

Causal Panel Analysis under Parallel Trends: Lessons from a Large Reanalysis Study

ALBERT CHIU *Stanford University, United States*

XINGCHEN LAN *New York University, United States*

ZIYI LIU *University of California, Berkeley, United States*

YIQING XU *Stanford University, United States*

Two-way fixed effects (TWFE) models are widely used in political science to establish causality, but recent methodological discussions highlight their limitations under heterogeneous treatment effects (HTE) and violations of the parallel trends (PT) assumption. This growing literature has introduced numerous new estimators and procedures, causing confusion among researchers about the reliability of existing results and best practices. To address these concerns, we replicated and reanalyzed 49 studies from leading journals that employ TWFE models for causal inference using observational panel data with binary treatments. Using six HTE-robust estimators, diagnostic tests, and sensitivity analyses, we find: (i) HTE-robust estimators yield qualitatively similar but highly variable results; (ii) while a few studies show clear signs of PT violations, many lack evidence to support this assumption; and (iii) many studies are underpowered when accounting for HTE and potential PT violations. We emphasize the importance of strong research designs and rigorous validation of key identifying assumptions.

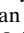
Over the past decade, political scientists have increasingly relied on panel data to draw causal conclusions (Xu 2023). A favored method for such analyses is the two-way fixed effects (TWFE) model because of its ability to control for unobserved time-invariant confounders and common time trends. In our survey of 102 articles published from 2017 to 2023 in three top political science journals using observational panel data with *binary* treatments, 64 studies (63%) assume a TWFE model with the following functional form or a close variant:¹


$$Y_{i,t} = \tau^{TWFE} D_{i,t} + X'_{i,t} \beta + \alpha_i + \zeta_t + \varepsilon_{i,t}, \quad \text{for all } i, t, \quad (1)$$

where $Y_{i,t}$ and $D_{i,t}$ are the outcome and treatment variables for unit i at time t ; $X_{i,t}$ is a vector of time-varying covariates; α_i and ζ_t are unit and time fixed effects; and $\varepsilon_{i,t}$ is idiosyncratic errors.² Researchers typically interpret τ^{TWFE} as the treatment effect and estimate the model using ordinary least squares. The resulting estimator for τ^{TWFE} is commonly known as the TWFE estimator. Moreover, researchers frequently conflate this model with a difference-in-differences (DID) design, and use the two terms interchangeably.³

Recent methodological discussions have raised concerns about the validity of TWFE models and the associated identifying assumptions, leaving many researchers in a quandary. First, existing findings based on the TWFE models may not hold given recent developments. Second, with the introduction of numerous new estimators and diagnostics, there is confusion about the current best practices. This article seeks to bridge this gap by

Albert Chiu , PhD student, Department of Political Science, Stanford University, United States, altchiu@stanford.edu.

Xingchen Lan , PhD student, Wilf Family Department of Politics, New York University, United States, xingchenlan@nyu.edu.

Ziyi Liu , PhD student, Haas School of Business, University of California, Berkeley, United States, zyliu2023@berkeley.edu.

Corresponding author: Yiqing Xu , Assistant Professor, Department of Political Science, Stanford University, United States, yiqingxu@stanford.edu.

Handling editor: Andrew Eggers.

Received: June 18, 2024; revised: November 29, 2024; accepted: April 17, 2025.

¹ The remaining 38 studies can be categorized into five groups: studies focusing on interaction effects (8 studies), studies using non-linear links such as logit and Poisson (5 studies), studies employing instrumental variables or regression discontinuity designs (8 studies), and studies using other linear specifications, such as only one-way fixed effects or lagged dependent variables (17 studies).

² In some studies classified as using TWFE models, “unit” fixed effects are specified at the group level g , where multiple units i are nested (e.g., county fixed effects when i indexes cities), or time fixed effects are at a higher level p (e.g., year fixed effects when t indexes days). For simplicity, we use the notation α_i and ζ_t rather than the more general α_g and ζ_p .

³ We use the phrase “DID design” in reference to DID research design, which differs from the typical usage in the statistics literature that refers to treatment assignment mechanism (Xu, Zhao, and Ding 2024).

reviewing new estimation, inference, and diagnostics methods from the methodological literature and by re-analyzing published studies using both new estimators and the TWFE estimator. Based on the findings, we offer several practical recommendations for researchers.

These criticisms of the use of TWFE models mainly come from two directions. First, causal identification using TWFE models requires the *strict exogeneity* assumption, which critics argue is stronger than many researchers realize and is often unrealistic in real-world settings (e.g., Imai and Kim 2019). Strict exogeneity states that

(Strict exogeneity)

$$\mathbb{E}[\varepsilon_{i,t} | \mathbf{D}_i, \mathbf{X}_i, \alpha_i, \zeta_t] = \mathbb{E}[\varepsilon_{i,t} | D_{i,t}, X_{i,t}, \alpha_i, \zeta_t] = 0, \quad \forall i, t,$$

in which $\mathbf{D}_i = \{D_{i,1}, D_{i,2}, \dots, D_{i,T}\}$ and $\mathbf{X}_i = \{X_{i,1}, X_{i,2}, \dots, X_{i,T}\}$. It means that once current treatment status, covariates, and fixed effects are accounted for, treatment status in any other periods has no additional effect on $Y_{i,t}$ (Wooldridge 2010, 253). Under Equation 1, strict exogeneity implies a parallel trends (PT) assumption:

(Parallel trends)

$$\begin{aligned} \mathbb{E}[Y_{i,t}(0) - Y_{i,s}(0) | D_{i,t} = 1, D_{i,s} = 0, X_{i,t} - X_{i,s} = x_0] \\ = \mathbb{E}[Y_{j,t}(0) - Y_{j,s}(0) | D_{j,t} = 0, D_{j,s} = 0, X_{j,t} - X_{j,s} = x_0], \end{aligned}$$

in which $Y_{i,t}(0) = Y_{i,t}(d_{i,t} = 0)$ represents the untreated potential outcome for unit i at time t . It states that the change in untreated potential outcomes between any two periods is mean independent of the change in observed treatment status during those periods, once changes in covariate values are controlled for. Threats to PT, such as the presence of time-varying confounders and feedback from past outcomes to current treatment assignment, also invalidate strict exogeneity. Therefore, throughout the rest of the article, we use the term “PT violations” to encompass violations of strict exogeneity.⁴

The second group of criticisms concerns the consequences of heterogeneous treatment effects (HTE), that is, τ^{TWFE} is not a constant (e.g., Athey and Imbens 2022; Borusyak, Jaravel, and Spiess 2024; Callaway and Sant’Anna 2021; de Chaisemartin and D’Haultfœuille 2020; Goodman-Bacon 2021; Strezhnev 2018; Sun and Abraham 2021). Researchers have shown that, under HTE, TWFE estimates in general do not converge to a convex combination of the individual treatment effects for observations under the treatment condition, even when the PT assumption is valid. The so-called “negative weighting” problem, as described in de Chaisemartin and D’Haultfœuille (2020), is an alarming theoretical result because it implies that a TWFE estimand can be negative (positive) even when all individual treatment effects are positive (negative). To address this issue, researchers have

proposed many new estimators that are “HTE-robust”—that is, estimators that converge to some convex combinations of individual treatment effects under their identifying assumptions.

This article thus pursues two goals. First, we explain and compare six recent proposals to amend TWFE models, including the interaction weighted (IW) estimator (Sun and Abraham 2021), stacked DID (Cengiz et al. 2019), CSDID (Callaway and Sant’Anna 2021), DID multiple (de Chaisemartin and D’Haultfœuille 2020; 2024), PanelMatch (Imai, Kim, and Wang 2023, hereafter IKW 2023), and the imputation method (Borusyak, Jaravel, and Spiess 2024, hereafter BJS 2024; Liu, Wang, and Xu 2024, hereafter LWX 2024). These estimators produce causally interpretable estimates under HTE and PT (or its variants). Second, we replicate and re-analyze 49 studies published in the *American Political Science Review* (APSR), *American Journal of Political Science* (AJPS), and *The Journal of Politics* (JOP) from 2017 to 2023 which rely on a TWFE model to draw causal conclusions.⁵ Our aim is to assess the consequences of using or not using HTE-robust estimators and shed light on the severity of PT violations in political science research.

Our reanalysis shows that, in most studies, the HTE-robust estimators yield qualitatively similar estimates to TWFE models. However, there is considerable variation in how closely these estimators align with TWFE. In three cases, at least one HTE-robust estimator produces an estimate with an opposite sign to the TWFE estimate; in one of these cases, the opposite-sign estimate is also statistically significant at the 5% level. There is also a more widespread problem of power: HTE-robust estimators tend to have larger measures of uncertainty, which, combined with even small fluctuations in point estimates, can weaken statistical confidence. This is especially relevant for results that originally teeter on the brink of significance.

The primary concern, however, is the validity of the PT assumption. While only a few studies show clear signs of PT violations, which likely lead to spurious findings, most studies lack the power to rule out that realistic PT violations could explain a nonzero estimated causal effect. In such cases, even mild PT violations (informed by pre-treatment estimates) prevent us from concluding that the original treatment effect is nonzero. This does not mean that these studies are wrong; rather, it indicates that the available data do not have sufficient power to reject the null hypothesis of no effects when the PT assumption is not perfectly met.

Overall, we find that a small minority of the studies in our sample meet our criteria of being highly credible. In these studies, we can statistically distinguish the treatment effect from zero using an HTE-robust estimator, even when allowing for mild PT violations benchmarked against placebo estimates using pre-treatment data. We recognize this as a high standard, as most researchers do not account for the power needed for

⁴ We discuss the relationship between strict exogeneity and PT, as well as other assumptions, under Equation 1 in Section A.1 of the Supplementary Material (SM). Note that the PT assumption invoked by many HTE-robust estimators does not depend on Equation 1.

⁵ Replication materials for this article are available for download at the APSR Harvard Dataverse (Chiu et al. 2025).

the sensitivity analysis we perform. To be clear, our intent is not to criticize the authors of the studies we have replicated, since many of the methods we used were not available at the time their studies were conducted. Our goal is to guide and improve future research.

In light of these findings, we urge researchers to prioritize a strong research design and sufficient power in causal panel studies. Credible observational studies should feature a well-defined treatment-outcome pair, shock-induced variation in treatment assignment, and sufficient power to ensure results are not undermined by small perturbations of key identifying assumptions. Research design has often been overlooked in causal panel analyses, likely because researchers have become accustomed to accepting the strong parametric and exogeneity assumptions behind TWFE models. Recent studies have emphasized the importance of (quasi-) randomness in treatment assignment for the robustness of DID findings (e.g., Roth and Sant'Anna 2023).

This article makes several contributions. First, we propose a typology of various estimators for causal panel analysis. Our typology is based on the settings in which an estimator can be used and how controls are chosen. We also provide a comprehensive comparison of these estimators and show how several proposals are equivalent in some circumstances. We hope this discussion helps researchers deepen their understanding of these estimators. Second, we adapt the robust confidence set approach for sensitivity analysis proposed by Rambachan and Roth (2023) to the setting of imputation estimators. We find it highly useful as it avoids the issue of conditional inference—where hypothesis testing conditional on passing a pretest (e.g., a pretrend test) can distort estimation and inference (Roth 2022). Third, our reanalysis instills confidence in existing political science research that uses TWFE models correctly while also cautioning against potential risks, such as the failure of the PT assumption and insufficient power. Based on these findings, we provide recommendations to improve practices, including the choice of estimators and the use of proper diagnostics. Finally, we contribute to the ongoing conversation on replication and reproducibility in political science (e.g., Eggers et al. 2015; Lall 2016; Hainmueller, Mummolo, and Xu 2019; Lal et al. 2024).

Our work is closely related to Baker, Larcker, and Wang (2022), who evaluate the credibility of a handful of studies with staggered adoption treatments in finance and accounting. It differs in that: (i) we use a wider range of estimators and diagnostic tests on a larger and more diverse set of empirical applications, many of which involve treatment reversals; (ii) our review suggests that while the weighting issue under HTE is important, the main threats to causal inference with panel data are PT violations and insufficient power. Our work also relates to Roth et al. (2023), Xu (2023), de Chaisemartin and D'Haultfoeuille (2023), Arkhangelsky and Imbens (2024), and Baker et al. (2025), who review and synthesize the recent methodological advancements in the DID literature. What sets this article apart is our application of these innovations to data, allowing us to evaluate the practical relevance of the theoretical critiques.

This research has a few limitations. First, we do not examine methods based on sequential ignorability, an alternative identification framework that assumes no unobserved confounding but allows for dynamic treatment selection up to the current time period. Second, our analysis does not encompass studies that use continuous treatments, which is common in political science research. Finally, as the methodological literature continues to evolve rapidly, our recommendations should be regarded as reflecting current best practices.

TWFE AND ITS PITFALLS

In this section, we review the pitfalls of TWFE models identified in the literature. In the classic two-group and two-period case, the TWFE estimator $\hat{\tau}^{TWFE}$ is equivalent to the DID estimator, which consistently estimates the average treatment effect on the treated (ATT) under no anticipation and PT even with HTE. These results do not hold more generally in more complex settings with differential treatment adoption times (known as staggered adoption) or treatment reversal, as we will discuss below.

