

Two-way fixed effects instrumental variable regressions in staggered DID-IV designs.*

Miyaji Sho[†]

May 28, 2024

Abstract

Many studies run two-way fixed effects instrumental variable (TWFEIV) regressions, leveraging variation in the timing of policy adoption across units as an instrument for treatment. This paper studies the properties of the TWFEIV estimator in staggered instrumented difference-in-differences (DID-IV) designs. We show that in settings with the staggered adoption of the instrument across units, the TWFEIV estimator can be decomposed into a weighted average of all possible two-group/two-period Wald-DID estimators. Under staggered DID-IV designs, a causal interpretation of the TWFEIV estimand hinges on the stable effects of the instrument on the treatment and the outcome over time. We illustrate the use of our decomposition theorem for the TWFEIV estimator through an empirical application.

arXiv:2405.16467v1 [econ.EM] 26 May 2024

*I am grateful to Daiji Kawaguchi and Ryo Okui for their continued guidance and support. All errors are my own.

[†]Graduate School of Economics, The University of Tokyo, 7-3-1 Hongo, Bunkyo, Tokyo 113-0033, Japan; Email: shomiyaji-apple@g.ecc.u-tokyo.ac.jp.

1 Introduction

Instrumented difference-in-differences (DID-IV) is a method to estimate the effect of a treatment on an outcome, exploiting variation in the timing of policy adoption across units as an instrument for the treatment. In a simple setting with two groups and two periods, some units become exposed to the policy shock in the second period (exposed group), whereas others are not over two periods (unexposed group). The estimator is constructed by running the following IV regression with the group and post-time dummies as included instruments and the interaction of the two as the excluded instrument (e.g., [Duflo \(2001\)](#), [Field \(2007\)](#)):

$$Y_{i,t} = \beta_0 + \beta_{i,\cdot}\text{Exposed}_i + \beta_{\cdot,t}\text{POST}_t + \beta_{IV}D_{i,t} + \epsilon_{i,t}.$$

The resulting IV estimand β_{IV} scales the DID estimand of the outcome by the DID estimand of the treatment, the so-called Wald-DID estimand ([de Chaisemartin and D’Haultfoeuille \(2018\)](#), [Miyaji \(2024\)](#)). In this two-group/two-period (2×2) setting, DID-IV designs mainly consist of a monotonicity assumption and parallel trends assumptions in the treatment and the outcome between the two groups, and allow for the Wald-DID estimand to capture the local average treatment effect on the treated (LATET) ([de Chaisemartin \(2010\)](#), [Hudson et al. \(2017\)](#), and [Miyaji \(2024\)](#)). DID-IV designs have gained popularity over DID designs in practice when there is no control group or the treatment adoption is potentially endogenous over time ([Miyaji \(2024\)](#)).

In reality, however, most DID-IV applications go beyond the canonical DID-IV set up, and leverage variation in the timing of policy adoption across units in more than two periods, instrumenting for the treatment with the natural variation. The instrument is constructed, for instance from the staggered adoption of school reforms across countries or municipalities (e.g. [Oreopoulos \(2006\)](#), [Lundborg et al. \(2014\)](#), [Meghir et al. \(2018\)](#)), the phase-in introduction of head starts across states (e.g. [Johnson and Jackson \(2019\)](#)), or the gradual adoption of broadband internet programs (e.g. [Akerman et al. \(2015\)](#), [Bhuller et al. \(2013\)](#)). These policy changes can be viewed as some natural experiments, but not randomized in reality.

Recently, [Miyaji \(2024\)](#) formalizes the underlying identification strategy as a staggered DID-IV design. In this design, the treatment adoption is allowed to be endogenous over time, while the instrument is required to be uncorrelated with time-varying unobservables in the treatment and the outcome; the assignment of the treatment can be non-staggered across units, while the assignment of the instrument is staggered across units: they are partitioned into mutually exclusive and exhaustive cohorts by the initial adoption date of the instrument. The target parameter is the cohort specific local average treatment effect on the treated (CLATT); this parameter measures the treatment effects among the units who belong to cohort e and are induced to the treatment by instrument in a given relative period l after the initial adoption of the instrument. The identifying assumptions are the natural generalization of those in 2×2 DID-IV designs.

In practice, empirical researchers commonly implement this design via linear instrumental variable regressions with time and unit fixed effects, the so-called two-way fixed effects instrumental variable (TWFEIV) regressions (e.g., [Black et al. \(2005\)](#), [Lundborg et al. \(2017\)](#), [Johnson and Jackson \(2019\)](#)):

$$Y_{i,t} = \phi_i + \lambda_t + \beta_{IV}D_{i,t} + v_{i,t}, \tag{1}$$

$$D_{i,t} = \gamma_i + \zeta_t + \pi Z_{i,t} + \eta_{i,t}. \tag{2}$$

In contrast to the canonical DID-IV set up, however, the validity of running TWFEIV regressions seems less clear under staggered DID-IV designs. The IV estimate is commonly interpreted

as measuring the local average treatment effect in the presence of heterogeneous treatment effects as in [Imbens and Angrist \(1994\)](#), whereas the target parameter is not stated formally. We know little about how the IV estimator is constructed by comparing the evolution of the treatment and the outcome across units and over time. Finally, we have no tools to illustrate the identifying variations in the IV estimate in a given application.

In this paper, we study the properties of two-way fixed effects instrumental variable estimators under staggered DID-IV designs. Specifically, we present the decomposition result for the TWFEIV estimator, and study the causal interpretation of the TWFEIV estimand under staggered DID-IV designs.

First, we derive the decomposition theorem for the TWFEIV estimator with settings of the staggered adoption of the instrument across units. We show that the TWFEIV estimator is equal to a weighted average of all possible 2×2 Wald-DID estimators arising from the three types of the DID-IV design. First, in an Unexposed/Exposed design, some units are never exposed to the instrument during the sample period (unexposed group), whereas some units start exposed at a particular date and remain exposed (exposed group). Second, in an Exposed/Not Yet Exposed design, some units start exposed earlier, whereas some units are not yet exposed during the design period (not yet exposed group). Finally, in an Exposed/Exposed Shift design, some units are already exposed, whereas some units start exposed later at a particular point during the design period (exposed shift group). The weight assigned to each Wald-DID estimator reflects all the identifying variations in each DID-IV design: the sample share, the variance of the instrument, and the DID estimator of the treatment between the two groups.

Built on the decomposition result, we next uncover the shortcomings of running TWFEIV regressions under staggered DID-IV designs. We show that the TWFEIV estimand potentially fails to summarize the treatment effects under staggered DID-IV designs due to negative weights. Specifically, we show that this estimand is equal to a weighted average of all possible cohort specific local average treatment effect on the treated (CLATT) parameters, but some weights can be negative. The negative weight problem potentially arises due to the "bad comparisons" (c.f. [Goodman-Bacon \(2021\)](#)) performed by TWFEIV regressions: the already exposed units play the role of controls in the Exposed/Exposed Shift design in the first stage and reduced form regressions. Given the negative result of using the TWFEIV estimand under staggered DID-IV designs, we also investigate the sufficient conditions for this estimand to attain its causal interpretation. We show that this estimand can be interpreted as causal only if the effects of the instrument on the treatment and the outcome are stable over time.

We extend our decomposition result in several directions. We first consider non-binary, ordered treatment. We also derive the decomposition result for the TWFEIV estimand in unbalanced panel settings. Lastly, we consider the case when the adoption date of the instrument is randomized across units. In all cases, we show that the TWFEIV estimand potentially fails to summarize the treatment effects under staggered DID-IV designs due to negative weights.

We illustrate our findings with the setting of [Miller and Segal \(2019\)](#) who estimate the effect of female police officers' share on intimate partner homicide rate, leveraging the timing variation of AA (affirmative action) plans across U.S. counties. In this application, we first assess the plausibility of the staggered DID-IV design implicitly imposed by [Miller and Segal \(2019\)](#) and confirm its validity. We then estimate TWFEIV regressions, slightly modifying the authors' setting, and apply our DID-IV decomposition theorem to the IV estimate. We find that the estimate assigns more weights to the Unexposed/Exposed design and less weights to the other two types of the DID-IV design. Despite the small weight on the Exposed/Exposed Shift design, we also find that the IV estimate suffers from the substantial downward bias

arising from the bad comparisons in the Exposed/Exposed Shift design.

Finally, we develop simple tools to examine how different specifications affect the change in TWFEIV estimates, and illustrate these by revisiting [Miller and Segal \(2019\)](#). In many empirical settings, researchers typically diverge from a simple TWFEIV regression as in equation (1) and estimate various specifications such as weighting or including time-varying covariates. We follow [Goodman-Bacon \(2021\)](#) and decompose the difference between the two specifications into the changes in Wald-DID estimates, the changes in weights, and the interaction of the two. This decomposition result enables the researchers to quantify the contribution of the changes in each term to the difference in the overall estimates. In addition, plotting the pairs of Wald-DID estimates and associated weights obtained from the two specifications allows the researchers to investigate which components have the significant impact on these contributions.

Overall, this paper shows the negative result of using TWFEIV estimators under staggered DID-IV designs in more than two periods, and provide tools to illustrate how serious that concern is in a given application. Specifically, our decomposition result for the TWFEIV estimator enables the researchers to quantify the bias term arising from the bad comparisons in Exposed/Exposed Shift designs in the data. Fortunately, [Miyaji \(2024\)](#) recently proposes the alternative estimation method in staggered DID-IV designs that is robust to treatment effect heterogeneity. Using such estimation method allows the practitioners to avoid the issue of TWFEIV estimators in practice, and facilitates the credibility of their empirical findings.

The rest of the paper is organized as follows. The next subsection discusses the related literature. Section 2 presents our decomposition theorem for the TWFEIV estimator. Section 3 formally introduces staggered instrumented difference-in-differences designs. Section 4 presents the pitfalls of running TWFEIV regressions under staggered DID-IV designs, and explores the sufficient conditions for the TWFEIV estimand to attain its causal interpretation. Section 5 describes some of the extensions. Section 6 presents our empirical application. Section 7 explain how different specifications affect the difference in estimates and Section 8 concludes. All proofs are given in the Appendix.

1.1 Related literature

Our paper is related to the recent DID-IV literature ([de Chaisemartin \(2010\)](#); [Hudson et al. \(2017\)](#); [de Chaisemartin and D’Haultfœuille \(2018\)](#); [Miyaji \(2024\)](#)). In this literature, [de Chaisemartin \(2010\)](#) first formalizes 2×2 DID-IV designs and shows that a Wald-DID estimand identifies the local average treatment effect on the treated (LATET) if the parallel trends assumptions in the treatment and the outcome, and a monotonicity assumption are satisfied. [Hudson et al. \(2017\)](#) also consider 2×2 DID-IV designs with non-binary, ordered treatment settings. Build on the work in [de Chaisemartin \(2010\)](#), however, [de Chaisemartin and D’Haultfœuille \(2018\)](#) formalize 2×2 DID-IV designs differently, and call them Fuzzy DID. [Miyaji \(2024\)](#) compares 2×2 DID-IV to Fuzzy DID designs and points out the issues embedded in Fuzzy DID designs, and extends 2×2 DID-IV design to multiple period settings with the staggered adoption of the instrument across units, which the author calls staggered DID-IV designs. [Miyaji \(2024\)](#) also provides a reliable estimation method in staggered DID-IV designs that is robust to treatment effect heterogeneity.

In this paper, we contribute to the literature by showing the properties of two-way fixed instrumental variable estimators in staggered DID-IV designs. In reality, when empirical researchers implicitly rely on the staggered DID-IV design, they commonly implement this design via TWFEIV regressions (e.g. [Black et al. \(2005\)](#), [Lundborg et al. \(2014\)](#), [Meghir et al. \(2018\)](#)). This paper presents the issues of the conventional approach, and provides the sufficient conditions for this estimand to attain its causal interpretation.

Our paper is also related to a recent DID literature on the causal interpretation of two-way fixed effects (TWFE) regressions and its dynamic specifications under heterogeneous treatment effects (Athey and Imbens (2022); Borusyak et al. (2021); de Chaisemartin and D’Haultfoeulle (2020); Goodman-Bacon (2021); Imai and Kim (2021); Sun and Abraham (2021)).

Specifically, this paper is closely connected to Goodman-Bacon (2021), who derives the decomposition theorem for the TWFE estimator with settings of the staggered adoption of the treatment across units. In this paper, we establish the decompose theorem for the TWFEIV estimator with settings of the staggered adoption of the instrument across units, which is a natural generalization of their theorem 1.

This paper is also closely connected to de Chaisemartin and D’Haultfoeulle (2020), who decompose the TWFE estimand and present the issue of using this estimand under DID designs: some weights assigned to the causal parameters in this estimand can be potentially negative. In their appendix, the authors also decompose the TWFEIV estimand and refer to the negative weight problem in this estimand. Specifically, they apply the decomposition theorem for the TWFE estimand to the numerator and denominator in the TWFEIV estimand respectively, and conclude that this estimand identifies the LATE as in Imbens and Angrist (1994) only if the effects of the instrument on the treatment and outcome are constant across groups and over time. However, their decomposition result for the TWFEIV estimand has some drawbacks. First, they do not formally state the target parameter and identifying assumptions in DID-IV designs. Second, their decomposition result is not based on the target parameter in DID-IV designs. Finally, the sufficient conditions for this estimand to be interpretable causal parameter are not well investigated.

In this paper, we investigate the causal interpretation of the TWFEIV estimand more clearly than that of de Chaisemartin and D’Haultfoeulle (2020). Specifically, we first decompose the TWFEIV estimator into all possible 2×2 Wald-DID estimators. We then formally introduce the target parameter and identifying assumptions in staggered DID-IV designs, built on the recent work in Miyaji (2024). This allows us to decompose the TWFEIV estimand into a weighted average of the target parameter in staggered DID-IV designs. Finally, we assess the causal interpretation of the TWFEIV estimand under a variety of restrictions on the effects of the instrument on the treatment and outcome, which clarifies the sufficient conditions for this estimand to attain its causal interpretation.

We note that this paper is distinct from the recent IV literature on the causal interpretation of two stage least square (TSLS) estimators with covariates under heterogeneous treatment effects (Słoczyński (2020), Blandhol et al. (2022)). These recent studies investigate the causal interpretation of the TSLS estimand with covariates under the random variation of the instrument conditional on covariates, and cast doubt on the LATE (or LATEs) interpretation of this estimand. In this literature, the identifying variations come from the assignment process of the instrument. In this paper, however, we investigate the causal interpretation of the TWFEIV estimand (where time and unit dummies can be viewed as covariates) under staggered DID-IV designs: our identifying variations mainly come from the parallel trends assumptions in the treatment and the outcome over time.

2 Instrumented difference-in-differences decomposition

In this section, we present a decomposition result for the two-way fixed effects instrumental variable (TWFEIV) estimator in multiple time period settings with the staggered adoption of the instrument across units.

2.1 Set up

We introduce the notation we use throughout this article. We consider a panel data setting with T periods and N units. For each $i \in \{1, \dots, N\}$ and $t \in \{1, \dots, T\}$, let $Y_{i,t}$ denote the outcome and $D_{i,t} \in \{0, 1\}$ denote the treatment status, and $Z_{i,t} \in \{0, 1\}$ denote the instrument status. Let $D_i = (D_{i,1}, \dots, D_{i,T})$ and $Z_i = (Z_{i,1}, \dots, Z_{i,T})$ denote the path of the treatment and the path of the instrument for unit i , respectively. Throughout this article, we assume that $\{Y_{i,t}, D_{i,t}, Z_{i,t}\}_{t=1}^T$ are independent and identically distributed (i.i.d).

We make the following assumption about the assignment process of the instrument.

Assumption 1 (Staggered adoption for $Z_{i,t}$). For $s < t$, $Z_{i,s} \leq Z_{i,t}$ where $s, t \in \{1, \dots, T\}$.

Assumption 1 requires that once units start exposed to the instrument, they remain exposed to that instrument afterward. In the DID literature, several recent papers impose this assumption on the adoption process of the treatment and sometimes call it the "staggered treatment adoption", see, e.g., [Athey and Imbens \(2022\)](#), [Callaway and Sant'Anna \(2021\)](#) and [Sun and Abraham \(2021\)](#).

Given Assumption 1, we can uniquely characterize the instrument path by the time period when unit i is first exposed to the instrument, denoted by $E_i = \min\{t : Z_{i,t} = 1\}$. If unit i is not exposed to the instrument for all time periods, we define $E_i = \infty$. Based on the initial exposure period E_i , we can uniquely partition units into mutually exclusive and exhaustive cohorts e for $e \in \{1, 2, \dots, T, \infty\}$: all the units in cohort e are first exposed to the instrument at time $E_i = e$. Hereafter, to ease the notation, we assume that the data contain K cohorts ($K \leq T$) where $e \in \{1, \dots, k, \dots, K\}$, and define U as the never exposed cohort $E_i = \infty$.

Let n_e be the relative sample share for cohort e and let \bar{Z}_e be the time share of the exposure to the instrument for cohort e :

$$n_e \equiv \frac{\sum_i \mathbf{1}\{E_i = e\}}{N}, \quad \bar{Z}_e \equiv \frac{\sum_t \mathbf{1}\{t \geq e\}}{T}.$$

We also define $n_{ab} \equiv \frac{n_a}{n_a + n_b}$ to be the relative sample share between cohort a and b .

In contrast to the staggered adoption of the instrument across units, we allow the general adoption process for the treatment: the treatment can potentially turn on/off repeatedly over time. [de Chaisemartin and D'Haultfœuille \(2020\)](#) and [Imai and Kim \(2021\)](#) consider the same setting in the recent DID literature.

The notations $PRE(a)$, $MID(a, b)$, and $POST(a)$ represent the corresponding time window, respectively: $PRE(a) \equiv [1, a)$, $MID(a, b) \equiv [a, b)$, and $POST(a) \equiv [a, T]$. Let $\bar{R}_e^{POST(a)}$ be the sample mean of the random variable $R_{i,t}$ in cohort e during the time window $POST(a)$:

$$\bar{R}_e^{POST(a)} \equiv \frac{1}{T - (a - 1)} \sum_a^T \left[\frac{\sum_i R_{i,t} \mathbf{1}\{E_i = e\}}{\sum_i \mathbf{1}\{E_i = e\}} \right].$$

We define $\bar{R}_e^{PRE(a)}$ and $\bar{R}_e^{MID(a,b)}$ analogously, representing the sample mean of the random variable $R_{i,t}$ in cohort e during the time window $PRE(a)$ and $MID(a, b)$ respectively.

