

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email info.bmjopen@bmj.com

BMJ Open

Problems with Evidence Assessment in COVID-19 Health Policy Impact Evaluation: A systematic review of study design and evidence strength

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-053820
Article Type:	Original research
Date Submitted by the Author:	26-May-2021
Complete List of Authors:	<p>Haber, Noah A.; Stanford University Clarke-Deelder, Emma; Harvard University T H Chan School of Public Health, Department of Global Health and Population Feller, Avi; University of California Berkeley, Goldman School of Public Policy Smith, Emily; George Washington University School of Public Health and Health Services, Department of Global Health Salomon, Joshua ; Stanford University MacCormack-Gelles, Benjamin; Harvard University T H Chan School of Public Health, Department of Global Health and Population Stone, Elizabeth M.; Johns Hopkins University Bloomberg School of Public Health, Department of Health Policy and Management Bolster-Foucault, Clara; McGill University, Epidemiology, Biostatistics, and Occupational Health Daw, Jamie R.; Columbia University Mailman School of Public Health, Health Policy and Management Fry, Carrie E.; Vanderbilt University, Department of Health Policy Boyer, Christopher B.; Harvard University T H Chan School of Public Health, Department of Epidemiology Ben-Michael, Eli; University of California Berkeley, Department of Statistics Joyce, Caroline M.; McGill University, Epidemiology, Biostatistics, and Occupational Health Linas, Beth S.; Johns Hopkins University Bloomberg School of Public Health, Department of Epidemiology; MITRE Corp Schmid, Ian; Johns Hopkins University Bloomberg School of Public Health, Department of Mental Health Au, Eric; The University of Sydney, School of Public Health Wieten, Sarah; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Axfors, Cathrine; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Nguyen, Van; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Bilinski, Alyssa; Harvard University Graduate School of Arts and Sciences Hatfield, L; Harvard Medical School, Biostatistics Jarrett, Brooke; Johns Hopkins University, Epidemiology Griffin, Bath; RAND Corp Stuart, Elizabeth A; Johns Hopkins University Bloomberg School of Public</p>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

	Health, Department of Mental Health
Keywords:	COVID-19, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, STATISTICS & RESEARCH METHODS





I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

Problems with Evidence Assessment in COVID-19 Health Policy Impact Evaluation: A systematic review of study design and evidence strength

Noah A. Haber, ScD¹, Emma Clarke-Deelder, MPhil², Avi Feller, PhD³, Emily R. Smith, ScD⁴, Joshua Salomon, PhD⁵, Benjamin MacCormack-Gelles, MS², Elizabeth M. Stone, MS⁶, Clara Bolster-Foucault, MScPH⁷, Jamie R. Daw, PhD⁸, Laura A. Hatfield, PhD⁹, Carrie E. Fry, PhD¹⁰, Christopher B. Boyer, MPH¹¹, Eli Ben-Michael, PhD¹², Caroline M. Joyce, MPH⁷, Beth S. Linas, PhD, MHS^{13,14}, Ian Schmid, ScM¹⁵, Eric H. Au, MPH¹⁶, Sarah E. Wieten, PhD¹, Brooke A Jarrett, MSPH¹³, Cathrine Axfors, MD, PhD¹, Van Thu Nguyen, PhD¹, Beth Ann Griffin, PhD¹⁷, Alyssa Bilinski, MS¹⁸, Elizabeth A. Stuart, PhD¹⁵

1. Meta-Research Innovation Center at Stanford (METRICS), Stanford University, Stanford, CA, USA
2. Department of Global Health and Population, Harvard T. H. Chan School of Public Health, Boston, MA, USA
3. Goldman School of Public Policy, UC Berkeley, Berkeley, CA, USA
4. Department of Global Health, Milken Institute School of Public Health, George Washington University, Washington, D.C, USA
5. Center for Health Policy and Center for Primary Care and Outcomes Research, Stanford University, Stanford, CA, USA
6. Department of Health Policy and Management, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
7. Epidemiology, Biostatistics, and Occupational Health, McGill University, Montreal, Canada
8. Health Policy and Management, Columbia University Mailman School of Public Health, New York, NY, USA
9. Department of Health Care Policy, Harvard Medical School, Boston, MA, USA
10. Department of Health Policy, Vanderbilt University, Nashville, TN, USA
11. Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA
12. Institute for Quantitative Social Science, Harvard University, Cambridge, MA, USA
13. Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
14. Clinical Quality and Informatics, MITRE Corp, McLean, VA, USA
15. Department of Mental Health, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
16. School of Public Health, University of Sydney, Sydney, Australia
17. RAND Corporation, Arlington, VA, USA
18. Interfaculty Initiative in Health Policy, Harvard Graduate School of Arts and Sciences, Cambridge, MA, USA

Corresponding author:

Noah A. Haber, ScD

noahhaber@stanford.edu

Meta Research Innovation Center at Stanford University

Stanford University

1265 Welch Rd

Palo Alto, CA 94305

(650) 497-0811

Word count: 4,894

Abstract

Introduction: Assessing the impact of COVID-19 policy is critical for informing future policies. However, there are concerns about the overall strength of COVID-19 impact evaluation studies given the circumstances for evaluation and concerns about the publication environment. This study systematically reviewed the strength of evidence in the published COVID-19 policy impact evaluation literature.

Methods: We included studies that were primarily designed to estimate the quantitative impact of one or more implemented COVID-19 policies on direct SARS-CoV-2 and COVID-19 outcomes published on November 26 or earlier and screening. The review tool was based on previously developed and released review guidance for COVID-19 policy impact evaluation, assessing what impact evaluation method was used, graphical display of outcomes data, functional form for the outcomes, timing between policy and impact, concurrent changes to the outcomes, and an overall rating.

Results: After 102 articles were identified as potentially meeting inclusion criteria, we identified 36 published articles that evaluated the quantitative impact of COVID-19 policies on direct COVID-19 outcomes. The majority (n=23/36) of studies in our sample examined the impact of stay-at-home requirements. Reviewers found that only four of the 36 identified published and peer-reviewed health policy impact evaluation studies passed a set of key design checks for identifying the causal impact of policies on COVID-19 outcomes. The most common issues were lack of functional form justification and/or failure to eliminate concurrent changes to the outcomes.

Discussion: The reviewed literature directly evaluating the impact of COVID-19 policies largely failed to meet key design criteria for useful inference. This was largely driven by the circumstances under which policies were passed making it difficult to attribute changes in COVID-19 outcomes to particular policies. More reliable evidence review is needed to both identify and produce policy-actionable evidence, alongside the recognition that actionable evidence is often unlikely to be feasible.

Strengths and limitations

- This study is based on previously released review guidance for discerning and evaluating critical minimal methodological design aspects of the COVID-19 health policy impact evaluation.
- The review tool assesses critical aspects of study design grounded in impact evaluation methods that must be true for the papers to provide useful policy impact evaluation, including what type of impact evaluation method was used, graphical display of outcomes data, functional form for the outcomes, timing between policy and impact, concurrent changes to the outcomes, and an overall rating.
- This study used a consensus reviewer model with three reviewers in order to obtain replicable results for study strength ratings.
- While the vast majority of studies in our sample received low ratings for useful causal policy impact evaluation, they may make other contributions to the literature.
- Because our review tool was limited to a very narrow - albeit critical - set of items, weaknesses in other aspects not reviewed (e.g. data quality or other aspects of statistical inference) may further weaken studies that were found to meet our criteria.

Introduction

Policy decisions to mitigate the impact of COVID-19 on morbidity and mortality are some of the most important issues policymakers have had to make since January 2020. Decisions regarding which policies are enacted depend in part on the evidence base for those policies, including understanding what impact past policies had on COVID-19 outcomes.[1,2] Unfortunately, there are substantial concerns that much of the existing literature may be methodologically flawed, which could render its conclusions unreliable for informing policy. The combination of circumstances being difficult for strong impact evaluation, the importance of the topic, and concerns over the publication environment may lead to the proliferation of low strength studies.

High-quality causal evidence requires a combination of rigorous methods, clear reporting, appropriate caveats, and the appropriate circumstances for the methods used.[3–6] Rigorous evidence is difficult in the best of circumstances, and the circumstances for evaluating non-pharmaceutical intervention (NPI) policy effects on COVID-19 are particularly challenging.[5] The global pandemic has yielded a combination of a large number of concurrent policy and non-policy changes, complex infectious disease dynamics, and unclear timing between policy implementation and impact; all of this makes isolating the causal impact of any particular policy or policies exceedingly difficult.[7]

1
2
3 The scientific literature on COVID-19 is exceptionally large and fast growing. Scientists
4 published more than 100,000 papers related to COVID-19 in 2020.[8] There is some general
5 concern that the volume and speed at which this work has been produced may result in a
6 literature that is overall low quality and unreliable.[9–15]
7

8
9 Given the importance of the topic, it is critical that decision-makers are able to understand what
10 is known and knowable from observational data in COVID-19 policy,[5,16] as well as what is
11 unknown and/or unknowable.
12

13
14 Motivated by concerns about the methodological strength of COVID-19 policy evaluations, we
15 set out to review the literature using a set of methodological design checks tailored to common
16 policy impact evaluation methods. Our primary objective was to evaluate each paper for
17 methodological strength and reporting, based on pre-existing review guidance developed for
18 this purpose.[17] As a secondary objective, we also studied our own process: examining the
19 consistency, ease of use, and clarity of this review guidance.
20
21

22
23 This protocol differs in several ways from more traditional systematic review protocols given the
24 atypical objectives and scope of the systematic review. First, this is a systematic review of
25 methodological strength of evidence for a given literature as opposed to a review summary of
26 the evidence of a particular topic. As such, we do not summarize and attempt to combine the
27 results for any of the literature. Second, rather than being a comprehensive review of every
28 possible aspect of what might be considered “quality,” this is a review of targeted critical design
29 features for useful inference for COVID-19 policy impact evaluation and methods. It is designed
30 to be a set of broad criteria for minimal plausibility of useful causal inference, where each of the
31 criteria is necessary but not sufficient for strong design. Issues in other domains (data, details of
32 the design, statistics, etc) further reduce overall usefulness and quality, and thorough review in
33 those domains is needed for any studies passing our basic minimal criteria. Third, because the
34 scope relies on guided, but difficult and subjective assessments of methodological
35 appropriateness, we utilize a discussion-based consensus process to arrive at consistent and
36 replicable results, rather than a more common model with two independent reviewers with
37 conflict resolution. The independent review serves primarily as a starting point for discussion,
38 but is neither designed nor expected to be a strong indicator of the overall consensus ratings of
39 the group.
40
41
42
43
44

45 Methods

46 Overview

47
48
49 This protocol and study was written and developed following the release of the review guidance
50 written by the author team in September 2020 on which the review tool is based. The protocol
51 for this study was pre-registered on OSF.io in November 2020 following PRISMA
52 guidelines.[18,19] Deviations from the original protocol are discussed in Appendix 1, and
53 consisted largely of language clarifications and error corrections for both the inclusion criteria
54
55
56
57
58
59

and review tool, an increase in the number of reviewers per fully reviewed article from two to three, and simplification of the statistical methods used to assess the data.

This systematic review of the strength of evidence took place in three phases: search, screening, and full review.

Eligibility criteria

The following eligibility criteria were used to determine the papers to include:

- The primary topic of the article must be evaluating one or more individual COVID-19 policies on direct COVID-19 outcomes
 - The primary exposure(s) must be a policy, defined as a government-issued order at any government level to address a directly COVID-19-related outcome (e.g., mask requirements, travel restrictions, etc).
 - COVID-19 outcomes may include cases detected, mortality, number of tests taken, test positivity rates, Rt, etc.
 - This may NOT include indirect impacts of COVID-19 on things such as income, childcare, trust in science, etc.
- The primary outcome being examined must be a COVID-19-specific outcome, as above.
- The study must be designed as an impact evaluation study from primary data (i.e., not primarily a predictive or simulation model or meta-analysis).
- The study must be peer reviewed, and published in a peer-reviewed journal indexed by PubMed.
- The study must have the title and abstract available via PubMed at the time of the study start date (November 26).
- The study must be written in English.

These eligibility criteria were designed to identify the literature primarily concerning the quantitative impact of one or more implemented COVID-19 policies on COVID-19 outcomes. Studies in which impact evaluation was secondary to another analysis (such as a hypothetical projection model) were eliminated because they were less relevant to our objectives and/or may not contain sufficient information for evaluation. Categories for types of policies were from the Oxford COVID-19 Government Response Tracker.[20]

Reviewer recruitment, training, and communication

Reviewers were recruited through personal contacts and postings on online media. All reviewers had experience in systematic review, quantitative causal inference, epidemiology, econometrics, public health, methods evaluation, or policy review. All reviewers participated in two meetings in which the procedures and the review tool were demonstrated. Screening reviewers participated in an additional meeting specific to the screening process. Throughout the main review process, reviewers communicated with the administrators and each other through Slack for any additional clarifications, questions, corrections, and procedures. The main administrator (NH), who was also a reviewer, was available to answer general questions and make clarifications, but did not answer questions specific to any given article.

Review phases and procedures

Search strategy

The search terms combined four Boolean-based search terms: a) COVID-19 research, b) regional government units (e.g., country, state, county, and specific country, state, or province, etc.), c) policy or policies, and d) impact or effect. The full search terms are available in Appendix 2.

Information Sources

The search was limited to published articles in peer-reviewed journals. This was largely to attempt to identify literature that was high quality, relevant, prominent, and most applicable to the review guidance. PubMed was chosen as the exclusive indexing source due to the prevalence and prominence of policy impact studies in the health and medical field. Preprints were excluded to limit the volume of studies to be screened and to ensure each had met the standards for publication through peer review. The search was conducted on November 26, 2020.

Study Selection

Eight reviewers screened the title and abstract of each article for the inclusion criteria. Two reviewers were randomly selected to screen each article for acceptance/rejection. In the case of a dispute, a third randomly selected reviewer decided on acceptance/rejection. Training consisted of a one-hour instruction meeting, a review of the first 50 items on each reviewers' list of assigned articles, and a brief asynchronous online discussion before conducting the full review.

Full article review

The full article review consisted of two sub-phases: the independent primary review phase, and a group consensus phase. The independent review phase was designed primarily for the purpose of supporting and facilitating useful discussion in the consensus discussion, rather than as high stakes definitive review data on its own. The consensus process was considered the primary way in which review data would be generated, rather than synthesis from the independent reviews.

Each article was randomly assigned to three of the 23 reviewers in our review pool. Each reviewer independently reviewed each article on their list, first for whether the study met the eligibility criteria, then responding to methods identification and guided strength of evidence questions using the review tool, as described below. Reviewers were able to recuse themselves for any reason, in which case another reviewer was randomly selected. Once all three reviewers had reviewed a given article, all articles that weren't unanimously determined to not meet the inclusion criteria underwent a consensus process.

1
2
3 During the consensus round, the three reviewers were given all three primary reviews for
4 reference, and were tasked with generating a consensus opinion among the group. One
5 randomly selected reviewer was tasked to act as the arbitrator. The arbitrator's primary task was
6 facilitative useful discussion and for moving the group toward establishing a consensus that
7 represented the collective subjective assessments of the group. If consensus could not be
8 reached, a fourth randomly selected reviewer was brought into the discussion to help resolve
9 disputes.
10
11

12 13 Review tool for data collection 14

15 This review tool and data collection process was an operationalized and lightly adapted version
16 of the COVID-19 health policy impact evaluation review guidance literature, written by the lead
17 authors of this study. The main adaptation was removing references to the COVID-19 literature.
18 All reviewers were instructed to read and refer to this guidance document to guide their
19 assessments. The full guidance manuscript contains additional explanation and rationale for all
20 parts of this review and the tool, and is available both in the adapted form as was provided to
21 the reviewers in a supplementary file "CHSPER review guidance refs removed.pdf" and in an
22 updated version in Haber et al., 2020.[17] The full review tool is attached as supplementary file
23 "review tool final.pdf".
24
25
26

27 The review tool consisted of two main parts: methods design categorization and full review. The
28 review tool and guidance categorizes policy causal inference designs based on the structure of
29 their assumed counterfactual. This is assessed through identifying the data structure and
30 comparison(s) being made. There are two main items for this determination: the number of pre-
31 period time points (if any) used to assess pre-policy outcome trends, and whether or not policy
32 regions were compared with non-policy regions. These, and other supporting questions, broadly
33 allowed categorization of methods into cross-sectional, pre/post, interrupted time series (ITS),
34 difference-in-differences (DiD), comparative interrupted time-series (CITS), (randomized) trials,
35 or other. Given that most papers have several analyses, reviewers were asked to focus
36 exclusively on the impact evaluation analysis that was used as the primary support for the main
37 conclusion of the article.
38
39
40
41

42 Studies categorized as cross-sectional, pre/post, randomized controlled trial designs, and other
43 were set aside for no further review for the purposes of this research. Cross-sectional and pre-
44 post designs were considered inappropriate for policy causal inference for COVID-19 due
45 largely to inability to account for a large number of potential issues, including confounding,
46 epidemic trends, and selection biases. Randomized controlled trials were assumed to broadly
47 meet key design checks. Studies categorized as "other" received no further review, as the
48 review guidance would be unable to assess them. Additional justification and explanation for
49 this decision is available in the review guidance.
50
51
52

53 For the methods receiving full review (ITS, DiD, and CITS), reviewers were asked to identify
54 potential issues and give a category-specific rating. The specific study designs triggered sub-
55
56
57
58
59
60

1
2
3 questions and/or slightly altered the language of the questions being asked, but all three of the
4 methods design categories shared these four key questions:
5

- 6 ● Graphical presentation: “Does the analysis provide graphical representation of the
7 outcome over time?”
- 8 ● Functional form: “Is the functional form of the counterfactual (e.g., linear) well-justified
9 and appropriate?”
- 10 ● Timing of policy impact: “Is the date or time threshold set to the appropriate date or time
11 (e.g., is there lag between the intervention and outcome)?”
- 12 ● Concurrent changes: “Is this policy the only uncontrolled or unadjusted-for way in which
13 the outcome could have changed during the measurement period [differently for policy
14 and non-policy regions]?”
15
16
17

18
19 For each of the four key questions, reviewers were given the option to select “No,” “Mostly no,”
20 “Mostly yes,” and “Yes” with justification text requested for all answers other than “Yes.” Each
21 question had additional prompts as guidance, and with much more detail provided in the full
22 guidance document.
23

24
25 The criteria were designed to establish minimal plausibility of useful and actionable evidence,
26 rather than certification of high quality. Graphical representation is included here primarily as a
27 key way to assess the plausibility and justification of key model assumptions, rather than being
28 necessary for validity by itself. For example, rather than having the “right” functional form or lag
29 structure, the review guidance asks whether the functional form and lags is discussed at all and
30 (if discussed) reasonable.
31

32
33 These four questions were selected and designed being critical to evaluating strength of study
34 design for policy impact evaluation in general, direct relevance for COVID-19 policy, feasibility
35 for use in guided review. These questions are designed as minimal and key criteria for plausibly
36 useful impact evaluation design for COVID-19 policy impact evaluation, rather than as a
37 comprehensive tool assessing overall quality. Thorough review of data quality, statistical
38 validity, and other issues are also critical points of potential weakness in study designs, and
39 would be needed in addition to these criteria, if these key design criteria are met. A thorough
40 justification and explanation of how and why these questions were selected is available in the
41 provided guidance document and in Haber et al., 2020.[17]
42
43
44

45 Finally, reviewers were asked a summary question:
46

- 47 ● Overall: “Do you believe that the design is appropriate for identifying the policy impact(s)
48 of interest?”
49
50

51 Reviewers were asked to consider the scale of this question to be both independent/not relative
52 to any other papers, and that any one substantial issue with the study design could render it a
53 “No” or “Mostly no.” Reviewers were asked to follow the guidance and their previous answers,
54 allowing for their own weighting of how important each issue was to the final result. A study
55
56
57
58
59

1
2
3 could be excellent on all dimensions except for one, and that one dimension could render it
4 inappropriate for causal inference. As such, in addition to the overall rating question, we also
5 generated a “weakest link” metric for overall assessment, representing the lowest rating among
6 the four key questions (graphical representation, functional form, timing of policy impact, and
7 concurrent changes). A “mostly yes” or “yes” is considered a passing rating, indicating that the
8 study was not found to be inappropriate on the specific dimension of interest.
9
10

11 A “yes” rating does not necessarily indicate that the study is strongly designed, conducted, or is
12 useful; it only means that it passes a series of key design checks for policy impact evaluation
13 and should be considered for further evaluation. The papers may contain any number of other
14 issues that were not reviewed (e.g., statistical issues, inappropriate comparisons,
15 generalizability, etc.). As such, this should only be considered an initial assessment of
16 plausibility that the study is well-designed, rather than confirmation that it is appropriate and
17 applicable.
18
19
20

21 The full review tool is available in the supplementary materials.
22
23

24 Statistical analysis

25
26 Statistics provided are nearly exclusively counts and percentages of the final dataset. Analyses
27 and graphics were performed in R.[21] Inter-rater reliability was assessed using Krippendorff's
28 alpha[22] using the IRR package.[23] Relative risks were estimated using the epitools
29 package.[24]
30
31

32 Citation counts for accepted articles were obtained through Google Scholar on January 11,
33 2021.[25] Journal impact factors were obtained from the 2019 Journal Citation Reports.[26]
34
35

36 Data and code

37
38 Data, code, the review tool, and the review guidance are stored and available here:
39 <https://osf.io/9xmke/files/>. The dataset includes full results from the search and screening and all
40 review tool responses from reviewers during the full review phase.
41
42
43

44 Patient and Public Involvement Statement

45
46 Patients or public stakeholders were not consulted in the design or conduct of this systematic
47 evaluation.
48
49

50 Results

51 Search and screening

52
53
54
55
56 *Figure 1: PRISMA diagram of systematic review process*
57
58
59
60

As shown in Figure 1, after search and screening of titles and abstracts, 102 articles were identified as likely or potentially meeting our inclusion criteria. Of those 102 articles, 36 studies met inclusion after independent review and deliberation in the consensus process. The most common reasons for rejection at this stage were that the study did not measure the quantitative direct impact of specific policies and/or that such an impact was not the main purpose of the study. Many of these studies implied that they measured policy impact in the abstract or introduction, but instead measured correlations with secondary outcomes (e.g., the effect of movement reductions, which are influenced by policy) and/or performed cursory policy impact evaluation secondary to projection modelling efforts.

Descriptive statistics

Figure 2: Descriptive sample statistics (n=36)

Publication information from our sample is shown in Figure 2. The articles in our sample were generally published in journals with high impact factors (median impact factor: 3.6) and have already been cited in the academic literature (median citation count: 5, on 1/11/21). The most commonly evaluated policy type was stay at home requirements (64% n=23/36). Reviewers noted that many articles referenced “lockdowns,” but did not define the specific policies to which this referred.

Reviewers most commonly selected interrupted time-series (39% n=14/36) as the methods design, followed by difference-in-differences (9% n=9/36) and pre-post (8% n=8/36). There were no randomized controlled trials of COVID-19 health policies identified (0% n=0/36), nor were any studies identified that reviewers could not categorize on the review guidance (0% n=0/36).

Table 1: Summary of articles reviewed and reviewer ratings for key and overall questions

Category ratings order		Legend for color coded ratings						
Graphical presentation	Timing of policy impact	N/A	Unclear	No*	No**	Mostly no	Mostly yes	Yes
Functional form	Concurrent changes	method determined to me inappropriate by: * guidance (cross sectional or pre/post) or ** reviewer consensus						
Citation	Title	Journal	Publication date	Methods design	Category ratings	Overall rating		
Cobb and Seale, 2020[27]	Examining the effect of social distancing on the compound growth rate of COVID-19 at the county level (United States) using statistical analyses and a random forest machine learning model.	Public Health	4/28/2020	Pre/post	Mostly no	No*		
Lyu and Wehby, 2020a[28]	Comparison of Estimated Rates of Coronavirus Disease 2019 (COVID-19) in Border Counties in Iowa Without a Stay-at-Home Order and Border Counties in Illinois With a Stay-at-Home Order.	JAMA Network Open	5/1/2020	Difference-in-differences	Mostly no	Mostly yes		
Tam et al., 2020[29]	Effect of mitigation measures on the spreading of COVID-19 in hard-hit states in the U.S.	PloS One	5/1/2020	Interrupted time-series	Mostly no	Mostly yes		
Courtemanche et al., 2020[30]	Strong Social Distancing Measures In The United States Reduced The COVID-19 Growth Rate.	Health Affairs	5/14/2020	Difference-in-differences	Mostly yes	Yes		
Crokidakis, 2020[31]	COVID-19 spreading in Rio de Janeiro, Brazil: Do the policies of social isolation really work?	Chaos, Solitons, and Fractals	5/23/2020	Interrupted time-series	Mostly yes	Yes		
Hyafil and Moríña, 2020[32]	Analysis of the impact of lockdown on the reproduction number of the SARS-Cov-2 in Spain.	Gaceta Aanitaria	5/23/2020	Pre/post	Mostly yes	No*		

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Castillo, et al., 2020[33]	The effect of state-level stay-at-home orders on COVID-19 infection rates.	American Journal of Infection Control	5/24/2020	Pre/post		
Alfano and Ercolano, 2020[34]	The Efficacy of Lockdown Against COVID-19: A Cross-Country Panel Analysis.	Applied Health Economics and Health Policy	6/3/2020	Difference-in-differences		
Lyu and Wehby, 2020b[35]	Community Use Of Face Masks And COVID-19: Evidence From A Natural Experiment Of State Mandates In The US.	Health Affairs	6/16/2020	Difference-in-differences		
Zhang, et al., 2020[36]	Identifying airborne transmission as the dominant route for the spread of COVID-19.	PNAS	6/30/2020	Interrupted time-series		
Xu et al., 2020[37]	Associations of Stay-at-Home Order and Face-Masking Recommendation with Trends in Daily New Cases and Deaths of Laboratory-Confirmed COVID-19 in the United States.	Exploratory research and hypothesis in medicine	7/8/2020	Interrupted time-series		
Lyu and Wehby, 2020c[38]	Shelter-In-Place Orders Reduced COVID-19 Mortality And Reduced The Rate Of Growth In Hospitalizations.	Health Affairs	7/9/2020	Difference-in-differences		
Wagner, et al., 2020[39]	Social distancing merely stabilized COVID-19 in the US.	Stat (International Statistical Institute)	7/13/2020	Interrupted time-series		
Di Bari et al., 2020[40]	Extensive Testing May Reduce COVID-19 Mortality: A Lesson From Northern Italy.	Frontiers in Medicine	7/14/2020	Comparative interrupted time-series		
Islam et al., 2020[41]	Physical distancing interventions and incidence of coronavirus disease 2019: natural experiment in 149 countries.	BMJ (Clinical research ed.)	7/15/2020	Interrupted time-series		
Wong et al., 2020[42]	Impact of National Containment Measures on Decelerating the Increase in Daily New Cases of COVID-19 in 54 Countries and 4 Epicenters of the Pandemic: Comparative Observational Study.	Journal of Medical Internet Research	7/22/2020	Pre/post		
Liang et al., 2020[43]	Effects of policies and containment measures on control of COVID-19 epidemic in Chongqing.	World Journal of Clinical Cases	7/26/2020	Pre/post		
Banerjee and Nayak, 2020[44]	U.S. county level analysis to determine if social distancing slowed the spread of COVID-19.	Pan American Journal of Public Health	7/31/2020	Difference-in-differences		
Dave et al., 2020a[45]	When Do Shelter-in-Place Orders Fight COVID-19 Best? Policy Heterogeneity Across States and Adoption Time.	Economic inquiry	8/3/2020	Difference-in-differences		
Hsiang et al., 2020[46]	The effect of large-scale anti-contagion policies on the COVID-19 pandemic.	Nature	8/22/2020	Interrupted time-series		
Lim et al., 2020[47]	Revealing regional disparities in the transmission potential of SARS-CoV-2 from interventions in Southeast Asia.	Proceedings. Biological sciences	8/26/2020	Interrupted time-series		
Arshed et al., 2020[48]	Empirical assessment of government policies and flattening of the COVID19 curve.	Journal of Public Affairs	8/27/2020	Cross-sectional analysis		
Wang et al., 2020[49]	Fangcang shelter hospitals are a One Health approach for responding to the COVID-19 outbreak in Wuhan, China.	One Health	8/29/2020	Interrupted time-series		
Kang et al., 2020[50]	The Effects of Border Shutdowns on the Spread of COVID-19.	Journal of Preventive Medicine and Public Health	8/30/2020	Comparative interrupted time-series		
Auger et al., 2020[51]	Association Between Statewide School Closure and COVID-19 Incidence and Mortality in the US.	JAMA	9/1/2020	Interrupted time-series		
Santamaria et al., 2020[52]	COVID-19 effective reproduction number dropped during Spain's nationwide dropdown, then spiked at lower-incidence regions.	The Science of the Total Environment	9/9/2020	Interrupted time-series		
Bennett, 2020[53]	All things equal? Heterogeneity in policy effectiveness against COVID-19 spread in Chile.	World Development	9/24/2020	Comparative interrupted time-series		
Yang et al., 2020[54]	Lessons Learnt from China: National Multidisciplinary Healthcare Assistance.	Risk Management and Healthcare Policy	9/30/2020	Difference-in-differences		
Padalabalanarayanan et al., 2020[55]	Association of State Stay-at-Home Orders and State-Level African American Population With COVID-19 Case Rates.	JAMA Network Open	10/1/2020	Comparative interrupted time-series		
Edelstein et al., 2020[56]	SARS-CoV-2 infection in London, England: changes to community point prevalence around lockdown time, March-May 2020.	Journal of Epidemiology and Community Health	10/1/2020	Pre/post		
Tsai et al., 2020[57]	COVID-19 transmission in the U.S. before vs. after relaxation of statewide social distancing measures.	Clinical Infectious Diseases	10/3/2020	Interrupted time-series		
Singh et al., 2020[58]	Public health interventions slowed but did not halt the spread of COVID-19 in India.	Transboundary and Emerging Diseases	10/4/2020	Pre/post		
Gallaway et al., 2020[59]	Trends in COVID-19 Incidence After Implementation of Mitigation Measures - Arizona, January 22-August 7, 2020.	Morbidity and Mortality Weekly Report	10/9/2020	Pre/post		
Castex et al., 2020[60]	COVID-19: The impact of social distancing policies, cross-country analysis.	Economics of Disasters and Climate Change	10/15/2020	Interrupted time-series		
Silva et al., 2020[61]	The effect of lockdown on the COVID-19 epidemic in Brazil: evidence from an interrupted time series design.	Cadernos de Saude Publica	10/19/2020	Interrupted time-series		
Dave et al., 2020b[62]	Were Urban Cowboys Enough to Control COVID-19? Local Shelter-in-Place Orders and Coronavirus Case Growth.	Journal of Urban Economics	11/6/2020	Difference-in-differences		

1
2
3
4
5
6 The identified articles and selected review results are summarized in Table 1.
7

8 9 Strength of methods assessment

10 *Figure 3: Main consensus results summary for key and overall questions*

11
12
13 Graphical representation of the outcome over time was relatively well-rated in our sample, with
14 74% (n=20/27) studies being given a “mostly yes” or “yes” rating for appropriateness. Reasons
15 cited for non-“yes” ratings included a lack of graphical representation of the data, alternative
16 scales used, and not showing the dates of policy implementation, as shown in Figure 3.
17

18
19 Functional form issues appear to have presented a major issue in these studies, with only 19%
20 receiving a “mostly yes” or “yes” rating, 78% (n=21/27) receiving a “no” rating, and 4% (n=1/27)
21 “unclear.” There were two common themes in this category: studies generally using scales that
22 were broadly considered inappropriate for infectious disease outcomes (e.g., linear counts),
23 and/or studies lacking stated justification for the scale used. Reviewers also noted disconnects
24 between clear curvature in the outcomes in the graphical representations and the analysis
25 models and outcome scales used (e.g., linear). In one case, reviewers could not identify the
26 functional form actually used in analysis.
27
28

29
30 Reviewers broadly found that these studies dealt with timing of policy impact (e.g., lags between
31 policy implementation and expected impact) relatively well, with 70% (n=19/27) rated “yes” or
32 “mostly yes.” Reasons for non-“yes” responses included not adjusting for lags and a lack of
33 justification for the specific lags used.
34

35
36 Concurrent changes were found to be a major issue in these studies, with only 11% (n=3/27)
37 studies receiving passing ratings (“yes” or “mostly yes”) with regard to uncontrolled concurrent
38 changes to the outcomes. Reviewers nearly ubiquitously noted that the articles failed to account
39 for the impact of other policies that could have impacted COVID-19 outcomes concurrent with
40 the policies of interest. Other issues cited were largely related to non-policy-induced behavioral
41 and societal changes.
42
43

44
45 When reviewers were asked if sensitivity analyses had been performed on key assumptions and
46 parameters, about half (56% n=15/27) answered “mostly yes” or “yes.” The most common
47 reason for non-“yes” ratings was that, while sensitivity analyses were performed, they did not
48 address the most substantial assumptions and issues.
49

50
51 Overall, reviewers rated only four studies (11%, n=4/36,) as being plausibly appropriate (“mostly
52 yes” or “yes”) for identifying the impact of specific policies on COVID-19 outcomes, as shown in
53 Figure 3. 25% (n=9/36) were automatically categorized as being inappropriate due to being
54 either cross-sectional or pre/post in design, 33% (n=12/36) of studies were given a “no” rating
55 for appropriateness, 31% “mostly no” (n=11/36), 8% “mostly yes” (n=3/36), and 3% “yes”
56
57
58
59

(n=1/36). The most common reason cited for non-“yes” overall ratings was failure to account for concurrent changes (particularly policy and societal changes).

Figure 4: Comparison of independent reviews, weakest link, and direct consensus review

As shown in Figure 4, the consensus overall proportion passing (“mostly yes” or “yes”) was a quarter of what it was from the initial independent reviews. 45% (n=34/75) of studies were rated as “yes” or “mostly yes” in the initial independent review, as compared to 11% (n=4/36) in the consensus round (RR 0.25, 95%CI 0.09:0.64). The issues identified and discussed in combination during consensus discussions, as well as additional clarity on the review process, resulted in reduced overall confidence in the findings. Increased clarity on the review guidance with experience and time may also have reduced these ratings further.

The large majority of studies had at least one “no” or “unclear” rating in one of the four categories (74% n=20/27), with only one study whose lowest rating was a “mostly yes,” no studies rated “yes” in all four categories. Only one study was found to pass design criteria in all four key questions categories, as shown in the “weakest link” column in Figure 4.

Review process assessment

During independent review, all three reviewers independently came to the same conclusions on the main methods design category for 33% (n=12/36) articles, two out of the three reviewers agreed for 44% (n=16/36) articles, and none of the reviewers agreed in 22% (n=8/36) cases. One major contributor to these discrepancies were the 31% (n=11/36) cases where one or more reviewers marked the study as not meeting eligibility criteria, 64% (n=7/11) of which the other two reviewers agreed on the methods design category.

Inter-rater reliability of the primary independent reviews was relatively low across the board for the key questions. For the overall scores, Krippendorff’s alpha was only 0.16 due to widely varying opinions between raters. The four key categorical questions had slightly better inter-rater reliability than the overall question, with Krippendorff’s alphas of 0.59 for graphical representation, 0.34 for functional form, 0.44 for timing of policy impact, and 0.15 for concurrent changes, respectively.

The consensus rating for overall strength was equal to the lowest rating among the independent reviews in 78% (n=21/27) of cases, and only one higher than the lowest in the remaining 22% (n=6/27). This strongly suggests that the multiple reviewer review, discussion, and consensus process better identifies issues than independent review alone. Differences in initial opinions between reviewers may be attributable to any number of factors, including true differences in opinion, misunderstandings/learning about the review tool and process, and expected reliance on the consensus process. Notably, there were two cases for which reviewers requested an additional fourth reviewer to help resolve standing issues for which the reviewers felt they were unable to come to consensus.

1
2
3
4 The most consistent point of feedback from reviewers was the value of having a three reviewer
5 team with whom to discuss and deliberate, rather than two as initially planned. This was
6 reported to help catch a larger number of issues and clarify both the papers and the
7 interpretation of the review tool questions. Reviewers also expressed that one of the most
8 difficult parts of this process was assessing the inclusion criteria, some of the implications of
9 which are discussed below.
10
11
12

13 Discussion

14
15
16 This systematic review of evidence strength found that only four (or only one by a stricter
17 standard) of the 36 identified published and peer-reviewed health policy impact evaluation
18 studies passed a set of key checks for identifying the causal impact of policies on COVID-19
19 outcomes. Because this systematic review examined a limited set of key study design features
20 and did not address more detailed aspects of study design, statistical issues, generalizability,
21 and any number of other issues, this result may be considered an upper bound on the overall
22 strength of evidence within this sample. Two major problems are nearly ubiquitous throughout
23 this literature: failure to isolate the impact of the policy(s) of interest from other changes that
24 were occurring contemporaneously, and failure to appropriately address the functional form of
25 infectious disease outcomes in a population setting. While policy decisions are being made on
26 the backs of high impact-factor papers, we find that the citation-based metrics do not
27 correspond to “high-quality” research as Yin et al., 2021 claims.[63] Similar to other areas in the
28 COVID-19 literature,[64] we found the current literature directly evaluating the impact of COVID-
29 19 policies largely fails to meet key design criteria for useful inference.
30
31
32
33

34 The framework for the review tool is based on the requirements and assumptions built into
35 policy evaluation methods. Quasi-experimental methods rely critically on the scenarios in which
36 the data are generated. These assumptions and the circumstances in which they are plausible
37 are well-documented and understood,[2,4–6,17,65] including one paper discussing application
38 of difference-in-differences methods specifically for COVID-19 health policy, released in May
39 2020.[5] While “no uncontrolled concurrent changes” is a difficult bar to clear, that bar is
40 fundamental to inference using these methods.
41
42
43

44 The circumstances of isolating the impact of policies in COVID-19 - including large numbers of
45 policies, infectious disease dynamics, and massive changes to social behaviors - make those
46 already difficult fundamental assumptions broadly much less likely to be met. Some of the
47 studies in our sample were nearly the best feasible studies that could be done given the
48 circumstances, but the best that can be done often yields little useful inference. The relative
49 paucity of strong studies does not in any way imply a lack of impact of those policies; only that
50 we lack the circumstances to have evaluated their effects.
51
52
53

54 Because the studies estimating the harms of policies share the same fundamental
55 circumstances, the evidence of COVID-19 policy harms is likely to be of similarly poor strength.
56 Identifying the effects of many of these policies, particularly for the spring of 2020, is likely to be
57
58
59
60

1
2
3 unknown and perhaps unknowable. However, there remains additional opportunities with more
4 favorable circumstances, such as measuring overall impact of NPIs as bundles, rather than
5 individual policies. Similarly, studies estimating the impact of re-opening policies or policy
6 cancellation are likely to have fewer concurrent changes to address.
7

8
9 The review process itself demonstrates how guided and targeted peer review can efficiently
10 evaluate studies in ways that the traditional peer review systems do not. The studies in our
11 sample had passed the full peer review process, were published in largely high-profile journals,
12 and are highly cited, but contained substantial flaws that rendered their inference utility
13 questionable. The relatively small number of studies included, as compared to the size of the
14 literature concerning itself with COVID-19 policy, may suggest that there was relative restraint
15 from journal editors and reviewers for publishing these types of studies. The large number of
16 models, but relatively small number of primary evaluation analyses is consistent with other
17 areas of COVID-19.[66,67] At minimum, the flaws and limitations in their inference could have
18 been communicated at the time of publication, when they are needed most. In other cases, it is
19 plausible that many of these studies would not have been published had a more thorough or
20 better targeted methodological review been performed.
21
22
23

24
25 This systematic review of evidence strength has limitations. The tool itself was limited to a very
26 narrow - albeit critical - set of items. Low ratings in our study should not be interpreted as being
27 overall poor studies, as they may make other contributions to the literature that we did not
28 evaluate. While the guidance provided a well-structured framework and our reviewer pool was
29 well-qualified, strength of evidence review is inherently subjective. It is plausible and likely that
30 other sets of reviewers would come to different conclusions. However, the consensus process
31 was designed with these issues subjectivity in mind, and demonstrates the value of consensus
32 processes for overcoming hurdles with subjective and difficult decisions.
33
34
35

36 Most importantly, this review does not cover all policy inference in the scientific literature. One
37 large literature from which there may be COVID-19 policy evaluation otherwise meeting our
38 inclusion criteria are pre-prints. Many pre-prints would likely fare well in our review process.
39 Higher strength papers often require more time for review and publication, and many high
40 quality papers may be in the publication pipeline now. Second, this review excluded studies that
41 had a quantitative impact evaluation as a secondary part of the study (e.g., to estimate
42 parameters for microsimulation or disease modeling). Not only are these assessments not the
43 primary purpose of those studies, they also typically lack the detail requisite to make a critical
44 assessment of the study design and methods used. Third, the review does not include policy
45 inference studies that do not measure the impact of a specific policy. For instance, there are
46 studies that estimate the impact of reduced mobility on COVID-19 outcomes but do not attribute
47 the reduced mobility to any specific policy change. Finally, a considerable number of studies
48 that present analyses of COVID-19 outcomes to inform policy are excluded because they do not
49 present a quantitative estimate of specific policies' treatment effects.
50
51
52

53
54 While COVID-19 policy is one of the most important problems of our time, the circumstances
55 under which those policies were enacted severely hamper our ability to study and understand
56
57
58
59

1
2
3 their effects. Claimed conclusions are only as valuable as the methods by which they are
4 produced. Replicable, rigorous, intense, and methodologically guided review is needed to both
5 communicate our limitations and make more useful inference. Weak, unreliable, and
6 overconfident evidence leads to poor decisions and undermines trust in science.[15,68] In the
7 case of COVID-19 health policy, a frank appraisal of the strength of the studies on which
8 policies are based is needed, alongside the understanding that we often must make decisions
9 when strong evidence is not feasible.[69]
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figures

Figure 1: PRISMA diagram of systematic review process

Caption: This chart shows the PRISMA diagram for the process of screening the literature from search to the full review phase.

Figure 2: Descriptive sample statistics (n=36)

Caption: This chart shows descriptive statistics of the 36 studies entered into our systematic evidence review.

Figure 3: Main consensus results summary for key and overall questions

Caption: This chart shows the final overall ratings (left) and the key design question ratings for the consensus review of the 36 included studies, answering the degree to which the articles met the given key design question criteria. The key design question ratings were not asked for the nine included articles which selected methods assumed by the guidance to be non-appropriate. The question prompt in the figure is shortened for clarity, where the full prompt for each key question is available in the Methods section.

Figure 4: Comparison of independent reviews, weakest link, and direct consensus review

Caption: This chart shows the final overall ratings by three different possible metrics. The first column contains all of the independent review ratings for the 27 studies which were eventually included in our sample, noting that reviewers who either selected them as not meeting inclusion criteria or selected a method that didn't receive the full review did not contribute. The middle column contains the final consensus reviews among the 27 articles which received full review. The last column contains the weakest link rating, as described in the methods section. The question prompt in the figure is shortened for clarity, where the full prompt for each key question is available in the Methods section.

Works cited

- 1 Fischhoff B. Making Decisions in a COVID-19 World. *JAMA* 2020;324:139.
doi:10.1001/jama.2020.10178
- 2 COVID-19 Statistics, Policy modeling, and Epidemiology Collective. Defining high-value information for COVID-19 decision-making. *Health Policy* 2020.
doi:10.1101/2020.04.06.20052506
- 3 Hernán MA, Robins JM. *Causal Inference: What If*. Boca Raton: : Chapman & Hall/CRC
- 4 Angrist J, Pischke J-S. *Mostly Harmless Econometrics: An Empiricist's Companion*. 1st ed. Princeton University Press 2009. <https://EconPapers.repec.org/RePEc:pup:pbooks:8769>
- 5 Goodman-Bacon A, Marcus J. Using Difference-in-Differences to Identify Causal Effects of COVID-19 Policies. SSRN Journal Published Online First: 2020. doi:10.2139/ssrn.3603970
- 6 Bärnighausen T, Oldenburg C, Tugwell P, et al. Quasi-experimental study designs series—paper 7: assessing the assumptions. *Journal of Clinical Epidemiology* 2017;89:53–66.
doi:10.1016/j.jclinepi.2017.02.017
- 7 Haushofer J, Metcalf CJE. Which interventions work best in a pandemic? *Science* 2020;368:1063–5. doi:10.1126/science.abb6144
- 8 Else H. How a torrent of COVID science changed research publishing — in seven charts. *Nature* 2020;588:553–553. doi:10.1038/d41586-020-03564-y
- 9 Palayew A, Norgaard O, Safreed-Harmon K, et al. Pandemic publishing poses a new COVID-19 challenge. *Nat Hum Behav* 2020;4:666–9. doi:10.1038/s41562-020-0911-0
- 10 Bagdasarian N, Cross GB, Fisher D. Rapid publications risk the integrity of science in the era of COVID-19. *BMC Med* 2020;18:192. doi:10.1186/s12916-020-01650-6
- 11 Yeo-Teh NSL, Tang BL. An alarming retraction rate for scientific publications on Coronavirus Disease 2019 (COVID-19). *Accountability in Research* 2020;0:1–7.
doi:10.1080/08989621.2020.1782203
- 12 Abritis A, Marcus A, Oransky I. An “alarming” and “exceptionally high” rate of COVID-19 retractions? *Accountability in Research* 2021;28:58–9. doi:10.1080/08989621.2020.1793675
- 13 Zdravkovic M, Berger-Estilita J, Zdravkovic B, et al. Scientific quality of COVID-19 and SARS CoV-2 publications in the highest impact medical journals during the early phase of the pandemic: A case control study. *PLoS ONE* 2020;15:e0241826.
doi:10.1371/journal.pone.0241826
- 14 Elgendy IY, Nimri N, Barakat AF, et al. A systematic bias assessment of top-cited full-length original clinical investigations related to COVID-19. *European Journal of Internal Medicine* 2021;:S0953620521000182. doi:10.1016/j.ejim.2021.01.018
- 15 Glasziou PP, Sanders S, Hoffmann T. Waste in covid-19 research. *BMJ* 2020;369:m1847. doi:10.1136/bmj.m1847
- 16 Powell M, Koenecke A, Byrd JB, et al. A how-to guide for conducting retrospective analyses: example COVID-19 study. *Open Science Framework* 2020. doi:10.31219/osf.io/3drch
- 17 Haber NA, Clarke-Deelder E, Salomon JA, et al. Policy evaluation in COVID-19: A guide to common design issues. *American Journal of Epidemiology*. [in print]
- 18 Haber N. Systematic review of COVID-19 policy evaluation methods and design. Published Online First: 26 November 2020.<https://osf.io/7nbk6> (accessed 15 Jan 2021).
- 19 PRISMA. <http://www.prisma-statement.org/PRISMAStatement/> (accessed 15 Jan 2021).

- 1
2
3 20 Petherick A, Kira B, Hale T, et al. Variation in Government Responses to COVID-19.
4 <https://www.bsg.ox.ac.uk/research/publications/variation-government-responses-covid-19>
5 (accessed 24 Nov 2020).
6
7 21 R Core Team. R: A language and environment for statistical computing. Vienna, Austria:
8 : R Foundation for Statistical Computing 2019. <https://www.R-project.org/>
9
10 22 Krippendorff KH. Content Analysis: An Introduction to Its Methodology. SAGE
11 Publications 1980.
12
13 23 Gamer M, Lemon J, Fellows I, et al. irr: Various Coefficients of Interrater Reliability and
14 Agreement. <https://cran.r-project.org/web/packages/irr/index.html>
15
16 24 Aragon TJ, Fay MP, Wollschlaeger D, et al. Epitools. CRAN: 2017. [https://cran.r-](https://cran.r-project.org/web/packages/epitools/epitools.pdf)
17 [project.org/web/packages/epitools/epitools.pdf](https://cran.r-project.org/web/packages/epitools/epitools.pdf)
18
19 25 About Google Scholar. <https://scholar.google.com/intl/en/scholar/about.html> (accessed
20 15 Jan 2021).
21
22 26 Clarivate Analytics. Journal Citation Reports. 2019.
23
24 27 Cobb JS, Seale MA. Examining the effect of social distancing on the compound growth
25 rate of COVID-19 at the county level (United States) using statistical analyses and a random
26 forest machine learning model. *Public Health* 2020;185:27–9. doi:10.1016/j.puhe.2020.04.016
27
28 28 Lyu W, Wehby GL. Comparison of Estimated Rates of Coronavirus Disease 2019
29 (COVID-19) in Border Counties in Iowa Without a Stay-at-Home Order and Border Counties in
30 Illinois With a Stay-at-Home Order. *JAMA Netw Open* 2020;3:e2011102.
31 doi:10.1001/jamanetworkopen.2020.11102
32
33 29 Tam K-M, Walker N, Moreno J. Effect of mitigation measures on the spreading of
34 COVID-19 in hard-hit states in the U.S. *PLoS One* 2020;15:e0240877.
35 doi:10.1371/journal.pone.0240877
36
37 30 Courtemanche C, Garuccio J, Le A, et al. Strong Social Distancing Measures In The
38 United States Reduced The COVID-19 Growth Rate: Study evaluates the impact of social
39 distancing measures on the growth rate of confirmed COVID-19 cases across the United States.
40 *Health Affairs* 2020;39:1237–46. doi:10.1377/hlthaff.2020.00608
41
42 31 Crokidakis N. COVID-19 spreading in Rio de Janeiro, Brazil: Do the policies of social
43 isolation really work? *Chaos Solitons Fractals* 2020;136:109930.
44 doi:10.1016/j.chaos.2020.109930
45
46 32 Hyafil A, Moriña D. Analysis of the impact of lockdown on the reproduction number of the
47 SARS-Cov-2 in Spain. *Gaceta Sanitaria* 2020;:S0213911120300984.
48 doi:10.1016/j.gaceta.2020.05.003
49
50 33 Castillo RC, Staguhn ED, Weston-Farber E. The effect of state-level stay-at-home
51 orders on COVID-19 infection rates. *American Journal of Infection Control* 2020;48:958–60.
52 doi:10.1016/j.ajic.2020.05.017
53
54 34 Alfano V, Ercolano S. The Efficacy of Lockdown Against COVID-19: A Cross-Country
55 Panel Analysis. *Appl Health Econ Health Policy* 2020;18:509–17. doi:10.1007/s40258-020-
56 00596-3
57
58 35 Lyu W, Wehby GL. Community Use Of Face Masks And COVID-19: Evidence From A
59 Natural Experiment Of State Mandates In The US: Study examines impact on COVID-19 growth
60 rates associated with state government mandates requiring face mask use in public. *Health
Affairs* 2020;39:1419–25. doi:10.1377/hlthaff.2020.00818

- 1
2
3 36 Zhang R, Li Y, Zhang AL, et al. Identifying airborne transmission as the dominant route
4 for the spread of COVID-19. *Proc Natl Acad Sci USA* 2020;117:14857–63.
5 doi:10.1073/pnas.2009637117
6
7 37 Xu J, Hussain S, Lu G, et al. Associations of Stay-at-Home Order and Face-Masking
8 Recommendation with Trends in Daily New Cases and Deaths of Laboratory-Confirmed
9 COVID-19 in the United States. *Explor Res Hypothesis Med* 2020;:1–10.
10 doi:10.14218/ERHM.2020.00045
11
12 38 Lyu W, Wehby GL. Shelter-In-Place Orders Reduced COVID-19 Mortality And Reduced
13 The Rate Of Growth In Hospitalizations. *Health Aff (Millwood)* 2020;39:1615–23.
14 doi:10.1377/hlthaff.2020.00719
15
16 39 Wagner AB, Hill EL, Ryan SE, et al. Social distancing merely stabilized COVID-19 in the
17 United States. *Stat* 2020;9. doi:10.1002/sta4.302
18
19 40 Di Bari M, Balzi D, Carreras G, et al. Extensive Testing May Reduce COVID-19
20 Mortality: A Lesson From Northern Italy. *Front Med* 2020;7:402. doi:10.3389/fmed.2020.00402
21
22 41 Islam N, Sharp SJ, Chowell G, et al. Physical distancing interventions and incidence of
23 coronavirus disease 2019: natural experiment in 149 countries. *BMJ* 2020;:m2743.
24 doi:10.1136/bmj.m2743
25
26 42 Wong LP, Alias H. Temporal changes in psychobehavioural responses during the early
27 phase of the COVID-19 pandemic in Malaysia. *J Behav Med* 2020;:1–11. doi:10.1007/s10865-
28 020-00172-z
29
30 43 Liang X-H, Tang X, Luo Y-T, et al. Effects of policies and containment measures on
31 control of COVID-19 epidemic in Chongqing. *WJCC* 2020;8:2959–76.
32 doi:10.12998/wjcc.v8.i14.2959
33
34 44 Banerjee T, Nayak A. U.S. county level analysis to determine If social distancing slowed
35 the spread of COVID-19. *Revista Panamericana de Salud Pública* 2020;44:1.
36 doi:10.26633/RPSP.2020.90
37
38 45 Dave D, Friedson AI, Matsuzawa K, et al. When Do Shelter-in-Place Orders Fight
39 COVID-19 Best? Policy Heterogeneity Across States and Adoption Time. *Econ Inq* Published
40 Online First: 3 August 2020. doi:10.1111/ecin.12944
41
42 46 Hsiang S, Allen D, Annan-Phan S, et al. The effect of large-scale anti-contagion policies
43 on the COVID-19 pandemic. *Nature* 2020;584:262–7. doi:10.1038/s41586-020-2404-8
44
45 47 Lim JT, Dickens BSL, Choo ELW, et al. Revealing regional disparities in the
46 transmission potential of SARS-CoV-2 from interventions in Southeast Asia. *Proc R Soc B*
47 2020;287:20201173. doi:10.1098/rspb.2020.1173
48
49 48 Arshed N, Meo MS, Farooq F. Empirical assessment of government policies and
50 flattening of the COVID 19 curve. *J Public Affairs* Published Online First: 27 August 2020.
51 doi:10.1002/pa.2333
52
53 49 Wang K-W, Gao J, Song X-X, et al. Fangcang shelter hospitals are a One Health
54 approach for responding to the COVID-19 outbreak in Wuhan, China. *One Health*
55 2020;10:100167. doi:10.1016/j.onehlt.2020.100167
56
57 50 Kang N, Kim B. The Effects of Border Shutdowns on the Spread of COVID-19. *J Prev*
58 *Med Public Health* 2020;53:293–301. doi:10.3961/jpmph.20.332
59
60

- 1
2
3 51 Auger KA, Shah SS, Richardson T, et al. Association Between Statewide School
4 Closure and COVID-19 Incidence and Mortality in the US. *JAMA* 2020;324:859.
5 doi:10.1001/jama.2020.14348
6
7 52 Santamaría L, Hortal J. COVID-19 effective reproduction number dropped during Spain's
8 nationwide dropdown, then spiked at lower-incidence regions. *Science of The Total*
9 *Environment* 2021;751:142257. doi:10.1016/j.scitotenv.2020.142257
10
11 53 Bennett M. All things equal? Heterogeneity in policy effectiveness against COVID-19
12 spread in Chile. *World Development* 2021;137:105208. doi:10.1016/j.worlddev.2020.105208
13
14 54 Yang T, Shi H, Liu J, et al. Lessons Learnt from China: National Multidisciplinary
15 Healthcare Assistance. *RMHP* 2020;Volume 13:1835–7. doi:10.2147/RMHP.S269523
16
17 55 Padalabalanarayanan S, Hanumanthu VS, Sen B. Association of State Stay-at-Home
18 Orders and State-Level African American Population With COVID-19 Case Rates. *JAMA Netw*
19 *Open* 2020;3:e2026010. doi:10.1001/jamanetworkopen.2020.26010
20
21 56 Edelstein M, Obi C, Chand M, et al. SARS-CoV-2 infection in London, England: changes
22 to community point prevalence around lockdown time, March–May 2020. *J Epidemiol*
23 *Community Health* 2020;;jech-2020-214730. doi:10.1136/jech-2020-214730
24
25 57 Tsai AC, Harling G, Reynolds Z, et al. COVID-19 transmission in the U.S. before vs.
26 after relaxation of statewide social distancing measures. *Clin Infect Dis Published Online First:*
27 *3 October 2020.* doi:10.1093/cid/ciaa1502
28
29 58 Singh BB, Lowerison M, Lewinson RT, et al. Public health interventions slowed but did
30 not halt the spread of COVID-19 in India. *Transbound Emerg Dis Published Online First:* 4
31 *October 2020.* doi:10.1111/tbed.13868
32
33 59 Gallaway MS, Rigler J, Robinson S, et al. Trends in COVID-19 Incidence After
34 Implementation of Mitigation Measures — Arizona, January 22–August 7, 2020. *MMWR Morb*
35 *Mortal Wkly Rep* 2020;69:1460–3. doi:10.15585/mmwr.mm6940e3
36
37 60 Castex G, Dechter E, Lorca M. COVID-19: The impact of social distancing policies,
38 cross-country analysis. *EconDisCliCha Published Online First:* 15 October 2020.
39 doi:10.1007/s41885-020-00076-x
40
41 61 Silva L, Figueiredo Filho D, Fernandes A. The effect of lockdown on the COVID-19
42 epidemic in Brazil: evidence from an interrupted time series design. *Cad Saúde Pública*
43 *2020;36:e00213920.* doi:10.1590/0102-311x00213920
44
45 62 Dave D, Friedson A, Matsuzawa K, et al. Were Urban Cowboys Enough to Control
46 COVID-19? Local Shelter-in-Place Orders and Coronavirus Case Growth. *J Urban Econ*
47 *2020;;103294.* doi:10.1016/j.jue.2020.103294
48
49 63 Yin Y, Gao J, Jones BF, et al. Coevolution of policy and science during the pandemic.
50 *Science* 2021;371:128–30. doi:10.1126/science.abe3084
51
52 64 Wynants L, Van Calster B, Collins GS, et al. Prediction models for diagnosis and
53 prognosis of covid-19: systematic review and critical appraisal. *BMJ* 2020;;m1328.
54 doi:10.1136/bmj.m1328
55
56 65 Clarke GM, Conti S, Wolters AT, et al. Evaluating the impact of healthcare interventions
57 using routine data. *BMJ* 2019;;l2239. doi:10.1136/bmj.l2239
58
59 66 Krishnaratne S, Pfadenhauer LM, Coenen M, et al. Measures implemented in the school
60 setting to contain the COVID-19 pandemic: a rapid scoping review. *Cochrane Database of*

1
2
3 Systematic Reviews Published Online First: 17 December 2020.