Our survey of the top journals reveals that the TWFE model under Equation 1 is the most commonly adopted approach for estimating causal effects using panel data in political science. Fixed effects models began their rise to prominence in political science in the early 2000s, and criticism promptly followed. In a debate with Green, Kim, and Yoon (2001), Beck and Katz (2001) and King (2001) argue that linear fixed effects models often lead to misleading findings because they throw away valuable information, ignore rich temporal dynamics, and are incapable of capturing complex time-varying heterogeneities. Moreover, since both treatment and outcome variables are often serially correlated in panel data, researchers have advised against using standard error (SE) estimators suitable for cross-sectional data, such as Huber-White robust SEs (Bertrand, Duflo, and Mullainathan 2004). Scholars also recommend using bootstrap procedures to more effectively control Type I error rates when the number of clusters (units) is small (Cameron, Gelbach, and Miller 2008).

Recent Criticisms

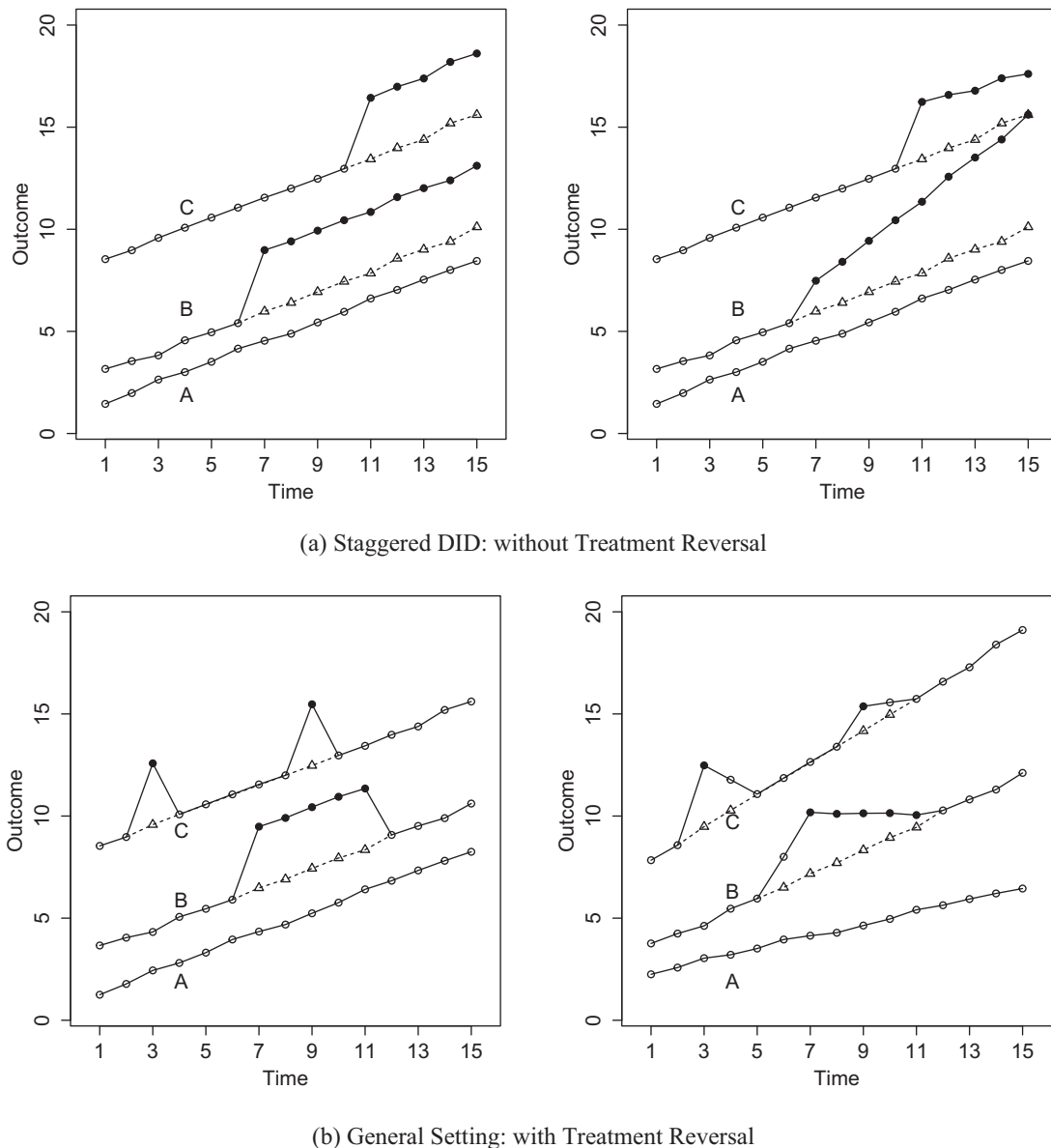
In the past few years, a surge of studies has renewed investigation into the properties of the TWFE estimator and the assumptions it requires to achieve causal identification. One group of work urges researchers to better understand TWFE models (and their assumptions) from a design-based perspective, with a focus on restrictions on treatment assignment mechanisms. For example, Imai and Kim (2019) point out that the strict exogeneity assumption not only implies the well-known no time-varying confounder requirement, but it also forbids a “feedback” effect from past outcomes to treatment assignment. Blackwell and Glynn (2018) clarify that such an assumption is closely

related to baseline randomization in which the treatment vector is generated prior to, or independent of, the realization of the outcome. Another body of work cautions researchers that the PT assumption, even in a simple 2×2 DID setting, is sensitive to functional form. For example, Kahn-Lang and Lang (2020) emphasize the implicit functional form restrictions imposed by PT, encouraging researchers to justify it from a (quasi-) experimental perspective and address pre-treatment covariate imbalance. Roth and Sant'Anna (2023) point out that strong assumptions are needed for PT to be scale-independent, ensuring that a monotonic

transformation of the outcome does not invalidate it. In practice, we find that many political science studies do not provide strong justification for strict exogeneity or PT.

A second body of research explores the implications of HTE with binary treatments within TWFE models (e.g., Athey and Imbens 2022; Borusyak, Jaravel, and Spiess 2024; Callaway and Sant'Anna 2021; de Chaisemartin and D'Haultfœuille 2020; Goodman-Bacon 2021; Strezhnev 2018). Most of this literature assumes staggered adoption, but the insights from that setting are still relevant when there are treatment reversals. In Figure 1,

FIGURE 1. Toy Examples: TWFE Assumptions Satisfied vs. Violated



Note: The above panels show outcome trajectories of units in a staggered adoption setting (a) and in a general setting (b). Solid and hollow circles represent observed outcomes under the treatment and control conditions during the current period, respectively, while triangles represent counterfactual outcomes (in the absence of the treatment across all periods), $Y_{i,t}(\mathbf{d}_i = \mathbf{0})$. The data on the *left* panels in both (a) and (b) are generated by DGPs that satisfy TWFE assumptions while the data on the *right* are not. The divergence between hollow circles and triangles in the right panel of (b), both of which are under the control condition, is caused by anticipation or carryover effects.

we present two simplified examples for the staggered adoption and general settings. The left panel of Figure 1a represents outcome trajectories in line with standard TWFE assumptions, which not only include PT but also require that the treatment effect be contemporaneous and unvarying across units and over time. In contrast, the right panel portrays a scenario where PT holds, but the constant treatment effect assumption is not met. Various decompositions by the aforementioned researchers reveal that even under PT, when treatments begin at different times (such as in staggered adoption) and treatment effects evolve over time and vary across units, the TWFE estimand is generally not a convex combination of the individual treatment effects for observations subjected to the treatment. The basic intuition behind this theoretical result is that TWFE models use post-treatment data from units who adopt treatment earlier in the panel as controls for those who adopt the treatment later (e.g., Goodman-Bacon 2021). HTE-robust estimators capitalize on this insight by avoiding these “invalid” comparisons between two treated observations.

A third limitation of the canonical TWFE specification is its presumption of no temporal and spatial interference. In most uses of TWFE models, researchers assume that there are no spatial spillovers and that treatment effects occur contemporaneously, hence no anticipation or carryover effects. No anticipation effects means that future treatments do not affect today’s potential outcomes, while no carryover effects means that today’s treatment does not affect future potential outcomes: For any i, t ,

(No anticipation effects)

$$Y_{i,t}(\mathbf{d}_i) = Y_{i,t}(d_{i,1}, d_{i,2}, \dots, d_{i,t}),$$

(No carryover effects)

$$Y_{i,t}(\mathbf{d}_i) = Y_{i,t}(d_{i,t}, d_{i,t+1}, \dots, d_{i,T}),$$

in which $\mathbf{d}_i = (d_{i,1}, d_{i,2}, \dots, d_{i,T})$ is unit i ’s vector of potential treatment conditions in all periods (from period 1 to T) and $Y_{i,t}(\mathbf{d}_i)$ is the potential outcome in period t given \mathbf{d}_i . The TWFE model specified in Equation 1 satisfies these two assumptions because $Y_{i,t}(\mathbf{d}_i) = Y_{i,t}(d_{i,t})$. These assumptions are obviously very strong, but they are rarely questioned or tested in practice (Athey and Imbens 2022; Imai and Kim 2019; Wang 2021). Although some recent methods permit arbitrary carryover effects in staggered adoption settings, these effects are not distinguishable from contemporaneous effects.⁶ This limitation becomes more complex when treatment reversal is possible, as demonstrated in Figure 1b. In Figure 1b, data in the left panel are consistent with TWFE assumptions, while the right panel illustrates deviations from PT, constant treatment effect, and the absence of anticipation or carryover effects. Real-world data often encounter the challenging scenarios depicted in the right panel

rather than the idealized conditions in the left. Scholars have proposed methods to handle limited carryover effects in the general setting (IKW 2023; LWX 2024). The challenge of addressing spatial spillover effects without strong structural assumptions still persists (Wang et al. 2025; Wang 2021), but its resolution is beyond the scope of this article.

Causal Estimands

To define the estimands clearly, consider the panel setting where multiple units $i \in \{1, \dots, N\}$ are observed at each time period $t \in \{1, \dots, T\}$. Each unit-time pair (i, t) uniquely identifies an observation. Define $E_{i,t}$ as unit i ’s “event time” at time t . For each i , let $E_{i,t} = \max\{t' : t' \leq t, D_{i,t'} = 1, D_{i,t'-1} = 0\}$ if $\exists s \leq t : D_{i,s} = 1$, and $E_{i,t} = \min\{t' : D_{i,t'} = 1, D_{i,t'-1} = 0\}$ otherwise. That is to say, $E_{i,t}$ is the most recent time at which unit i switched into treatment or, if i has not yet been treated at any point up until time t , the first time i switches into treatment. If i is never treated, we let $E_{i,t} = \infty$. In the staggered setting, the event time for each unit is constant, $E_i = E_{i,t}$, and $D_{i,t} = \mathbf{1}\{t \geq E_{i,t}\}$, where $\mathbf{1}\{\cdot\}$ is the indicator function. In such settings, we can partition units into distinct “cohorts” $g \in \{1, \dots, G\}$ according to the timing of treatment adoption E_i . Units transitioning to treatment at period g ($i : E_{i,t} = g$) form cohort g , whereas units that never undergo treatment belong to the “never-treated” cohort ($i : E_{i,t} = \infty$). $Z_{i,t}$ ($Z_{i,g,t}$) represents the variable Z for unit i (part of cohort g) at time t . We use $Y_{i,t}(1)$ and $Y_{i,t}(0)$ to denote the potential outcomes under treatment and control, respectively, and $Y_{i,t} = D_{i,t}Y_{i,t}(1) + (1-D_{i,t})Y_{i,t}(0)$ to denote the observed outcome.⁷

The finest estimand is the individual treatment effect, $\tau_{i,t} = Y_{i,t}(1) - Y_{i,t}(0)$, of which there exists one for each observation (i, t) .⁸ Most political science research, however, typically focuses on estimating a single summary statistic. Commonly, this is the ATT, which represents individual treatment effects averaged over all observations exposed to the treatment condition. In between these extremes of granularity and coarseness are time-varying dynamic treatment effects, which are across-unit averages of individual treatment effects at each time period relative to treatment adoption. In the staggered adoption setting, we can further subdivide by cohort. We use τ_l ($\tau_{g,l}$) to denote the dynamic treatment effect l periods after treatment adoption (for treatment

⁷ The current notation will not cause confusion because we do not allow feedback or temporal spillover. In some of the studies we refer to, potential outcomes are defined in terms of treatment history, as opposed to current treatment status. We adopt similar notations for these frameworks. For instance, we use $Y_{i,t}(D_{i,t} = 1, \{D_{i,s}\}_{s < t} = 0)$ to represent the potential outcome under the specified treatment history.

⁸ This is without loss of generality when feedback and interference are excluded. In staggered DID designs, carryover effects are permissible. When potential outcomes are defined in terms of treatment history, $\tau_{i,t}$ is defined as $Y_{i,t}(1) - Y_{i,t}(\infty)$, where $Y_{i,t}(\infty)$ signifies the untreated potential outcome when unit i never undergoes treatment.