2.2 Decomposing the TWFEIV estimator

We consider a TWFEIV regression in multiple time period settings with the staggered adoption of the instrument across units:

$$Y_{i,t} = \phi_i + \lambda_t + \beta_{IV} D_{i,t} + v_{i,t}, \quad (3)$$

$$D_{i,t} = \gamma_i + \zeta_t + \pi Z_{i,t} + \eta_{i,t}. \quad (4)$$

By substituting the first stage regression (4) into the structural equation (3), we obtain the reduced form regression:

$$Y_{i,t} = \phi_i + \lambda_t + \alpha Z_{i,t} + v_{i,t}. \quad (5)$$

The ratio between the first stage coefficient $\hat{\pi}$ and the reduced form coefficient $\hat{\alpha}$ yields the TWFEIV estimator $\hat{\beta}_{IV}$. By the Frisch-Waugh-Lovell theorem, the IV estimator $\hat{\beta}_{IV}$ is equal to the ratio between the coefficient from regressing $Y_{i,t}$ on the double-demeaning variable $\tilde{Z}_{i,t}$ and the coefficient from regressing $D_{i,t}$ on the same variable:

$$\hat{\beta}_{IV} = \frac{\frac{1}{NT} \sum_i \sum_t \tilde{Z}_{i,t} Y_{i,t}}{\frac{1}{NT} \sum_i \sum_t \tilde{Z}_{i,t} D_{i,t}}, \quad (6)$$

where $\tilde{Z}_{i,t}$ is the double demeaning variable defined below:

$$\begin{aligned} \tilde{Z}_{i,t} &= Z_{i,t} - \frac{1}{T} \sum_{t=1}^T Z_{i,t} - \frac{1}{N} \sum_{i=1}^N Z_{i,t} + \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N Z_{i,t} \\ &\equiv (Z_{i,t} - \bar{Z}_i) - (\bar{Z}_t - \bar{\bar{Z}}). \end{aligned}$$

Note that the TWFEIV regression runs the two-way fixed effects (TWFE) regression twice, as can be seen in equations (4) and (5). Because we assume the staggered assignment of the instrument across units, if we focus on the TWFE coefficient on $Z_{i,t}$ in the first stage or the reduced form regression, we can show that it is equal to a weighted average of all possible 2×2 DID estimators of the treatment or the outcome from the decomposition result for the TWFE estimator shown by [Goodman-Bacon \(2021\)](#).

Consider the simple setting where we have only two periods and two cohorts: one cohort is not exposed to the instrument during the two periods ($E_i = U$), whereas the other cohort starts exposed to the instrument in the second period ($E_i = 2$). In this setting, the TWFEIV estimator takes the following form, the so-called Wald-DID estimator ([de Chaisemartin and D'Haultfoeuille \(2018\)](#), [Miyaji \(2024\)](#)):

$$\hat{\beta}_{IV} = \frac{\bar{Y}_{2,2} - \bar{Y}_{2,1} - (\bar{Y}_{U,2} - \bar{Y}_{U,1})}{\bar{D}_{2,2} - \bar{D}_{2,1} - (\bar{D}_{U,2} - \bar{D}_{U,1})},$$

where $\bar{R}_{a,t}$ is the sample mean of the random variable $R_{i,t}$ for cohort $E_i = a$ in time t . This estimator scales the DID estimator of the outcome by the DID estimator of the treatment between cohort $E_i = U$ and $E_i = 2$.

The above observations bring us the intuition about how we can decompose the TWFEIV estimator with settings of the staggered adoption of the instrument across units; we expect that the TWFEIV estimator can be decomposed into a weighted average of all possible 2×2 Wald-DID estimators (instead of DID-estimators).

To clarify this intuition, assume for now that we have only three cohorts, an early exposed cohort k , a middle exposed cohort l ($k < l$), and a never exposed cohort U ($E_i = \infty$). [Figure 1](#) plots the simulated data for the time trends of the average treatment (first stage) and the average outcome (reduced form) in three cohorts.

From the data structure, we can construct the Wald-DID estimator in three ways. First, we can compare the evolution of the treatment and the outcome between exposed cohort $j = k, l$

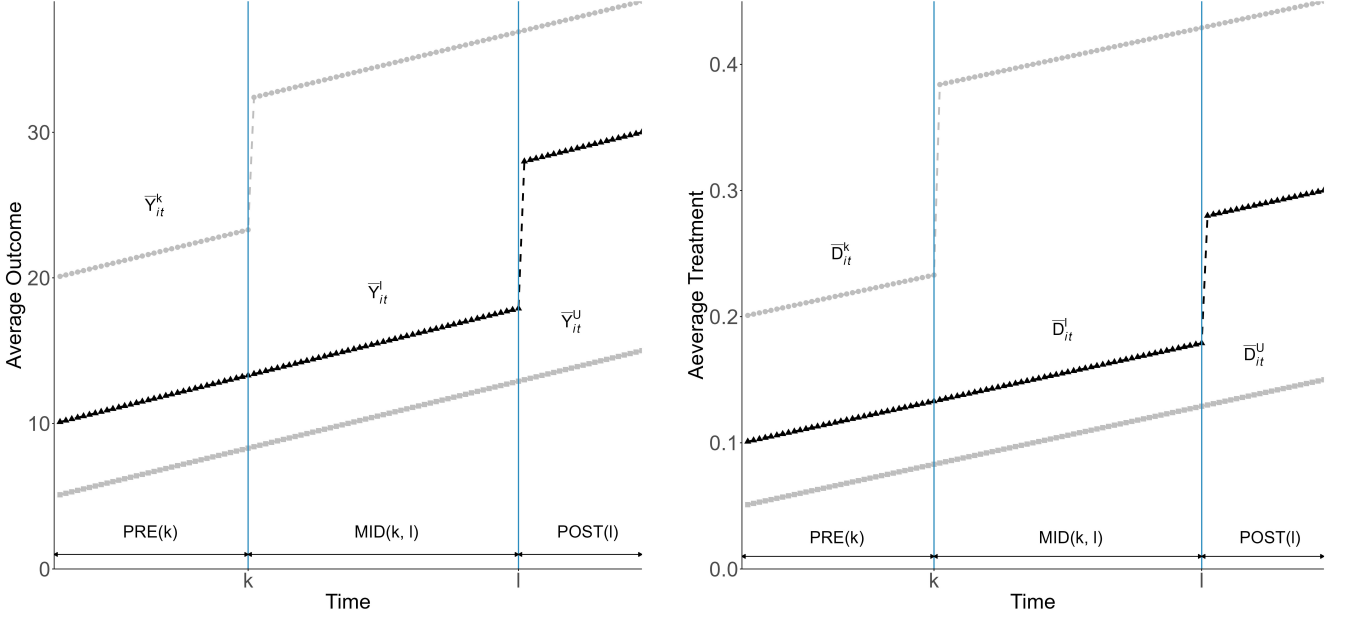


Fig. 1. Instrumented difference-in-differences with three cohorts. *Notes:* This figure plots the simulated data for the time trends of the average treatment (first stage) and the average outcome (reduced form) with time length $T = 100$ in three cohorts: an early exposed cohort k , which is exposed to the instrument at $k = \frac{34}{100}T$; a middle exposed cohort l , which is exposed to the instrument at $l = \frac{80}{100}T$; a never exposed cohort, U . The x -axis consists of three time windows: the pre-exposed period for cohort k , $[1, k - 1]$, denoted by $PRE(k)$; the middle exposed period when the cohort k is already exposed but cohort l is not yet exposed, $[k, l - 1]$, denoted by $MID(k, l)$; and post-exposed period when cohort l is already exposed, $[l, T]$, denoted by $POST(l)$. The effects of the instrument on the treatment and the outcome are 0.15 and 9 in cohort k respectively; 0.1 and 10 in cohort l respectively.

and never exposed cohort U , exploiting the time window $POST(j)$ and $PRE(j)$, which we call an Unexposed/Exposed design:

$$\begin{aligned} \hat{\beta}_{IV,jU}^{2 \times 2} &\equiv \frac{\left(\bar{y}_j^{POST(j)} - \bar{y}_j^{PRE(j)}\right) - \left(\bar{y}_U^{POST(j)} - \bar{y}_U^{PRE(j)}\right)}{\left(\bar{D}_j^{POST(j)} - \bar{D}_j^{PRE(j)}\right) - \left(\bar{D}_U^{POST(j)} - \bar{D}_U^{PRE(j)}\right)}, \quad j = k, l, \\ &\equiv \frac{\hat{\beta}_{jU}^{2 \times 2}}{\hat{D}_{jU}^{2 \times 2}}, \quad j = k, l. \end{aligned} \quad (7)$$

Second, we can construct the Wald-DID estimator, leveraging variation in the timing of the initial exposure to the instrument between exposed cohorts. Consider an early exposed cohort k and a middle exposed cohort l . Before period l , the early exposed cohort k is already exposed to the instrument, while the middle exposed cohort l is not yet exposed to the instrument. In this setting, we can view that the middle exposed cohort l plays the role of the control group in both the first stage and the reduced form. From this observation, we can compare the evolution of the treatment and the outcome between the early exposed cohort k and middle exposed cohort l , exploiting the time window $MID(k, l)$ and $PRE(k)$, which we call an Exposed/Not Yet

Exposed design:

$$\begin{aligned}\hat{\beta}_{IV,kl}^{2 \times 2,k} &\equiv \frac{\left(\bar{y}_k^{MID(k,l)} - \bar{y}_k^{PRE(k)}\right) - \left(\bar{y}_l^{MID(k,l)} - \bar{y}_l^{PRE(k)}\right)}{\left(\bar{D}_k^{MID(k,l)} - \bar{D}_k^{PRE(k)}\right) - \left(\bar{D}_l^{MID(k,l)} - \bar{D}_l^{PRE(k)}\right)} \\ &\equiv \frac{\hat{\beta}_{kl}^{2 \times 2,k}}{\hat{D}_{kl}^{2 \times 2,k}}.\end{aligned}\tag{8}$$

Finally, if we focus on the middle exposed cohort l , which changes the exposure status from being unexposed to being exposed at time l , we can regard the early exposed cohort k as the control group after time l because this cohort is already exposed to the instrument at time l . We can compare the evolution of the treatment and the outcome between early exposed cohort k and middle exposed cohort l , exploiting the time window $MID(k, l)$ and $POST(l)$, which we call an Exposed/Exposed Shift design:

$$\begin{aligned}\hat{\beta}_{IV,kl}^{2 \times 2,l} &\equiv \frac{\left(\bar{y}_l^{POST(l)} - \bar{y}_l^{MID(k,l)}\right) - \left(\bar{y}_k^{POST(l)} - \bar{y}_k^{MID(k,l)}\right)}{\left(\bar{D}_l^{POST(l)} - \bar{D}_l^{MID(k,l)}\right) - \left(\bar{D}_k^{POST(l)} - \bar{D}_k^{MID(k,l)}\right)} \\ &\equiv \frac{\hat{\beta}_{kl}^{2 \times 2,l}}{\hat{D}_{kl}^{2 \times 2,l}}.\end{aligned}\tag{9}$$

In each type of the DID-IV design, we have three sources of variation. First, each design exploits the subsample from all NT observations. The Unexposed/Exposed DID-IV design in (7) uses two cohorts and all time periods, indicating that the relative sample share is $n_k + n_u$. The Exposed/Not Yet Exposed DID-IV design in (8) uses two cohorts but exploits only the time periods before period l , so the relative sample share is $(1 - \bar{Z}_l)(n_k + n_l)$. The Exposed/Exposed Shift DID-IV design in (9) uses two cohorts but exploits only the time periods after period k , so the relative sample share is $\bar{Z}_k(n_k + n_l)$.

Second, the variation in each type of the DID-IV design partly comes from the variation of the instrument in its subsample. It is equal to the variance of the double demeaning variable $\tilde{Z}_{i,t}$ in each design:

$$\hat{V}_{jU}^Z \equiv n_{jU}(1 - n_{jU})\bar{Z}_j(1 - \bar{Z}_j), \quad j = k, l,\tag{10}$$

$$\hat{V}_{kl}^{Z,k} \equiv n_{kl}(1 - n_{kl}) \left(\frac{\bar{Z}_k - \bar{Z}_l}{1 - \bar{Z}_l}\right) \left(\frac{1 - \bar{Z}_k}{1 - \bar{Z}_l}\right),\tag{11}$$

$$\hat{V}_{kl}^{Z,l} \equiv n_{kl}(1 - n_{kl}) \left(\frac{\bar{Z}_l}{\bar{Z}_k}\right) \left(\frac{\bar{Z}_k - \bar{Z}_l}{\bar{Z}_k}\right),\tag{12}$$

where the \hat{V}_{jU}^Z , $\hat{V}_{kl}^{Z,k}$ and $\hat{V}_{kl}^{Z,l}$ represent the variance of the double demeaning variable $\tilde{Z}_{i,t}$ in Unexposed/Exposed, Exposed/Not Yet Exposed, and Exposed/Exposed Shift DID-IV designs, respectively. In the staggered DID set up, [Goodman-Bacon \(2021\)](#) also describes the two variations, that is, the relative sample share and the variance of the double demeaning treatment variable in each type of the DID designs.

Unlike the staggered DID set up, however, each DID-IV design has an additional source of the variation; the effect of the instrument on the treatment in the first stage. This comes from the fact that each DID-IV design allows the noncompliance of receiving the treatment when units are exposed to the instrument. The amount of this variation is equal to the 2×2 DID

estimator of the treatment in each DID-IV design:

$$\begin{aligned}\hat{D}_{jU}^{2 \times 2} &\equiv \left(\bar{D}_j^{POST(j)} - \bar{D}_j^{PRE(j)} \right) - \left(\bar{D}_U^{POST(j)} - \bar{D}_U^{PRE(j)} \right) \quad j = k, l, \\ \hat{D}_{kl}^{2 \times 2, k} &\equiv \left(\bar{D}_k^{MID(k, l)} - \bar{D}_k^{PRE(k)} \right) - \left(\bar{D}_l^{MID(k, l)} - \bar{D}_l^{PRE(k)} \right), \\ \hat{D}_{kl}^{2 \times 2, l} &\equiv \left(\bar{D}_l^{POST(l)} - \bar{D}_l^{MID(k, l)} \right) - \left(\bar{D}_k^{POST(l)} - \bar{D}_k^{MID(k, l)} \right).\end{aligned}$$

Note that the denominator of the TWFEIV estimator $\hat{\beta}_{IV}$ in (6), which we denote $\hat{C}^{D, Z}$ hereafter, measures the covariance between the instrument $Z_{i,t}$ and the treatment $D_{i,t}$ in whole samples. By some calculations (see the proof of Theorem 1 below), one can show that $\hat{C}^{D, Z}$ is equal to a weighted average of all possible 2×2 DID estimators of the treatment in each DID-IV design:

$$\hat{C}^{D, Z} = \sum_{k \neq U} \hat{w}_{kU} \hat{D}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} [\hat{w}_{kl}^k \hat{D}_{kl}^{2 \times 2, k} + \hat{w}_{kl}^l \hat{D}_{kl}^{2 \times 2, l}],$$

where the weights are:

$$\begin{aligned}\hat{w}_{kU} &= (n_k + n_u)^2 \hat{V}_{kU}^Z, \\ \hat{w}_{kl}^k &= ((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z, k}, \\ \hat{w}_{kl}^l &= ((n_k + n_l)\bar{Z}_k)^2 \hat{V}_{kl}^{Z, l}.\end{aligned}$$

Hereafter, we refer to \hat{w}_{kU} , \hat{w}_{kl}^k , and \hat{w}_{kl}^l as the first stage weights. This decomposition result for $\hat{C}^{D, Z}$ is almost identical to that of [Goodman-Bacon \(2021\)](#) for the TWFE estimator under staggered DID designs, but the slight difference here is that each weight is not scaled by the variance of the double demeaning variable \tilde{Z}_{it} in whole samples.

We now present the decomposition theorem for the TWFEIV estimator under the staggered assignment of the instrument across units. Theorem 1 below is a generalization of the decomposition result for the TWFE estimator with settings of the staggered assignment of the treatment across units in [Goodman-Bacon \(2021\)](#).

Theorem 1 (Instrumented Difference-in-Differences Decomposition Theorem). Suppose that there exist K cohorts, $e = 1, \dots, k, \dots, K$. The data may also contain a never exposed cohort U . Then, the two-way fixed effects instrumental variable estimator $\hat{\beta}_{IV}$ in (6) is a weighted average of all possible 2×2 Wald-DID estimators.

$$\hat{\beta}_{IV} = \left[\sum_{k \neq U} \hat{w}_{IV, kU} \hat{\beta}_{IV, kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} \hat{w}_{IV, kl}^k \hat{\beta}_{IV, kl}^{2 \times 2, k} + \hat{w}_{IV, kl}^l \hat{\beta}_{IV, kl}^{2 \times 2, l} \right].$$

The 2×2 Wald-DID estimators are:

$$\begin{aligned}\hat{\beta}_{IV, kU}^{2 \times 2} &\equiv \frac{\left(\bar{y}_k^{POST(k)} - \bar{y}_k^{PRE(k)} \right) - \left(\bar{y}_U^{POST(k)} - \bar{y}_U^{PRE(k)} \right)}{\left(\bar{D}_k^{POST(k)} - \bar{D}_k^{PRE(k)} \right) - \left(\bar{D}_U^{POST(k)} - \bar{D}_U^{PRE(k)} \right)}, \\ \hat{\beta}_{IV, kl}^{2 \times 2, k} &\equiv \frac{\left(\bar{y}_k^{MID(k, l)} - \bar{y}_k^{PRE(k)} \right) - \left(\bar{y}_l^{MID(k, l)} - \bar{y}_l^{PRE(k)} \right)}{\left(\bar{D}_k^{MID(k, l)} - \bar{D}_k^{PRE(k)} \right) - \left(\bar{D}_l^{MID(k, l)} - \bar{D}_l^{PRE(k)} \right)}, \\ \hat{\beta}_{IV, kl}^{2 \times 2, l} &\equiv \frac{\left(\bar{y}_l^{POST(l)} - \bar{y}_l^{MID(k, l)} \right) - \left(\bar{y}_k^{POST(l)} - \bar{y}_k^{MID(k, l)} \right)}{\left(\bar{D}_l^{POST(l)} - \bar{D}_l^{MID(k, l)} \right) - \left(\bar{D}_k^{POST(l)} - \bar{D}_k^{MID(k, l)} \right)}.\end{aligned}$$

The weights are:

$$\begin{aligned}\hat{w}_{IV,kU} &= \frac{\hat{w}_{kU} \hat{D}_{kU}^{2 \times 2}}{\hat{C}^{D,Z}} \\ \hat{w}_{IV,kl}^k &= \frac{\hat{w}_{kl}^k \hat{D}_{kl}^{2 \times 2,k}}{\hat{C}^{D,Z}} \\ \hat{w}_{IV,kl}^l &= \frac{\hat{w}_{kl}^l \hat{D}_{kl}^{2 \times 2,l}}{\hat{C}^{D,Z}}.\end{aligned}$$

and sum to one, that is, we have $\sum_{k \neq U} w_{IV,kU} + \sum_{k \neq U} \sum_{l > k} [w_{IV,kl}^k + w_{IV,kl}^l] = 1$.