4 doi:10.1002/14651858.CD013812

5 67 Raynaud M, Zhang H, Louis K, et al. COVID-19-related medical research: a meta-
6 research and critical appraisal. BMC Med Res Methodol 2021;21:1. doi:10.1186/s12874-020-
7 01190-w

8 68 Casigliani V, De Nard F, De Vita E, et al. Too much information, too little evidence: is
9 waste in research fuelling the covid-19 infodemic? BMJ 2020;:m2672. doi:10.1136/bmj.m2672

10 69 Greenhalgh T. Will COVID-19 be evidence-based medicine's nemesis? PLoS Med
11 2020;17:e1003266. doi:10.1371/journal.pmed.1003266
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Transparency declaration

The lead author (NH) affirms that this manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been explained.

Transparency declaration

The lead author (NH) affirms that this manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been explained.

Acknowledgements

We would like to thank Dr. Steven Goodman and Dr. John Ioannidis for their support during the development of this study, and Dr. Lars Hemkins and Dr. Mario Malicki for helpful comments in the protocol development.

Author roles

Screening reviewers: CJ, SW, CB, CA, NH, CBF, VN, and Keletso Makofane

Full article reviewers: NH, ECD, AF, BMG, ES, CBF, JD, LH, CG, CB, EBM, CJ, BL, IS, EA, SW, BJ, CA, VN, BG, AB, ES

Protocol development: NH, EC, JS, AF, ES

Administration, primary manuscript writing, data management, and analysis: NH

All authors participated in manuscript editing.

Funding

No funding was provided specifically for this research, and no funder was involved in the study design, data collection, analysis, interpretation, or writing of this study. All researchers were acting independently from funders. All authors had full access to all of the in the study and can take responsibility for the integrity of the data and the accuracy of the data analysis is also required.

Elizabeth Stone receives funding under the National Institutes of Health grant T32MH109436.

Ian Schmid receives funding under the National Institutes of Health grant T32MH122357.

Brooke Jarrett receives funding under the National Institutes of Health grant MH121128.

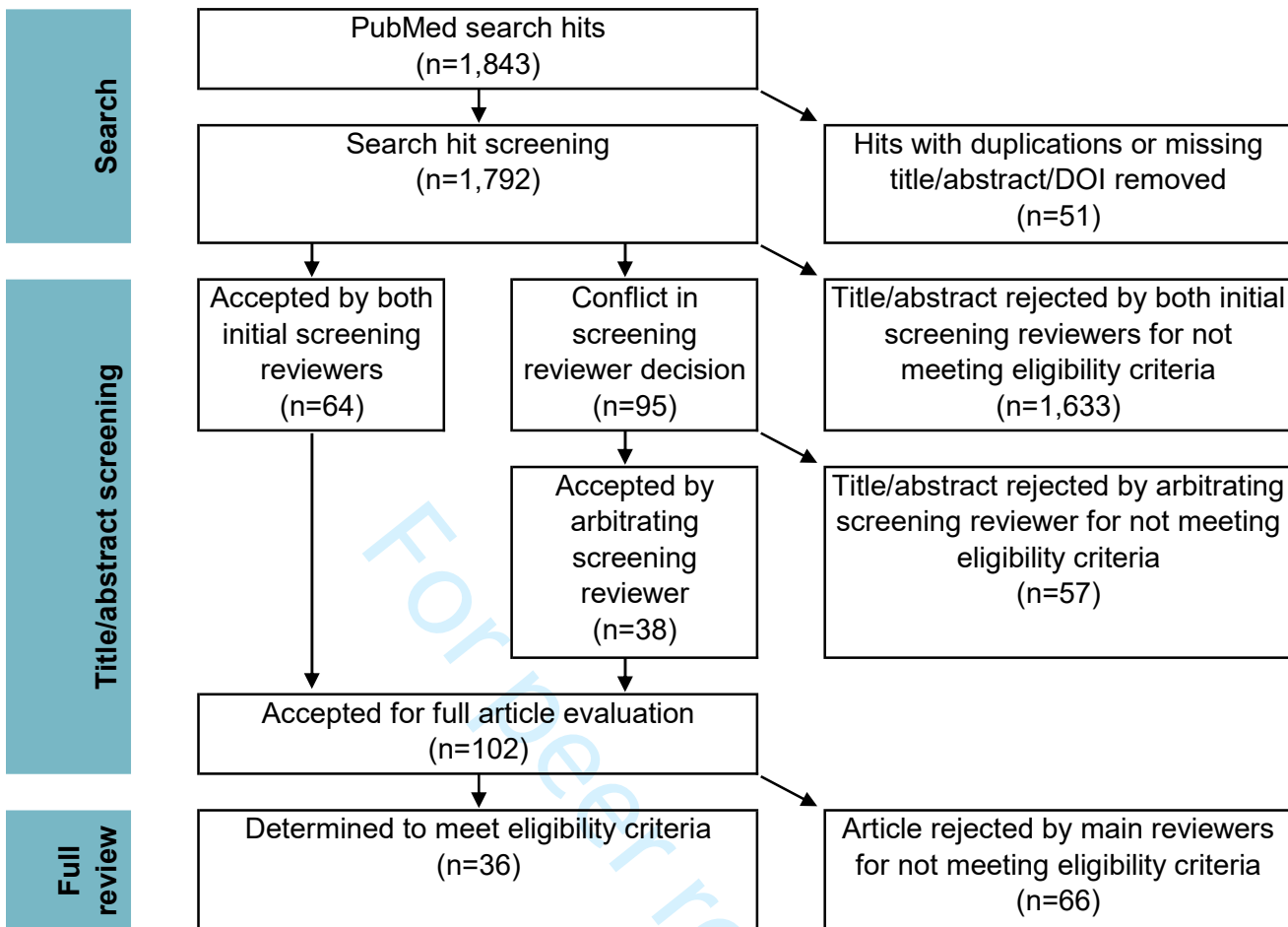
Christopher Boyer receives funding under the National Institutes of Health grant T32HL098048

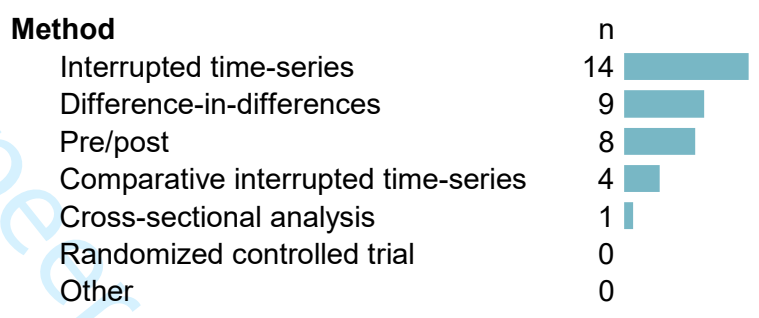
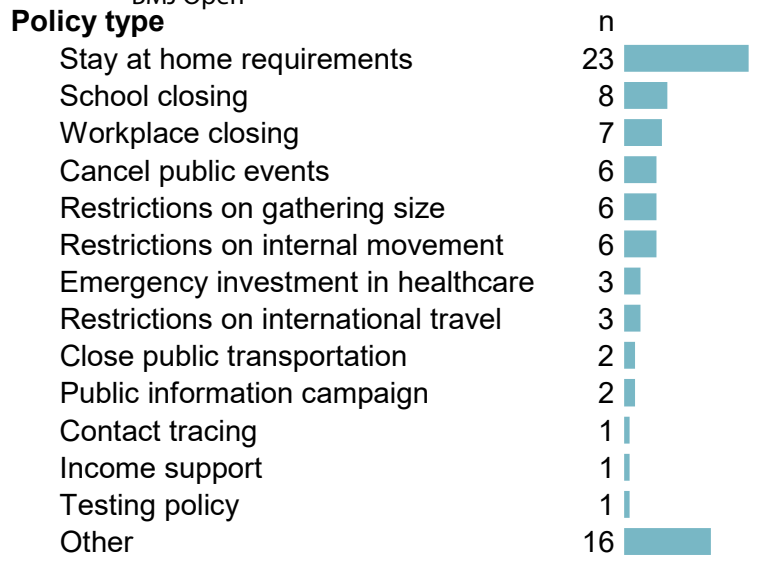
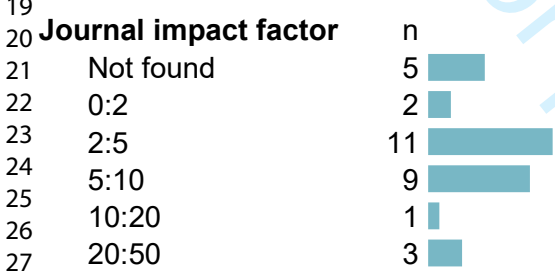
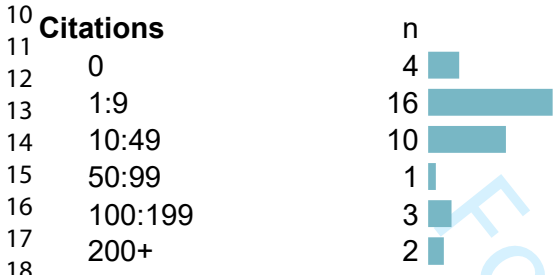
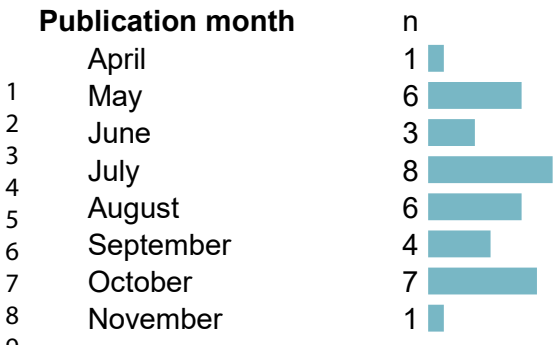
1
2
3 Cathrine Axfors receives funding from the Knut and Alice Wallenberg Foundation, grant KAW
4 2019.0561.

5 Beth Ann Griffin and Elizabeth Stuart were supported by award number P50DA046351 from the
6 National Institute on Drug Abuse. Elizabeth Stuart's time was also supported by the Bloomberg
7 American Health Initiative. Caroline Joyce receives funding from the Ferring Foundation. Meta-
8 Research Innovation Center at Stanford (METRICS), Stanford University is supported by a grant
9 from the Laura and John Arnold Foundation
10
11

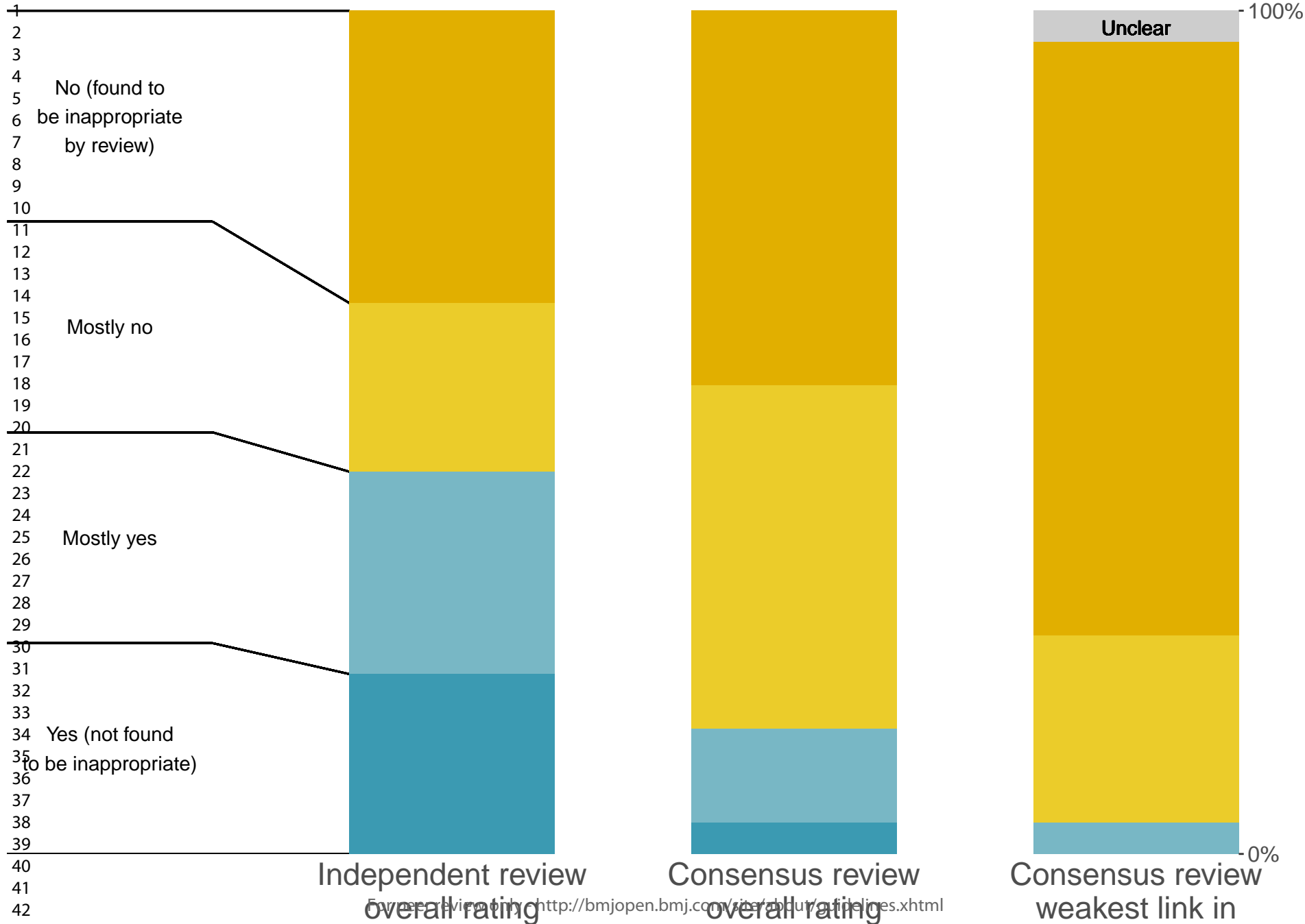
12 13 Conflicts of interest disclosure

14
15 The authors have no financial or social conflicts or competing interests to declare.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60





Did the study meet design criteria?



Independent review overall rating

Consensus review overall rating

Consensus review weakest link in key questions

For more information on this article, please visit <http://bmjopen.bmj.com/content/10/2/e20150111>

40
41
42
43
44

COVID-19 Health Policy Impact Evaluation Review

Start of Block: Main form

Q10 Administrative information

Q8 Study DOI

Q3 Reviewer number

Q54 Review type/round

The first round (Primary/independent review round) is for the independent first reviews of every article; the second (Secondary/consensus round) is for the second round of review for each article.

- Primary/independent review round (1)
- Secondary/consensus round (2)

Q50 Screening

1
2
3 Q52 Do you wish to recuse yourself from reviewing this study for any reason (e.g. social or
4 professional relationship with the authors, financial conflict of interest, etc)?
5

- 6
7 No, I do not wish to recuse myself. (1)
8
9 Yes, I recuse myself from reviewing this paper. (2)
10

11
12 *Skip To: End of Survey If Q52 = Yes, I recuse myself from reviewing this paper.*
13
14

15
16 Q51 Do you believe that this study meets the inclusion criteria for this research?
17

18 The inclusion criteria are: The primary topic of the article must be evaluating one or more
19 individual COVID-19 policies on direct COVID-19 outcomes The primary
20 exposure(s) must be a policy, defined as a government-issued order at any government level to
21 address a directly COVID-19-related outcome (e.g. mask requirements, travel restrictions, etc).
22

23 COVID-19 outcomes may include cases detected, mortality, number of tests taken, test
24 positivity rates, Rt, etc. This may NOT include indirect impacts of COVID-19 on
25 things such as income, childcare, trust in science, etc. The primary outcome
26 being examined must be a COVID-19-specific outcome, as above. The study must be
27 designed as an impact evaluation study from primary data (i.e. not primarily a predictive or
28 simulation model or meta-analysis) The study must be peer reviewed, and published in a peer-
29 reviewed journal indexed by PubMed The study must have the title and abstract available
30 via PubMed at the time of the study start date The study must be written in English
31
32
33

- 34 Yes (1)
35
36 No (2) _____
37
38

39
40 *Skip To: End of Survey If Q51 = No*
41
42

43 Q7 Study topic information

44
45 Please consult review guidance ([available here](#)) for additional guidance on answering these
46 questions.
47
48
49

50
51
52 Q6 Main impact sentence
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Copy and paste the sentence from the abstract that best describes the main claim of the study (e.g. "Policy X had a positive impact on outcome Y")

For peer review only

1
2
3 Q9 Main COVID-19 policy type evaluated
45 Select all that apply. Note: categorization from the Oxford Government Response Tracker
6

- 7
-
- 8
-
- School closing (1)
-
- 9
-
- 10
-
- Workplace closing (2)
-
- 11
-
- 12
-
- Cancel public events (3)
-
- 13
-
- 14
-
- Restrictions on gathering size (4)
-
- 15
-
- 16
-
- Close public transportation (5)
-
- 17
-
- 18
-
- Stay at home requirements (6)
-
- 19
-
- 20
-
- Restrictions on internal movement (7)
-
- 21
-
- 22
-
- Restrictions on international travel (8)
-
- 23
-
- 24
-
- Income support (9)
-
- 25
-
- 26
-
- Debt/contract relief for household (10)
-
- 27
-
- 28
-
- Fiscal measures (11)
-
- 29
-
- 30
-
- Giving international support (12)
-
- 31
-
- 32
-
- Public information campaign (13)
-
- 33
-
- 34
-
- Testing policy (14)
-
- 35
-
- 36
-
- Contact tracing (15)
-
- 37
-
- 38
-
- Emergency investment in healthcare (16)
-
- 39
-
- 40
-
- Investment in COVID-19 vaccines (17)
-
- 41
-
- 42
-
- 43
-
- 44
-
- 45
-
- 46
-
- 47
-
- 48
-
- 49
-
- 50
-
- 51
-
- 52
-
- 53
-
- 54
-
- 55
-
- 56
-
- 57
-
- 58
-
- 59
-
- 60

Other policy response (fill in) (18)

Q12 Main COVID-19 outcome type evaluated

Select all that apply

COVID-19 cases (1)

COVID-19 test positivity (2)

COVID-19 deaths (3)

COVID-19 hospitalizations (4)

SARS-CoV-2 infections and infection rate (e.g. effective R) (8)

Other (fill in) (9) _____

Q13 **Method(s) identification**

For this section, consider only the data structure as it enters into the main statistical model. In other words, if the original dataset is of individuals at many time points, but the main statistical model uses a regional-level aggregated count of cases, the data as it enters into the main statistical model is a regional aggregate at one time point.

Q14 What is the level of aggregation for the main outcome data?

Individual level (1)

Regional aggregate (e.g. count, mean, etc.) (2)

1
2
3
4 Q16 How many regional units included in the main statistical model received the policy of
5 interest?
6

7
8 If 2-20, enter the number of regional units analyzed which received the policy of interest.
9

10 One (1) (1)

11
12 Two through twenty (2-20) (2)
13
14 _____

15
16 More than twenty (21+) (3)
17

18 Unclear or N/A (4) _____
19
20

21
22
23
24 Q17 How many regional units were included which did NOT receive any form of the policy of
25 interest?
26

27
28 If 2-20, enter the number of regional units analyzed which did not receive the policy of interest.
29

30 Zero (0) (1)

31
32 One (1) (2)

33
34 Two through twenty (2-20) (3)
35
36 _____

37
38 More than twenty (21+) (4)
39

40 Unclear or N/A (5) _____
41
42

43
44
45 *Display This Question:*

46 *If Q17 = Zero (0)*
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q25 Did different regions receive different intensities of the policy of interest for comparison?
4
5

6 For example, the study might compare places with more intense versions of policy or policies
7 vs. places with less intense versions of policy or policies, rather than just places with and
8 without the policy or policies.
9

- 10 Yes (regions with more intense policy were compared with regions with less intense
11 policy) (1)
12
13 No (2)
14
15 Unclear or N/A (3)
16
17
18
-

19
20
21
22 Q18 For each regional unit, how many time point observations were in the model *before* the
23 policy was enacted?
24

- 25 None (0) (1)
26
27 One (1) (2)
28
29 More than one (2+) (3)
30
31 Unclear or N/A (4) _____
32
33
34
-

35
36
37
38 Q19 For each regional unit, how many time point observations were in the model *after* the policy
39 was enacted?
40

- 41 None (0) (1)
42
43 One (1) (2)
44
45 More than one (2+) (3)
46
47 Unclear or N/A (4) _____
48
49
50
51
-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q19 = One (1)

And Q18 = One (1)

Or If

Q19 = More than one (2+)

Or If

Q18 = More than one (2+)

Q20 How would you describe the time intervals between observations?

- Days (1-5 days between observations) (1)
- Weeks (about 5-10 days between observations) (2)
- Multiple weeks (11-25 days between observations) (3)
- Monthly (26 or more days between observations) (4)

Display This Question:

If Q17 = Zero (0)

Q21 Did the pre-policy period for any region act as a "control" for different region post-policy enactment?

In other words, was there any pre-period in one or more region's being used to control or compare for the trends of any one or more *different* regions' post-period?

- No (pre-periods were treated as controls only within-region) (1)
- Yes (pre-periods were treated as controls with other regions) (2)
- Unclear or N/A (3)
-

Q22 Was any unit assigned the policy or the timing of the policy externally (i.e. as an experiment/trial)?

- No (observational data only) (1)
- Yes (treatment assigned as part of research or evaluation) (2)
- Unclear or N/A (3)

Display This Question:

If Q22 = Yes (treatment assigned as part of research or evaluation)

Q23 Was the assignment randomized?

- Yes (1)
- No (2)

Q27 Based on your answers above and the guidance document, please select the type of study that best resembles the design of the main analysis.

Please note that the design(s) named in the paper may not match with the method described below, nor is this the actual exact design that was used. If you believe that the design used differs from the choices below in a way that makes this choice impossible, please contact the study administrator before selecting "other."

	Design	
	Units (e.g., regions of comparison)	
	Time points measured per unit	
	Assumed counterfactual.	
	"If not for the intervention, ___"	
	Without intervention	With intervention
	Before intervention	
	After intervention	
	Cross-sectional	
	At least one	
At least one		N/A
	One time point	
	Outcome in intervention units would have been the same	

as the outcome in the non-intervention units.

Pre/post

At least one

None

At least one (typically one)

At least one (typically one)

Outcome would have stayed the same from the pre period to the post period.

Interrupted

time-series

(ITS)

At least one

None

More than one

At least one (typically several)

Outcome slope and level* would have continued along the same modelled trajectory from the pre-period to the post period.

Difference-in-differences

(DiD)

At least one

At least one†

At least one (typically one)

At least one (typically one)

Outcome in intervention units would have changed as much as (or in parallel with) the outcome in the non-intervention units.

Comparative interrupted time series (CITS)

At least one

At least one†

More than one (typically several)

At least one (typically several)

Outcome slope and level* would have changed as much as non-intervention group's slope and level* changed.

* Assessing both slope and level only applicable if there are multiple data points during the post period
 † Units without the intervention may be the pre-period of a different unit that eventually receives the intervention.

Cross-sectional analysis (1)

Non-randomized experiment/trial (2)

Randomized controlled trial (3)

- 1
2
3
4 Pre/post (4)
5
6 Interrupted time-series (5)
7
8
9 Difference-in-differences (6)
10
11 Comparative interrupted time-series (7)
12
13
14 Other (please contact administrator before selecting) (8)
15 _____
16

17
18
19
20 **Q49 Design evaluation**
21
22

23
24 *Display This Question:*

25 *If Q27 = Interrupted time-series*

26 *Or Q27 = Difference-in-differences*

27 *Or Q27 = Comparative interrupted time-series*
28
29
30

31 **Q29 Does the analysis provide graphical representation of the outcome over time?**
32
33

34
35 If not "Yes" please describe (three short sentences max).
36

37 -Check for a chart that shows the outcome over time, with the dates of interest, separated by
38 policy/non policy groups if applicable. Outcomes may be aggregated for clarity (e.g. means and
39 CIs at discrete time points).
40

- 41
42 Yes (1)
43
44 Mostly yes (2) _____
45
46 Mostly no (3) _____
47
48 No (4) _____
49
50 Unclear (5) _____
51
52
53
54

1
2
3 *Display This Question:*

4 *If Q27 = Interrupted time-series*

5 *Or Q27 = Difference-in-differences*

6 *Or Q27 = Comparative interrupted time-series*
7
8
9

10 Q30 Is there sufficient pre-intervention data to characterize pre-trends in the data?
11
12

13
14 If not "Yes" please describe (three short sentences max).
15

16 -Check the chart(s) to see if there are several time points over a reasonable period of time over
17 which to establish stability and curvature in the pre-trends.
18

19
20 Yes (1)
21

22 Mostly yes (2) _____
23

24 Mostly no (3) _____
25

26 No (4) _____
27

28 Unclear (5) _____
29
30
31

32
33
34 *Display This Question:*

35 *If Q27 = Interrupted time-series*

36 *Or Q27 = Difference-in-differences*

37 *Or Q27 = Comparative interrupted time-series*
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q32 Is the pre-trend stable?
4
5

6
7 If not "Yes" please describe (three short sentences max).
8

9 -Check if there are sufficient data to reasonably determine a stable functional form for the pre-
10 trends, and that they follow a modelable functional form.
11

- 12
13 Yes (1)
14
15 Mostly yes (2) _____
16
17 Mostly no (3) _____
18
19 No (4) _____
20
21 Unclear (5) _____
22
23
24

25
26
27 *Display This Question:*

28 *If Q27 = Interrupted time-series*

29 *Or Q27 = Comparative interrupted time-series*
30
31

32 Q31 Is there sufficient post-intervention data to observe post trends in the data?
33
34

35
36 If not "Yes" please describe (three short sentences max).
37

38 -Check the chart(s) to see if there are several time points over a reasonable period of time over
39 which to establish stability and curvature in the post- trends.
40
41

- 42 Yes (1)
43
44 Mostly yes (2) _____
45
46 Mostly no (3) _____
47
48 No (4) _____
49
50 Unclear (5) _____
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q33 Is the functional form of the counterfactual (e.g. linear) well-justified and appropriate?

If not "Yes" please describe (three short sentences max).

- Check whether the authors explain and justify their choice of functional form.
- Check if there is any curvature in the pre-trend.
- Consider the nature of the outcome. Is it sensible to measure the trend of this outcome on the scale and form used? Note: infectious disease dynamics are rarely linear.
- Consider that while pre-trend fit is a necessary condition for an appropriate linear counterfactual model, it is not sufficient. Check if the authors provide justification for the functional form to continue to be of the same functional form (e.g. linear).

- Yes (1)
- Mostly yes (2) _____
- Mostly no (3) _____
- No (4) _____
- Unclear (5) _____

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q34 Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?

If not "Yes" please describe (three short sentences max).

- Check whether the authors justify the use of the date threshold relative to the date of the intervention.
- Trace the process between the intervention being put in place to when observable effects in

1
2
3 the outcome might appear over time.

4 -Consider whether there are anticipation effects (e.g. do people change behaviors before the
5 date when the intervention begins?)

6
7 -Consider whether there are lag effects. (e.g. does it take time for behaviors to change,
8 behavior change to impact infections, infections to impact testing, and data to be collected, etc?)

9 -Check if authors appropriately and directly account for these time effects.
10

11 Yes (1)

12 Mostly yes (2) _____
13

14 Mostly no (3) _____
15

16 No (4) _____
17

18 Unclear (5) _____
19
20
21
22
23

24 -----
25 *Display This Question:*

26 *If Q27 = Interrupted time-series*
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q36 Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have
4 changed during the measurement period?
5

6
7 If not "Yes" please describe (three short sentences max).
8

- 9 -Consider other policies or interventions which could impact the outcome during this time.
10 -Consider social behaviors changed which could meaningfully impact the outcome during this
11 time.
12 -Consider economic conditions changed which could meaningfully impact the outcome during
13 this time.
14 -Note that the actual concurrent changes do not need to happen during the period of
15 measurement, just their effects.
16
17

- 18
19 Yes (1)
20
21 Mostly yes (2) _____
22
23 Mostly no (3) _____
24
25 No (4) _____
26
27 Unclear (5) _____
28
29
30
31

32
33 *Display This Question:*

34 *If Q27 = Difference-in-differences*

35 *Or Q27 = Comparative interrupted time-series*
36
37

38
39 Q53

40 Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have
41 changed during the measurement period, differently for policy and non-policy regions?
42

43 If not "Yes" please describe (three short sentences max).
44

45 -Consider any uncontrolled factor which could have influenced the outcome differently in policy
46 and non-policy regions.
47

48 -This may include (but is not limited to)

- 49 -Other policies
50 -Social behaviors
51 -Economic conditions
52

53 -Are these factors justified as having negligible impact?

54 -If justified, is the argument that these have negligible impact convincing?
55
56
57
58
59
60

1
2
3 -Note that the actual concurrent changes do not need to happen during the period of
4 measurement, just their effects.
5

- 6
7 Yes (1)
8
9 Mostly yes (2) _____
10
11 Mostly no (3) _____
12
13 No (4) _____
14
15 Unclear (5) _____
16
17
18
19

20
21 *Display This Question:*

22 *If Q27 = Interrupted time-series*

23 *Or Q27 = Difference-in-differences*

24 *Or Q27 = Comparative interrupted time-series*
25
26
27

28 **Q38**

29 Did authors provide diagnostics or show robustness and/or sensitivity of results to alternative
30 model choices?
31

32
33
34 If not "Yes" please describe (three short sentences max).
35

- 36 Yes (1)
37
38 Mostly yes (2) _____
39
40 Mostly no (3) _____
41
42 No (4) _____
43
44 Unclear (5) _____
45
46
47
48
49

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q39 Given the above, do you believe that the design is appropriate for identifying the policy impact(s) of interest?

This should be taken as independent of what you believe about other studies, and/or the feasibility of other designs.

If not "Yes" please describe (three short sentences max).

Yes (1)

Mostly yes (2) _____

Mostly no (3) _____

No (4) _____

Unclear (5) _____

Display This Question:

If Q54 = Secondary/consensus round

Q55 General and/or additional comments on this paper from consensus discussion. This may include any additional information worth commenting on regarding the paper, difficulties encountered evaluating it, etc.

(three short sentences max)

End of Block: Main form

Review version with references removed; NOT FOR DISTRIBUTION

Policy evaluation in COVID-19: A guide to common design issues

Noah A Haber, Emma Clarke-Deelder, Joshua A Salomon, Avi Feller, Elizabeth A Stuart

Noah A Haber, ScD*

noahhaber@stanford.edu

Meta Research Innovation Center at Stanford University

Stanford University

1265 Welch Rd

Palo Alto, CA 94305

(650) 497-0811

Emma Clarke-Deelder, MPhil

Department of Global Health & Population

Harvard T. H. Chan School of Public Health

665 Huntington Avenue

Building 1, room 1104

Boston, Massachusetts 02115

Joshua A Salomon, PhD

Department of Medicine

Center for Health Policy and Center for Primary Care and Outcomes Research

Stanford University School of Medicine

Encina Commons, Room 118

615 Crothers Way

Stanford, CA 94305-6019

Avi Feller, PhD

Goldman School of Public Policy

University of California, Berkeley

2607 Hearst Avenue

Room 309

Berkeley, CA 94720

Elizabeth A Stuart, PhD

Department of Mental Health

Johns Hopkins Bloomberg School of Public Health

624 N. Broadway

Hampton House 839

Baltimore, MD 21205

* corresponding author

Review version with references removed; NOT FOR DISTRIBUTION

Abstract

Policy responses to COVID-19, particularly those related to non-pharmaceutical interventions, are unprecedented in scale and scope. Researchers and policymakers are striving to understand the impact of these policies on a variety of outcomes. Policy impact evaluations always require a complex combination of circumstance, study design, data, statistics, and analysis. Beyond the issues that are faced for any policy, evaluation of COVID-19 policies is complicated by additional challenges related to infectious disease dynamics and lags, lack of direct observation of key outcomes, and a multiplicity of interventions occurring on an accelerated time scale.

In this paper, we (1) introduce the basic suite of policy impact evaluation designs for observational data, including cross-sectional analyses, pre/post, interrupted time-series, and difference-in-differences analysis, (2) demonstrate key ways in which the requirements and assumptions underlying these designs are often violated in the context of COVID-19, and (3) provide decision-makers and reviewers a conceptual and graphical guide to identifying these key violations. The overall goal of this paper is to help policy-makers, journal editors, journalists, researchers, and other research consumers understand and weigh the strengths and limitations of evidence that is essential to decision-making.

Introduction

The response to the global COVID-19 pandemic has demanded urgent decision making in the face of substantial uncertainties. Policies to arrest transmission, including stay-at-home orders and other non-pharmaceutical interventions (NPIs), have wide-reaching consequences that touch many aspects of well being. Decision-making in the public interest requires evaluating and weighing the evidence on both intended and unintended consequences in order to best predict outcomes. The wide range of policy interventions implemented by different jurisdictions may yield opportunities for learning from what has already happened to inform future policymaking, and we have observed a proliferation of studies aimed at such policy evaluations. However, policy evaluation requires a complex combination of circumstance, data, study design, analysis, and interpretation in order to be informative.

Policy impact evaluation aims to answer questions about the extent to which the realized outcomes given a particular policy would have been different in the absence of that policy. Estimating the causal impact of the policy with observational data is challenging because what would have happened in the absence of the policy change (the “counterfactual”) is, by definition, unobserved. Randomized controlled trials (RCTs) of policies related to COVID-19 interventions may not always be practical or ethical. In this context, a large and growing number of studies have attempted to evaluate the impact of COVID-19 policies using observational data. There

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 are many potential pitfalls in the use of observational data for evaluation generally, and some
4 additional methodological design challenges relating to COVID-19 policies in particular.
5

6
7 This paper provides a graphical guide to policy impact evaluations for COVID-19, targeted to
8 decision-makers, researchers and evidence curators. Our aim is to provide a coherent
9 framework for conceptualizing and identifying common pitfalls in COVID-19 policy evaluation.
10 Importantly, this should not be taken either as a comprehensive guide to policy evaluation more
11 broadly or as guidance on performing analysis, which may be found elsewhere. Rather, we
12 review relevant study designs for policy evaluations — including pre/post, interrupted time
13 series, and difference-in-difference approaches — and provide guidance and tools for
14 identifying key issues with each type of study as they relate to NPIs and other COVID-19 policy
15 interventions. Improving our ability to identify key pitfalls will enhance our ability to identify and
16 produce valid and useful evidence for informing policymaking.
17
18
19
20

21 Common policy evaluation designs and their pitfalls 22 in COVID-19 23 24

25 Identifying the type of design 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

Review version with references removed; NOT FOR DISTRIBUTION

Table 1: Summary definitions of policy impact evaluation designs commonly used for COVID-19

Design	Units (e.g., regions of comparison)		Time points measured per unit		Assumed counterfactual. “If not for the intervention, _____”
	With intervention	Without intervention	Before intervention	After intervention	
Cross-sectional	At least one	At least one	N/A	One time point	Outcome in intervention units would have been the same as the outcome in the non-intervention units.
Pre/post Figure 1A	At least one	None	At least one (typically one)	At least one (typically one)	Outcome would have stayed the same from the pre period to the post period.
Interrupted time-series (ITS) Figure 1B	At least one	None	More than one	At least one (typically several)	Outcome slope and level* would have continued along the same modelled trajectory from the pre-period to the post period.
Difference-in-differences (DiD) Figure 1C	At least one	At least one [†]	At least one (typically one)	At least one (typically one)	Outcome in intervention units would have changed as much as (or in parallel with) the outcome in the non-intervention units.
Comparative interrupted time series (CITS) Figure 1D	At least one	At least one [†]	More than one (typically several)	At least one (typically several)	Outcome slope and level* would have changed as much as non-intervention group's slope and level* changed.
* Assessing both slope and level only applicable if there are multiple data points during the post period † Units without the intervention may be the pre-period of a different unit that eventually receives the intervention.					

Identifying the underlying design in a given analysis often requires using a combination of the methods as reported and evaluating the data structure that is used for the main analysis, as shown in Table 1. COVID-19-related policy evaluation analyses typically fall under these categories. In most cases, the design can be categorized using a combination of whether there are also units that did not receive the treatment (columns 2-3) and whether there are time points both before and after intervention for those units (columns 4-5). The final column describes the implied counterfactual, discussed further in subsequent sections. Cross sectional designs typically compare units with vs without the treatment at single time points. Pre/post studies typically compare within units who received the intervention at two points: before and after a policy. Interrupted time-series analyses compare outcomes within units within units who received the intervention at greater than two time points before the intervention vs with at least one (typically multiple) after the intervention. Difference-in-differences analysis compares the outcome change in units which received the intervention with those that did not (or have not yet), with at least one point before and one after the intervention. In cases with multiple periods, that may involve a comparison with the pre-policy period of one region with the post-period of a different region, even though all regions eventually receive the intervention.

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

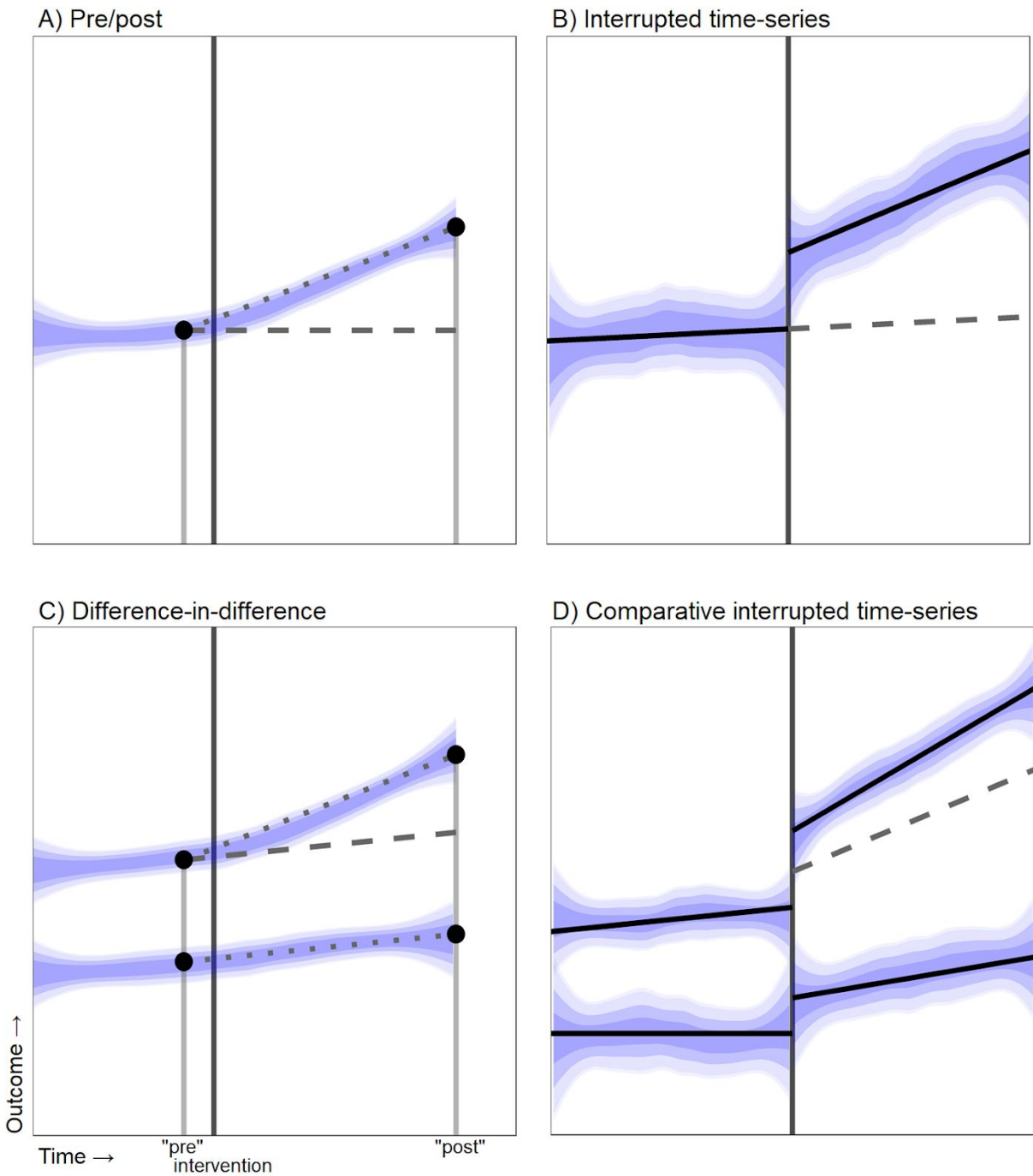
Methods descriptions may not always provide a precise or reliable guide to which of the design approaches has been used. Some studies do not explicitly name these designs (or may classify them differently); and these are only a small fraction of designs and frameworks that are possible to use for policy evaluation. Studies may have data at multiple time points but are effectively cross-sectional. DiD, ITS, and CITS designs based on repeated cross-sectional data are sometimes described as “cross-sectional” instead of longitudinal. The term “event study” is often used to refer to studies with a single unit and one change over time resembling ITS, but may refer to other designs. Although ITS is often used to describe changes in one unit, it may also refer to settings in which many treated units adopt an intervention over time. Studies will also frequently employ multiple designs, while others use more complex methods of generating counterfactuals. Definitions of these terms vary widely, and the definitions above should be considered as guidance only.

Policy impact evaluation design foundations for COVID-19

The simplest design is the cross-sectional analysis, which compares COVID-19 outcomes between units of observation (e.g., cities) at a single calendar time or time since an event, typically post-intervention. These studies are unlikely to be appropriate for COVID-19-related policy evaluations, but provide a useful starting point for reasoning about different designs. Just as with comparisons of non-randomized medical treatments, the localities that adopt a particular policy likely differ substantially from those that don't on both observed and unobserved characteristics on a number of dimensions, including epidemic status and timing.

Figure 1: Longitudinal designs overview

Review version with references removed; NOT FOR DISTRIBUTION



This chart shows four canonical longitudinal designs. In all cases: the blue shading represents the underlying data trends, the solid vertical grey line represents the time of intervention, the grey dashed lines represent the assumed counterfactual in the absence of the intervention, as discussed in the text. The impact estimate is obtained by comparing the outcomes observed for the treated unit in the post period (the solid line) with the implied counterfactual line (the dashed line). In the case of the pre/post and

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 difference-in-differences panels the large black dots represent the time of measurement,
4 connected by the grey dotted lines.
5

6
7 Given the challenges in a simple cross-sectional comparison, which compare post-intervention
8 outcomes, it is important to consider longitudinal designs, which instead look at differences or
9 trends across time, as summarized in Figure 1. These can be distinguished by the data used
10 and the construction of the counterfactual. Pre/post, for example, has only one unit, measured
11 at two time points. Two common strategies expand on the logic and data requirements of the
12 pre/post design. Interrupted time series designs (Figure 1B) incorporate multiple time points
13 before the intervention, and usually multiple time points after the intervention, to enable a more
14 complete view on changes in levels and trends that are temporally related to the intervention.
15 Difference-in-difference designs (Figure 1C) add a set of comparison points from a group or
16 location that did not have the intervention. Another related design (comparative interrupted
17 time-series, Figure 1D, discussed only briefly here), uses both aspects — a change over time
18 and a comparison group — to compare the observed change in slopes for the intervention
19 group with the change in slope for the comparison group.
20
21
22
23

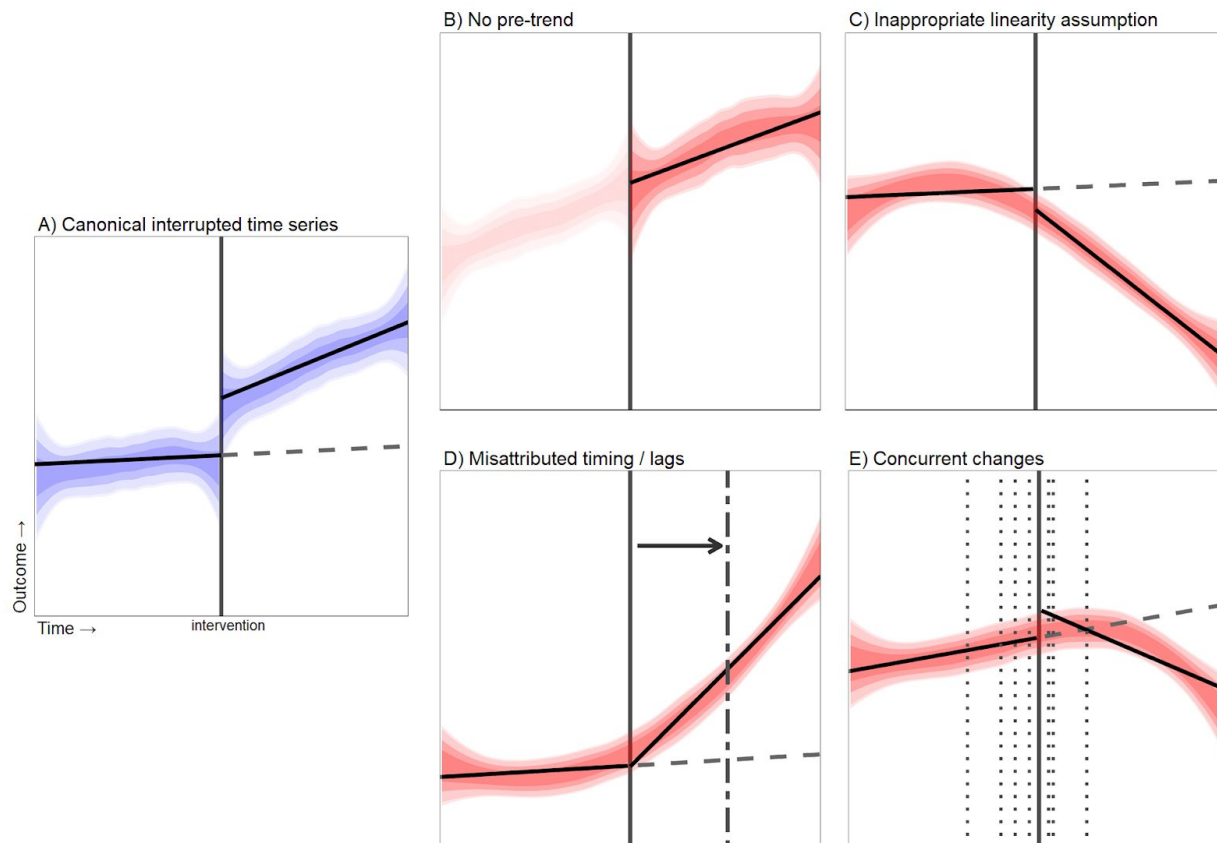
24 Pre/post studies

25
26 The simplest longitudinal design is a pre/post analysis, where some outcome is observed before
27 policy implementation, and again after, in a single group (Figure 1A). Pre/post studies are
28 analogous to a single arm trial with no control and only a single follow-up observation after
29 treatment.. This effectively imposes the assumption that the counterfactual trend is completely
30 flat (i.e., that the outcome in the post-period in the absence of the policy change is the same as
31 the value of the outcome before the policy change) without accounting for pre-existing
32 underlying trends, and attributing all outcome changes completely to the intervention of interest.
33 Just as the outcomes for an individual patient might be expected to change before and after
34 treatment, for reasons unrelated to the treatment, outcomes related to policy interventions will
35 change for reasons not caused by the policy. Infection rates, for example, would not be
36 expected to remain stationary except in very specific circumstances, but a pre/post
37 measurement would assume that any changes in infection rates are attributable to the policy.
38
39
40
41
42

43 Interrupted time-series

44
45 Figure 2: Interrupted time-series graphical guidance for identifying common pitfalls
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION



This chart shows one canonical design for ITS (blue, Panel A) and four panels demonstrating common issues with ITS analysis (red, panels B-E) discussed in the text. In all cases: the lag/red shading represents the underlying data trends, the vertical grey line represents the time of intervention, the grey dashed lines represent the assumed counterfactual in the absence of the intervention. In panel D) the dash-dot line represents the time at which the policy is expected to impact the outcome. In panel E), the vertical dotted lines represent concurrent events and changes.

Interrupted time-series (ITS) is a strategy that uses a projection of the pre-policy outcome trend as a counterfactual for how the outcome would have changed if the policy had not been introduced. In other words, in the absence of the policy change, ITS assumes the outcome would have continued on its pre-policy trend during the study period. ITS can be a useful tool in policy evaluation because it allows researchers to account for underlying trends in the outcome and, by comparing the treated unit (or location) to itself; it can therefore eliminate some of the confounding concerns that arise in cross-sectional or pre-post studies.

However, the validity of ITS depends critically on how well counterfactual trends in the outcome are modelled, and whether the policy of interest is the only relevant change during the study period. In the canonical setting (Figure 2A), the pre-policy trend is stable and can be feasibly modelled with the available data; the researcher appropriately models the timing of the change in the slope and/or level of the outcome; the researcher has sufficient information to conclude

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 that there were no other changes during the study period that would be expected to influence
4 the outcome. These elements are largely not satisfied in studies of COVID-related policy, as
5 described below.
6
7

8 ITS relies critically on modelled trends of the outcome over time. Key components of ITS
9 analyses include both visual and statistical examination of trends, preferentially alongside a
10 theoretical justification of the model used. At a minimum, analyses should provide graphical
11 representation of the data and model over time to examine whether pre-trend outcomes are
12 stable, all trends are well-fit to the data, “interrupted” at the appropriate time point, and sensibly
13 modelled (Figure 2B). In the case where an ITS includes a large number of units (e.g. states), it
14 can be difficult to display this information graphically.
15
16
17

18 One common pitfall in ITS is adoption of inappropriate assumptions on the outcome trend
19 (Figure 2C). The estimate of policy impact will be biased if a linear trend is assumed but the
20 outcome and response to interventions instead follow nonlinear trends (either before or after the
21 policy). In some cases, transformation of the outcome, for example using a log scale, may
22 improve the suitability of a linear model. Imposing linearity inappropriately is a serious risk in the
23 context of COVID-19, as trends in infectious disease dynamics are inherently non-linear. For
24 intuition, terms such as “exponential growth,” “flattening,” and “s-curves” all refer to non-linear
25 infectious disease trends. Depending on the particular situation, non-linearity or other modelled
26 trends can have complicated and counterintuitive impact on policy impact. Apparent linearity
27 may also be temporary and an artifact of testing, which may give a misleading impression that
28 linear models for infectious disease trends are appropriate indefinitely. While some use linear
29 progression in order to avoid more complex infectious disease models, in fact, linear projections
30 impose strict and often unrealistic models, generally resulting in an inappropriate counterfactual.
31
32
33
34
35

36 Researchers can easily misattribute the timing of the policy impact, resulting in spurious
37 inference and bias (Figure 2D). Some public health policies can be expected to translate into
38 immediate results (e.g., smoking bans and acute coronary events). In contrast, nearly every
39 outcome of interest in COVID-19 exhibits complex and difficult to infer time lags typically in the
40 realm of many weeks. The time between policy implementation and expected effect in the data
41 can be large and highly variable. For example, in order to see the impact of a mask order, first
42 the mask order takes effect, then people change their behaviors over time to comply with the
43 order (or sometimes the reverse in the case of anticipation effects), mask use behavior
44 produces changes in infections, then infections later result in symptoms, symptoms induce
45 people to seek testing, the tests must then be processed in labs, and then finally the results get
46 reported in data monitoring efforts. Selection of lead/lag time should be justifiable *a priori* or
47 external data. Selecting a lag based on the data risks issues comparable to p-hacking.
48
49
50
51

52 Finally, and perhaps most concerningly in the context of COVID-19, ITS fails when the policy of
53 interest coincides in time with other changes that affect the outcome (Figure 2E). For example, if
54 both mask and bar closure orders are rolled out together as a package, ITS cannot isolate the
55 impact of bar closures specifically. These changes do not need to have taken place exactly
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

concurrently with the policy implementation date of interest; they merely need to have some effect within the time period of measurement to result in potentially serious bias in effect estimates if unaddressed. ITS will also likely be biased if, during the study period, there is a change in the way the outcome data is collected or measured. This might occur if the introduction of a COVID-19 control policy is combined with an effort to collect better data on infection or mortality cases. Analogously, if an RCT involves randomizing people to a group receiving both A and B vs. control, we typically can't disentangle the effects of A from the effects of B, unless we also have separate A- and B-only arms. Ultimately, if multiple things are changing at the same time, ITS may not be an appropriate design for policy evaluation.

COVID-19 policies rarely arrive alone; they are typically created alongside other policies, unofficial action, and large scale behavior changes which themselves impact COVID-19-related outcomes. In some cases, anticipation of a policy may induce behavior change before the actual policy takes effect. The policies themselves may have been chosen due to the expectation of change in disease outcomes, which introduces additional biases related to “reverse” causality.

Table 2: Checklist for identifying common pitfalls for ITS to evaluate COVID-19 policy

Key design questions. If any answer is “no,” this analysis is unlikely to be appropriate or useful for estimating the impact of the intervention of interest.	Details and suggestions for identifying issues:
Does the analysis provide graphical representation of the outcome over time?	-Check for a chart that shows the outcome over time, with the dates of interest. Outcomes may be aggregated for clarity (e.g. means and CIs at discrete time points).
Is there sufficient pre-intervention data to characterize pre-trends in the data?	-Check the chart(s) to see if there are several time points over a reasonable period of time over which to establish stability and curvature in the pre-trends.
Is the pre-trend stable?	-Check if there are sufficient data to reasonably determine a stable functional form for the pre-trends, and that they follow a modelable functional form.
Is the functional form of the counterfactual (e.g. linear) well-justified and appropriate?	-Check whether the authors explain and justify their choice of functional form. -Check if there is any curvature in the pre-trend. -Consider the nature of the outcome. Is it sensible to measure the trend of this outcome on the scale and form used? Note: infectious disease dynamics are rarely linear. -Consider that while pre-trend fit is a necessary condition for an appropriate linear counterfactual model, it is not sufficient. Check if the authors provide justification for the functional form to continue to be of the same functional form (e.g. linear).
Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?	-Check whether the authors justify the use of the date threshold relative to the date of the intervention. -Trace the process between the intervention being put in place to when observable effects in the outcome might appear over time. -Consider whether there are anticipation effects (e.g. do people change behaviors before the date when the intervention begins?) -Consider whether there are lag effects. (e.g. does it take time for behaviors to change, behavior change to impact infections, infections to impact testing, and data to be collected, etc?)

Review version with references removed; NOT FOR DISTRIBUTION

	-Check if authors appropriately and directly account for these time effects.
Is this policy the only thing to happen which could have impacted the outcome during the measurement period, differently for policy and non-policy regions??	<ul style="list-style-type: none"> -Consider other policies or interventions which could impact the outcome during this time. -Consider social behaviors changed which could meaningfully impact the outcome during this time. -Consider economic conditions changed which could meaningfully impact the outcome during this time. -Note that the actual concurrent changes do not need to happen during the period of measurement, just their effects.