⁶ See Section A.4 of the SM in Liu, Wang and Xu (2024) for more details.

cohort g), with $l = 1$ representing the period immediately after treatment adoption.⁹ $\tau_{g,l}$ is also what some authors refer to as a cohort average treatment effect on the treated (Strezhnev 2018; Sun and Abraham 2021) or group-time average treatment effect (Callaway and Sant'Anna 2021).

Each of the estimators we discuss can be used to estimate τ_l . The outcome model analogous to TWFE for estimating dynamic effects is a lags-and-leads specification. For simplicity, we first describe the staggered setting. Let $K_{i,t} = (t - E_{i,t} + 1)$ be the number of periods until (when $K_{i,t} \leq 0$) or since unit i 's event time at time t (e.g., $K_{i,t} = 1$ if unit i switches into treatment at time t). Consider a regression based on the following specification:

$$Y_{i,t} = \alpha_i + \zeta_t + X'_{i,t}\beta + \sum_{\substack{l=-a \\ l \neq 0}}^b \tau_l^{TWFE} \cdot \mathbf{1}\{K_{i,t} = l\} + \tau_{b+}^{TWFE} \mathbf{1}\{K_{i,t} > b\} \cdot D_{i,t} + \varepsilon_{i,t}, \quad (2)$$

where a and b are the number of lag and lead terms (BJS 2024). In the social science literature, the typical practice is to exclude $l = 0$, which corresponds to the time period immediately before the transition into the treatment phase, and use it as a reference period (Roth 2022). Conventionally, $\hat{\tau}_l^{TWFE}$ is interpreted as an estimate of τ_l or as a meaningful weighted average of pertinent individual treatment effects. Meanwhile, $\hat{\tau}_{b+}^{TWFE}$ is viewed as an estimate for the long-term effect.

HTE-ROBUST ESTIMATORS

In this section, we offer a brief overview and comparison of several recently introduced HTE-robust estimators. We use the term ‘‘HTE’’ to refer to individual treatment effects that are arbitrarily heterogeneous, that is, $\tau_{i,t} \neq \tau_{j,s}$ for some i, j, s, t . HTE-robust estimators are defined as those that produce causally interpretable estimates under their respective identifying assumptions. For a more comprehensive discussion on these estimators, please refer to the SM.

Summary of HTE-Robust Estimators

Table 1 summarizes the estimators we discuss in this article. The primary difference resides in the mechanics of their estimation strategies: there are methods based on canonical DIDs and methods based on imputation. We refer to the former as *DID extensions* and the latter as *imputation methods*. DID extensions use dynamic treatment effects, estimated from local, 2×2 DIDs between treated and control observations, as building blocks. Imputation methods use individual treatment

effects, estimated as the difference between an imputed outcome under control and the observed outcome (under treatment), as building blocks. The imputation estimator we use in this article employs TWFE, fitted with all observations under the control condition, to impute treated counterfactuals. Different strategies also entail different assumptions. Each DID extension, for example, relies on a particular type of PT assumption, whereas imputation methods presuppose a TWFE model for untreated potential outcomes and require a zero mean for the error terms, which is implied by strict exogeneity.

Another noteworthy difference lies in the settings in which these estimators are applicable: some estimators can only be used in settings with staggered treatment adoption, while others can accommodate treatment reversals. In the latter setting, all estimators we discuss also require no anticipation and no or limited carryover effects. Furthermore, the estimators diverge in terms of (1) how they select untreated observations as controls for treated units, (2) how they incorporate pre-treatment or exogenous covariates, and (3) the choice of the reference period. We discuss these details further below and in Section A.1 of the SM.

DID Extensions

DID extensions are all built from local, 2×2 DID estimates. The overarching strategy for these estimators is to estimate the dynamic treatment effect, τ_l (or $\tau_{g,l}$ for each cohort g in the staggered setting), for each period since the most recent initiation of treatment, l , using one or more *valid* 2×2 DIDs. By ‘‘valid,’’ we mean that the DID includes (1) a pre-period and a post-period and (2) a treated group and a comparison group. The pre-period is such that all observations in both groups are in control, whereas the post-period is such that observations from the treated group are in treatment and those from the comparison group are in control. The choice of the comparison group is the primary distinction between estimators in this category. To obtain higher-level averages such as the ATT, we then average over our estimates of τ_l (or $\tau_{g,l}$), typically employing appropriate, convex weights.

Three estimators in this category are appropriate only for the staggered setting. Sun and Abraham (2021) propose an interaction-weighted (IW) estimator, which is a weighted average of $\tau_{g,l}$ estimates obtained from a TWFE regression with cohort dummies fully interacted with indicators of relative time to the onset of treatment. They demonstrate that, in a balanced panel, each resulting estimate of $\tau_{g,l}$ can be characterized as a difference in the change in average outcome from a fixed pre-period $s < g$ to a post-period l periods since g between the treated cohort g and the comparison cohort(s) in some set \mathcal{C} . The authors recommend using $\mathcal{C} = \sup_i E_i$, which is either the never-treated cohort or, if no such cohort exists, the last-treated cohort. By default, IW uses $l = 0$ as the reference period and can accommodate covariates in the TWFE regression.

⁹ Some of the authors we reference denote this first post-treatment period with $l = 0$.

TABLE 1. Summary of HTE-Robust Estimators							
Type	DID extensions: uses 2 × 2 DIDs as building blocks					Imputation methods	
Setting	Staggered: treatment reversals not allowed			General: treatment reversals allowed			
Research article	Sun and Abraham (2021)	Callaway and Sant’Anna (2021)	Cengiz et al. (2019)	de Chaisemartin and D’Haultfœuille (2020; 2024)	IKW (2023)	BJS (2024)	LWX (2024)
Method known as	IW	CSDID	StackedDID	did_multitple	PanelMatch	DID _{impute}	FEct
Key ID assumption	Parallel trends	Parallel trends	Parallel trends	Parallel trends	Parallel trends	Strict exogeneity	Strict exogeneity
Finest estimand	$\tau_{g,l}$	$\tau_{g,l}$	τ_l^{YW}	$\tau_{g,l}$	τ_l	$\tau_{i,t}$	$\tau_{i,t}$
Common aggregated estimand	ATT	ATT	ATT	ATT for switchers for l periods		ATT	ATT
Comparison group	Never-treated or last-treated	Never-treated or not-yet-treated	Never-treated	Matched stable group (not-yet-treated)	Matched stable group (not-yet-treated)	Imputed counterfactuals (not-yet-treated)	Imputed counterfactuals (not-yet-treated)
Reference period(s)	Period 0	An arbitrary pre-treatment period	Period 0	Period 0	Period 0	All pre-treatment periods	All pre-treatment periods
Covariate adjustment	Possible extension	Outcome and propensity score modeling	Outcome modeling	Possible extension	Refined matched set and outcome modeling	Outcome modeling	Outcome modeling

Employing the same general approach, Callaway and Sant’Anna (2021) propose two estimators, one of which uses never-treated units ($\hat{\tau}_{nev}^{CS}$) and the other not-yet-treated units ($\hat{\tau}_{ny}^{CS}$) as the comparison group. We label these estimators collectively as CSDID. Note that $\hat{\tau}_{nev}^{CS}$ uses the same comparison group as IW when a never-treated cohort exists,¹⁰ whereas $\hat{\tau}_{ny}^{CS}$ uses all untreated observations of not only never-treated units but also later adopters as controls for earlier adopters. Besides the choice of comparison cohort, CSDID estimators differ from IW in that they allow users to condition on pre-treatment covariates using both an explicit outcome model and inverse probability weighting simultaneously; consistency of the estimators requires at least one of these to be correct. By default, both IW and CSDID use the period immediately before the treatment’s onset as the reference period for estimating the ATT.

Stacked DID, first formally introduced by Cengiz et al. (2019), is another related estimator sometimes used to address HTE concerns. As described by Baker, Larcker, and Wang (2022), it involves creating separate sub-datasets for each treated cohort by combining data from that cohort (around treatment adoption) and data from the never-treated cohort from the same periods. These cohort-specific datasets are then “stacked” to form a single dataset. An event-study regression akin to Equation 2 with the addition of sub-dataset specific unit and time dummies is then run. This method uses the same comparison group as IW and the never-treated version of CSDID without covariates, but stacked DID estimates a single dynamic treatment effect for a given relative period rather than separate estimates for each cohort. Essentially, stacked DID is a special case of IW that uses immutable weights selected by OLS. We denote the corresponding estimand τ_l^{yw} to reflect the fact that it is variance-weighted. These weights are generally neither proportional to cohort sizes nor guaranteed to sum to one (Wing, Freedman, and Hollingsworth 2024). Thus, while stacked DID avoids the “negative weighting” problem and meets our criteria for HTE-robustness, its estimands are not the same as τ_l or the ATT.

In settings with treatment reversals, separate groups of researchers have converged on the same strategy for choosing a comparison group: matching treated and control observations that belong to units with identical treatment histories. IKW (2023) suggest one such estimator, PanelMatch, which begins by constructing a “matched set” for each observation (i, t) such that unit i transitions into treatment at time t . This matched set includes units that both (1) are not under treatment at time t and (2) share the same treatment history as i for a fixed number of periods leading up to the treatment

onset. For each treated observation (i, t) and for every post-period $(t + l - 1)$ such that unit i is still under treatment, it then estimates a local DID using the same pre-period $s < t$. The treatment “group” comprises solely observation (i, t) , and the members of the matched set for (i, t) that are still under control during $t + l - 1$ serve as the comparison group. To obtain τ_l for a given l , it then averages over the corresponding local DID estimates from all treated observations. IKW (2023) propose incorporating covariates by “refining” matched sets and use $l = 0$ as the reference period.

Using a similar strategy, de Chaisemartin and D’Haultfœuille (2020) propose a “multiple DID” estimator, DID_multiple. A notable difference is that they include local DIDs for units leaving the treatment and not just those joining the treatment; when there are no treatment reversals or covariates, DID_multiple is a special case of PanelMatch. The original proposal for DID_multiple also only considers the case where we match on a single period and where $l = 1$, but since it has been extended (de Chaisemartin and D’Haultfœuille 2024). Consequently, the target estimand is not the ATT but rather an average of the contemporaneous effects of “switching” (i.e., the effect of joining or the negative of the effect of leaving at the time of doing so). In the staggered setting, the PanelMatch estimator aligns with the not-yet-treated version of CSDID (without covariate adjustment). We delve into details on the connections between these three estimators in the SM.

All DID extensions are built using local, 2×2 DIDs, and their assumptions reflect this. Specifically, they each rely on a form of the PT assumption—that is, the expected changes in untreated potential outcomes from one period to the other are equal between the treated and the chosen comparison groups. We defer readers to the SM for a fuller account of each method’s assumptions.

The Imputation Method

Imputation estimators do not explicitly estimate local DIDs. Instead, they take the difference between the observed outcome and an imputed counterfactual outcome for each treated observation. The connection to the TWFE model is in the functional form assumption used to impute counterfactual outcomes. Specifically, an imputation estimator first fits a parametric model for the potential outcome under control $Y_{i,t}(0)$ —in our case, $Y_{i,t}(0) = X'_{i,t}\beta + \alpha_t + \zeta_t + \varepsilon_{i,t}$ —using only control observations $\{(i, t) : D_{i,t} = 0\}$. It is also through this outcome model that one can adjust for time-varying covariates. Then, it imputes $Y_{i,t}(0)$ for all treated observations $\{(i, t) : D_{i,t} = 1\}$ using the estimated parameters. Finally, it estimates the individual treatment effect, $\tau_{i,t}$, for each treated observation (i, t) by calculating the difference between the observation’s observed outcome $Y_{i,t} = Y_{i,t}(1)$ and its imputed counterfactual outcome $\hat{Y}_{i,t}(0)$. Inference for the estimated $\hat{\tau}_{i,t}$ is possible, although uncertainty estimates need to be adjusted to account for the presence of

¹⁰ This equivalence holds when there are no missing data. Otherwise, IW from the saturated regression differs from one that directly estimates local DIDs, including the never-treated version of CSDID. These estimates are typically close but can differ substantially, as in Kuipers and Sahn (2023).

idiosyncratic errors (e.g., Bai and Ng 2021). BJS (2024) and LWX (2024) each propose estimators in this category. Each article proposes a more general framework that nests many models, including TWFE. The latter also introduces several specific imputation estimators, one of which uses the TWFE model, and the authors refer to the resulting estimator as the fixed effects counterfactual estimator, or FEct.

Although DID extensions and imputation methods rely on slightly different identification assumptions, these assumptions usually lead to similar observable implications. Researchers commonly use the presence or absence of pretrends to judge how plausible the PT assumption is. In the classic two-group setting, if there are data from multiple pre-treatment periods, researchers can plot the time series of average outcomes of each group and visually inspect whether they indeed trend together. The intuition is that if PT holds and the average outcome trends of the treated and control groups are indeed parallel in pre-treatment periods when $Y(0)$'s are observed for all units, then it is plausible that PT also holds in the post-treatment periods, when $Y(0)$'s are no longer observable for units in the treatment group. Conversely, differential trends in the pre-treatment periods should make us suspicious of PT. In more complex settings or where we wish to control for observed confounders, researchers often use dynamic estimates before and after the onset of treatment, τ_t , to construct so-called “event-study plots” to judge the presence of pretrends. If PT holds, then pre-treatment dynamic estimates should be around zero. We provide a more thorough discussion and an example of the event-study plot in the next section when we introduce our procedure.