Proof. See Appendix A. □

Theorem 1 shows that when the assignment of the instrument is staggered across units, the TWFEIV estimator is a weighted average of all possible 2×2 Wald-DID estimators. If there exist K cohorts in the data, we have $K^2 - K$ Wald-DID estimators, which come from either Exposed/Not Yet Exposed designs as in (8) or Exposed/Exposed shift designs as in (9). If the data contains a never exposed cohort U , we have additionally K Wald-DID estimators, which come from Unexposed/Exposed designs as in (7). If both situations occur, the TWFEIV estimator equals a weighted average of K^2 Wald-DID estimators.

The weight assigned to each Wald-DID estimator consists of three parts: the relative sample share squared, the variance of the double demeaning variable $\tilde{Z}_{i,t}$, and the DID estimator of the treatment in each DID-IV design. The first part depends on the sample share of two cohorts and the timing of the initial exposure date. The second part reflects the variation of the instrument in the subsample, represented by (10)-(12), and depends on the relative sample share between two cohorts and the timing of the initial exposure date. Finally, the remaining part reflects variation in the evolution of the treatment between the two cohorts. Note that the weight is not guaranteed to be non-negative in finite sample settings: although the first and second parts are always non-negative, the DID estimator of the treatment can be potentially negative in the data.

Theorem 1 also shows that if we subset the data containing only two cohorts (cohorts k and l), the TWFEIV estimator $\beta_{IV,kl}^{2 \times 2}$ in the subsample can be written as:

$$\beta_{IV,kl}^{2 \times 2} = \frac{\hat{w}_{kl}^k \hat{D}_{kl}^{2 \times 2,k}}{\hat{w}_{kl}^k \hat{D}_{kl}^{2 \times 2,k} + \hat{w}_{kl}^l \hat{D}_{kl}^{2 \times 2,l}} \beta_{IV,kl}^{2 \times 2,k} + \frac{\hat{w}_{kl}^l \hat{D}_{kl}^{2 \times 2,l}}{\hat{w}_{kl}^k \hat{D}_{kl}^{2 \times 2,k} + \hat{w}_{kl}^l \hat{D}_{kl}^{2 \times 2,l}} \beta_{IV,kl}^{2 \times 2,l}.$$

The TWFEIV estimator $\beta_{IV,kl}^{2 \times 2}$ is a weighted average of the Wald-DID estimators which come from either Exposed/Not Yet Exposed design or Exposed/Exposed Shift design, and the weight assigned to each Wald-DID estimator reflects the first stage weight and the DID estimator of the treatment in each DID-IV design.

To make the DID-IV decomposition theorem concrete, we provide a simple numerical example. Suppose we have three cohorts with equal sample size, as shown in Figure 1. In this figure, we set an early exposed period k and a middle exposed period l such that $\tilde{Z}_k = 0.67$ and $\tilde{Z}_l = 0.21$. We assume that the effect of the instrument on the treatment is 0.15 in cohort k and 0.1 in cohort l over time. This means that the units in cohort k are more induced to the treatment by the instrument than those in cohort l and the effects are stable in both cohorts. The DID estimates of the treatment are $\{\hat{D}_{kU}^{2 \times 2}, \hat{D}_{lU}^{2 \times 2}, \hat{D}_{kl}^{2 \times 2,k}, \hat{D}_{kl}^{2 \times 2,l}\} = \{0.15, 0.1, 0.15, 0.1\}$. We also assume that the effect of the instrument on the outcome through treatment is 9 in cohort k and 10 in cohort l over time. The DID estimates of the outcome are $\{\hat{Y}_{kU}^{2 \times 2}, \hat{Y}_{lU}^{2 \times 2}, \hat{Y}_{kl}^{2 \times 2,k}, \hat{Y}_{kl}^{2 \times 2,l}\} =$

$\{9, 10, 9, 10\}$. Dividing the DID estimate of the treatment by the DID estimate of the outcome yields the Wald-DID estimate: $\{\hat{\beta}_{kU}^{2 \times 2}, \hat{\beta}_{lU}^{2 \times 2}, \hat{\beta}_{kl}^{2 \times 2, k}, \hat{\beta}_{kl}^{2 \times 2, l}\} = \{60, 100, 60, 100\}$. The Wald-DID estimate is larger in cohort l than that of cohort k , though as we already noted, the effect of the instrument on the treatment is larger in cohort k than that of cohort l .

The DID estimates of the treatment and the exposure timing determine the amount of the weight assigned to each Wald-DID estimate, holding the sample size equal across cohorts. In the above setting, the resulting weights are $\{\hat{w}_{IV, kU}, \hat{w}_{IV, lU}, \hat{w}_{IV, kl}^k, \hat{w}_{IV, kl}^l\} = \{0.28, 0.12, 0.40, 0.20\}$. In Unexposed/Exposed designs, we have $\hat{w}_{IV, kU} > \hat{w}_{IV, lU}$ for two reasons. First, the DID estimate of the treatment is larger in cohort k than that of cohort l , that is, we have $\hat{D}_{kU}^{2 \times 2} = 0.15 > 0.1 = \hat{D}_{lU}^{2 \times 2}$. Second, the time period k is closer to the middle in the whole period than the time period l , that is, we have $\bar{Z}_k(1 - \bar{Z}_k) = 0.22 > 0.17 = \bar{Z}_l(1 - \bar{Z}_l)$, which implies $\hat{w}_{kU} > \hat{w}_{lU}$ in the first stage weight. By the similar argument, we have $\hat{w}_{IV, kl}^k > \hat{w}_{IV, kl}^l$ between Exposed/Not Yet Exposed and Exposed/Exposed Shift designs: we have $\hat{D}_{kl}^{2 \times 2, k} = 0.15 > 0.1 = \hat{D}_{kl}^{2 \times 2, l}$ and $\hat{w}_{kl}^k > \hat{w}_{kl}^l$ in the first stage weight. If the DID estimates of the treatment are equal between the two designs, the exposure timing matters: we have $\hat{w}_{IV, kU} < \hat{w}_{IV, kl}^k$ and $\hat{w}_{IV, lU} < \hat{w}_{IV, kl}^l$. The DD estimates are the same in each comparison, that is, we have $\hat{D}_{kU}^{2 \times 2} = \hat{D}_{kl}^{2 \times 2, k}$ and $\hat{D}_{lU}^{2 \times 2} = \hat{D}_{kl}^{2 \times 2, l}$. However, the different initial exposure date yields different weights in the first stage, that is, we have $\hat{w}_{kU} < \hat{w}_{kl}^k$ and $\hat{w}_{lU} < \hat{w}_{kl}^l$, which make the difference above the two comparisons.

In this numerical example, the simple average of the Wald-DID estimates is 80 and the weighted average is $100 \times \frac{3}{5} + 60 \times \frac{2}{5} = 84$ where the weight assigned to the Wald-DID estimate reflects the relative amount of the DID estimate of the treatment. The TWFEIV estimate, however, is $\hat{\beta}_{IV} = 60 \times (0.28 + 0.40) + 100 \times (0.12 + 0.20) = 72.8$ because it assigns more weights on the smaller Wald-DID estimate.

Theorem 1 is a decomposition result for the TWFEIV estimator and not for the estimand. Related to the work in this paper, [de Chaisemartin and D'Haultfoeuille \(2020\)](#) decompose the TWFE estimand and present the issue regarding the use of this estimand under DID designs: some weights assigned to the causal parameters in this estimand can be potentially negative. In their appendix, the authors also decompose the TWFEIV estimand, and refer to the negative weight problem in this estimand. Specifically, they apply their decomposition theorem for the TWFE estimand to the numerator and the denominator of the TWFEIV estimand respectively, and conclude that this estimand identifies the local average treatment effect as in [Imbens and Angrist \(1994\)](#) only if the effects of the instrument on the treatment and the outcome are homogeneous across groups and over time. In fact, the population coefficients on the instrument in the first stage and the reduced form regressions take the form of the TWFE estimand and their decomposition theorem for the TWFE estimand is also applicable to the analysis of the TWFEIV estimand. However, the way of their decomposition for the TWFEIV estimand has some drawbacks. First, they do not formally state the target parameter and identifying assumptions in DID-IV designs. Second, their decomposition for the TWFEIV estimand is not based on the target parameter in DID-IV designs. Finally, the sufficient conditions for this estimand to have its causal interpretation are not well explored.

In the following section, we explore the causal interpretation of the TWFEIV estimand under staggered DID-IV designs. In section 3, we first define the target parameter and identifying assumptions in staggered DID-IV designs. In section 4, based on the decomposition theorem for the TWFEIV estimator, we then provide the causal interpretation of the TWFEIV estimand under staggered DID-IV designs. Finally, we investigate the sufficient conditions for this estimand to attain its causal interpretation under staggered DID-IV designs.

3 Staggered instrumented difference-in-differences

In this section, we formalize the staggered instrumented difference-in-differences (DID-IV), built on the recent work in [Miyaji \(2024\)](#). We first introduce the additional notation. We then define the target parameter and identifying assumptions in staggered DID-IV designs.

3.1 Notation

First, we introduce the potential outcomes framework. Let $Y_{i,t}(d, z)$ denote the potential outcome in period t when unit i receives the treatment path $d \in \mathcal{S}(D)$ and the instrument path $z \in \mathcal{S}(Z)$. Similarly, let $D_{i,t}(z)$ denote the potential treatment status in period t when unit i receives the instrument path $z \in \mathcal{S}(Z)$.

Assumption 1 allows us to rewrite $D_{i,t}(z)$ by the initial adoption date $E_i = e$. Let $D_{i,t}^e$ denote the potential treatment status in period t if unit i is first exposed to the instrument in period e . Let $D_{i,t}^\infty$ denote the potential treatment status in period t if unit i is never exposed to the instrument. Hereafter, we call $D_{i,t}^\infty$ the "never exposed treatment". Since the adoption date of the instrument uniquely pins down one's instrument path, we can write the observed treatment status $D_{i,t}$ for unit i at time t as

$$D_{i,t} = D_{i,t}^\infty + \sum_{1 \leq e \leq T} (D_{i,t}^e - D_{i,t}^\infty) \cdot \mathbf{1}\{E_i = e\}.$$

We define $D_{i,t} - D_{i,t}^\infty$ to be the effect of the instrument on the treatment for unit i at time t , which is the difference between the observed treatment status $D_{i,t}$ to the never exposed treatment status $D_{i,t}^\infty$. Hereafter, we refer to $D_{i,t} - D_{i,t}^\infty$ as the individual exposed effect in the first stage. In the DID literature, [Callaway and Sant'Anna \(2021\)](#) and [Sun and Abraham \(2021\)](#) define the effect of the treatment on the outcome in the same fashion.

Next, we introduce the group variable which describes the type of unit i at time t , based on the reaction of potential treatment choices at time t to the instrument path z . Let $G_{i,e,t} \equiv (D_{i,t}^\infty, D_{i,t}^e)(t \geq e)$ be the group variable at time t for unit i and the initial exposure date e . Specifically, the first element $D_{i,t}^\infty$ represents the treatment status at time t if unit i is never exposed to the instrument $E_i = \infty$ and the second element $D_{i,t}^e$ represents the treatment status at time t if unit i starts exposed to the instrument at $E_i = e$. Following to the terminology in [Imbens and Angrist \(1994\)](#), we define $G_{i,e,t} = (0, 0) \equiv NT_{e,t}$ to be the never-takers, $G_{i,e,t} = (1, 1) \equiv AT_{e,t}$ to be the always-takers, $G_{i,e,t} = (0, 1) \equiv CM_{e,t}$ to be the compliers and $G_{i,e,t} = (1, 0) \equiv DF_{e,t}$ to be the defiers at time t and the initial exposure date e .

Finally, we make a no carryover assumption on potential outcomes $Y_{i,t}(d, z)$.

Assumption 2 (No carryover assumption).

$$\forall z \in \mathcal{S}(Z), \forall d \in \mathcal{S}(D), \forall t \in \{1, \dots, T\}, Y_{i,t}(d, z) = Y_{i,t}(d_t, z),$$

where $d = (d_1, \dots, d_T)$ is the generic element of the treatment path D_i .

This assumption requires that potential outcomes $Y_{i,t}(d, z)$ depend only on the current treatment status d_t and the instrument path z . In the DID literature, several recent papers impose this assumption with settings of a non-staggered treatment; see, e.g., [de Chaisemartin and D'Haultfoeulle \(2020\)](#) and [Imai and Kim \(2021\)](#). Although it can be possible to weaken this assumption by introducing the treatment path d in potential outcomes $Y_{i,t}(d, z)$, this requires the cumbersome notation and complicates the definition of our target parameter, thus is beyond the scope of this paper.

Henceforth, we keep Assumption 1 and 2. In the next section, we define the target parameter in staggered DID-IV designs.

3.2 Target parameter in staggered DID-IV designs

Our target parameter in staggered DID-IV designs is the cohort specific local average treatment effect on the treated (CLATT) defined below.

Def. The cohort specific local average treatment effect on the treated (CLATT) at a given relative period l from the initial adoption of the instrument is

$$\begin{aligned} CLATT_{e,l} &= E[Y_{i,e+l}(1) - Y_{i,e+l}(0) | E_i = e, D_{i,e+l}^e > D_{i,e+l}^\infty] \\ &= E[Y_{i,e+l}(1) - Y_{i,e+l}(0) | E_i = e, CM_{e,e+l}]. \end{aligned}$$

This parameter measures the treatment effects at a given relative period l from the initial instrument adoption date $E_i = e$, for those who belong to cohort e , and are the compliers $CM_{e,e+l}$, that is, who are induced to treatment by instrument at time $e + l$. Each $CLATT_{e,l}$ can potentially vary across cohorts and over time, as it depends on cohort e , relative period l , and the compliers $CM_{e,e+l}$.

3.3 Identifying assumptions in staggered DID-IV designs

In this section, we state the identifying assumptions in staggered DID-IV designs based on Miyaji (2024).

Assumption 3 (Exclusion Restriction in multiple time periods).

$$\forall z \in \mathcal{S}(Z), \forall d_t \in \mathcal{S}(D_t), \forall t \in \{1, \dots, T\}, Y_{i,t}(d, z) = Y_{i,t}(d) \quad a.s.$$

Assumption 3 requires that the path of the instrument does not directly affect the potential outcome for all time periods and its effects are only through treatment. Given Assumption 2 and Assumption 3, we can write the potential outcome $Y_{i,t}(d, z)$ as $Y_{i,t}(d_t) = D_{i,t}Y_{i,t}(1) + (1 - D_{i,t})Y_{i,t}(0)$.

Here, we introduce the potential outcomes at time t if unit i is assigned to the instrument path $z \in \mathcal{S}(Z)$:

$$Y_{i,t}(D_{i,t}(z)) \equiv D_{i,t}(z)Y_{i,t}(1) + (1 - D_{i,t}(z))Y_{i,t}(0).$$

Since the exposure timing E_i completely determines the path of the instrument, we can write the potential outcomes for cohort e and cohort ∞ as $Y_{i,t}(D_{i,t}^e)$ and $Y_{i,t}(D_{i,t}^\infty)$, respectively. The potential outcome $Y_{i,t}(D_{i,t}^e)$ represents the outcome status at time t if unit i is first exposed to the instrument at time e and the potential outcome $Y_{i,t}(D_{i,t}^\infty)$ represents the outcome status at time t if unit i is never exposed to the instrument. Hereafter, we refer to $Y_{i,t}(D_{i,t}^\infty)$ as the "never exposed outcome".

Assumption 4 (Monotonicity Assumption in multiple time periods).

$$Pr(D_{i,e+l}^e \geq D_{i,e+l}^\infty) = 1 \quad \text{or} \quad Pr(D_{i,e+l}^e \leq D_{i,e+l}^\infty) = 1 \quad \text{for all } e \in \mathcal{S}(E_i) \text{ and for all } l \geq 0.$$

This assumption requires that the instrument path affects the treatment adoption behavior in a monotone way for all relative periods after the initial exposure. Recall that we define $D_{i,t} - D_{i,t}^\infty$ to be the effect of the instrument on the treatment for unit i at time t . Assumption 4 requires that the individual exposed effect in the first stage is non-negative (or non-positive) for all i and all the time periods after the initial exposure. This assumption implies that the group variable $G_{i,e,t} \equiv (D_{i,t}^\infty, D_{i,t}^e)$ can take three values with non-zero probability for all e and all $t \geq e$. Hereafter, we consider the type of the monotonicity assumption that rules out the existence of the defiers $DF_{e,t}$ for all $t \geq e$ in any cohort e .

Assumption 5 (No anticipation in the first stage).

$$D_{i,e+l}^e = D_{i,e+l}^\infty \text{ a.s. for all units } i, \text{ for all } e \in \mathcal{S}(E_i) \text{ and for all } l < 0.$$

Assumption 5 requires that the potential treatment choice for the treatment in any l period before the initial exposure to the instrument is equal to the never exposed treatment. This assumption restricts the anticipatory behavior before the initial exposure in the first stage.

Assumption 6 (Parallel Trends Assumption in the treatment in multiple time periods).

$$\text{For all } s \neq t, E[D_{i,t}^\infty - D_{i,s}^\infty | E_i = e] \text{ is same for all } e \in \mathcal{S}(E_i).$$

Assumption 6 is a parallel trends assumption in the treatment in multiple periods and multiple cohorts. This assumption requires that the trends of the treatment across cohorts would have followed the same path, on average, if there is no exposure to the instrument. Assumption 6 is analogous to that of Callaway and Sant’Anna (2021) and Sun and Abraham (2021) in DID designs: both papers impose the same type of the parallel trends assumption on untreated outcomes with settings of multiple periods and multiple cohorts.

Assumption 7 (Parallel Trends Assumption in the outcome in multiple time periods).

$$\text{For all } s < t, E[Y_{i,t}(D_{i,t}^\infty) - Y_{i,s}(D_{i,s}^\infty) | E_i = e] \text{ is same for all } e \in \mathcal{S}(E_i).$$

Assumption 7 is a parallel trends assumption in the outcome with settings of multiple periods and multiple cohorts. This assumption requires that the expectation of the never exposed outcome across cohorts would have followed the same evolution if the assignment of the instrument had not occurred. From the discussions in Miyaji (2024), we can interpret that this assumption requires the same expected time gain across cohorts and over time: the effects of time on outcome through treatment are the same on average across cohorts and over time.