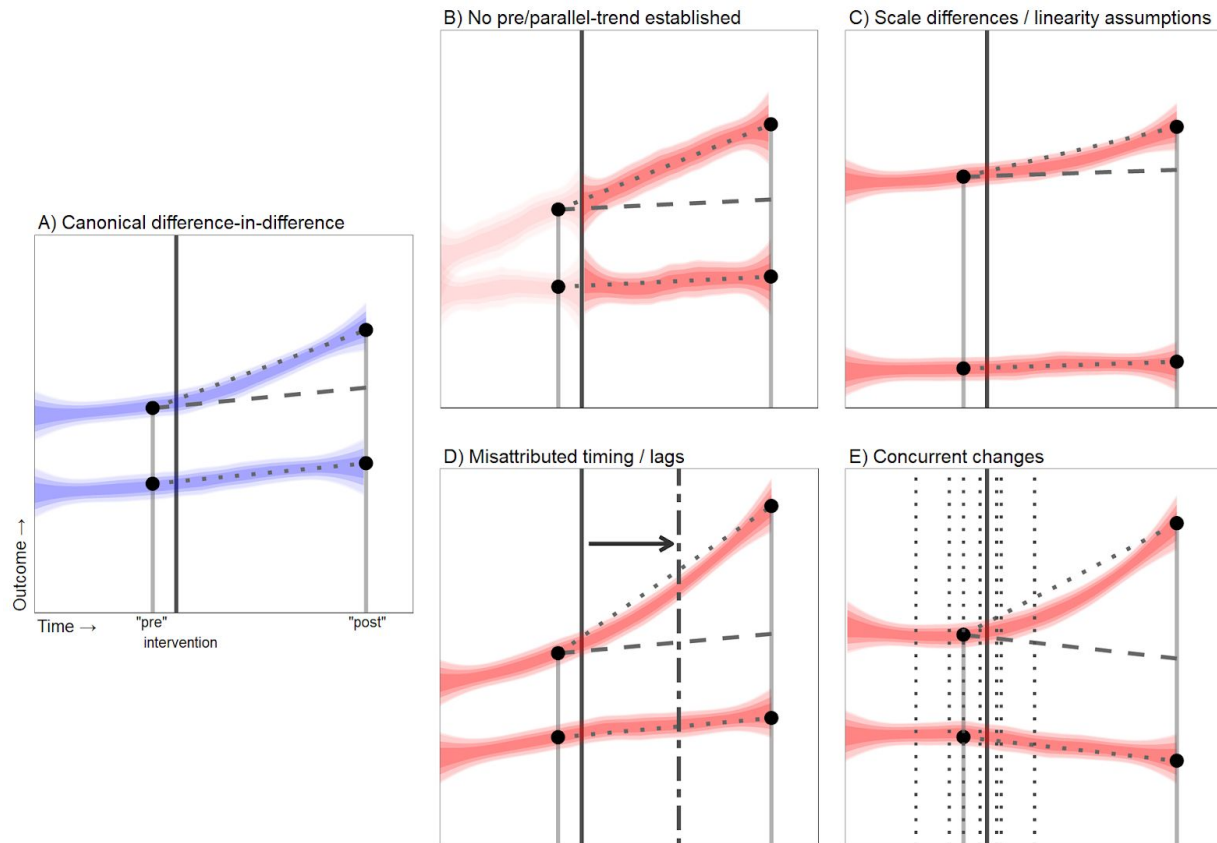
These issues are summarized as a checklist of questions to identify common pitfalls in Table 2.

Difference-in-differences

The difference-in-difference (DiD) approach uses concurrent non-intervention groups as a counterfactual. Typically, this consists of one set of units (e.g., regions) that had the intervention and one set that did not, with each measured before and after the intervention took place. DiD is more directly analogous to a non-randomized medical study with at least one treatment and control group but limited observation before and after treatment. In contrast to ITS, which compares a unit with itself over time, DiD compares differences between treatment arms or units at two observation points. In many analyses, a DiD approach is implied by comparing regions over time, without formally naming or modelling it. Other DiD approaches use interventions implemented at multiple time points.

Figure 3: Difference-in-differences graphical guidance for identifying common pitfalls

Review version with references removed; NOT FOR DISTRIBUTION



This chart shows one canonical design for DiD (blue, Panel A) and four panels demonstrating common issues with DiD analysis (red, panels B-E). In all cases: the blue/red shading represents the underlying data trends, the vertical grey line represents the time of intervention, the grey dashed lines represent the assumed counterfactual in the absence of the intervention. In panel D) the dash-dot line represents the time at which the policy is expected to impact the outcome. In panel E), the vertical dotted lines represent concurrent events and changes.

One key component of the standard DiD approach is the parallel counterfactual trends assumption: that the intervention and comparison groups would have had parallel trends over time in the absence of the intervention. In some cases, the parallel trends assumption may be referenced or examined implicitly but not named.

Ideally, pre-intervention trends would be shown to be clearly identifiable, stable, of a similar level, and parallel between groups. With only one observation before and only one after the intervention, assessment of the plausibility of the parallel counterfactual trends assumption is not possible. Absent this confirmation the evaluation runs the risk of biased estimation due to differential pre-trends (Figure 3B). Pre-trends approaching the ceiling or floor may also not be informative about stable and parallel pre-trends. Empirical assessment of whether pre-intervention trends were parallel and stable between groups is possible when multiple observations are available at multiple time points before the intervention, noting that this can

Review version with references removed; NOT FOR DISTRIBUTION

begin to resemble a CITS design. In this scenario, pre-trend data should be visually and statistically established and documented. While parallel trends before intervention (which we can observe and may be testable) do not guarantee parallel *counterfactual* trends in the post-intervention period (which we cannot observe and are generally untestable), examining pre-intervention parallel trends is a minimal requirement for DiD reliability.

It is also important to consider the scale and level on which the outcome is measured (Figure 3C). As with ITS, if the outcomes in the treatment and comparison groups are moving in parallel on a logged scale, they will not be moving in parallel on a natural scale. Level differences by themselves may be a problem for COVID-19 outcomes, as infectious disease transmission dynamics dictate that infection risks are related to the prevalence of infected people in a population, i.e. the rate of change is linked intrinsically to the level. A population with an extremely low prevalence will tend to have an inherently slower rise in infection rates than an otherwise identical population with merely a low prevalence. Just as importantly, large level differences in the outcome between intervention and comparison groups is often indicative of other important differences between comparators, which may result in other assumptions being violated.

While DiD is in some ways more robust to very specific kinds of timing effects (Figure 3D) and concurrent changes (Figure 3E), it also introduces additional risks. DiD effectively doubles the opportunity for concurrent changes to spuriously impact results, since they can occur in the treatment or comparison groups. As above, this can become even more problematic for DiD in the typical case where intervention groups enact more or very contextually different policies than non-intervention groups. Even cases where concurrent changes happen equally in both treatment and comparison groups can lead to overwhelming bias, particularly when approaching the maximum or minimum levels of the outcome. If either the treatment or control group is approaching the floor (e.g. 0% prevalence) or ceiling for an outcome of interest due to other policies concurrent in both places (e.g. national lockdowns, but region-level differences in mask policy), this can lead to bias when comparing changes between the two groups.

Table 3: Checklist for identifying common pitfalls for DiD to evaluate COVID-19 policy

Key design questions. If any answer is “no,” this analysis is unlikely to be appropriate or useful for estimating the impact of the intervention of interest	Details and suggestions for inspection:
Does the analysis provide graphical representation of the outcome over time?	-Check for a graph that shows the outcome over time for all groups, with the dates of interest. Outcomes may be aggregated for clarity (e.g. mean and CI at discrete time points).
Is there sufficient pre-intervention data to observe both pre and post trends in the data?	-Check the chart(s) to see if there are several time points over a reasonable period of time over which to establish stability and curvature in the pre- and post- trends.
Are the pre-trends stable?	-Check if there are sufficient graphical data to reasonably determine a stable functional form for the pre-trends, and that they follow a modelable functional form.

Review version with references removed; NOT FOR DISTRIBUTION

Are the pre-trends parallel?	-Observe if the trends in the intervention and comparison groups appear to move together at the same rate at the same time.
Are the pre-trends at a similar level?	-Check if the trends in the intervention and comparison groups are at similar levels. -Note that non-level trends exacerbates other problems with the analysis, including linearity assumptions
Are intervention and non-groups broadly comparable?	-Consider areas where comparison groups may be dissimilar for comparison beyond just the level of the outcome.
Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?	-Check whether the authors justify the use of the date threshold relative to the date of the intervention. -Trace the process between the intervention being put in place to when observable effects in the outcome might appear over time. -Consider whether there are anticipation effects (e.g. do people change behaviors before the date when the intervention begins?) -Consider whether there are lag effects. (e.g. does it take time for behaviors to change, behavior change to impact infections, infections to impact testing, and data to be collected, etc?) -Check if authors appropriately and directly account for these time effects.
Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have changed during the measurement period, differently for policy and non-policy regions?	-Consider any uncontrolled factor which could have influenced the outcome differently in policy and non-policy regions. -This may include (but is not limited to) -Other policies -Social behaviors -Economic conditions -Are these factors justified as having negligible impact? -If justified, is the argument that these have negligible impact convincing? -Note that the actual concurrent changes do not need to happen during the period of measurement, just their effects.

Similarly to the ITS section, these issues are summarized as a checklist of questions to identify common pitfalls in Table 3.

Discussion

In recent months, there has been a proliferation of research evaluating policies related to the COVID-19 pandemic. As with other areas of COVID-19 research, quality has been highly variable, with low quality studies resulting in poorly or mis-informed policy decisions, wasted resources, and undermined trust in research. To support high quality policy evaluations, in this paper we describe common approaches to evaluating policies using observational data, and describe key issues that can arise in applying these approaches. We hope that this guidance can help support researchers, editors, reviewers, and decision-makers in conducting high quality policy evaluations and in assessing the strength of the evidence that has already been published.

Policy evaluation — far from a simple task in normal circumstances — is particularly challenging during a pandemic. Cross-sectional comparisons of states or countries are likely to be biased by selection into treatment: for example, countries with worse outbreaks may be more likely to

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 implement policies such as mask requirements. In analyses of changes over time – such as
4 single-unit studies using interrupted time-series or multi-unit comparisons using
5 difference-in-differences or comparative interrupted time-series – it may not be possible to parse
6 apart the effects of different policies implemented around the same time, such as mask
7 mandates paired with limits on social gatherings. Analyses of changes over time may also be
8 biased if disease or human behavioral dynamics are not modelled appropriately. This can be
9 challenging because case counts typically do not grow linearly and there is often a lag between
10 a policy change and a behavioral response.
11
12
13

14 This guidance should be considered minimal screening to identify low quality policy impact
15 evaluation in COVID-19, but is in no way sufficient to identify high quality evidence or
16 actionability. Decision-makers and researchers should pay particular attention to the relevance
17 of the intervention as it was evaluated to relevant decisions being made. The evaluated impact
18 of a program encouraging mask use through messages might not be informative about mask
19 requirement orders. Differences in level of aggregation may be important, such as ecological
20 fallacy arising from a situation in which areas with higher overall mask use have higher
21 transmission, but transmission is actually lower for individuals wearing masks. Policy impact
22 evaluation is only as useful as the question it asks, data it uses, and the way it is analyzed.
23 Problems with measurement, spillover effects, generalizability, changes in measurement
24 overtime (e.g. varying test availability), statistics, testing robustness to alternative assumptions,
25 and many issues can undermine an otherwise robust evaluation, and are not discussed here.
26
27
28
29

30 While this guidance is not comprehensive, it may help inform study designs not covered here.
31 Issues with comparative interrupted time-series and synthetic control methods, for example, are
32 broadly similar to the issues with difference-in-differences analyses we discuss here. Other
33 approaches may include adjustment and matching based observational causal inference
34 designs, instrumental variables and related quasi-experimental approaches, and randomized
35 controlled trials. Each has its own set of practical, ethical, and inferential limitations.
36
37
38

39 In the face of these challenges, we recommend careful scrutiny and attention to potential
40 sources of bias in COVID-19-related policy evaluations, but we remain optimistic about the
41 potential for robust evaluations to inform decision-making. Researchers and decision-makers
42 should triangulate across a large variety of approaches from theory to evidence, invest in better
43 data and more reliable and useful evidence wherever feasible, clearly acknowledge limitations
44 and potential sources of bias, and acknowledge when actionable evidence is not feasible. We
45 anticipate increasing opportunities for better examining policies moving forward, particularly if
46 policies and interventions are designed with policy impact evaluation and data collection in
47 mind.
48
49
50

51 The COVID-19 pandemic requires urgent decisions about policies that affect millions of people's
52 lives in significant ways. High quality evidence on the effects of these policies is critical to
53 informing decision-making, but is very hard to generate. Evidence-based decision-making
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 depends on research that carefully considers potential sources of bias, and clearly
4 communicates underlying assumptions and sources of uncertainty.
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Appendix 1: Changes from pre-registered protocol and justifications

The full, original pre-registered protocol is available here: <https://osf.io/7nbk6>

Inclusion criteria

Minor language edits were made to the inclusion criteria to improve clarity and fix grammatical and typographical errors. This largely centered around improving clarity that a study must estimate the quantitative impact of policies that had already been enacted. The word “quantitative” was not explicitly stated in the original version.

Procedures

The original protocol specified that each article would receive two independent reviewers. This was increased to three reviewers per article once it became clear both that the number of articles which would be accepted for full review was lower than expectations, and that there would be substantial differences in opinion between reviewers.

Statistical analysis

Firstly, the original protocol specified that 95% confidence intervals would be calculated. However, after further discussion and review, we determined that sampling-based confidence intervals were not appropriate. Our results are not indicative nor intended to be representative of any super- or target-population, and as such sampling-based error is not an appropriate metric for the conclusions of this study.

Secondly, the original protocol specified Kappa-based interrater reliability statistics. However, using three reviewers, rather than the originally registered two, meant that most Kappa statistics would not be appropriate for our review process. Given the three-rater, four-level ordinal scale used, we opted instead to use Krippendorff’s Alpha.

Review tool

A number of changes were made to the review tool during the course of the review process. While the original protocol included logic to allow pre/post for review in some of the key questions, this was removed for consistency with the guidance document.

The remaining changes to the review tool were error corrections and clarifications (e.g. correcting the text for the concurrent changes sections in difference-in-differences so that it

1
2
3 stated “uncontrolled” concurrent changes, and distinguishing the DiD/CITS requirements from
4 the ITS requirements to emphasize differential concurrent changes).
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Appendix 2: Full search terms

Note: The search filter for COVID-19 and SARS-CoV-2 were the exact search terms used for the National Library of Medicine one-click search option at the time of the protocol development and when the search took place. This reflects that some of the early literature referred to Wuhan specifically (both in geographic reference for where the SARS-CoV-2 was initially found, and unfortunately also early naming of the virus/disease) before official naming conventions became ubiquitous in the literature. In order to comprehensively capture the literature and use searching best practices, we used the most standard and recommended terms.

```
((((wuhan[All Fields] AND ("coronavirus"[MeSH Terms] OR "coronavirus"[All Fields])) AND 2019/12[PDAT] : 2030[PDAT]) OR 2019-nCoV[All Fields] OR 2019nCoV[All Fields] OR COVID-19[All Fields] OR SARS-CoV-2[All Fields])
```

```
AND ("impact"[TIAB] OR "effect"[TIAB])
```

```
AND ("policy"[TIAB] OR "policies"[TIAB] OR "order"[TIAB] OR "mandate"[TIAB])
```

```
AND ("countries"[TIAB] OR "country"[TIAB] OR "state"[TIAB] OR "provinc"[TIAB] OR "county"[TIAB] OR "parish"[TIAB] OR "region"[TIAB] OR "city"[TIAB] OR "cities"[TIAB] OR "continent"[TIAB] "Asia"[TIAB] OR "Europe"[TIAB] OR "Africa"[TIAB] OR "America"[TIAB] OR "Australia"[TIAB] OR "Antarctica"[TIAB] OR "Afghanistan"[TIAB] OR "Aland Islands"[TIAB] OR "Åland Islands"[TIAB] OR "Albania"[TIAB] OR "Algeria"[TIAB] OR "American Samoa"[TIAB] OR "Andorra"[TIAB] OR "Angola"[TIAB] OR "Anguilla"[TIAB] OR "Antarctica"[TIAB] OR "Antigua"[TIAB] OR "Argentina"[TIAB] OR "Armenia"[TIAB] OR "Aruba"[TIAB] OR "Australia"[TIAB] OR "Austria"[TIAB] OR "Azerbaijan"[TIAB] OR "Bahamas"[TIAB] OR "Bahrain"[TIAB] OR "Bangladesh"[TIAB] OR "Barbados"[TIAB] OR "Barbuda"[TIAB] OR "Belarus"[TIAB] OR "Belgium"[TIAB] OR "Belize"[TIAB] OR "Benin"[TIAB] OR "Bermuda"[TIAB] OR "Bhutan"[TIAB] OR "Bolivia"[TIAB] OR "Bonaire"[TIAB] OR "Bosnia"[TIAB] OR "Botswana"[TIAB] OR "Bouvet Island"[TIAB] OR "Brazil"[TIAB] OR "British Indian Ocean Territory"[TIAB] OR "Brunei"[TIAB] OR "Bulgaria"[TIAB] OR "Burkina Faso"[TIAB] OR "Burundi"[TIAB] OR "Cabo Verde"[TIAB] OR "Cambodia"[TIAB] OR "Cameroon"[TIAB] OR "Canada"[TIAB] OR "Cayman Islands"[TIAB] OR "Central African Republic"[TIAB] OR "Chad"[TIAB] OR "Chile"[TIAB] OR "China"[TIAB] OR "Christmas Island"[TIAB] OR "Cocos Islands"[TIAB] OR "Colombia"[TIAB] OR "Comoros"[TIAB] OR "Congo"[TIAB] OR "Congo"[TIAB] OR "Cook Islands"[TIAB] OR "Costa Rica"[TIAB] OR "Côte d'Ivoire"[TIAB] OR "Croatia"[TIAB] OR "Cuba"[TIAB] OR "Curaçao"[TIAB] OR "Cyprus"[TIAB] OR "Czechia"[TIAB] OR "Denmark"[TIAB] OR "Djibouti"[TIAB] OR "Dominica"[TIAB] OR "Dominican Republic"[TIAB] OR "Ecuador"[TIAB] OR "Egypt"[TIAB] OR "El Salvador"[TIAB] OR "Equatorial Guinea"[TIAB] OR "Eritrea"[TIAB] OR "Estonia"[TIAB] OR "Eswatini"[TIAB] OR "Ethiopia"[TIAB] OR "Falkland Islands"[TIAB] OR "Faroe Islands"[TIAB] OR "Fiji"[TIAB] OR "Finland"[TIAB] OR "France"[TIAB] OR "French Guiana"[TIAB] OR "French Polynesia"[TIAB] OR "French Southern
```

1
2
3 Territories[TIAB] OR "Futuna"[TIAB] OR "Gabon"[TIAB] OR "Gambia"[TIAB] OR
4 "Georgia"[TIAB] OR "Germany"[TIAB] OR "Ghana"[TIAB] OR "Gibraltar"[TIAB] OR
5 "Greece"[TIAB] OR "Greenland"[TIAB] OR "Grenada"[TIAB] OR "Grenadines"[TIAB] OR
6 "Guadeloupe"[TIAB] OR "Guam"[TIAB] OR "Guatemala"[TIAB] OR "Guernsey"[TIAB] OR
7 "Guinea"[TIAB] OR "Guinea-Bissau"[TIAB] OR "Guyana"[TIAB] OR "Haiti"[TIAB] OR "Heard
8 Island"[TIAB] OR "Herzegovina"[TIAB] OR "Holy See"[TIAB] OR "Honduras"[TIAB] OR "Hong
9 Kong"[TIAB] OR "Hungary"[TIAB] OR "Iceland"[TIAB] OR "India"[TIAB] OR "Indonesia"[TIAB]
10 OR "Iran"[TIAB] OR "Iraq"[TIAB] OR "Ireland"[TIAB] OR "Isle of Man"[TIAB] OR "Israel"[TIAB]
11 OR "Italy"[TIAB] OR "Jamaica"[TIAB] OR "Jan Mayen Islands"[TIAB] OR "Japan"[TIAB] OR
12 "Jersey"[TIAB] OR "Jordan"[TIAB] OR "Kazakhstan"[TIAB] OR "Keeling Islands"[TIAB] OR
13 "Kenya"[TIAB] OR "Kiribati"[TIAB] OR "Korea"[TIAB] OR "Korea"[TIAB] OR "Kuwait"[TIAB] OR
14 "Kyrgyzstan"[TIAB] OR "Lao People's Democratic Republic"[TIAB] OR "Laos"[TIAB] OR
15 "Latvia"[TIAB] OR "Lebanon"[TIAB] OR "Lesotho"[TIAB] OR "Liberia"[TIAB] OR "Libya"[TIAB]
16 OR "Liechtenstein"[TIAB] OR "Lithuania"[TIAB] OR "Luxembourg"[TIAB] OR "Macao"[TIAB] OR
17 "Madagascar"[TIAB] OR "Malawi"[TIAB] OR "Malaysia"[TIAB] OR "Maldives"[TIAB] OR
18 "Mali"[TIAB] OR "Malta"[TIAB] OR "Malvinas"[TIAB] OR "Marshall Islands"[TIAB] OR
19 "Martinique"[TIAB] OR "Mauritania"[TIAB] OR "Mauritius"[TIAB] OR "Mayotte"[TIAB] OR
20 "McDonald Islands"[TIAB] OR "Mexico"[TIAB] OR "Micronesia"[TIAB] OR "Moldova"[TIAB] OR
21 "Monaco"[TIAB] OR "Mongolia"[TIAB] OR "Montenegro"[TIAB] OR "Montserrat"[TIAB] OR
22 "Morocco"[TIAB] OR "Mozambique"[TIAB] OR "Myanmar"[TIAB] OR "Namibia"[TIAB] OR
23 "Nauru"[TIAB] OR "Nepal"[TIAB] OR "Netherlands"[TIAB] OR "Nevis"[TIAB] OR "New
24 Caledonia"[TIAB] OR "New Zealand"[TIAB] OR "Nicaragua"[TIAB] OR "Niger"[TIAB] OR
25 "Nigeria"[TIAB] OR "Niue"[TIAB] OR "Norfolk Island"[TIAB] OR "North Macedonia"[TIAB] OR
26 "Northern Mariana Islands"[TIAB] OR "Norway"[TIAB] OR "Oman"[TIAB] OR "Pakistan"[TIAB]
27 OR "Palau"[TIAB] OR "Panama"[TIAB] OR "Papua New Guinea"[TIAB] OR "Paraguay"[TIAB]
28 OR "Peru"[TIAB] OR "Philippines"[TIAB] OR "Pitcairn"[TIAB] OR "Poland"[TIAB] OR
29 "Portugal"[TIAB] OR "Principe"[TIAB] OR "Puerto Rico"[TIAB] OR "Qatar"[TIAB] OR
30 "Réunion"[TIAB] OR "Romania"[TIAB] OR "Russian Federation"[TIAB] OR "Rwanda"[TIAB] OR
31 "Saba"[TIAB] OR "Saint Barthélemy"[TIAB] OR "Saint Helena"[TIAB] OR "Saint Kitts"[TIAB] OR
32 "Saint Lucia"[TIAB] OR "Saint Martin"[TIAB] OR "Saint Pierre and Miquelon"[TIAB] OR "Saint
33 Vincent"[TIAB] OR "Samoa"[TIAB] OR "San Marino"[TIAB] OR "Sao Tome"[TIAB] OR
34 "Sark"[TIAB] OR "Saudi Arabia"[TIAB] OR "Senegal"[TIAB] OR "Serbia"[TIAB] OR
35 "Seychelles"[TIAB] OR "Sierra Leone"[TIAB] OR "Singapore"[TIAB] OR "Sint Eustatius"[TIAB]
36 OR "Sint Maarten"[TIAB] OR "Slovakia"[TIAB] OR "Slovenia"[TIAB] OR "Solomon
37 Islands"[TIAB] OR "Somalia"[TIAB] OR "South Africa"[TIAB] OR "South Georgia"[TIAB] OR
38 "South Sandwich Islands"[TIAB] OR "South Sudan"[TIAB] OR "Spain"[TIAB] OR "Sri
39 Lanka"[TIAB] OR "State of Palestine"[TIAB] OR "Sudan"[TIAB] OR "Suriname"[TIAB] OR
40 "Svalbard"[TIAB] OR "Sweden"[TIAB] OR "Switzerland"[TIAB] OR "Syria"[TIAB] OR "Syrian
41 Arab Republic"[TIAB] OR "Tajikistan"[TIAB] OR "Thailand"[TIAB] OR "Timor-Leste"[TIAB] OR
42 "Tobago"[TIAB] OR "Togo"[TIAB] OR "Tokelau"[TIAB] OR "Tonga"[TIAB] OR "Trinidad"[TIAB]
43 OR "Tunisia"[TIAB] OR "Turkey"[TIAB] OR "Turkmenistan"[TIAB] OR "Turks and Caicos"[TIAB]
44 OR "Tuvalu"[TIAB] OR "Uganda"[TIAB] OR "UK"[TIAB] OR "Ukraine"[TIAB] OR "United Arab
45 Emirates"[TIAB] OR "United Kingdom"[TIAB] OR "United Republic of Tanzania"[TIAB] OR
46 "United States Minor Outlying Islands"[TIAB] OR "United States of America"[TIAB] OR
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 "Uruguay"[TIAB] OR "USA"[TIAB] OR "Uzbekistan"[TIAB] OR "Vanuatu"[TIAB] OR
4 "Venezuela"[TIAB] OR "Viet Nam"[TIAB] OR "Vietnam"[TIAB] OR "Virgin Islands"[TIAB] OR
5 "Virgin Islands"[TIAB] OR "Wallis"[TIAB] OR "Western Sahara"[TIAB] OR "Yemen"[TIAB] OR
6 "Zambia"[TIAB] OR "Zimbabwe"[TIAB] OR "Alabama"[TIAB] OR "Alaska"[TIAB] OR
7 "Arizona"[TIAB] OR "Arkansas"[TIAB] OR "California"[TIAB] OR "Colorado"[TIAB] OR
8 "Connecticut"[TIAB] OR "Delaware"[TIAB] OR "Florida"[TIAB] OR "Georgia"[TIAB] OR
9 "Hawaii"[TIAB] OR "Idaho"[TIAB] OR "Illinois"[TIAB] OR "Indiana"[TIAB] OR "Iowa"[TIAB] OR
10 "Kansas"[TIAB] OR "Kentucky"[TIAB] OR "Louisiana"[TIAB] OR "Maine"[TIAB] OR
11 "Maryland"[TIAB] OR "Massachusetts"[TIAB] OR "Michigan"[TIAB] OR "Minnesota"[TIAB] OR
12 "Mississippi"[TIAB] OR "Missouri"[TIAB] OR "Montana"[TIAB] OR "Nebraska"[TIAB] OR
13 "Nevada"[TIAB] OR "New Hampshire"[TIAB] OR "New Jersey"[TIAB] OR "New Mexico"[TIAB]
14 OR "New York"[TIAB] OR "North Carolina"[TIAB] OR "North Dakota"[TIAB] OR "Ohio"[TIAB] OR
15 "Oklahoma"[TIAB] OR "Oregon"[TIAB] OR "Pennsylvania"[TIAB] OR "Rhode Island"[TIAB] OR
16 "South Carolina"[TIAB] OR "South Dakota"[TIAB] OR "Tennessee"[TIAB] OR "Texas"[TIAB] OR
17 "Utah"[TIAB] OR "Vermont"[TIAB] OR "Virginia"[TIAB] OR "Washington"[TIAB] OR "West
18 Virginia"[TIAB] OR "Wisconsin"[TIAB] OR "Wyoming"[TIAB] OR "Ontario"[TIAB] OR
19 "Quebec"[TIAB] OR "Nova Scotia"[TIAB] OR "New Brunswick"[TIAB] OR "Manitoba"[TIAB] OR
20 "British Columbia"[TIAB] OR "Prince Edward Island"[TIAB] OR "Saskatchewan"[TIAB] OR
21 "Alberta"[TIAB] OR "Newfoundland"[TIAB] OR "Labrador"[TIAB])
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

BMJ Open

Problems with Evidence Assessment in COVID-19 Health Policy Impact Evaluation: A systematic review of study design and evidence strength

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-053820.R1
Article Type:	Original research
Date Submitted by the Author:	14-Sep-2021
Complete List of Authors:	<p>Haber, Noah ; Stanford University, Clarke-Deelder, Emma; Harvard University T H Chan School of Public Health, Department of Global Health and Population Feller, Avi; University of California Berkeley, Goldman School of Public Policy Smith, Emily; George Washington University School of Public Health and Health Services, Department of Global Health Salomon, Joshua ; Stanford University MacCormack-Gelles, Benjamin; Harvard University T H Chan School of Public Health, Department of Global Health and Population Stone, Elizabeth M.; Johns Hopkins University Bloomberg School of Public Health, Department of Health Policy and Management Bolster-Foucault, Clara; McGill University, Epidemiology, Biostatistics, and Occupational Health Daw, Jamie R.; Columbia University Mailman School of Public Health, Health Policy and Management Hatfield, Laura; Harvard Medical School, Biostatistics Fry, Carrie E.; Vanderbilt University, Department of Health Policy Boyer, Christopher B.; Harvard University T H Chan School of Public Health, Department of Epidemiology Ben-Michael, Eli; University of California Berkeley, Department of Statistics Joyce, Caroline M.; McGill University, Epidemiology, Biostatistics, and Occupational Health Linas, Beth S.; Johns Hopkins University Bloomberg School of Public Health, Department of Epidemiology; MITRE Corp Schmid, Ian; Johns Hopkins University Bloomberg School of Public Health, Department of Mental Health Au, Eric; The University of Sydney, School of Public Health Wieten, Sarah; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Jarrett, Brooke; Johns Hopkins University, Epidemiology Axfors, Cathrine; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Nguyen, Van; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Griffin, Beth; RAND Corp, Bilinski, Alyssa; Harvard University Graduate School of Arts and Sciences Stuart, Elizabeth A; Johns Hopkins University Bloomberg School of Public</p>

	Health, Department of Mental Health
Primary Subject Heading :	Health policy
Secondary Subject Heading :	Research methods, Public health, Epidemiology, Global health
Keywords :	COVID-19, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

Problems with Evidence Assessment in COVID-19 Health Policy Impact Evaluation: A systematic review of study design and evidence strength

Noah A. Haber, ScD¹, Emma Clarke-Deelder, MPhil², Avi Feller, PhD³, Emily R. Smith, ScD⁴, Joshua Salomon, PhD⁵, Benjamin MacCormack-Gelles, MS², Elizabeth M. Stone, MS⁶, Clara Bolster-Foucalt, MScPH⁷, Jamie R. Daw, PhD⁸, Laura A. Hatfield, PhD⁹, Carrie E. Fry, PhD¹⁰, Christopher B. Boyer, MPH¹¹, Eli Ben-Michael, PhD¹², Caroline M. Joyce, MPH⁷, Beth S. Linas, PhD, MHS^{13,14}, Ian Schmid, ScM¹⁵, Eric H. Au, MPH¹⁶, Sarah E. Wieten, PhD¹, Brooke A Jarrett, MSPH¹³, Cathrine Axfors, MD, PhD¹, Van Thu Nguyen, PhD¹, Beth Ann Griffin, PhD¹⁷, Alyssa Bilinski, MS¹⁸, Elizabeth A. Stuart, PhD¹⁵

1. Meta-Research Innovation Center at Stanford (METRICS), Stanford University, Stanford, CA, USA
2. Department of Global Health and Population, Harvard T. H. Chan School of Public Health, Boston, MA, USA
3. Goldman School of Public Policy, UC Berkeley, Berkeley, CA, USA
4. Department of Global Health, Milken Institute School of Public Health, George Washington University, Washington, D.C, USA
5. Center for Health Policy and Center for Primary Care and Outcomes Research, Stanford University, Stanford, CA, USA
6. Department of Health Policy and Management, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
7. Epidemiology, Biostatistics, and Occupational Health, McGill University, Montreal, Canada
8. Health Policy and Management, Columbia University Mailman School of Public Health, New York, NY, USA
9. Department of Health Care Policy, Harvard Medical School, Boston, MA, USA
10. Department of Health Policy, Vanderbilt University, Nashville, TN, USA
11. Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA
12. Institute for Quantitative Social Science, Harvard University, Cambridge, MA, USA
13. Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
14. Clinical Quality and Informatics, MITRE Corp, McLean, VA, USA
15. Department of Mental Health, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
16. School of Public Health, University of Sydney, Sydney, Australia
17. RAND Corporation, Arlington, VA, USA
18. Interfaculty Initiative in Health Policy, Harvard Graduate School of Arts and Sciences, Cambridge, MA, USA

Corresponding author:

Noah A. Haber, ScD

noahhaber@stanford.edu

Meta Research Innovation Center at Stanford University

Stanford University

1265 Welch Rd

Palo Alto, CA 94305

(650) 497-0811

Abstract

Introduction: Assessing the impact of COVID-19 policy is critical for informing future policies. However, there are concerns about the overall strength of COVID-19 impact evaluation studies given the circumstances for evaluation and concerns about the publication environment. This study systematically reviewed the strength of evidence in the published COVID-19 policy impact evaluation literature.

Methods: We included studies that were primarily designed to estimate the quantitative impact of one or more implemented COVID-19 policies on direct SARS-CoV-2 and COVID-19 outcomes. After searching PubMed for peer-reviewed articles published on November 26, 2020 or earlier and screening, all studies were reviewed by three reviewers first independently and then to consensus. The review tool was based on previously developed and released review guidance for COVID-19 policy impact evaluation, assessing what impact evaluation method was used, graphical display of outcomes data, functional form for the outcomes, timing between policy and impact, concurrent changes to the outcomes, and an overall rating.

Results: After 102 articles were identified as potentially meeting inclusion criteria, we identified 36 published articles that evaluated the quantitative impact of COVID-19 policies on direct COVID-19 outcomes. The majority (n=23/36) of studies in our sample examined the impact of stay-at-home requirements. Nine studies were set aside because the study design was considered inappropriate for COVID-19 policy impact evaluation (n=8 pre/post; n=1 cross-section), and 27 articles were given a full consensus assessment. 20/27 met criteria for graphical display of data, 5/27 for functional form, 19/27 for timing between policy implementation and impact, and only 3/27 for concurrent changes to the outcomes. Only 1/27 studies passed all of the above checks, and 4/27 were rated as overall appropriate. Including the 9 studies set aside, reviewers found that only four of the 36 identified published and peer-reviewed health policy impact evaluation studies passed a set of key design checks for identifying the causal impact of policies on COVID-19 outcomes.

Discussion: The reviewed literature directly evaluating the impact of COVID-19 policies largely failed to meet key design criteria for inference of sufficient rigor to be actionable by policy-makers. This was largely driven by the circumstances under which policies were passed making it difficult to attribute changes in COVID-19 outcomes to particular policies. More reliable evidence review is needed to both identify and produce policy-actionable evidence, alongside the recognition that actionable evidence is often unlikely to be feasible.

Strengths and limitations

- This study is based on previously released review guidance for discerning and evaluating critical minimal methodological design aspects of the COVID-19 health policy impact evaluation.
- The review tool assesses critical aspects of study design grounded in impact evaluation methods that must be true for the papers to provide useful policy impact evaluation, including what type of impact evaluation method was used, graphical display of outcomes data, functional form for the outcomes, timing between policy and impact, concurrent changes to the outcomes, and an overall rating.
- This study used a consensus reviewer model with three reviewers in order to obtain replicable results for study strength ratings, noting that ratings were based on subjective assessments of key strength of evidence criteria, which may change between reviewers.
- While the vast majority of studies in our sample received low ratings for useful causal policy impact evaluation, they may make other contributions to the literature.
- Because our review tool was limited to a very narrow - albeit critical - set of items, weaknesses in other aspects not reviewed (e.g. data quality or other aspects of statistical inference) may further weaken studies that were found to meet our criteria.

Introduction

Policy decisions to mitigate the impact of COVID-19 on morbidity and mortality are some of the most important issues policymakers have had to make since January 2020. Decisions regarding which policies are enacted depend in part on the evidence base for those policies, including understanding what impact past policies had on COVID-19 outcomes[1,2] Unfortunately, there are substantial concerns that much of the existing literature may be methodologically flawed, which could render its conclusions unreliable for informing policy. The combination of circumstances being difficult for strong impact evaluation, the importance of the topic, and concerns over the publication environment may lead to the proliferation of low strength studies.

High-quality causal evidence requires a combination of rigorous methods, clear reporting, appropriate caveats, and the appropriate circumstances for the methods used.[3–6] Rigorous evidence is difficult in the best of circumstances, and the circumstances for evaluating non-pharmaceutical intervention (NPI) policy effects on COVID-19 are particularly challenging.[5] The global pandemic has yielded a combination of a large number of concurrent policy and non-policy changes, complex infectious disease dynamics, and unclear timing between policy implementation and impact; all of this makes isolating the causal impact of any particular policy or policies exceedingly difficult.[7]

The scientific literature on COVID-19 is exceptionally large and fast growing. Scientists published more than 100,000 papers related to COVID-19 in 2020.[8] There is some general concern that the volume and speed[9,10] at which this work has been produced may result in a literature that is overall low quality and unreliable.[11–15]

1
2
3
4 Given the importance of the topic, it is critical that decision-makers are able to understand what
5 is known and knowable[5,16] from observational data in COVID-19 policy, as well as what is
6 unknown and/or unknowable.
7

8
9 Motivated by concerns about the methodological strength of COVID-19 policy evaluations, we
10 set out to review the literature using a set of methodological design checks tailored to common
11 policy impact evaluation methods. Our primary objective was to evaluate each paper for
12 methodological strength and reporting, based on pre-existing review guidance developed for
13 this purpose.[17] As a secondary objective, we also studied our own process: examining the
14 consistency, ease of use, and clarity of this review guidance.
15
16

17
18 This protocol differs in several ways from more traditional systematic review protocols given the
19 atypical objectives and scope of the systematic review. First, this is a systematic review of
20 methodological strength of evidence for a given literature as opposed to a review summary of
21 the evidence of a particular topic. As such, we do not summarize and attempt to combine the
22 results for any of the literature. Second, rather than being a comprehensive review of every
23 possible aspect of what might be considered “quality,” this is a review of targeted critical design
24 features for actionable inference for COVID-19 policy impact evaluation and methods. It is
25 designed to be a set of broad criteria for minimal plausibility of actionable causal inference,
26 where each of the criteria is necessary but not sufficient for strong design. Issues in other
27 domains (data, details of the design, statistics, etc) further reduce overall actionability and
28 quality, and thorough review in those domains is needed for any studies passing our basic
29 minimal criteria. Third, because the scope relies on guided, but difficult and subjective
30 assessments of methodological appropriateness, we utilize a discussion-based consensus
31 process to arrive at consistent and replicable results, rather than a more common model with
32 two independent reviewers with conflict resolution. The independent review serves primarily as
33 a starting point for discussion, but is neither designed nor expected to be a strong indicator of
34 the overall consensus ratings of the group.
35
36
37
38
39

40 Methods

41 Overview

42
43
44 This protocol and study was written and developed following the release of the review guidance
45 written by the author team in September 2020 on which the review tool is based. The protocol
46 for this study was pre-registered on OSF.io[18] in November 2020 following PRISMA
47 guidelines.[19] Deviations from the original protocol are discussed in Appendix 1, and consisted
48 largely of language clarifications and error corrections for both the inclusion criteria and review
49 tool, an increase in the number of reviewers per fully reviewed article from two to three, and
50 simplification of the statistical methods used to assess the data.
51
52
53
54
55
56
57
58
59
60

1
2
3 For this study, we ascertain minimal criteria for studies to be able to plausibly identify causal
4 effects of policies, which is the information of greatest interest to inform policy decisions. The
5 causal estimand is something that, if known, would definitely help policy makers decide what to
6 do (e.g., whether to implement or discontinue a policy). The study estimates that target causal
7 quantity with a rigorous design and appropriate data in a relevant population/sample. For
8 shorthand, we refer to this as minimal properties of “actionable” evidence.
9
10

11 This systematic review of the strength of evidence took place in three phases: search,
12 screening, and full review.
13
14

15 Eligibility criteria

16 The following eligibility criteria were used to determine the papers to include:
17

- 18 ● The primary topic of the article must be evaluating one or more individual COVID-19 or SARS-CoV-2
19 policies on direct COVID-19 or SARS-CoV-2 outcomes
 - 20 ○ The primary exposure(s) must be a policy, defined as a government-issued order at any
21 government level to address a directly COVID-19-related outcome (e.g., mask requirements, travel
22 restrictions, etc).
 - 23 ○ *Direct COVID-19 or SARS-CoV-2 outcomes are those that are specific to disease and health*
24 *outcomes* may include cases detected, mortality, number of tests taken, test positivity rates, Rt, etc.
 - 25 ○ This may NOT include indirect impacts of COVID-19 on items that are not direct *COVID-19 or*
26 *SARS-CoV-2 impacts* such as income, childcare, *economic impacts, beliefs and attitudes, etc.*
- 27 ● The primary outcome being examined must be a COVID-19-specific outcome, as above.
- 28 ● The study must be designed as an impact evaluation study from primary data (i.e., not primarily a predictive
29 or simulation model or meta-analysis).
- 30 ● The study must be peer reviewed, and published in a peer-reviewed journal indexed by PubMed.
- 31 ● The study must have the title and abstract available via PubMed at the time of the study start date
32 (November 26).
- 33 ● The study must be written in English.
34
35
36

37 These eligibility criteria were designed to identify the literature primarily concerning the
38 quantitative impact of one or more implemented COVID-19 policies on COVID-19 outcomes.
39 Studies in which impact evaluation was secondary to another analysis (such as a hypothetical
40 projection model) were eliminated because they were less relevant to our objectives and/or may
41 not contain sufficient information for evaluation. Categories for types of policies were from the
42 Oxford COVID-19 Government Response Tracker.[20]
43
44
45

46 Reviewer recruitment, training, and communication

47 Reviewers were recruited through personal contacts and postings on online media. All
48 reviewers had experience in systematic review, quantitative causal inference, epidemiology,
49 econometrics, public health, methods evaluation, or policy review. All reviewers participated in
50 two meetings in which the procedures and the review tool were demonstrated. Screening
51 reviewers participated in an additional meeting specific to the screening process. Throughout
52 the main review process, reviewers communicated with the administrators and each other
53 through Slack for any additional clarifications, questions, corrections, and procedures. The main
54
55
56
57
58
59
60

1
2
3 administrator (NH), who was also a reviewer, was available to answer general questions and
4 make clarifications, but did not answer questions specific to any given article.
5
6

7 Review phases and procedures

10 Search strategy

11 The search terms combined four Boolean-based search terms: a) COVID-19 research, 17 b)
12 regional government units (e.g., country, state, county, and specific country, state, or province,
13 etc.), c) policy or policies, and d) impact or effect. The full search terms are available in
14 Appendix 2.
15
16

18 Information Sources

19 The search was limited to published articles in peer-reviewed journals. This was largely to
20 attempt to identify literature that was high quality, relevant, prominent, and most applicable to
21 the review guidance. PubMed was chosen as the exclusive indexing source due to the
22 prevalence and prominence of policy impact studies in the health and medical field. Preprints
23 were excluded to limit the volume of studies to be screened and to ensure each had met the
24 standards for publication through peer review. The search was conducted on November 26,
25 2020.
26
27
28

30 Study Selection

31 Two reviewers were randomly selected to screen the title and abstract of each article for the
32 inclusion criteria. In the case of a dispute, a third randomly selected reviewer decided on
33 acceptance/rejection. Eight reviewers participated in the screening. Training consisted of a one-
34 hour instruction meeting, a review of the first 50 items on each reviewers' list of assigned
35 articles, and a brief asynchronous online discussion before conducting the full review.
36
37
38

39 Full article review

40 The full article review consisted of two sub-phases: the independent primary review phase, and
41 a group consensus phase. The independent review phase was designed primarily for the
42 purpose of supporting and facilitating discussion in the consensus discussion, rather than as
43 high stakes definitive review data on its own. The consensus process was considered the
44 primary way in which review data would be generated, rather than synthesis from the
45 independent reviews. A flow diagram of the review process is available in Appendix 3
46
47
48

49 Each article was randomly assigned to three of the 23 reviewers in our review pool. Each
50 reviewer independently reviewed each article on their list, first for whether the study met the
51 eligibility criteria, then responding to methods identification and guided strength of evidence
52 questions using the review tool, as described below. Reviewers were able to recuse themselves
53 for any reason, in which case another reviewer was randomly selected. Once all three reviewers
54
55
56
57
58
59
60

1
2
3 had reviewed a given article, all articles that weren't unanimously determined to not meet the
4 inclusion criteria underwent a consensus process.
5

6
7 During the consensus round, the three reviewers were given all three primary reviews for
8 reference, and were tasked with generating a consensus opinion among the group. One
9 randomly selected reviewer was tasked to act as the arbitrator. The arbitrator's primary task was
10 facilitating discussion and for moving the group toward establishing a consensus that
11 represented the collective subjective assessments of the group. If consensus could not be
12 reached, a fourth randomly selected reviewer was brought into the discussion to help resolve
13 disputes.
14
15

16 17 Review tool for data collection 18

19 This review tool and data collection process was an operationalized and lightly adapted version
20 of the COVID-19 health policy impact evaluation review guidance literature, written by the lead
21 authors of this study and released in September 2020 as a pre-print.[21] The main adaptation
22 was removing references to the COVID-19 literature. All reviewers were instructed to read and
23 refer to this guidance document to guide their assessments. The full guidance manuscript
24 contains additional explanation and rationale for all parts of this review and the tool, and is
25 available both in the adapted form as was provided to the reviewers in a supplementary file
26 "CHSPER review guidance refs removed.pdf" and in an updated version in Haber et al.,
27 2020.[17] The full review tool is attached as supplementary file "review tool final.pdf".
28
29
30

31 The review tool consisted of two main parts: methods design categorization and full review. The
32 review tool and guidance categorizes policy causal inference designs based on the structure of
33 their assumed counterfactual. This is assessed through identifying the data structure and
34 comparison(s) being made. There are two main items for this determination: the number of pre-
35 period time points (if any) used to assess pre-policy outcome trends, and whether or not policy
36 regions were compared with non-policy regions. These, and other supporting questions, broadly
37 allowed categorization of methods into cross-sectional, pre/post, interrupted time series (ITS),
38 difference-in-differences (DiD), comparative interrupted time-series (CITS), (randomized) trials,
39 or other. Given that most papers have several analyses, reviewers were asked to focus
40 exclusively on the impact evaluation analysis that was used as the primary support for the main
41 conclusion of the article.
42
43
44

45 Studies categorized as cross-sectional, pre/post, randomized controlled trial designs, and other
46 were included in our sample, but set aside for no further review for the purposes of this
47 research. Cross-sectional and pre/post studies are not considered sufficient to yield well-
48 identified causal inference in the specific context of COVID-19 policy impact evaluation, as
49 explained in the policy impact evaluation guidance documentation. Cross-sectional and pre-post
50 designs were considered inappropriate for policy causal inference for COVID-19 due largely to
51 inability to account for a large number of potential issues, including confounding, epidemic
52 trends, and selection biases. Randomized controlled trials were assumed to broadly meet key
53 design checks. Studies categorized as "other" received no further review, as the review
54
55
56
57
58
59
60

1
2
3 guidance would be unable to assess them. Additional justification and explanation for this
4 decision is available in the review guidance.
5

6
7 For the methods receiving full review (ITS, DiD, and CITS), reviewers were asked to identify
8 potential issues and give a category-specific rating. The specific study designs triggered sub-
9 questions and/or slightly altered the language of the questions being asked, but all three of the
10 methods design categories shared these four key questions:
11

- 12
13 ● Graphical presentation: “Does the analysis provide graphical representation of the
14 outcome over time?”
 - 15 ○ Graphical presentation refers to how the authors present the data underlying
16 their impact evaluation method. This is a critical criteria for assessing the
17 potential validity of the assumed model. The key questions here are whether any
18 chart shows the outcome over time and the assumed models of the
19 counterfactuals. To meet a high degree of confidence in this category, graphical
20 displays must show the outcome and connect to the counterfactual construction
21 method.
22
- 23
24 ● Functional form: “Is the functional form of the model used for the trend in counterfactual
25 infectious disease outcomes (e.g., linear, non-parametric, exponential, logarithmic, etc.)
26 well-justified and appropriate?”
 - 27 ○ Functional form refers to the statistical functional form of the trend in
28 counterfactual infectious disease outcomes (i.e. the assumptions used to
29 construct counterfactual outcomes). This may be a linear function, non-
30 parametric, exponential or logarithmic function, infectious disease model
31 projection, or any other functional form. The key criteria here are whether this is
32 discussed and justified in the manuscript, and if so, is it a plausibly appropriate
33 choice given infectious disease outcomes.
34
- 35
36 ● Timing of policy impact: “Is the date or time threshold set to the appropriate date or time
37 (e.g., is there lag between the intervention and outcome)?”
 - 38 ○ Timing of policy impact refers to assumptions about when we would expect to
39 see an impact from the policy vis-a-vis the timing of the policy introduction. This
40 would typically be modelled with leads and lags. The impact of policy can occur
41 before enactment (e.g., in cases where behavior change after policy is
42 announced, but before it takes place in anticipation) or long after the policy is
43 enacted (e.g., in cases where it takes time to ramp up policy implementation or
44 impacts). The key criteria here are whether this is discussed and justified in the
45 manuscript, and if so, whether it is a plausibly appropriate choice given the policy
46 and outcome.
47
- 48
49 ● Concurrent changes: “Is this policy the only uncontrolled or unadjusted-for way in which
50 the outcome could have changed during the measurement period [differently for policy
51 and non-policy regions]?”
 - 52 ○ Concurrent changes refers to the presence of uncontrolled other events and
53 changes that may influence outcomes at the same time as the policy would
54 impact outcomes. In order to assess the impact of one policy or set of policies,
55
56
57
58
59

1
2
3 the impact of all other forces that differentially impact the outcome must either be
4 negligible or controlled for. The key criteria here are whether it is likely that there
5 are substantial other uncontrolled forces (e.g. policies, behavioral changes, etc)
6 which may be differentially impacting outcomes at the same time as the policy of
7 interest.
8
9

10 For each of the four key questions, reviewers were given the option to select “No,” “Mostly no,”
11 “Mostly yes,” and “Yes” with justification text requested for all answers other than “Yes.” Each
12 question had additional prompts as guidance, and with much more detail provided in the full
13 guidance document. Ratings are, by design, subjective assessments of the category according
14 to the guidance. We do not use numerical scoring, for similar reasons as Cochrane suggests
15 that the algorithms for summary judgements for the RoB2 tool are merely “proposed”
16 assessments, which reviewers should change as they believe appropriate.[22] It is entirely
17 plausible, for example, for a study to meet all but one criteria but for the one remaining to be
18 sufficiently violated that the entire collective category is compromised. Alternatively, there could
19 be many minor violations of all of the criteria, but that they were collectively not sufficiently
20 problematic to impact overall ratings. Further, reviewers were also tasked with considering room
21 for doubt in cases where answers to these questions were unclear.
22
23
24
25

26 The criteria were designed to establish minimal plausibility of actionable evidence, rather than
27 certification of high quality. Graphical representation is included here primarily as a key way to
28 assess the plausibility and justification of key model assumptions, rather than being necessary
29 for validity by itself. For example, rather than having the “right” functional form or lag structure,
30 the review guidance asks whether the functional form and lags is discussed at all and (if
31 discussed) reasonable.
32
33
34

35 These four questions were selected and designed being critical to evaluating strength of study
36 design for policy impact evaluation in general, direct relevance for COVID-19 policy, feasibility
37 for use in guided review. These questions are designed as minimal and key criteria for plausibly
38 actionable impact evaluation design for COVID-19 policy impact evaluation, rather than as a
39 comprehensive tool assessing overall quality. Thorough review of data quality, statistical
40 validity, and other issues are also critical points of potential weakness in study designs, and
41 would be needed in addition to these criteria, if these key design criteria are met. A thorough
42 justification and explanation of how and why these questions were selected is available in the
43 provided guidance document and in Haber et al., 2020.[17]
44
45
46

47 Finally, reviewers were asked a summary question:

- 48
49 ● Overall: “Do you believe that the design is appropriate for identifying the policy impact(s)
50 of interest?”
51

52 Reviewers were asked to consider the scale of this question to be both independent/not relative
53 to any other papers, and that any one substantial issue with the study design could render it a
54 “No” or “Mostly no.” Reviewers were asked to follow the guidance and their previous answers,
55
56
57
58
59
60

1
2
3 allowing for their own weighting of how important each issue was to the final result. A study
4 could be excellent on all dimensions except for one, and that one dimension could render it
5 inappropriate for causal inference. As such, in addition to the overall rating question, we also
6 generated a “weakest link” metric for overall assessment, representing the lowest rating among
7 the four key questions (graphical representation, functional form, timing of policy impact, and
8 concurrent changes). A “mostly yes” or “yes” is considered a passing rating, indicating that the
9 study was not found to be inappropriate on the specific dimension of interest.
10
11

12
13 A “yes” rating does not necessarily indicate that the study is strongly designed, conducted, or is
14 actionable; it only means that it passes a series of key design checks for policy impact
15 evaluation and should be considered for further evaluation. The papers may contain any
16 number of other issues that were not reviewed (e.g., statistical issues, inappropriate
17 comparisons, generalizability, etc.). As such, this should only be considered an initial
18 assessment of plausibility that the study is well-designed, rather than confirmation that it is
19 appropriate and applicable.
20
21

22 The full review tool is available in the supplementary materials.
23
24

25 Heterogeneity

26
27 Inter-rater reliability (IRR) was assessed using Krippendorff’s alpha.[23,24] Rather than more
28 typical uses intended as an examination of the “validity” of ratings, the IRR statistic in this case
29 is being used as a heuristic indicator of heterogeneity between reviewers during the
30 independent phase, where heterogeneity is both expected and not necessarily undesirable. As
31 a second examination of reviewer heterogeneity, we also show the distribution of category
32 differences between primary reviewers within a study (e.g. if primary reviewers rated “Yes,”
33 “Mostly no,” and “Mostly yes” there are two pairs of answers that were one category different,
34 and one pair that was two categories different).
35
36
37

38 Statistical analysis

39
40 Statistics provided are nearly exclusively counts and percentages of the final dataset. Analyses
41 and graphics were performed in R.[25] Krippendorff’s alpha was calculated using the IRR
42 package.[26] Relative risks were estimated using the epitools package.[27]
43
44

45 Citation counts for accepted articles were obtained through Google Scholar[28] on January 11,
46 2021. Journal impact factors were obtained from the 2019 Journal Citation Reports.[29]
47
48

49 Data and code

50
51 Data, code, the review tool, and the review guidance are stored and available here:
52 <https://osf.io/9xmke/files/>. The dataset includes full results from the search and screening and all
53 review tool responses from reviewers during the full review phase.
54
55
56
57
58
59

Patient and Public Involvement Statement

Patients or public stakeholders were not consulted in the design or conduct of this systematic evaluation.

Results

Search and screening

Figure 1: PRISMA diagram of systematic review process

After search and screening of titles and abstracts, 102 articles were identified as likely or potentially meeting our inclusion criteria. Of those 102 articles, 36 studies met inclusion after independent review and deliberation in the consensus process. The most common reasons for rejection at this stage were that the study did not measure the quantitative direct impact of specific policies and/or that such an impact was not the main purpose of the study. Many of these studies implied that they measured policy impact in the abstract or introduction, but instead measured correlations with secondary outcomes (e.g., the effect of movement reductions, which are influenced by policy) and/or performed cursory policy impact evaluation secondary to projection modelling efforts.

Descriptive statistics

Figure 2: Descriptive sample statistics (n=36)

Publication information from our sample is shown in Figure 2. The articles in our sample were generally published in journals with high impact factors (median impact factor: 3.6, 25th percentile: 2.3, 75th percentile: 5.3 IQR: 3.0) and have already been cited in the academic literature (median citation count: 5.0, 25th percentile: 2.0, 75th percentile: 26.8, IQR 24.8, on 1/11/21). The most commonly evaluated policy type was stay at home requirements (64% n=23/36). Reviewers noted that many articles referenced “lockdowns,” but did not define the specific policies to which this referred. Reviewers specified mask mandates for 3 of the studies, and noted either a combination of many interventions or unspecified specific policies in 7 cases.

Reviewers most commonly selected interrupted time-series (39% n=14/36) as the methods design, followed by difference-in-differences (9% n=9/36) and pre-post (8% n=8/36). There were no randomized controlled trials of COVID-19 health policies identified (0% n=0/36), nor were any studies identified that reviewers could not categorize based on the review guidance (0% n=0/36).