Choice of Estimators

In general, we believe the credibility of identifying assumptions is more important than the choice of estimator. After all, in the staggered setting, when assumptions hold, IW, CSDID, DID_multiple, Panel-Match, and the imputation estimator all converge to the same or a similar estimand. However, there are a few reasons to favor the imputation estimator. First, it can handle complex settings, including those with treatment reversal—which account for over half of the studies in our sample—and can accommodate time-varying covariates, additional fixed effects, and unit- or group-specific time trends commonly seen in social science research. The imputation estimator connects to TWFE through a shared outcome model, and thus any of the aforementioned modifications to the outcome model can be directly mirrored. DID extensions relate to TWFE through their shared connection to DID in the two-group, two-period setting. Classic DID's inability to naturally accommodate these complexities limits DID extensions on this front. Just like TWFE, the imputation estimator risks misspecification bias, and adding more redundant terms may significantly increase variance. However, we still

consider the added flexibility to be an advantage. Second, imputation estimators are the most efficient under homoskedastic errors (BJS 2024).¹¹ Moreover, by using the average of all pre-treatment periods as the reference point rather than a single pre-period, as the default in DID extensions, they provide greater power in hypothesis testing for pretrends. The main drawback of imputation estimators is that their current implementations (either FEct or DID_impute) do not allow for automated adjustment of time-invariant covariates, an advantage offered by CSDID and Panel-Match. Adjusting for pre-treatment characteristics can improve credibility of research, as conditional PT may be more plausible than the unconditional one (Sant'Anna and Zhao 2020).

DATA AND PROCEDURE

Next, we assess the robustness of empirical findings from causal panel analyses in political science and compare results obtained using the different methods we have discussed. We will explain our sample selection rules, describe standard practices in the field, and outline our reanalysis approach. Readers can find a more detailed explanation of our sample selection criteria and replication and reanalysis procedure in Section A.2 of the SM.

Data

Our replication sample comprises studies from three leading political science journals, *APSR*, *AJPS*, and *JOP*, published over a recent 7-year span from 2017 to 2023. We initially include all studies, including both long and short articles, that employ panel data analyses with a binary treatment as a key component of their causal argument, resulting in a total of 102 studies. After a careful review, as explained in footnote 1, we find that 64 studies employ a TWFE model similar to Equation 1. We then attempt to replicate the main results of these 64 studies and are successful in 49 cases (76.6%). Though a significant proportion of studies failed to replicate, we note that the success rate is still higher than that of Hainmueller, Mummolo and Xu (2019) at 55% and Lal et al. (2024) at 67%. We credit this to the new replicability standards set by journals. A detailed explanation of how we select the “main model” is provided in Section A.3 of the SM. Table 2 depicts the distribution of successful replications, along with reasons for replication failures, across the various journals.

¹¹ In Section A.3 of the SM, we demonstrate that the imputation estimator tends to yield larger z -scores, based on data from our sample. In the majority of cases, the imputation estimator has a smaller standard error, and the difference can be especially dramatic for IW and the never-treated version of CSDID, which often discard the vast majority of untreated observations.

TABLE 2. Sample Selection and Replicability of Qualified Studies

Journal	All	TWFE (attempted)	Incomplete data	Replication error	Replicable	Success rate%
<i>APSR</i>	22	13	2	1	10	76.9
<i>AJPS</i>	31	21	3	3	15	71.4
<i>JOP</i>	49	30	6	0	24	80.0
Total	102	64	11	4	49	76.6

TABLE 3. Settings and Common Practice

<i>Motivations for TWFE</i>				<i>Variance estimator</i>		
“Difference-in-differences”	33	67.3%		Cluster-robust SE or PCSE	48	98.0%
“Within” variations	16	32.7%		Cluster-bootstrapped procedures	8	16.3%
<i>Treatment setting</i>				<i>Variants in specifications</i>		
Classic 2x2 DID	3	6.1%		Lagged dependent variables	8	16.3%
Multi-period block DID	6	12.2%		Higher-than-unit-level time trends	1	2.0%
Staggered DID	13	26.5%		Unit-specific linear time trends	15	30.6%
General	27	55.1%		<i>Data visualization</i>		
<i>Outcome variable</i>				Group average outcomes	19	38.8%
Continuous	44	89.8%		Event-study plots	23	46.9%
Binary	5	10.2%		Neither	17	34.7%

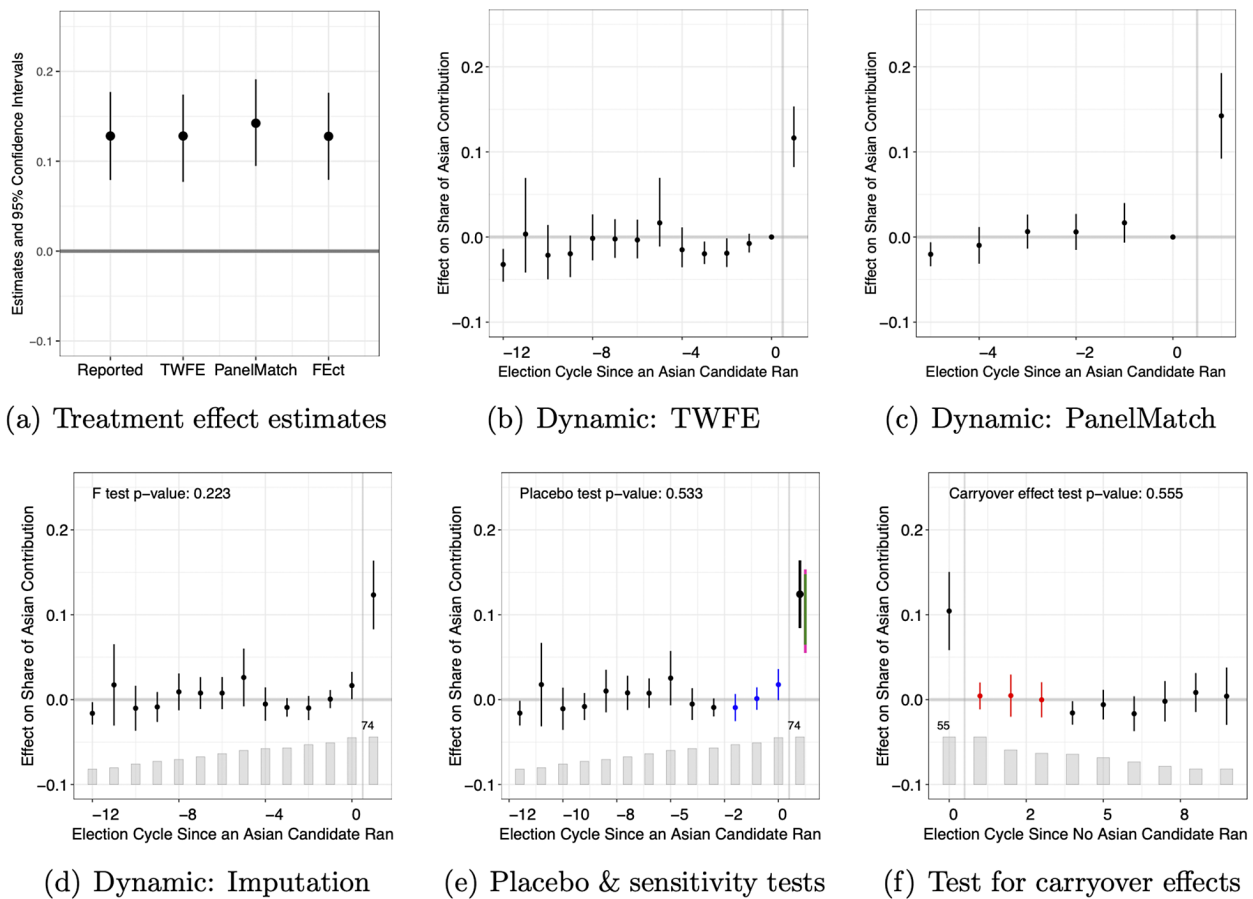
Settings and Common Practices

Table 3 presents an overview of the standard practices and settings in the studies that we successfully replicated. The majority of studies in our sample (67.3%) use the DID design/method/approach to justify the use of the TWFE model, while the remaining studies advocate for the model’s ability to exploit “within” variations in the data. Out of the 49 studies, nine (18.4%) employ a classic block DID setting, which includes two-group, two-period designs (three studies) and multi-period block DID designs (six studies). Thirteen studies (26.5%) use a staggered DID design, while the remaining 27 studies (55.1%) fall into the “general” category, meaning they allow for treatment reversals. Except for five, all studies feature a continuous outcome of interest. Most use cluster-robust SEs or panel-corrected SEs (Beck and Katz 1995), and eight studies employ bootstrap procedures for estimating uncertainties. A subset of studies explore alternative model specifications by adding lagged dependent variables (eight studies), unit-specific linear time trends (fifteen studies), and higher-than-unit-level time trends (one study). Notably, 32 studies use some type of visual inspection—either average outcomes over time, event-study plots, or both—to evaluate the plausibility of PT. Four studies published in 2023 (33%) employ HTE-robust estimators, compared to none before 2023, indicating rapid adoption of these methods. Of these, two use CSDID, one PanelMatch, and one the imputation estimator.

Procedure

We use data from Grumbach and Sahn (2020) to illustrate our process for replication and reanalysis. The authors investigate coethnic mobilization by minority candidates during U.S. congressional elections. To simplify our analysis, we focus on the impact of the presence of an Asian candidate on the proportion of general election contributions from Asian donors. To begin, we aim to understand the research setting and data structure. We visualize the patterns of treatment and outcome variables using plots, which are shown in the SM. In this application, treatment reversals clearly take place. Some data are missing (due to redistricting), but the issue does not seem to be severe. We record important details such as the number of observations, units, and time periods, the type of variance estimator, and other specifics of the main model. Next, we replicate the main finding, employing both the original variance estimator and a cluster-bootstrap procedure.

We then re-estimate the ATT and dynamic treatment effects using estimators discussed in the previous section. For staggered adoption treatment cases, we apply seven estimators: TWFE (with always treated units removed for easier comparisons with other estimators), the imputation estimator (FEct), PanelMatch, DID_multiple, StackedDID, IW, and CSDID (both not-yet-treated and never-treated versions). For applications with treatment reversals like Grumbach and Sahn (2020), we use the first three

FIGURE 2. Reanalysis of Grumbach and Sahn (2020)

Note: Reanalysis of data from Grumbach and Sahn (2020). (a) Treatment effect or ATT estimates from multiple methods. (b–d) Event-study plots using TWFE, PanelMatch, and the imputation estimator (FEct). (e,f) Results from the placebo test (and robust confidence set) and test for carryover effects using FEct—the blue points in (e) and red points in (f) represent the holdout periods in the respective tests. In (e), the green and pink bars represent the 95% robust confidence sets when $\bar{M} = 0$ and $\bar{M} = 0.5$, respectively. CIs in all subfigures—excepted for the reported estimate in (a)—are produced by bootstrap percentile methods.

estimators only.¹² The comparison between the TWFE estimate and the other estimates sheds light on whether original findings are sensitive to relaxing the constant treatment effect assumption. Figures 2a–d show the main results from this example. The similarity between estimates for the ATT in panel a suggests that the original finding is robust to the choice of estimators. The event-study plots from HTE-robust estimators in panels c and d are broadly consistent with the event-study plot from TWFE in panel b.

Next, we conduct diagnostic tests based on the imputation estimator, including the F test and the placebo test, to further assess the plausibility of PT and, in applications with treatment reversal, the no-carryover-

effect assumption. We use the imputation estimator because it is applicable across all studies in our replication sample, can incorporate time-varying covariates, and remains highly efficient. Figures 2d–f show the results from the F test, placebo test, and test for no carryover effects on our running example, respectively. Both a visual inspection and the formal tests suggest that PT and no-carryover-effect assumptions are quite plausible.

Finally, we compute the robust confidence sets proposed by Rambachan and Roth (2023), which account for potential PT violations when testing the null hypothesis of no post-treatment effect. Specifically, we employ the relative magnitude restriction, with two modifications to accommodate the imputation method. First, we use estimates from the placebo test to ensure that benchmark pre-treatment estimates are obtained using the same approach as post-treatment ATT estimates. This alignment prevents potential asymmetry in testing and treatment effect estimation (Roth 2024). Second, since the imputation method does not rely on a single reference period, we explicitly

¹² A recent development of DID_multiple, DIDmultiplegDYN, allows for the estimation of dynamic and long-term effects (de Chaisemartin and d'Haultfoeuille 2024). It defines an estimand that accounts for carryover effects and targets the first onset of the treatment. To remain consistent with the imputation framework, which allows multiple treatment onsets, we apply DID_multiple only to cases featuring staggered treatment timing.

incorporate the placebo estimate from the last pre-treatment period ($\hat{\delta}_0$) to account for deviations of post-treatment estimates from earlier reference periods. Mathematically, we decompose each dynamic estimate, μ_t , into the true treatment effect, τ_t , and a trend (bias) component, δ_t , such that: $\mu_t = \tau_t + \delta_t$. Our modified relative magnitude restriction then requires that, for all $t \geq 0$,

$$|\delta_{t+1} - \delta_t| \leq \bar{M} \cdot \max_{s \in \mathcal{P} \setminus \{0\}} |\delta_{s+1} - \delta_s|, \quad (3)$$

where \mathcal{P} is the set of placebo periods. In our application, we set $\mathcal{P} = \{-2, -1, 0\}$, so the maximum violation among placebo periods is: $\max \{|\delta_0 - \delta_{-1}|, |\delta_{-1} - \delta_{-2}|\}$.