4 Causal interpretation of the TWFEIV estimand

In this section, we explore the causal interpretation of the TWFEIV estimand under staggered DID-IV designs. In section 4.1, we first define the main building block parameter in the first stage and reduced form regressions, respectively. In section 4.2, we then interpret the TWFEIV estimand under staggered DID-IV designs, and show that this estimand potentially fails to summarize the treatment effects. In section 4.3, given the negative result of using the TWFEIV estimand under staggered DID-IV designs, we describe the various restrictions on main building block parameter in each stage regression. In section 4.4, as a preparation, we then describe the causal interpretation of the denominator in the TWFEIV estimand under these restrictions. In section 4.5, we finally investigate the sufficient conditions for the TWFEIV estimand to attain its causal interpretation.

4.1 Main building block parameter in each stage regression

As we already mentioned in section 2, the TWFEIV regression employs the TWFEIV regression twice in the first stage and reduced form regressions. In this section, we define the main building block parameter in each stage regression.

In the first stage regression, our building block parameter is the average of individual exposed effect at a given relative period l from the initial exposure to the instrument in cohort e . We call this the cohort specific average exposed effect on the treated in the first stage ($\text{CAET}_{e,l}^1$) defined below.

Def. The cohort specific average exposed effect on the treated in the first stage ($CAET^1$) at a given relative period l from the initial adoption of the instrument is

$$CAET_{e,l}^1 = E[D_{i,e+l} - D_{i,e+l}^\infty | E_i = e].$$

We use the superscript 1 to make it clear that we define this parameter for the first stage regression. In the recent DID literature, [Sun and Abraham \(2021\)](#) define their main building block parameter in staggered DID designs in a similar fashion and call it the cohort specific average treatment effect on the treated. [Callaway and Sant'Anna \(2021\)](#) call the same parameter the group-time average treatment effect.

If the treatment is binary and monotonicity assumption (Assumption 4) holds, the $CAET_{e,l}^1$ is equal to the share of the compliers $CM_{e,e+l}$ in cohort e at period $e + l$:

$$\begin{aligned} CAET_{e,l}^1 &= E[D_{i,e+l}^e - D_{i,e+l}^\infty | E_i = e] \\ &= Pr(CM_{e,e+l} | E_i = e). \end{aligned}$$

In the reduced form regression, our building block parameter is the average of individual effect of the instrument on the outcome through treatment at a given relative period l from the initial exposure to the instrument in cohort e . We call this the cohort specific average intention to exposed effect on the treated in the reduced form ($CAIET_{e,l}$) defined below.

Def. The cohort specific average intention to exposed effect on the treated in the reduced form ($CAIET$) at a given relative period l from the initial adoption of the instrument is

$$CAIET_{e,l} = E[Y_{i,e+l}(D_{i,e+l}) - Y_{i,e+l}(D_{i,e+l}^\infty) | E_i = e].$$

If we assume the identifying assumptions in staggered DID-IV designs (Assumptions 1 to 7), this parameter is equal to a product of the $CLATT_{e,l}$ and $CAET_{e,l}^1$:

$$\begin{aligned} CAIET_{e,l} &= E[Y_{i,e+l}(D_{i,e+l}^e) - Y_{i,e+l}(D_{i,e+l}^\infty) | E_i = e] \\ &= E[(D_{i,e+l}^e - D_{i,e+l}^\infty)(Y_{i,e+l}(1) - Y_{i,e+l}(0)) | E_i = e] \\ &= E[Y_{i,e+l}(1) - Y_{i,e+l}(0) | E_i = e, CM_{e,e+l}] \cdot Pr(CM_{e,e+l} | E_i = e) \\ &= CLATT_{e,l} \cdot CAET_{e,l}^1. \end{aligned} \tag{13}$$

In other words, if we scale the $CAIET_{e,l}$ in the reduced form by the $CAET_{e,l}^1$ in the first stage, we obtain the $CLATT_{e,l}$, which is the reason why we call this the cohort specific average "intention to exposed effect" on the treated in the reduced form.

4.2 Interpreting the TWFEIV estimand under staggered DID-IV designs

We now interpret the TWFEIV estimand under staggered DID-IV designs based on the DID-IV decomposition theorem derived in section 2 and the main building block parameters defined in the previous section. This section presumes the monotonicity assumption (Assumption 4) to clarify the interpretation of each notation defined below.

First, we introduce the additional notation. Let $CLATT_k^{CM}(W)$ denote a weighted average of each $CLATT_{k,t}$ in the time window W (with T_W periods) where the weight reflects the relative amount of the exposed effect in the first stage in cohort k at period t :

$$\begin{aligned} CLATT_k^{CM}(W) &\equiv \sum_{t \in W} \frac{CAET_{k,t}^1}{\sum_{t \in W} CAET_{k,t}^1} CLATT_{k,t} \\ &= \sum_{t \in W} \frac{Pr(CM_{k,t} | E_i = k)}{\sum_{t \in W} Pr(CM_{k,t} | E_i = k)} CLATT_{k,t}. \end{aligned}$$

The first equality holds because we have a binary treatment and assume the monotonicity assumption (Assumption 4). Each weight assigned to each $CLATT_{k,t}$ reflects the relative share of the compliers at period t in cohort k during the time window W . We call this the compliers weighted scheme. This would be one of the reasonable weighting schemes for two reasons. First, the weight is designed to be larger in the period when the proportion of the compliers is higher in cohort k . Second, the sum of the weight is one by construction: the proportion of the compliers in each period in cohort k is divided by the total amount of the compliers in the time window W in cohort k .

We also define the similar notation $CLATT_k(W)$, in which the proportion of the compliers in cohort k at period t is divided by the time length T_W :

$$\begin{aligned} CLATT_k(W) &\equiv \frac{1}{T_W} \sum_{t \in W} CAET_{k,t}^1 CLATT_{k,t} \\ &= \frac{1}{T_W} \sum_{t \in W} Pr(CM_{k,t} | E_i = k) CLATT_{k,t}. \end{aligned}$$

We call this the time-corrected weighting scheme. In contrast to $CLATT_k^{CM}(W)$, the weight assigned to each $CLATT_{k,t}$ can be inappropriate: each weight does not reflect the relative share of the compliers in cohort k at period t . In addition, the sum of each weight is not equal to one in general.

Theorem 2 below shows the probability limit of the TWFEIV estimator $\hat{\beta}_{IV}$ under staggered DID-IV designs (Assumptions 1-7).

Theorem 2. Suppose Assumptions 1-7 hold. Then, the TWFEIV estimand β_{IV} consists of two terms:

$$\begin{aligned} \hat{\beta}_{IV} &= \left[\sum_{k \neq U} \hat{w}_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} \hat{w}_{IV,kl}^k \hat{\beta}_{IV,kl}^{2 \times 2,k} + \hat{w}_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2,l} \right] \\ &\xrightarrow{p} WCLATT - \Delta CLATT. \end{aligned}$$

where we define:

$$\begin{aligned} WCLATT &\equiv \sum_{k \neq U} w_{IV,kU} CLATT_k^{CM}(POST(k)) + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^k CLATT_k^{CM}(MID(k,l)) \\ &\quad + \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot CLATT_l(POST(l)) \\ \Delta CLATT &\equiv \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot [CLATT_k(POST(l)) - CLATT_k(MID(k,l))]. \end{aligned}$$

The weights $w_{IV,kU}$ and $w_{IV,kl}^k$ are the probability limit of $\hat{w}_{IV,kU}$ and $\hat{w}_{IV,kl}^k$, respectively. The weight $\sigma_{IV,kl}^l$ is the probability limit of $\frac{\hat{w}_{kl}^l}{\hat{C}_{D,Z}^l} \neq \frac{\hat{w}_{kl}^l \hat{D}_{kl}^{2 \times 2,l}}{\hat{C}_{D,Z}^l} = \hat{w}_{IV,kl}^l$. The specific expressions for each weight are shown in equations (49), (50), and (51) in Appendix B.

Proof. See Appendix B. □

Theorem 2 shows that the TWFEIV estimand β_{IV} consists of two terms ($WCLATT$ and $\Delta CLATT$) and potentially fails to aggregate the treatment effects under staggered DID-IV designs.

The first term $WCLATT$ is a positively weighted average of each $CLATT_{k,t}$ for the post-exposed period in cohort k . We call this a weighted average cohort specific local average treatment effect on the treated ($WCLATT$) parameter. The first and the second terms in the $WCLATT$ use the compliers weighted scheme, but the third term in $WCLATT$ uses the time-corrected one.

Although $WCLATT$ can be a causal parameter, the amount of this parameter may be difficult to interpret in practice for two reasons. First, the weight $\sigma_{IV,kl}^l$ assigned to $CLATT_l(POST(l))$ reflects only the sample share and the variation of the instrument, and does not reflect the variation of the treatment $D_{kl}^{2 \times 2, l}$ in the first stage. Because the other weights, $w_{IV,kU}$ and $w_{IV,kl}^k$ precisely reflect all the variations in each DID-IV design, this asymmetry can break the implication of the magnitude of this parameter in a given application. Second, the $CLATT_l(POST(l))$ in the third term is a weighted average of $CLATT_{k,t}$ for the post exposed periods in cohort k , but the weight assigned to each $CLATT_{k,t}$ seems not reasonable: it does not reflect the relative share of the compliers in period t in cohort k and the sum of the weight is not equal to one.

The problem of the $WCLATT$ is due to the "bad comparisons" in the first stage TWFE regression: when we compare the evolution of the treatment in Exposed/Exposed Shift designs, we use already exposed cohorts as controls. In these comparisons, we should offset the DID estimator of the treatment in each weight in Exposed/Exposed Shift designs by the one appeared in the denominator of the corresponding Wald-DID estimator, which produces the weight $\sigma_{IV,kl}^l$ and $CLATT_l(POST(l))$ in the third term.

The second term $\Delta CLATT$ is a weighted sum of the differences in the positively weighted average of each $CLATT_{k,t}$ from the exposed period k to before period l ($k < l$) and after period l in the already exposed cohort k . This term fails to properly aggregate the treatment effects because the $CLATT_k(POST(l))$ is canceled out by the $CLATT_k(MID(k,l))$ in each cohort k . This problem arises due to the "bad comparisons" in the reduced form TWFE regression: when we compare the evolution of the outcome in Exposed/Exposed Shift designs, we use already exposed cohorts as controls. In these comparisons, we subtract their expected trends of unexposed potential outcomes and average intention exposed effects, which yields the $\Delta CLATT$.

Overall, this section shows that the TWFEIV estimand potentially fails to summarize the treatment effects under staggered DID-IV designs. In the next section, we first describe various restrictions on main building block parameters in the first stage and the reduced form regressions. Given these restrictions on exposed effect heterogeneity, we then explore the sufficient conditions for the TWFEIV estimand to be causally interpretable parameter.

4.3 Restrictions on exposed effect heterogeneity

First, we describe the restrictions on the $CAET_{e,l}^1$ in the first stage regression.

Assumption 8 (Exposed effect homogeneity across cohorts in the first stage). For each relative period l , $CAET_{e,l}^1$ does not depend on cohort e and is equal to AET_l^1 .

Assumption 8 requires that the exposed effects in the first stage depend on only the relative time period l after the initial exposure to the instrument and do not depend on the cohort e . This assumption does not exclude the dynamic effects of the instrument on the treatment, but requires that the exposed effects are the same across cohorts for all relative periods.

Assumption 9 (Stable exposed effect over time within cohort in the first stage). For each cohort e , $CAET_{e,l}^1$ does not depend on the relative time period l and is equal to $CAET_e^1$.

Assumption 9 rules out the dynamic effects of the instrument on the treatment within cohort e in the first stage regression. Assumption 9 permits the heterogeneous exposed effects across cohort e , but requires the homogeneous exposed effects over time after the initial adoption of the instrument within cohort e .

The recent DID literature imposes the similar restrictions as in Assumption 8 and Assumption 9 on treatment effects. Sun and Abraham (2021) assume that "each cohort experiences the same path of treatment effects", which is in line with Assumption 8. Goodman-Bacon (2021) requires heterogeneous treatment effects to either be "constant over time but vary across units" or "vary over time but not across units". The former corresponds to Assumption 9 and the latter corresponds to Assumption 8.

Next, we describe the restrictions on the $CAIET_{e,l}$ in the reduced form regression. Following to Assumption 8 and Assumption 9 on the $CAET_{e,l}^1$, we consider Assumption 10 and Assumption 11 below.

Assumption 10 (Exposed effect homogeneity across cohorts in the reduced form). For each relative period l , $CAIET_{e,l}$ does not depend on cohort e and is equal to $AIET_l$.

Assumption 11 (Stable exposed effect over time within cohort in the reduced form). For each cohort e , $CAIET_{e,l}$ does not depend on the relative time period l and is equal to $CAIET_e$.

Assumption 10 requires that the evolution of the average intention to exposed effect after the initial exposure is the same across cohorts. Assumption 11 requires that the average intention to exposed effects are stable over time in all relative periods within cohort e .

Note that given Assumption 8 and Assumption 10, we have the following restriction on the $CLATT_{e,l}$, which follows from equation (13) in section 4.1.

Assumption 12 (Treatment effect homogeneity across cohorts for $CLATT_{e,l}$). For each relative period l , $CLATT_{e,l}$ does not depend on cohort e and is equal to $LATT_l$.

Similarly, given Assumption 9 and Assumption 11, we have the following restriction on the $CLATT_{e,l}$.

Assumption 13 (Stable treatment effect over time within cohort for $CLATT_{e,l}$). For each cohort e , $CLATT_{e,l}$ does not depend on the relative time period l and is equal to $CLATT_e$.

4.4 The denominator in the TWFEIV estimand

In this section, we first interpret the denominator in the TWFEIV estimand under various restrictions considered in section 4.3. This section is a preparation for the next section, in which we analyze the TWFEIV estimand itself.

As we already noted, the denominator in the TWFEIV estimator (see equation (6)), $\hat{C}^{D,Z}$ can be decomposed into a weighted average of all possible 2×2 DID estimators of the treatment. In the following discussion, we show that this estimand can potentially fail to aggregate the effects of the instrument on the treatment in the first stage regression without additional restrictions. We then briefly describe the interpretation of this estimand by imposing Assumption 8 or Assumption 9, and state the implications.

First, we introduce the additional notation. Let $CAET_k^1(W)$ denote an equally weighted average of the $CAET_{k,t}^1$ in the time window W (with T_W period length):

$$CAET_k^1(W) \equiv \frac{1}{T_W} \sum_{t \in W} CAET_{k,t}^1.$$

If we assume Assumption 4 (monotonicity assumption), $CAET_k^1(W)$ is an equally weighted average of the fraction of the compliers in cohort k in the time window W . For instance, the $CAET_k^1(POST(k))$ is an equally weighted average of the $CAET_k^1$ during the periods after the initial exposure date k and rewritten as

$$\begin{aligned} CAET_k^1(POST(k)) &= \frac{1}{T - (k - 1)} \sum_{t=k}^T CAET_{k,t}^1 \\ &= \frac{1}{T - (k - 1)} \sum_{t=k}^T Pr(CM_{k,t} | E_i = k). \end{aligned}$$

Lemma 1 below shows the probability limit of the denominator $\hat{C}^{D,Z}$ under staggered DID-IV designs. This lemma is mainly based on the result of Goodman-Bacon (2021), who shows the probability limit of the two-way fixed effects estimator under staggered DID designs. The slight difference here is that each weight assigned to each $CAET_k^1(W)$ in $C^{D,Z}$ is not divided by the probability limit of the grand mean $\frac{1}{NT} \sum_i \sum_t \tilde{Z}_{it}$.

Lemma 1. Suppose Assumptions 1-7 hold. Then, the probability limit of the denominator of the TWFEIV estimator, $C^{D,Z}$ consists of two terms:

$$\begin{aligned} \hat{C}^{D,Z} &= \sum_{k \neq U} \hat{w}_{kU} \hat{D}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} [\hat{w}_{kl}^k \hat{D}_{kl}^{2 \times 2, k} + \hat{w}_{kl}^l \hat{D}_{kl}^{2 \times 2, l}] \\ &\xrightarrow{p} WCAET - \Delta CAET^1. \end{aligned}$$

where we define:

$$\begin{aligned} WCAET &\equiv \sum_{k \neq U} w_{kU} CAET_k^1(POST(k)) + \sum_{k \neq U} \sum_{l > k} w_{kl}^k CAET_k^1(MID(k, l)) + w_{kl}^l CAET_l^1(POST(l)), \\ \Delta CAET^1 &\equiv \sum_{k \neq U} \sum_{l > k} w_{kl}^l [CAET_k^1(POST(l)) - CAET_k^1(MID(k, l))]. \end{aligned}$$

The weights w_{kU}, w_{kl}^k and w_{kl}^l are the probability limit of $\hat{w}_{kU}, \hat{w}_{kl}^k$ and \hat{w}_{kl}^l defined in section 3 respectively, and are non-negative. The specific expressions in each weight are shown in equations (35)-(37) in Appendix B.

Proof. See Appendix B. □

Lemma 1 shows that we can decompose $C^{D,Z}$ into two terms. The first term is a positively weighted average of each $CAET_{k,t}^1$ during the periods after the initial exposure in exposed cohorts, allowing for its causal interpretation. Following the terminology in Goodman-Bacon (2021), we call this a weighted average cohort specific exposed effect on the treated (WCAET) parameter.

The second term $\Delta CAET^1$ is equal to the sum of the difference in the positively weighted average of exposed effect $CAET_{k,t}^1$ from the exposed period k to before period l ($k < l$) and after period l in the already exposed cohort k . This term fails to properly aggregate the causal parameter in the first stage because some exposed effects are canceled out by other exposed effects.

Lemma 1 implies that if we assume only Assumptions 1-7, the probability limit of the denominator in the TWFEIV estimand, $C^{D,Z}$ generally fails to properly summarize the exposed effects in the first stage due to the second term $\Delta CAET^1$. This problem arises from the "bad

comparisons" performed by the TWFE regression in the first stage: we treat the already exposed cohorts as control groups in the Exposed/ Exposed Shift designs. In these comparisons, we should subtract their expected trends of unexposed potential treatment choices and their expected exposed effects, which yields the second term $\Delta CAET^1$. In the DID literature, [Borusyak et al. \(2021\)](#), [de Chaisemartin and D'Haultfœuille \(2020\)](#), and [Goodman-Bacon \(2021\)](#) point out the same issue for the TWFE estimand in staggered DID designs.

Based on the negative result shown in Lemma 1, we consider the restrictions on exposed effect heterogeneity in the first stage regression. The conclusion here is that $C^{D,Z}$ properly aggregates each $CAET_{k,t}^1$ only if Assumption 9 holds, that is, the exposed effects are stable over time within cohort e . Because [Goodman-Bacon \(2021\)](#) have already made the same point for the TWFE estimand, we briefly summarize the interpretation of $C^{D,Z}$ under Assumption 8 or Assumption 9 in the following. For the more detailed discussions, see section 3.1 in [Goodman-Bacon \(2021\)](#).