Table 1: Summary of articles reviewed and reviewer ratings for key and overall questions

Category ratings order

Legend for color coded ratings



method determined to me inappropriate by: * guidance (cross sectional or pre/post) or ** reviewer consensus

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Citation	Title	Journal	Publication date	Methods design	Category ratings	Overall rating
Cobb and Seale, 2020[30]	Examining the effect of social distancing on the compound growth rate of COVID-19 at the county level (United States) using statistical analyses and a random forest machine learning model.	Public Health	4/28/2020	Pre/post		Red
Lyu and Wehby, 2020a[31]	Comparison of Estimated Rates of Coronavirus Disease 2019 (COVID-19) in Border Counties in Iowa Without a Stay-at-Home Order and Border Counties in Illinois With a Stay-at-Home Order.	JAMA Network Open	5/1/2020	Difference-in-differences	Light blue, Yellow	Yellow
Tam et al., 2020[32]	Effect of mitigation measures on the spreading of COVID-19 in hard-hit states in the U.S.	PloS One	5/1/2020	Interrupted time-series	Light blue, Orange	Orange
Courtemanche et al., 2020[33]	Strong Social Distancing Measures In The United States Reduced The COVID-19 Growth Rate.	Health Affairs	5/14/2020	Difference-in-differences	Light blue, Yellow	Yellow
Crokidakis, 2020[34]	COVID-19 spreading in Rio de Janeiro, Brazil: Do the policies of social isolation really work?	Chaos, Solitons, and Fractals	5/23/2020	Interrupted time-series	Light blue, Orange	Orange
Hyafil and Morifa, 2020[35]	Analysis of the impact of lockdown on the reproduction number of the SARS-Cov-2 in Spain.	Gaceta Aanitaria	5/23/2020	Pre/post		Red
Castillo, et al., 2020[36]	The effect of state-level stay-at-home orders on COVID-19 infection rates.	American Journal of Infection Control	5/24/2020	Pre/post		Red
Alfano and Ercolano, 2020[37]	The Efficacy of Lockdown Against COVID-19: A Cross-Country Panel Analysis.	Applied Health Economics and Health Policy	6/3/2020	Difference-in-differences	Orange, Light blue	Orange
Lyu and Wehby, 2020b[38]	Community Use Of Face Masks And COVID-19: Evidence From A Natural Experiment Of State Mandates In The US.	Health Affairs	6/16/2020	Difference-in-differences	Yellow, Light blue	Yellow
Zhang, et al., 2020[39]	Identifying airborne transmission as the dominant route for the spread of COVID-19.	PNAS	6/30/2020	Interrupted time-series	Light blue, Orange	Orange
Xu et al., 2020[40]	Associations of Stay-at-Home Order and Face-Masking Recommendation with Trends in Daily New Cases and Deaths of Laboratory-Confirmed COVID-19 in the United States.	Exploratory research and hypothesis in medicine	7/8/2020	Interrupted time-series	Light blue, Orange	Orange
Lyu and Wehby, 2020c[41]	Shelter-In-Place Orders Reduced COVID-19 Mortality And Reduced The Rate Of Growth In Hospitalizations.	Health Affairs	7/9/2020	Difference-in-differences	Light blue, Orange	Yellow
Wagner, et al., 2020[42]	Social distancing merely stabilized COVID-19 in the US.	Stat (International Statistical Institute)	7/13/2020	Interrupted time-series	Light blue, Yellow	Yellow
Di Bari et al., 2020[43]	Extensive Testing May Reduce COVID-19 Mortality: A Lesson From Northern Italy.	Frontiers in Medicine	7/14/2020	Comparative interrupted time-series	Yellow, Light blue	Orange
Islam et al., 2020[44]	Physical distancing interventions and incidence of coronavirus disease 2019: natural experiment in 149 countries.	BMJ (Clinical research ed.)	7/15/2020	Interrupted time-series	Yellow, Light blue	Orange
Wong et al., 2020[45]	Impact of National Containment Measures on Decelerating the Increase in Daily New Cases of COVID-19 in 54 Countries and 4 Epicenters of the Pandemic: Comparative Observational Study.	Journal of Medical Internet Research	7/22/2020	Pre/post		Red
Liang et al., 2020[46]	Effects of policies and containment measures on control of COVID-19 epidemic in Chongqing.	World Journal of Clinical Cases	7/26/2020	Pre/post		Red
Banerjee and Nayak, 2020[47]	U.S. county level analysis to determine If social distancing slowed the spread of COVID-19.	Pan American Journal of Public Health	7/31/2020	Difference-in-differences	Yellow, Orange	Orange
Dave et al., 2020a[48]	When Do Shelter-in-Place Orders Fight COVID-19 Best? Policy Heterogeneity Across States and Adoption Time.	Economic inquiry	8/3/2020	Difference-in-differences	Light blue, Orange	Light blue
Hsiang et al., 2020[49]	The effect of large-scale anti-contagion policies on the COVID-19 pandemic.	Nature	8/22/2020	Interrupted time-series	Light blue, Yellow	Yellow
Lim et al., 2020[50]	Revealing regional disparities in the transmission potential of SARS-CoV-2 from interventions in Southeast Asia.	Proceedings. Biological sciences	8/26/2020	Interrupted time-series	Light blue, Grey, Yellow	Orange
Arshed et al., 2020[51]	Empirical assessment of government policies and flattening of the COVID19 curve.	Journal of Public Affairs	8/27/2020	Cross-sectional analysis		Red
Wang et al., 2020[52]	Fangcang shelter hospitals are a One Health approach for responding to the COVID-19 outbreak in Wuhan, China.	One Health	8/29/2020	Interrupted time-series	Light blue, Orange	Orange
Kang et al., 2020[53]	The Effects of Border Shutdowns on the Spread of COVID-19.	Journal of Preventive Medicine and Public Health	8/30/2020	Comparative interrupted time-series	Light blue, Yellow	Yellow
Auger et al., 2020[54]	Association Between Statewide School Closure and COVID-19 Incidence and Mortality in the US.	JAMA	9/1/2020	Interrupted time-series	Light blue, Orange	Yellow
Santamaria et al., 2020[55]	COVID-19 effective reproduction number dropped during Spain's nationwide dropdown, then spiked at lower-incidence regions.	The Science of the Total Environment	9/9/2020	Interrupted time-series	Light blue, Orange	Light blue
Bennett, 2020[56]	All things equal? Heterogeneity in policy effectiveness against COVID-19 spread in Chile.	World Development	9/24/2020	Comparative interrupted time-series	Light blue, Orange	Light blue
Yang et al.,	Lessons Learnt from China: National Multidisciplinary Healthcare	Risk Management and	9/30/2020	Difference-in-	Orange	Orange

2020[57]	Assistance.	Healthcare Policy			differences		
Padalabalanar ayanan et al., 2020[58]	Association of State Stay-at-Home Orders and State-Level African American Population With COVID-19 Case Rates.	JAMA Network Open	10/1/2020		Comparative interrupted time-series		
Edelstein et al., 2020[59]	SARS-CoV-2 infection in London, England: changes to community point prevalence around lockdown time, March-May 2020.	Journal of Epidemiology and Community Health	10/1/2020		Pre/post		
Tsai et al., 2020[60]	COVID-19 transmission in the U.S. before vs. after relaxation of statewide social distancing measures.	Clinical Infectious Diseases	10/3/2020		Interrupted time-series		
Singh et al., 2020[61]	Public health interventions slowed but did not halt the spread of COVID-19 in India.	Transboundary and Emerging Diseases	10/4/2020		Pre/post		
Galloway et al., 2020[62]	Trends in COVID-19 Incidence After Implementation of Mitigation Measures - Arizona, January 22-August 7, 2020.	Morbidity and Mortality Weekly Report	10/9/2020		Pre/post		
Castex et al., 2020[63]	COVID-19: The impact of social distancing policies, cross-country analysis.	Economics of Disasters and Climate Change	10/15/2020		Interrupted time-series		
Silva et al., 2020[64]	The effect of lockdown on the COVID-19 epidemic in Brazil: evidence from an interrupted time series design.	Cadernos de Saude Publica	10/19/2020		Interrupted time-series		
Dave et al., 2020b[65]	Were Urban Cowboys Enough to Control COVID-19? Local Shelter-in-Place Orders and Coronavirus Case Growth.	Journal of Urban Economics	11/6/2020		Difference-in-differences		

The identified articles and selected review results are summarized in Table 1.

Strength of methods assessment

Figure 3: Main consensus results summary for key and overall questions

Graphical representation of the outcome over time was relatively well-rated in our sample, with 74% (n=20/27) studies being given a “mostly yes” or “yes” rating for appropriateness. Reasons cited for non-“yes” ratings included a lack of graphical representation of the data, alternative scales used, and not showing the dates of policy implementation.

Functional form issues appear to have presented a major issue in these studies, with only 19% receiving a “mostly yes” or “yes” rating, 78% (n=21/27) receiving a “no” rating, and 4% (n=1/27) “unclear.” There were two common themes in this category: studies generally using scales that were broadly considered inappropriate for infectious disease outcomes (e.g., linear counts), and/or studies lacking stated justification for the scale used. Reviewers also noted disconnects between clear curvature in the outcomes in the graphical representations and the analysis models and outcome scales used (e.g., linear). In one case, reviewers could not identify the functional form actually used in analysis.

Reviewers broadly found that these studies dealt with timing of policy impact (e.g., lags between policy implementation and expected impact) relatively well, with 70% (n=19/27) rated “yes” or “mostly yes.” Reasons for non-“yes” responses included not adjusting for lags and a lack of justification for the specific lags used.

Concurrent changes were found to be a major issue in these studies, with only 11% (n=3/27) studies receiving passing ratings (“yes” or “mostly yes”) with regard to uncontrolled concurrent changes to the outcomes. Reviewers nearly ubiquitously noted that the articles failed to account for the impact of other policies that could have impacted COVID-19 outcomes concurrent with

1
2
3 the policies of interest. Other issues cited were largely related to non-policy-induced behavioral
4 and societal changes.
5

6
7 When reviewers were asked if sensitivity analyses had been performed on key assumptions and
8 parameters, about half (56% n=15/27) answered “mostly yes” or “yes.” The most common
9 reason for non-“yes” ratings was that, while sensitivity analyses were performed, they did not
10 address the most substantial assumptions and issues.
11

12
13 Overall, reviewers rated only four studies (11%, n=4/36,) as being plausibly appropriate (“mostly
14 yes” or “yes”) for identifying the impact of specific policies on COVID-19 outcomes, as shown in
15 Figure 3. 25% (n=9/36) were automatically categorized as being inappropriate due to being
16 either cross-sectional or pre/post in design, 33% (n=12/36) of studies were given a “no” rating
17 for appropriateness, 31% “mostly no” (n=11/36), 8% “mostly yes” (n=3/36), and 3% “yes”
18 (n=1/36). The most common reason cited for non-“yes” overall ratings was failure to account for
19 concurrent changes (particularly policy and societal changes).
20
21

22
23 Figure 4: Comparison of independent reviews, weakest link, and direct consensus review
24

25
26 As shown in Figure 4, the consensus overall proportion passing (“mostly yes” or “yes”) was a
27 quarter of what it was from the initial independent reviews. 45% (n=34/75) of studies were rated
28 as “yes” or “mostly yes” in the initial independent review, as compared to 11% (n=4/36) in the
29 consensus round (RR 0.25, 95%CI 0.09:0.64). The issues identified and discussed in
30 combination during consensus discussions, as well as additional clarity on the review process,
31 resulted in reduced overall confidence in the findings. Increased clarity on the review guidance
32 with experience and time may also have reduced these ratings further.
33

34
35 The large majority of studies had at least one “no” or “unclear” rating in one of the four
36 categories (74% n=20/27), with only one study whose lowest rating was a “mostly yes,” no
37 studies rated “yes” in all four categories. Only one study was found to pass design criteria in all
38 four key questions categories, as shown in the “weakest link” column in Figure 4.
39
40

41 Review process assessment

42
43 During independent review, all three reviewers independently came to the same conclusions on
44 the main methods design category for 33% (n=12/36) articles, two out of the three reviewers
45 agreed for 44% (n=16/36) articles, and none of the reviewers agreed in 22% (n=8/36) cases.
46 One major contributor to these discrepancies were the 31% (n=11/36) cases where one or more
47 reviewers marked the study as not meeting eligibility criteria, 64% (n=7/11) of which the other
48 two reviewers agreed on the methods design category.
49
50

51
52 Reviewers’ initial independent reviews were heterogeneous for key rating questions. For the
53 overall scores, Krippendorff’s alpha was only 0.16 due to widely varying opinions between
54 raters. The four key categorical questions had slightly better inter-rater reliability than the overall
55 question, with Krippendorff’s alphas of 0.59 for graphical representation, 0.34 for functional form,
56
57
58
59

0.44 for timing of policy impact, and 0.15 for concurrent changes, respectively. For the main summary rating, primary reviewers within each study agreed in 26% of cases (n=16), were one category different in 45% (n=46), two categories different in 19% (n=12), and three categories (i.e. the maximum distance, “Yes” vs “No”) in 10% of cases (n=6).

The consensus rating for overall strength was equal to the lowest rating among the independent reviews in 78% (n=21/27) of cases, and only one higher than the lowest in the remaining 22% (n=6/27). This strongly suggests that the multiple reviewer review, discussion, and consensus process more thoroughly identifies issues than independent review alone. There were two cases for which reviewers requested an additional fourth reviewer to help resolve standing issues for which the reviewers felt they were unable to come to consensus.

The most consistent point of feedback from reviewers was the value of having a three reviewer team with whom to discuss and deliberate, rather than two as initially planned. This was reported to help catch a larger number of issues and clarify both the papers and the interpretation of the review tool questions. Reviewers also expressed that one of the most difficult parts of this process was assessing the inclusion criteria, some of the implications of which are discussed below.

Discussion

This systematic review of evidence strength found that only four (or only one by a stricter standard) of the 36 identified published and peer-reviewed health policy impact evaluation studies passed a set of key checks for identifying the causal impact of policies on COVID-19 outcomes. Because this systematic review examined a limited set of key study design features and did not address more detailed aspects of study design, statistical issues, generalizability, and any number of other issues, this result may be considered an upper bound on the overall strength of evidence within this sample. Two major problems are nearly ubiquitous throughout this literature: failure to isolate the impact of the policy(s) of interest from other changes that were occurring contemporaneously, and failure to appropriately address the functional form of infectious disease outcomes in a population setting. While policy decisions are being made on the backs of high impact-factor papers, we find that the citation-based metrics do not correspond to “quality” research as used by Yin et al., 2021.[66] Similar to other areas in the COVID-19 literature,[67] we found the current literature directly evaluating the impact of COVID-19 policies largely fails to meet key design criteria for actionable inference to inform policy decisions.

The framework for the review tool is based on the requirements and assumptions built into policy evaluation methods. Quasi-experimental methods rely critically on the scenarios in which the data are generated. These assumptions and the circumstances in which they are plausible are well-documented and understood,[2,4–6,17,68] including one paper discussing application of difference-in-differences methods specifically for COVID-19 health policy, released in May 2020.[5] While “no uncontrolled concurrent changes” is a difficult bar to clear, that bar is fundamental to inference using these methods.

1
2
3
4 The circumstances of isolating the impact of policies in COVID-19 - including large numbers of
5 policies, infectious disease dynamics, and massive changes to social behaviors - make those
6 already difficult fundamental assumptions broadly much less likely to be met. Some of the
7 studies in our sample were nearly the best feasible studies that could be done given the
8 circumstances, but the best that can be done often yields little actionable inference. The relative
9 paucity of strong studies does not in any way imply a lack of impact of those policies; only that
10 we lack the circumstances to have evaluated their effects.
11
12

13
14 Because the studies estimating the harms of policies share the same fundamental
15 circumstances, the evidence of COVID-19 policy harms is likely to be of similarly poor strength.
16 Identifying the effects of many of these policies, particularly for the spring of 2020, is likely to be
17 unknown and perhaps unknowable. However, there remains additional opportunities with more
18 favorable circumstances, such as measuring overall impact of NPIs as bundles, rather than
19 individual policies. Similarly, studies estimating the impact of re-opening policies or policy
20 cancellation are likely to have fewer concurrent changes to address.
21
22

23
24 The review process itself demonstrates how guided and targeted peer review can efficiently
25 evaluate studies in ways that the traditional peer review systems do not. The studies in our
26 sample had passed the full peer review process, were published in largely high-profile journals,
27 and are highly cited, but contained substantial flaws that rendered their inference utility
28 questionable. The relatively small number of studies included, as compared to the size of the
29 literature concerning itself with COVID-19 policy, may suggest that there was relative restraint
30 from journal editors and reviewers for publishing these types of studies. The large number of
31 models, but relatively small number of primary evaluation analyses is consistent with other
32 areas of COVID-19.[69,70] At minimum, the flaws and limitations in their inference could have
33 been communicated at the time of publication, when they are needed most. In other cases, it is
34 plausible that many of these studies would not have been published had a more thorough or
35 more targeted methodological review been performed.
36
37
38

39
40 This systematic review of evidence strength has limitations. The tool itself was limited to a very
41 narrow - albeit critical - set of items. Low ratings in our study should not be interpreted as being
42 overall poor studies, as they may make other contributions to the literature that we did not
43 evaluate. While the guidance and tool provided a well-structured framework and our reviewer
44 pool was well-qualified, strength of evidence review is inherently subjective. It is plausible and
45 likely that other sets of reviewers would come to different conclusions for each study, but
46 unlikely that the overall conclusions of our assessment would change substantially. However,
47 the consensus process was designed with these issues subjectivity in mind, and demonstrates
48 the value of consensus processes for overcoming hurdles with subjective and difficult decisions.
49
50

51
52 While subjective assessments are inherently subject to the technical expertise, experiences,
53 and opinions of reviewers, we argue they are both appropriate and necessary to reliably assess
54 strength of evidence based on theoretical methodological issues. With the exception of the
55 graphical assessment, proper assessment of the core methodological issues requires that
56
57
58
59

1
2
3 reviewers are able to weigh the evidence as they see fit. Much like standard institutional peer
4 review, reviewers independently had highly heterogeneous opinions, attributable to differences
5 in opinion or training, misunderstandings/learning about the review tool and process, and
6 expected reliance on the consensus process. Unlike traditional peer review, there was subject-
7 matter-specific guidance and a process to consolidate and discuss those heterogenous initial
8 opinions. The reduction in ratings from the initial highly heterogeneous ratings to a lower
9 heterogeneity in ratings indicates that reviewers had initially identified issues differently, but that
10 the discussion and consensus process helped elucidate the extent of the different issues that
11 each reviewer detected and brought to discussion. This also reflects reviewer learning over
12 time, where reviewers were better able to identify issues at the consensus phase than earlier. It
13 is plausible that stronger opinions had more weight, but we expect that this was largely
14 mitigated by the random assignment of the arbitrator, and reviewer experiences did not indicate
15 this as an issue.
16
17
18
19

20 Most importantly, this review does not cover all policy inference in the scientific literature. One
21 large literature from which there may be COVID-19 policy evaluation otherwise meeting our
22 inclusion criteria are pre-prints. Many pre-prints would likely fare well in our review process.
23 Higher strength papers often require more time for review and publication, and many high
24 quality papers may be in the publication pipeline now. Second, this review excluded studies that
25 had a quantitative impact evaluation as a secondary part of the study (e.g., to estimate
26 parameters for microsimulation or disease modeling). Third, the review does not include policy
27 inference studies that do not measure the impact of a specific policy. For instance, there are
28 studies that estimate the impact of reduced mobility on COVID-19 outcomes but do not attribute
29 the reduced mobility to any specific policy change. Finally, a considerable number of studies
30 that present analyses of COVID-19 outcomes to inform policy are excluded because they do not
31 present a quantitative estimate of specific policies' treatment effects.
32
33
34
35

36 While COVID-19 policy is one of the most important problems of our time, the circumstances
37 under which those policies were enacted severely hamper our ability to study and understand
38 their effects. Claimed conclusions are only as valuable as the methods by which they are
39 produced. Replicable, rigorous, intense, and methodologically guided review is needed to both
40 communicate our limitations and make more actionable inference. Weak, unreliable, and
41 overconfident evidence leads to poor decisions and undermines trust in science.[15,71] In the
42 case of COVID-19 health policy, a frank appraisal of the strength of the studies on which
43 policies are based is needed, alongside the understanding that we often must make decisions
44 when strong evidence is not feasible.[72]
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Ethical approval statement

Not applicable

Works cited

- 1 Fischhoff B. Making Decisions in a COVID-19 World. *JAMA* 2020;324:139. doi:10.1001/jama.2020.10178
- 2 COVID-19 Statistics, Policy modeling, and Epidemiology Collective. Defining high-value information for COVID-19 decision-making. *Health Policy* 2020. doi:10.1101/2020.04.06.20052506
- 3 Hernán MA, Robins JM. *Causal Inference: What If*. Boca Raton: : Chapman & Hall/CRC
- 4 Angrist J, Pischke J-S. *Mostly Harmless Econometrics: An Empiricist's Companion*. 1st ed. Princeton University Press 2009. <https://EconPapers.repec.org/RePEc:pup:pbooks:8769>
- 5 Goodman-Bacon A, Marcus J. Using Difference-in-Differences to Identify Causal Effects of COVID-19 Policies. SSRN Journal Published Online First: 2020. doi:10.2139/ssrn.3603970
- 6 Bärnighausen T, Oldenburg C, Tugwell P, et al. Quasi-experimental study designs series—paper 7: assessing the assumptions. *Journal of Clinical Epidemiology* 2017;89:53–66. doi:10.1016/j.jclinepi.2017.02.017
- 7 Haushofer J, Metcalf CJE. Which interventions work best in a pandemic? *Science* 2020;368:1063–5. doi:10.1126/science.abb6144
- 8 Else H. How a torrent of COVID science changed research publishing — in seven charts. *Nature* 2020;588:553–553. doi:10.1038/d41586-020-03564-y
- 9 Palayew A, Norgaard O, Safreed-Harmon K, et al. Pandemic publishing poses a new COVID-19 challenge. *Nat Hum Behav* 2020;4:666–9. doi:10.1038/s41562-020-0911-0
- 10 Bagdasarian N, Cross GB, Fisher D. Rapid publications risk the integrity of science in the era of COVID-19. *BMC Med* 2020;18:192. doi:10.1186/s12916-020-01650-6
- 11 Yeo-Teh NSL, Tang BL. An alarming retraction rate for scientific publications on Coronavirus Disease 2019 (COVID-19). *Accountability in Research* 2020;0:1–7. doi:10.1080/08989621.2020.1782203
- 12 Abritis A, Marcus A, Oransky I. An “alarming” and “exceptionally high” rate of COVID-19 retractions? *Accountability in Research* 2021;28:58–9. doi:10.1080/08989621.2020.1793675
- 13 Zdravkovic M, Berger-Estilita J, Zdravkovic B, et al. Scientific quality of COVID-19 and SARS CoV-2 publications in the highest impact medical journals during the early phase of the pandemic: A case control study. *PLoS ONE* 2020;15:e0241826. doi:10.1371/journal.pone.0241826
- 14 Elgendy IY, Nimri N, Barakat AF, et al. A systematic bias assessment of top-cited full-length original clinical investigations related to COVID-19. *European Journal of Internal Medicine* 2021;:S0953620521000182. doi:10.1016/j.ejim.2021.01.018
- 15 Glasziou PP, Sanders S, Hoffmann T. Waste in covid-19 research. *BMJ* 2020;369:m1847. doi:10.1136/bmj.m1847
- 16 Powell M, Koenecke A, Byrd JB, et al. A how-to guide for conducting retrospective analyses: example COVID-19 study. *Open Science Framework* 2020. doi:10.31219/osf.io/3drch
- 17 Haber NA, Clarke-Deelder E, Salomon JA, et al. COVID-19 Policy Impact Evaluation: A guide to common design issues. *American Journal of Epidemiology* 2021;:kwab185. doi:10.1093/aje/kwab185
- 18 Haber N. Systematic review of COVID-19 policy evaluation methods and design. Published

- 1
2
3 Online First: 26 November 2020. <https://osf.io/7nbk6> (accessed 15 Jan 2021).
- 4 19 PRISMA. <http://www.prisma-statement.org/PRISMAStatement/> (accessed 15 Jan 2021).
- 5 20 Petherick A, Kira B, Hale T, et al. Variation in Government Responses to COVID-19.
6 <https://www.bsg.ox.ac.uk/research/publications/variation-government-responses-covid-19>
7 (accessed 24 Nov 2020).
- 8 21 Haber NA, Clarke-Deelder E, Salomon JA, et al. Policy evaluation in COVID-19: A guide to
9 common design issues. arXiv:200901940 [stat] Published Online First: 31 December
10 2020. <http://arxiv.org/abs/2009.01940> (accessed 15 Jan 2021).
- 11 22 Chapter 8: Assessing risk of bias in a randomized trial.
12 <https://training.cochrane.org/handbook/current/chapter-08> (accessed 8 Sep 2021).
- 13 23 Krippendorff KH. Content Analysis: An Introduction to Its Methodology. SAGE Publications
14 1980.
- 15 24 Zhao X, Liu JS, Deng K. Assumptions behind Intercoder Reliability Indices. *Annals of the*
16 *International Communication Association* 2013;36:419–80.
17 doi:10.1080/23808985.2013.11679142
- 18 25 R Core Team. R: A language and environment for statistical computing. Vienna, Austria: R
19 Foundation for Statistical Computing 2019. <https://www.R-project.org/>
- 20 26 Gamer M, Lemon J, Fellows I, et al. irr: Various Coefficients of Interrater Reliability and
21 Agreement. <https://cran.r-project.org/web/packages/irr/index.html>
- 22 27 Aragon TJ, Fay MP, Wollschlaeger D, et al. Epitools. CRAN: 2017. [https://cran.r-](https://cran.r-project.org/web/packages/epitools/epitools.pdf)
23 [project.org/web/packages/epitools/epitools.pdf](https://cran.r-project.org/web/packages/epitools/epitools.pdf)
- 24 28 About Google Scholar. <https://scholar.google.com/intl/en/scholar/about.html> (accessed 15
25 Jan 2021).
- 26 29 Clarivate Analytics. Journal Citation Reports. 2019.
- 27 30 Cobb JS, Seale MA. Examining the effect of social distancing on the compound growth rate
28 of COVID-19 at the county level (United States) using statistical analyses and a random
29 forest machine learning model. *Public Health* 2020;185:27–9.
30 doi:10.1016/j.puhe.2020.04.016
- 31 31 Lyu W, Wehby GL. Comparison of Estimated Rates of Coronavirus Disease 2019 (COVID-
32 19) in Border Counties in Iowa Without a Stay-at-Home Order and Border Counties in Illinois
33 With a Stay-at-Home Order. *JAMA Netw Open* 2020;3:e2011102.
34 doi:10.1001/jamanetworkopen.2020.11102
- 35 32 Tam K-M, Walker N, Moreno J. Effect of mitigation measures on the spreading of COVID-19
36 in hard-hit states in the U.S. *PLoS One* 2020;15:e0240877.
37 doi:10.1371/journal.pone.0240877
- 38 33 Courtemanche C, Garuccio J, Le A, et al. Strong Social Distancing Measures In The United
39 States Reduced The COVID-19 Growth Rate: Study evaluates the impact of social
40 distancing measures on the growth rate of confirmed COVID-19 cases across the United
41 States. *Health Affairs* 2020;39:1237–46. doi:10.1377/hlthaff.2020.00608
- 42 34 Crokidakis N. COVID-19 spreading in Rio de Janeiro, Brazil: Do the policies of social
43 isolation really work? *Chaos Solitons Fractals* 2020;136:109930.
44 doi:10.1016/j.chaos.2020.109930
- 45 35 Hyafil A, Morfiña D. Analysis of the impact of lockdown on the reproduction number of the
46 SARS-Cov-2 in Spain. *Gaceta Sanitaria* 2020;;S0213911120300984.
47 doi:10.1016/j.gaceta.2020.05.003
- 48 36 Castillo RC, Staguñ ED, Weston-Farber E. The effect of state-level stay-at-home orders on
49 COVID-19 infection rates. *American Journal of Infection Control* 2020;48:958–60.
50 doi:10.1016/j.ajic.2020.05.017
- 51 37 Alfano V, Ercolano S. The Efficacy of Lockdown Against COVID-19: A Cross-Country Panel
52 Analysis. *Appl Health Econ Health Policy* 2020;18:509–17. doi:10.1007/s40258-020-00596-
53 3
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- 38 Lyu W, Wehby GL. Community Use Of Face Masks And COVID-19: Evidence From A Natural Experiment Of State Mandates In The US: Study examines impact on COVID-19 growth rates associated with state government mandates requiring face mask use in public. *Health Affairs* 2020;39:1419–25. doi:10.1377/hlthaff.2020.00818
- 39 Zhang R, Li Y, Zhang AL, et al. Identifying airborne transmission as the dominant route for the spread of COVID-19. *Proc Natl Acad Sci USA* 2020;117:14857–63. doi:10.1073/pnas.2009637117
- 40 Xu J, Hussain S, Lu G, et al. Associations of Stay-at-Home Order and Face-Masking Recommendation with Trends in Daily New Cases and Deaths of Laboratory-Confirmed COVID-19 in the United States. *Explor Res Hypothesis Med* 2020;:1–10. doi:10.14218/ERHM.2020.00045
- 41 Lyu W, Wehby GL. Shelter-In-Place Orders Reduced COVID-19 Mortality And Reduced The Rate Of Growth In Hospitalizations. *Health Aff (Millwood)* 2020;39:1615–23. doi:10.1377/hlthaff.2020.00719
- 42 Wagner AB, Hill EL, Ryan SE, et al. Social distancing merely stabilized COVID-19 in the United States. *Stat* 2020;9. doi:10.1002/sta4.302
- 43 Di Bari M, Balzi D, Carreras G, et al. Extensive Testing May Reduce COVID-19 Mortality: A Lesson From Northern Italy. *Front Med* 2020;7:402. doi:10.3389/fmed.2020.00402
- 44 Islam N, Sharp SJ, Chowell G, et al. Physical distancing interventions and incidence of coronavirus disease 2019: natural experiment in 149 countries. *BMJ* 2020;:m2743. doi:10.1136/bmj.m2743
- 45 Wong LP, Alias H. Temporal changes in psychobehavioural responses during the early phase of the COVID-19 pandemic in Malaysia. *J Behav Med* 2020;:1–11. doi:10.1007/s10865-020-00172-z
- 46 Liang X-H, Tang X, Luo Y-T, et al. Effects of policies and containment measures on control of COVID-19 epidemic in Chongqing. *WJCC* 2020;8:2959–76. doi:10.12998/wjcc.v8.i14.2959
- 47 Banerjee T, Nayak A. U.S. county level analysis to determine If social distancing slowed the spread of COVID-19. *Revista Panamericana de Salud Pública* 2020;44:1. doi:10.26633/RPSP.2020.90
- 48 Dave D, Friedson AI, Matsuzawa K, et al. When Do Shelter-in-Place Orders Fight COVID-19 Best? Policy Heterogeneity Across States and Adoption Time. *Econ Inq* Published Online First: 3 August 2020. doi:10.1111/ecin.12944
- 49 Hsiang S, Allen D, Annan-Phan S, et al. The effect of large-scale anti-contagion policies on the COVID-19 pandemic. *Nature* 2020;584:262–7. doi:10.1038/s41586-020-2404-8
- 50 Lim JT, Dickens BSL, Choo ELW, et al. Revealing regional disparities in the transmission potential of SARS-CoV-2 from interventions in Southeast Asia. *Proc R Soc B* 2020;287:20201173. doi:10.1098/rspb.2020.1173
- 51 Arshed N, Meo MS, Farooq F. Empirical assessment of government policies and flattening of the COVID 19 curve. *J Public Affairs* Published Online First: 27 August 2020. doi:10.1002/pa.2333
- 52 Wang K-W, Gao J, Song X-X, et al. Fangcang shelter hospitals are a One Health approach for responding to the COVID-19 outbreak in Wuhan, China. *One Health* 2020;10:100167. doi:10.1016/j.onehlt.2020.100167
- 53 Kang N, Kim B. The Effects of Border Shutdowns on the Spread of COVID-19. *J Prev Med Public Health* 2020;53:293–301. doi:10.3961/jpmph.20.332
- 54 Auger KA, Shah SS, Richardson T, et al. Association Between Statewide School Closure and COVID-19 Incidence and Mortality in the US. *JAMA* 2020;324:859. doi:10.1001/jama.2020.14348
- 55 Santamaría L, Hortal J. COVID-19 effective reproduction number dropped during Spain's nationwide dropdown, then spiked at lower-incidence regions. *Science of The Total*

- Environment 2021;751:142257. doi:10.1016/j.scitotenv.2020.142257
- 56 Bennett M. All things equal? Heterogeneity in policy effectiveness against COVID-19 spread in Chile. *World Development* 2021;137:105208. doi:10.1016/j.worlddev.2020.105208
- 57 Yang T, Shi H, Liu J, et al. Lessons Learnt from China: National Multidisciplinary Healthcare Assistance. *RMHP* 2020;Volume 13:1835–7. doi:10.2147/RMHP.S269523
- 58 Padalabalanarayanan S, Hanumanthu VS, Sen B. Association of State Stay-at-Home Orders and State-Level African American Population With COVID-19 Case Rates. *JAMA Netw Open* 2020;3:e2026010. doi:10.1001/jamanetworkopen.2020.26010
- 59 Edelstein M, Obi C, Chand M, et al. SARS-CoV-2 infection in London, England: changes to community point prevalence around lockdown time, March–May 2020. *J Epidemiol Community Health* 2020;:jech-2020-214730. doi:10.1136/jech-2020-214730
- 60 Tsai AC, Harling G, Reynolds Z, et al. COVID-19 transmission in the U.S. before vs. after relaxation of statewide social distancing measures. *Clin Infect Dis Published Online First*: 3 October 2020. doi:10.1093/cid/ciaa1502
- 61 Singh BB, Lowerison M, Lewinson RT, et al. Public health interventions slowed but did not halt the spread of COVID-19 in India. *Transbound Emerg Dis Published Online First*: 4 October 2020. doi:10.1111/tbed.13868
- 62 Gallaway MS, Rigler J, Robinson S, et al. Trends in COVID-19 Incidence After Implementation of Mitigation Measures — Arizona, January 22–August 7, 2020. *MMWR Morb Mortal Wkly Rep* 2020;69:1460–3. doi:10.15585/mmwr.mm6940e3
- 63 Castex G, Dechter E, Lorca M. COVID-19: The impact of social distancing policies, cross-country analysis. *EconDisCliCha Published Online First*: 15 October 2020. doi:10.1007/s41885-020-00076-x
- 64 Silva L, Figueiredo Filho D, Fernandes A. The effect of lockdown on the COVID-19 epidemic in Brazil: evidence from an interrupted time series design. *Cad Saúde Pública* 2020;36:e00213920. doi:10.1590/0102-311x00213920
- 65 Dave D, Friedson A, Matsuzawa K, et al. Were Urban Cowboys Enough to Control COVID-19? Local Shelter-in-Place Orders and Coronavirus Case Growth. *J Urban Econ* 2020;:103294. doi:10.1016/j.jue.2020.103294
- 66 Yin Y, Gao J, Jones BF, et al. Coevolution of policy and science during the pandemic. *Science* 2021;371:128–30. doi:10.1126/science.abe3084
- 67 Wynants L, Van Calster B, Collins GS, et al. Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. *BMJ* 2020;:m1328. doi:10.1136/bmj.m1328
- 68 Clarke GM, Conti S, Wolters AT, et al. Evaluating the impact of healthcare interventions using routine data. *BMJ* 2019;:l2239. doi:10.1136/bmj.l2239
- 69 Krishnaratne S, Pfadenhauer LM, Coenen M, et al. Measures implemented in the school setting to contain the COVID-19 pandemic: a rapid scoping review. *Cochrane Database of Systematic Reviews Published Online First*: 17 December 2020. doi:10.1002/14651858.CD013812
- 70 Raynaud M, Zhang H, Louis K, et al. COVID-19-related medical research: a meta-research and critical appraisal. *BMC Med Res Methodol* 2021;21:1. doi:10.1186/s12874-020-01190-w
- 71 Casigliani V, De Nard F, De Vita E, et al. Too much information, too little evidence: is waste in research fuelling the covid-19 infodemic? *BMJ* 2020;:m2672. doi:10.1136/bmj.m2672
- 72 Greenhalgh T. Will COVID-19 be evidence-based medicine's nemesis? *PLoS Med* 2020;17:e1003266. doi:10.1371/journal.pmed.1003266

Acknowledgements

We would like to thank Keletso Makofane for assisting with the screening, Dr. Steven Goodman and Dr. John Ioannidis for their support during the development of this study, and Dr. Lars Hemkins and Dr. Mario Malicki for helpful comments in the protocol development.

Transparency declaration

The lead author (NH) affirms that this manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been explained.

Author roles

NH led the protocol development, study design, administration, data curation, data management, statistical analysis, graphical design, manuscript writing, and manuscript editing, and serves as the primary guarantor of the study.

NH, ECD, JS, AF, and ESt co-wrote the review guidance on which the design of the study review tool is based

NH, ECD, JS, AF, ESm and ESt designed, wrote, and supported the pre-registered protocol.

NH, CJ, SW, CB, CA, CBF, VN, and Keletso Makofane were the screening reviewers for this study, analysing the abstracts and titles for inclusion criteria.

NH, ECD, AF, BMG, ES, CBF, JD, LH, CG, CB, EBM, CJ, BL, IS, EA, SW, BJ, CA, VN, BG, AB, ES were the main reviewers for this study, and contributed to the analysis and evaluation of the studies entering into the main review phase.

NH, ECD, JA, AF, BMG, ESm, CBF, JD, LH, CG, CB, EBM, CJ, BL, IS, EA, SW, BJ, CA, VN, BG, AB, and ESt all contributed to and supported the manuscript editing.

Funding

No funding was provided specifically for this research.

Elizabeth Stone receives funding under the National Institutes of Health grant T32MH109436.

Ian Schmid receives funding under the National Institutes of Health grant T32MH122357.

Brooke Jarrett receives funding under the National Institutes of Health grant MH121128.

Christopher Boyer receives funding under the National Institutes of Health grant T32HL098048

Cathrine Axfors receives funding from the Knut and Alice Wallenberg Foundation, grant KAW 2019.0561.

Beth Ann Griffin and Elizabeth Stuart were supported by award number P50DA046351 from the National Institute on Drug Abuse. Elizabeth Stuart's time was also supported by the Bloomberg

American Health Initiative. Caroline Joyce receives funding from the Ferring Foundation. Meta-Research Innovation Center at Stanford (METRICS), Stanford University is supported by a grant from the Laura and John Arnold Foundation

Conflicts of interest disclosure

The authors have no financial or social conflicts of interest to declare.

Figures

Figure 1: PRISMA diagram of systematic review process

Caption: This chart shows the PRISMA diagram for the process of screening the literature from search to the full review phase.

Figure 2: Descriptive sample statistics (n=36)

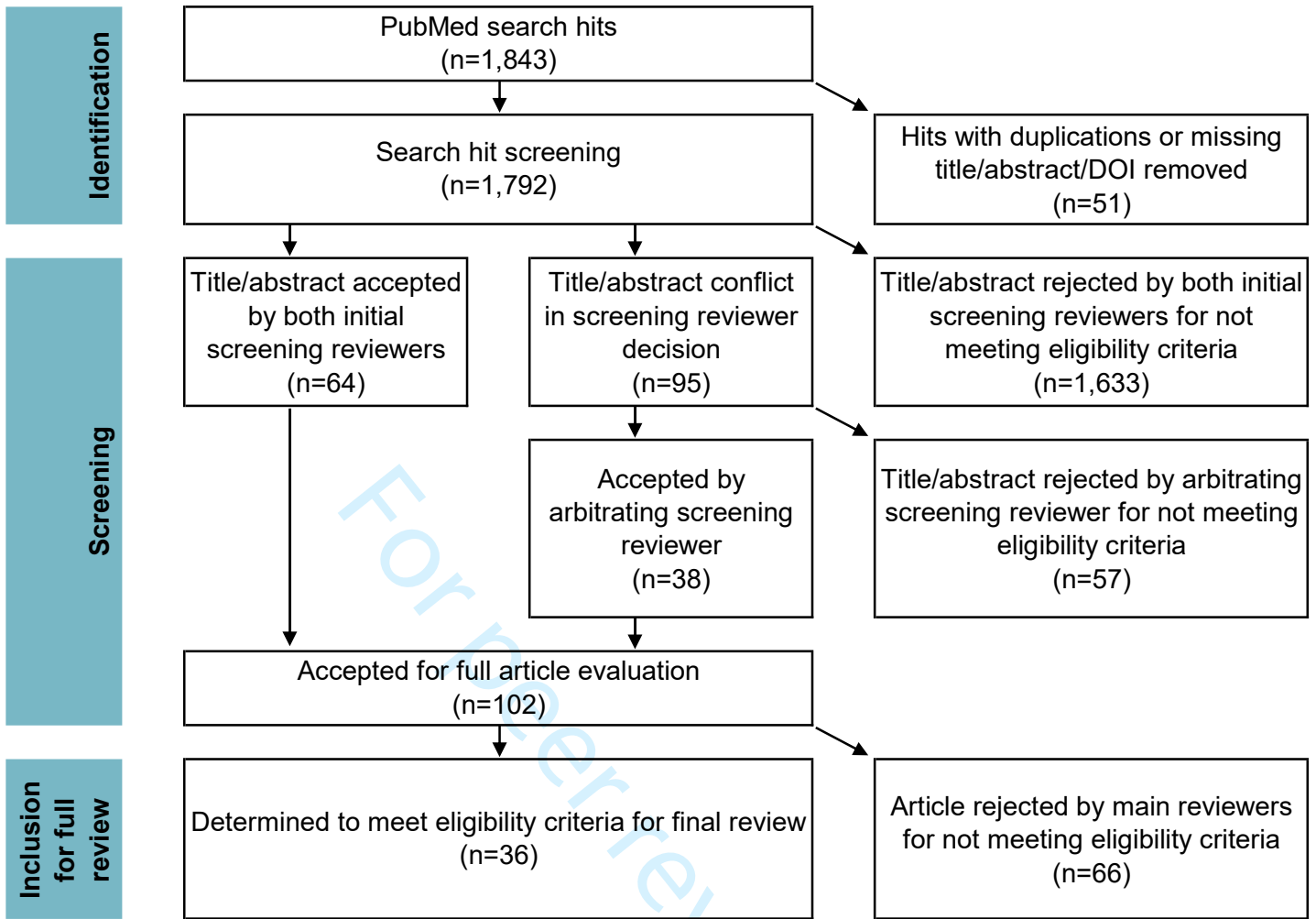
Caption: This chart shows descriptive statistics of the 36 studies entered into our systematic evidence review.

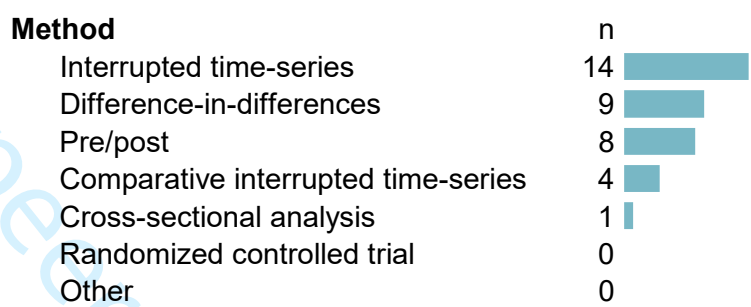
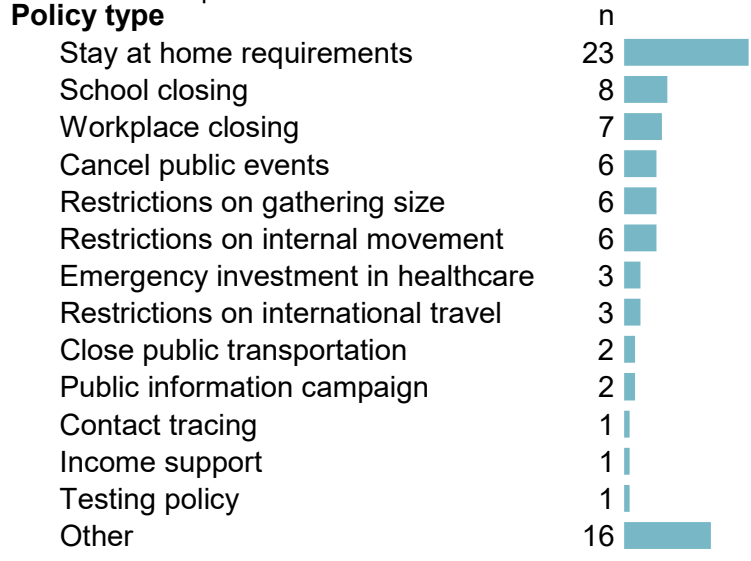
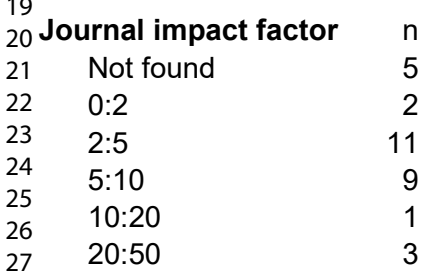
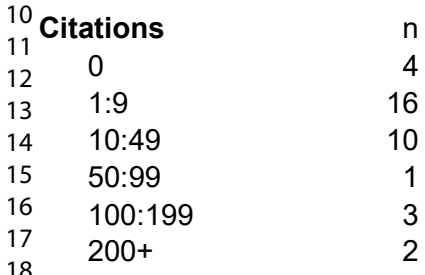
Figure 3: Main consensus results summary for key and overall questions

Caption: This chart shows the final overall ratings (left) and the key design question ratings for the consensus review of the 36 included studies, answering the degree to which the articles met the given key design question criteria. The key design question ratings were not asked for the nine included articles which selected methods assumed by the guidance to be non-appropriate. The question prompt in the figure is shortened for clarity, where the full prompt for each key question is available in the Methods section.

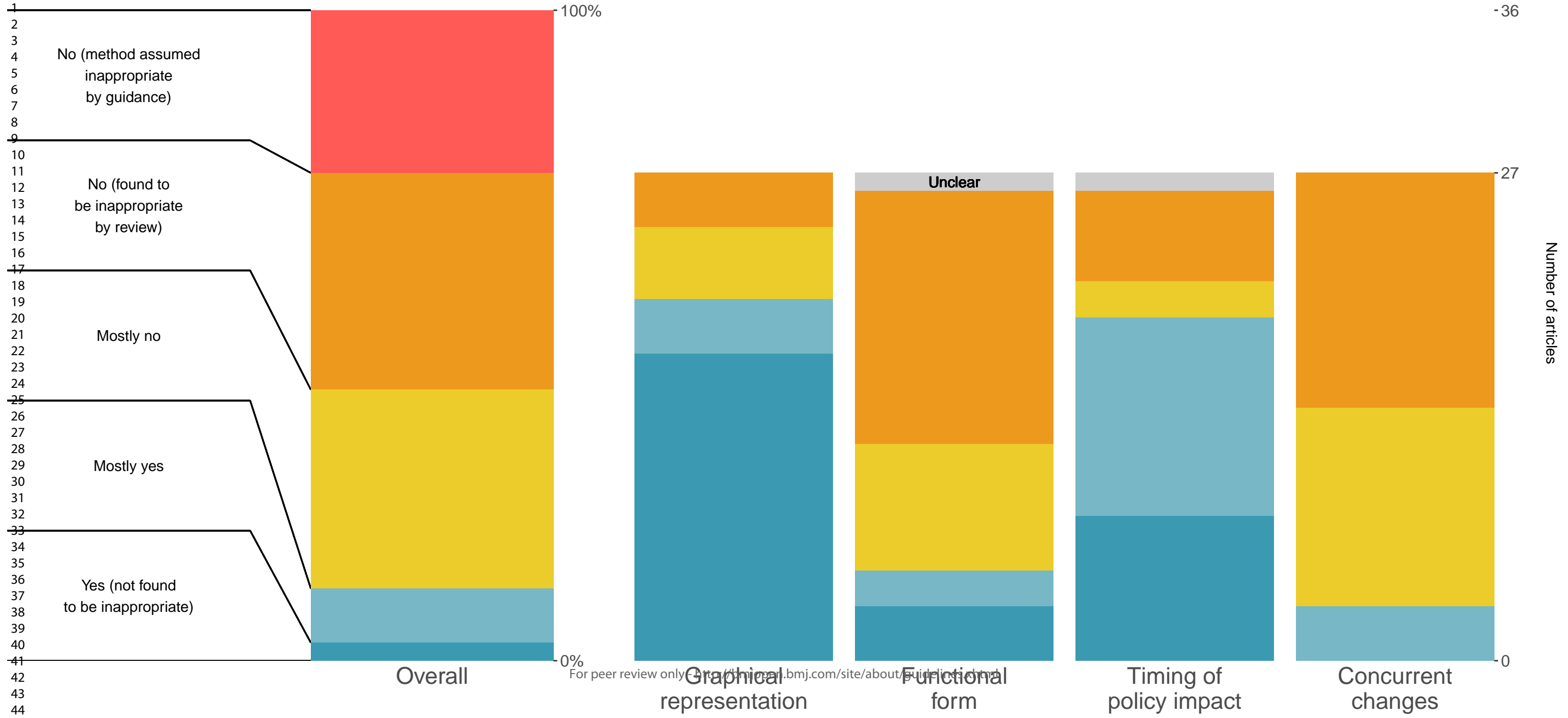
Figure 4: Comparison of independent reviews, weakest link, and direct consensus review

Caption: This chart shows the final overall ratings by three different possible metrics. The first column contains all of the independent review ratings for the 27 studies which were eventually included in our sample, noting that reviewers who either selected them as not meeting inclusion criteria or selected a method that didn't receive the full review did not contribute. The middle column contains the final consensus reviews among the 27 articles which received full review. The last column contains the weakest link rating, as described in the methods section. The question prompt in the figure is shortened for clarity, where the full prompt for each key question is available in the Methods section.





Did the study meet design criteria?



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

- 36

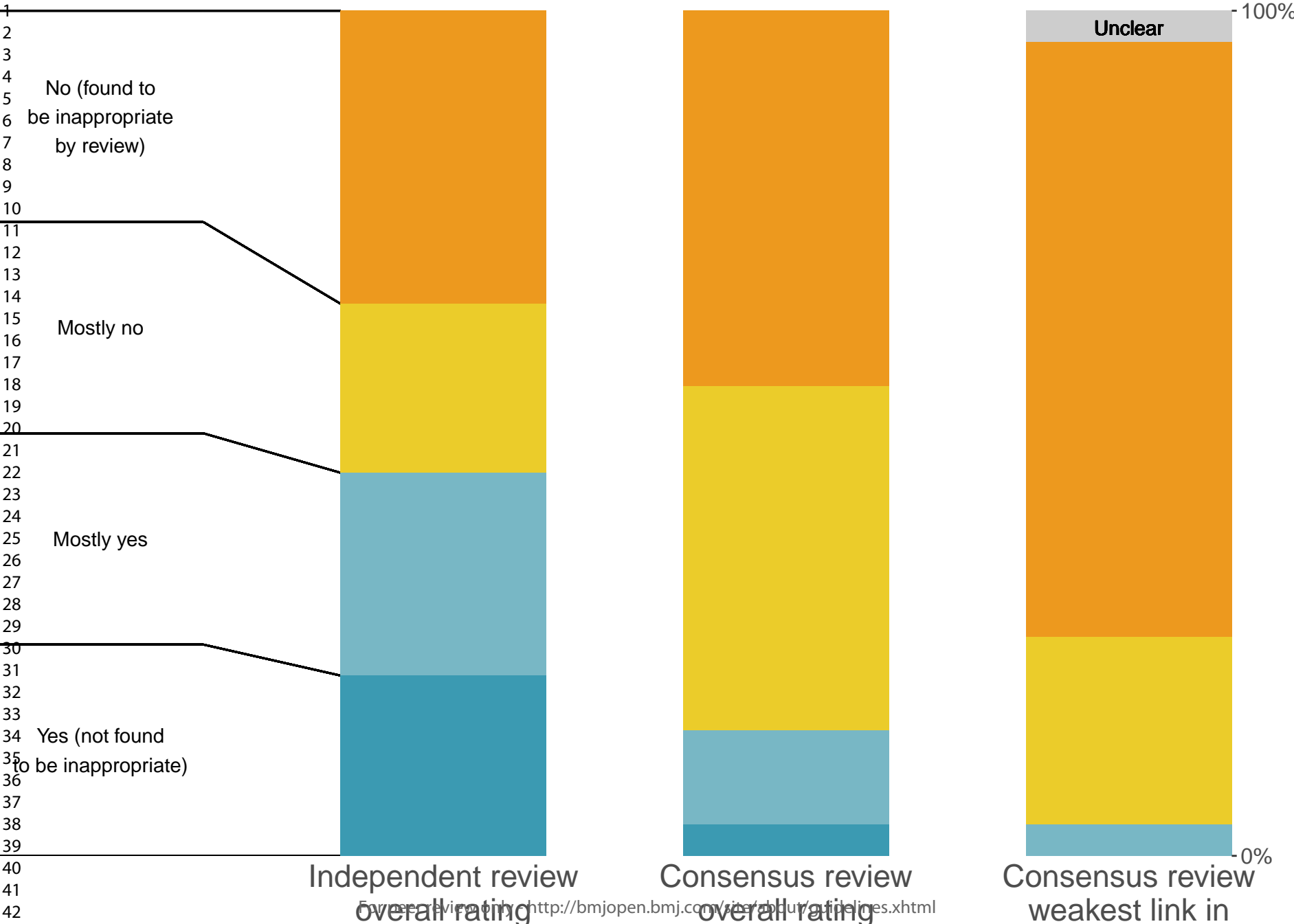
27

Number of articles

0

For peer review only - <https://bmjopen.bmj.com/site/about/guidelines.xhtml>

Did the study meet design criteria?



COVID-19 Health Policy Impact Evaluation Review

Start of Block: Main form

Q10 Administrative information

Q8 Study DOI

Q3 Reviewer number

Q54 Review type/round

The first round (Primary/independent review round) is for the independent first reviews of every article; the second (Secondary/consensus round) is for the second round of review for each article.

- Primary/independent review round (1)
- Secondary/consensus round (2)

Q50 Screening

1
2
3 Q52 Do you wish to recuse yourself from reviewing this study for any reason (e.g. social or
4 professional relationship with the authors, financial conflict of interest, etc)?
5

- 6
7 No, I do not wish to recuse myself. (1)
8
9 Yes, I recuse myself from reviewing this paper. (2)
10

11
12 *Skip To: End of Survey If Q52 = Yes, I recuse myself from reviewing this paper.*
13
14

15
16 Q51 Do you believe that this study meets the inclusion criteria for this research?
17

18 The inclusion criteria are: The primary topic of the article must be evaluating one or more
19 individual COVID-19 policies on direct COVID-19 outcomes The primary
20 exposure(s) must be a policy, defined as a government-issued order at any government level to
21 address a directly COVID-19-related outcome (e.g. mask requirements, travel restrictions, etc).
22

23 COVID-19 outcomes may include cases detected, mortality, number of tests taken, test
24 positivity rates, Rt, etc. This may NOT include indirect impacts of COVID-19 on
25 things such as income, childcare, trust in science, etc. The primary outcome
26 being examined must be a COVID-19-specific outcome, as above. The study must be
27 designed as an impact evaluation study from primary data (i.e. not primarily a predictive or
28 simulation model or meta-analysis) The study must be peer reviewed, and published in a peer-
29 reviewed journal indexed by PubMed The study must have the title and abstract available
30 via PubMed at the time of the study start date The study must be written in English
31
32
33

- 34 Yes (1)
35
36 No (2) _____
37
38

39
40 *Skip To: End of Survey If Q51 = No*
41
42

43 **Q7 Study topic information**
44

45
46 Please consult review guidance ([available here](#)) for additional guidance on answering these
47 questions.
48
49

50
51
52 Q6 Main impact sentence
53
54
55
56
57
58
59
60

1
2
3 Copy and paste the sentence from the abstract that best describes the main claim of the study
4 (e.g. "Policy X had a positive impact on outcome Y")
5
6
7

8
9

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Q9 Main COVID-19 policy type evaluated

Select all that apply. Note: categorization from the Oxford Government Response Tracker

- School closing (1)
- Workplace closing (2)
- Cancel public events (3)
- Restrictions on gathering size (4)
- Close public transportation (5)
- Stay at home requirements (6)
- Restrictions on internal movement (7)
- Restrictions on international travel (8)
- Income support (9)
- Debt/contract relief for household (10)
- Fiscal measures (11)
- Giving international support (12)
- Public information campaign (13)
- Testing policy (14)
- Contact tracing (15)
- Emergency investment in healthcare (16)
- Investment in COVID-19 vaccines (17)

Other policy response (fill in) (18)

Q12 Main COVID-19 outcome type evaluated

Select all that apply

- COVID-19 cases (1)
- COVID-19 test positivity (2)
- COVID-19 deaths (3)
- COVID-19 hospitalizations (4)
- SARS-CoV-2 infections and infection rate (e.g. effective R) (8)
- Other (fill in) (9) _____

Q13 Method(s) identification

For this section, consider only the data structure as it enters into the main statistical model. In other words, if the original dataset is of individuals at many time points, but the main statistical model uses a regional-level aggregated count of cases, the data as it enters into the main statistical model is a regional aggregate at one time point.

Q14 What is the level of aggregation for the main outcome data?

- Individual level (1)
- Regional aggregate (e.g. count, mean, etc.) (2)

1
2
3
4 Q16 How many regional units included in the main statistical model received the policy of
5 interest?
6

7
8 If 2-20, enter the number of regional units analyzed which received the policy of interest.
9

10 One (1) (1)

11
12 Two through twenty (2-20) (2)
13
14 _____
15

16 More than twenty (21+) (3)

17
18 Unclear or N/A (4) _____
19
20
21

22
23
24 Q17 How many regional units were included which did NOT receive any form of the policy of
25 interest?
26

27
28 If 2-20, enter the number of regional units analyzed which did not receive the policy of interest.
29

30 Zero (0) (1)

31
32 One (1) (2)

33
34 Two through twenty (2-20) (3)
35
36 _____
37

38 More than twenty (21+) (4)

39
40 Unclear or N/A (5) _____
41
42
43

44
45 *Display This Question:*

46 *If Q17 = Zero (0)*
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q25 Did different regions receive different intensities of the policy of interest for comparison?
4

5 For example, the study might compare places with more intense versions of policy or policies
6 vs. places with less intense versions of policy or policies, rather than just places with and
7 without the policy or policies.
8
9

- 10 Yes (regions with more intense policy were compared with regions with less intense
11 policy) (1)
12
13 No (2)
14
15 Unclear or N/A (3)
16
17
18
-

19
20
21
22 Q18 For each regional unit, how many time point observations were in the model *before* the
23 policy was enacted?
24

- 25 None (0) (1)
26
27 One (1) (2)
28
29 More than one (2+) (3)
30
31 Unclear or N/A (4) _____
32
33
34
-

35
36
37
38 Q19 For each regional unit, how many time point observations were in the model *after* the policy
39 was enacted?
40

- 41 None (0) (1)
42
43 One (1) (2)
44
45 More than one (2+) (3)
46
47 Unclear or N/A (4) _____
48
49
50
51
-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q19 = One (1)

And Q18 = One (1)

Or If

Q19 = More than one (2+)

Or If

Q18 = More than one (2+)

Q20 How would you describe the time intervals between observations?