When $\bar{M} = 0$, the relative magnitude restriction reported in Equation 3 implies that $\delta_t = \delta_0$ for all $t > 0$, meaning that the PT violation remains fixed at the same level as in the last pre-treatment placebo period. In this case, the robust confidence set obtained at $\bar{M} = 0$ acts as a debiased confidence interval, using the placebo estimate from the last pre-treatment period as the benchmark for bias. Allowing $\bar{M} > 0$ permits PT violations to vary over time, but constrains the change in magnitude of violations between consecutive post-treatment periods to remain within \bar{M} times the largest consecutive discrepancy observed during the placebo periods.

Rambachan and Roth (2023, 2653) suggest using $\bar{M} = 1$ as a “natural benchmark” when the number of placebo periods is roughly equal to the number of post-treatment periods, treating any potential PT violations as no worse than those already observed.¹³ In our reanalysis, we first construct 95% robust confidence sets for each post-treatment dynamic effect and the ATT at $\bar{M} = 0$ and $\bar{M} = 0.5$. Figure 2e illustrates these robust confidence sets for the estimated ATT using the imputation method. The center of the robust confidence sets is smaller than the point estimate because $\hat{\delta}_0 > 0$. If the confidence sets for $\bar{M} = 0$ does not include zero, as in this case, we conduct a sensitivity analysis by varying \bar{M} over a wider range to determine the “breakdown value” \bar{M} , which is the smallest value of \bar{M} at which the robust confidence set first includes zero. In the case of Grumbach and Sahn (2020), the breakdown value is $\bar{M} = 2.5$, which means that the estimated coethnic mobilization effect remains statistically distinguishable from zero unless PT violations are more than 2.5 times the largest discrepancy observed during the placebo periods.

Overall, the results from Grumbach and Sahn (2020) appear highly robust, regardless of the choice of point and variance estimators. The PT and no-carryover-effect assumptions seem plausible. The study also has sufficient power to distinguish the ATT from zero, even under potential, realistic PT violations.

SYSTEMATIC ASSESSMENT

We perform the replication and reanalysis procedure described above for all 49 studies in our sample. This section offers a summary of our findings, with complete results for each article available in the SM. We organize our results around two main questions: (1) Are existing empirical findings based on TWFE models robust to HTE-robust estimators? (2) Is the PT assumption plausible, and do original findings remain robust to mild PT violations informed by pretrends? We also discuss other issues observed in the replicated studies, including the presence of carryover effects and sensitivity to model specifications.

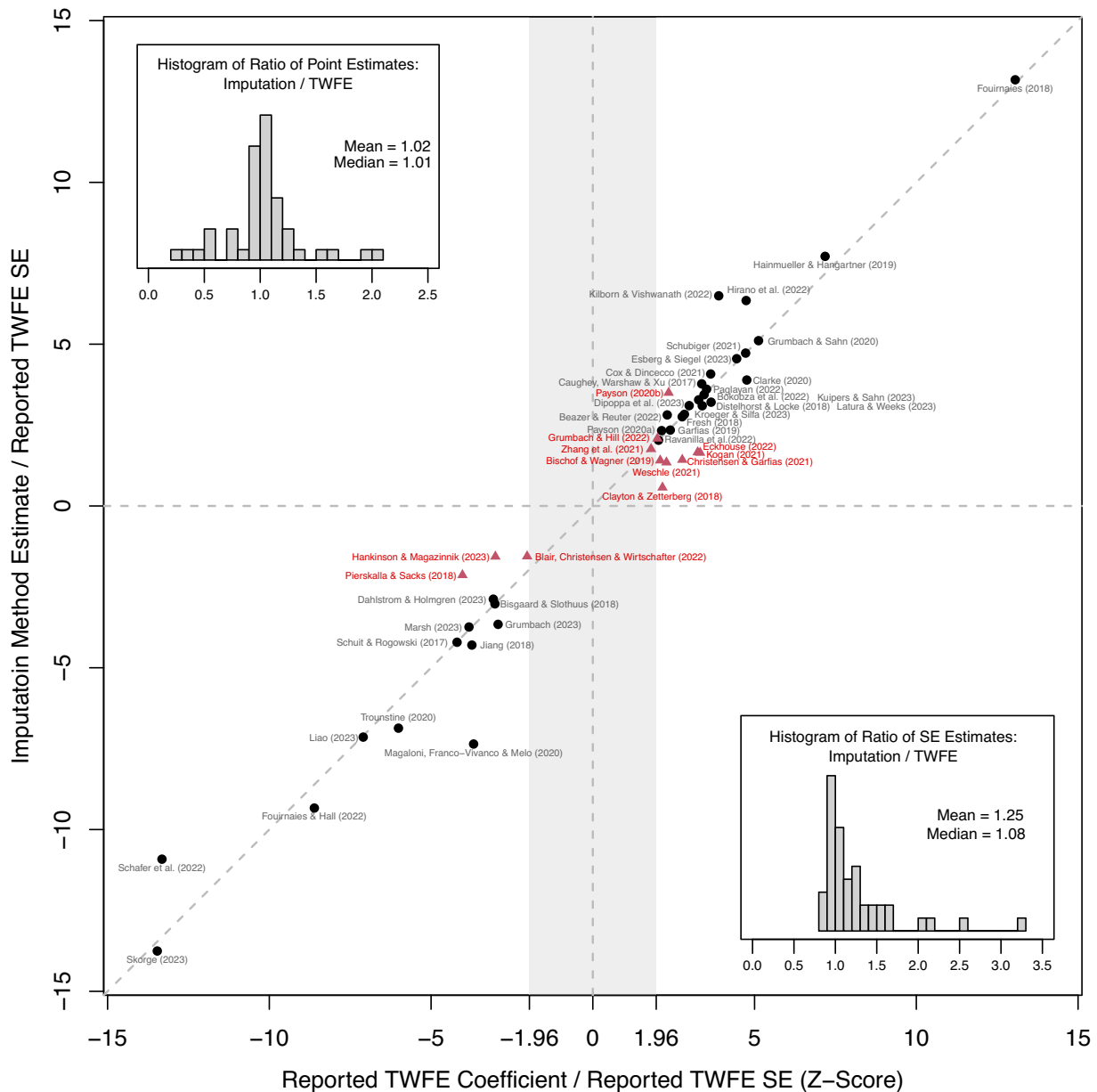
HTE-Robust Estimators Yield Qualitatively Similar but Highly Variable Estimates

To examine the impact of the weighting problem caused by HTE associated with TWFE models, we first compare the estimates obtained from the imputation estimator, FEct, for all studies to those originally reported. We choose the imputation estimator for the reason mentioned earlier. Most importantly, it is applicable to all studies in our sample, including those with treatment reversals and those with additional time trends. Figure 3 plots the comparison. The horizontal axis represents the originally reported TWFE estimates, and the vertical axis represents FEct estimates, both normalized using the same set of originally reported SEs. If the point estimates are identical, then the corresponding point should lie exactly on the 45-degree line. Red triangles represent studies where the imputation estimates are statistically insignificant at the 5% level, based on cluster-bootstrapped SEs.

We observe several patterns. First, TWFE coefficients are statistically significant at the 5% level in all but one study, and the absolute values of z-scores for a significant minority of studies cluster around 1.96, indicating possibly a file-drawer problem and potential publication bias. Second, the points largely follow the 45-degree line, with the imputation estimates *always* having the same sign as the original estimates. This suggests that while scenarios where accounting for HTE completely reverses the empirical findings are theoretically possible, they are rare.

However, results sometimes deviate significantly. In the top left corner of Figure 3, we plot the histogram of the ratio of imputation to TWFE estimates. Although the mean and median of the ratio are close to one, at the extremes, we observe imputation estimates as small as one-fourth or as large as more than double the TWFE estimates. The most consequential deviations occur in studies that were originally near the margins of statistical significance. Additionally, we plot the ratio of SE estimates from the imputation method to TWFE in the bottom right corner of Figure 3. The median is 1.08, meaning that in the majority of cases the SE estimate from the imputation method is at least 8% larger than the SE from TWFE. The mean is 1.25, and the distribution is right-skewed; in the extreme, the ratio was almost as high as three. Combined, drops in point

¹³ In our setting, the number of post-treatment periods typically exceeds the number of placebo periods, which means the criterion is even more lenient than the authors have suggested.

FIGURE 3. TWFE vs. The Imputation Estimator: All Cases

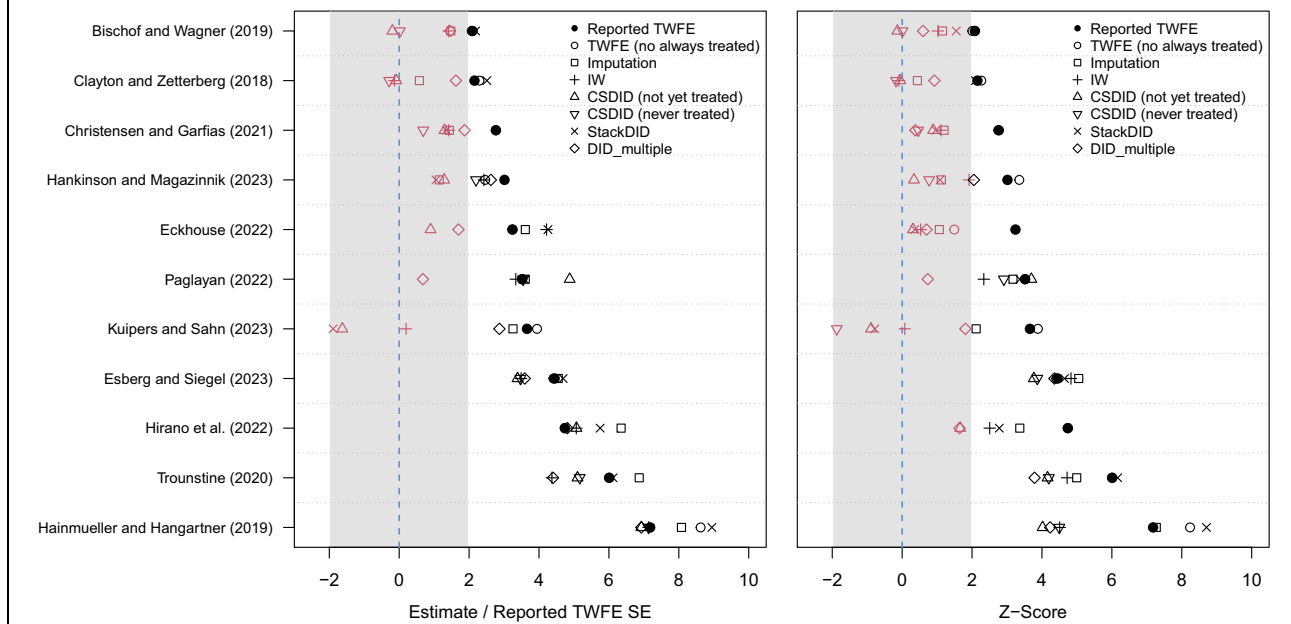
Note: The above figure compares reported TWFE coefficients with imputation method (FEct) estimates. Both estimates for each application are normalized by the same reported TWFE SE. Fourniaies and Hall (2018) and Hall and Yoder (2022) are close to the 45-degree line but are not included in the figure as their TWFE z-scores exceed 15. Black circles (red triangles) represent studies whose imputation method estimates for the ATT are statistically significant (insignificant) at the 5% level, based on cluster-bootstrapped SEs. The top-left (bottom-right) corners display histograms of the ratio of point (SE) estimates based on the imputation method and TWFE. These plots show that changes in point estimates, combined with the efficiency loss from using the imputation method, contribute to the loss of statistical significance in some studies.

estimates and increases in uncertainty lead to the third pattern: When we switch from TWFE to the imputation estimator, the number of studies that are statistically insignificant at the 5% level increases from one to twelve (24%).

If we restrict our attention to eleven studies with staggered treatments, we can broaden our comparison set to include more HTE-robust estimators

discussed earlier.¹⁴ Figure 4 visually compares the points estimates (left panel) and z-scores (right

¹⁴ Kogan (2021) and Magaloni, Franco-Vivanco, and Melo (2020) are excluded because the original specifications include additional time trends, which are not supported by HTE-robust estimators except the imputation estimator.