Interpreting $C^{D,Z}$ under Assumption 8 only

Even when Assumption 8 holds, that is, the exposed effects are the same across cohorts but vary over time in the first stage, we have $\Delta CAET^1 \neq 0$ in general. This implies that if we impose only Assumption 8, we cannot generally interpret the $C^{D,Z}$ as measuring the positively weighted average of exposed effects in the first stage.

Interpreting $C^{D,Z}$ under Assumption 9 only

If Assumption 9 holds, that is, the exposed effects are stable over time within cohort e in the first stage, we have $CAET_k^1(W) = CAET_k^1$. This implies that the second term $\Delta CAET^1$ is equal to zero:

$$\begin{aligned} \Delta CAET^1 &= \sum_{k \neq U} \sum_{l > k} w_{kl}^l [CAET_k^1 - CAET_k^1] \\ &= 0. \end{aligned}$$

Thus, $C^{D,Z}$ simplifies to:

$$\begin{aligned} C^{D,Z} &= WCAET \\ &= \sum_{k \neq U} CAET_k^1 \underbrace{\left[w_{kU} + \sum_{j=1}^{k-1} w_{jk}^k + \sum_{j=k+1}^K w_{kj}^k \right]}_{\equiv w_k}. \end{aligned}$$

$C^{D,Z}$ weights each $CAET_k^1$ positively across cohorts under Assumption 9 only. We note, however, that each weight assigned to each $CAET_k^1$, w_k is not equal to the sample share in cohort k , but is a function of the sample share and the timing of the initial exposure date.

In this section, we have considered whether the denominator in the TWFEIV estimand properly aggregates the exposed effects in the first stage. We have two implications. First, if we do not impose Assumption 9, the weight assigned to each 2×2 Wald-DID in the TWFEIV estimand may not be properly normalized because the numerator in each weight is divided by $C^{D,Z}$, and the denominator potentially fail to aggregate the exposed effects in the first stage. Second, if we do not impose Assumption 9, some weights assigned to 2×2 Wald-DID estimands

can be potentially negative. This is because the DID estimand of the treatment forms the part of each weight and can be negative due to the "bad comparisons" in the first stage regression.

From the discussion so far, hereafter, we impose Assumption 9 when we consider the restrictions on exposed effect heterogeneity in the first stage.

4.5 Interpreting the TWFEIV estimand under additional restrictions

We now describe the interpretation of the TWFEIV estimand under additional restrictions.

Interpretation under Assumption 9 only

First, we consider imposing Assumption 9 only, that is, we assume only the stable exposed effect over time in the first stage. If Assumption 9 holds, the $CLATT_k^{CM}(W)$ simplifies to an equally weighted average of $CLATT_{k,t}$:

$$\begin{aligned} CLATT_k^{CM}(W) &= \sum_{t \in W} \frac{Pr(CM_{k,t} | E_i = k)}{\sum_{t \in W} Pr(CM_{k,t} | E_i = k)} CLATT_{k,t} \\ &= \frac{1}{T_W} \sum_{t \in W} CLATT_{k,t} \\ &\equiv CLATT_k^{eq}(W). \end{aligned}$$

The $CLATT_k^{eq}(W)$ weights each $CLATT_{k,t}$ equally in the time window W and the weight sum to one by construction. We call this an equal weighting scheme.

Lemma 2 presents the interpretation of the TWFEIV estimand under staggered DID-IV designs and Assumption 9.

Lemma 2. Suppose Assumptions 1-7 hold. If Assumption 9 holds additionally, the TWFEIV estimand β_{IV} consists of two terms:

$$\beta_{IV} = WCLATT - \Delta CLATT.$$

where we define:

$$\begin{aligned} WCLATT &\equiv \sum_{k \neq U} w_{IV,kU} CLATT_k^{eq}(POST(k)) + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^k CLATT_k^{eq}(MID(k, l)) \\ &\quad + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^l CLATT_l^{eq}(POST(l)), \\ \Delta CLATT &\equiv \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot [CLATT_k(POST(l)) - CLATT_k(MID(k, l))]. \end{aligned}$$

The weights $w_{IV,kU}$, $w_{IV,kl}^k$ and $w_{IV,kl}^l$ are the probability limit of $\hat{w}_{IV,kU}$, $\hat{w}_{IV,kl}^k$ and $\hat{w}_{IV,kl}^l$ respectively, and are non-negative. The specific expressions for these weights are shown in equations (54), (55), and (59) in Appendix B. The weight $\sigma_{IV,kl}^l$ is already defined in Theorem 2.

Proof. See Appendix B. □

Lemma 2 shows that Assumption 9 is not sufficient for the TWFEIV estimand to attain its causal interpretation. If the exposed effects in the first stage are stable over time, we can

interpret the first term $WCLATT$ causally and its interpretation seems clear: this parameter is a positively weighted average of each $CLATT_k^{eq}(W)$ and each weight assigned to each $CLATT_k^{eq}(W)$ reflects all the variations in each DID-IV design. However, the second term $\Delta CLATT$ still remains, which contaminates the causal interpretation of the TWFEIV estimand.

Interpretation under Assumption 9 and Assumption 10

Next, we assume Assumption 9 and Assumption 10 additionally. Even in this case, we still have the second term $\Delta CLATT \neq 0$ in general. This implies that the TWFEIV estimand identifies $WCLATT - \Delta CLATT$, that is, this estimand does not generally attain its causal interpretation.

Interpretation under Assumption 9 and Assumption 11

As we already noted in section 4.2, if we assume Assumption 9 and Assumption 11 additionally, we have Assumption 13, that is, $CLATT_{e,t} = CLATT_e$ holds. Then, we obtain the following Lemma.

Lemma 3. Suppose Assumptions 1-7 hold. In addition, if Assumption 9 and Assumption 11 hold, the TWFEIV estimand β_{IV} is:

$$\beta_{IV} = \sum_{k \neq U} CLATT_k \underbrace{\left[w_{IV,kU} + \sum_{j=1}^{k-1} w_{IV,jk}^k + \sum_{j=k+1}^K w_{IV,kj}^k \right]}_{\equiv w_{k,IV}}.$$

where the weights $w_{IV,kU}$, $w_{IV,kj}^k$ and $w_{IV,jk}^k$ are the probability limit of $\hat{w}_{IV,kU}$, $\hat{w}_{IV,kj}^k$ and $\hat{w}_{IV,jk}^k$ respectively.

Proof. See Appendix B. □

If Assumption 9 and Assumption 11 are satisfied, the TWFEIV estimand is a positively weighted average of each $CLATT_k$ across exposed cohorts, which implies that we can interpret this estimand causally. However, at the same time, we also note that the weight $w_{k,IV}$ assigned to each $CLATT_k$ does not reflect only the cohort share and the fraction of the compliers, but is a function of the cohort share, the fraction of the compliers, and the timing of the initial exposure to the instrument.

5 Extensions

This section briefly describes the extensions in section 4. We consider a non-binary, ordered treatment and unbalanced panel settings. It also includes the case when the adoption date of the instrument is randomized across units. For the proofs and the specific discussions, see Appendix C.

Non-binary, ordered treatment

Up to now, we have considered only the case of a binary treatment. When treatment takes a finite number of ordered values, $D_{i,t} \in \{0, 1, \dots, J\}$, our target parameter in staggered DID-IV design is the cohort specific average causal response on the treated (CACRT) defined below.

Def. The cohort specific average causal response on the treated (CACRT) at a given relative period l from the initial adoption of the instrument is

$$CACRT_{e,l} \equiv \sum_{j=1}^J w_{e+l,j}^e \cdot E[Y_{i,e+l}(j) - Y_{i,e+l}(j-1) | E_i = e, D_{i,e+l}^e \geq j > D_{i,e+l}^\infty]$$

where the weights $w_{e+l,j}^e$ are:

$$w_{e+l,j}^e = \frac{\Pr(D_{i,e+l}^e \geq j > D_{i,e+l}^\infty | E_i = e)}{\sum_{j=1}^J \Pr(D_{i,e+l}^e \geq j > D_{i,e+l}^\infty | E_i = e)}.$$

The CACRT is a weighted average of the effect of a unit increase in treatment on outcome, for those who are in cohort e and induced to increase treatment by instrument at a relative period l after the initial exposure. This parameter is similar to the average causal response (ACR) considered in Angrist and Imbens (1995), but the difference here is that there exist dynamic effects in the first stage, and each weight $w_{e+l,j}^e$ and the associated causal parameters in CACRT are conditioned on $E_i = 1$.

If we have a non-binary, ordered treatment, one can show that we have Theorem 2 and Lemmas 2-3 in section 4, which replace $CLATT_{e,k}$ with $CACRT_{e,k}$. Note that our decomposition result for the TWFEIV estimator is unchanged under non-binary, ordered treatment settings.

Unbalanced panel case

Throughout sections 2 to 4, we have considered a balanced panel setting. If we assume an unbalanced panel (or repeated cross section) setting, we obtain the following theorem.

Theorem 3. Suppose Assumptions 1-7 hold. If we assume a binary treatment and an unbalanced panel setting, the population regression coefficient β_{IV} is a weighted average of each $CLATT_{e,t}$ in all relative periods after the initial exposure across cohorts with potentially some negative weights:

$$\beta_{IV} = \sum_e \sum_{t \geq e} w_{e,t} \cdot CLATT_{e,t}.$$

where the weight $w_{e,t}$ is:

$$w_{e,t} = \frac{E[\hat{Z}_{i,t} | E_i = e] \cdot n_{e,t} \cdot CAET_{e,t}^1}{\sum_e \sum_{t \geq e} E[\hat{Z}_{i,t} | E_i = e] \cdot n_{e,t} \cdot CAET_{e,t}^1},$$

where $E[\hat{Z}_{i,t} | E_i = e]$ is the population residuals from regression $Z_{i,t}$ on unit and time fixed effects in cohort e and $n_{e,t}$ is the population share for cohort e at time t . The weights sum to one.

Proof. See Appendix C. □

Theorem 3 shows that the population regression coefficient β_{IV} is a weighted average of all possible $CLATT_{e,t}$ across cohorts, but some weights can be negative. Theorem 3 is related to de Chaisemartin and D'Haultfœuille (2020), who show the decomposition theorem for the TWFEIV estimand when the assignment of the instrument is non-staggered and a no carry over assumption is satisfied in the first stage regression. Theorem 3 instead considers the case

when the assignment of the instrument is staggered and there exist dynamic effects in the first stage. Theorem 3 assumes a binary treatment, but a non-binary, ordered treatment case is easy to extend: one can obtain the theorem which replaces $CLATT_{e,t}$ with $CACRT_{e,t}$.

If one wants to check the validity of the TWFEIV estimator in a given application, one can estimate each weight by constructing the consistent estimator for $CAET_{e,t}^1$. If there does not exist a never exposed cohort, however, it is not feasible to obtain the consistent estimator for $CAET_{l,t}^1$ in the last exposed cohort $l = \max\{E_i\}$. In Appendix C, we provide another representation of the decomposition theorem, in which we can estimate each weight consistently and quantify the bias term arising from the bad comparisons performed by TWFEIV regressions.

Random assignment of the adoption date

In practice, researchers may use the TWFEIV regression when the adoption date of the instrument is randomized across units (e.g., Randomized control trial). In Appendix C, we consider the causal interpretation of the TWFEIV estimand under the random assignment assumption. In the DID literature, a similar issue is analyzed in Athey and Imbens (2022): they investigate the causal interpretation of the TWFE estimand when the adoption date of the treatment is randomized across units.

First, we define the random assignment assumption of the adoption date E_i .

Assumption 14 (Random assignment assumption of adoption date E_i). For all $t \in \{1, \dots, T\}$ and all $z \in \mathcal{S}(Z)$, E_i is independent of potential outcomes:

$$(Y_{i,t}(1), Y_{i,t}(0), D_{i,t}(z)) \perp\!\!\!\perp E_i.$$

When the assignment of the adoption date is totally randomized, our target parameter is the local average treatment effect (LATE) defined below.

Def. The local average treatment effect (LATE) at a given relative period l from the initial adoption of the instrument is

$$LATE_{e,l} = E[Y_{i,e+l}(1) - Y_{i,e+l}(0) | CM_{e,e+l}].$$

Unlike the CLATT, this parameter is not conditioned on the adoption date E_i due to the independence assumption. The causal parameter in the first stage, $CAET_{e,l}^1$, is also simplified to the average exposed effect ($AE_{e,l}^1$) defined below:

$$\begin{aligned} CAET_{e,l}^1 &= E[D_{i,e+l} - D_{i,e+l}^\infty] \\ &\equiv AE_{e,l}^1. \end{aligned}$$

If Assumptions 1- 7 and Assumption 14 hold, one can obtain the theorem and lemmas in section 4, which replace $CAET_{k,t}^1$ and $CLATT_{k,t}$ with $AE_{k,t}^1$ and $LATE_{k,t}$, respectively. This implies that even when the adoption date of the instrument is randomized, we cannot interpret the TWFEIV estimand causally in general, and the causal interpretation requires the stable exposed assumptions in both the first stage and reduced form regressions.

6 Application

In this section, we illustrate our DID-IV decomposition theorem in the setting of Miller and Segal (2019). We first explain our dataset. We then assess the plausibility of the staggered DID-IV identification strategy implicitly imposed by Miller and Segal (2019). Finally, we present the DID-IV decomposition result and state the implication.

Miller and Segal (2019) study the effect of an increase in the share of female police officers on intimate partner homicide (IPH) rates among women in the United States between 1977 and 1991. The increase was in line with a shift in gender norms during these periods and there was growing interest in whether the female integration improved police quality in addressing violence against women.

To establish the causal relationship, Miller and Segal (2019) first regress the IPH rates on the lagged female officers' share with county and year fixed effects. In the second part of their analysis, Miller and Segal (2019) exploit "plausibly exogenous variation in female integration from externally imposed AA (affirmative action) following employment discrimination cases against particular departments in different years" across 255 counties. Specifically, Miller and Segal (2019) use the two-way fixed effects instrumental variable regression, instrumenting the lagged female officers' share with the exposure years of AA plans.

Miller and Segal (2019) implicitly rely on staggered DID-IV designs to estimate the causal effects: Miller and Segal (2019) concern that "AA itself might have occurred following increasing trends" in the share of female officers or the IPH rates. To address this concern, Miller and Segal (2019) check the trends of these variables before AA introduction using event study regressions in the first stage and reduced form.

In this application, we slightly modify the authors' setting for simplicity. Specifically, unlike Miller and Segal (2019), we use the staggered adoption of AA plans as our instrument instead of the exposure years. In the authors' setting, AA plans were terminated in some counties during the sample period, which is probably the reason why Miller and Segal (2019) use the exposure years of AA plans as their instrument. We instead drop such counties from our sample and make the instrument assignment staggered. Although it reduces our sample size, it allows us to have a clearer staggered DID-IV identification strategy. In addition, it enables us to apply our DID-IV decomposition theorem to the TWFEIV estimate in the authors' setting.

Data

The data come from Miller and Segal (2019). Our final sample differs from their main analysis sample in two ways. First, unlike Miller and Segal (2019), we only include the counties whose variables are observable for all sample periods. This restriction excludes 20 counties and allows us to create the balanced panel data set. Second, as we already noted, we construct an instrument that takes one after the AA introduction. Miller and Segal (2019) use data on AA plans from Miller and Segal (2012) and define the instrument as the difference between the current year and the start year of AA introduction¹; see Miller and Segal (2012), Miller and Segal (2019) for details. We identify the initial year of AA plans in each county, and discard the counties whose AA plans ended between 1976 and 1990 (8 counties dropped) and whose AA plans were already implemented before 1976 (23 counties dropped). Table 1 shows the timing of AA adoption across 199 counties between 1976 and 1990.

Summary statistics for county characteristics are reported in Table 2. We have a smaller sample size, but otherwise have a similar sample to that of Miller and Segal (2019). Counties are separated into exposed and unexposed counties based on whether the county experienced AA introduction. In both types of counties, the lagged female officers' share increased over time. However, it increases more in counties who are exposed to AA plans during sample periods. The IPH rates had downward trends in all counties, but it seems that there are no

¹As one can see in this construction, Miller and Segal (2019) create the lagged instrument in line with the lagged female officers' share. Therefore, we construct the lagged staggered instrument instead of the current one.

Table 1

The staggered AA introduction: exposure year, cohort sizes.

Start year of AA plans	Number of counties
Unexposed counties	159
1976	6
1977	3
1978	3
1979	4
1980	3
1981	5
1982	4
1983	3
1984	2
1985	1
1986	1
1987	3
1988	1
1990	1

Notes: This table presents the initial exposure year of AA plans and the number of counties in each year in our final sample.

Table 2

Summary statistics

		All counties	Unexposed counties	Exposed counties
IPH per 100000 population	1977-91	0.544	0.521	0.638
	1977	0.549	0.526	0.641
	1991	0.489	0.461	0.599
Lagged female officer share	1977-91	0.053	0.050	0.066
	1977	0.033	0.033	0.032
	1991	0.077	0.071	0.101
Counties		199	159	40
Observations		2985	2385	600

Notes: This table presents summary statistics on our final sample from 1977 to 1991. The sample consists of 199 counties.

systematic differences in the trends between exposed and unexposed counties.

Assessing the identifying assumptions in staggered DID-IV design

In this section, we discuss the validity of the staggered DID-IV identification strategy implicitly imposed by [Miller and Segal \(2019\)](#). Note that in the authors' setting, our target parameter is the cohort specific average causal response on the treated (CACRT) as female officer share is a non-binary, ordered treatment. We therefore expect that we can identify each CACRT if the underlying staggered DID-IV identification strategy seems plausible, which we will check below. Here, we presume the no carry over assumption ([Assumption 2](#)).