- Days (1-5 days between observations) (1)
- Weeks (about 5-10 days between observations) (2)
- Multiple weeks (11-25 days between observations) (3)
- Monthly (26 or more days between observations) (4)
-

Display This Question:

If Q17 = Zero (0)

Q21 Did the pre-policy period for any region act as a "control" for different region post-policy enactment?

In other words, was there any pre-period in one or more region's being used to control or compare for the trends of any one or more *different* regions' post-period?

- No (pre-periods were treated as controls only within-region) (1)
- Yes (pre-periods were treated as controls with other regions) (2)
- Unclear or N/A (3)
-

Q22 Was any unit assigned the policy or the timing of the policy externally (i.e. as an experiment/trial)?

- No (observational data only) (1)
- Yes (treatment assigned as part of research or evaluation) (2)
- Unclear or N/A (3)

Display This Question:

If Q22 = Yes (treatment assigned as part of research or evaluation)

Q23 Was the assignment randomized?

- Yes (1)
- No (2)

Q27 Based on your answers above and the guidance document, please select the type of study that best resembles the design of the main analysis.

Please note that the design(s) named in the paper may not match with the method described below, nor is this the actual exact design that was used. If you believe that the design used differs from the choices below in a way that makes this choice impossible, please contact the study administrator before selecting "other."

	Design	
	Units (e.g., regions of comparison)	
	Time points measured per unit	
	Assumed counterfactual.	
	“If not for the intervention, ___”	
	Without intervention	With intervention
	Before intervention	
	After intervention	
	Cross-sectional	
	At least one	
At least one		N/A
	One time point	
	Outcome in intervention units would have been the same	

1
2
3 as the outcome in the non-intervention units.

4 Pre/post

5
6
7 At least one

8 None

9 At least one (typically one)

10 At least one (typically one)

11 Outcome would have stayed the same from the pre period to the post period.

12 Interrupted

13 time-series

14 (ITS)

15
16
17 At least one

18 None

19 More than one

20 At least one (typically several)

21 Outcome slope and level* would have continued along the same modelled trajectory
22 from the pre-period to the post period.

23 Difference-in-differences

(DiD)

24
25
26 At least one

27 At least one†

28 At least one (typically one)

29 At least one (typically one)

30 Outcome in intervention units would have changed as much as (or in parallel with) the
31 outcome in the non-intervention units.

32 Comparative interrupted time series (CITS)

33
34
35 At least one

36 At least one†

37 More than one (typically several)

38 At least one (typically several)

39 Outcome slope and level* would have changed as much as non-
40 intervention group's slope and level* changed.

41 * Assessing both slope and level only applicable if
42 there are multiple data points during the post period † Units without the
43 intervention may be the pre-period of a different unit that eventually receives the intervention.

44
45
46 Cross-sectional analysis (1)

47
48
49 Non-randomized experiment/trial (2)

50
51
52 Randomized controlled trial (3)

- 1
2
3
4 Pre/post (4)
5
6
7 Interrupted time-series (5)
8
9 Difference-in-differences (6)
10
11 Comparative interrupted time-series (7)
12
13
14 Other (please contact administrator before selecting) (8)
15 _____
16

17
18
19
20 **Q49 Design evaluation**
21
22

23
24 *Display This Question:*

25 *If Q27 = Interrupted time-series*

26 *Or Q27 = Difference-in-differences*

27 *Or Q27 = Comparative interrupted time-series*
28
29

30
31 **Q29 Does the analysis provide graphical representation of the outcome over time?**
32
33

34
35 If not "Yes" please describe (three short sentences max).
36

37 -Check for a chart that shows the outcome over time, with the dates of interest, separated by
38 policy/non policy groups if applicable. Outcomes may be aggregated for clarity (e.g. means and
39 CIs at discrete time points).
40

- 41
42 Yes (1)
43
44 Mostly yes (2) _____
45
46 Mostly no (3) _____
47
48 No (4) _____
49
50 Unclear (5) _____
51
52
53
54

1
2
3 *Display This Question:*

4 *If Q27 = Interrupted time-series*

5 *Or Q27 = Difference-in-differences*

6 *Or Q27 = Comparative interrupted time-series*
7
8
9

10 Q30 Is there sufficient pre-intervention data to characterize pre-trends in the data?
11
12

13
14 If not "Yes" please describe (three short sentences max).
15

16 -Check the chart(s) to see if there are several time points over a reasonable period of time over
17 which to establish stability and curvature in the pre-trends.
18

19
20 Yes (1)
21

22 Mostly yes (2) _____
23

24 Mostly no (3) _____
25

26 No (4) _____
27

28 Unclear (5) _____
29
30
31

32
33
34 *Display This Question:*

35 *If Q27 = Interrupted time-series*

36 *Or Q27 = Difference-in-differences*

37 *Or Q27 = Comparative interrupted time-series*
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q32 Is the pre-trend stable?
4
5

6
7 If not "Yes" please describe (three short sentences max).
8

9 -Check if there are sufficient data to reasonably determine a stable functional form for the pre-
10 trends, and that they follow a modelable functional form.
11

- 12
13 Yes (1)
14
15 Mostly yes (2) _____
16
17 Mostly no (3) _____
18
19 No (4) _____
20
21 Unclear (5) _____
22
23
24

25
26
27 *Display This Question:*

28 *If Q27 = Interrupted time-series*

29 *Or Q27 = Comparative interrupted time-series*
30

31
32 Q31 Is there sufficient post-intervention data to observe post trends in the data?
33
34

35
36 If not "Yes" please describe (three short sentences max).
37

38 -Check the chart(s) to see if there are several time points over a reasonable period of time over
39 which to establish stability and curvature in the post- trends.
40

- 41
42 Yes (1)
43
44 Mostly yes (2) _____
45
46 Mostly no (3) _____
47
48 No (4) _____
49
50 Unclear (5) _____
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q33 Is the functional form of the counterfactual (e.g. linear) well-justified and appropriate?

If not "Yes" please describe (three short sentences max).

- Check whether the authors explain and justify their choice of functional form.
- Check if there is any curvature in the pre-trend.
- Consider the nature of the outcome. Is it sensible to measure the trend of this outcome on the scale and form used? Note: infectious disease dynamics are rarely linear.
- Consider that while pre-trend fit is a necessary condition for an appropriate linear counterfactual model, it is not sufficient. Check if the authors provide justification for the functional form to continue to be of the same functional form (e.g. linear).

- Yes (1)
- Mostly yes (2) _____
- Mostly no (3) _____
- No (4) _____
- Unclear (5) _____

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q34 Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?

If not "Yes" please describe (three short sentences max).

- Check whether the authors justify the use of the date threshold relative to the date of the intervention.
- Trace the process between the intervention being put in place to when observable effects in

1
2
3 the outcome might appear over time.

4 -Consider whether there are anticipation effects (e.g. do people change behaviors before the
5 date when the intervention begins?)

6
7 -Consider whether there are lag effects. (e.g. does it take time for behaviors to change,
8 behavior change to impact infections, infections to impact testing, and data to be collected, etc?)

9 -Check if authors appropriately and directly account for these time effects.
10

11 Yes (1)

12 Mostly yes (2) _____
13

14 Mostly no (3) _____
15

16 No (4) _____
17

18 Unclear (5) _____
19
20
21
22
23

24 -----
25 *Display This Question:*

26 *If Q27 = Interrupted time-series*
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q36 Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have
4 changed during the measurement period?
5

6
7 If not "Yes" please describe (three short sentences max).
8

- 9 -Consider other policies or interventions which could impact the outcome during this time.
10 -Consider social behaviors changed which could meaningfully impact the outcome during this
11 time.
12 -Consider economic conditions changed which could meaningfully impact the outcome during
13 this time.
14 -Note that the actual concurrent changes do not need to happen during the period of
15 measurement, just their effects.
16
17

- 18
19 Yes (1)
20
21 Mostly yes (2) _____
22
23 Mostly no (3) _____
24
25 No (4) _____
26
27 Unclear (5) _____
28
29
30
31

32
33 *Display This Question:*

34 *If Q27 = Difference-in-differences*

35 *Or Q27 = Comparative interrupted time-series*
36
37

38 Q53

39 Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have
40 changed during the measurement period, differently for policy and non-policy regions?
41

42 If not "Yes" please describe (three short sentences max).
43

44 -Consider any uncontrolled factor which could have influenced the outcome differently in policy
45 and non-policy regions.
46

47 -This may include (but is not limited to)

- 48 -Other policies
49 -Social behaviors
50 -Economic conditions
51

52 -Are these factors justified as having negligible impact?
53

54 -If justified, is the argument that these have negligible impact convincing?
55
56
57
58
59

1
2
3 -Note that the actual concurrent changes do not need to happen during the period of
4 measurement, just their effects.
5

- 6
7 Yes (1)
8
9 Mostly yes (2) _____
10
11 Mostly no (3) _____
12
13 No (4) _____
14
15 Unclear (5) _____
16
17
18
19

20
21 *Display This Question:*

22 *If Q27 = Interrupted time-series*

23 *Or Q27 = Difference-in-differences*

24 *Or Q27 = Comparative interrupted time-series*
25
26
27

28 Q38

29 Did authors provide diagnostics or show robustness and/or sensitivity of results to alternative
30 model choices?
31

32
33
34 If not "Yes" please describe (three short sentences max).
35

- 36 Yes (1)
37
38 Mostly yes (2) _____
39
40 Mostly no (3) _____
41
42 No (4) _____
43
44 Unclear (5) _____
45
46
47
48
49

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q39 Given the above, do you believe that the design is appropriate for identifying the policy impact(s) of interest?

This should be taken as independent of what you believe about other studies, and/or the feasibility of other designs.

If not "Yes" please describe (three short sentences max).

Yes (1)

Mostly yes (2) _____

Mostly no (3) _____

No (4) _____

Unclear (5) _____

Display This Question:

If Q54 = Secondary/consensus round

Q55 General and/or additional comments on this paper from consensus discussion. This may include any additional information worth commenting on regarding the paper, difficulties encountered evaluating it, etc.

(three short sentences max)

End of Block: Main form

Appendix 1: Changes from pre-registered protocol and justifications

The full, original pre-registered protocol is available here: <https://osf.io/7nbk6>

Inclusion criteria

Minor language edits were made to the inclusion criteria to improve clarity and fix grammatical and typographical errors. This largely centered around improving clarity that a study must estimate the quantitative impact of policies that had already been enacted. The word “quantitative” was not explicitly stated in the original version.

Procedures

The original protocol specified that each article would receive two independent reviewers. This was increased to three reviewers per article once it became clear both that the number of articles which would be accepted for full review was lower than expectations, and that there would be substantial differences in opinion between reviewers.

Statistical analysis

Firstly, the original protocol specified that 95% confidence intervals would be calculated. However, after further discussion and review, we determined that sampling-based confidence intervals were not appropriate. Our results are not indicative nor intended to be representative of any super- or target-population, and as such sampling-based error is not an appropriate metric for the conclusions of this study.

Secondly, the original protocol specified Kappa-based interrater reliability statistics. However, using three reviewers, rather than the originally registered two, meant that most Kappa statistics would not be appropriate for our review process. Given the three-rater, four-level ordinal scale used, we opted instead to use Krippendorff’s Alpha.

Review tool

A number of changes were made to the review tool during the course of the review process. While the original protocol included logic to allow pre/post for review in some of the key questions, this was removed for consistency with the guidance document.

The remaining changes to the review tool were error corrections and clarifications (e.g. correcting the text for the concurrent changes sections in difference-in-differences so that it

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

stated “uncontrolled” concurrent changes, and distinguishing the DiD/CITS requirements from the ITS requirements to emphasize differential concurrent changes).

Appendix 2: Full search terms

Note: The search filter for COVID-19 and SARS-CoV-2 were the exact search terms used for the National Library of Medicine one-click search option at the time of the protocol development and when the search took place. This reflects that some of the early literature referred to Wuhan specifically (both in geographic reference for where the SARS-CoV-2 was initially found, and unfortunately also early naming of the virus/disease) before official naming conventions became ubiquitous in the literature. In order to comprehensively capture the literature and use searching best practices, we used the most standard and recommended terms.

(((((wuhan[All Fields] AND ("coronavirus"[MeSH Terms] OR "coronavirus"[All Fields])) AND 2019/12[PDAT] : 2030[PDAT]) OR 2019-nCoV[All Fields] OR 2019nCoV[All Fields] OR COVID-19[All Fields] OR SARS-CoV-2[All Fields])

AND ("impact"[TIAB] OR "effect"[TIAB])

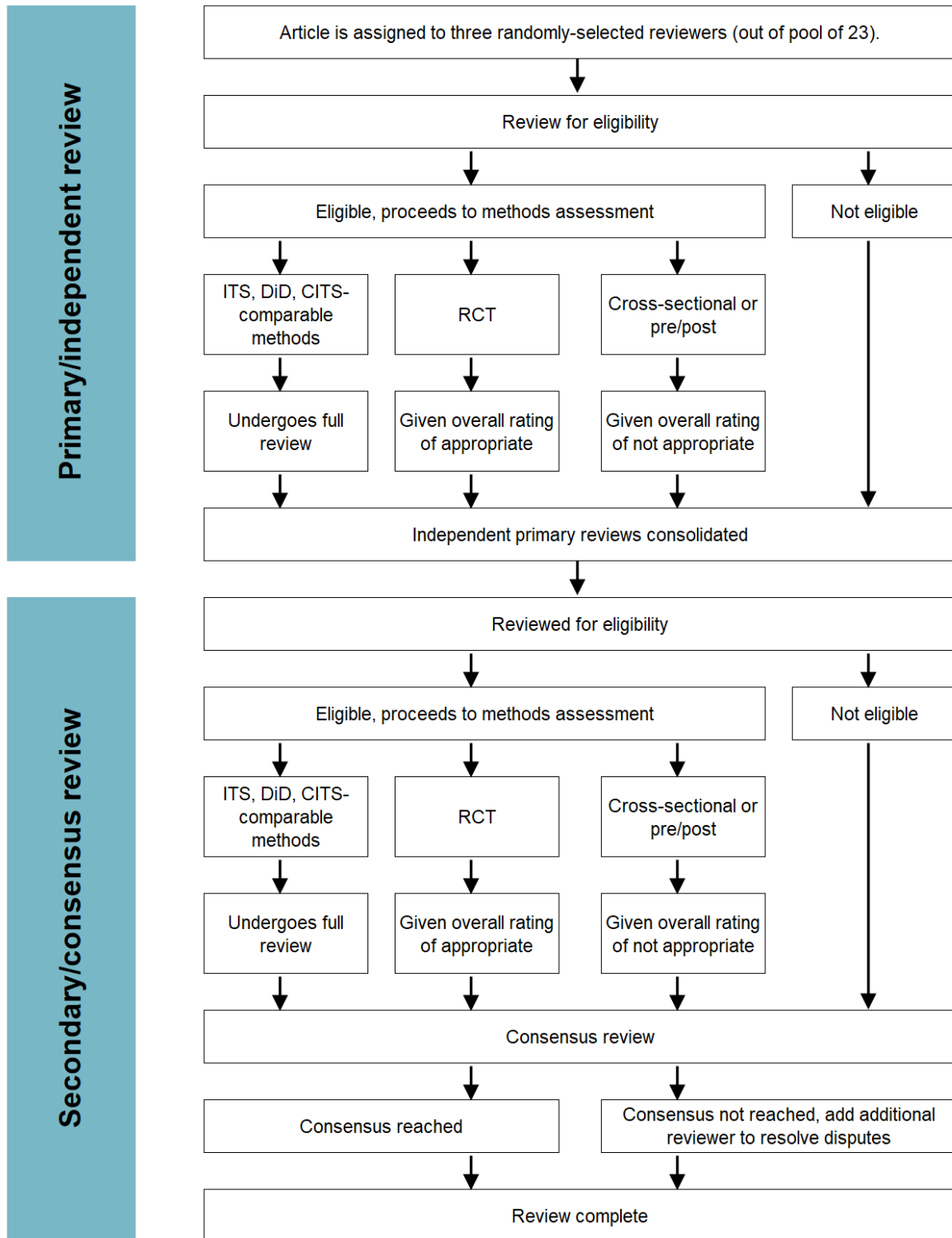
AND ("policy"[TIAB] OR "policies"[TIAB] OR "order"[TIAB] OR "mandate"[TIAB])

AND ("countries"[TIAB] OR "country"[TIAB] OR "state"[TIAB] OR "provinc"[TIAB] OR "county"[TIAB] OR "parish"[TIAB] OR "region"[TIAB] OR "city"[TIAB] OR "cities"[TIAB] OR "continent"[TIAB] "Asia"[TIAB] OR "Europe"[TIAB] OR "Africa"[TIAB] OR "America"[TIAB] OR "Australia"[TIAB] OR "Antarctica"[TIAB] OR "Afghanistan"[TIAB] OR "Aland Islands"[TIAB] OR "Åland Islands"[TIAB] OR "Albania"[TIAB] OR "Algeria"[TIAB] OR "American Samoa"[TIAB] OR "Andorra"[TIAB] OR "Angola"[TIAB] OR "Anguilla"[TIAB] OR "Antarctica"[TIAB] OR "Antigua"[TIAB] OR "Argentina"[TIAB] OR "Armenia"[TIAB] OR "Aruba"[TIAB] OR "Australia"[TIAB] OR "Austria"[TIAB] OR "Azerbaijan"[TIAB] OR "Bahamas"[TIAB] OR "Bahrain"[TIAB] OR "Bangladesh"[TIAB] OR "Barbados"[TIAB] OR "Barbuda"[TIAB] OR "Belarus"[TIAB] OR "Belgium"[TIAB] OR "Belize"[TIAB] OR "Benin"[TIAB] OR "Bermuda"[TIAB] OR "Bhutan"[TIAB] OR "Bolivia"[TIAB] OR "Bonaire"[TIAB] OR "Bosnia"[TIAB] OR "Botswana"[TIAB] OR "Bouvet Island"[TIAB] OR "Brazil"[TIAB] OR "British Indian Ocean Territory"[TIAB] OR "Brunei"[TIAB] OR "Bulgaria"[TIAB] OR "Burkina Faso"[TIAB] OR "Burundi"[TIAB] OR "Cabo Verde"[TIAB] OR "Cambodia"[TIAB] OR "Cameroon"[TIAB] OR "Canada"[TIAB] OR "Cayman Islands"[TIAB] OR "Central African Republic"[TIAB] OR "Chad"[TIAB] OR "Chile"[TIAB] OR "China"[TIAB] OR "Christmas Island"[TIAB] OR "Cocos Islands"[TIAB] OR "Colombia"[TIAB] OR "Comoros"[TIAB] OR "Congo"[TIAB] OR "Congo"[TIAB] OR "Cook Islands"[TIAB] OR "Costa Rica"[TIAB] OR "Côte d'Ivoire"[TIAB] OR "Croatia"[TIAB] OR "Cuba"[TIAB] OR "Curaçao"[TIAB] OR "Cyprus"[TIAB] OR "Czechia"[TIAB] OR "Denmark"[TIAB] OR "Djibouti"[TIAB] OR "Dominica"[TIAB] OR "Dominican Republic"[TIAB] OR "Ecuador"[TIAB] OR "Egypt"[TIAB] OR "El Salvador"[TIAB] OR "Equatorial Guinea"[TIAB] OR "Eritrea"[TIAB] OR "Estonia"[TIAB] OR "Eswatini"[TIAB] OR "Ethiopia"[TIAB] OR "Falkland Islands"[TIAB] OR "Faroe Islands"[TIAB] OR "Fiji"[TIAB] OR "Finland"[TIAB] OR "France"[TIAB] OR "French Guiana"[TIAB] OR "French Polynesia"[TIAB] OR "French Southern

1
2
3 Territories[TIAB] OR "Futuna"[TIAB] OR "Gabon"[TIAB] OR "Gambia"[TIAB] OR
4 "Georgia"[TIAB] OR "Germany"[TIAB] OR "Ghana"[TIAB] OR "Gibraltar"[TIAB] OR
5 "Greece"[TIAB] OR "Greenland"[TIAB] OR "Grenada"[TIAB] OR "Grenadines"[TIAB] OR
6 "Guadeloupe"[TIAB] OR "Guam"[TIAB] OR "Guatemala"[TIAB] OR "Guernsey"[TIAB] OR
7 "Guinea"[TIAB] OR "Guinea-Bissau"[TIAB] OR "Guyana"[TIAB] OR "Haiti"[TIAB] OR "Heard
8 Island"[TIAB] OR "Herzegovina"[TIAB] OR "Holy See"[TIAB] OR "Honduras"[TIAB] OR "Hong
9 Kong"[TIAB] OR "Hungary"[TIAB] OR "Iceland"[TIAB] OR "India"[TIAB] OR "Indonesia"[TIAB]
10 OR "Iran"[TIAB] OR "Iraq"[TIAB] OR "Ireland"[TIAB] OR "Isle of Man"[TIAB] OR "Israel"[TIAB]
11 OR "Italy"[TIAB] OR "Jamaica"[TIAB] OR "Jan Mayen Islands"[TIAB] OR "Japan"[TIAB] OR
12 "Jersey"[TIAB] OR "Jordan"[TIAB] OR "Kazakhstan"[TIAB] OR "Keeling Islands"[TIAB] OR
13 "Kenya"[TIAB] OR "Kiribati"[TIAB] OR "Korea"[TIAB] OR "Korea"[TIAB] OR "Kuwait"[TIAB] OR
14 "Kyrgyzstan"[TIAB] OR "Lao People's Democratic Republic"[TIAB] OR "Laos"[TIAB] OR
15 "Latvia"[TIAB] OR "Lebanon"[TIAB] OR "Lesotho"[TIAB] OR "Liberia"[TIAB] OR "Libya"[TIAB]
16 OR "Liechtenstein"[TIAB] OR "Lithuania"[TIAB] OR "Luxembourg"[TIAB] OR "Macao"[TIAB] OR
17 "Madagascar"[TIAB] OR "Malawi"[TIAB] OR "Malaysia"[TIAB] OR "Maldives"[TIAB] OR
18 "Mali"[TIAB] OR "Malta"[TIAB] OR "Malvinas"[TIAB] OR "Marshall Islands"[TIAB] OR
19 "Martinique"[TIAB] OR "Mauritania"[TIAB] OR "Mauritius"[TIAB] OR "Mayotte"[TIAB] OR
20 "McDonald Islands"[TIAB] OR "Mexico"[TIAB] OR "Micronesia"[TIAB] OR "Moldova"[TIAB] OR
21 "Monaco"[TIAB] OR "Mongolia"[TIAB] OR "Montenegro"[TIAB] OR "Montserrat"[TIAB] OR
22 "Morocco"[TIAB] OR "Mozambique"[TIAB] OR "Myanmar"[TIAB] OR "Namibia"[TIAB] OR
23 "Nauru"[TIAB] OR "Nepal"[TIAB] OR "Netherlands"[TIAB] OR "Nevis"[TIAB] OR "New
24 Caledonia"[TIAB] OR "New Zealand"[TIAB] OR "Nicaragua"[TIAB] OR "Niger"[TIAB] OR
25 "Nigeria"[TIAB] OR "Niue"[TIAB] OR "Norfolk Island"[TIAB] OR "North Macedonia"[TIAB] OR
26 "Northern Mariana Islands"[TIAB] OR "Norway"[TIAB] OR "Oman"[TIAB] OR "Pakistan"[TIAB]
27 OR "Palau"[TIAB] OR "Panama"[TIAB] OR "Papua New Guinea"[TIAB] OR "Paraguay"[TIAB]
28 OR "Peru"[TIAB] OR "Philippines"[TIAB] OR "Pitcairn"[TIAB] OR "Poland"[TIAB] OR
29 "Portugal"[TIAB] OR "Principe"[TIAB] OR "Puerto Rico"[TIAB] OR "Qatar"[TIAB] OR
30 "Réunion"[TIAB] OR "Romania"[TIAB] OR "Russian Federation"[TIAB] OR "Rwanda"[TIAB] OR
31 "Saba"[TIAB] OR "Saint Barthélemy"[TIAB] OR "Saint Helena"[TIAB] OR "Saint Kitts"[TIAB] OR
32 "Saint Lucia"[TIAB] OR "Saint Martin"[TIAB] OR "Saint Pierre and Miquelon"[TIAB] OR "Saint
33 Vincent"[TIAB] OR "Samoa"[TIAB] OR "San Marino"[TIAB] OR "Sao Tome"[TIAB] OR
34 "Sark"[TIAB] OR "Saudi Arabia"[TIAB] OR "Senegal"[TIAB] OR "Serbia"[TIAB] OR
35 "Seychelles"[TIAB] OR "Sierra Leone"[TIAB] OR "Singapore"[TIAB] OR "Sint Eustatius"[TIAB]
36 OR "Sint Maarten"[TIAB] OR "Slovakia"[TIAB] OR "Slovenia"[TIAB] OR "Solomon
37 Islands"[TIAB] OR "Somalia"[TIAB] OR "South Africa"[TIAB] OR "South Georgia"[TIAB] OR
38 "South Sandwich Islands"[TIAB] OR "South Sudan"[TIAB] OR "Spain"[TIAB] OR "Sri
39 Lanka"[TIAB] OR "State of Palestine"[TIAB] OR "Sudan"[TIAB] OR "Suriname"[TIAB] OR
40 "Svalbard"[TIAB] OR "Sweden"[TIAB] OR "Switzerland"[TIAB] OR "Syria"[TIAB] OR "Syrian
41 Arab Republic"[TIAB] OR "Tajikistan"[TIAB] OR "Thailand"[TIAB] OR "Timor-Leste"[TIAB] OR
42 "Tobago"[TIAB] OR "Togo"[TIAB] OR "Tokelau"[TIAB] OR "Tonga"[TIAB] OR "Trinidad"[TIAB]
43 OR "Tunisia"[TIAB] OR "Turkey"[TIAB] OR "Turkmenistan"[TIAB] OR "Turks and Caicos"[TIAB]
44 OR "Tuvalu"[TIAB] OR "Uganda"[TIAB] OR "UK"[TIAB] OR "Ukraine"[TIAB] OR "United Arab
45 Emirates"[TIAB] OR "United Kingdom"[TIAB] OR "United Republic of Tanzania"[TIAB] OR
46 "United States Minor Outlying Islands"[TIAB] OR "United States of America"[TIAB] OR
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 "Uruguay"[TIAB] OR "USA"[TIAB] OR "Uzbekistan"[TIAB] OR "Vanuatu"[TIAB] OR
4 "Venezuela"[TIAB] OR "Viet Nam"[TIAB] OR "Vietnam"[TIAB] OR "Virgin Islands"[TIAB] OR
5 "Virgin Islands"[TIAB] OR "Wallis"[TIAB] OR "Western Sahara"[TIAB] OR "Yemen"[TIAB] OR
6 "Zambia"[TIAB] OR "Zimbabwe"[TIAB] OR "Alabama"[TIAB] OR "Alaska"[TIAB] OR
7 "Arizona"[TIAB] OR "Arkansas"[TIAB] OR "California"[TIAB] OR "Colorado"[TIAB] OR
8 "Connecticut"[TIAB] OR "Delaware"[TIAB] OR "Florida"[TIAB] OR "Georgia"[TIAB] OR
9 "Hawaii"[TIAB] OR "Idaho"[TIAB] OR "Illinois"[TIAB] OR "Indiana"[TIAB] OR "Iowa"[TIAB] OR
10 "Kansas"[TIAB] OR "Kentucky"[TIAB] OR "Louisiana"[TIAB] OR "Maine"[TIAB] OR
11 "Maryland"[TIAB] OR "Massachusetts"[TIAB] OR "Michigan"[TIAB] OR "Minnesota"[TIAB] OR
12 "Mississippi"[TIAB] OR "Missouri"[TIAB] OR "Montana"[TIAB] OR "Nebraska"[TIAB] OR
13 "Nevada"[TIAB] OR "New Hampshire"[TIAB] OR "New Jersey"[TIAB] OR "New Mexico"[TIAB]
14 OR "New York"[TIAB] OR "North Carolina"[TIAB] OR "North Dakota"[TIAB] OR "Ohio"[TIAB] OR
15 "Oklahoma"[TIAB] OR "Oregon"[TIAB] OR "Pennsylvania"[TIAB] OR "Rhode Island"[TIAB] OR
16 "South Carolina"[TIAB] OR "South Dakota"[TIAB] OR "Tennessee"[TIAB] OR "Texas"[TIAB] OR
17 "Utah"[TIAB] OR "Vermont"[TIAB] OR "Virginia"[TIAB] OR "Washington"[TIAB] OR "West
18 Virginia"[TIAB] OR "Wisconsin"[TIAB] OR "Wyoming"[TIAB] OR "Ontario"[TIAB] OR
19 "Quebec"[TIAB] OR "Nova Scotia"[TIAB] OR "New Brunswick"[TIAB] OR "Manitoba"[TIAB] OR
20 "British Columbia"[TIAB] OR "Prince Edward Island"[TIAB] OR "Saskatchewan"[TIAB] OR
21 "Alberta"[TIAB] OR "Newfoundland"[TIAB] OR "Labrador"[TIAB])
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Appendix 3: Article review flow diagram



Review version with references removed; NOT FOR DISTRIBUTION

Policy evaluation in COVID-19: A guide to common design issues

Noah A Haber, Emma Clarke-Deelder, Joshua A Salomon, Avi Feller, Elizabeth A Stuart

Noah A Haber, ScD*

noahhaber@stanford.edu

Meta Research Innovation Center at Stanford University

Stanford University

1265 Welch Rd

Palo Alto, CA 94305

(650) 497-0811

Emma Clarke-Deelder, MPhil

Department of Global Health & Population

Harvard T. H. Chan School of Public Health

665 Huntington Avenue

Building 1, room 1104

Boston, Massachusetts 02115

Joshua A Salomon, PhD

Department of Medicine

Center for Health Policy and Center for Primary Care and Outcomes Research

Stanford University School of Medicine

Encina Commons, Room 118

615 Crothers Way

Stanford, CA 94305-6019

Avi Feller, PhD

Goldman School of Public Policy

University of California, Berkeley

2607 Hearst Avenue

Room 309

Berkeley, CA 94720

Elizabeth A Stuart, PhD

Department of Mental Health

Johns Hopkins Bloomberg School of Public Health

624 N. Broadway

Hampton House 839

Baltimore, MD 21205

* corresponding author

Review version with references removed; NOT FOR DISTRIBUTION

Abstract

Policy responses to COVID-19, particularly those related to non-pharmaceutical interventions, are unprecedented in scale and scope. Researchers and policymakers are striving to understand the impact of these policies on a variety of outcomes. Policy impact evaluations always require a complex combination of circumstance, study design, data, statistics, and analysis. Beyond the issues that are faced for any policy, evaluation of COVID-19 policies is complicated by additional challenges related to infectious disease dynamics and lags, lack of direct observation of key outcomes, and a multiplicity of interventions occurring on an accelerated time scale.

In this paper, we (1) introduce the basic suite of policy impact evaluation designs for observational data, including cross-sectional analyses, pre/post, interrupted time-series, and difference-in-differences analysis, (2) demonstrate key ways in which the requirements and assumptions underlying these designs are often violated in the context of COVID-19, and (3) provide decision-makers and reviewers a conceptual and graphical guide to identifying these key violations. The overall goal of this paper is to help policy-makers, journal editors, journalists, researchers, and other research consumers understand and weigh the strengths and limitations of evidence that is essential to decision-making.

Introduction

The response to the global COVID-19 pandemic has demanded urgent decision making in the face of substantial uncertainties. Policies to arrest transmission, including stay-at-home orders and other non-pharmaceutical interventions (NPIs), have wide-reaching consequences that touch many aspects of well being. Decision-making in the public interest requires evaluating and weighing the evidence on both intended and unintended consequences in order to best predict outcomes. The wide range of policy interventions implemented by different jurisdictions may yield opportunities for learning from what has already happened to inform future policymaking, and we have observed a proliferation of studies aimed at such policy evaluations. However, policy evaluation requires a complex combination of circumstance, data, study design, analysis, and interpretation in order to be informative.

Policy impact evaluation aims to answer questions about the extent to which the realized outcomes given a particular policy would have been different in the absence of that policy. Estimating the causal impact of the policy with observational data is challenging because what would have happened in the absence of the policy change (the “counterfactual”) is, by definition, unobserved. Randomized controlled trials (RCTs) of policies related to COVID-19 interventions may not always be practical or ethical. In this context, a large and growing number of studies have attempted to evaluate the impact of COVID-19 policies using observational data. There

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 are many potential pitfalls in the use of observational data for evaluation generally, and some
4 additional methodological design challenges relating to COVID-19 policies in particular.
5

6
7 This paper provides a graphical guide to policy impact evaluations for COVID-19, targeted to
8 decision-makers, researchers and evidence curators. Our aim is to provide a coherent
9 framework for conceptualizing and identifying common pitfalls in COVID-19 policy evaluation.
10 Importantly, this should not be taken either as a comprehensive guide to policy evaluation more
11 broadly or as guidance on performing analysis, which may be found elsewhere. Rather, we
12 review relevant study designs for policy evaluations — including pre/post, interrupted time
13 series, and difference-in-difference approaches — and provide guidance and tools for
14 identifying key issues with each type of study as they relate to NPIs and other COVID-19 policy
15 interventions. Improving our ability to identify key pitfalls will enhance our ability to identify and
16 produce valid and useful evidence for informing policymaking.
17
18
19
20

21 Common policy evaluation designs and their pitfalls 22 in COVID-19 23 24

25 Identifying the type of design 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

Review version with references removed; NOT FOR DISTRIBUTION

Table 1: Summary definitions of policy impact evaluation designs commonly used for COVID-19

Design	Units (e.g., regions of comparison)		Time points measured per unit		Assumed counterfactual. “If not for the intervention, _____”
	With intervention	Without intervention	Before intervention	After intervention	
Cross-sectional	At least one	At least one	N/A	One time point	Outcome in intervention units would have been the same as the outcome in the non-intervention units.
Pre/post Figure 1A	At least one	None	At least one (typically one)	At least one (typically one)	Outcome would have stayed the same from the pre period to the post period.
Interrupted time-series (ITS) Figure 1B	At least one	None	More than one	At least one (typically several)	Outcome slope and level* would have continued along the same modelled trajectory from the pre-period to the post period.
Difference-in-differences (DiD) Figure 1C	At least one	At least one [†]	At least one (typically one)	At least one (typically one)	Outcome in intervention units would have changed as much as (or in parallel with) the outcome in the non-intervention units.
Comparative interrupted time series (CITS) Figure 1D	At least one	At least one [†]	More than one (typically several)	At least one (typically several)	Outcome slope and level* would have changed as much as non-intervention group's slope and level* changed.

* Assessing both slope and level only applicable if there are multiple data points during the post period
[†] Units without the intervention may be the pre-period of a different unit that eventually receives the intervention.

Identifying the underlying design in a given analysis often requires using a combination of the methods as reported and evaluating the data structure that is used for the main analysis, as shown in Table 1. COVID-19-related policy evaluation analyses typically fall under these categories. In most cases, the design can be categorized using a combination of whether there are also units that did not receive the treatment (columns 2-3) and whether there are time points both before and after intervention for those units (columns 4-5). The final column describes the implied counterfactual, discussed further in subsequent sections. Cross sectional designs typically compare units with vs without the treatment at single time points. Pre/post studies typically compare within units who received the intervention at two points: before and after a policy. Interrupted time-series analyses compare outcomes within units within units who received the intervention at greater than two time points before the intervention vs with at least one (typically multiple) after the intervention. Difference-in-differences analysis compares the outcome change in units which received the intervention with those that did not (or have not yet), with at least one point before and one after the intervention. In cases with multiple periods, that may involve a comparison with the pre-policy period of one region with the post-period of a different region, even though all regions eventually receive the intervention.

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

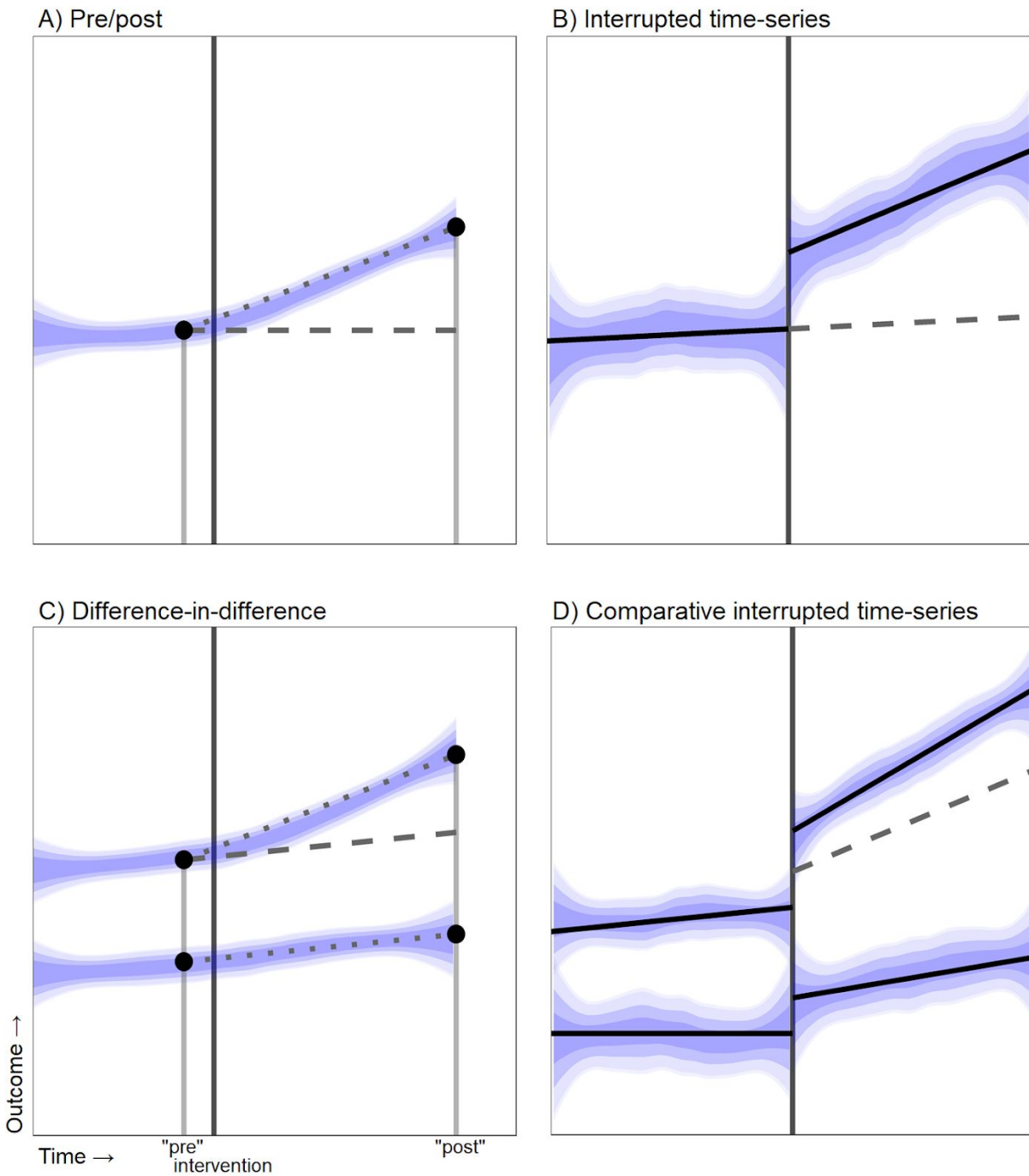
Methods descriptions may not always provide a precise or reliable guide to which of the design approaches has been used. Some studies do not explicitly name these designs (or may classify them differently); and these are only a small fraction of designs and frameworks that are possible to use for policy evaluation. Studies may have data at multiple time points but are effectively cross-sectional. DiD, ITS, and CITS designs based on repeated cross-sectional data are sometimes described as “cross-sectional” instead of longitudinal. The term “event study” is often used to refer to studies with a single unit and one change over time resembling ITS, but may refer to other designs. Although ITS is often used to describe changes in one unit, it may also refer to settings in which many treated units adopt an intervention over time. Studies will also frequently employ multiple designs, while others use more complex methods of generating counterfactuals. Definitions of these terms vary widely, and the definitions above should be considered as guidance only.

Policy impact evaluation design foundations for COVID-19

The simplest design is the cross-sectional analysis, which compares COVID-19 outcomes between units of observation (e.g., cities) at a single calendar time or time since an event, typically post-intervention. These studies are unlikely to be appropriate for COVID-19-related policy evaluations, but provide a useful starting point for reasoning about different designs. Just as with comparisons of non-randomized medical treatments, the localities that adopt a particular policy likely differ substantially from those that don't on both observed and unobserved characteristics on a number of dimensions, including epidemic status and timing.

Figure 1: Longitudinal designs overview

Review version with references removed; NOT FOR DISTRIBUTION



48 This chart shows four canonical longitudinal designs. In all cases: the blue shading
49 represents the underlying data trends, the solid vertical grey line represents the time of
50 intervention, the grey dashed lines represent the assumed counterfactual in the absence of
51 the intervention, as discussed in the text. The impact estimate is obtained by comparing the
52 outcomes observed for the treated unit in the post period (the solid line) with the implied
53 counterfactual line (the dashed line). In the case of the pre/post and
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 difference-in-differences panels the large black dots represent the time of measurement,
4 connected by the grey dotted lines.
5
6

7 Given the challenges in a simple cross-sectional comparison, which compare post-intervention
8 outcomes, it is important to consider longitudinal designs, which instead look at differences or
9 trends across time, as summarized in Figure 1. These can be distinguished by the data used
10 and the construction of the counterfactual. Pre/post, for example, has only one unit, measured
11 at two time points. Two common strategies expand on the logic and data requirements of the
12 pre/post design. Interrupted time series designs (Figure 1B) incorporate multiple time points
13 before the intervention, and usually multiple time points after the intervention, to enable a more
14 complete view on changes in levels and trends that are temporally related to the intervention.
15 Difference-in-difference designs (Figure 1C) add a set of comparison points from a group or
16 location that did not have the intervention. Another related design (comparative interrupted
17 time-series, Figure 1D, discussed only briefly here), uses both aspects — a change over time
18 and a comparison group — to compare the observed change in slopes for the intervention
19 group with the change in slope for the comparison group.
20
21
22
23

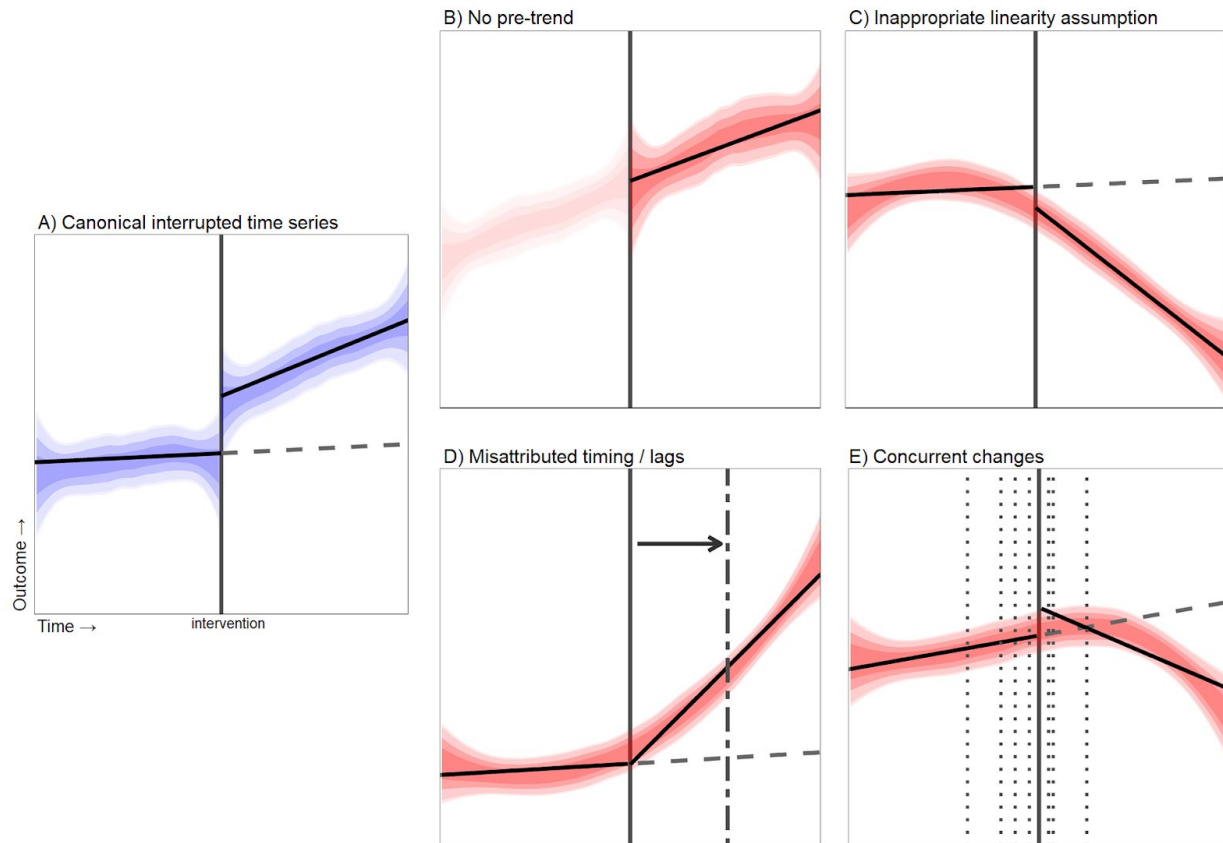
24 Pre/post studies

25
26 The simplest longitudinal design is a pre/post analysis, where some outcome is observed before
27 policy implementation, and again after, in a single group (Figure 1A). Pre/post studies are
28 analogous to a single arm trial with no control and only a single follow-up observation after
29 treatment.. This effectively imposes the assumption that the counterfactual trend is completely
30 flat (i.e., that the outcome in the post-period in the absence of the policy change is the same as
31 the value of the outcome before the policy change) without accounting for pre-existing
32 underlying trends, and attributing all outcome changes completely to the intervention of interest.
33 Just as the outcomes for an individual patient might be expected to change before and after
34 treatment, for reasons unrelated to the treatment, outcomes related to policy interventions will
35 change for reasons not caused by the policy. Infection rates, for example, would not be
36 expected to remain stationary except in very specific circumstances, but a pre/post
37 measurement would assume that any changes in infection rates are attributable to the policy.
38
39
40
41
42

43 Interrupted time-series

44
45 Figure 2: Interrupted time-series graphical guidance for identifying common pitfalls
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION



This chart shows one canonical design for ITS (blue, Panel A) and four panels demonstrating common issues with ITS analysis (red, panels B-E) discussed in the text. In all cases: the lag/red shading represents the underlying data trends, the vertical grey line represents the time of intervention, the grey dashed lines represent the assumed counterfactual in the absence of the intervention. In panel D) the dash-dot line represents the time at which the policy is expected to impact the outcome. In panel E), the vertical dotted lines represent concurrent events and changes.

Interrupted time-series (ITS) is a strategy that uses a projection of the pre-policy outcome trend as a counterfactual for how the outcome would have changed if the policy had not been introduced. In other words, in the absence of the policy change, ITS assumes the outcome would have continued on its pre-policy trend during the study period. ITS can be a useful tool in policy evaluation because it allows researchers to account for underlying trends in the outcome and, by comparing the treated unit (or location) to itself; it can therefore eliminate some of the confounding concerns that arise in cross-sectional or pre-post studies.

However, the validity of ITS depends critically on how well counterfactual trends in the outcome are modelled, and whether the policy of interest is the only relevant change during the study period. In the canonical setting (Figure 2A), the pre-policy trend is stable and can be feasibly modelled with the available data; the researcher appropriately models the timing of the change in the slope and/or level of the outcome; the researcher has sufficient information to conclude

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 that there were no other changes during the study period that would be expected to influence
4 the outcome. These elements are largely not satisfied in studies of COVID-related policy, as
5 described below.
6
7

8 ITS relies critically on modelled trends of the outcome over time. Key components of ITS
9 analyses include both visual and statistical examination of trends, preferentially alongside a
10 theoretical justification of the model used. At a minimum, analyses should provide graphical
11 representation of the data and model over time to examine whether pre-trend outcomes are
12 stable, all trends are well-fit to the data, “interrupted” at the appropriate time point, and sensibly
13 modelled (Figure 2B). In the case where an ITS includes a large number of units (e.g. states), it
14 can be difficult to display this information graphically.
15
16
17

18 One common pitfall in ITS is adoption of inappropriate assumptions on the outcome trend
19 (Figure 2C). The estimate of policy impact will be biased if a linear trend is assumed but the
20 outcome and response to interventions instead follow nonlinear trends (either before or after the
21 policy). In some cases, transformation of the outcome, for example using a log scale, may
22 improve the suitability of a linear model. Imposing linearity inappropriately is a serious risk in the
23 context of COVID-19, as trends in infectious disease dynamics are inherently non-linear. For
24 intuition, terms such as “exponential growth,” “flattening,” and “s-curves” all refer to non-linear
25 infectious disease trends. Depending on the particular situation, non-linearity or other modelled
26 trends can have complicated and counterintuitive impact on policy impact. Apparent linearity
27 may also be temporary and an artifact of testing, which may give a misleading impression that
28 linear models for infectious disease trends are appropriate indefinitely. While some use linear
29 progression in order to avoid more complex infectious disease models, in fact, linear projections
30 impose strict and often unrealistic models, generally resulting in an inappropriate counterfactual.
31
32
33
34
35

36 Researchers can easily misattribute the timing of the policy impact, resulting in spurious
37 inference and bias (Figure 2D). Some public health policies can be expected to translate into
38 immediate results (e.g., smoking bans and acute coronary events). In contrast, nearly every
39 outcome of interest in COVID-19 exhibits complex and difficult to infer time lags typically in the
40 realm of many weeks. The time between policy implementation and expected effect in the data
41 can be large and highly variable. For example, in order to see the impact of a mask order, first
42 the mask order takes effect, then people change their behaviors over time to comply with the
43 order (or sometimes the reverse in the case of anticipation effects), mask use behavior
44 produces changes in infections, then infections later result in symptoms, symptoms induce
45 people to seek testing, the tests must then be processed in labs, and then finally the results get
46 reported in data monitoring efforts. Selection of lead/lag time should be justifiable *a priori* or
47 external data. Selecting a lag based on the data risks issues comparable to p-hacking.
48
49
50
51

52 Finally, and perhaps most concerningly in the context of COVID-19, ITS fails when the policy of
53 interest coincides in time with other changes that affect the outcome (Figure 2E). For example, if
54 both mask and bar closure orders are rolled out together as a package, ITS cannot isolate the
55 impact of bar closures specifically. These changes do not need to have taken place exactly
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

concurrently with the policy implementation date of interest; they merely need to have some effect within the time period of measurement to result in potentially serious bias in effect estimates if unaddressed. ITS will also likely be biased if, during the study period, there is a change in the way the outcome data is collected or measured. This might occur if the introduction of a COVID-19 control policy is combined with an effort to collect better data on infection or mortality cases. Analogously, if an RCT involves randomizing people to a group receiving both A and B vs. control, we typically can't disentangle the effects of A from the effects of B, unless we also have separate A- and B-only arms. Ultimately, if multiple things are changing at the same time, ITS may not be an appropriate design for policy evaluation.

COVID-19 policies rarely arrive alone; they are typically created alongside other policies, unofficial action, and large scale behavior changes which themselves impact COVID-19-related outcomes. In some cases, anticipation of a policy may induce behavior change before the actual policy takes effect. The policies themselves may have been chosen due to the expectation of change in disease outcomes, which introduces additional biases related to "reverse" causality.

Table 2: Checklist for identifying common pitfalls for ITS to evaluate COVID-19 policy

Key design questions. If any answer is "no," this analysis is unlikely to be appropriate or useful for estimating the impact of the intervention of interest.	Details and suggestions for identifying issues:
Does the analysis provide graphical representation of the outcome over time?	-Check for a chart that shows the outcome over time, with the dates of interest. Outcomes may be aggregated for clarity (e.g. means and CIs at discrete time points).
Is there sufficient pre-intervention data to characterize pre-trends in the data?	-Check the chart(s) to see if there are several time points over a reasonable period of time over which to establish stability and curvature in the pre-trends.
Is the pre-trend stable?	-Check if there are sufficient data to reasonably determine a stable functional form for the pre-trends, and that they follow a modelable functional form.
Is the functional form of the counterfactual (e.g. linear) well-justified and appropriate?	-Check whether the authors explain and justify their choice of functional form. -Check if there is any curvature in the pre-trend. -Consider the nature of the outcome. Is it sensible to measure the trend of this outcome on the scale and form used? Note: infectious disease dynamics are rarely linear. -Consider that while pre-trend fit is a necessary condition for an appropriate linear counterfactual model, it is not sufficient. Check if the authors provide justification for the functional form to continue to be of the same functional form (e.g. linear).
Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?	-Check whether the authors justify the use of the date threshold relative to the date of the intervention. -Trace the process between the intervention being put in place to when observable effects in the outcome might appear over time. -Consider whether there are anticipation effects (e.g. do people change behaviors before the date when the intervention begins?) -Consider whether there are lag effects. (e.g. does it take time for behaviors to change, behavior change to impact infections, infections to impact testing, and data to be collected, etc?)

Review version with references removed; NOT FOR DISTRIBUTION

	-Check if authors appropriately and directly account for these time effects.
Is this policy the only thing to happen which could have impacted the outcome during the measurement period, differently for policy and non-policy regions??	<ul style="list-style-type: none"> -Consider other policies or interventions which could impact the outcome during this time. -Consider social behaviors changed which could meaningfully impact the outcome during this time. -Consider economic conditions changed which could meaningfully impact the outcome during this time. -Note that the actual concurrent changes do not need to happen during the period of measurement, just their effects.

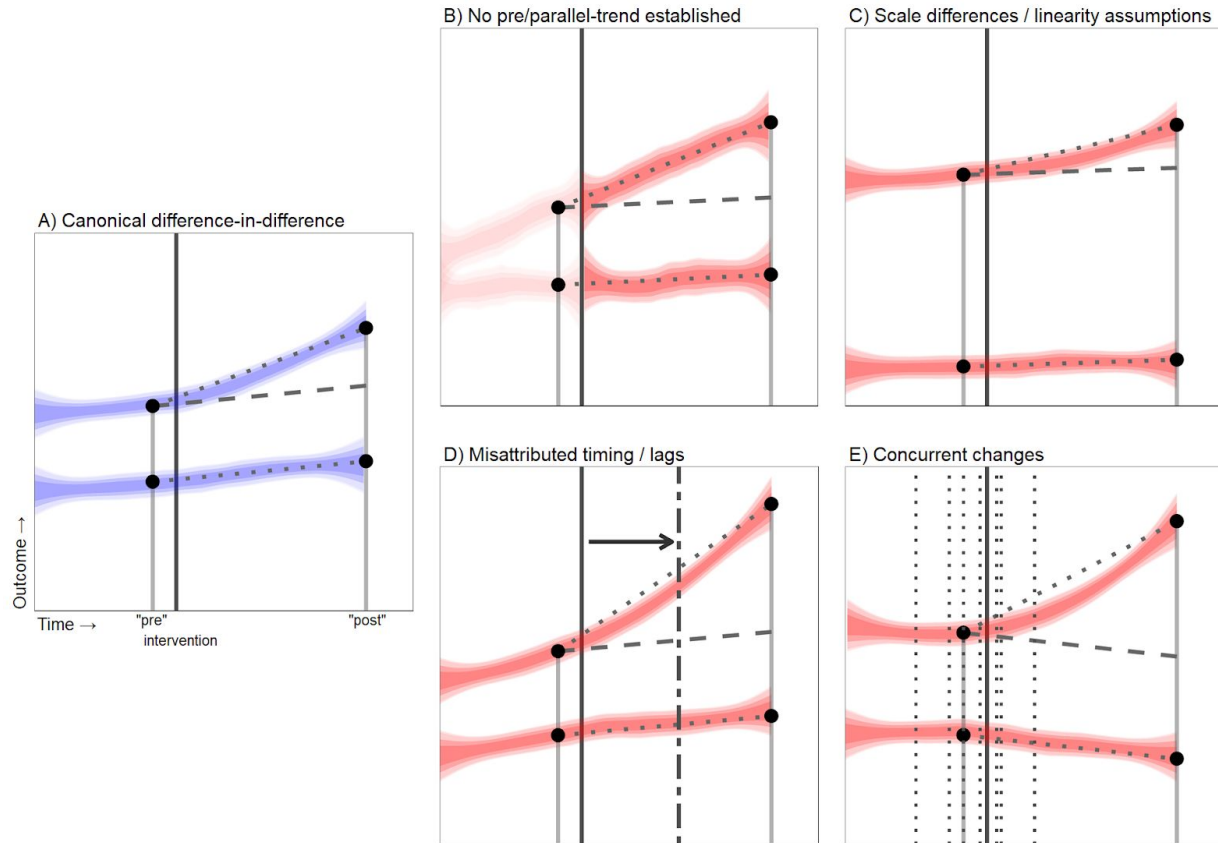
These issues are summarized as a checklist of questions to identify common pitfalls in Table 2.