FIGURE 4. Comparison of Estimates: The Staggered Setting

Note: The above figures compare reported TWFE coefficients with estimates from various alternative estimators. In the left panel, all eight estimates for each application are normalized by the same reported SE to highlight changes resulting from the use of alternative estimators. In the right panel, the estimates are divided by their respective bootstrapped SEs. To facilitate visualization, we multiply all estimates by the sign of the reported coefficient. In both figures, black (red) symbols represent estimates that would be statistically significant (insignificant) at the 5% level, assuming they were treated as z-scores. The normalized CSDID (never treated) estimate for Kuipers and Sahn (2023), -5.36 , falls out of plotting area and is therefore not shown in the left panel. Kogan (2021) and Magaloni, Franco-Vivanco, and Melo (2020) are excluded because the authors' original TWFE specifications include unit-specific linear time trends, which are not supported by most HTE-robust estimators except the imputation estimator. Some estimates are missing because of too few never-treated units. PanelMatch is excluded because it targets a different estimand.

panel). In the left panel, all point estimates are normalized by the same set of reported TWFE SEs for each study. In all but three studies, the estimates from all HTE-robust estimators share the same sign, though there is a noticeable amount of variation in the estimated effect size.¹⁵

As in Figure 3, TWFE does not appear to be systematically upward or downward biased compared to HTE-robust estimators. Another observation that carries over is that HTE-robust estimators generally require more power to reject the null hypothesis of no effect. In five of the eleven studies, at least four HTE-robust estimates per study are statistically insignificant. The left panel shows that the changes in point estimate alone are often sufficient to render the results statistically insignificant. The comparison of z-scores in the right panel highlights that increased uncertainties can be sub-

stantial. Combined with earlier evidence, these findings from the staggered cases suggest that the HTE issue regarding TWFE is empirically significant and warrants careful consideration by researchers.

It is worth noting that while only a small fraction of studies in our sample (eight studies, 16.3%) employ a bootstrap procedure to estimate SEs or CIs, the widely practiced cluster-robust SE typically performs adequately. This is because the number of units (clusters) is generally large, with a median of 317. However, ten studies have fewer than 50 units; among them, two studies that were significant at the 5% level using cluster-robust SEs fell below this threshold when using cluster-bootstrapped SEs, both of which were already marginally significant. We provide comparisons of reported, cluster-robust, and cluster-bootstrapped SEs in the SM.

Relaxing PT Renders Most Studies Unable to Reject the Null

Although the recent methodological literature has heavily focused on HTE, PT violations—long recognized as a potential pitfall—remain a primary concern in practice. In Figure 5, we present the event-study plot based on estimates from the imputation estimator for each study in our sample. We also report the ATT estimates and their bootstrapped SEs. Due to space

¹⁵ Although a sign change is observed in Clayton and Zetterberg (2018) when using CSDID, IW, and DID_multiple, these estimates are negligibly small and statistically insignificant. Similarly, the not-yet-treated version of CSDID is the opposite sign but miniscule and statistically insignificant in Bischof and Wagner (2019). The estimates from IW and CSDID are also of the opposite sign for Kuipers and Sahn (2023), and they are of a much larger magnitude. The never-treated version of CSDID is also statistically significant.

FIGURE 5. Event-Study Plots w/ Imputation Estimator

Note: We report the estimated ATT and corresponding bootstrap SEs (in parentheses) using FEct. For Skorge (2023), we use an “exit” plot because all treated units receives the treatment in the first period.

limitations, we present the event-study plots from other estimators, as well as results from the placebo tests, robust confidence sets, and sensitivity analyses, in the SM.

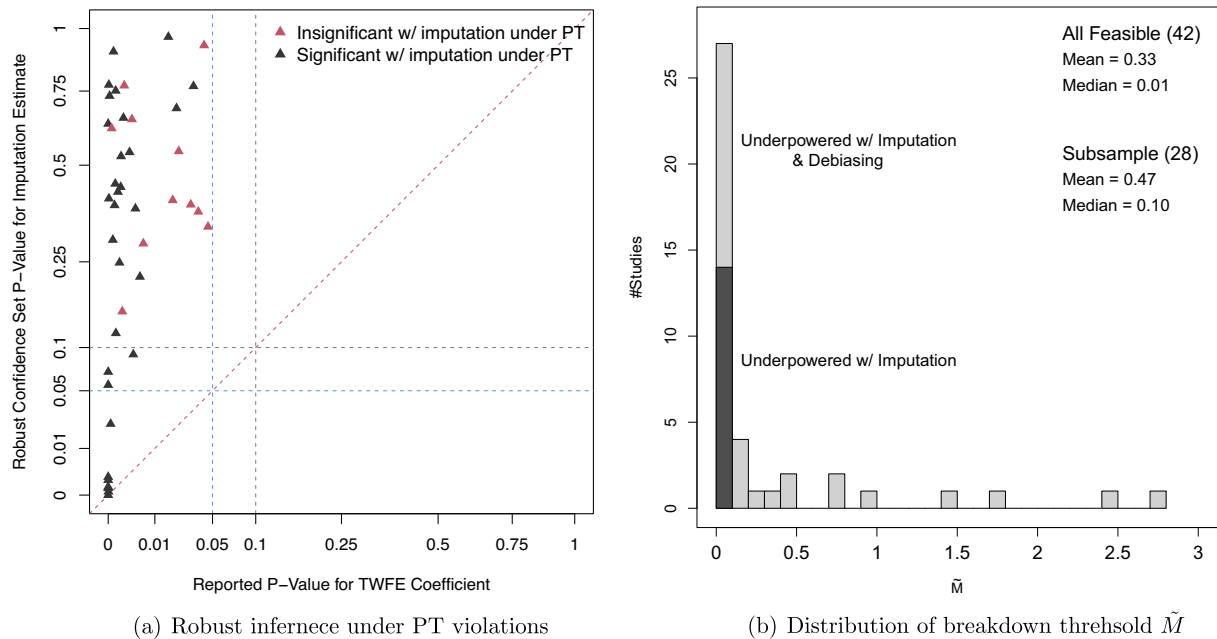
What is concerning is that, in at least six studies (12%), the PT assumption seems highly implausible. In these studies, the dynamic estimates in the pre-treatment periods deviate substantially from zero compared with

FIGURE 5. Continued

those in the post-treatment period, and the F test rejects the null.¹⁶ For the remaining studies, the CIs of pre-

treatment estimates mostly cover zero, and the F test and placebo test do not reject. However, this could be simply due to a lack of statistical power. Therefore, to assess the robustness of the findings, we need additional tools that simultaneously account for both the estimated pretrend and statistical power. The sensitivity analysis proposed by Rambachan and Roth (2023) addresses this issue. We conduct such an analysis with a modified

¹⁶ There are other cases where the F test rejects, but we do not consider them highly problematic because, with a large sample size, a confounder inconsequential to the ATT estimate can still produce a small p -value in the F test. This is why the sensitivity analysis approach is particularly useful.

FIGURE 6. Allowing PT Violations with Robust Confidence Sets

Note: The above figures present findings from the sensitivity analysis. The sample consists of 42 studies with more than three pre-treatment periods, allowing for such an analysis. Subfigure (a) displays the p -values of partially identified ATT estimates using the imputation method under restricted relative magnitude PT violations with $\bar{M} = 0.5$, compared to reported p -values for TWFE coefficients assuming PT. A square root scale is used to facilitate visualization. Black (red) triangles represent studies that are statistically significant (insignificant) at the 5% level when using the imputation method under PT. Subfigure (b) shows a histogram of \bar{M} , the breakdown values of \bar{M} ; the dark gray bar represents studies whose ATT estimates are statistically insignificant at the 5% level when using the imputation method.

relative magnitude restriction from Equation 3 for 42 studies that have at least three pre-treatment periods.¹⁷

Figure 6a shows that when we construct a robust confidence set with $\bar{M} = 0.5$, the null hypothesis of no effect is rejected at the 5% level in only eight (19%) of the 42 studies. Many studies that were robust to HTE now appear underpowered. Figure 6b displays the distribution of breakdown values \bar{M} . The spike at 0 consists of two types of studies: those that are statistically insignificant with the imputation estimator when three placebo periods are considered (dark gray) and those that become statistically insignificant due to debiasing using the placebo estimate $\hat{\delta}_0$. In other words, many results are not robust when we take into account $\hat{\delta}_0$ even without a relative magnitude shift. Among the 42 studies, the median is close to 0 and the mean is 0.33. Focusing only on the studies that remain statistically significant with the imputation estimator, the median and mean are still as low as 0.10 and 0.47, respectively.

These results suggest that although most statistically significant findings are robust to the imputation estimator, in the vast majority of these studies, accounting for potential PT violations—even very mild ones based on estimated pretrends during the placebo periods—prevents us from rejecting the null. In other words, in

most studies, the ATT estimates are not substantial enough to be differentiated from realistic PT violations or from estimation and sampling errors.

Other Issues

Our reanalysis highlights several additional issues. First, the presence of missing values is widespread. Although most methodological work presumes balanced panels without missing data, in reality, many empirical studies encounter varying degrees of data missingness. Substantial differences in results for estimators that are numerically identical in balanced panels suggest that such violations may have important implications. During replication, we generate plots that display the patterns of treatment status for each study (Mou, Liu, and Xu 2023). Based on these plots, we also see that in some studies the pattern of missingness is either seemingly nonrandom or extremely prevalent, which weakens our confidence in the respective empirical findings.

Second, we perform tests for carryover effects for all studies with treatment reversals. If this test is rejected, it suggests that the treatment effects persist beyond the treatment periods. Among 27 studies, five reject the null hypothesis of no carryover effects at 5%. Part of the concern is that the imputation method and DID extensions will use control observations from previously treated units to fit the potential outcome model or as

¹⁷ For studies with only three pre-treatment periods, we set the number of placebo periods in their placebo tests to 2.

TABLE 4. Summary of Findings

	<i>APSR</i>	<i>AJPS</i>	<i>JOP</i>	All	<i>n</i>
Imputation to TWFE ratio < 0.8	0.00	0.13	0.30	0.18	49
Imputation to TWFE ratio < 0.5	0.00	0.00	0.13	0.06	49
<i>F</i> test reject null	0.30	0.21	0.30	0.27	44
Placebo test reject null	0.40	0.29	0.15	0.25	44
TWFE <i>not</i> reject null	0.00	0.00	0.04	0.02	49
Imputation <i>not</i> reject null	0.00	0.13	0.44	0.25	49
Imputation <i>not</i> reject null w/ $\bar{M} = 0$	0.40	0.23	0.74	0.50	42
Imputation <i>not</i> reject null w/ $\bar{M} = 0.5$	0.70	0.77	0.90	0.81	42

Note: Entries (except in the last column) are proportions of studies satisfying each set of conditions. A null is deemed rejected if $p < 0.05$. Five studies with a single pre-treatment period are not included in the summary statistics for the *F* test and placebo test. Seven studies are not included in the last two rows due to too few pre-treatment periods.

comparisons for treated observations, and if there are carryover effects, then the comparisons will become tainted. LWX (2024) note that the presence of carryover effects for a limited number of periods is less concerning, as researchers can recode treatment to persist for some time after a unit transitions out of treatment. Despite its prevalence, we observe that carryover effects do not substantially alter the findings in most studies. Specifically, in six studies that reject the null of no carryover effects, when we exclude two periods after the treatment switches back to zero, the ATT estimates remain similar in magnitude, and statistically significant results remain significant (Figure A8 in the SM). Nevertheless, we recommend that researchers make it a practice to check for potential carryover effects, considering the low cost of conducting such tests and adjustments.

Finally, many findings are sensitive to model specifications. Some studies that we exclude from our sample employ one-way fixed effects or fixed effects at a level different from that at which treatment is assigned. Many of these findings do not hold when we reanalyze them using a TWFE model. We should clarify that this does not imply that the original results are wrong; rather, these models implicitly operate under different identifying assumptions, and there is substantial variation in how much consideration authors give to this point. Some studies do not provide a rationale for their choice to use one-way fixed effects, while others explicitly outline the type of unobserved confounders they intend to control for. In one instance, the authors inaccurately label their specification as a DID design. The TWFE and DID estimators are generally not equivalent, and we emphasize that this difference becomes even more pronounced when fixed effects are not assigned at the level of treatment. In such cases, a TWFE specification does not correspond to even a broadly defined DID design.

Summary

In Table 4, we summarize the main findings of our reanalysis. The numbers represent the proportion of studies in a given journal or across all journals that fall

under the respective category (hence, a smaller number is better). The first two rows relate to the ratio of the imputation estimate to the TWFE estimate, which proxies the consequences of the weighting problem due to HTE. Across all journals, this ratio is less than 0.8 in 18% of studies and less than 0.5 in 6% of studies, indicating that using an HTE-robust estimator does sometimes have a substantively significant impact and is important to use at least as a robustness check.

The third and fourth rows indicate the proportion of studies with evidence of PT violations based on the *F* test and placebo tests, respectively. Across all journals, around a quarter of studies reject the null hypothesis at the 5% level using the *F* test, with slightly fewer rejecting when using the placebo test due to the loss in power from excluding data from the placebo periods. Note, though, that failure to reject can result from insufficient power and is not sufficient to support that the PT holds.

The fifth and sixth rows display the proportion of studies in which the null hypothesis of no effect is not rejected using the TWFE and imputation estimators, respectively. When TWFE is employed, this occurs in only one study (2%); however, when the imputation estimator is used, this number increases to 24%, suggesting that many studies in our sample are potentially underpowered with an HTE-robust estimator. The last two rows show that in 50% and 81% of the studies, the robust confidence sets for the ATT include zero with $\bar{M} = 0$ (debiasing only) and $\bar{M} = 0.5$ (debiasing plus restricted relative magnitude in biases), respectively. The comparison of results in the last four rows highlights that the main source of fragility in the existing literature is potential PT violations, rather than concerns of HTE.