Exclusion restriction ([Assumption 3](#)). It would be plausible, given that the AA plans (instrument) did not affect IPH rates other than by increasing the female officers' share. This assumption may be violated for instance if the AA plans increased both the black and female

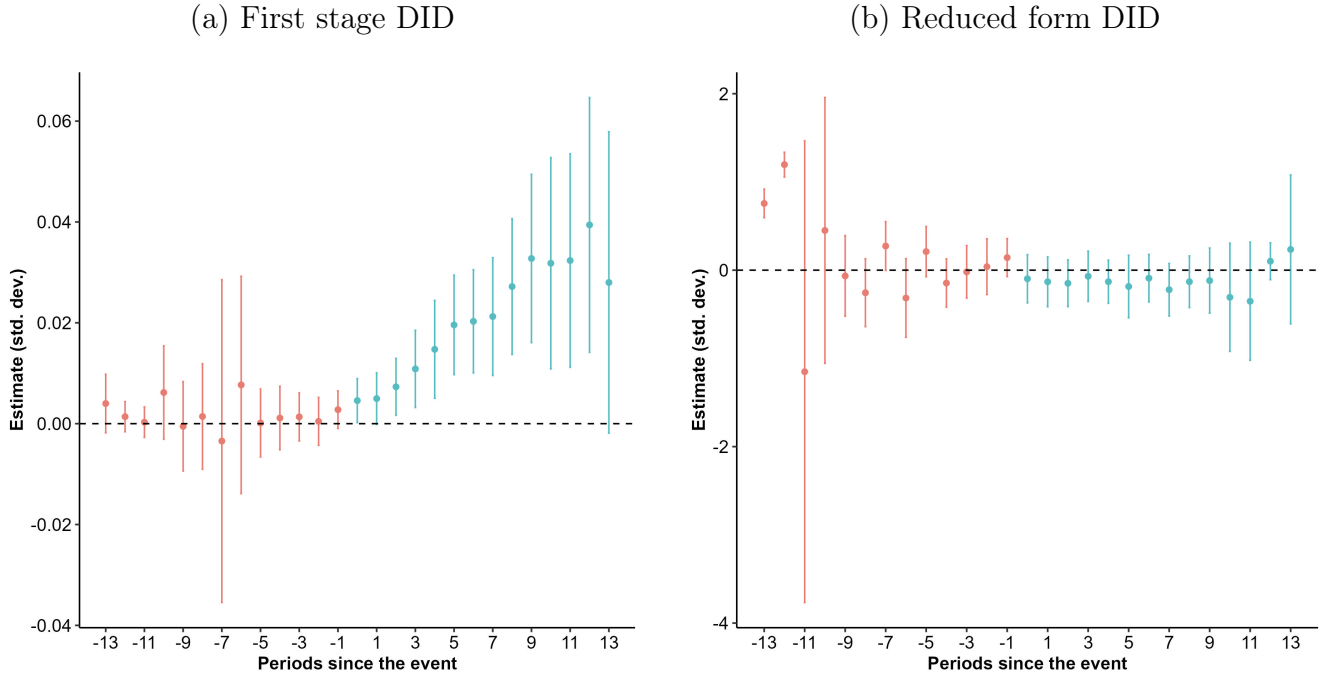


Fig. 2. Weighted average of the effects of AA introduction on female officer share and IPH rates in [Miller and Segal \(2019\)](#). *Notes:* The results for the effects of AA plans on the female officers share (Panel (a)) and on the IPH rates (Panel (b)) under the staggered DID-IV identification strategy. The red line represents the weighted average of the estimates with simultaneous 95% confidence intervals for pre-exposed periods in both panels where the weight reflects cohort size in each period. The control group is a never exposed cohort and the reference period is $t - 1$ for period t estimate. These should be equal to zero under the null hypothesis that parallel trends assumptions in the treatment and the outcome hold. The blue line represents the weighted average of the estimates with simultaneous 95% confidence intervals for post exposed periods in both panels where the weight reflects cohort size in each period. The control group is a never exposed cohort and the reference period is $t = -1$ for all post-exposed period estimates. All the standard errors are clustered at county level.

officer shares and changes in IPH rates reflect both effects. [Miller and Segal \(2019\)](#) conduct the robustness check and confirm that this is not the case; see footnote 42 in [Miller and Segal \(2019\)](#) for details.

Monotonicity assumption (Assumption 4). It would be automatically satisfied in the authors' setting: the AA plans (instrument) were imposed on departments with the intent to increase the share of female police officers. This ensures that the dynamic effects of the instrument on female police officers should be non-negative after the AA introduction.

No anticipation in the first stage (Assumption 5). It would be plausible that there is no anticipatory behavior, given that the treatment status, i.e., the female officers share before the AA plans is equal to the one in the absence of the AA introduction across counties. This assumption may be violated if the police departments in some counties had private knowledge about the probability of the AA introduction and manipulated their treatment status before the implementation.

Next, we assess the plausibility of the parallel trends assumptions in the treatment and the outcome. To do so, we apply the method proposed by [Callaway and Sant'Anna \(2021\)](#) to the

first stage and reduced form, respectively². Specifically, we estimate the weighted average of the effects of the instrument on the treatment and outcome in each relative period where the weight reflects the cohort size. We depict the results in Figure 2. The plots report estimates for the effects before and after AA plans with a simultaneous 95% confidence interval in each stage. The confidence intervals account for clustering at the county level.

Parallel trends assumption in the treatment (Assumption 6). It requires that if the AA plans had not occurred, the average time trends of the female officers share would have been the same across counties and over time. The pre-exposed estimates in Panel (a) in Figure 2 seem consistent with the parallel trends assumption in the treatment: the pre-exposed estimates around AA plans are not significantly different from zero.

Parallel trends assumption in the outcome (Assumption 7). It would be plausible if the AA plans had not been implemented, the average time trends of the IPH rates would have been the same across counties and over time. Panel (b) in Figure 2 presents that the pre-exposed estimates around AA introduction are not significantly different from zero, which indicates that the parallel trends assumption in the outcome is also plausible.

Figure 2 also sheds light on the dynamic effects of the AA plans on the female officer share and IPH rates during the post-exposed periods. The figure indicates that the effect of the AA plans on the female officer share increases over time, whereas the effect on IPH rates through the female officer share has downward trends during the post-exposed periods. We note that the estimated effects in the reduced form are not scaled by the ones in the first stage, i.e., these estimates do not capture each CACRT after the AA shock.

Illustrating the weights in TWFEIV regression

First, we estimate the two-way fixed effects instrumental variable regression in the authors' setting. To clearly illustrate the shortcomings of the TWFEIV regression, we modify the authors' specification in two ways: Miller and Segal (2019) include some covariates and weight their regression with county population, whereas we exclude such covariates and do not apply their weights to our regression.

The result is shown in Table 3. The two-way fixed effects instrumental variable estimate is -0.646 and it is not significantly different from zero³. However, as we already noted in section 4, we cannot generally interpret the IV estimate as measuring a properly weighted average of each CACRT if the effect of the AA introduction on female officer share or IPH rates is not stable over time.

Our DID-IV decomposition theorem (Theorem 1) allows us to visualize the source of variations in the three types of the DID-IV design: Unexposed/Exposed, Exposed/Not Yet Exposed, and Exposed/Exposed Shift designs. Panel (a) in Figure 3 plots the weights and the corresponding Wald-DID estimates for all designs and Panels (b), (c), and (d) in Figure 3 plot them for each type of the DID-IV design, respectively. Table 4 reports the total weight, total Wald-DID estimate, and weighted average of Wald-DID estimates in each type of the DID-IV design.

²Unfortunately, in the presence of heterogeneous treatment effects, the coefficients on event study regression face a contamination bias shown by Sun and Abraham (2021).

³Although Miller and Segal (2019) do not report the TWFEIV estimate without weights and covariates, when we run such a TWFEIV regression in their final analysis sample, the IV estimate is -1.445 and is not significantly different from zero. This implies that we reach the same conclusion as in Miller and Segal (2019) in our data.

The total weight and total Wald-DID estimate are calculated by summing the weights and Wald-DID estimates respectively, and the weighted average of Wald-DID estimates is calculated by summing the products of the weight and the associated Wald-DID estimate. Summing all the weighted average of Wald-DID estimates yields the two-way fixed instrumental variable estimate (-0.646).

Panel (a) in Figure 3 shows that the weights are heavily assigned to the Wald-DID estimates in Unexposed/Exposed designs. This is due to the large sample size of the unexposed cohort in the authors' setting. Panels (b), (c), and (d) in Figure 3 highlight that some weights in each type can be negative: 2 out of 14 weights are negative in Unexposed/Exposed designs, 29 out of 91 weights are negative in Exposed/Not Yet Exposed designs and 50 out of 91 weights are negative in Exposed/Exposed Shift designs. The negative weights arise because some DID estimates of the treatment in the first stage are negative in each type of the DID-IV design.

The TWFEIV estimate suffers from a downward bias due to the bad comparisons arising from the Exposed/Exposed shift designs. As we already mentioned in section 4, the TWFEIV estimand potentially fails to summarize the causal effects if the effect of the instrument on the treatment or the outcome evolves over time. Table 4 indicates that the estimated bias occurring from the Exposed/Exposed shift designs is quantitatively not negligible: the weighted average of the Wald-DID estimates in the Exposed/Exposed shift designs is -0.093 , which accounts for one-seventh of our IV estimate.

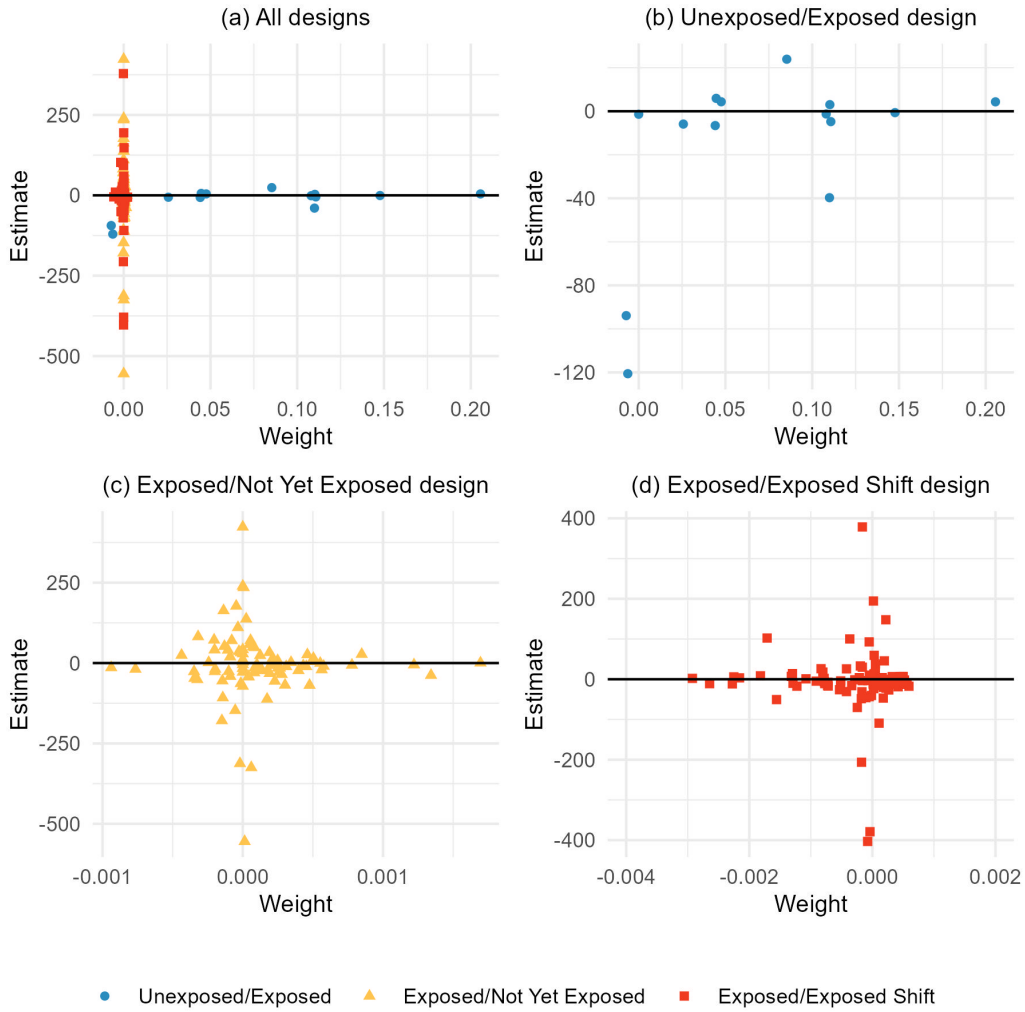


Fig. 3. Instrumented difference-in-differences decomposition result in the setting of Miller and Segal (2019). *Notes:* Panel (a) plots the weights and the corresponding Wald-DID estimates for all DID-IV designs and Panels (b), (c), and (d) plot them for each type of the DID-IV design, respectively. Unexposed/Exposed designs yield blue circles, Exposed/Not Yet Exposed designs yield yellow triangles and Exposed/Exposed Shift designs yield red squares.

Table 3

Estimate for the effect of female officers share on IPH rates.

	Estimate	Standard Error	95% CI
TSLS with fixed effects	-0.646	3.284	[-7.594, 6.301]

Notes: Sample consists of 199 counties. Confidence intervals account for clustering at the county level.

Table 4

Total weight, Total and Weighted WDD estimates in each type of the DID-IV design.

	Total weight	Total WDD estimate	Weighted WDD estimate
Unexposed/Exposed	1.026	-233.198	-0.399
Exposed/Not Yet Exposed	0.010	-490.665	-0.154
Exposed/Exposed Shift	-0.036	-569.426	-0.093

Notes: This table presents the total weight, total Wald-DID estimate, and weighted average of Wald-DID estimates in each type of the DID-IV design: Unexposed/Exposed, Exposed/Not Yet Exposed, and Exposed/Exposed Shift designs.

7 Alternative specifications

So far, we have considered simple TWFEIV regressions as in equation (1). However, many studies routinely estimate various specifications, such as weighting or introducing covariates, to check the robustness of their findings. In this section, we extend our DID-IV decomposition theorem to the settings with weighting and covariates, and provide simple tools to examine how different specifications affect differences in estimates. We illustrate these by revisiting [Miller and Segal \(2019\)](#).

The tools we provide here are based on [Goodman-Bacon \(2021\)](#). Recall that our DID-IV decomposition theorem shows that the TWFEIV estimator can be written as the product of a vector of 2×2 Wald-DID estimators ($\hat{\beta}_{IV}^{2 \times 2}$) and a vector of weights (\mathbf{s}), that is, $\hat{\beta}^{IV} = \mathbf{s}' \hat{\beta}_{IV}^{2 \times 2}$. When a TWFEIV estimator generated from different specification ($\hat{\beta}_{IV,alt}$) can also be written as the product of a vector of 2×2 Wald-DID estimators ($\hat{\beta}_{IV,alt}^{2 \times 2}$) and a vector of their associated weights (\mathbf{s}_{alt}), one can decompose the difference between the two specifications as

$$\hat{\beta}_{IV,alt} - \hat{\beta}_{IV} = \underbrace{\mathbf{s}'(\hat{\beta}_{IV,alt}^{2 \times 2} - \hat{\beta}_{IV}^{2 \times 2})}_{\text{Due to } 2 \times 2 \text{ Wald-DIDs}} + \underbrace{(\mathbf{s}'_{alt} - \mathbf{s}')\hat{\beta}_{IV}^{2 \times 2}}_{\text{Due to } 2 \times 2 \text{ weights}} + \underbrace{(\mathbf{s}'_{alt} - \mathbf{s}')(\hat{\beta}_{IV,alt}^{2 \times 2} - \hat{\beta}_{IV}^{2 \times 2})}_{\text{Due to the interaction of the two}}.$$

It takes the form of a Oaxaca-Blinder-Kitagawa decomposition ([Oaxaca \(1973\)](#), [Blinder \(1973\)](#), [Kitagawa \(1955\)](#)) and indicates that the difference comes from changes in 2×2 Wald-DID estimators, changes in weights, and the interaction of the two. Dividing both sides by $\hat{\beta}_{IV,alt} - \hat{\beta}_{IV}$, one can measure the proportional contribution of each term on the difference. Plotting each pair in $(\hat{\beta}_{IV,alt}^{2 \times 2}, \hat{\beta}_{IV}^{2 \times 2})$ and $(\mathbf{s}'_{alt}, \mathbf{s}')$, one can also examine which elements in each term have a significant impact on the difference.

7.1 Weighted TWFEIV regression

When researchers use weighted TWFEIV regression instead of unweighted one, it potentially changes the influence of Wald-DID estimators ($\hat{\beta}_{IV,WLS}^{2 \times 2}$) by replacing the DIDs of the treatment and the outcome with the weighted ones. It also potentially change the influence of weights (\mathbf{s}'_{WLS}) by replacing the sample share with the relative amount of the specified weight and the DIDs of the treatment with the weighted ones. Table 5 shows the result of our TWFEIV regression weighted by county population in [Miller and Segal \(2019\)](#): the estimate changes from -0.646 to -0.386 . The decomposition result indicates that the contribution of the changes in 2×2 Wald-DIDs is negative, whereas the contributions of the changes in weights and the interaction are positive.

Figure 4 plots the 2×2 Wald-DIDs and the associated weights in WLS against those in OLS. Panel (a) shows that most comparisons of the Wald-DID between OLS and WLS are located at the 45-degree line, but some comparisons generated from Exposed/Not Yet Exposed and Exposed/Exposed Shift designs are away from the 45-degree line. In addition, this figure indicates that the Wald-DID generated from the comparison between 1978 and 1991 counties (1991 counties are the controls) is much more negative in WLS than in OLS, which drives the overall negative impact of the changes in 2×2 Wald-DIDs on the difference between the two specifications. Panel (b) shows that most comparisons of the decomposition weight between OLS and WLS are near the 45-degree line and the origin, but some comparisons generated from Unexposed/Exposed designs are away from the 45-degree line and the origin. This figure also indicates that the decomposition weight generated from the comparison between 1982 and unexposed counties is much more positive in WLS than in OLS, which causes the overall positive impact of the changes in weights on the difference between the two specifications.

Table 5

Estimate for the effect of female officers share on IPH rates.

	(1)	(2)	(3)
	Baseline	WLS	Covariates
Estimate	-0.646	-0.386	-0.868
Standard Error	3.284	2.452	3.968
Difference from baseline		0.260	-0.222
Difference comes from:			
2×2 Wald-DIDs		-4.048	0.370
Weights		2.341	16.503
Interaction		1.966	-17.107
Within term		0	0.012

Notes: This table presents TWFEIV estimates in the setting of [Miller and Segal \(2019\)](#). Column (1) is a simple TWFEIV estimate from Eq. (1). Column (2) is a TWFEIV estimate weighted by county population in 1977. Column (3) is a TWFEIV estimate with time-varying covariates which include the lagged local area controls, the county's non-IPH rate, and the state-level crack cocaine index. All the standard errors are clustered at county level.

7.2 TWFEIV regression with time-varying covariates

In most applications of the DID-IV method, researchers typically estimate TWFEIV models that include time-varying covariates, in addition to the simple ones, based on the belief that it enhances the validity of the parallel trends assumptions in the first stage and reduced form regressions:

$$Y_{i,t} = \phi_i + \lambda_t + \beta_{IV}^X D_{i,t} + \psi X_{i,t} + v_{i,t}, \quad (14)$$

$$D_{i,t} = \gamma_i + \zeta_t + \pi^X Z_{i,t} + \tilde{\psi} X_{i,t} + \eta_{i,t}. \quad (15)$$

In this section, we derive a DID-IV decomposition result for the case when we introduce the time-varying covariates into TWFEIV regressions. Our decomposition result in this section is based on [Goodman-Bacon \(2021\)](#), who decomposes TWFE estimators with time-varying covariates. Appendix D further considers the causal interpretation of the covariate-adjusted TWFEIV estimand under additional conditions.