Difference-in-differences

The difference-in-difference (DiD) approach uses concurrent non-intervention groups as a counterfactual. Typically, this consists of one set of units (e.g., regions) that had the intervention and one set that did not, with each measured before and after the intervention took place. DiD is more directly analogous to a non-randomized medical study with at least one treatment and control group but limited observation before and after treatment. In contrast to ITS, which compares a unit with itself over time, DiD compares differences between treatment arms or units at two observation points. In many analyses, a DiD approach is implied by comparing regions over time, without formally naming or modelling it. Other DiD approaches use interventions implemented at multiple time points.

Figure 3: Difference-in-differences graphical guidance for identifying common pitfalls

Review version with references removed; NOT FOR DISTRIBUTION



This chart shows one canonical design for DiD (blue, Panel A) and four panels demonstrating common issues with DiD analysis (red, panels B-E). In all cases: the blue/red shading represents the underlying data trends, the vertical grey line represents the time of intervention, the grey dashed lines represent the assumed counterfactual in the absence of the intervention. In panel D) the dash-dot line represents the time at which the policy is expected to impact the outcome. In panel E), the vertical dotted lines represent concurrent events and changes.

One key component of the standard DiD approach is the parallel counterfactual trends assumption: that the intervention and comparison groups would have had parallel trends over time in the absence of the intervention. In some cases, the parallel trends assumption may be referenced or examined implicitly but not named.

Ideally, pre-intervention trends would be shown to be clearly identifiable, stable, of a similar level, and parallel between groups. With only one observation before and only one after the intervention, assessment of the plausibility of the parallel counterfactual trends assumption is not possible. Absent this confirmation the evaluation runs the risk of biased estimation due to differential pre-trends (Figure 3B). Pre-trends approaching the ceiling or floor may also not be informative about stable and parallel pre-trends. Empirical assessment of whether pre-intervention trends were parallel and stable between groups is possible when multiple observations are available at multiple time points before the intervention, noting that this can

Review version with references removed; NOT FOR DISTRIBUTION

begin to resemble a CITS design. In this scenario, pre-trend data should be visually and statistically established and documented. While parallel trends before intervention (which we can observe and may be testable) do not guarantee parallel *counterfactual* trends in the post-intervention period (which we cannot observe and are generally untestable), examining pre-intervention parallel trends is a minimal requirement for DiD reliability.

It is also important to consider the scale and level on which the outcome is measured (Figure 3C). As with ITS, if the outcomes in the treatment and comparison groups are moving in parallel on a logged scale, they will not be moving in parallel on a natural scale. Level differences by themselves may be a problem for COVID-19 outcomes, as infectious disease transmission dynamics dictate that infection risks are related to the prevalence of infected people in a population, i.e. the rate of change is linked intrinsically to the level. A population with an extremely low prevalence will tend to have an inherently slower rise in infection rates than an otherwise identical population with merely a low prevalence. Just as importantly, large level differences in the outcome between intervention and comparison groups is often indicative of other important differences between comparators, which may result in other assumptions being violated.

While DiD is in some ways more robust to very specific kinds of timing effects (Figure 3D) and concurrent changes (Figure 3E), it also introduces additional risks. DiD effectively doubles the opportunity for concurrent changes to spuriously impact results, since they can occur in the treatment or comparison groups. As above, this can become even more problematic for DiD in the typical case where intervention groups enact more or very contextually different policies than non-intervention groups. Even cases where concurrent changes happen equally in both treatment and comparison groups can lead to overwhelming bias, particularly when approaching the maximum or minimum levels of the outcome. If either the treatment or control group is approaching the floor (e.g. 0% prevalence) or ceiling for an outcome of interest due to other policies concurrent in both places (e.g. national lockdowns, but region-level differences in mask policy), this can lead to bias when comparing changes between the two groups.

Table 3: Checklist for identifying common pitfalls for DiD to evaluate COVID-19 policy

Key design questions. If any answer is “no,” this analysis is unlikely to be appropriate or useful for estimating the impact of the intervention of interest	Details and suggestions for inspection:
Does the analysis provide graphical representation of the outcome over time?	-Check for a graph that shows the outcome over time for all groups, with the dates of interest. Outcomes may be aggregated for clarity (e.g. mean and CI at discrete time points).
Is there sufficient pre-intervention data to observe both pre and post trends in the data?	-Check the chart(s) to see if there are several time points over a reasonable period of time over which to establish stability and curvature in the pre- and post- trends.
Are the pre-trends stable?	-Check if there are sufficient graphical data to reasonably determine a stable functional form for the pre-trends, and that they follow a modelable functional form.

Review version with references removed; NOT FOR DISTRIBUTION

Are the pre-trends parallel?	-Observe if the trends in the intervention and comparison groups appear to move together at the same rate at the same time.
Are the pre-trends at a similar level?	-Check if the trends in the intervention and comparison groups are at similar levels. -Note that non-level trends exacerbates other problems with the analysis, including linearity assumptions
Are intervention and non-groups broadly comparable?	-Consider areas where comparison groups may be dissimilar for comparison beyond just the level of the outcome.
Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?	-Check whether the authors justify the use of the date threshold relative to the date of the intervention. -Trace the process between the intervention being put in place to when observable effects in the outcome might appear over time. -Consider whether there are anticipation effects (e.g. do people change behaviors before the date when the intervention begins?) -Consider whether there are lag effects. (e.g. does it take time for behaviors to change, behavior change to impact infections, infections to impact testing, and data to be collected, etc?) -Check if authors appropriately and directly account for these time effects.
Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have changed during the measurement period, differently for policy and non-policy regions?	-Consider any uncontrolled factor which could have influenced the outcome differently in policy and non-policy regions. -This may include (but is not limited to) -Other policies -Social behaviors -Economic conditions -Are these factors justified as having negligible impact? -If justified, is the argument that these have negligible impact convincing? -Note that the actual concurrent changes do not need to happen during the period of measurement, just their effects.

Similarly to the ITS section, these issues are summarized as a checklist of questions to identify common pitfalls in Table 3.

Discussion

In recent months, there has been a proliferation of research evaluating policies related to the COVID-19 pandemic. As with other areas of COVID-19 research, quality has been highly variable, with low quality studies resulting in poorly or mis-informed policy decisions, wasted resources, and undermined trust in research. To support high quality policy evaluations, in this paper we describe common approaches to evaluating policies using observational data, and describe key issues that can arise in applying these approaches. We hope that this guidance can help support researchers, editors, reviewers, and decision-makers in conducting high quality policy evaluations and in assessing the strength of the evidence that has already been published.

Policy evaluation — far from a simple task in normal circumstances — is particularly challenging during a pandemic. Cross-sectional comparisons of states or countries are likely to be biased by selection into treatment: for example, countries with worse outbreaks may be more likely to

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 implement policies such as mask requirements. In analyses of changes over time – such as
4 single-unit studies using interrupted time-series or multi-unit comparisons using
5 difference-in-differences or comparative interrupted time-series – it may not be possible to parse
6 apart the effects of different policies implemented around the same time, such as mask
7 mandates paired with limits on social gatherings. Analyses of changes over time may also be
8 biased if disease or human behavioral dynamics are not modelled appropriately. This can be
9 challenging because case counts typically do not grow linearly and there is often a lag between
10 a policy change and a behavioral response.
11
12
13

14 This guidance should be considered minimal screening to identify low quality policy impact
15 evaluation in COVID-19, but is in no way sufficient to identify high quality evidence or
16 actionability. Decision-makers and researchers should pay particular attention to the relevance
17 of the intervention as it was evaluated to relevant decisions being made. The evaluated impact
18 of a program encouraging mask use through messages might not be informative about mask
19 requirement orders. Differences in level of aggregation may be important, such as ecological
20 fallacy arising from a situation in which areas with higher overall mask use have higher
21 transmission, but transmission is actually lower for individuals wearing masks. Policy impact
22 evaluation is only as useful as the question it asks, data it uses, and the way it is analyzed.
23 Problems with measurement, spillover effects, generalizability, changes in measurement
24 overtime (e.g. varying test availability), statistics, testing robustness to alternative assumptions,
25 and many issues can undermine an otherwise robust evaluation, and are not discussed here.
26
27
28
29

30 While this guidance is not comprehensive, it may help inform study designs not covered here.
31 Issues with comparative interrupted time-series and synthetic control methods, for example, are
32 broadly similar to the issues with difference-in-differences analyses we discuss here. Other
33 approaches may include adjustment and matching based observational causal inference
34 designs, instrumental variables and related quasi-experimental approaches, and randomized
35 controlled trials. Each has its own set of practical, ethical, and inferential limitations.
36
37
38

39 In the face of these challenges, we recommend careful scrutiny and attention to potential
40 sources of bias in COVID-19-related policy evaluations, but we remain optimistic about the
41 potential for robust evaluations to inform decision-making. Researchers and decision-makers
42 should triangulate across a large variety of approaches from theory to evidence, invest in better
43 data and more reliable and useful evidence wherever feasible, clearly acknowledge limitations
44 and potential sources of bias, and acknowledge when actionable evidence is not feasible. We
45 anticipate increasing opportunities for better examining policies moving forward, particularly if
46 policies and interventions are designed with policy impact evaluation and data collection in
47 mind.
48
49
50

51 The COVID-19 pandemic requires urgent decisions about policies that affect millions of people's
52 lives in significant ways. High quality evidence on the effects of these policies is critical to
53 informing decision-making, but is very hard to generate. Evidence-based decision-making
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

depends on research that carefully considers potential sources of bias, and clearly communicates underlying assumptions and sources of uncertainty.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



PRISMA 2009 Checklist

Section/topic	#	Checklist item	Reported on page #
TITLE			
Title	1	Identify the report as a systematic review, meta-analysis, or both.	1
ABSTRACT			
Structured summary	2	Provide a structured summary including, as applicable: background; objectives; data sources; study eligibility criteria, participants, and interventions; study appraisal and synthesis methods; results; limitations; conclusions and implications of key findings; systematic review registration number.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known.	3
Objectives	4	Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS).	3
METHODS			
Protocol and registration	5	Indicate if a review protocol exists, if and where it can be accessed (e.g., Web address), and, if available, provide registration information including registration number.	4
Eligibility criteria	6	Specify study characteristics (e.g., PICOS, length of follow-up) and report characteristics (e.g., years considered, language, publication status) used as criteria for eligibility, giving rationale.	4
Information sources	7	Describe all information sources (e.g., databases with dates of coverage, contact with study authors to identify additional studies) in the search and date last searched.	5
Search	8	Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated.	5 + appendix 1
Study selection	9	State the process for selecting studies (i.e., screening, eligibility, included in systematic review, and, if applicable, included in the meta-analysis).	5
Data collection process	10	Describe method of data extraction from reports (e.g., piloted forms, independently, in duplicate) and any processes for obtaining and confirming data from investigators.	5
Data items	11	List and define all variables for which data were sought (e.g., PICOS, funding sources) and any assumptions and simplifications made.	6-7
Risk of bias in individual studies	12	Describe methods used for assessing risk of bias of individual studies (including specification of whether this was done at the study or outcome level), and how this information is to be used in any data synthesis.	6-7
Summary measures	13	State the principal summary measures (e.g., risk ratio, difference in means).	8



PRISMA 2009 Checklist

Synthesis of results	14	Describe the methods of handling data and combining results of studies, if done, including measures of consistency (e.g., I^2) for each meta-analysis.	8
----------------------	----	---	---

Page 1 of 2

Section/topic	#	Checklist item	Reported on page #
Risk of bias across studies	15	Specify any assessment of risk of bias that may affect the cumulative evidence (e.g., publication bias, selective reporting within studies).	8
Additional analyses	16	Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified.	8
RESULTS			
Study selection	17	Give numbers of studies screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally with a flow diagram.	8
Study characteristics	18	For each study, present characteristics for which data were extracted (e.g., study size, PICOS, follow-up period) and provide the citations.	9
Risk of bias within studies	19	Present data on risk of bias of each study and, if available, any outcome level assessment (see item 12).	9-13
Results of individual studies	20	For all outcomes considered (benefits or harms), present, for each study: (a) simple summary data for each intervention group (b) effect estimates and confidence intervals, ideally with a forest plot.	9-11
Synthesis of results	21	Present results of each meta-analysis done, including confidence intervals and measures of consistency.	10
Risk of bias across studies	22	Present results of any assessment of risk of bias across studies (see Item 15).	9-12
Additional analysis	23	Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]).	10-12
DISCUSSION			
Summary of evidence	24	Summarize the main findings including the strength of evidence for each main outcome; consider their relevance to key groups (e.g., healthcare providers, users, and policy makers).	12
Limitations	25	Discuss limitations at study and outcome level (e.g., risk of bias), and at review-level (e.g., incomplete retrieval of identified research, reporting bias).	13
Conclusions	26	Provide a general interpretation of the results in the context of other evidence, and implications for future research.	14
FUNDING			
Funding	27	Describe sources of funding for the systematic review and other support (e.g., supply of data); role of funders for the systematic review.	23

From: Moher D, Liberati A, Tetzlaff J, Altman DG, The PRISMA Group (2009). Preferred Reporting Items for Systematic Reviews and Meta-Analyses: The PRISMA Statement. PLoS Med 6(7): e1000097. doi:10.1371/journal.pmed1000097

For peer review only - <http://bmjopen.bmj.com/site/about/guidelines.xhtml>
For more information, visit: www.prisma-statement.org



PRISMA 2009 Checklist

For peer review only

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47

BMJ Open

Problems with Evidence Assessment in COVID-19 Health Policy Impact Evaluation: A systematic review of study design and evidence strength

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-053820.R2
Article Type:	Original research
Date Submitted by the Author:	10-Nov-2021
Complete List of Authors:	<p>Haber, Noah ; Stanford University, Clarke-Deelder, Emma; Harvard University T H Chan School of Public Health, Department of Global Health and Population Feller, Avi; University of California Berkeley, Goldman School of Public Policy Smith, Emily; George Washington University School of Public Health and Health Services, Department of Global Health Salomon, Joshua ; Stanford University MacCormack-Gelles, Benjamin; Harvard University T H Chan School of Public Health, Department of Global Health and Population Stone, Elizabeth M.; Johns Hopkins University Bloomberg School of Public Health, Department of Health Policy and Management Bolster-Foucault, Clara; McGill University, Epidemiology, Biostatistics, and Occupational Health Daw, Jamie R.; Columbia University Mailman School of Public Health, Health Policy and Management Hatfield, Laura; Harvard Medical School, Biostatistics Fry, Carrie E.; Vanderbilt University, Department of Health Policy Boyer, Christopher B.; Harvard University T H Chan School of Public Health, Department of Epidemiology Ben-Michael, Eli; University of California Berkeley, Department of Statistics Joyce, Caroline M.; McGill University, Epidemiology, Biostatistics, and Occupational Health Linas, Beth S.; Johns Hopkins University Bloomberg School of Public Health, Department of Epidemiology; MITRE Corp Schmid, Ian; Johns Hopkins University Bloomberg School of Public Health, Department of Mental Health Au, Eric; The University of Sydney, School of Public Health Wieten, Sarah; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Jarrett, Brooke; Johns Hopkins University, Epidemiology Axfors, Cathrine; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Nguyen, Van; Stanford University, Meta Research Innovation Center at Stanford University (METRICS) Griffin, Beth; RAND Corp, Bilinski, Alyssa; Harvard University Graduate School of Arts and Sciences Stuart, Elizabeth A; Johns Hopkins University Bloomberg School of Public</p>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

	Health, Department of Mental Health
Primary Subject Heading:	Health policy
Secondary Subject Heading:	Research methods, Public health, Epidemiology, Global health
Keywords:	COVID-19, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, STATISTICS & RESEARCH METHODS





I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

Problems with Evidence Assessment in COVID-19 Health Policy Impact Evaluation: A systematic review of study design and evidence strength

Noah A. Haber, ScD¹, Emma Clarke-Deelder, MPhil², Avi Feller, PhD³, Emily R. Smith, ScD⁴, Joshua Salomon, PhD⁵, Benjamin MacCormack-Gelles, MS², Elizabeth M. Stone, MS⁶, Clara Bolster-Foucalt, MScPH⁷, Jamie R. Daw, PhD⁸, Laura A. Hatfield, PhD⁹, Carrie E. Fry, PhD¹⁰, Christopher B. Boyer, MPH¹¹, Eli Ben-Michael, PhD¹², Caroline M. Joyce, MPH⁷, Beth S. Linas, PhD, MHS^{13,14}, Ian Schmid, ScM¹⁵, Eric H. Au, MPH¹⁶, Sarah E. Wieten, PhD¹, Brooke A Jarrett, MSPH¹³, Cathrine Axfors, MD, PhD¹, Van Thu Nguyen, PhD¹, Beth Ann Griffin, PhD¹⁷, Alyssa Bilinski, MS¹⁸, Elizabeth A. Stuart, PhD¹⁵

1. Meta-Research Innovation Center at Stanford (METRICS), Stanford University, Stanford, CA, USA
2. Department of Global Health and Population, Harvard T. H. Chan School of Public Health, Boston, MA, USA
3. Goldman School of Public Policy, UC Berkeley, Berkeley, CA, USA
4. Department of Global Health, Milken Institute School of Public Health, George Washington University, Washington, D.C, USA
5. Center for Health Policy and Center for Primary Care and Outcomes Research, Stanford University, Stanford, CA, USA
6. Department of Health Policy and Management, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
7. Epidemiology, Biostatistics, and Occupational Health, McGill University, Montreal, Canada
8. Health Policy and Management, Columbia University Mailman School of Public Health, New York, NY, USA
9. Department of Health Care Policy, Harvard Medical School, Boston, MA, USA
10. Department of Health Policy, Vanderbilt University, Nashville, TN, USA
11. Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA
12. Institute for Quantitative Social Science, Harvard University, Cambridge, MA, USA
13. Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
14. Clinical Quality and Informatics, MITRE Corp, McLean, VA, USA
15. Department of Mental Health, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
16. School of Public Health, University of Sydney, Sydney, Australia
17. RAND Corporation, Arlington, VA, USA
18. Interfaculty Initiative in Health Policy, Harvard Graduate School of Arts and Sciences, Cambridge, MA, USA

Corresponding author:

Noah A. Haber, ScD

noahhaber@stanford.edu

Meta Research Innovation Center at Stanford University

Stanford University

1265 Welch Rd

Palo Alto, CA 94305

(650) 497-0811

Abstract

Introduction: Assessing the impact of COVID-19 policy is critical for informing future policies. However, there are concerns about the overall strength of COVID-19 impact evaluation studies given the circumstances for evaluation and concerns about the publication environment.

Methods: We included studies that were primarily designed to estimate the quantitative impact of one or more implemented COVID-19 policies on direct SARS-CoV-2 and COVID-19 outcomes. After searching PubMed for peer-reviewed articles published on November 26, 2020 or earlier and screening, all studies were reviewed by three reviewers first independently and then to consensus. The review tool was based on previously developed and released review guidance for COVID-19 policy impact evaluation.

Results: After 102 articles were identified as potentially meeting inclusion criteria, we identified 36 published articles that evaluated the quantitative impact of COVID-19 policies on direct COVID-19 outcomes. Nine studies were set aside because the study design was considered inappropriate for COVID-19 policy impact evaluation (n=8 pre/post; n=1 cross-sectional), and 27 articles were given a full consensus assessment. 20/27 met criteria for graphical display of data, 5/27 for functional form, 19/27 for timing between policy implementation and impact, and only 3/27 for concurrent changes to the outcomes. Only 4/27 were rated as overall appropriate. Including the 9 studies set aside, reviewers found that only four of the 36 identified published and peer-reviewed health policy impact evaluation studies passed a set of key design checks for identifying the causal impact of policies on COVID-19 outcomes.

Discussion: The reviewed literature directly evaluating the impact of COVID-19 policies largely failed to meet key design criteria for inference of sufficient rigor to be actionable by policy-makers. More reliable evidence review is needed to both identify and produce policy-actionable evidence, alongside the recognition that actionable evidence is often unlikely to be feasible.

Strengths and limitations

- This study is based on previously released review guidance for discerning and evaluating critical minimal methodological design aspects of the COVID-19 health policy impact evaluation.
- The review tool assesses critical aspects of study design grounded in impact evaluation methods that must be true for the papers to provide useful policy impact evaluation, including what type of impact evaluation method was used, graphical display of outcomes data, functional form for the outcomes, timing between policy and impact, concurrent changes to the outcomes, and an overall rating.
- This study used a consensus reviewer model with three reviewers in order to obtain replicable results for study strength ratings.
- While the vast majority of studies in our sample received low ratings for useful causal policy impact evaluation, they may make other contributions to the literature.
- Because our review tool was limited to a very narrow - albeit critical - set of items, weaknesses in other aspects not reviewed (e.g. data quality or other aspects of statistical inference) may further weaken studies that were found to meet our criteria.

Introduction

Policy decisions to mitigate the impact of COVID-19 on morbidity and mortality are some of the most important issues policymakers have had to make since January 2020. Decisions regarding which policies are enacted depend in part on the evidence base for those policies, including understanding what impact past policies had on COVID-19 outcomes[1,2] Unfortunately, there are substantial concerns that much of the existing literature may be methodologically flawed, which could render its conclusions unreliable for informing policy. The combination of circumstances being difficult for strong impact evaluation, the importance of the topic, and concerns over the publication environment may lead to the proliferation of low strength studies.

High-quality causal evidence requires a combination of rigorous methods, clear reporting, appropriate caveats, and the appropriate circumstances for the methods used.[3–6] Rigorous evidence is difficult in the best of circumstances, and the circumstances for evaluating non-pharmaceutical intervention (NPI) policy effects on COVID-19 are particularly challenging.[5] The global pandemic has yielded a combination of a large number of concurrent policy and non-policy changes, complex infectious disease dynamics, and unclear timing between policy implementation and impact; all of this makes isolating the causal impact of any particular policy or policies exceedingly difficult.[7]

The scientific literature on COVID-19 is exceptionally large and fast growing. Scientists published more than 100,000 papers related to COVID-19 in 2020.[8] There is some general concern that the volume and speed[9,10] at which this work has been produced may result in a literature that is overall low quality and unreliable.[11–15]

Given the importance of the topic, it is critical that decision-makers are able to understand what is known and knowable[5,16] from observational data in COVID-19 policy, as well as what is unknown and/or unknowable.

Motivated by concerns about the methodological strength of COVID-19 policy evaluations, we set out to review the literature using a set of methodological design checks tailored to common policy impact evaluation methods. Our primary objective was to evaluate each paper for methodological strength and reporting, based on pre-existing review guidance developed for this purpose.[17] As a secondary objective, we also studied our own process: examining the consistency, ease of use, and clarity of this review guidance.

This protocol differs in several ways from more traditional systematic review protocols given the atypical objectives and scope of the systematic review. First, this is a systematic review of methodological strength of evidence for a given literature as opposed to a review summary of the evidence of a particular topic. As such, we do not summarize and attempt to combine the results for any of the literature. Second, rather than being a comprehensive review of every

possible aspect of what might be considered “quality,” this is a review of targeted critical design features for actionable inference for COVID-19 policy impact evaluation and methods. It is designed to be a set of broad criteria for minimal plausibility of actionable causal inference, where each of the criteria is necessary but not sufficient for strong design. Issues in other domains (data, details of the design, statistics, etc) further reduce overall actionability and quality, and thorough review in those domains is needed for any studies passing our basic minimal criteria. Third, because the scope relies on guided, but difficult and subjective assessments of methodological appropriateness, we utilize a discussion-based consensus process to arrive at consistent and replicable results, rather than a more common model with two independent reviewers with conflict resolution. The independent review serves primarily as a starting point for discussion, but is neither designed nor expected to be a strong indicator of the overall consensus ratings of the group.

Methods

Overview

This protocol and study was written and developed following the release of the review guidance written by the author team in September 2020 on which the review tool is based. The protocol for this study was pre-registered on OSF.io[18] in November 2020 following PRISMA guidelines.[19] Deviations from the original protocol are discussed in Appendix 1, and consisted largely of language clarifications and error corrections for both the inclusion criteria and review tool, an increase in the number of reviewers per fully reviewed article from two to three, and simplification of the statistical methods used to assess the data.

For this study, we ascertain minimal criteria for studies to be able to plausibly identify causal effects of policies, which is the information of greatest interest to inform policy decisions. The causal estimand is something that, if known, would definitely help policy makers decide what to do (e.g., whether to implement or discontinue a policy). The study estimates that target causal quantity with a rigorous design and appropriate data in a relevant population/sample. For shorthand, we refer to this as minimal properties of “actionable” evidence.

This systematic review of the strength of evidence took place in three phases: search, screening, and full review.

Eligibility criteria

The following eligibility criteria were used to determine the papers to include:

- The primary topic of the article must be evaluating one or more individual COVID-19 or SARS-CoV-2 policies on direct COVID-19 or SARS-CoV-2 outcomes
 - The primary exposure(s) must be a policy, defined as a government-issued order at any government level to address a directly COVID-19-related outcome (e.g., mask requirements, travel restrictions, etc).

- Direct COVID-19 or SARS-CoV-2 outcomes are those that are specific to disease and health outcomes may include cases detected, mortality, number of tests taken, test positivity rates, Rt, etc.
- This may NOT include indirect impacts of COVID-19 on items that are not direct COVID-19 or SARS-CoV-2 impacts such as income, childcare, economic impacts, beliefs and attitudes, etc.
- The primary outcome being examined must be a COVID-19-specific outcome, as above.
- The study must be designed as an impact evaluation study from primary data (i.e., not primarily a predictive or simulation model or meta-analysis).
- The study must be peer reviewed, and published in a peer-reviewed journal indexed by PubMed.
- The study must have the title and abstract available via PubMed at the time of the study start date (November 26).
- The study must be written in English.

These eligibility criteria were designed to identify the literature primarily concerning the quantitative impact of one or more implemented COVID-19 policies on COVID-19 outcomes. Studies in which impact evaluation was secondary to another analysis (such as a hypothetical projection model) were eliminated because they were less relevant to our objectives and/or may not contain sufficient information for evaluation. Categories for types of policies were from the Oxford COVID-19 Government Response Tracker.[20]

Reviewer recruitment, training, and communication

Reviewers were recruited through personal contacts and postings on online media. All reviewers had experience in systematic review, quantitative causal inference, epidemiology, econometrics, public health, methods evaluation, or policy review. All reviewers participated in two meetings in which the procedures and the review tool were demonstrated. Screening reviewers participated in an additional meeting specific to the screening process. Throughout the main review process, reviewers communicated with the administrators and each other through Slack for any additional clarifications, questions, corrections, and procedures. The main administrator (NH), who was also a reviewer, was available to answer general questions and make clarifications, but did not answer questions specific to any given article.

Review phases and procedures

Search strategy

The search terms combined four Boolean-based search terms: a) COVID-19 research, b) regional government units (e.g., country, state, county, and specific country, state, or province, etc.), c) policy or policies, and d) impact or effect. The full search terms are available in Appendix 2.

Information Sources

The search was limited to published articles in peer-reviewed journals. This was largely to attempt to identify literature that was high quality, relevant, prominent, and most applicable to the review guidance. PubMed was chosen as the exclusive indexing source due to the

1
2
3 prevalence and prominence of policy impact studies in the health and medical field. Preprints
4 were excluded to limit the volume of studies to be screened and to ensure each had met the
5 standards for publication through peer review. The search was conducted on November 26,
6 2020.
7

8 9 Study Selection

10
11 Two reviewers were randomly selected to screen the title and abstract of each article for the
12 inclusion criteria. In the case of a dispute, a third randomly selected reviewer decided on
13 acceptance/rejection. Eight reviewers participated in the screening. Training consisted of a one-
14 hour instruction meeting, a review of the first 50 items on each reviewers' list of assigned
15 articles, and a brief asynchronous online discussion before conducting the full review.
16
17

18 Full article review

19
20 The full article review consisted of two sub-phases: the independent primary review phase, and
21 a group consensus phase. The independent review phase was designed primarily for the
22 purpose of supporting and facilitating discussion in the consensus discussion, rather than as
23 high stakes definitive review data on its own. The consensus process was considered the
24 primary way in which review data would be generated, rather than synthesis from the
25 independent reviews. A flow diagram of the review process is available in Appendix 3
26
27

28
29 Each article was randomly assigned to three of the 23 reviewers in our review pool. Each
30 reviewer independently reviewed each article on their list, first for whether the study met the
31 eligibility criteria, then responding to methods identification and guided strength of evidence
32 questions using the review tool, as described below. Reviewers were able to recuse themselves
33 for any reason, in which case another reviewer was randomly selected. Once all three reviewers
34 had reviewed a given article, all articles that weren't unanimously determined to not meet the
35 inclusion criteria underwent a consensus process.
36
37

38
39 During the consensus round, the three reviewers were given all three primary reviews for
40 reference, and were tasked with generating a consensus opinion among the group. One
41 randomly selected reviewer was tasked to act as the arbitrator. The arbitrator's primary task was
42 facilitating discussion and for moving the group toward establishing a consensus that
43 represented the collective subjective assessments of the group. If consensus could not be
44 reached, a fourth randomly selected reviewer was brought into the discussion to help resolve
45 disputes.
46
47

48 Review tool for data collection

49
50 This review tool and data collection process was an operationalized and lightly adapted version
51 of the COVID-19 health policy impact evaluation review guidance literature, written by the lead
52 authors of this study and released in September 2020 as a pre-print.[21] The main adaptation
53 was removing references to the COVID-19 literature. All reviewers were instructed to read and
54 refer to this guidance document to guide their assessments. The full guidance manuscript
55
56
57
58
59

1
2
3 contains additional explanation and rationale for all parts of this review and the tool, and is
4 available both in the adapted form as was provided to the reviewers in a supplementary file
5 “CHSPER review guidance refs removed.pdf” and in an updated version in Haber et al.,
6 2021.[17] The full review tool is attached as supplementary file “review tool final.pdf”.
7
8

9 The review tool consisted of two main parts: methods design categorization and full review. The
10 review tool and guidance categorizes policy causal inference designs based on the structure of
11 their assumed counterfactual. This is assessed through identifying the data structure and
12 comparison(s) being made. There are two main items for this determination: the number of pre-
13 period time points (if any) used to assess pre-policy outcome trends, and whether or not policy
14 regions were compared with non-policy regions. These, and other supporting questions, broadly
15 allowed categorization of methods into cross-sectional, pre/post, interrupted time series (ITS),
16 difference-in-differences (DiD), comparative interrupted time-series (CITS), (randomized) trials,
17 or other. Given that most papers have several analyses, reviewers were asked to focus
18 exclusively on the impact evaluation analysis that was used as the primary support for the main
19 conclusion of the article.
20
21
22

23 Studies categorized as cross-sectional, pre/post, randomized controlled trial designs, and other
24 were included in our sample, but set aside for no further review for the purposes of this
25 research. Cross-sectional and pre/post studies are not considered sufficient to yield well-
26 identified causal inference in the specific context of COVID-19 policy impact evaluation, as
27 explained in the policy impact evaluation guidance documentation. Cross-sectional and pre-post
28 designs were considered inappropriate for policy causal inference for COVID-19 due largely to
29 inability to account for a large number of potential issues, including confounding, epidemic
30 trends, and selection biases. Randomized controlled trials were assumed to broadly meet key
31 design checks. Studies categorized as “other” received no further review, as the review
32 guidance would be unable to assess them. Additional justification and explanation for this
33 decision is available in the review guidance.
34
35
36
37

38 For the methods receiving full review (ITS, DiD, and CITS), reviewers were asked to identify
39 potential issues and give a category-specific rating. The specific study designs triggered sub-
40 questions and/or slightly altered the language of the questions being asked, but all three of the
41 methods design categories shared these four key questions:
42
43

- 44 ● Graphical presentation: “Does the analysis provide graphical representation of the
45 outcome over time?”
 - 46 ○ Graphical presentation refers to how the authors present the data underlying
47 their impact evaluation method. This is a critical criteria for assessing the
48 potential validity of the assumed model. The key questions here are whether any
49 chart shows the outcome over time and the assumed models of the
50 counterfactuals. To meet a high degree of confidence in this category, graphical
51 displays must show the outcome and connect to the counterfactual construction
52 method.
53
54
55
56
57
58
59
60

- 1
- 2
- 3 ● Functional form: “Is the functional form of the model used for the trend in counterfactual
- 4 infectious disease outcomes (e.g., linear, non-parametric, exponential, logarithmic, etc.)
- 5 well-justified and appropriate?”
- 6 ○ Functional form refers to the statistical functional form of the trend in
- 7 counterfactual infectious disease outcomes (i.e. the assumptions used to
- 8 construct counterfactual outcomes). This may be a linear function, non-
- 9 parametric, exponential or logarithmic function, infectious disease model
- 10 projection, or any other functional form. The key criteria here are whether this is
- 11 discussed and justified in the manuscript, and if so, is it a plausibly appropriate
- 12 choice given infectious disease outcomes.
- 13
- 14
- 15 ● Timing of policy impact: “Is the date or time threshold set to the appropriate date or time
- 16 (e.g., is there lag between the intervention and outcome)?”
- 17 ○ Timing of policy impact refers to assumptions about when we would expect to
- 18 see an impact from the policy vis-a-vis the timing of the policy introduction. This
- 19 would typically be modelled with leads and lags. The impact of policy can occur
- 20 before enactment (e.g., in cases where behavior change after policy is
- 21 announced, but before it takes place in anticipation) or long after the policy is
- 22 enacted (e.g., in cases where it takes time to ramp up policy implementation or
- 23 impacts). The key criteria here are whether this is discussed and justified in the
- 24 manuscript, and if so, whether it is a plausibly appropriate choice given the policy
- 25 and outcome.
- 26
- 27
- 28 ● Concurrent changes: “Is this policy the only uncontrolled or unadjusted-for way in which
- 29 the outcome could have changed during the measurement period [differently for policy
- 30 and non-policy regions]?”
- 31 ○ Concurrent changes refers to the presence of uncontrolled other events and
- 32 changes that may influence outcomes at the same time as the policy would
- 33 impact outcomes. In order to assess the impact of one policy or set of policies,
- 34 the impact of all other forces that differentially impact the outcome must either be
- 35 negligible or controlled for. The key criteria here are whether it is likely that there
- 36 are substantial other uncontrolled forces (e.g. policies, behavioral changes, etc)
- 37 which may be differentially impacting outcomes at the same time as the policy of
- 38 interest.
- 39
- 40
- 41
- 42

43 For each of the four key questions, reviewers were given the option to select “No,” “Mostly no,”
44 “Mostly yes,” and “Yes” with justification text requested for all answers other than “Yes.” Each
45 question had additional prompts as guidance, and with much more detail provided in the full
46 guidance document. Ratings are, by design, subjective assessments of the category according
47 to the guidance. We do not use numerical scoring, for similar reasons as Cochrane suggests
48 that the algorithms for summary judgements for the RoB2 tool are merely “proposed”
49 assessments, which reviewers should change as they believe appropriate.[22] It is entirely
50 plausible, for example, for a study to meet all but one criteria but for the one remaining to be
51 sufficiently violated that the entire collective category is compromised. Alternatively, there could
52 be many minor violations of all of the criteria, but that they were collectively not sufficiently
53
54
55
56
57
58
59

1
2
3 problematic to impact overall ratings. Further, reviewers were also tasked with considering room
4 for doubt in cases where answers to these questions were unclear.
5

6
7 The criteria were designed to establish minimal plausibility of actionable evidence, rather than
8 certification of high quality. Graphical representation is included here primarily as a key way to
9 assess the plausibility and justification of key model assumptions, rather than being necessary
10 for validity by itself. For example, rather than having the “right” functional form or lag structure,
11 the review guidance asks whether the functional form and lags is discussed at all and (if
12 discussed) reasonable.
13

14
15 These four questions were selected and designed being critical to evaluating strength of study
16 design for policy impact evaluation in general, direct relevance for COVID-19 policy, feasibility
17 for use in guided review. These questions are designed as minimal and key criteria for plausibly
18 actionable impact evaluation design for COVID-19 policy impact evaluation, rather than as a
19 comprehensive tool assessing overall quality. Thorough review of data quality, statistical
20 validity, and other issues are also critical points of potential weakness in study designs, and
21 would be needed in addition to these criteria, if these key design criteria are met. A thorough
22 justification and explanation of how and why these questions were selected is available in the
23 provided guidance document and in Haber et al., 2020.[17]
24
25

26
27 Finally, reviewers were asked a summary question:
28

- 29
30 ● Overall: “Do you believe that the design is appropriate for identifying the policy impact(s)
31 of interest?”
32

33 Reviewers were asked to consider the scale of this question to be both independent/not relative
34 to any other papers, and that any one substantial issue with the study design could render it a
35 “No” or “Mostly no.” Reviewers were asked to follow the guidance and their previous answers,
36 allowing for their own weighting of how important each issue was to the final result. A study
37 could be excellent on all dimensions except for one, and that one dimension could render it
38 inappropriate for causal inference. As such, in addition to the overall rating question, we also
39 generated a “weakest link” metric for overall assessment, representing the lowest rating among
40 the four key questions (graphical representation, functional form, timing of policy impact, and
41 concurrent changes). A “mostly yes” or “yes” is considered a passing rating, indicating that the
42 study was not found to be inappropriate on the specific dimension of interest.
43
44
45

46
47 A “yes” rating does not necessarily indicate that the study is strongly designed, conducted, or is
48 actionable; it only means that it passes a series of key design checks for policy impact
49 evaluation and should be considered for further evaluation. The papers may contain any
50 number of other issues that were not reviewed (e.g., statistical issues, inappropriate
51 comparisons, generalizability, etc.). As such, this should only be considered an initial
52 assessment of plausibility that the study is well-designed, rather than confirmation that it is
53 appropriate and applicable.
54
55
56
57
58
59

Heterogeneity

Inter-rater reliability (IRR) was assessed using Krippendorff's alpha.[23,24] Rather than more typical uses intended as an examination of the "validity" of ratings, the IRR statistic in this case is being used as a heuristic indicator of heterogeneity between reviewers during the independent phase, where heterogeneity is both expected and not necessarily undesirable. As a second examination of reviewer heterogeneity, we also show the distribution of category differences between primary reviewers within a study (e.g. if primary reviewers rated "Yes," "Mostly no," and "Mostly yes" there are two pairs of answers that were one category different, and one pair that was two categories different).

Statistical analysis

Statistics provided are nearly exclusively counts and percentages of the final dataset. Analyses and graphics were performed in R.[25] Krippendorff's alpha was calculated using the IRR package.[26] Relative risks were estimated using the epitools package.[27]

Citation counts for accepted articles were obtained through Google Scholar[28] on January 11, 2021. Journal impact factors were obtained from the 2019 Journal Citation Reports.[29]

Data sharing

Data, code, the review tool, and the review guidance are stored and available at the OSF.io repository for this study[30] here: <https://osf.io/9xmke/files/>. The dataset includes full results from the search and screening and all review tool responses from reviewers during the full review phase.

Patient and Public Involvement Statement

Patients or public stakeholders were not consulted in the design or conduct of this systematic evaluation.

Results

Search and screening

<Figure 1: PRISMA diagram of systematic review process>

After search and screening of titles and abstracts, 102 articles were identified as likely or potentially meeting our inclusion criteria (Figure 1). Of those 102 articles, 36 studies met inclusion after independent review and deliberation in the consensus process. The most common reasons for rejection at this stage were that the study did not measure the quantitative direct impact of specific policies and/or that such an impact was not the main purpose of the

study. Many of these studies implied that they measured policy impact in the abstract or introduction, but instead measured correlations with secondary outcomes (e.g., the effect of movement reductions, which are influenced by policy) and/or performed cursory policy impact evaluation secondary to projection modelling efforts.

Descriptive statistics

<Figure 2: Descriptive sample statistics (n=36)>

Publication information from our sample is shown in Figure 2. The articles in our sample were generally published in journals with high impact factors (median impact factor: 3.6, 25th percentile: 2.3, 75th percentile: 5.3 IQR: 3.0) and have already been cited in the academic literature (median citation count: 5.0, 25th percentile: 2.0, 75th percentile: 26.8, IQR 24.8, on 1/11/21). The most commonly evaluated policy type was stay at home requirements (64% n=23/36). Reviewers noted that many articles referenced “lockdowns,” but did not define the specific policies to which this referred. Reviewers specified mask mandates for 3 of the studies, and noted either a combination of many interventions or unspecified specific policies in 7 cases.

Reviewers most commonly selected interrupted time-series (39% n=14/36) as the methods design, followed by difference-in-differences (9% n=9/36) and pre-post (8% n=8/36). There were no randomized controlled trials of COVID-19 health policies identified (0% n=0/36), nor were any studies identified that reviewers could not categorize based on the review guidance (0% n=0/36).

Table 1: Summary of articles reviewed and reviewer ratings for key and overall questions

Category ratings order		Legend for color coded ratings						
Graphical presentation Functional form	Timing of policy impact Concurrent changes	N/A	Unclear	No*	No**	Mostly no	Mostly yes	Yes
method determined to me inappropriate by: * guidance (cross sectional or pre/post) or ** reviewer consensus								
Citation	Title	Journal	Publication date	Methods design	Category ratings	Overall rating		
Cobb and Seale, 2020[31]	Examining the effect of social distancing on the compound growth rate of COVID-19 at the county level (United States) using statistical analyses and a random forest machine learning model.	Public Health	4/28/2020	Pre/post	Yes	Yes		
Lyu and Wehby, 2020a[32]	Comparison of Estimated Rates of Coronavirus Disease 2019 (COVID-19) in Border Counties in Iowa Without a Stay-at-Home Order and Border Counties in Illinois With a Stay-at-Home Order.	JAMA Network Open	5/1/2020	Difference-in-differences	Mostly yes	Mostly yes		
Tam et al., 2020[33]	Effect of mitigation measures on the spreading of COVID-19 in hard-hit states in the U.S.	PloS One	5/1/2020	Interrupted time-series	Mostly yes	Mostly no		
Courtemanche et al., 2020[34]	Strong Social Distancing Measures In The United States Reduced The COVID-19 Growth Rate.	Health Affairs	5/14/2020	Difference-in-differences	Mostly yes	Mostly yes		
Crokidakis, 2020[35]	COVID-19 spreading in Rio de Janeiro, Brazil: Do the policies of social isolation really work?	Chaos, Solitons, and Fractals	5/23/2020	Interrupted time-series	Mostly yes	Mostly yes		
Hyafil and Morifa, 2020[36]	Analysis of the impact of lockdown on the reproduction number of the SARS-Cov-2 in Spain.	Gaceta Anitaria	5/23/2020	Pre/post	Mostly yes	Mostly no		
Castillo, et al., 2020[37]	The effect of state-level stay-at-home orders on COVID-19 infection rates.	American Journal of infection control	5/24/2020	Pre/post	Mostly yes	Mostly no		
Alfano and Ercolano, 2020[38]	The Efficacy of Lockdown Against COVID-19: A Cross-Country Panel Analysis.	Applied Health Economics and Health Policy	6/3/2020	Difference-in-differences	Mostly yes	Mostly no		
Lyu and Wehby, 2020b[39]	Community Use Of Face Masks And COVID-19: Evidence From A Natural Experiment Of State Mandates In The US.	Health Affairs	6/16/2020	Difference-in-differences	Mostly yes	Mostly yes		

1								
2								
3	Zhang, et al., 2020[40]	Identifying airborne transmission as the dominant route for the spread of COVID-19.	PNAS	6/30/2020	Interrupted time-series			
4								
5	Xu et al., 2020[41]	Associations of Stay-at-Home Order and Face-Masking Recommendation with Trends in Daily New Cases and Deaths of Laboratory-Confirmed COVID-19 in the United States.	Exploratory research and hypothesis in medicine	7/8/2020	Interrupted time-series			
6								
7	Lyu and Wehby, 2020[42]	Shelter-In-Place Orders Reduced COVID-19 Mortality And Reduced The Rate Of Growth In Hospitalizations.	Health Affairs	7/9/2020	Difference-in-differences			
8								
9	Wagner, et al., 2020[43]	Social distancing merely stabilized COVID-19 in the US.	Stat (International Statistical Institute)	7/13/2020	Interrupted time-series			
10								
11	Di Bari et al., 2020[44]	Extensive Testing May Reduce COVID-19 Mortality: A Lesson From Northern Italy.	Frontiers in Medicine	7/14/2020	Comparative interrupted time-series			
12								
13	Islam et al., 2020[45]	Physical distancing interventions and incidence of coronavirus disease 2019: natural experiment in 149 countries.	BMJ (Clinical research ed.)	7/15/2020	Interrupted time-series			
14								
15	Wong et al., 2020[46]	Impact of National Containment Measures on Decelerating the Increase in Daily New Cases of COVID-19 in 54 Countries and 4 Epicenters of the Pandemic: Comparative Observational Study.	Journal of Medical Internet Research	7/22/2020	Pre/post			
16								
17	Liang et al., 2020[47]	Effects of policies and containment measures on control of COVID-19 epidemic in Chongqing.	World Journal of Clinical Cases	7/26/2020	Pre/post			
18								
19	Banerjee and Nayak, 2020[48]	U.S. county level analysis to determine if social distancing slowed the spread of COVID-19.	Pan American Journal of Public Health	7/31/2020	Difference-in-differences			
20								
21	Dave et al., 2020a[49]	When Do Shelter-in-Place Orders Fight COVID-19 Best? Policy Heterogeneity Across States and Adoption Time.	Economic inquiry	8/3/2020	Difference-in-differences			
22								
23	Hsiang et al., 2020[50]	The effect of large-scale anti-contagion policies on the COVID-19 pandemic.	Nature	8/22/2020	Interrupted time-series			
24								
25	Lim et al., 2020[51]	Revealing regional disparities in the transmission potential of SARS-CoV-2 from interventions in Southeast Asia.	Proceedings. Biological sciences	8/26/2020	Interrupted time-series			
26								
27	Arshed et al., 2020[52]	Empirical assessment of government policies and flattening of the COVID19 curve.	Journal of Public Affairs	8/27/2020	Cross-sectional analysis			
28								
29	Wang et al., 2020[53]	Fangcang shelter hospitals are a One Health approach for responding to the COVID-19 outbreak in Wuhan, China.	One Health	8/29/2020	Interrupted time-series			
30								
31	Kang et al., 2020[54]	The Effects of Border Shutdowns on the Spread of COVID-19.	Journal of Preventive Medicine and Public Health	8/30/2020	Comparative interrupted time-series			
32								
33	Auger et al., 2020[55]	Association Between Statewide School Closure and COVID-19 Incidence and Mortality in the US.	JAMA	9/1/2020	Interrupted time-series			
34								
35	Santamaria et al., 2020[56]	COVID-19 effective reproduction number dropped during Spain's nationwide dropdown, then spiked at lower-incidence regions.	The Science of the Total Environment	9/9/2020	Interrupted time-series			
36								
37	Bennett, 2020[57]	All things equal? Heterogeneity in policy effectiveness against COVID-19 spread in Chile.	World Development	9/24/2020	Comparative interrupted time-series			
38								
39	Yang et al., 2020[58]	Lessons Learnt from China: National Multidisciplinary Healthcare Assistance.	Risk Management and Healthcare Policy	9/30/2020	Difference-in-differences			
40								
41	Padalabalanarayanan et al., 2020[59]	Association of State Stay-at-Home Orders and State-Level African American Population With COVID-19 Case Rates.	JAMA Network Open	10/1/2020	Comparative interrupted time-series			
42								
43	Edelstein et al., 2020[60]	SARS-CoV-2 infection in London, England: changes to community point prevalence around lockdown time, March-May 2020.	Journal of Epidemiology and Community Health	10/1/2020	Pre/post			
44								
45	Tsai et al., 2020[61]	COVID-19 transmission in the U.S. before vs. after relaxation of statewide social distancing measures.	Clinical Infectious Diseases	10/3/2020	Interrupted time-series			
46								
47	Singh et al., 2020[62]	Public health interventions slowed but did not halt the spread of COVID-19 in India.	Transboundary and Emerging Diseases	10/4/2020	Pre/post			
48								
49	Gallaway et al., 2020[63]	Trends in COVID-19 Incidence After Implementation of Mitigation Measures - Arizona, January 22-August 7, 2020.	Morbidity and Mortality Weekly Report	10/9/2020	Pre/post			
50								
51	Castex et al., 2020[64]	COVID-19: The impact of social distancing policies, cross-country analysis.	Economics of Disasters and Climate Change	10/15/2020	Interrupted time-series			
52								
53	Silva et al., 2020[65]	The effect of lockdown on the COVID-19 epidemic in Brazil: evidence from an interrupted time series design.	Cadernos de Saude Publica	10/19/2020	Interrupted time-series			
54								
55	Dave et al., 2020b[66]	Were Urban Cowboys Enough to Control COVID-19? Local Shelter-in-Place Orders and Coronavirus Case Growth.	Journal of Urban Economics	11/6/2020	Difference-in-differences			
56								
57								
58								
59								
60								

The identified articles and selected review results are summarized in Table 1.

Strength of methods assessment

<Figure 3: Main consensus results summary for key and overall questions>

Graphical representation of the outcome over time was relatively well-rated in our sample, with 74% (n=20/27) studies being given a “mostly yes” or “yes” rating for appropriateness. Reasons cited for non-“yes” ratings included a lack of graphical representation of the data, alternative scales used, and not showing the dates of policy implementation.

Functional form issues appear to have presented a major issue in these studies, with only 19% receiving a “mostly yes” or “yes” rating, 78% (n=21/27) receiving a “no” rating, and 4% (n=1/27) “unclear.” There were two common themes in this category: studies generally using scales that were broadly considered inappropriate for infectious disease outcomes (e.g., linear counts), and/or studies lacking stated justification for the scale used. Reviewers also noted disconnects between clear curvature in the outcomes in the graphical representations and the analysis models and outcome scales used (e.g., linear). In one case, reviewers could not identify the functional form actually used in analysis.

Reviewers broadly found that these studies dealt with timing of policy impact (e.g., lags between policy implementation and expected impact) relatively well, with 70% (n=19/27) rated “yes” or “mostly yes.” Reasons for non-“yes” responses included not adjusting for lags and a lack of justification for the specific lags used.

Concurrent changes were found to be a major issue in these studies, with only 11% (n=3/27) studies receiving passing ratings (“yes” or “mostly yes”) with regard to uncontrolled concurrent changes to the outcomes. Reviewers nearly ubiquitously noted that the articles failed to account for the impact of other policies that could have impacted COVID-19 outcomes concurrent with the policies of interest. Other issues cited were largely related to non-policy-induced behavioral and societal changes.

When reviewers were asked if sensitivity analyses had been performed on key assumptions and parameters, about half (56% n=15/27) answered “mostly yes” or “yes.” The most common reason for non-“yes” ratings was that, while sensitivity analyses were performed, they did not address the most substantial assumptions and issues.

Overall, reviewers rated only four studies (11%, n=4/36,) as being plausibly appropriate (“mostly yes” or “yes”) for identifying the impact of specific policies on COVID-19 outcomes, as shown in Figure 3. 25% (n=9/36) were automatically categorized as being inappropriate due to being either cross-sectional or pre/post in design, 33% (n=12/36) of studies were given a “no” rating for appropriateness, 31% “mostly no” (n=11/36), 8% “mostly yes” (n=3/36), and 3% “yes” (n=1/36). The most common reason cited for non-“yes” overall ratings was failure to account for concurrent changes (particularly policy and societal changes).

<Figure 4: Comparison of independent reviews, weakest link, and direct consensus review>

1
2
3 As shown in Figure 4, the consensus overall proportion passing (“mostly yes” or “yes”) was a
4 quarter of what it was from the initial independent reviews. 45% (n=34/75) of studies were rated
5 as “yes” or “mostly yes” in the initial independent review, as compared to 11% (n=4/36) in the
6 consensus round (RR 0.25, 95%CI 0.09:0.64). The issues identified and discussed in
7 combination during consensus discussions, as well as additional clarity on the review process,
8 resulted in reduced overall confidence in the findings. Increased clarity on the review guidance
9 with experience and time may also have reduced these ratings further.
10
11

12
13 The large majority of studies had at least one “no” or “unclear” rating in one of the four
14 categories (74% n=20/27), with only one study whose lowest rating was a “mostly yes,” no
15 studies rated “yes” in all four categories. Only one study was found to pass design criteria in all
16 four key questions categories, as shown in the “weakest link” column in Figure 4.
17
18

19 Review process assessment

20
21 During independent review, all three reviewers independently came to the same conclusions on
22 the main methods design category for 33% (n=12/36) articles, two out of the three reviewers
23 agreed for 44% (n=16/36) articles, and none of the reviewers agreed in 22% (n=8/36) cases.
24 One major contributor to these discrepancies were the 31% (n=11/36) cases where one or more
25 reviewers marked the study as not meeting eligibility criteria, 64% (n=7/11) of which the other
26 two reviewers agreed on the methods design category.
27
28

29
30 Reviewers’ initial independent reviews were heterogeneous for key rating questions. For the
31 overall scores, Krippendorff’s alpha was only 0.16 due to widely varying opinions between
32 raters. The four key categorical questions had slightly better inter-rater reliability than the overall
33 question, with Krippendorff’s alphas of 0.59 for graphical representation, 0.34 for functional form,
34 0.44 for timing of policy impact, and 0.15 for concurrent changes, respectively. For the main
35 summary rating, primary reviewers within each study agreed in 26% of cases (n=16), were one
36 category different in 45% (n=46), two categories different in 19% (n=12), and three categories
37 (i.e. the maximum distance, “Yes” vs “No”) in 10% of cases (n=6).
38
39

40
41 The consensus rating for overall strength was equal to the lowest rating among the independent
42 reviews in 78% (n=21/27) of cases, and only one higher than the lowest in the remaining 22%
43 (n=6/27). This strongly suggests that the multiple reviewer review, discussion, and consensus
44 process more thoroughly identifies issues than independent review alone. There were two
45 cases for which reviewers requested an additional fourth reviewer to help resolve standing
46 issues for which the reviewers felt they were unable to come to consensus.
47
48

49
50 The most consistent point of feedback from reviewers was the value of having a three reviewer
51 team with whom to discuss and deliberate, rather than two as initially planned. This was
52 reported to help catch a larger number of issues and clarify both the papers and the
53 interpretation of the review tool questions. Reviewers also expressed that one of the most
54 difficult parts of this process was assessing the inclusion criteria, some of the implications of
55 which are discussed below.
56
57
58
59
60

Discussion

This systematic review of evidence strength found that only four (or only one by a stricter standard) of the 36 identified published and peer-reviewed health policy impact evaluation studies passed a set of key checks for identifying the causal impact of policies on COVID-19 outcomes. Because this systematic review examined a limited set of key study design features and did not address more detailed aspects of study design, statistical issues, generalizability, and any number of other issues, this result may be considered an upper bound on the overall strength of evidence within this sample. Two major problems are nearly ubiquitous throughout this literature: failure to isolate the impact of the policy(s) of interest from other changes that were occurring contemporaneously, and failure to appropriately address the functional form of infectious disease outcomes in a population setting. While policy decisions are being made on the backs of high impact-factor papers, we find that the citation-based metrics do not correspond to “quality” research as used by Yin et al., 2021.[67] Similar to other areas in the COVID-19 literature,[68] we found the current literature directly evaluating the impact of COVID-19 policies largely fails to meet key design criteria for actionable inference to inform policy decisions.

The framework for the review tool is based on the requirements and assumptions built into policy evaluation methods. Quasi-experimental methods rely critically on the scenarios in which the data are generated. These assumptions and the circumstances in which they are plausible are well-documented and understood,[2,4–6,17,69] including one paper discussing application of difference-in-differences methods specifically for COVID-19 health policy, released in May 2020.[5] While “no uncontrolled concurrent changes” is a difficult bar to clear, that bar is fundamental to inference using these methods.

The circumstances of isolating the impact of policies in COVID-19 - including large numbers of policies, infectious disease dynamics, and massive changes to social behaviors - make those already difficult fundamental assumptions broadly much less likely to be met. Some of the studies in our sample were nearly the best feasible studies that could be done given the circumstances, but the best that can be done often yields little actionable inference. The relative paucity of strong studies does not in any way imply a lack of impact of those policies; only that we lack the circumstances to have evaluated their effects.

Because the studies estimating the harms of policies share the same fundamental circumstances, the evidence of COVID-19 policy harms is likely to be of similarly poor strength. Identifying the effects of many of these policies, particularly for the spring of 2020, is likely to be unknown and perhaps unknowable. However, there remains additional opportunities with more favorable circumstances, such as measuring overall impact of NPIs as bundles, rather than individual policies. Similarly, studies estimating the impact of re-opening policies or policy cancellation are likely to have fewer concurrent changes to address.