Our reanalysis is not meant to criticize existing studies, many of which were conducted before recent methodological advances. In fact, we have observed rapid adoption of these new methods and greater statistical power in more recently published studies. However, we want to emphasize two key points: (1) many studies do not adequately assess the plausibility of PT, the key identifying assumption of their research designs; and (2) given recent methodological developments, causal panel analysis under PT requires significantly more

statistical power to account for HTE and potential PT violations than previously believed.

RECOMMENDATIONS

Based on the findings of the reanalysis, we provide the following recommendations.

Research Design Is the Key

An important component of a strong research design is a clear understanding of how treatment is assigned. The PT assumption, which both TWFE models and most modern methods rely on, is silent on the assignment mechanism. As a result, many researchers assume that these methods can be applied when the assignment mechanism is unknown and that the absence of a pre-trend is sufficient to make PT credible. Another perspective suggests that (quasi-)randomization is required for PT to be credible (e.g., Kahn-Lang and Lang 2020). Our view aligns more closely with the latter, but is less stringent. We argue that for PT to be credible, researchers must justify the following assumption:

$$\Delta_{s,t}Y_{i,t}(0) \perp\!\!\!\perp D_{i,t}, \quad \forall s, t,$$

where $\Delta_{s,t}Y_{i,t}(0) = Y_{i,t}(0) - Y_{i,s}(0)$ is the before-and-after difference in untreated potential outcomes across any two periods (assuming no covariates, though this naturally extends to include them). This assumption is slightly stronger than PT, which only assumes mean independence, but is more intuitive. It demonstrates both the strength and limitation of causal panel analysis under PT: while the panel structure helps account for time-invariant unobserved confounding through before-and-after differences, the introduction of the treatment must act as a shock in so much as it is orthogonal to the evolution of untreated potential outcomes; hence, any dynamic relationships between past outcomes or covariates and treatment are ruled out. In other words, under PT, a strong causal panel design still requires some (quasi-)random element in treatment assignment. Importantly, the research design issues cannot be resolved simply by applying the novel estimators surveyed in this article.

Inspecting Raw Data Helps Spot Obvious Issues

The research design phase should also include inspection of the raw data. We encourage researchers to plot the data at hand to better understand patterns of the treatment status, missingness, and outliers (Mou, Liu, and Xu 2023). Treatment status should vary both by unit and time. If the majority of variation occurs over time (across units) with little or no variation between units at any given time period (or across time within a given unit), the TWFE estimand will be likely dominated by impermissible comparisons and thus susceptible to larger biases. Moreover, HTE-robust estimators will estimate the

treatment effect using very little data and thus be underpowered. Equally important is the need for researchers to understand the degree and possible origins of data missingness prior to initiating statistical analysis. If missingness does not seem to be random, or if it is too prevalent, leaving insufficient meaningful variation in the data, researchers should consider halting the analysis at this stage. Just as in the cross-sectional case, plotting the raw data can also help researchers to spot outliers and highly skewed distributions, which may require additional pre-processing. At this stage, researchers can also trim the data to make the units in comparison more similar in terms of pre-treatment or time-invariant covariates (Sant'Anna and Zhao 2020).

When Estimates Diverge, Understand Why

At the estimation stage, we recommend using at least one HTE-robust estimator alongside a benchmark TWFE model. While TWFE is often more efficient, its constant treatment effects assumption is too restrictive in many contexts. As shown in this article, TWFE does not produce systematic upward or downward bias compared to the imputation estimator, but it can severely bias causal estimates in individual cases. We recommend the imputation estimator in most settings primarily for its flexibility, though other estimators also have their advantages. The more critical question, however, is why results differ between estimators when they do. If TWFE and an HTE-robust estimator diverge, it could be due to HTE or PT violations. If HTE-robust estimators themselves diverge, it is often because the data are too sparse, leading to high variability, or the PT assumption fails differentially, causing estimators weighting control units differently to produce varying estimates. Plotting raw data and using diagnostics (such as the Goodman-Bacon decomposition) typically clarifies these issues. We also recommend keeping the benchmark TWFE model for its transparency.

For uncertainty estimates, researchers should use cluster-robust SEs when the number of clusters is large (e.g., exceeds 50) and opt for cluster-bootstrap or cluster-jackknife procedures when the number of clusters is relatively small. The clustering level should match the higher of either the time-series units or the level of treatment assignment. This follows the rule of thumb to cluster at the level of potential outcome input index, taking into account both treatment assignment and potential temporal spillover (Fu, Samii, and Wang 2024). For novel HTE-robust estimators, cluster-bootstrap or jackknife is generally safer than relying on various analytical SEs. In the SM, we show that cluster-bootstrapped SEs are typically larger than analytically derived SEs for five HTE-robust estimators using data from our sample.

Conduct Diagnostics to Assess Key Assumptions and Robustness of Findings

With a clear research design, researchers should critically evaluate key identification assumptions and test the robustness of findings when these assumptions are violated. Event-study plots, available for TWFE and most

HTE-robust estimators, are valuable tools for assessing whether the no-anticipation and PT assumptions hold. We recommend creating an event-study plot using the chosen estimator(s), followed by both visual inspection and statistical tests to assess how plausible the PT assumption is. Importantly, the absence of statistical significance in pretrend coefficients should not be taken as conclusive evidence for the validity of PT. To avoid the conditional inference problem, we recommend performing a sensitivity analysis with robust confidence sets across different values of \bar{M} , regardless of the pretest results. As shown in this article, such tests require more statistical power than rejecting the null that the average treatment effect is zero. Researchers should, therefore, allocate sufficient statistical power for these diagnostic tests during the research design phase.

Panel data provide valuable opportunities for social scientists to tackle complex causal questions; however, these data, especially when analyzed under PT, present distinct challenges. Our findings are not intended to dissuade researchers from employing PT-based research designs in causal panel analysis. Rather, our aim is to guide researchers in conducting their analyses more transparently and credibly. To facilitate this, we have integrated all the procedures described in this article into the open-source R package `fect`, and we offer detailed tutorials for these methods.

SUPPLEMENTARY MATERIAL

To view supplementary material for this article, please visit <https://doi.org/10.1017/S0003055425000243>.

DATA AVAILABILITY STATEMENT

Research documentation and data that support the findings of this study are openly available at the American Political Science Review Dataverse: <https://doi.org/10.7910/DVN/9RJFZF>.

ACKNOWLEDGMENTS

We thank Susan Athey, Quintin Beazer, Kirill Borusyak, Kirk Bansak, Clément de Chaisemartin, Xavier d'Hautfoeuille, Gary Cox, Avi Feller, Anthony Fowler, Francisco Garfias, Justin Grimm, Jens Hainmueller, Erin Hartman, Guido Imbens, Julia Payson, Annamaria Prati, Jonathan Roth, Pedro Sant'anna, Ye Wang, Luwei Ying, and participants at the 39th PolMeth, the 11th Asian PolMeth, Seminars at UCSD, UCLA, UC Berkeley, and University of Chicago for helpful comments and suggestions. We are also grateful to Neil Malhotra, Andrew Baker, Anton Strezhnev, and the Alethia platform for providing pre-publication reviews of this paper. Further thanks go to two anonymous reviewers and the APSR editor, Andrew Eggers, for their helpful comments. We appreciate the authors of the studies we reviewed for their

constructive feedback and for making their data publicly available. We thank Anran Liu, Tianzhu Qin, and Jinwen Wu for excellent research assistance.

CONFLICT OF INTERESTS

The authors declare no ethical issues or conflict of interests in this research.

ETHICAL STANDARDS

The authors affirm this research did not involve human participants.

REFERENCES

- Arkhangelsky, Dmitry, and Guido Imbens. 2024. "Causal Models for Longitudinal and Panel Data: A Survey." *The Econometrics Journal* 27 (3): C1–61.
- Athey, Susan, and Guido W. Imbens. 2022. "Design-Based Analysis in Difference-in-Differences Settings with Staggered Adoption." *Journal of Econometrics* 226 (1): 62–79.
- Bai, Jushan, and Serena Ng. 2021. "Matrix Completion, Counterfactuals, and Factor Analysis of Missing Data." *Journal of the American Statistical Association* 116 (536): 1746–63.
- Baker, Andrew, Brantly Callaway, Scott Cunningham, Andrew Goodman-Bacon, and Pedro H.C. Sant'Anna. 2025. "Difference-in-Differences Designs: A Practitioner's Guide." *arXiv*. [arXiv: 2503.13323](https://arxiv.org/abs/2503.13323).
- Baker, Andrew C., David F. Larcker, and Charles C.Y. Wang. 2022. "How Much Should We Trust Staggered Difference-in-Differences Estimates?" *Journal of Financial Economics* 144 (2): 370–95.
- Beazer, Quintin H., and Ora John Reuter. 2022. "Do Authoritarian Elections Help the Poor? Evidence from Russian Cities." *The Journal of Politics* 84 (1): 437–54.
- Beck, Nathaniel, and Jonathan N. Katz. 1995. "What to Do (and Not to Do) with Time-Series Cross-Section Data." *American Political Science Review* 89 (3): 634–47.
- Beck, Nathaniel, and Jonathan N. Katz. 2001. "Throwing Out the Baby with the Bath Water: A Comment on Green, Kim, and Yoon." *International Organization* 55 (2): 487–95.
- Bertrand, Marianne, Esther Duflo, and Sendhil Mullainathan. 2004. "How Much Should We Trust Differences-in-Differences Estimates?" *Quarterly Journal of Economics* 119 (1): 249–75.
- Bischof, Daniel, and Markus Wagner. 2019. "Do Voters Polarize When Radical Parties Enter Parliament?" *American Journal of Political Science* 63 (4): 888–904.
- Bisgaard, Martin, and Rune Slothuus. 2018. "Partisan Elites as Culprits? How Party Cues Shape Partisan Perceptual Gaps." *American Journal of Political Science* 62 (2): 456–69.
- Blackwell, Matthew, and Adam Glynn. 2018. "How to Make Causal Inferences with Time-Series Cross-Sectional Data under Selection on Observables." *American Political Science Review* 112 (2): 1067–82.
- Blair, Graeme, Darin Christensen, and Valerie Wirtschafter. 2022. "How Does Armed Conflict Shape Investment? Evidence from the Mining Sector." *The Journal of Politics* 84 (1): 116–33.
- Bokobza, Laure, Suthan Krishnarajan, Jacob Nyrop, Casper Sakstrup, and Lasse Aaskoven. 2022. "The Morning After: Cabinet Instability and the Purging of Ministers after Failed Coup Attempts in Autocracies." *The Journal of Politics* 84 (3): 1437–52.
- Borusyak, Kirill, Xavier Jaravel, and Jann Spiess. 2024. "Revisiting Event-Study Designs: Robust and Efficient Estimation." *Review of Economic Studies* 91 (6): 3253–85.
- Callaway, Brantly, and Pedro H.C. Sant'Anna. 2021. "Difference-in-Differences with Multiple Time Periods." *Journal of Econometrics* 225 (2): 200–30.