First, consider the coefficient on instrument (α^X) in the reduced form regression:

$$Y_{i,t} = \phi_i + \lambda_t + \alpha^X Z_{i,t} + \xi X_{i,t} + v_{i,t}. \quad (16)$$

Let $\tilde{Z}_{i,t}$ and $\tilde{X}_{i,t}$ denote the double demeaning variables of $Z_{i,t}$ and $X_{i,t}$ respectively, obtained from regressing $Z_{i,t}$ and $X_{i,t}$ on time and unit fixed effects. Let $\tilde{z}_{i,t}$ denote the residuals obtained from regressing $\tilde{Z}_{i,t}$ on $\tilde{X}_{i,t}$:

$$\tilde{Z}_{i,t} = \hat{\Gamma} \tilde{X}_{i,t} + \tilde{z}_{i,t},$$

Here, we define the linear projection as $\tilde{p}_{i,t} \equiv \hat{\Gamma} \tilde{X}_{i,t}$. The specific expression for $\tilde{z}_{i,t}$ is:

$$\begin{aligned} \tilde{z}_{i,t} &= [(Z_{i,t} - \bar{Z}_i) - (\hat{\Gamma} X_{i,t} - \hat{\Gamma} \bar{X}_i)] - [(\bar{Z}_t - \bar{Z}) - (\hat{\Gamma} \bar{X}_t - \hat{\Gamma} \bar{X})] \\ &\equiv (z_{i,t} - \bar{z}_i) - (\bar{z}_t - \bar{z}). \end{aligned}$$

By the FWL theorem, we then obtain the following expression for $\hat{\alpha}^X$:

$$\hat{\alpha}^X = \frac{\hat{C}(Y_{i,t}, \tilde{z}_{i,t})}{\hat{V}^{\tilde{z}}} = \frac{\hat{C}(Y_{i,t}, \tilde{Z}_{i,t} - \tilde{p}_{i,t})}{\hat{V}^{\tilde{z}}},$$

where $\hat{V}^{\tilde{z}}$ is the variance of $\tilde{z}_{i,t}$. By symmetry, we can also express the first stage coefficient on instrument $\hat{\pi}^X$ as follows:

$$\hat{\pi}^X = \frac{\hat{C}(D_{i,t}, \tilde{z}_{i,t})}{\hat{V}^{\tilde{z}}} = \frac{\hat{C}(D_{i,t}, \tilde{Z}_{i,t} - \tilde{p}_{i,t})}{\hat{V}^{\tilde{z}}}.$$

Because the IV estimator $\hat{\beta}_{IV}^X$ is the ration between the first stage coefficient $\hat{\pi}^X$ and the reduced form coefficient $\hat{\alpha}^X$, we obtain the following expression for $\hat{\beta}_{IV}^X$:

$$\hat{\beta}_{IV}^X = \frac{\hat{C}(Y_{i,t}, \tilde{z}_{i,t})}{\hat{C}(D_{i,t}, \tilde{z}_{i,t})} = \frac{\hat{C}(Y_{i,t}, \tilde{Z}_{i,t} - \tilde{p}_{i,t})}{\hat{C}(D_{i,t}, \tilde{Z}_{i,t} - \tilde{p}_{i,t})}. \quad (17)$$

In contrast to the unconditional TWFEIV estimator $\hat{\beta}_{IV}$, the covariate-adjusted TWFEIV estimator exploits the variation in both $\tilde{Z}_{i,t}$ and $\tilde{p}_{i,t}$. $\tilde{Z}_{i,t}$ varies at cohort and time level, but $\tilde{p}_{i,t}$ varies at unit and time level because $X_{i,t}$ varies at unit and time level.

To decompose the covariate-adjusted TWFEIV estimator $\hat{\beta}_{IV}^X$, we first partition $\tilde{z}_{i,t}$ into "within" and "between" terms as in [Goodman-Bacon \(2021\)](#). Let $\bar{z}_{k,t} - \bar{z}_k = (\bar{Z}_{k,t} - \bar{Z}_k) - (\hat{\Gamma}\bar{X}_{k,t} - \hat{\Gamma}\bar{X}_k)$ be the average of $z_{i,t} - \bar{z}_i$ in cohort k . By adding and subtracting $\bar{z}_{k,t} - \bar{z}_k$, we can decompose $\tilde{z}_{i,t}$ into two terms:

$$\tilde{z}_{i,t} = \underbrace{[(z_{i,t} - \bar{z}_i) - (\bar{z}_{k,t} - \bar{z}_k)]}_{\tilde{z}_{i(k),t}} + \underbrace{[(\bar{z}_{k,t} - \bar{z}_k) - (\bar{z}_t - \bar{z})]}_{\tilde{z}_{k,t}}. \quad (18)$$

The first term $\tilde{z}_{i(k),t}$ measures the deviation of $z_{i,t} - \bar{z}_i$ from the average $\bar{z}_{k,t} - \bar{z}_k$ in cohort k , which we call the within term of $\tilde{z}_{i,t}$. The second term $\tilde{z}_{k,t}$ measures the deviation of $\bar{z}_{k,t} - \bar{z}_k$ from the average $\bar{z}_t - \bar{z}$ in whole sample, which we call the between term of $\tilde{z}_{i,t}$. The within term $\tilde{z}_{i(k),t}$ varies at unit and time level because of $\tilde{p}_{i,t}$, whereas the between term $\tilde{z}_{k,t}$ varies at cohort and time level.

By substituting (18) into (17), we obtain

$$\hat{\beta}_{IV}^X = \frac{\hat{C}(Y_{i,t}, \tilde{z}_{i(k),t}) + \hat{C}(Y_{i,t}, \tilde{z}_{k,t})}{\hat{C}(D_{i,t}, \tilde{z}_{i(k),t}) + \hat{C}(D_{i,t}, \tilde{z}_{k,t})} = \frac{\hat{V}_w^z \hat{\beta}_w^{p,y} + \hat{V}_b^z \hat{\beta}_b^{z,y}}{\hat{V}_w^z \hat{\beta}_w^{p,d} + \hat{V}_b^z \hat{\beta}_b^{z,d}} \quad (19)$$

$$= \underbrace{\frac{\hat{C}_w^{D,\tilde{z}}}{\hat{C}_w^{D,\tilde{z}} + \hat{C}_b^{D,\tilde{z}}}}_{\Omega} \cdot \underbrace{\frac{\hat{\beta}_w^{p,y}}{\hat{\beta}_w^{p,d}}}_{\hat{\beta}_{w,IV}^p} + \underbrace{\frac{\hat{C}_b^{D,\tilde{z}}}{\hat{C}_w^{D,\tilde{z}} + \hat{C}_b^{D,\tilde{z}}}}_{1-\Omega} \cdot \underbrace{\frac{\hat{\beta}_b^{z,y}}{\hat{\beta}_b^{z,d}}}_{\hat{\beta}_{b,IV}^{z,y}}. \quad (20)$$

We use the subscript w to denote within components and the subscript b to denote between components. \hat{V}_w^z and \hat{V}_b^z are the variances of $\tilde{z}_{i(k),t}$ and $\tilde{z}_{k,t}$, respectively. $\hat{C}_w^{D,\tilde{z}}$ is the covariance between $D_{i,t}$ and $\tilde{z}_{i(k),t}$, the within term of $\tilde{z}_{i,t}$. $\hat{C}_b^{D,\tilde{z}}$ is the covariance between $D_{i,t}$ and $\tilde{z}_{k,t}$, the between term of $\tilde{z}_{i,t}$. The weight $\Omega = \frac{\hat{C}_w^{D,\tilde{z}}}{\hat{C}_w^{D,\tilde{z}} + \hat{C}_b^{D,\tilde{z}}}$ measures the relative amount of the within covariance $\hat{C}_w^{D,\tilde{z}}$.

$\hat{\beta}_w^{p,y} \equiv \frac{\hat{C}(Y_{i,t}, \tilde{z}_{i(k),t})}{\hat{V}_w^z}$ measures the relationship between $Y_{i,t}$ and $\tilde{z}_{i(k),t}$. Similarly, $\hat{\beta}_w^{p,d} \equiv \frac{\hat{C}(D_{i,t}, \tilde{z}_{i(k),t})}{\hat{V}_w^z}$ measures the relationship between $D_{i,t}$ and $\tilde{z}_{i(k),t}$. We call these the within coefficients in the first stage and reduced form regressions. $\hat{\beta}_{w,IV}^p \equiv \frac{\hat{\beta}_w^{p,y}}{\hat{\beta}_w^{p,d}}$ scales the within coefficient in the reduced form regression by the one in the first stage regression. We call this the within

IV coefficient⁴. This IV coefficient arises because $\tilde{z}_{i(k),t}$ varies at unit and time level. Similar to what [Goodman-Bacon \(2021\)](#) points out for the covariate-adjusted TWFE estimator, time-varying covariates bring a new source of identifying variation in the TWFEIV estimator, within variation of $X_{i,t}$ in each cohort.

$\hat{\beta}_b^{z,y} \equiv \frac{\hat{C}(Y_{i,t}, \tilde{z}_{k,t})}{\hat{V}_b^z}$ measures the relationship between $Y_{i,t}$ and $\tilde{z}_{k,t}$. Similarly $\hat{\beta}_b^{z,d} \equiv \frac{\hat{C}(D_{i,t}, \tilde{z}_{k,t})}{\hat{V}_b^z}$ measures the relationship between $D_{i,t}$ and $\tilde{z}_{k,t}$. We call these the between coefficients in the first stage and reduced form regressions. $\hat{\beta}_{b,IV}^z \equiv \frac{\hat{\beta}_b^{z,y}}{\hat{\beta}_b^{z,d}}$ divides the between coefficient in the reduced form regression by the one in the first stage regression, and have the following specific expression:

$$\hat{\beta}_{b,IV}^z = \frac{\hat{C}^{D,Z} \hat{\beta}_{IV} - \hat{C}_b^p \hat{\beta}_{b,IV}^p}{\hat{C}_b^{D,\tilde{z}}}. \quad (21)$$

$\hat{C}^{D,Z}$ and $\hat{\beta}_{IV}$ are already defined in section 3. \hat{C}_b^p is the covariance between $D_{i,t}$ and $\tilde{p}_{k,t}$ (the between term of $\tilde{p}_{i,t}$). $\hat{\beta}_{b,IV}^p$ is the estimator, obtained from an IV regression of $Y_{i,t}$ on $D_{i,t}$ with $\tilde{p}_{k,t}$ as the excluded instrument. We call $\hat{\beta}_{b,IV}^z$ the between IV coefficient, which exploits the cohort and time level variation in $\tilde{z}_{k,t}$. This IV coefficient is not equal to the unconditional TWFEIV coefficient $\hat{\beta}_{IV}$: $\hat{\beta}_{b,IV}^z$ subtracts the influence of $\hat{\beta}_{b,IV}^p$ from the unconditional IV estimator $\hat{\beta}_{IV}$. This indicates that time-varying covariates $X_{i,t}$ changes the identifying variation at cohort and time level through $\tilde{p}_{k,t}$, the between term of the linear projection $\tilde{p}_{i,t}$.

We can further decompose the between IV coefficient as follows:

$$\hat{\beta}_{b,IV}^z = \sum_k \sum_{l>k} \underbrace{(n_k + n_l)^2}_{s_{b,kl}} \frac{\hat{C}_{b,kl}^{D,\tilde{z}}}{\hat{C}_b^{D,\tilde{z}}} \underbrace{\left[\frac{\hat{C}_{kl}^{D,Z} \hat{\beta}_{IV,kl}^{2 \times 2} - \hat{C}_{b,kl}^p \hat{\beta}_{b,IV,kl}^p}{\hat{C}_{b,kl}^{D,\tilde{z}}} \right]}_{\hat{\beta}_{b,IV,kl}^z}. \quad (22)$$

The proof is given in Appendix D. Each notation is similarly defined in (k, l) cell subsamples. $\hat{\beta}_{b,IV,kl}^z$ and $s_{b,kl}$ are the between IV coefficient and the corresponding weight in (k, l) cell subsamples. Equation (22) indicates that time-varying covariates $X_{i,t}$ affect the between IV coefficient $\hat{\beta}_{b,IV}^z$ by changing both the 2×2 between IV coefficient and the associated weight in each (k, l) cell.

To sum up, combining (22) with (19), we can decompose the covariate-adjusted TWFEIV estimator $\hat{\beta}_{IV}^X$ as

$$\hat{\beta}_{IV}^X = \Omega \hat{\beta}_{w,IV}^p + (1 - \Omega) \underbrace{\sum_k \sum_{l>k} s_{b,kl} \hat{\beta}_{b,IV,kl}^z}_{\hat{\beta}_{b,IV}^z}.$$

The weight Ω is assigned to the within IV coefficient $\hat{\beta}_{w,IV}^p$ and the weight $1 - \Omega$ is assigned to the between IV coefficient $\hat{\beta}_{b,IV}^z$, which is equal to a weighted average of all possible 2×2 between IV coefficients $\hat{\beta}_{b,IV,kl}^z$ as in Theorem 1.

Table 5 presents the result of our TWFEIV regression with time-varying covariates in [Miller and Segal \(2019\)](#). We follow [Miller and Segal \(2019\)](#) and include the lagged local area controls,

⁴One can obtain this coefficient by running an IV regression of the outcome on the treatment with $\tilde{z}_{i(k),t}$ as the excluded instrument.

the county's non-IPH rate, and the state-level crack cocaine index; see [Miller and Segal \(2019\)](#) for details. The estimate changes from -0.646 to -0.868 . The decomposition result shows that the contribution of the within term is positive but negligible, whereas the contribution of the between term is negative and substantial. Specifically, in the between term, the contribution of the changes in 2×2 Wald-DIDs and weights are positive, but these are offset by the negative contribution of the interaction. This result indicates that in [Miller and Segal \(2019\)](#), the time-varying covariates affect the IV estimate mainly through the identifying variation in cohort and time level, that is, the between term of the linear projection $\tilde{p}_{i,t}$.

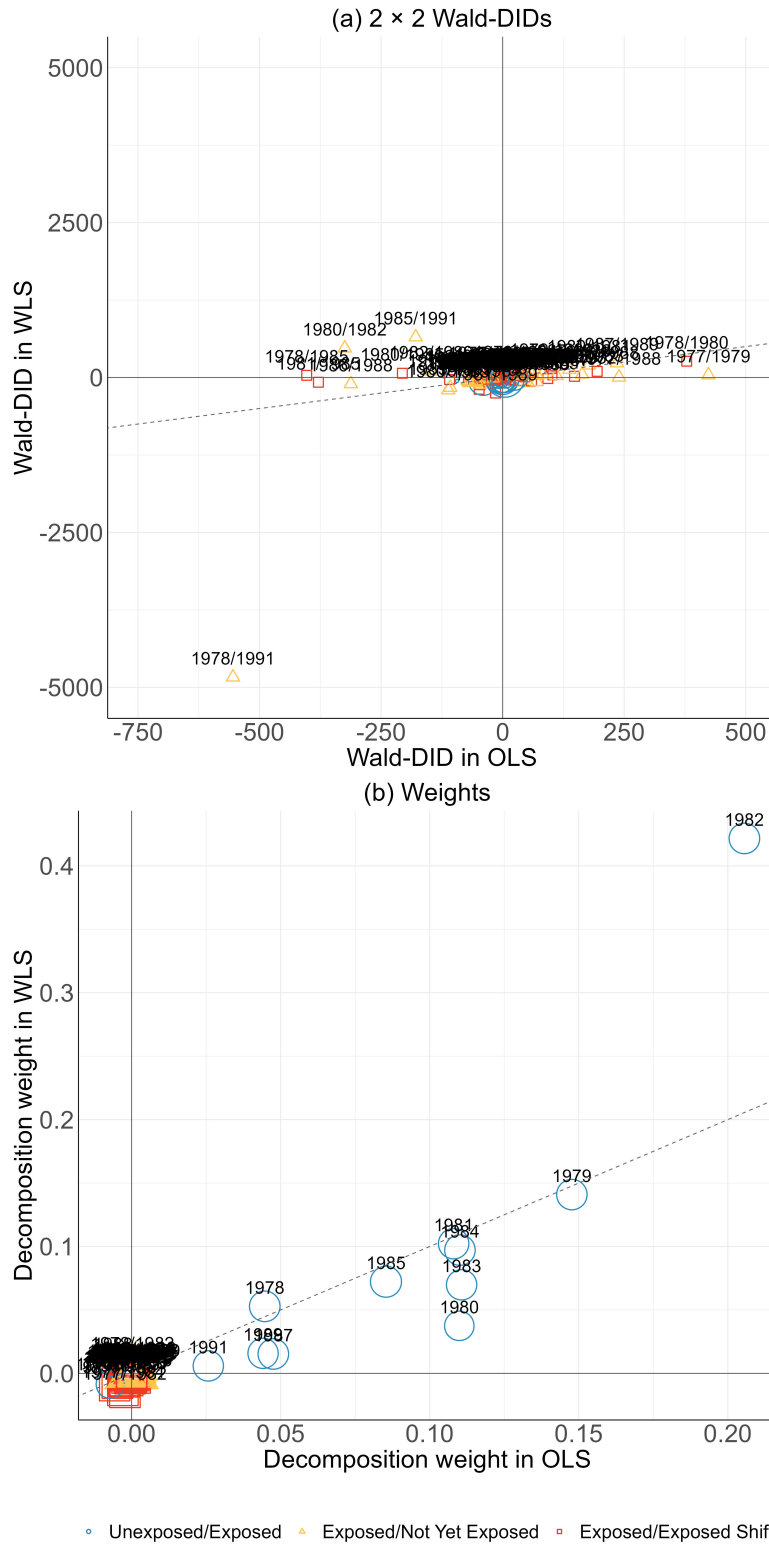


Fig. 4. Comparisons of 2×2 WaldDIDs and weights between OLS and WLS in the setting of Miller and Segal (2019). *Notes:* Panel (a) plots the 2×2 Wald-DIDs in WLS against those in OLS for all DID-IV designs. Panel (b) plots the decomposition weights in WLS against those in OLS for all DID-IV designs. In Panel (a), the size of each point is proportional to the corresponding weight in OLS. In Panel (b), the size of each point is proportional to the corresponding Wald-DID estimate in OLS. In both panels, the dotted lines represent 45-degree lines. In both panels, Unexposed/Exposed designs yield blue circles, Exposed/Not Yet Exposed designs yield yellow triangles, and Exposed/Exposed Shift designs yield red squares. In both panels, the dotted lines represent 45-degree lines.