The review process itself demonstrates how guided and targeted peer review can efficiently evaluate studies in ways that the traditional peer review systems do not. The studies in our

1
2
3 sample had passed the full peer review process, were published in largely high-profile journals,
4 and are highly cited, but contained substantial flaws that rendered their inference utility
5 questionable. The relatively small number of studies included, as compared to the size of the
6 literature concerning itself with COVID-19 policy, may suggest that there was relative restraint
7 from journal editors and reviewers for publishing these types of studies. The large number of
8 models, but relatively small number of primary evaluation analyses is consistent with other
9 areas of COVID-19.[70,71] At minimum, the flaws and limitations in their inference could have
10 been communicated at the time of publication, when they are needed most. In other cases, it is
11 plausible that many of these studies would not have been published had a more thorough or
12 more targeted methodological review been performed.
13
14
15

16 This systematic review of evidence strength has limitations. The tool itself was limited to a very
17 narrow - albeit critical - set of items. Low ratings in our study should not be interpreted as being
18 overall poor studies, as they may make other contributions to the literature that we did not
19 evaluate. While the guidance and tool provided a well-structured framework and our reviewer
20 pool was well-qualified, strength of evidence review is inherently subjective. It is plausible and
21 likely that other sets of reviewers would come to different conclusions for each study, but
22 unlikely that the overall conclusions of our assessment would change substantially. However,
23 the consensus process was designed with subjectivity in mind, and demonstrates the value of
24 consensus processes for overcoming hurdles with subjective and difficult decisions.
25
26
27

28 While subjective assessments are inherently subject to the technical expertise, experiences,
29 and opinions of reviewers, we argue they are both appropriate and necessary to reliably assess
30 strength of evidence based on theoretical methodological issues. With the exception of the
31 graphical assessment, proper assessment of the core methodological issues requires that
32 reviewers are able to weigh the evidence as they see fit. Much like standard institutional peer
33 review, reviewers independently had highly heterogeneous opinions, attributable to differences
34 in opinion or training, misunderstandings/learning about the review tool and process, and
35 expected reliance on the consensus process. Unlike traditional peer review, there was subject-
36 matter-specific guidance and a process to consolidate and discuss those heterogenous initial
37 opinions. The reduction in ratings from the initial highly heterogeneous ratings to a lower
38 heterogeneity in ratings indicates that reviewers had initially identified issues differently, but that
39 the discussion and consensus process helped elucidate the extent of the different issues that
40 each reviewer detected and brought to discussion. This also reflects reviewer learning over
41 time, where reviewers were better able to identify issues at the consensus phase than earlier. It
42 is plausible that stronger opinions had more weight, but we expect that this was largely
43 mitigated by the random assignment of the arbitrator, and reviewer experiences did not indicate
44 this as an issue.
45
46
47
48
49

50 Most importantly, this review does not cover all policy inference in the scientific literature. One
51 large literature from which there may be COVID-19 policy evaluation otherwise meeting our
52 inclusion criteria are pre-prints. Many pre-prints would likely fare well in our review process.
53 Higher strength papers often require more time for review and publication, and many high
54 quality papers may be in the publication pipeline now. Second, this review excluded studies that
55
56
57
58
59

1
2
3 had a quantitative impact evaluation as a secondary part of the study (e.g., to estimate
4 parameters for microsimulation or disease modeling). Third, the review does not include policy
5 inference studies that do not measure the impact of a specific policy. For instance, there are
6 studies that estimate the impact of reduced mobility on COVID-19 outcomes but do not attribute
7 the reduced mobility to any specific policy change. A considerable number of studies that
8 present analyses of COVID-19 outcomes to inform policy are excluded because they do not
9 present a quantitative estimate of specific policies' treatment effects. Importantly, this study was
10 designed to assess a minimal set of criteria critical to the design of impact evaluation studies of
11 COVID-19 policies. Studies found meeting these criteria would require further and more
12 comprehensive review for assessing overall quality and actionability. Unfortunately, exceedingly
13 few studies we reviewed, taken largely from the high-profile literature, were found to meet these
14 minimal criteria.
15
16
17
18

19 While COVID-19 policy is one of the most important problems of our time, the circumstances
20 under which those policies were enacted severely hamper our ability to study and understand
21 their effects. Claimed conclusions are only as valuable as the methods by which they are
22 produced. Replicable, rigorous, intense, and methodologically guided review is needed to both
23 communicate our limitations and make more actionable inference. Weak, unreliable, and
24 overconfident evidence leads to poor decisions and undermines trust in science.[15,72] In the
25 case of COVID-19 health policy, a frank appraisal of the strength of the studies on which
26 policies are based is needed, alongside the understanding that we often must make decisions
27 when strong evidence is not feasible.[73]
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Ethical approval statement

Not applicable

Works cited

- 1 Fischhoff B. Making Decisions in a COVID-19 World. *JAMA* 2020;**324**:139. doi:10.1001/jama.2020.10178
- 2 COVID-19 Statistics, Policy modeling, and Epidemiology Collective. Defining high-value information for COVID-19 decision-making. *Health Policy* 2020. doi:10.1101/2020.04.06.20052506
- 3 Hernán MA, Robins JM. *Causal Inference: What If*. Boca Raton: : Chapman & Hall/CRC
- 4 Angrist J, Pischke J-S. *Mostly Harmless Econometrics: An Empiricist's Companion*. 1st ed. Princeton University Press 2009. <https://EconPapers.repec.org/RePEc:pup:pbooks:8769>
- 5 Goodman-Bacon A, Marcus J. Using Difference-in-Differences to Identify Causal Effects of COVID-19 Policies. *SSRN Journal* Published Online First: 2020. doi:10.2139/ssrn.3603970
- 6 Bärnighausen T, Oldenburg C, Tugwell P, *et al*. Quasi-experimental study designs series—paper 7: assessing the assumptions. *Journal of Clinical Epidemiology* 2017;**89**:53–66. doi:10.1016/j.jclinepi.2017.02.017
- 7 Haushofer J, Metcalf CJE. Which interventions work best in a pandemic? *Science* 2020;**368**:1063–5. doi:10.1126/science.abb6144
- 8 Else H. How a torrent of COVID science changed research publishing — in seven charts. *Nature* 2020;**588**:553–553. doi:10.1038/d41586-020-03564-y
- 9 Palayew A, Norgaard O, Safreed-Harmon K, *et al*. Pandemic publishing poses a new COVID-19 challenge. *Nat Hum Behav* 2020;**4**:666–9. doi:10.1038/s41562-020-0911-0
- 10 Bagdasarian N, Cross GB, Fisher D. Rapid publications risk the integrity of science in the era of COVID-19. *BMC Med* 2020;**18**:192. doi:10.1186/s12916-020-01650-6
- 11 Yeo-Teh NSL, Tang BL. An alarming retraction rate for scientific publications on Coronavirus Disease 2019 (COVID-19). *Accountability in Research* 2020;**0**:1–7. doi:10.1080/08989621.2020.1782203
- 12 Abritis A, Marcus A, Oransky I. An “alarming” and “exceptionally high” rate of COVID-19 retractions? *Accountability in Research* 2021;**28**:58–9. doi:10.1080/08989621.2020.1793675
- 13 Zdravkovic M, Berger-Estilita J, Zdravkovic B, *et al*. Scientific quality of COVID-19 and SARS CoV-2 publications in the highest impact medical journals during the early phase of the pandemic: A case control study. *PLoS ONE* 2020;**15**:e0241826. doi:10.1371/journal.pone.0241826
- 14 Elgendy IY, Nimri N, Barakat AF, *et al*. A systematic bias assessment of top-cited full-length original clinical investigations related to COVID-19. *European Journal of Internal Medicine* 2021;:S0953620521000182. doi:10.1016/j.ejim.2021.01.018
- 15 Glasziou PP, Sanders S, Hoffmann T. Waste in covid-19 research. *BMJ* 2020;**369**:m1847. doi:10.1136/bmj.m1847
- 16 Powell M, Koenecke A, Byrd JB, *et al*. A how-to guide for conducting retrospective analyses: example COVID-19 study. *Open Science Framework* 2020. doi:10.31219/osf.io/3drch
- 17 Haber NA, Clarke-Deelder E, Salomon JA, *et al*. COVID-19 Policy Impact Evaluation: A guide to common design issues. *American Journal of Epidemiology* 2021;:kwab185. doi:10.1093/aje/kwab185

- 1
2
3 18 Haber N. Systematic review of COVID-19 policy evaluation methods and design. Published
4 Online First: 26 November 2020. <https://osf.io/7nbk6> (accessed 15 Jan 2021).
- 5 19 PRISMA. <http://www.prisma-statement.org/PRISMAStatement/> (accessed 15 Jan 2021).
- 6 20 Petherick A, Kira B, Hale T, *et al.* Variation in Government Responses to COVID-19.
7 <https://www.bsg.ox.ac.uk/research/publications/variation-government-responses-covid-19>
8 (accessed 24 Nov 2020).
- 9 21 Haber NA, Clarke-Deelder E, Salomon JA, *et al.* Policy evaluation in COVID-19: A guide to
10 common design issues. *arXiv:200901940 [stat]* Published Online First: 31 December
11 2020. <http://arxiv.org/abs/2009.01940> (accessed 15 Jan 2021).
- 12 22 Chapter 8: Assessing risk of bias in a randomized trial.
13 <https://training.cochrane.org/handbook/current/chapter-08> (accessed 8 Sep 2021).
- 14 23 Krippendorff KH. *Content Analysis: An Introduction to Its Methodology*. SAGE Publications
15 1980.
- 16 24 Zhao X, Liu JS, Deng K. Assumptions behind Intercoder Reliability Indices. *Annals of the*
17 *International Communication Association* 2013;**36**:419–80.
18 doi:10.1080/23808985.2013.11679142
- 19 25 R Core Team. *R: A language and environment for statistical computing*. Vienna, Austria: R
20 Foundation for Statistical Computing 2019. <https://www.R-project.org/>
- 21 26 Gamer M, Lemon J, Fellows I, *et al.* *irr: Various Coefficients of Interrater Reliability and*
22 *Agreement*. <https://cran.r-project.org/web/packages/irr/index.html>
- 23 27 Aragon TJ, Fay MP, Wollschlaeger D, *et al.* *Epitools*. CRAN: 2017. [https://cran.r-](https://cran.r-project.org/web/packages/epitools/epitools.pdf)
24 [project.org/web/packages/epitools/epitools.pdf](https://cran.r-project.org/web/packages/epitools/epitools.pdf)
- 25 28 About Google Scholar. <https://scholar.google.com/intl/en/scholar/about.html> (accessed 15
26 Jan 2021).
- 27 29 Clarivate Analytics. Journal Citation Reports. 2019.
- 28 [dataset] 30 Haber N. Data repository for Systematic review of COVID-19 policy evaluation
29 methods and design. 2020. <https://osf.io/9xmke/files> (accessed 9 Nov 2021).
- 30 31 Cobb JS, Seale MA. Examining the effect of social distancing on the compound growth rate
31 of COVID-19 at the county level (United States) using statistical analyses and a random
32 forest machine learning model. *Public Health* 2020;**185**:27–9.
33 doi:10.1016/j.puhe.2020.04.016
- 34 32 Lyu W, Wehby GL. Comparison of Estimated Rates of Coronavirus Disease 2019 (COVID-
35 19) in Border Counties in Iowa Without a Stay-at-Home Order and Border Counties in Illinois
36 With a Stay-at-Home Order. *JAMA Netw Open* 2020;**3**:e2011102.
37 doi:10.1001/jamanetworkopen.2020.11102
- 38 33 Tam K-M, Walker N, Moreno J. Effect of mitigation measures on the spreading of COVID-19
39 in hard-hit states in the U.S. *PLoS One* 2020;**15**:e0240877.
40 doi:10.1371/journal.pone.0240877
- 41 34 Courtemanche C, Garuccio J, Le A, *et al.* Strong Social Distancing Measures In The United
42 States Reduced The COVID-19 Growth Rate: Study evaluates the impact of social
43 distancing measures on the growth rate of confirmed COVID-19 cases across the United
44 States. *Health Affairs* 2020;**39**:1237–46. doi:10.1377/hlthaff.2020.00608
- 45 35 Crokidakis N. COVID-19 spreading in Rio de Janeiro, Brazil: Do the policies of social
46 isolation really work? *Chaos Solitons Fractals* 2020;**136**:109930.
47 doi:10.1016/j.chaos.2020.109930
- 48 36 Hyafil A, Moríña D. Analysis of the impact of lockdown on the reproduction number of the
49 SARS-Cov-2 in Spain. *Gaceta Sanitaria* 2020;:S0213911120300984.
50 doi:10.1016/j.gaceta.2020.05.003
- 51 37 Castillo RC, Staguhn ED, Weston-Farber E. The effect of state-level stay-at-home orders on
52 COVID-19 infection rates. *American Journal of Infection Control* 2020;**48**:958–60.
53 doi:10.1016/j.ajic.2020.05.017
- 54
55
56
57
58
59
60

- 1
2
3 38 Alfano V, Ercolano S. The Efficacy of Lockdown Against COVID-19: A Cross-Country Panel
4 Analysis. *Appl Health Econ Health Policy* 2020;**18**:509–17. doi:10.1007/s40258-020-00596-
5 3
- 6 39 Lyu W, Wehby GL. Community Use Of Face Masks And COVID-19: Evidence From A
7 Natural Experiment Of State Mandates In The US: Study examines impact on COVID-19
8 growth rates associated with state government mandates requiring face mask use in public.
9 *Health Affairs* 2020;**39**:1419–25. doi:10.1377/hlthaff.2020.00818
- 10 40 Zhang R, Li Y, Zhang AL, *et al*. Identifying airborne transmission as the dominant route for
11 the spread of COVID-19. *Proc Natl Acad Sci USA* 2020;**117**:14857–63.
12 doi:10.1073/pnas.2009637117
- 13 41 Xu J, Hussain S, Lu G, *et al*. Associations of Stay-at-Home Order and Face-Masking
14 Recommendation with Trends in Daily New Cases and Deaths of Laboratory-Confirmed
15 COVID-19 in the United States. *Explor Res Hypothesis Med* 2020;:1–10.
16 doi:10.14218/ERHM.2020.00045
- 17 42 Lyu W, Wehby GL. Shelter-In-Place Orders Reduced COVID-19 Mortality And Reduced The
18 Rate Of Growth In Hospitalizations. *Health Aff (Millwood)* 2020;**39**:1615–23.
19 doi:10.1377/hlthaff.2020.00719
- 20 43 Wagner AB, Hill EL, Ryan SE, *et al*. Social distancing merely stabilized COVID-19 in the
21 United States. *Stat* 2020;**9**. doi:10.1002/sta4.302
- 22 44 Di Bari M, Balzi D, Carreras G, *et al*. Extensive Testing May Reduce COVID-19 Mortality: A
23 Lesson From Northern Italy. *Front Med* 2020;**7**:402. doi:10.3389/fmed.2020.00402
- 24 45 Islam N, Sharp SJ, Chowell G, *et al*. Physical distancing interventions and incidence of
25 coronavirus disease 2019: natural experiment in 149 countries. *BMJ* 2020;:m2743.
26 doi:10.1136/bmj.m2743
- 27 46 Wong LP, Alias H. Temporal changes in psychobehavioural responses during the early
28 phase of the COVID-19 pandemic in Malaysia. *J Behav Med* 2020;:1–11.
29 doi:10.1007/s10865-020-00172-z
- 30 47 Liang X-H, Tang X, Luo Y-T, *et al*. Effects of policies and containment measures on control
31 of COVID-19 epidemic in Chongqing. *WJCC* 2020;**8**:2959–76.
32 doi:10.12998/wjcc.v8.i14.2959
- 33 48 Banerjee T, Nayak A. U.S. county level analysis to determine If social distancing slowed the
34 spread of COVID-19. *Revista Panamericana de Salud Pública* 2020;**44**:1.
35 doi:10.26633/RPSP.2020.90
- 36 49 Dave D, Friedson AI, Matsuzawa K, *et al*. When Do Shelter-in-Place Orders Fight COVID-
37 19 Best? Policy Heterogeneity Across States and Adoption Time. *Econ Inq* Published
38 Online First: 3 August 2020. doi:10.1111/ecin.12944
- 39 50 Hsiang S, Allen D, Annan-Phan S, *et al*. The effect of large-scale anti-contagion policies on
40 the COVID-19 pandemic. *Nature* 2020;**584**:262–7. doi:10.1038/s41586-020-2404-8
- 41 51 Lim JT, Dickens BSL, Choo ELW, *et al*. Revealing regional disparities in the transmission
42 potential of SARS-CoV-2 from interventions in Southeast Asia. *Proc R Soc B*
43 2020;**287**:20201173. doi:10.1098/rspb.2020.1173
- 44 52 Arshed N, Meo MS, Farooq F. Empirical assessment of government policies and flattening
45 of the COVID 19 curve. *J Public Affairs* Published Online First: 27 August 2020.
46 doi:10.1002/pa.2333
- 47 53 Wang K-W, Gao J, Song X-X, *et al*. Fangcang shelter hospitals are a One Health approach
48 for responding to the COVID-19 outbreak in Wuhan, China. *One Health* 2020;**10**:100167.
49 doi:10.1016/j.onehlt.2020.100167
- 50 54 Kang N, Kim B. The Effects of Border Shutdowns on the Spread of COVID-19. *J Prev Med*
51 *Public Health* 2020;**53**:293–301. doi:10.3961/jpmph.20.332
- 52 55 Auger KA, Shah SS, Richardson T, *et al*. Association Between Statewide School Closure
53 and COVID-19 Incidence and Mortality in the US. *JAMA* 2020;**324**:859.
54
55
56
57
58
59
60

- doi:10.1001/jama.2020.14348
- 56 Santamaría L, Hortal J. COVID-19 effective reproduction number dropped during Spain's nationwide dropdown, then spiked at lower-incidence regions. *Science of The Total Environment* 2021;**751**:142257. doi:10.1016/j.scitotenv.2020.142257
- 57 Bennett M. All things equal? Heterogeneity in policy effectiveness against COVID-19 spread in Chile. *World Development* 2021;**137**:105208. doi:10.1016/j.worlddev.2020.105208
- 58 Yang T, Shi H, Liu J, *et al.* Lessons Learnt from China: National Multidisciplinary Healthcare Assistance. *RMHP* 2020;**Volume 13**:1835–7. doi:10.2147/RMHP.S269523
- 59 Padalabalanarayanan S, Hanumanthu VS, Sen B. Association of State Stay-at-Home Orders and State-Level African American Population With COVID-19 Case Rates. *JAMA Netw Open* 2020;**3**:e2026010. doi:10.1001/jamanetworkopen.2020.26010
- 60 Edelstein M, Obi C, Chand M, *et al.* SARS-CoV-2 infection in London, England: changes to community point prevalence around lockdown time, March–May 2020. *J Epidemiol Community Health* 2020;:jech-2020-214730. doi:10.1136/jech-2020-214730
- 61 Tsai AC, Harling G, Reynolds Z, *et al.* COVID-19 transmission in the U.S. before vs. after relaxation of statewide social distancing measures. *Clin Infect Dis* Published Online First: 3 October 2020. doi:10.1093/cid/ciaa1502
- 62 Singh BB, Lowerison M, Lewinson RT, *et al.* Public health interventions slowed but did not halt the spread of COVID-19 in India. *Transbound Emerg Dis* Published Online First: 4 October 2020. doi:10.1111/tbed.13868
- 63 Gallaway MS, Rigler J, Robinson S, *et al.* Trends in COVID-19 Incidence After Implementation of Mitigation Measures — Arizona, January 22–August 7, 2020. *MMWR Morb Mortal Wkly Rep* 2020;**69**:1460–3. doi:10.15585/mmwr.mm6940e3
- 64 Castex G, Dechter E, Lorca M. COVID-19: The impact of social distancing policies, cross-country analysis. *EconDisCiiCha* Published Online First: 15 October 2020. doi:10.1007/s41885-020-00076-x
- 65 Silva L, Figueiredo Filho D, Fernandes A. The effect of lockdown on the COVID-19 epidemic in Brazil: evidence from an interrupted time series design. *Cad Saúde Pública* 2020;**36**:e00213920. doi:10.1590/0102-311x00213920
- 66 Dave D, Friedson A, Matsuzawa K, *et al.* Were Urban Cowboys Enough to Control COVID-19? Local Shelter-in-Place Orders and Coronavirus Case Growth. *J Urban Econ* 2020;:103294. doi:10.1016/j.jue.2020.103294
- 67 Yin Y, Gao J, Jones BF, *et al.* Coevolution of policy and science during the pandemic. *Science* 2021;**371**:128–30. doi:10.1126/science.abe3084
- 68 Wynants L, Van Calster B, Collins GS, *et al.* Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. *BMJ* 2020;:m1328. doi:10.1136/bmj.m1328
- 69 Clarke GM, Conti S, Wolters AT, *et al.* Evaluating the impact of healthcare interventions using routine data. *BMJ* 2019;:l2239. doi:10.1136/bmj.l2239
- 70 Krishnaratne S, Pfadenhauer LM, Coenen M, *et al.* Measures implemented in the school setting to contain the COVID-19 pandemic: a rapid scoping review. *Cochrane Database of Systematic Reviews* Published Online First: 17 December 2020. doi:10.1002/14651858.CD013812
- 71 Raynaud M, Zhang H, Louis K, *et al.* COVID-19-related medical research: a meta-research and critical appraisal. *BMC Med Res Methodol* 2021;**21**:1. doi:10.1186/s12874-020-01190-w
- 72 Casigliani V, De Nard F, De Vita E, *et al.* Too much information, too little evidence: is waste in research fuelling the covid-19 infodemic? *BMJ* 2020;:m2672. doi:10.1136/bmj.m2672
- 73 Greenhalgh T. Will COVID-19 be evidence-based medicine's nemesis? *PLoS Med* 2020;**17**:e1003266. doi:10.1371/journal.pmed.1003266

Acknowledgements

We would like to thank Keletso Makofane for assisting with the screening, Dr. Steven Goodman and Dr. John Ioannidis for their support during the development of this study, and Dr. Lars Hemkins and Dr. Mario Malicki for helpful comments in the protocol development.

Contributorship statement

NH led the protocol development, study design, administration, data curation, data management, statistical analysis, graphical design, manuscript writing, and manuscript editing, and serves as the primary guarantor of the study.

NH, ECD, JS, AF, and ESt co-wrote the review guidance on which the design of the study review tool is based

NH, ECD, JS, AF, ESm and ESt designed, wrote, and supported the pre-registered protocol.

NH, CJ, SW, CB, CA, CBF, VN, and Keletso Makofane were the screening reviewers for this study, analysing the abstracts and titles for inclusion criteria.

NH, ECD, AF, BMG, ES, CBF, JD, LH, CF, CB, EBM, CJ, BL, IS, EA, SW, BJ, CA, VN, BG, AB, ES were the main reviewers for this study, and contributed to the analysis and evaluation of the studies entering into the main review phase.

NH, ECD, JA, AF, BMG, ESm, CBF, JD, LH, CF, CB, EBM, CJ, BL, IS, EA, SW, BJ, CA, VN, BG, AB, and ESt all contributed to and supported the manuscript editing.

Funding

This research received no specific grant from any funding agency in the public, commercial or not-for-profit sectors.

Elizabeth Stone receives funding under the National Institutes of Health grant T32MH109436.

Ian Schmid receives funding under the National Institutes of Health grant T32MH122357.

Brooke Jarrett receives funding under the National Institutes of Health grant MH121128.

Christopher Boyer receives funding under the National Institutes of Health grant T32HL098048

Cathrine Axfors receives funding from the Knut and Alice Wallenberg Foundation, grant KAW 2019.0561.

Beth Ann Griffin and Elizabeth Stuart were supported by award number P50DA046351 from the National Institute on Drug Abuse. Elizabeth Stuart's time was also supported by the Bloomberg American Health Initiative. Caroline Joyce receives funding from the Ferring Foundation. Meta-Research Innovation Center at Stanford (METRICS), Stanford University is supported by a grant from the Laura and John Arnold Foundation

Data sharing

Data, code, the review tool, and the review guidance are stored and available at the OSF.io repository for this study here: <https://osf.io/9xmke/files/>. The dataset includes full results from the search and screening and all review tool responses from reviewers during the full review phase.

Competing interests

The authors have no competing financial or social conflicts of interest to declare.

Figures

Figure 1: PRISMA diagram of systematic review process

Caption: This chart shows the PRISMA diagram for the process of screening the literature from search to the full review phase.

Figure 2: Descriptive sample statistics (n=36)

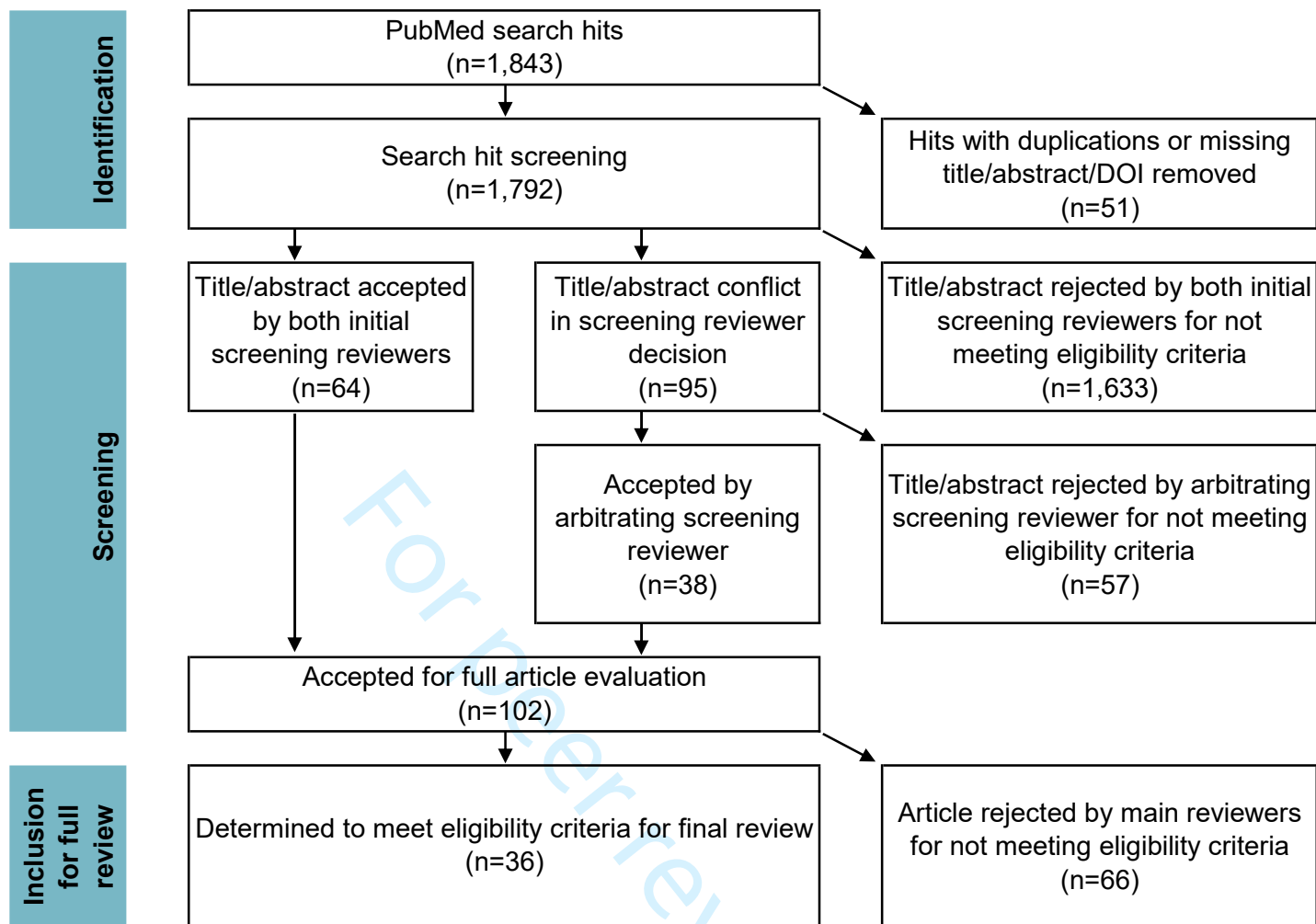
Caption: This chart shows descriptive statistics of the 36 studies entered into our systematic evidence review.

Figure 3: Main consensus results summary for key and overall questions

Caption: This chart shows the final overall ratings (left) and the key design question ratings for the consensus review of the 36 included studies, answering the degree to which the articles met the given key design question criteria. The key design question ratings were not asked for the nine included articles which selected methods assumed by the guidance to be non-appropriate. The question prompt in the figure is shortened for clarity, where the full prompt for each key question is available in the Methods section.

Figure 4: Comparison of independent reviews, weakest link, and direct consensus review

Caption: This chart shows the final overall ratings by three different possible metrics. The first column contains all of the independent review ratings for the 27 studies which were eventually included in our sample, noting that reviewers who either selected them as not meeting inclusion criteria or selected a method that didn't receive the full review did not contribute. The middle column contains the final consensus reviews among the 27 articles which received full review. The last column contains the weakest link rating, as described in the methods section. The question prompt in the figure is shortened for clarity, where the full prompt for each key question is available in the Methods section.



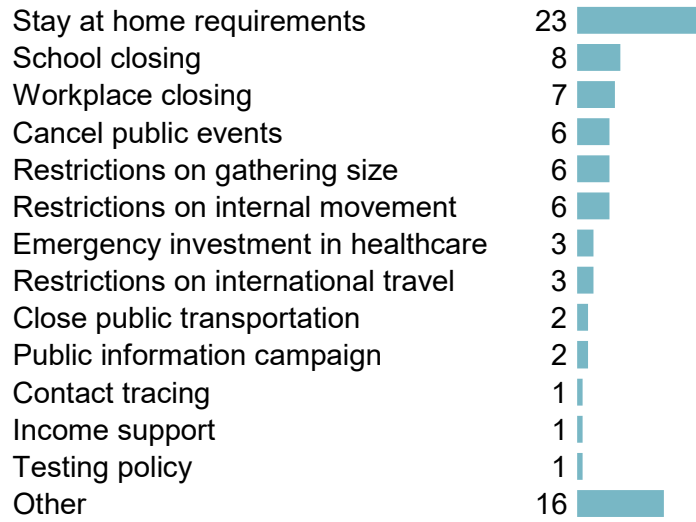
Publication month

n



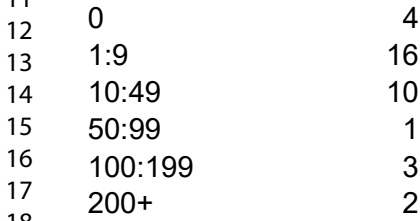
Policy type

n



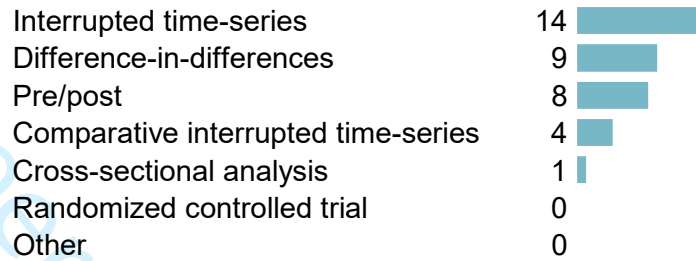
Citations

n



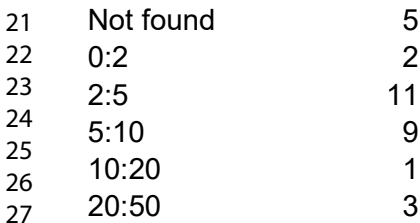
Method

n

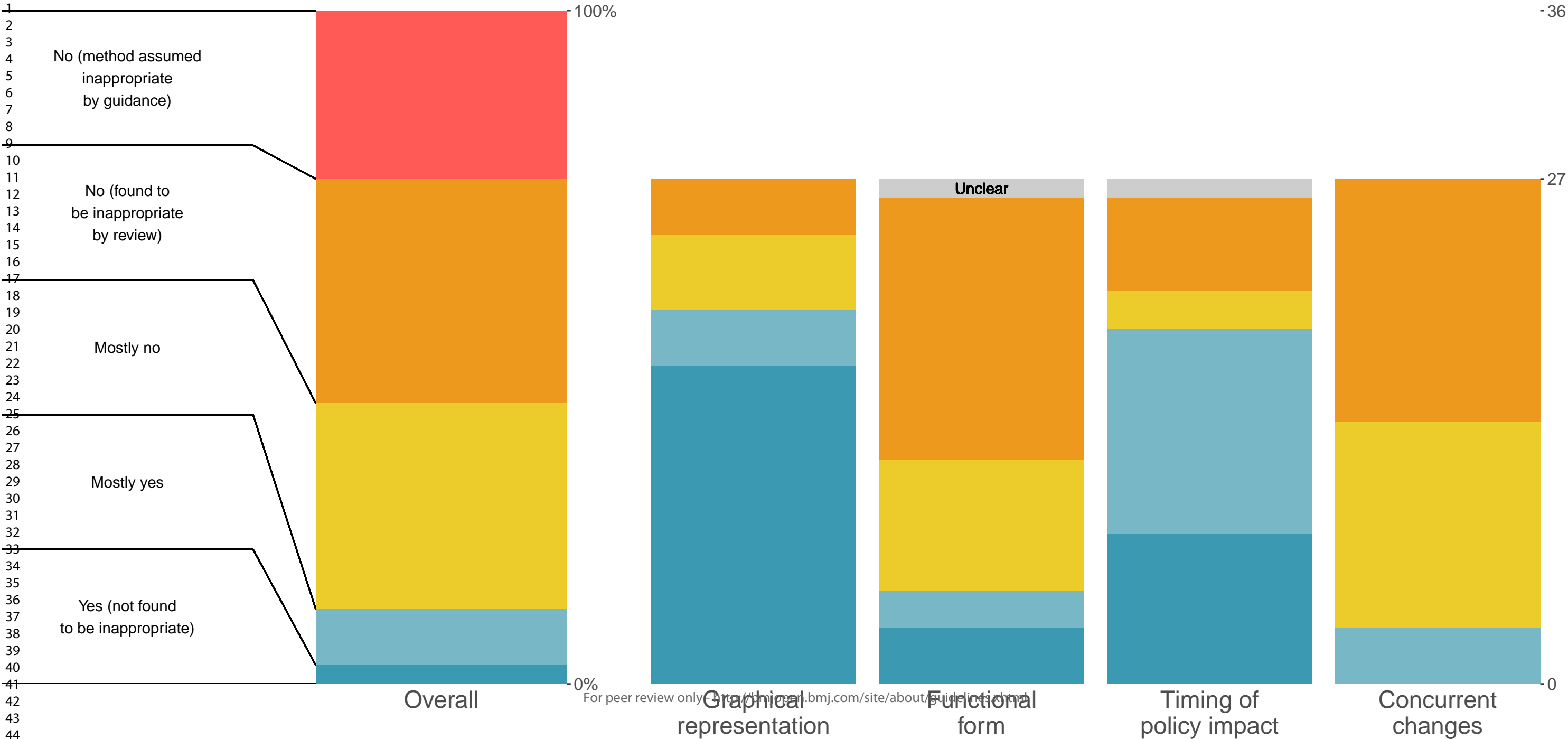


Journal impact factor

n



Did the study meet design criteria?



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

No (method assumed inappropriate by guidance)

No (found to be inappropriate by review)

Mostly no

Mostly yes

Yes (not found to be inappropriate)

Overall

Graphical representation

Functional form

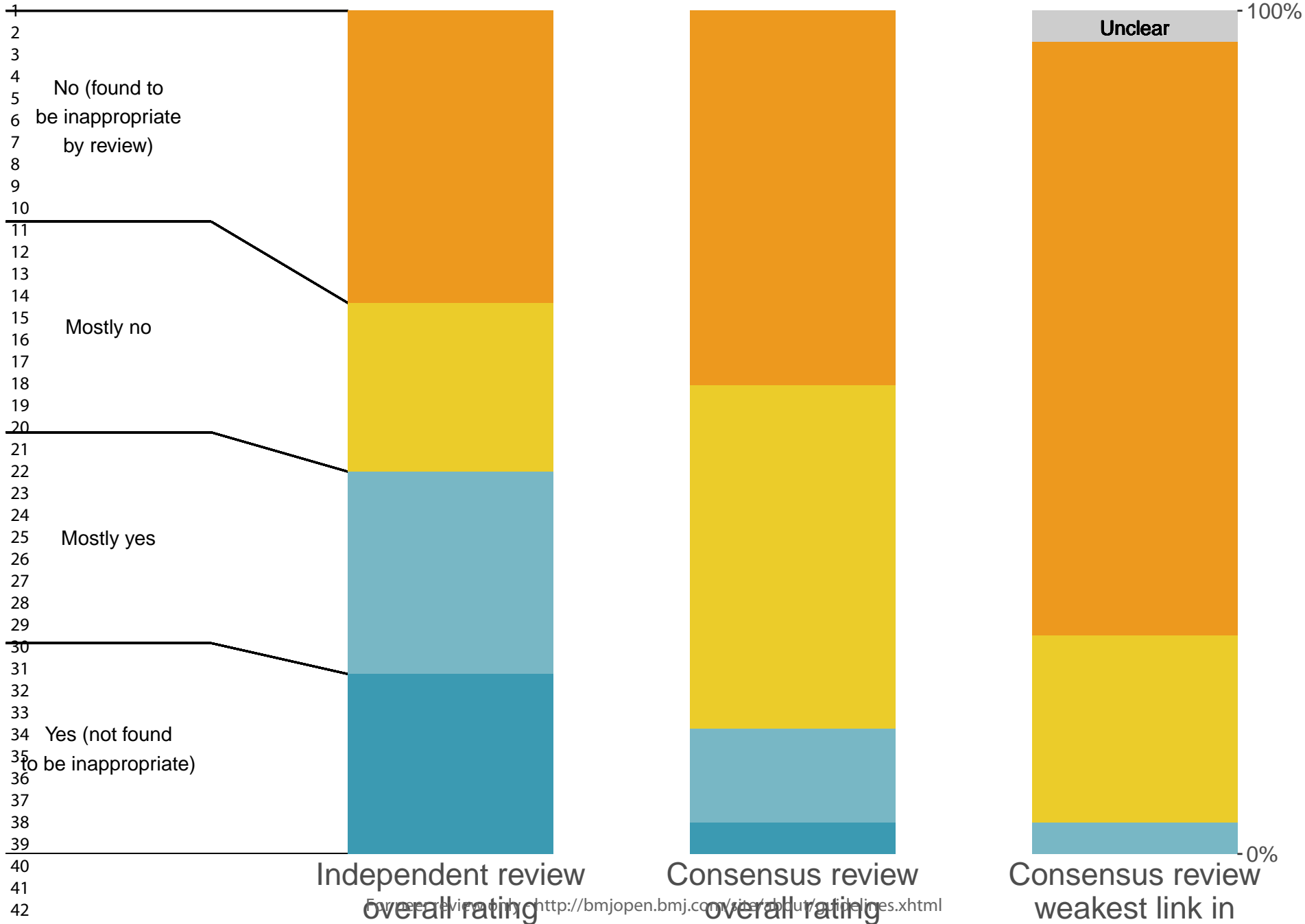
Timing of policy impact

Concurrent changes

Number of articles

For peer review only: <https://bmjopen.bmj.com/site/about/guidelines.xhtml>

Did the study meet design criteria?



Appendix 1: Changes from pre-registered protocol and justifications

The full, original pre-registered protocol is available here: <https://osf.io/7nbk6>

Inclusion criteria

Minor language edits were made to the inclusion criteria to improve clarity and fix grammatical and typographical errors. This largely centered around improving clarity that a study must estimate the quantitative impact of policies that had already been enacted. The word “quantitative” was not explicitly stated in the original version.

Procedures

The original protocol specified that each article would receive two independent reviewers. This was increased to three reviewers per article once it became clear both that the number of articles which would be accepted for full review was lower than expectations, and that there would be substantial differences in opinion between reviewers.

Statistical analysis

Firstly, the original protocol specified that 95% confidence intervals would be calculated. However, after further discussion and review, we determined that sampling-based confidence intervals were not appropriate. Our results are not indicative nor intended to be representative of any super- or target-population, and as such sampling-based error is not an appropriate metric for the conclusions of this study.

Secondly, the original protocol specified Kappa-based interrater reliability statistics. However, using three reviewers, rather than the originally registered two, meant that most Kappa statistics would not be appropriate for our review process. Given the three-rater, four-level ordinal scale used, we opted instead to use Krippendorff's Alpha.

Review tool

A number of changes were made to the review tool during the course of the review process. While the original protocol included logic to allow pre/post for review in some of the key questions, this was removed for consistency with the guidance document.

The remaining changes to the review tool were error corrections and clarifications (e.g. correcting the text for the concurrent changes sections in difference-in-differences so that it

1
2
3 stated “uncontrolled” concurrent changes, and distinguishing the DiD/CITS requirements from
4 the ITS requirements to emphasize differential concurrent changes).
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Appendix 2: Full search terms

Note: The search filter for COVID-19 and SARS-CoV-2 were the exact search terms used for the National Library of Medicine one-click search option at the time of the protocol development and when the search took place. This reflects that some of the early literature referred to Wuhan specifically (both in geographic reference for where the SARS-CoV-2 was initially found, and unfortunately also early naming of the virus/disease) before official naming conventions became ubiquitous in the literature. In order to comprehensively capture the literature and use searching best practices, we used the most standard and recommended terms.

(((((wuhan[All Fields] AND ("coronavirus"[MeSH Terms] OR "coronavirus"[All Fields])) AND 2019/12[PDAT] : 2030[PDAT])) OR 2019-nCoV[All Fields] OR 2019nCoV[All Fields] OR COVID-19[All Fields] OR SARS-CoV-2[All Fields])

AND ("impact"[TIAB] OR "effect"[TIAB])

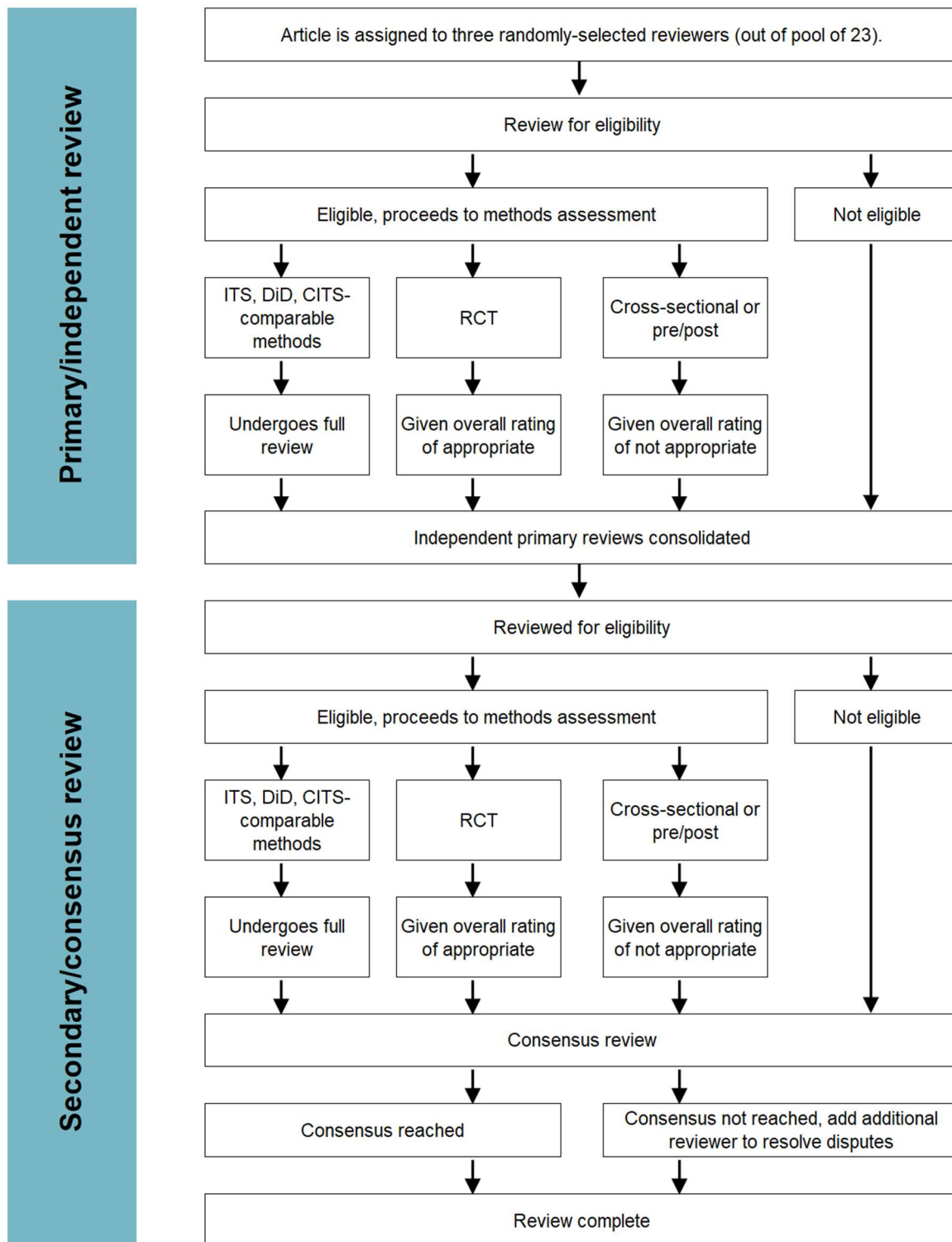
AND ("policy"[TIAB] OR "policies"[TIAB] OR "order"[TIAB] OR "mandate"[TIAB])

AND ("countries"[TIAB] OR "country"[TIAB] OR "state"[TIAB] OR "provinc"[TIAB] OR "county"[TIAB] OR "parish"[TIAB] OR "region"[TIAB] OR "city"[TIAB] OR "cities"[TIAB] OR "continent"[TIAB] "Asia"[TIAB] OR "Europe"[TIAB] OR "Africa"[TIAB] OR "America"[TIAB] OR "Australia"[TIAB] OR "Antarctica"[TIAB] OR "Afghanistan"[TIAB] OR "Aland Islands"[TIAB] OR "Åland Islands"[TIAB] OR "Albania"[TIAB] OR "Algeria"[TIAB] OR "American Samoa"[TIAB] OR "Andorra"[TIAB] OR "Angola"[TIAB] OR "Anguilla"[TIAB] OR "Antarctica"[TIAB] OR "Antigua"[TIAB] OR "Argentina"[TIAB] OR "Armenia"[TIAB] OR "Aruba"[TIAB] OR "Australia"[TIAB] OR "Austria"[TIAB] OR "Azerbaijan"[TIAB] OR "Bahamas"[TIAB] OR "Bahrain"[TIAB] OR "Bangladesh"[TIAB] OR "Barbados"[TIAB] OR "Barbuda"[TIAB] OR "Belarus"[TIAB] OR "Belgium"[TIAB] OR "Belize"[TIAB] OR "Benin"[TIAB] OR "Bermuda"[TIAB] OR "Bhutan"[TIAB] OR "Bolivia"[TIAB] OR "Bonaire"[TIAB] OR "Bosnia"[TIAB] OR "Botswana"[TIAB] OR "Bouvet Island"[TIAB] OR "Brazil"[TIAB] OR "British Indian Ocean Territory"[TIAB] OR "Brunei"[TIAB] OR "Bulgaria"[TIAB] OR "Burkina Faso"[TIAB] OR "Burundi"[TIAB] OR "Cabo Verde"[TIAB] OR "Cambodia"[TIAB] OR "Cameroon"[TIAB] OR "Canada"[TIAB] OR "Cayman Islands"[TIAB] OR "Central African Republic"[TIAB] OR "Chad"[TIAB] OR "Chile"[TIAB] OR "China"[TIAB] OR "Christmas Island"[TIAB] OR "Cocos Islands"[TIAB] OR "Colombia"[TIAB] OR "Comoros"[TIAB] OR "Congo"[TIAB] OR "Congo"[TIAB] OR "Cook Islands"[TIAB] OR "Costa Rica"[TIAB] OR "Côte d'Ivoire"[TIAB] OR "Croatia"[TIAB] OR "Cuba"[TIAB] OR "Curaçao"[TIAB] OR "Cyprus"[TIAB] OR "Czechia"[TIAB] OR "Denmark"[TIAB] OR "Djibouti"[TIAB] OR "Dominica"[TIAB] OR "Dominican Republic"[TIAB] OR "Ecuador"[TIAB] OR "Egypt"[TIAB] OR "El Salvador"[TIAB] OR "Equatorial Guinea"[TIAB] OR "Eritrea"[TIAB] OR "Estonia"[TIAB] OR "Eswatini"[TIAB] OR "Ethiopia"[TIAB] OR "Falkland Islands"[TIAB] OR "Faroe Islands"[TIAB] OR "Fiji"[TIAB] OR "Finland"[TIAB] OR "France"[TIAB] OR "French Guiana"[TIAB] OR "French Polynesia"[TIAB] OR "French Southern

1
2
3 Territories"[TIAB] OR "Futuna"[TIAB] OR "Gabon"[TIAB] OR "Gambia"[TIAB] OR
4 "Georgia"[TIAB] OR "Germany"[TIAB] OR "Ghana"[TIAB] OR "Gibraltar"[TIAB] OR
5 "Greece"[TIAB] OR "Greenland"[TIAB] OR "Grenada"[TIAB] OR "Grenadines"[TIAB] OR
6 "Guadeloupe"[TIAB] OR "Guam"[TIAB] OR "Guatemala"[TIAB] OR "Guernsey"[TIAB] OR
7 "Guinea"[TIAB] OR "Guinea-Bissau"[TIAB] OR "Guyana"[TIAB] OR "Haiti"[TIAB] OR "Heard
8 Island"[TIAB] OR "Herzegovina"[TIAB] OR "Holy See"[TIAB] OR "Honduras"[TIAB] OR "Hong
9 Kong"[TIAB] OR "Hungary"[TIAB] OR "Iceland"[TIAB] OR "India"[TIAB] OR "Indonesia"[TIAB]
10 OR "Iran"[TIAB] OR "Iraq"[TIAB] OR "Ireland"[TIAB] OR "Isle of Man"[TIAB] OR "Israel"[TIAB]
11 OR "Italy"[TIAB] OR "Jamaica"[TIAB] OR "Jan Mayen Islands"[TIAB] OR "Japan"[TIAB] OR
12 "Jersey"[TIAB] OR "Jordan"[TIAB] OR "Kazakhstan"[TIAB] OR "Keeling Islands"[TIAB] OR
13 "Kenya"[TIAB] OR "Kiribati"[TIAB] OR "Korea"[TIAB] OR "Korea"[TIAB] OR "Kuwait"[TIAB] OR
14 "Kyrgyzstan"[TIAB] OR "Lao People's Democratic Republic"[TIAB] OR "Laos"[TIAB] OR
15 "Latvia"[TIAB] OR "Lebanon"[TIAB] OR "Lesotho"[TIAB] OR "Liberia"[TIAB] OR "Libya"[TIAB]
16 OR "Liechtenstein"[TIAB] OR "Lithuania"[TIAB] OR "Luxembourg"[TIAB] OR "Macao"[TIAB] OR
17 "Madagascar"[TIAB] OR "Malawi"[TIAB] OR "Malaysia"[TIAB] OR "Maldives"[TIAB] OR
18 "Mali"[TIAB] OR "Malta"[TIAB] OR "Malvinas"[TIAB] OR "Marshall Islands"[TIAB] OR
19 "Martinique"[TIAB] OR "Mauritania"[TIAB] OR "Mauritius"[TIAB] OR "Mayotte"[TIAB] OR
20 "McDonald Islands"[TIAB] OR "Mexico"[TIAB] OR "Micronesia"[TIAB] OR "Moldova"[TIAB] OR
21 "Monaco"[TIAB] OR "Mongolia"[TIAB] OR "Montenegro"[TIAB] OR "Montserrat"[TIAB] OR
22 "Morocco"[TIAB] OR "Mozambique"[TIAB] OR "Myanmar"[TIAB] OR "Namibia"[TIAB] OR
23 "Nauru"[TIAB] OR "Nepal"[TIAB] OR "Netherlands"[TIAB] OR "Nevis"[TIAB] OR "New
24 Caledonia"[TIAB] OR "New Zealand"[TIAB] OR "Nicaragua"[TIAB] OR "Niger"[TIAB] OR
25 "Nigeria"[TIAB] OR "Niue"[TIAB] OR "Norfolk Island"[TIAB] OR "North Macedonia"[TIAB] OR
26 "Northern Mariana Islands"[TIAB] OR "Norway"[TIAB] OR "Oman"[TIAB] OR "Pakistan"[TIAB]
27 OR "Palau"[TIAB] OR "Panama"[TIAB] OR "Papua New Guinea"[TIAB] OR "Paraguay"[TIAB]
28 OR "Peru"[TIAB] OR "Philippines"[TIAB] OR "Pitcairn"[TIAB] OR "Poland"[TIAB] OR
29 "Portugal"[TIAB] OR "Principe"[TIAB] OR "Puerto Rico"[TIAB] OR "Qatar"[TIAB] OR
30 "Réunion"[TIAB] OR "Romania"[TIAB] OR "Russian Federation"[TIAB] OR "Rwanda"[TIAB] OR
31 "Saba"[TIAB] OR "Saint Barthélemy"[TIAB] OR "Saint Helena"[TIAB] OR "Saint Kitts"[TIAB] OR
32 "Saint Lucia"[TIAB] OR "Saint Martin"[TIAB] OR "Saint Pierre and Miquelon"[TIAB] OR "Saint
33 Vincent"[TIAB] OR "Samoa"[TIAB] OR "San Marino"[TIAB] OR "Sao Tome"[TIAB] OR
34 "Sark"[TIAB] OR "Saudi Arabia"[TIAB] OR "Senegal"[TIAB] OR "Serbia"[TIAB] OR
35 "Seychelles"[TIAB] OR "Sierra Leone"[TIAB] OR "Singapore"[TIAB] OR "Sint Eustatius"[TIAB]
36 OR "Sint Maarten"[TIAB] OR "Slovakia"[TIAB] OR "Slovenia"[TIAB] OR "Solomon
37 Islands"[TIAB] OR "Somalia"[TIAB] OR "South Africa"[TIAB] OR "South Georgia"[TIAB] OR
38 "South Sandwich Islands"[TIAB] OR "South Sudan"[TIAB] OR "Spain"[TIAB] OR "Sri
39 Lanka"[TIAB] OR "State of Palestine"[TIAB] OR "Sudan"[TIAB] OR "Suriname"[TIAB] OR
40 "Svalbard"[TIAB] OR "Sweden"[TIAB] OR "Switzerland"[TIAB] OR "Syria"[TIAB] OR "Syrian
41 Arab Republic"[TIAB] OR "Tajikistan"[TIAB] OR "Thailand"[TIAB] OR "Timor-Leste"[TIAB] OR
42 "Tobago"[TIAB] OR "Togo"[TIAB] OR "Tokelau"[TIAB] OR "Tonga"[TIAB] OR "Trinidad"[TIAB]
43 OR "Tunisia"[TIAB] OR "Turkey"[TIAB] OR "Turkmenistan"[TIAB] OR "Turks and Caicos"[TIAB]
44 OR "Tuvalu"[TIAB] OR "Uganda"[TIAB] OR "UK"[TIAB] OR "Ukraine"[TIAB] OR "United Arab
45 Emirates"[TIAB] OR "United Kingdom"[TIAB] OR "United Republic of Tanzania"[TIAB] OR
46 "United States Minor Outlying Islands"[TIAB] OR "United States of America"[TIAB] OR
47
48
49
50
51
52
53
54
55
56
57
58
59

1
2
3 "Uruguay"[TIAB] OR "USA"[TIAB] OR "Uzbekistan"[TIAB] OR "Vanuatu"[TIAB] OR
4 "Venezuela"[TIAB] OR "Viet Nam"[TIAB] OR "Vietnam"[TIAB] OR "Virgin Islands"[TIAB] OR
5 "Virgin Islands"[TIAB] OR "Wallis"[TIAB] OR "Western Sahara"[TIAB] OR "Yemen"[TIAB] OR
6 "Zambia"[TIAB] OR "Zimbabwe"[TIAB] OR "Alabama"[TIAB] OR "Alaska"[TIAB] OR
7 "Arizona"[TIAB] OR "Arkansas"[TIAB] OR "California"[TIAB] OR "Colorado"[TIAB] OR
8 "Connecticut"[TIAB] OR "Delaware"[TIAB] OR "Florida"[TIAB] OR "Georgia"[TIAB] OR
9 "Hawaii"[TIAB] OR "Idaho"[TIAB] OR "Illinois"[TIAB] OR "Indiana"[TIAB] OR "Iowa"[TIAB] OR
10 "Kansas"[TIAB] OR "Kentucky"[TIAB] OR "Louisiana"[TIAB] OR "Maine"[TIAB] OR
11 "Maryland"[TIAB] OR "Massachusetts"[TIAB] OR "Michigan"[TIAB] OR "Minnesota"[TIAB] OR
12 "Mississippi"[TIAB] OR "Missouri"[TIAB] OR "Montana"[TIAB] OR "Nebraska"[TIAB] OR
13 "Nevada"[TIAB] OR "New Hampshire"[TIAB] OR "New Jersey"[TIAB] OR "New Mexico"[TIAB]
14 OR "New York"[TIAB] OR "North Carolina"[TIAB] OR "North Dakota"[TIAB] OR "Ohio"[TIAB] OR
15 "Oklahoma"[TIAB] OR "Oregon"[TIAB] OR "Pennsylvania"[TIAB] OR "Rhode Island"[TIAB] OR
16 "South Carolina"[TIAB] OR "South Dakota"[TIAB] OR "Tennessee"[TIAB] OR "Texas"[TIAB] OR
17 "Utah"[TIAB] OR "Vermont"[TIAB] OR "Virginia"[TIAB] OR "Washington"[TIAB] OR "West
18 Virginia"[TIAB] OR "Wisconsin"[TIAB] OR "Wyoming"[TIAB] OR "Ontario"[TIAB] OR
19 "Quebec"[TIAB] OR "Nova Scotia"[TIAB] OR "New Brunswick"[TIAB] OR "Manitoba"[TIAB] OR
20 "British Columbia"[TIAB] OR "Prince Edward Island"[TIAB] OR "Saskatchewan"[TIAB] OR
21 "Alberta"[TIAB] OR "Newfoundland"[TIAB] OR "Labrador"[TIAB])
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Appendix 3: Article review flow diagram



COVID-19 Health Policy Impact Evaluation Review

Start of Block: Main form

Q10 Administrative information

Q8 Study DOI

Q3 Reviewer number

Q54 Review type/round

The first round (Primary/independent review round) is for the independent first reviews of every article; the second (Secondary/consensus round) is for the second round of review for each article.