- Cameron, A. Colin, Jonah B. Gelbach, and Douglas L. Miller. 2008. "Bootstrap-Based Improvements for Inference with Clustered Errors." *Review of Economics and Statistics* 90 (3): 414–27.
- Caughey, Devin, Christopher Warshaw, and Yiqing Xu. 2017. "Incremental Democracy: The Policy Effects of Partisan Control of State Government." *The Journal of Politics* 79 (4): 1342–58.
- Cengiz, Doruk, Arindrajit Dube, Attila Lindner, and Ben Zipperer. 2019. "The Effect of Minimum Wages on Low-Wage Jobs." *The Quarterly Journal of Economics* 134 (3): 1405–54.
- Chiu, Albert, Xingchen Lan, Ziyi Liu, and Yiqing Xu. 2025. "Replication Data for: Causal Panel Analysis under Parallel Trends: Lessons from a Large Reanalysis Study." Harvard Dataverse. Dataset. <https://doi.org/10.7910/DVN/9RJFZF>.
- Christensen, Darin, and Francisco Garfias. 2021. "The Politics of Property Taxation: Fiscal Infrastructure and Electoral Incentives in Brazil." *The Journal of Politics* 83 (4): 1399–416.
- Clarke, Andrew J. 2020. "Party Sub-Brands and American Party Factions." *American Journal of Political Science* 64 (3): 452–70.
- Clayton, Amanda, and Pär Zetterberg. 2018. "Quota Shocks: Electoral Gender Quotas and Government Spending Priorities Worldwide." *The Journal of Politics* 80 (3): 916–32.
- Cox, Gary W., and Mark Dincecco. 2021. "The Budgetary Origins of Fiscal-Military Prowess." *The Journal of Politics* 83 (3): 851–66.
- Dahlström, Carl, and Mikael Holmgren. 2023. "Loyal Leaders, Affluent Agencies: The Budgetary Implications of Political Appointments in the Executive Branch." *The Journal of Politics* 85 (2): 640–53.
- de Chaisemartin, Clément, and Xavier D'Haultfœuille. 2023. "Credible Answers to Hard Questions: Differences-in-Differences for Natural Experiments." Working Paper, SSRN.
- de Chaisemartin, Clément, and Xavier D'Haultfœuille. 2020. "Two-Way Fixed Effects Estimators with Heterogeneous Treatment Effects." *American Economic Review* 110 (9): 2964–96.
- de Chaisemartin, Clément, and Xavier d'Haultfœuille. 2024. "Difference-in-Differences Estimators of Intertemporal Treatment Effects." *Review of Economics and Statistics*, 1–45. https://doi.org/10.1162/rest_a_01414.
- Dippoppa, Gemma, Guy Grossman, and Stephanie Zonszein. 2023. "Locked Down, Lashing Out: COVID-19 Effects on Asian Hate Crimes in Italy." *The Journal of Politics* 85 (2): 389–404.
- Distelhorst, Greg, and Richard M. Locke. 2018. "Does Compliance Pay? Social Standards and Firm-Level Trade." *American Journal of Political Science* 62 (3): 695–711.
- Eckhouse, Laurel. 2022. "Metrics Management and Bureaucratic Accountability: Evidence from Policing." *American Journal of Political Science* 66 (2): 385–401.
- Eggers, Andrew C., Anthony Fowler, Jens Hainmueller, Andrew B. Hall, and James M. Snyder Jr. 2015. "On the Validity of the Regression Discontinuity Design for Estimating Electoral Effects: New Evidence from Over 40,000 Close Races." *American Journal of Political Science* 59 (1): 259–74.
- Esberg, Jane, and Alexandra A. Siegel. 2023. "How Exile Shapes Online Opposition: Evidence from Venezuela." *American Political Science Review* 117 (4): 1361–78.
- Fouirnaies, Alexander. 2018. "When Are Agenda Setters Valuable?" *American Journal of Political Science* 62 (1): 176–91.
- Fouirnaies, Alexander, and Andrew B. Hall. 2018. "How Do Interest Groups Seek Access to Committees?" *American Journal of Political Science* 62 (1): 132–47.
- Fouirnaies, Alexander, and Andrew B. Hall. 2022. "How Do Electoral Incentives Affect Legislator Behavior? Evidence from U.S. State Legislatures." *American Political Science Review* 116 (2): 662–76.
- Fresh, Adriane. 2018. "The Effect of the Voting Rights Act on Enfranchisement: Evidence from North Carolina." *The Journal of Politics* 80 (2): 713–18.
- Fu, Jiawei, Cyrus Samii, and Ye Wang. 2024. "Inference for Group Interaction Experiments." Paper presented at the Annual Meeting of the American Political Science Association.
- Garfias, Francisco. 2019. "Elite Coalitions, Limited Government, and Fiscal Capacity Development: Evidence from Bourbon Mexico." *The Journal of Politics* 81 (1): 95–111.
- Goodman-Bacon, Andrew. 2021. "Difference-in-Differences with Variation in Treatment Timing." *Journal of Econometrics* 225 (2): 254–77.
- Green, Donald P., Soo Yeon Kim, and David H. Yoon. 2001. "Dirty Pool." *International Organization* 55 (2): 441–68.
- Grumbach, Jacob M. 2023. "Laboratories of Democratic Backsliding." *American Political Science Review* 117 (3): 967–84.
- Grumbach, Jacob M., and Alexander Sahn. 2020. "Race and Representation in Campaign Finance." *American Political Science Review* 114 (1): 206–21.
- Grumbach, Jacob M., and Charlotte Hill. 2022. "Rock the Registration: Same Day Registration Increases Turnout of Young Voters." *The Journal of Politics* 84 (1): 405–17.
- Hainmueller, Jens, and Dominik Hangartner. 2019. "Does Direct Democracy Hurt Immigrant Minorities? Evidence from Naturalization Decisions in Switzerland." *American Journal of Political Science* 63 (3): 530–47.
- Hainmueller, Jens, Jonathan Mummolo, and Yiqing Xu. 2019. "How Much Should We Trust Estimates from Multiplicative Interaction Models? Simple Tools to Improve Empirical Practice." *Political Analysis* 27 (2): 163–92.
- Hall, Andrew B., and Jesse Yoder. 2022. "Does Homeownership Influence Political Behavior? Evidence from Administrative Data." *The Journal of Politics* 84 (1): 351–66.
- Hankinson, Michael, and Asya Magazinnik. 2023. "The Supply-Equity Trade-Off: The Effect of Spatial Representation on the Local Housing Supply." *The Journal of Politics* 85 (3): 1033–47.
- Hirano, Shigeo, Jaclyn Kaslovsy, Michael P. Olson, and James M. Snyder. 2022. "The Growth of Campaign Advertising in the United States, 1880–1930." *The Journal of Politics* 84 (3): 1482–96.
- Imai, Kosuke, and In Song Kim. 2019. "When Should We Use Unit Fixed Effects Regression Models for Causal Inference with Longitudinal Data?" *American Journal of Political Science* 63 (2): 467–90.
- Imai, Kosuke, In Song Kim, and Erik H. Wang. 2023. "Matching Methods for Causal Inference with Time-Series Cross-Sectional Data." *American Journal of Political Science* 67 (3): 587–605.
- Jiang, Junyan. 2018. "Making Bureaucracy Work: Patronage Networks, Performance Incentives, and Economic Development in China." *American Journal of Political Science* 62 (4): 982–99.
- Kahn-Lang, Ariella, and Kevin Lang. 2020. "The Promise and Pitfalls of Differences-in-Differences: Reflections on 16 and Pregnant and Other Applications." *Journal of Business & Economic Statistics* 38 (3): 613–20.
- Kilborn, Mitchell, and Arjun Vishwanath. 2022. "Public Money Talks Too: How Public Campaign Financing Degrades Representation." *American Journal of Political Science* 66 (3): 730–44.
- King, Gary. 2001. "Proper Nouns and Methodological Propriety: Pooling Dyads in International Relations Data." *International Organization* 55 (2): 497–507.
- Kogan, Vladimir. 2021. "Do Welfare Benefits Pay Electoral Dividends? Evidence from the National Food Stamp Program Rollout." *The Journal of Politics* 83 (1): 20–70.
- Kroeger, Mary, and Maria Silfa. 2023. "Motivated Corporate Political Action: Evidence from an SEC Experiment." *The Journal of Politics* 85 (3): 1139–44.
- Kuipers, Nicholas, and Alexander Sahn. 2023. "The Representational Consequences of Municipal Civil Service Reform." *American Political Science Review* 117 (1): 200–16.
- Lal, Apoorva, Mackenzie Lockhart, Yiqing Xu, and Ziwen Zu. 2024. "How Much Should We Trust Instrumental Variable Estimates in Political Science? Practical Advice Based on 67 Replicated Studies." *Political Analysis* 32 (4): 521–40.
- Lall, Ranjit. 2016. "How Multiple Imputation Makes a Difference." *Political Analysis* 24 (4): 414–33.
- Latura, Audrey, and Ana Catalano Weeks. 2023. "Corporate Board Quotas and Gender Equality Policies in the Workplace." *American Journal of Political Science* 67 (3): 606–22.
- Liao, Steven. 2023. "The Effect of Firm Lobbying on High-Skilled Visa Adjudication." *The Journal of Politics* 85 (4): 1416–29.
- Liu, Licheng, Ye Wang, and Yiqing Xu. 2024. "A Practical Guide to Counterfactual Estimators for Causal Inference with Time-Series Cross-Sectional Data." *American Journal of Political Science* 68 (1): 160–76.
- Magaloni, Beatriz, Edgar Franco-Vivanco, and Vanessa Melo. 2020. "Killing in the Slums: Social Order, Criminal Governance, and Police Violence in Rio de Janeiro." *American Political Science Review* 114 (2): 552–72.

- Marsh, Wayne Z.C. 2023. "Trauma and Turnout: The Political Consequences of Traumatic Events." *American Political Science Review* 117 (3): 1036–52.
- Mou, Hongyu, Licheng Liu, and Yiqing Xu. 2023. "Panel Data Visualization in R (panelView) and Stata (panelview)." *Journal of Statistical Software* 107 (7): 1–20.
- Paglayan, Agustina S. 2022. "Education or Indoctrination? The Violent Origins of Public School Systems in an Era of State-Building." *American Political Science Review* 116 (4): 1242–57.
- Payson, Julia A. 2020a. "The Partisan Logic of City Mobilization: Evidence from State Lobbying Disclosures." *American Political Science Review* 114 (3): 677–90.
- Payson, Julia A. 2020b. "Cities in the Statehouse: How Local Governments Use Lobbyists to Secure State Funding." *The Journal of Politics* 82 (2): 403–17.
- Pierskalla, Jan H., and Audrey Sacks. 2018. "Unpaved Road Ahead: The Consequences of Election Cycles for Capital Expenditures." *The Journal of Politics* 80 (2): 510–24.
- Rambachan, Ashesh, and Jonathan Roth. 2023. "A More Credible Approach to Parallel Trends." *The Review of Economic Studies* 90 (5): 2555–91.
- Ravanilla, Nico, Renard Sexton, and Dotan Haim. 2022. "Deadly Populism: How Local Political Outsiders Drive Duterte's War on Drugs in the Philippines." *The Journal of Politics* 84 (2): 1035–56.
- Roth, Jonathan. 2022. "Pretest with Caution: Event-Study Estimates after Testing for Parallel Trends." *American Economic Review: Insights* 4 (3): 305–22.
- Roth, Jonathan. 2024. "Interpreting Event-Studies from Recent Difference-in-Differences Methods." Working Paper.
- Roth, Jonathan, and Pedro H.C. Sant'Anna. 2023. "When is Parallel Trends Sensitive to Functional Form?" *Econometrica* 91 (2): 737–47.
- Roth, Jonathan, Pedro H.C. Sant'Anna, Alyssa Bilinski, and John Poe. 2023. "What's Trending in Difference-in-Differences? A Synthesis of the Recent Econometrics Literature." *Journal of Econometrics* 235 (2): 2218–44.
- Sanford, Luke. 2023. "Democratization, Elections, and Public Goods: The Evidence from Deforestation." *American Journal of Political Science* 67 (3): 748–63.
- Sant'Anna, Pedro H.C., and Jun B. Zhao. 2020. "Doubly Robust Difference-in-Differences Estimators." *Journal of Econometrics* 219 (1): 101–22.
- Schafer, Jerome, Enrico Cantoni, Giorgio Bellettini, and Carlotta Berti Ceroni. 2022. "Making Unequal Democracy Work? The Effects of Income on Voter Turnout in Northern Italy." *American Journal of Political Science* 66 (3): 745–61.
- Schubiger, Livia Isabella. 2021. "State Violence and Wartime Civilian Agency: Evidence from Peru." *The Journal of Politics* 83 (4): 1383–98.
- Schuit, Sophie, and Jon C. Rogowski. 2017. "Race, Representation, and the Voting Rights Act." *American Journal of Political Science* 61 (3): 513–26.
- Skorge, Øyvind Søråas. 2023. "Mobilizing the Underrepresented: Electoral Systems and Gender Inequality in Political Participation." *American Journal of Political Science* 67 (3): 538–52.
- Strezhnev, Anton. 2018. "Semiparametric Weighting Estimators for Multi-Period Difference-in-Differences Designs." Working Paper, New York University.
- Sun, Liyang, and Sarah Abraham. 2021. "Estimating Dynamic Treatment Effects in Event Studies with Heterogeneous Treatment Effects." *Journal of Econometrics* 225 (2): 175–99.
- Trounstein, Jessica. 2020. "The Geography of Inequality: How Land Use Regulation Produces Segregation." *American Political Science Review* 114 (2): 443–55.
- Wang, Ye. 2021. "Causal Inference under Temporal and Spatial Interference." Working Paper, New York University.
- Wang, Ye, Cyrus Samii, Haoge Chang, and P. M. Aronow. 2025. "Design-Based Inference for Spatial Experiments under Unknown Interference." *The Annals of Applied Statistics* 19 (1): 744–68.
- Weschle, Simon. 2021. "Parliamentary Positions and Politicians' Private Sector Earnings: Evidence from the UK House of Commons." *The Journal of Politics* 83 (2): 706–21.
- Wing, Coady, Seth M. Freedman, and Alex Hollingsworth. 2024. "Stacked Difference-in-Differences." Working Paper, NBER.
- Wooldridge, Jeffrey M. 2010. *Econometric Analysis of Cross Section and Panel Data*. Cambridge, MA: MIT Press.
- Xu, Yiqing. 2023. "Causal Inference with Time-Series Cross-Sectional Data: A Reflection." In *The Oxford Handbook of Engaged Methodological Pluralism in Political Science*, eds. Janet M. Box-Steffensmeier, Dino P. Christenson, and Valeria Sinclair-Chapman, Oxford: Oxford University Press. <https://doi.org/10.1093/oxfordhb/9780192868282.013.30>.
- Xu, Yiqing, Anqi Zhao, and Peng Ding. 2024. "Factorial Difference-in-Differences." *arXiv*, [arXiv:2407.11937](https://arxiv.org/abs/2407.11937).
- Zhang, Qi, Dong Zhang, Mingxing Liu, and Victor Shih. 2021. "Elite Cleavage and the Rise of Capitalism under Authoritarianism: A Tale of Two Provinces in China." *The Journal of Politics* 83 (3): 1010–23.