- Primary/independent review round (1)
- Secondary/consensus round (2)

Q50 Screening

1
2
3 Q52 Do you wish to recuse yourself from reviewing this study for any reason (e.g. social or
4 professional relationship with the authors, financial conflict of interest, etc)?
5

- 6
7 No, I do not wish to recuse myself. (1)
8
9 Yes, I recuse myself from reviewing this paper. (2)
10

11
12 *Skip To: End of Survey If Q52 = Yes, I recuse myself from reviewing this paper.*
13
14

15
16 Q51 Do you believe that this study meets the inclusion criteria for this research?
17

18 The inclusion criteria are: The primary topic of the article must be evaluating one or more
19 individual COVID-19 policies on direct COVID-19 outcomes The primary
20 exposure(s) must be a policy, defined as a government-issued order at any government level to
21 address a directly COVID-19-related outcome (e.g. mask requirements, travel restrictions, etc).
22

23 COVID-19 outcomes may include cases detected, mortality, number of tests taken, test
24 positivity rates, Rt, etc. This may NOT include indirect impacts of COVID-19 on
25 things such as income, childcare, trust in science, etc. The primary outcome
26 being examined must be a COVID-19-specific outcome, as above. The study must be
27 designed as an impact evaluation study from primary data (i.e. not primarily a predictive or
28 simulation model or meta-analysis) The study must be peer reviewed, and published in a peer-
29 reviewed journal indexed by PubMed The study must have the title and abstract available
30 via PubMed at the time of the study start date The study must be written in English
31
32
33

- 34 Yes (1)
35
36 No (2) _____
37
38

39
40 *Skip To: End of Survey If Q51 = No*
41
42

43 Q7 Study topic information

44
45
46 Please consult review guidance ([available here](#)) for additional guidance on answering these
47 questions.
48
49

50
51
52 Q6 Main impact sentence
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Copy and paste the sentence from the abstract that best describes the main claim of the study (e.g. "Policy X had a positive impact on outcome Y")

For peer review only

1
2
3 Q9 Main COVID-19 policy type evaluated
45 Select all that apply. Note: categorization from the Oxford Government Response Tracker
6

- 7
-
- 8
-
- School closing (1)
-
- 9
-
- 10
-
- Workplace closing (2)
-
- 11
-
- 12
-
- Cancel public events (3)
-
- 13
-
- 14
-
- Restrictions on gathering size (4)
-
- 15
-
- 16
-
- Close public transportation (5)
-
- 17
-
- 18
-
- Stay at home requirements (6)
-
- 19
-
- 20
-
- Restrictions on internal movement (7)
-
- 21
-
- 22
-
- Restrictions on international travel (8)
-
- 23
-
- 24
-
- Income support (9)
-
- 25
-
- 26
-
- Debt/contract relief for household (10)
-
- 27
-
- 28
-
- Fiscal measures (11)
-
- 29
-
- 30
-
- Giving international support (12)
-
- 31
-
- 32
-
- Public information campaign (13)
-
- 33
-
- 34
-
- Testing policy (14)
-
- 35
-
- 36
-
- Contact tracing (15)
-
- 37
-
- 38
-
- Emergency investment in healthcare (16)
-
- 39
-
- 40
-
- Investment in COVID-19 vaccines (17)
-
- 41
-
- 42
-
- 43
-
- 44
-
- 45
-
- 46
-
- 47
-
- 48
-
- 49
-
- 50
-
- 51
-
- 52
-
- 53
-
- 54
-
- 55
-
- 56
-
- 57
-
- 58
-
- 59
-
- 60

Other policy response (fill in) (18)

Q12 Main COVID-19 outcome type evaluated

Select all that apply

- COVID-19 cases (1)
- COVID-19 test positivity (2)
- COVID-19 deaths (3)
- COVID-19 hospitalizations (4)
- SARS-CoV-2 infections and infection rate (e.g. effective R) (8)
- Other (fill in) (9) _____

Q13 **Method(s) identification**

For this section, consider only the data structure as it enters into the main statistical model. In other words, if the original dataset is of individuals at many time points, but the main statistical model uses a regional-level aggregated count of cases, the data as it enters into the main statistical model is a regional aggregate at one time point.

Q14 What is the level of aggregation for the main outcome data?

- Individual level (1)
- Regional aggregate (e.g. count, mean, etc.) (2)

1
2
3
4 Q16 How many regional units included in the main statistical model received the policy of
5 interest?
6

7
8 If 2-20, enter the number of regional units analyzed which received the policy of interest.
9

10 One (1) (1)

11
12 Two through twenty (2-20) (2)
13
14 _____
15

16 More than twenty (21+) (3)

17
18 Unclear or N/A (4) _____
19
20
21

22
23
24 Q17 How many regional units were included which did NOT receive any form of the policy of
25 interest?
26

27
28 If 2-20, enter the number of regional units analyzed which did not receive the policy of interest.
29

30 Zero (0) (1)

31
32 One (1) (2)

33
34 Two through twenty (2-20) (3)
35
36 _____
37

38 More than twenty (21+) (4)

39
40 Unclear or N/A (5) _____
41
42
43
44

45 *Display This Question:*

46 *If Q17 = Zero (0)*
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q25 Did different regions receive different intensities of the policy of interest for comparison?
4

5 For example, the study might compare places with more intense versions of policy or policies
6 vs. places with less intense versions of policy or policies, rather than just places with and
7 without the policy or policies.
8
9

- 10 Yes (regions with more intense policy were compared with regions with less intense
11 policy) (1)
12
13 No (2)
14
15 Unclear or N/A (3)
16
17
18
-

19
20
21
22 Q18 For each regional unit, how many time point observations were in the model *before* the
23 policy was enacted?
24

- 25 None (0) (1)
26
27 One (1) (2)
28
29 More than one (2+) (3)
30
31 Unclear or N/A (4) _____
32
33
34
-

35
36
37
38 Q19 For each regional unit, how many time point observations were in the model *after* the policy
39 was enacted?
40

- 41 None (0) (1)
42
43 One (1) (2)
44
45 More than one (2+) (3)
46
47 Unclear or N/A (4) _____
48
49
50
51
-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q19 = One (1)

And Q18 = One (1)

Or If

Q19 = More than one (2+)

Or If

Q18 = More than one (2+)

Q20 How would you describe the time intervals between observations?

- Days (1-5 days between observations) (1)
- Weeks (about 5-10 days between observations) (2)
- Multiple weeks (11-25 days between observations) (3)
- Monthly (26 or more days between observations) (4)

Display This Question:

If Q17 = Zero (0)

Q21 Did the pre-policy period for any region act as a "control" for different region post-policy enactment?

In other words, was there any pre-period in one or more region's being used to control or compare for the trends of any one or more *different* regions' post-period?

- No (pre-periods were treated as controls only within-region) (1)
- Yes (pre-periods were treated as controls with other regions) (2)
- Unclear or N/A (3)
-

Q22 Was any unit assigned the policy or the timing of the policy externally (i.e. as an experiment/trial)?

- No (observational data only) (1)
- Yes (treatment assigned as part of research or evaluation) (2)
- Unclear or N/A (3)

Display This Question:

If Q22 = Yes (treatment assigned as part of research or evaluation)

Q23 Was the assignment randomized?

- Yes (1)
- No (2)

Q27 Based on your answers above and the guidance document, please select the type of study that best resembles the design of the main analysis.

Please note that the design(s) named in the paper may not match with the method described below, nor is this the actual exact design that was used. If you believe that the design used differs from the choices below in a way that makes this choice impossible, please contact the study administrator before selecting "other."

	Design	
	Units (e.g., regions of comparison)	
	Time points measured per unit	
	Assumed counterfactual.	
	"If not for the intervention, ___"	
	Without intervention	With intervention
	Before intervention	
	After intervention	
	Cross-sectional	
	At least one	
At least one		N/A
	One time point	
	Outcome in intervention units would have been the same	

1
2
3 as the outcome in the non-intervention units.

4 Pre/post

6 At least one

8 None

9 At least one (typically one)

10 At least one (typically one)

11 Outcome would have stayed the same from the pre period to the post period.

13 Interrupted

14 time-series

15 (ITS)

18 At least one

19 None

20 More than one

21 At least one (typically several)

22 Outcome slope and level* would have continued along the same modelled trajectory
23 from the pre-period to the post period.

25 Difference-in-differences

(DiD)

28 At least one

29 At least one†

30 At least one (typically one)

31 At least one (typically one)

32 Outcome in intervention units would have changed as much as (or in parallel with) the
33 outcome in the non-intervention units.

35 Comparative interrupted time series (CITS)

37 At least one

38 At least one†

39 More than one (typically several)

40 At least one (typically several)

41 Outcome slope and level* would have changed as much as non-
42 intervention group's slope and level* changed.

44 * Assessing both slope and level only applicable if
45 there are multiple data points during the post period
46 † Units without the
47 intervention may be the pre-period of a different unit that eventually receives the intervention.

50 Cross-sectional analysis (1)

52 Non-randomized experiment/trial (2)

54 Randomized controlled trial (3)

- 1
2
3
4 Pre/post (4)
5
6 Interrupted time-series (5)
7
8
9 Difference-in-differences (6)
10
11 Comparative interrupted time-series (7)
12
13
14 Other (please contact administrator before selecting) (8)
15 _____
16

17
18
19
20 **Q49 Design evaluation**
21
22

23
24 *Display This Question:*

25 *If Q27 = Interrupted time-series*

26 *Or Q27 = Difference-in-differences*

27 *Or Q27 = Comparative interrupted time-series*
28
29
30

31 **Q29 Does the analysis provide graphical representation of the outcome over time?**
32
33

34
35 If not "Yes" please describe (three short sentences max).
36

37 -Check for a chart that shows the outcome over time, with the dates of interest, separated by
38 policy/non policy groups if applicable. Outcomes may be aggregated for clarity (e.g. means and
39 CIs at discrete time points).
40

- 41
42 Yes (1)
43
44 Mostly yes (2) _____
45
46 Mostly no (3) _____
47
48 No (4) _____
49
50 Unclear (5) _____
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q30 Is there sufficient pre-intervention data to characterize pre-trends in the data?

If not "Yes" please describe (three short sentences max).

-Check the chart(s) to see if there are several time points over a reasonable period of time over which to establish stability and curvature in the pre-trends.

Yes (1)

Mostly yes (2) _____

Mostly no (3) _____

No (4) _____

Unclear (5) _____

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

1
2
3 Q32 Is the pre-trend stable?
4
5

6
7 If not "Yes" please describe (three short sentences max).
8

9 -Check if there are sufficient data to reasonably determine a stable functional form for the pre-
10 trends, and that they follow a modelable functional form.
11

- 12
13 Yes (1)
14
15 Mostly yes (2) _____
16
17 Mostly no (3) _____
18
19 No (4) _____
20
21 Unclear (5) _____
22
23
24

25
26
27 *Display This Question:*

28 *If Q27 = Interrupted time-series*

29 *Or Q27 = Comparative interrupted time-series*
30

31
32 Q31 Is there sufficient post-intervention data to observe post trends in the data?
33
34

35
36 If not "Yes" please describe (three short sentences max).
37

38 -Check the chart(s) to see if there are several time points over a reasonable period of time over
39 which to establish stability and curvature in the post- trends.
40

- 41
42 Yes (1)
43
44 Mostly yes (2) _____
45
46 Mostly no (3) _____
47
48 No (4) _____
49
50 Unclear (5) _____
51
52
53
54

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q33 Is the functional form of the counterfactual (e.g. linear) well-justified and appropriate?

If not "Yes" please describe (three short sentences max).

- Check whether the authors explain and justify their choice of functional form.
- Check if there is any curvature in the pre-trend.
- Consider the nature of the outcome. Is it sensible to measure the trend of this outcome on the scale and form used? Note: infectious disease dynamics are rarely linear.
- Consider that while pre-trend fit is a necessary condition for an appropriate linear counterfactual model, it is not sufficient. Check if the authors provide justification for the functional form to continue to be of the same functional form (e.g. linear).

- Yes (1)
- Mostly yes (2) _____
- Mostly no (3) _____
- No (4) _____
- Unclear (5) _____

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q34 Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?

If not "Yes" please describe (three short sentences max).

- Check whether the authors justify the use of the date threshold relative to the date of the intervention.
- Trace the process between the intervention being put in place to when observable effects in

1
2
3 the outcome might appear over time.

4 -Consider whether there are anticipation effects (e.g. do people change behaviors before the
5 date when the intervention begins?)

6
7 -Consider whether there are lag effects. (e.g. does it take time for behaviors to change,
8 behavior change to impact infections, infections to impact testing, and data to be collected, etc?)

9 -Check if authors appropriately and directly account for these time effects.
10

11 Yes (1)

12 Mostly yes (2) _____
13

14 Mostly no (3) _____
15

16 No (4) _____
17

18 Unclear (5) _____
19
20
21
22
23

24 -----
25 *Display This Question:*

26 *If Q27 = Interrupted time-series*
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Q36 Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have
4 changed during the measurement period?
5

6
7 If not "Yes" please describe (three short sentences max).
8

- 9 -Consider other policies or interventions which could impact the outcome during this time.
10 -Consider social behaviors changed which could meaningfully impact the outcome during this
11 time.
12 -Consider economic conditions changed which could meaningfully impact the outcome during
13 this time.
14 -Note that the actual concurrent changes do not need to happen during the period of
15 measurement, just their effects.
16
17

- 18
19 Yes (1)
20
21 Mostly yes (2) _____
22
23 Mostly no (3) _____
24
25 No (4) _____
26
27 Unclear (5) _____
28
29
30
31

32
33 *Display This Question:*

34 *If Q27 = Difference-in-differences*

35 *Or Q27 = Comparative interrupted time-series*
36
37

38
39 Q53

40 Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have
41 changed during the measurement period, differently for policy and non-policy regions?
42

43 If not "Yes" please describe (three short sentences max).
44

45 -Consider any uncontrolled factor which could have influenced the outcome differently in policy
46 and non-policy regions.
47

48 -This may include (but is not limited to)

- 49 -Other policies
50 -Social behaviors
51 -Economic conditions
52

53 -Are these factors justified as having negligible impact?

54 -If justified, is the argument that these have negligible impact convincing?
55
56
57
58
59
60

1
2
3 -Note that the actual concurrent changes do not need to happen during the period of
4 measurement, just their effects.
5

- 6
7 Yes (1)
8
9 Mostly yes (2) _____
10
11 Mostly no (3) _____
12
13 No (4) _____
14
15 Unclear (5) _____
16
17
18

19
20
21 *Display This Question:*

22 *If Q27 = Interrupted time-series*

23 *Or Q27 = Difference-in-differences*

24 *Or Q27 = Comparative interrupted time-series*
25
26
27

28 Q38

29 Did authors provide diagnostics or show robustness and/or sensitivity of results to alternative
30 model choices?
31

32
33
34 If not "Yes" please describe (three short sentences max).
35

- 36 Yes (1)
37
38 Mostly yes (2) _____
39
40 Mostly no (3) _____
41
42 No (4) _____
43
44 Unclear (5) _____
45
46
47
48
49

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Display This Question:

If Q27 = Interrupted time-series

Or Q27 = Difference-in-differences

Or Q27 = Comparative interrupted time-series

Q39 Given the above, do you believe that the design is appropriate for identifying the policy impact(s) of interest?

This should be taken as independent of what you believe about other studies, and/or the feasibility of other designs.

If not "Yes" please describe (three short sentences max).

Yes (1)

Mostly yes (2) _____

Mostly no (3) _____

No (4) _____

Unclear (5) _____

Display This Question:

If Q54 = Secondary/consensus round

Q55 General and/or additional comments on this paper from consensus discussion. This may include any additional information worth commenting on regarding the paper, difficulties encountered evaluating it, etc.

(three short sentences max)

End of Block: Main form

Review version with references removed; NOT FOR DISTRIBUTION

Policy evaluation in COVID-19: A guide to common design issues

Noah A Haber, Emma Clarke-Deelder, Joshua A Salomon, Avi Feller, Elizabeth A Stuart

Noah A Haber, ScD*

noahhaber@stanford.edu

Meta Research Innovation Center at Stanford University

Stanford University

1265 Welch Rd

Palo Alto, CA 94305

(650) 497-0811

Emma Clarke-Deelder, MPhil

Department of Global Health & Population

Harvard T. H. Chan School of Public Health

665 Huntington Avenue

Building 1, room 1104

Boston, Massachusetts 02115

Joshua A Salomon, PhD

Department of Medicine

Center for Health Policy and Center for Primary Care and Outcomes Research

Stanford University School of Medicine

Encina Commons, Room 118

615 Crothers Way

Stanford, CA 94305-6019

Avi Feller, PhD

Goldman School of Public Policy

University of California, Berkeley

2607 Hearst Avenue

Room 309

Berkeley, CA 94720

Elizabeth A Stuart, PhD

Department of Mental Health

Johns Hopkins Bloomberg School of Public Health

624 N. Broadway

Hampton House 839

Baltimore, MD 21205

* corresponding author

Review version with references removed; NOT FOR DISTRIBUTION

Abstract

Policy responses to COVID-19, particularly those related to non-pharmaceutical interventions, are unprecedented in scale and scope. Researchers and policymakers are striving to understand the impact of these policies on a variety of outcomes. Policy impact evaluations always require a complex combination of circumstance, study design, data, statistics, and analysis. Beyond the issues that are faced for any policy, evaluation of COVID-19 policies is complicated by additional challenges related to infectious disease dynamics and lags, lack of direct observation of key outcomes, and a multiplicity of interventions occurring on an accelerated time scale.

In this paper, we (1) introduce the basic suite of policy impact evaluation designs for observational data, including cross-sectional analyses, pre/post, interrupted time-series, and difference-in-differences analysis, (2) demonstrate key ways in which the requirements and assumptions underlying these designs are often violated in the context of COVID-19, and (3) provide decision-makers and reviewers a conceptual and graphical guide to identifying these key violations. The overall goal of this paper is to help policy-makers, journal editors, journalists, researchers, and other research consumers understand and weigh the strengths and limitations of evidence that is essential to decision-making.

Introduction

The response to the global COVID-19 pandemic has demanded urgent decision making in the face of substantial uncertainties. Policies to arrest transmission, including stay-at-home orders and other non-pharmaceutical interventions (NPIs), have wide-reaching consequences that touch many aspects of well being. Decision-making in the public interest requires evaluating and weighing the evidence on both intended and unintended consequences in order to best predict outcomes. The wide range of policy interventions implemented by different jurisdictions may yield opportunities for learning from what has already happened to inform future policymaking, and we have observed a proliferation of studies aimed at such policy evaluations. However, policy evaluation requires a complex combination of circumstance, data, study design, analysis, and interpretation in order to be informative.

Policy impact evaluation aims to answer questions about the extent to which the realized outcomes given a particular policy would have been different in the absence of that policy. Estimating the causal impact of the policy with observational data is challenging because what would have happened in the absence of the policy change (the “counterfactual”) is, by definition, unobserved. Randomized controlled trials (RCTs) of policies related to COVID-19 interventions may not always be practical or ethical. In this context, a large and growing number of studies have attempted to evaluate the impact of COVID-19 policies using observational data. There

Review version with references removed; NOT FOR DISTRIBUTION

are many potential pitfalls in the use of observational data for evaluation generally, and some additional methodological design challenges relating to COVID-19 policies in particular.

This paper provides a graphical guide to policy impact evaluations for COVID-19, targeted to decision-makers, researchers and evidence curators. Our aim is to provide a coherent framework for conceptualizing and identifying common pitfalls in COVID-19 policy evaluation. Importantly, this should not be taken either as a comprehensive guide to policy evaluation more broadly or as guidance on performing analysis, which may be found elsewhere. Rather, we review relevant study designs for policy evaluations — including pre/post, interrupted time series, and difference-in-difference approaches — and provide guidance and tools for identifying key issues with each type of study as they relate to NPIs and other COVID-19 policy interventions. Improving our ability to identify key pitfalls will enhance our ability to identify and produce valid and useful evidence for informing policymaking.

Common policy evaluation designs and their pitfalls in COVID-19

Identifying the type of design

Review version with references removed; NOT FOR DISTRIBUTION

Table 1: Summary definitions of policy impact evaluation designs commonly used for COVID-19

Design	Units (e.g., regions of comparison)		Time points measured per unit		Assumed counterfactual. “If not for the intervention, _____”
	With intervention	Without intervention	Before intervention	After intervention	
Cross-sectional	At least one	At least one	N/A	One time point	Outcome in intervention units would have been the same as the outcome in the non-intervention units.
Pre/post Figure 1A	At least one	None	At least one (typically one)	At least one (typically one)	Outcome would have stayed the same from the pre period to the post period.
Interrupted time-series (ITS) Figure 1B	At least one	None	More than one	At least one (typically several)	Outcome slope and level* would have continued along the same modelled trajectory from the pre-period to the post period.
Difference-in-differences (DiD) Figure 1C	At least one	At least one [†]	At least one (typically one)	At least one (typically one)	Outcome in intervention units would have changed as much as (or in parallel with) the outcome in the non-intervention units.
Comparative interrupted time series (CITS) Figure 1D	At least one	At least one [†]	More than one (typically several)	At least one (typically several)	Outcome slope and level* would have changed as much as non-intervention group's slope and level* changed.

* Assessing both slope and level only applicable if there are multiple data points during the post period
[†] Units without the intervention may be the pre-period of a different unit that eventually receives the intervention.

Identifying the underlying design in a given analysis often requires using a combination of the methods as reported and evaluating the data structure that is used for the main analysis, as shown in Table 1. COVID-19-related policy evaluation analyses typically fall under these categories. In most cases, the design can be categorized using a combination of whether there are also units that did not receive the treatment (columns 2-3) and whether there are time points both before and after intervention for those units (columns 4-5). The final column describes the implied counterfactual, discussed further in subsequent sections. Cross sectional designs typically compare units with vs without the treatment at single time points. Pre/post studies typically compare within units who received the intervention at two points: before and after a policy. Interrupted time-series analyses compare outcomes within units within units who received the intervention at greater than two time points before the intervention vs with at least one (typically multiple) after the intervention. Difference-in-differences analysis compares the outcome change in units which received the intervention with those that did not (or have not yet), with at least one point before and one after the intervention. In cases with multiple periods, that may involve a comparison with the pre-policy period of one region with the post-period of a different region, even though all regions eventually receive the intervention.

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

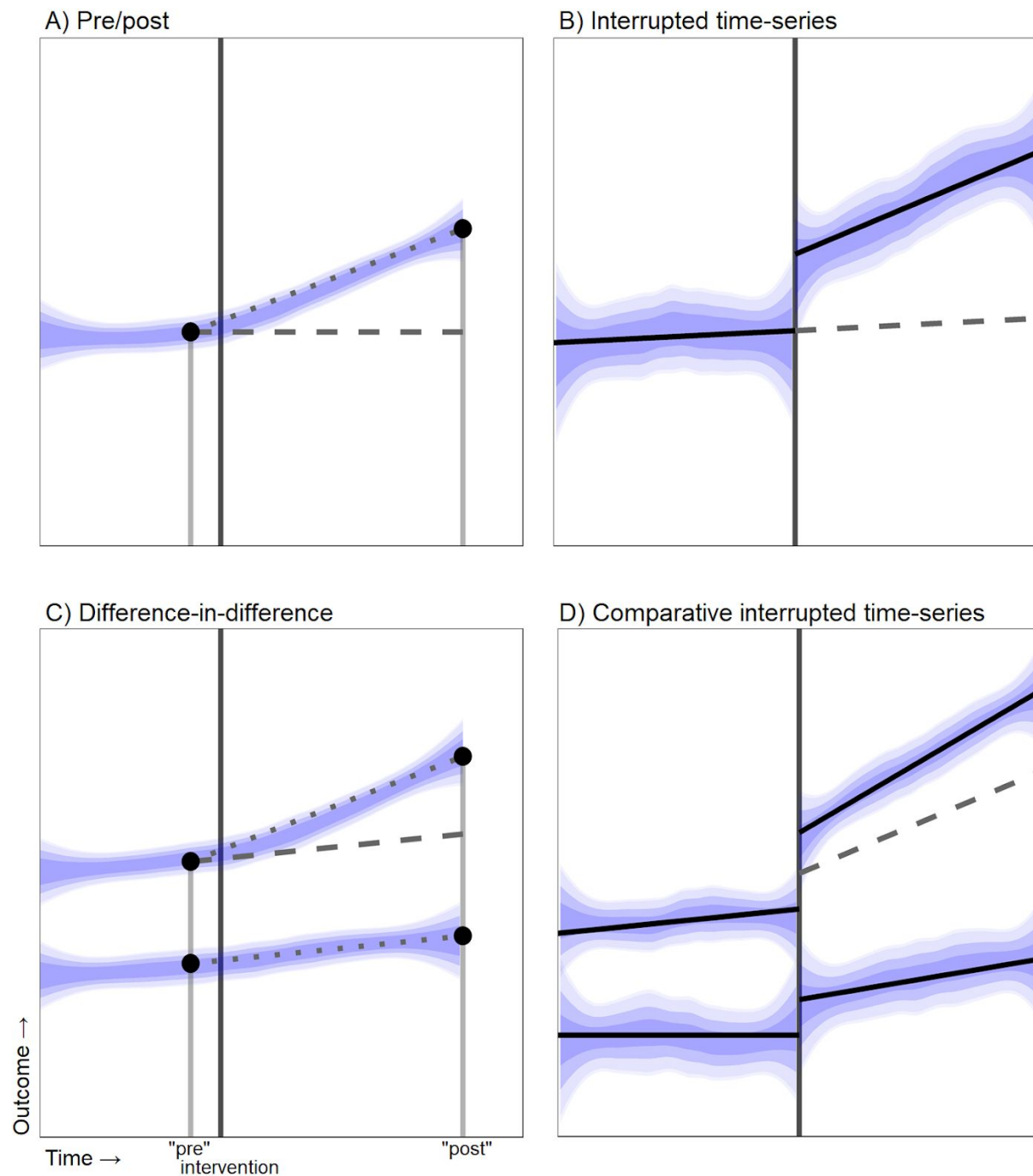
Methods descriptions may not always provide a precise or reliable guide to which of the design approaches has been used. Some studies do not explicitly name these designs (or may classify them differently); and these are only a small fraction of designs and frameworks that are possible to use for policy evaluation. Studies may have data at multiple time points but are effectively cross-sectional. DiD, ITS, and CITS designs based on repeated cross-sectional data are sometimes described as “cross-sectional” instead of longitudinal. The term “event study” is often used to refer to studies with a single unit and one change over time resembling ITS, but may refer to other designs. Although ITS is often used to describe changes in one unit, it may also refer to settings in which many treated units adopt an intervention over time. Studies will also frequently employ multiple designs, while others use more complex methods of generating counterfactuals. Definitions of these terms vary widely, and the definitions above should be considered as guidance only.

Policy impact evaluation design foundations for COVID-19

The simplest design is the cross-sectional analysis, which compares COVID-19 outcomes between units of observation (e.g., cities) at a single calendar time or time since an event, typically post-intervention. These studies are unlikely to be appropriate for COVID-19-related policy evaluations, but provide a useful starting point for reasoning about different designs. Just as with comparisons of non-randomized medical treatments, the localities that adopt a particular policy likely differ substantially from those that don't on both observed and unobserved characteristics on a number of dimensions, including epidemic status and timing.

Figure 1: Longitudinal designs overview

Review version with references removed; NOT FOR DISTRIBUTION



48 This chart shows four canonical longitudinal designs. In all cases: the blue shading
49 represents the underlying data trends, the solid vertical grey line represents the time of
50 intervention, the grey dashed lines represent the assumed counterfactual in the absence of
51 the intervention, as discussed in the text. The impact estimate is obtained by comparing the
52 outcomes observed for the treated unit in the post period (the solid line) with the implied
53 counterfactual line (the dashed line). In the case of the pre/post and
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 difference-in-differences panels the large black dots represent the time of measurement,
4 connected by the grey dotted lines.
5

6
7 Given the challenges in a simple cross-sectional comparison, which compare post-intervention
8 outcomes, it is important to consider longitudinal designs, which instead look at differences or
9 trends across time, as summarized in Figure 1. These can be distinguished by the data used
10 and the construction of the counterfactual. Pre/post, for example, has only one unit, measured
11 at two time points. Two common strategies expand on the logic and data requirements of the
12 pre/post design. Interrupted time series designs (Figure 1B) incorporate multiple time points
13 before the intervention, and usually multiple time points after the intervention, to enable a more
14 complete view on changes in levels and trends that are temporally related to the intervention.
15 Difference-in-difference designs (Figure 1C) add a set of comparison points from a group or
16 location that did not have the intervention. Another related design (comparative interrupted
17 time-series, Figure 1D, discussed only briefly here), uses both aspects — a change over time
18 and a comparison group — to compare the observed change in slopes for the intervention
19 group with the change in slope for the comparison group.
20
21
22
23

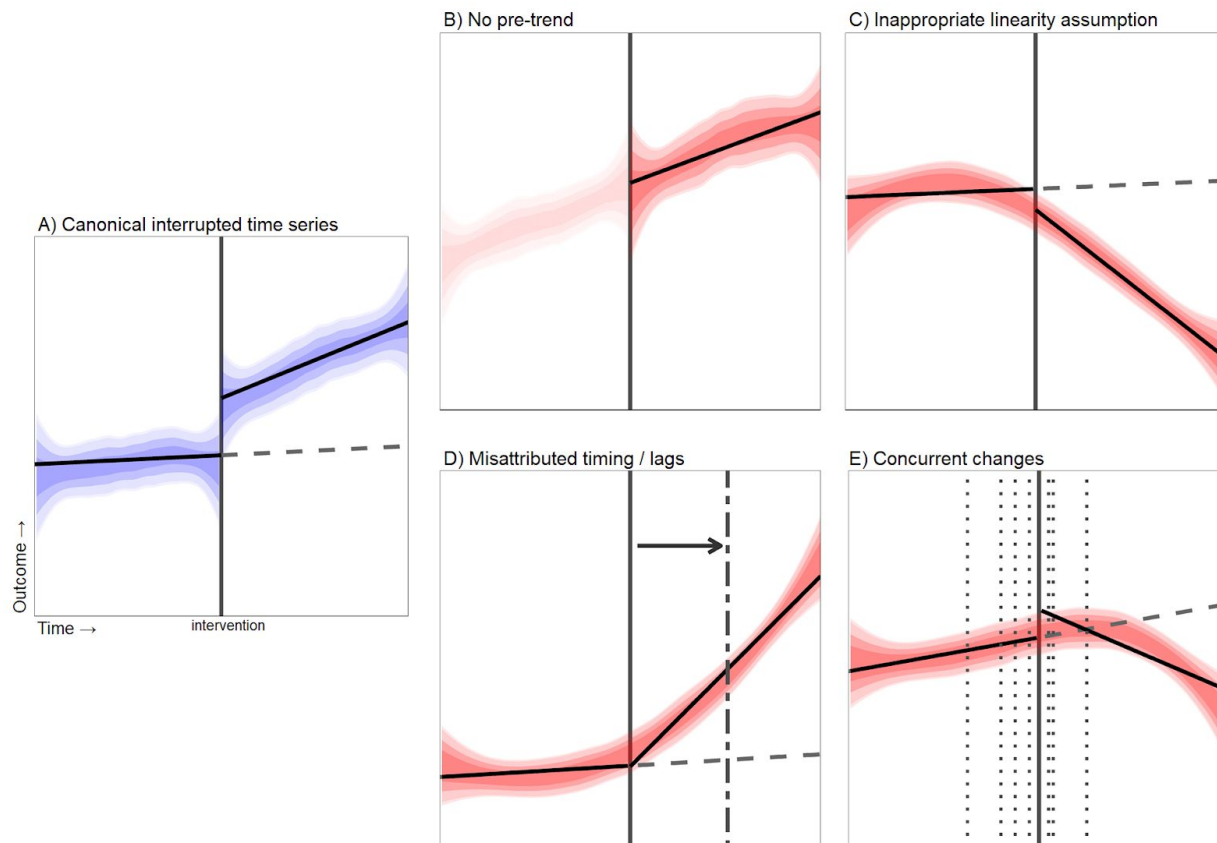
24 Pre/post studies

25
26 The simplest longitudinal design is a pre/post analysis, where some outcome is observed before
27 policy implementation, and again after, in a single group (Figure 1A). Pre/post studies are
28 analogous to a single arm trial with no control and only a single follow-up observation after
29 treatment.. This effectively imposes the assumption that the counterfactual trend is completely
30 flat (i.e., that the outcome in the post-period in the absence of the policy change is the same as
31 the value of the outcome before the policy change) without accounting for pre-existing
32 underlying trends, and attributing all outcome changes completely to the intervention of interest.
33 Just as the outcomes for an individual patient might be expected to change before and after
34 treatment, for reasons unrelated to the treatment, outcomes related to policy interventions will
35 change for reasons not caused by the policy. Infection rates, for example, would not be
36 expected to remain stationary except in very specific circumstances, but a pre/post
37 measurement would assume that any changes in infection rates are attributable to the policy.
38
39
40
41
42

43 Interrupted time-series

44
45 Figure 2: Interrupted time-series graphical guidance for identifying common pitfalls
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION



This chart shows one canonical design for ITS (blue, Panel A) and four panels demonstrating common issues with ITS analysis (red, panels B-E) discussed in the text. In all cases: the lag/red shading represents the underlying data trends, the vertical grey line represents the time of intervention, the grey dashed lines represent the assumed counterfactual in the absence of the intervention. In panel D) the dash-dot line represents the time at which the policy is expected to impact the outcome. In panel E), the vertical dotted lines represent concurrent events and changes.

Interrupted time-series (ITS) is a strategy that uses a projection of the pre-policy outcome trend as a counterfactual for how the outcome would have changed if the policy had not been introduced. In other words, in the absence of the policy change, ITS assumes the outcome would have continued on its pre-policy trend during the study period. ITS can be a useful tool in policy evaluation because it allows researchers to account for underlying trends in the outcome and, by comparing the treated unit (or location) to itself; it can therefore eliminate some of the confounding concerns that arise in cross-sectional or pre-post studies.

However, the validity of ITS depends critically on how well counterfactual trends in the outcome are modelled, and whether the policy of interest is the only relevant change during the study period. In the canonical setting (Figure 2A), the pre-policy trend is stable and can be feasibly modelled with the available data; the researcher appropriately models the timing of the change in the slope and/or level of the outcome; the researcher has sufficient information to conclude

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 that there were no other changes during the study period that would be expected to influence
4 the outcome. These elements are largely not satisfied in studies of COVID-related policy, as
5 described below.
6
7

8 ITS relies critically on modelled trends of the outcome over time. Key components of ITS
9 analyses include both visual and statistical examination of trends, preferentially alongside a
10 theoretical justification of the model used. At a minimum, analyses should provide graphical
11 representation of the data and model over time to examine whether pre-trend outcomes are
12 stable, all trends are well-fit to the data, “interrupted” at the appropriate time point, and sensibly
13 modelled (Figure 2B). In the case where an ITS includes a large number of units (e.g. states), it
14 can be difficult to display this information graphically.
15
16
17

18 One common pitfall in ITS is adoption of inappropriate assumptions on the outcome trend
19 (Figure 2C). The estimate of policy impact will be biased if a linear trend is assumed but the
20 outcome and response to interventions instead follow nonlinear trends (either before or after the
21 policy). In some cases, transformation of the outcome, for example using a log scale, may
22 improve the suitability of a linear model. Imposing linearity inappropriately is a serious risk in the
23 context of COVID-19, as trends in infectious disease dynamics are inherently non-linear. For
24 intuition, terms such as “exponential growth,” “flattening,” and “s-curves” all refer to non-linear
25 infectious disease trends. Depending on the particular situation, non-linearity or other modelled
26 trends can have complicated and counterintuitive impact on policy impact. Apparent linearity
27 may also be temporary and an artifact of testing, which may give a misleading impression that
28 linear models for infectious disease trends are appropriate indefinitely. While some use linear
29 progression in order to avoid more complex infectious disease models, in fact, linear projections
30 impose strict and often unrealistic models, generally resulting in an inappropriate counterfactual.
31
32
33
34
35

36 Researchers can easily misattribute the timing of the policy impact, resulting in spurious
37 inference and bias (Figure 2D). Some public health policies can be expected to translate into
38 immediate results (e.g., smoking bans and acute coronary events). In contrast, nearly every
39 outcome of interest in COVID-19 exhibits complex and difficult to infer time lags typically in the
40 realm of many weeks. The time between policy implementation and expected effect in the data
41 can be large and highly variable. For example, in order to see the impact of a mask order, first
42 the mask order takes effect, then people change their behaviors over time to comply with the
43 order (or sometimes the reverse in the case of anticipation effects), mask use behavior
44 produces changes in infections, then infections later result in symptoms, symptoms induce
45 people to seek testing, the tests must then be processed in labs, and then finally the results get
46 reported in data monitoring efforts. Selection of lead/lag time should be justifiable *a priori* or
47 external data. Selecting a lag based on the data risks issues comparable to p-hacking.
48
49
50
51

52 Finally, and perhaps most concerningly in the context of COVID-19, ITS fails when the policy of
53 interest coincides in time with other changes that affect the outcome (Figure 2E). For example, if
54 both mask and bar closure orders are rolled out together as a package, ITS cannot isolate the
55 impact of bar closures specifically. These changes do not need to have taken place exactly
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

concurrently with the policy implementation date of interest; they merely need to have some effect within the time period of measurement to result in potentially serious bias in effect estimates if unaddressed. ITS will also likely be biased if, during the study period, there is a change in the way the outcome data is collected or measured. This might occur if the introduction of a COVID-19 control policy is combined with an effort to collect better data on infection or mortality cases. Analogously, if an RCT involves randomizing people to a group receiving both A and B vs. control, we typically can't disentangle the effects of A from the effects of B, unless we also have separate A- and B-only arms. Ultimately, if multiple things are changing at the same time, ITS may not be an appropriate design for policy evaluation.

COVID-19 policies rarely arrive alone; they are typically created alongside other policies, unofficial action, and large scale behavior changes which themselves impact COVID-19-related outcomes. In some cases, anticipation of a policy may induce behavior change before the actual policy takes effect. The policies themselves may have been chosen due to the expectation of change in disease outcomes, which introduces additional biases related to “reverse” causality.

Table 2: Checklist for identifying common pitfalls for ITS to evaluate COVID-19 policy

Key design questions. If any answer is “no,” this analysis is unlikely to be appropriate or useful for estimating the impact of the intervention of interest.	Details and suggestions for identifying issues:
Does the analysis provide graphical representation of the outcome over time?	-Check for a chart that shows the outcome over time, with the dates of interest. Outcomes may be aggregated for clarity (e.g. means and CIs at discrete time points).
Is there sufficient pre-intervention data to characterize pre-trends in the data?	-Check the chart(s) to see if there are several time points over a reasonable period of time over which to establish stability and curvature in the pre-trends.
Is the pre-trend stable?	-Check if there are sufficient data to reasonably determine a stable functional form for the pre-trends, and that they follow a modelable functional form.
Is the functional form of the counterfactual (e.g. linear) well-justified and appropriate?	-Check whether the authors explain and justify their choice of functional form. -Check if there is any curvature in the pre-trend. -Consider the nature of the outcome. Is it sensible to measure the trend of this outcome on the scale and form used? Note: infectious disease dynamics are rarely linear. -Consider that while pre-trend fit is a necessary condition for an appropriate linear counterfactual model, it is not sufficient. Check if the authors provide justification for the functional form to continue to be of the same functional form (e.g. linear).
Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?	-Check whether the authors justify the use of the date threshold relative to the date of the intervention. -Trace the process between the intervention being put in place to when observable effects in the outcome might appear over time. -Consider whether there are anticipation effects (e.g. do people change behaviors before the date when the intervention begins?) -Consider whether there are lag effects. (e.g. does it take time for behaviors to change, behavior change to impact infections, infections to impact testing, and data to be collected, etc?)

Review version with references removed; NOT FOR DISTRIBUTION

	-Check if authors appropriately and directly account for these time effects.
Is this policy the only thing to happen which could have impacted the outcome during the measurement period, differently for policy and non-policy regions??	<ul style="list-style-type: none"> -Consider other policies or interventions which could impact the outcome during this time. -Consider social behaviors changed which could meaningfully impact the outcome during this time. -Consider economic conditions changed which could meaningfully impact the outcome during this time. -Note that the actual concurrent changes do not need to happen during the period of measurement, just their effects.

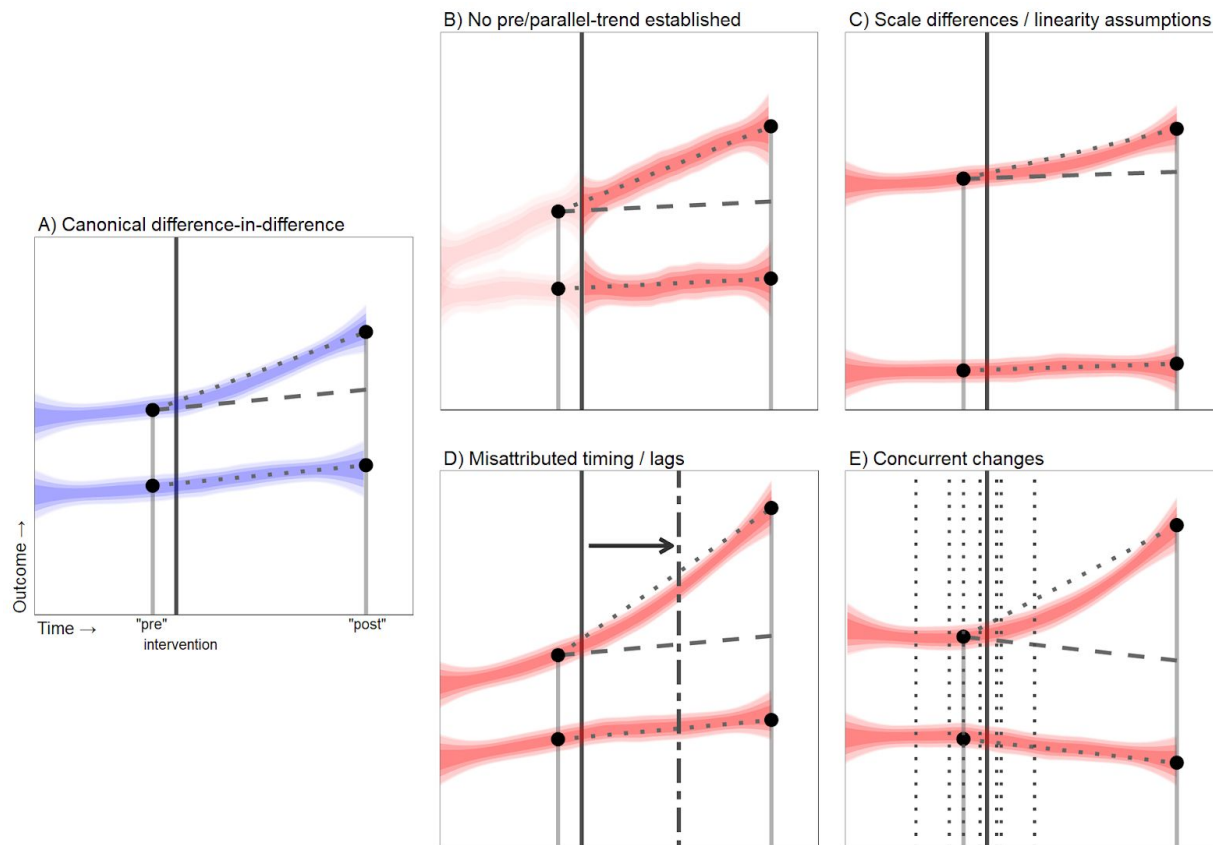
These issues are summarized as a checklist of questions to identify common pitfalls in Table 2.

Difference-in-differences

The difference-in-difference (DiD) approach uses concurrent non-intervention groups as a counterfactual. Typically, this consists of one set of units (e.g., regions) that had the intervention and one set that did not, with each measured before and after the intervention took place. DiD is more directly analogous to a non-randomized medical study with at least one treatment and control group but limited observation before and after treatment. In contrast to ITS, which compares a unit with itself over time, DiD compares differences between treatment arms or units at two observation points. In many analyses, a DiD approach is implied by comparing regions over time, without formally naming or modelling it. Other DiD approaches use interventions implemented at multiple time points.

Figure 3: Difference-in-differences graphical guidance for identifying common pitfalls

Review version with references removed; NOT FOR DISTRIBUTION



This chart shows one canonical design for DiD (blue, Panel A) and four panels demonstrating common issues with DiD analysis (red, panels B-E). In all cases: the blue/red shading represents the underlying data trends, the vertical grey line represents the time of intervention, the grey dashed lines represent the assumed counterfactual in the absence of the intervention. In panel D) the dash-dot line represents the time at which the policy is expected to impact the outcome. In panel E), the vertical dotted lines represent concurrent events and changes.

One key component of the standard DiD approach is the parallel counterfactual trends assumption: that the intervention and comparison groups would have had parallel trends over time in the absence of the intervention. In some cases, the parallel trends assumption may be referenced or examined implicitly but not named.

Ideally, pre-intervention trends would be shown to be clearly identifiable, stable, of a similar level, and parallel between groups. With only one observation before and only one after the intervention, assessment of the plausibility of the parallel counterfactual trends assumption is not possible. Absent this confirmation the evaluation runs the risk of biased estimation due to differential pre-trends (Figure 3B). Pre-trends approaching the ceiling or floor may also not be informative about stable and parallel pre-trends. Empirical assessment of whether pre-intervention trends were parallel and stable between groups is possible when multiple observations are available at multiple time points before the intervention, noting that this can

Review version with references removed; NOT FOR DISTRIBUTION

begin to resemble a CITS design. In this scenario, pre-trend data should be visually and statistically established and documented. While parallel trends before intervention (which we can observe and may be testable) do not guarantee parallel *counterfactual* trends in the post-intervention period (which we cannot observe and are generally untestable), examining pre-intervention parallel trends is a minimal requirement for DiD reliability.

It is also important to consider the scale and level on which the outcome is measured (Figure 3C). As with ITS, if the outcomes in the treatment and comparison groups are moving in parallel on a logged scale, they will not be moving in parallel on a natural scale. Level differences by themselves may be a problem for COVID-19 outcomes, as infectious disease transmission dynamics dictate that infection risks are related to the prevalence of infected people in a population, i.e. the rate of change is linked intrinsically to the level. A population with an extremely low prevalence will tend to have an inherently slower rise in infection rates than an otherwise identical population with merely a low prevalence. Just as importantly, large level differences in the outcome between intervention and comparison groups is often indicative of other important differences between comparators, which may result in other assumptions being violated.

While DiD is in some ways more robust to very specific kinds of timing effects (Figure 3D) and concurrent changes (Figure 3E), it also introduces additional risks. DiD effectively doubles the opportunity for concurrent changes to spuriously impact results, since they can occur in the treatment or comparison groups. As above, this can become even more problematic for DiD in the typical case where intervention groups enact more or very contextually different policies than non-intervention groups. Even cases where concurrent changes happen equally in both treatment and comparison groups can lead to overwhelming bias, particularly when approaching the maximum or minimum levels of the outcome. If either the treatment or control group is approaching the floor (e.g. 0% prevalence) or ceiling for an outcome of interest due to other policies concurrent in both places (e.g. national lockdowns, but region-level differences in mask policy), this can lead to bias when comparing changes between the two groups.

Table 3: Checklist for identifying common pitfalls for DiD to evaluate COVID-19 policy

Key design questions. If any answer is “no,” this analysis is unlikely to be appropriate or useful for estimating the impact of the intervention of interest	Details and suggestions for inspection:
Does the analysis provide graphical representation of the outcome over time?	-Check for a graph that shows the outcome over time for all groups, with the dates of interest. Outcomes may be aggregated for clarity (e.g. mean and CI at discrete time points).
Is there sufficient pre-intervention data to observe both pre and post trends in the data?	-Check the chart(s) to see if there are several time points over a reasonable period of time over which to establish stability and curvature in the pre- and post- trends.
Are the pre-trends stable?	-Check if there are sufficient graphical data to reasonably determine a stable functional form for the pre-trends, and that they follow a modelable functional form.

Review version with references removed; NOT FOR DISTRIBUTION

Are the pre-trends parallel?	-Observe if the trends in the intervention and comparison groups appear to move together at the same rate at the same time.
Are the pre-trends at a similar level?	-Check if the trends in the intervention and comparison groups are at similar levels. -Note that non-level trends exacerbates other problems with the analysis, including linearity assumptions
Are intervention and non-groups broadly comparable?	-Consider areas where comparison groups may be dissimilar for comparison beyond just the level of the outcome.
Is the date or time threshold set to the appropriate date or time (e.g. is there lag between the intervention and outcome)?	-Check whether the authors justify the use of the date threshold relative to the date of the intervention. -Trace the process between the intervention being put in place to when observable effects in the outcome might appear over time. -Consider whether there are anticipation effects (e.g. do people change behaviors before the date when the intervention begins?) -Consider whether there are lag effects. (e.g. does it take time for behaviors to change, behavior change to impact infections, infections to impact testing, and data to be collected, etc?) -Check if authors appropriately and directly account for these time effects.
Is this policy the only uncontrolled or unadjusted-for way in which the outcome could have changed during the measurement period, differently for policy and non-policy regions?	-Consider any uncontrolled factor which could have influenced the outcome differently in policy and non-policy regions. -This may include (but is not limited to) -Other policies -Social behaviors -Economic conditions -Are these factors justified as having negligible impact? -If justified, is the argument that these have negligible impact convincing? -Note that the actual concurrent changes do not need to happen during the period of measurement, just their effects.

Similarly to the ITS section, these issues are summarized as a checklist of questions to identify common pitfalls in Table 3.

Discussion

In recent months, there has been a proliferation of research evaluating policies related to the COVID-19 pandemic. As with other areas of COVID-19 research, quality has been highly variable, with low quality studies resulting in poorly or mis-informed policy decisions, wasted resources, and undermined trust in research. To support high quality policy evaluations, in this paper we describe common approaches to evaluating policies using observational data, and describe key issues that can arise in applying these approaches. We hope that this guidance can help support researchers, editors, reviewers, and decision-makers in conducting high quality policy evaluations and in assessing the strength of the evidence that has already been published.

Policy evaluation — far from a simple task in normal circumstances — is particularly challenging during a pandemic. Cross-sectional comparisons of states or countries are likely to be biased by selection into treatment: for example, countries with worse outbreaks may be more likely to

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 implement policies such as mask requirements. In analyses of changes over time – such as
4 single-unit studies using interrupted time-series or multi-unit comparisons using
5 difference-in-differences or comparative interrupted time-series – it may not be possible to parse
6 apart the effects of different policies implemented around the same time, such as mask
7 mandates paired with limits on social gatherings. Analyses of changes over time may also be
8 biased if disease or human behavioral dynamics are not modelled appropriately. This can be
9 challenging because case counts typically do not grow linearly and there is often a lag between
10 a policy change and a behavioral response.
11
12
13

14 This guidance should be considered minimal screening to identify low quality policy impact
15 evaluation in COVID-19, but is in no way sufficient to identify high quality evidence or
16 actionability. Decision-makers and researchers should pay particular attention to the relevance
17 of the intervention as it was evaluated to relevant decisions being made. The evaluated impact
18 of a program encouraging mask use through messages might not be informative about mask
19 requirement orders. Differences in level of aggregation may be important, such as ecological
20 fallacy arising from a situation in which areas with higher overall mask use have higher
21 transmission, but transmission is actually lower for individuals wearing masks. Policy impact
22 evaluation is only as useful as the question it asks, data it uses, and the way it is analyzed.
23 Problems with measurement, spillover effects, generalizability, changes in measurement
24 overtime (e.g. varying test availability), statistics, testing robustness to alternative assumptions,
25 and many issues can undermine an otherwise robust evaluation, and are not discussed here.
26
27
28
29

30 While this guidance is not comprehensive, it may help inform study designs not covered here.
31 Issues with comparative interrupted time-series and synthetic control methods, for example, are
32 broadly similar to the issues with difference-in-differences analyses we discuss here. Other
33 approaches may include adjustment and matching based observational causal inference
34 designs, instrumental variables and related quasi-experimental approaches, and randomized
35 controlled trials. Each has its own set of practical, ethical, and inferential limitations.
36
37
38

39 In the face of these challenges, we recommend careful scrutiny and attention to potential
40 sources of bias in COVID-19-related policy evaluations, but we remain optimistic about the
41 potential for robust evaluations to inform decision-making. Researchers and decision-makers
42 should triangulate across a large variety of approaches from theory to evidence, invest in better
43 data and more reliable and useful evidence wherever feasible, clearly acknowledge limitations
44 and potential sources of bias, and acknowledge when actionable evidence is not feasible. We
45 anticipate increasing opportunities for better examining policies moving forward, particularly if
46 policies and interventions are designed with policy impact evaluation and data collection in
47 mind.
48
49
50

51 The COVID-19 pandemic requires urgent decisions about policies that affect millions of people's
52 lives in significant ways. High quality evidence on the effects of these policies is critical to
53 informing decision-making, but is very hard to generate. Evidence-based decision-making
54
55
56
57
58
59
60

Review version with references removed; NOT FOR DISTRIBUTION

1
2
3 depends on research that carefully considers potential sources of bias, and clearly
4 communicates underlying assumptions and sources of uncertainty.
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



PRISMA 2009 Checklist

Section/topic	#	Checklist item	Reported on page #
TITLE			
Title	1	Identify the report as a systematic review, meta-analysis, or both.	1
ABSTRACT			
Structured summary	2	Provide a structured summary including, as applicable: background; objectives; data sources; study eligibility criteria, participants, and interventions; study appraisal and synthesis methods; results; limitations; conclusions and implications of key findings; systematic review registration number.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known.	3
Objectives	4	Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS).	3
METHODS			
Protocol and registration	5	Indicate if a review protocol exists, if and where it can be accessed (e.g., Web address), and, if available, provide registration information including registration number.	4
Eligibility criteria	6	Specify study characteristics (e.g., PICOS, length of follow-up) and report characteristics (e.g., years considered, language, publication status) used as criteria for eligibility, giving rationale.	4
Information sources	7	Describe all information sources (e.g., databases with dates of coverage, contact with study authors to identify additional studies) in the search and date last searched.	5
Search	8	Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated.	5 + appendix 1
Study selection	9	State the process for selecting studies (i.e., screening, eligibility, included in systematic review, and, if applicable, included in the meta-analysis).	5
Data collection process	10	Describe method of data extraction from reports (e.g., piloted forms, independently, in duplicate) and any processes for obtaining and confirming data from investigators.	5
Data items	11	List and define all variables for which data were sought (e.g., PICOS, funding sources) and any assumptions and simplifications made.	6-7
Risk of bias in individual studies	12	Describe methods used for assessing risk of bias of individual studies (including specification of whether this was done at the study or outcome level), and how this information is to be used in any data synthesis.	6-7
Summary measures	13	State the principal summary measures (e.g., risk ratio, difference in means).	8



PRISMA 2009 Checklist

Synthesis of results	14	Describe the methods of handling data and combining results of studies, if done, including measures of consistency (e.g., I^2) for each meta-analysis.	8
----------------------	----	---	---

Page 1 of 2

Section/topic	#	Checklist item	Reported on page #
Risk of bias across studies	15	Specify any assessment of risk of bias that may affect the cumulative evidence (e.g., publication bias, selective reporting within studies).	8
Additional analyses	16	Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified.	8
RESULTS			
Study selection	17	Give numbers of studies screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally with a flow diagram.	8
Study characteristics	18	For each study, present characteristics for which data were extracted (e.g., study size, PICOS, follow-up period) and provide the citations.	9
Risk of bias within studies	19	Present data on risk of bias of each study and, if available, any outcome level assessment (see item 12).	9-13
Results of individual studies	20	For all outcomes considered (benefits or harms), present, for each study: (a) simple summary data for each intervention group (b) effect estimates and confidence intervals, ideally with a forest plot.	9-11
Synthesis of results	21	Present results of each meta-analysis done, including confidence intervals and measures of consistency.	10
Risk of bias across studies	22	Present results of any assessment of risk of bias across studies (see Item 15).	9-12
Additional analysis	23	Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]).	10-12
DISCUSSION			
Summary of evidence	24	Summarize the main findings including the strength of evidence for each main outcome; consider their relevance to key groups (e.g., healthcare providers, users, and policy makers).	12
Limitations	25	Discuss limitations at study and outcome level (e.g., risk of bias), and at review-level (e.g., incomplete retrieval of identified research, reporting bias).	13
Conclusions	26	Provide a general interpretation of the results in the context of other evidence, and implications for future research.	14
FUNDING			
Funding	27	Describe sources of funding for the systematic review and other support (e.g., supply of data); role of funders for the systematic review.	23

From: Moher D, Liberati A, Tetzlaff J, Altman DG, The PRISMA Group (2009). Preferred Reporting Items for Systematic Reviews and Meta-Analyses: The PRISMA Statement. PLoS Med 6(7): e1000097. doi:10.1371/journal.pmed1000097

For peer review only - <http://bmjopen.bmj.com/site/about/guidelines.xhtml>
For more information, visit: www.prisma-statement.org



PRISMA 2009 Checklist

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47

For peer review only