8 Conclusion

Many studies run two-way fixed effects instrumental variable (TWFEIV) regressions, leveraging variation occurring from the different timing of policy adoption across units as an instrument for the treatment. In this paper, we study the causal interpretation of the TWFEIV estimator in staggered DID-IV designs. We first show that in settings with the staggered adoption of the instrument across units, the TWFEIV estimator is equal to a weighted average of all possible 2×2 Wald-DID estimators arising from the three types of the DID-IV design: Unexposed/Exposed, Exposed/Not Yet Exposed, and Exposed/Exposed Shift designs. The weight assigned to each Wald-DID estimator is a function of the sample share, the variance of the instrument, and the DID estimator of the treatment in each DID-IV design.

Based on the decomposition result, we then show that in staggered DID-IV designs, the TWFEIV estimand is equal to a weighted average of all possible cohort specific local average treatment effect on the treated parameters, but some weights can be negative. The negative weight problem arises due to the bad comparisons in the first and reduced form regressions: we use the already exposed units as controls. The TWFEIV estimand attains its causal interpretation if the effects of the instrument on the treatment and outcome are stable over time. The resulting causal parameter is a positively weighted average cohort specific local average treatment effect on the treated parameter.

Finally, we illustrate our findings with the setting of [Miller and Segal \(2019\)](#) who estimate the effect of female officers' share on the IPH rate, exploiting the timing variation of AA introduction across U.S. counties. We first assess the underlying staggered DID-IV identification strategy implicitly imposed by [Miller and Segal \(2019\)](#) and confirm its validity. We then apply our DID-IV decomposition theorem to the TWFEIV estimate, and find that the estimate suffers from the substantial downward bias arising from the bad comparisons in Exposed/Exposed shift DID-IV designs. We also decompose the difference between the two specifications and illustrate how different specifications affect the overall estimates in [Miller and Segal \(2019\)](#).

Overall, this paper shows the negative result of using TWFEIV estimators in the presence of heterogeneous treatment effects in staggered DID-IV designs in more than two periods. This paper provides simple tools to evaluate how serious that concern is in a given application. Specifically, we demonstrate that the TWFEIV estimator is not robust to the time-varying exposed effects in the first stage and reduced form regressions. Our DID-IV decomposition theorem allows the empirical researchers to assess the impact of the bias term arising from the bad comparisons on their TWFEIV estimate. Recently, [Miyaji \(2024\)](#) developed an alternative estimation method that is robust to treatment effects heterogeneity and proposes a weighting scheme to construct various summary measures in staggered DID-IV designs. Further developing alternative approaches and diagnostic tools will be a promising area for future work, facilitating the credibility of DID-IV design in practice.

References

- Akerman, A., Gaarder, I., Mogstad, M., 2015. The Skill Complementarity of Broadband Internet. *Q. J. Econ.* 130 (4), 1781–1824.
- Angrist, J. D., Imbens, G. W., 1995. Two-Stage Least Squares Estimation of Average Causal Effects in Models with Variable Treatment Intensity. *J. Am. Stat. Assoc.* 90 (430), 431–442.
- Athey, S., Imbens, G. W., 2022. Design-based analysis in Difference-In-Differences settings with staggered adoption. *J. Econom.* 226 (1), 62–79.
- Bhuller, M., Havnes, T., Leuven, E., Mogstad, M., 2013. Broadband Internet: An Information Superhighway to Sex Crime? *Rev. Econ. Stud.* 80 (4), 1237–1266.
- Black, S. E., Devereux, P. J., Salvanes, K. G., 2005. Why the Apple Doesn’t Fall Far: Understanding Intergenerational Transmission of Human Capital. *Am. Econ. Rev.* 95 (1), 437–449.
- Blandhol, C., Bonney, J., Mogstad, M., Torgovitsky, A., 2022. When is TSLS Actually LATE? February.
- Blinder, A. S., 1973. Wage Discrimination: Reduced Form and Structural Estimates. *J. Hum. Resour.* 8 (4), 436–455.
- Borusyak, K., Jaravel, X., Spiess, J., 2021. Revisiting Event Study Designs: Robust and Efficient Estimation.
- Callaway, B., Sant’Anna, P. H. C., 2021. Difference-in-Differences with multiple time periods. *J. Econom.* 225 (2), 200–230.
- de Chaisemartin, 2010. A note on instrumented difference in differences. Unpublished Manuscript.
- de Chaisemartin, C., D’Haultfœuille, X., 2018. Fuzzy Differences-in-Differences. *Rev. Econ. Stud.* 85 (2 (303)), 999–1028.
- de Chaisemartin, C., D’Haultfœuille, X., 2020. Two-Way Fixed Effects Estimators with Heterogeneous Treatment Effects. *Am. Econ. Rev.* 110 (9), 2964–2996.
- Duflo, E., 2001. Schooling and Labor Market Consequences of School Construction in Indonesia: Evidence from an Unusual Policy Experiment. *Am. Econ. Rev.* 91 (4), 795–813.
- Field, E., 2007. Entitled to Work: Urban Property Rights and Labor Supply in Peru. *Q. J. Econ.* 122 (4), 1561–1602.
- Goodman-Bacon, A., 2021. Difference-in-differences with variation in treatment timing. *J. Econom.* 225 (2), 254–277.
- Hudson, S., Hull, P., Liebersohn, J., 2017. Interpreting Instrumented Difference-in-Differences. Available at <http://www.mit.edu/~liebers/DDIV.pdf>.
- Imai, K., Kim, I. S., 2021. On the Use of Two-Way Fixed Effects Regression Models for Causal Inference with Panel Data. *Polit. Anal.* 29 (3), 405–415.

- Imbens, G. W., Angrist, J. D., 1994. Identification and Estimation of Local Average Treatment Effects. *Econometrica*. 62 (2), 467–475.
- Johnson, R. C., Jackson, C. K., 2019. Reducing Inequality through Dynamic Complementarity: Evidence from Head Start and Public School Spending. *Am. Econ. J. Econ. Policy*. 11 (4), 310–349.
- Kitagawa, E. M., 1955. Components of a Difference Between Two Rates. *J. Am. Stat. Assoc.* 50 (272), 1168–1194.
- Lundborg, P., Nilsson, A., Rooth, D.-O., 2014. Parental education and offspring outcomes: Evidence from the Swedish compulsory school reform. *Am. Econ. J. Appl. Econ.* 6 (1), 253–278.
- Lundborg, P., Plug, E., Rasmussen, A. W., 2017. Can Women Have Children and a Career? IV Evidence from IVF Treatments. *Am. Econ. Rev.* 107 (6), 1611–1637.
- Meghir, C., Palme, M., Simeonova, E., 2018. Education and mortality: Evidence from a social experiment. *Am. Econ. J. Appl. Econ.* 10 (2), 234–256.
- Miller, A. R., Segal, C., 2012. Does temporary affirmative action produce persistent effects? A study of black and female employment in law enforcement. *Rev. Econ. Stat.* 94 (4), 1107–1125.
- Miller, A. R., Segal, C., 2019. Do female officers improve law enforcement quality? Effects on crime reporting and domestic violence. *Rev. Econ. Stud.* 86 (5), 2220–2247.
- Miyaji, S., 2024. Instrumented Difference-in-Differences with heterogeneous treatment effects. Unpublished Manuscript.
- Oaxaca, R., 1973. Male-Female Wage Differentials in Urban Labor Markets. *Int. Econ. Rev.* 14 (3), 693–709.
- Oreopoulos, P., 2006. Estimating Average and Local Average Treatment Effects of Education when Compulsory Schooling Laws Really Matter. *Am. Econ. Rev.* 96 (1), 152–175.
- Słoczyński, T., 2020. When Should We (Not) Interpret Linear IV Estimands as LATE?
- Sun, L., Abraham, S., 2021. Estimating dynamic treatment effects in event studies with heterogeneous treatment effects. *J. Econom.* 225 (2), 175–199.

A Proof of the theorem in section 2

Before we proceed the proof of Theorem 1, we provide Lemma 4 below. This lemma is shown by Goodman-Bacon (2021).

Lemma 4 (Lemma 1 in Goodman-Bacon (2021)). The sample covariance between a cohort and time specific variable z_{kt} and a double demeaning variable $\tilde{x}_{kt} = (x_{kt} - \bar{x}_k) - (\bar{x}_t - \bar{x})$ is equal to a sum over every pair of observations of the period-by-period products of differences between cohorts in z_{kt} and \tilde{x}_{kt} .

$$\begin{aligned} & \sum_k n_k \frac{1}{T} \sum_t z_{kt} [(x_{kt} - \bar{x}_k) - (\bar{x}_t - \bar{x})] \\ &= \sum_k \sum_{l>k} n_l n_k \frac{1}{T} \sum_t (z_{kt} - z_{lt}) [(x_{kt} - \bar{x}_k) - (x_{lt} - \bar{x}_l)] \end{aligned} \quad (23)$$

Proof. See the proof of Lemma 1 in Goodman-Bacon (2021). \square

A.1 Proof of Theorem 1

Proof. From the FWL theorem, the TWFEIV estimator $\hat{\beta}_{IV}$ is:

$$\hat{\beta}_{IV} = \frac{\frac{1}{NT} \sum_i \sum_t \tilde{Z}_{i,t} Y_{i,t}}{\frac{1}{NT} \sum_i \sum_t \tilde{Z}_{i,t} D_{i,t}} \quad (24)$$

where $\tilde{Z}_{i,t}$ is a double-demeaning variable.

First, we consider the numerator of (24). In the following, we use $k(i)$ to express that unit i belongs to cohort k . We define $\bar{R}_{k(i),t}$ to be the sample mean of the random variable $R_{i,t}$ in cohort k at time t , and define $\bar{R}_{k(i)}$ to the average of $\bar{R}_{k(i),t}$ over time:

$$\bar{R}_{k(i),t} = \frac{\sum_i R_{i,t} \mathbf{1}\{E_i = k\}}{\sum_i \mathbf{1}\{E_i = k\}} \quad \text{and} \quad \bar{R}_{k(i)} = \frac{1}{T} \sum_{t=1}^T \bar{R}_{k(i),t}.$$

For the numerator of (24), by adding and subtracting $(\bar{Z}_{k(i),t} - \bar{Z}_{k(i)})$, we obtain

$$\begin{aligned} & \frac{1}{NT} \sum_i \sum_t \tilde{Z}_{i,t} Y_{i,t} \\ &= \frac{1}{NT} \sum_i \sum_t Y_{i,t} [(Z_{i,t} - \bar{Z}_i) - (\bar{Z}_t - \bar{Z})] \\ &= \frac{1}{NT} \sum_i \sum_t Y_{i,t} \left[\underbrace{(Z_{i,t} - \bar{Z}_i) - (\bar{Z}_{k(i),t} - \bar{Z}_{k(i)})}_{=0} + (\bar{Z}_{k(i),t} - \bar{Z}_{k(i)}) - (\bar{Z}_t - \bar{Z}) \right] \\ &= \frac{1}{NT} \sum_i \sum_t Y_{i,t} [(\bar{Z}_{k(i),t} - \bar{Z}_{k(i)}) - (\bar{Z}_t - \bar{Z})] \\ &= \sum_k n_k \frac{1}{T} \sum_t \bar{Y}_{k,t} [(\bar{Z}_{k,t} - \bar{Z}_k) - (\bar{Z}_t - \bar{Z})], \end{aligned}$$

where the third equality follows from the fact that $Z_{i,t} = \bar{Z}_{k(i),t}$ and $\bar{Z}_i = \bar{Z}_{k(i)}$ because all the units in cohort k have the same assignment of the instrument. The fourth equality follows because the expression only depends on cohort k and time t .

To further develop the expression, we use Lemma 4:

$$\begin{aligned}
& \sum_k n_k \frac{1}{T} \sum_t \bar{Y}_{kt} [(\bar{Z}_{kt} - \bar{Z}_k) - (\bar{Z}_t - \bar{Z})] \\
&= \sum_k \sum_{l>k} n_l n_k \frac{1}{T} \sum_t (\bar{Y}_{kt} - \bar{Y}_{lt}) [(\bar{Z}_{kt} - \bar{Z}_k) - (\bar{Z}_{lt} - \bar{Z}_l)]. \tag{25}
\end{aligned}$$

Next, we consider all possible expressions of (25). When $e = U$, that is, cohort e is never exposed cohort, we have $\bar{Z}_{Ut} - \bar{Z}_U = 0$. From this observation, for the pair (k, U) , we have:

$$\begin{aligned}
& \frac{1}{T} \sum_t (\bar{Y}_{kt} - \bar{Y}_{Ut}) [(\bar{Z}_{kt} - \bar{Z}_k) - (\bar{Z}_{Ut} - \bar{Z}_U)] \\
&= -\frac{1}{T} \sum_{t<k} (\bar{Y}_{kt} - \bar{Y}_{Ut}) \bar{Z}_k + \frac{1}{T} \sum_{t \geq k} (\bar{Y}_{kt} - \bar{Y}_{Ut}) (1 - \bar{Z}_k) \\
&= \left[(\bar{Y}_{kt}^{POST(k)} - \bar{Y}_{kt}^{PRE(k)}) - (\bar{Y}_{Ut}^{POST(k)} - \bar{Y}_{Ut}^{PRE(k)}) \right] \bar{Z}_k (1 - \bar{Z}_k).
\end{aligned}$$

By the similar argument, for the pair (k, l) where $k < l < T$, we obtain

$$\begin{aligned}
&= -\frac{1}{T} \sum_{t<k} (\bar{Y}_{kt} - \bar{Y}_{lt}) (\bar{Z}_k - \bar{Z}_l) + \frac{1}{T} \sum_{t \in [k,l]} (\bar{Y}_{kt} - \bar{Y}_{lt}) (1 - \bar{Z}_k + \bar{Z}_l) - \frac{1}{T} \sum_{t \geq l} (\bar{Y}_{kt} - \bar{Y}_{lt}) (\bar{Z}_k - \bar{Z}_l) \\
&= -\left[(\bar{Y}_{kt}^{PRE(k)} - \bar{Y}_{lt}^{PRE(k)}) \right] (\bar{Z}_k - \bar{Z}_l) (1 - \bar{Z}_k) + \left[(\bar{Y}_{kt}^{MID(k,l)} - \bar{Y}_{lt}^{MID(k,l)}) \right] (\bar{Z}_k - \bar{Z}_l) (1 - \bar{Z}_k + \bar{Z}_l) \\
&\quad - \left[(\bar{Y}_{kt}^{POST(l)} - \bar{Y}_{lt}^{POST(l)}) \right] (\bar{Z}_k - \bar{Z}_l) \bar{Z}_l \\
&= \left[(\bar{Y}_{kt}^{MID(k,l)} - \bar{Y}_{kt}^{PRE(k)}) - (\bar{Y}_{lt}^{MID(k,l)} - \bar{Y}_{lt}^{PRE(k)}) \right] (\bar{Z}_k - \bar{Z}_l) (1 - \bar{Z}_k) \\
&\quad + \left[(\bar{Y}_{lt}^{POST(l)} - \bar{Y}_{lt}^{MID(k,l)}) - (\bar{Y}_{kt}^{POST(l)} - \bar{Y}_{kt}^{MID(k,l)}) \right] (\bar{Z}_k - \bar{Z}_l) \bar{Z}_l.
\end{aligned}$$

To sum up, for the numerator of (24), we have

$$\begin{aligned}
& \frac{1}{NT} \sum_i \sum_t \tilde{Z}_{it} Y_{it} \\
&= \left[\sum_{k \neq U} (n_k + n_u)^2 n_{kU} (1 - n_{kU}) \bar{Z}_k (1 - \bar{Z}_k) \hat{\beta}_{kU}^{2 \times 2} \right. \\
&\quad + \sum_{k \neq U} \sum_{l>k} [(n_k + n_l) (1 - \bar{Z}_l)]^2 n_{kl} (1 - n_{kl}) \left(\frac{\bar{Z}_k - \bar{Z}_l}{1 - \bar{Z}_l} \right) \left(\frac{1 - \bar{Z}_k}{1 - \bar{Z}_l} \right) \hat{\beta}_{kl}^{2 \times 2, k} \\
&\quad \left. + ((n_k + n_l) \bar{Z}_k)^2 n_{kl} (1 - n_{kl}) \left(\frac{\bar{Z}_l}{\bar{Z}_k} \right) \left(\frac{\bar{Z}_k - \bar{Z}_l}{\bar{Z}_k} \right) \hat{\beta}_{kl}^{2 \times 2, l} \right] \\
&= \left[\sum_{k \neq U} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{\beta}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l>k} [(n_k + n_l) (1 - \bar{Z}_l)]^2 \hat{V}_{kl}^{Z,k} \hat{\beta}_{kl}^{2 \times 2, k} + ((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{\beta}_{kl}^{2 \times 2, l} \right]. \tag{26}
\end{aligned}$$

Next, we consider the denominator of (24). We note that the structure of the denominator is completely same as the one of the numerator in (24). Therefore, by the completely same

calculations, we obtain

$$\begin{aligned}
& \frac{1}{NT} \sum_i \sum_t \tilde{Z}_{it} D_{it} \\
& \equiv \hat{C}^{D,Z} \\
& = \left[\sum_{k \neq U} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} [((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{D}_{kl}^{2 \times 2,k} + ((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2,l}] \right].
\end{aligned} \tag{27}$$

Combining (26) with (27), we obtain

$$\begin{aligned}
& \hat{\beta}_{IV} \\
& = \frac{\left[\sum_{k \neq U} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{\beta}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} [((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{\beta}_{kl}^{2 \times 2,k} + ((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{\beta}_{kl}^{2 \times 2,l}] \right]}{\left[\sum_{k \neq U} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} [((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{D}_{kl}^{2 \times 2,k} + ((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2,l}] \right]} \\
& = \frac{\left[\sum_{k \neq U} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{\beta}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} [((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{\beta}_{kl}^{2 \times 2,k} + ((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{\beta}_{kl}^{2 \times 2,l}] \right]}{\hat{C}^{D,Z}} \\
& = \left[\sum_{k \neq U} \hat{w}_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} \hat{w}_{IV,kl}^k \hat{\beta}_{IV,kl}^{2 \times 2,k} + \hat{w}_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2,l} \right],
\end{aligned}$$

where the weights are:

$$\begin{aligned}
\hat{w}_{IV,kU} &= \frac{(n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2}}{\hat{C}^{D,Z}}, \\
\hat{w}_{IV,kl}^k &= \frac{((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{D}_{kl}^{2 \times 2,k}}{\hat{C}^{D,Z}}, \\
\hat{w}_{IV,kl}^l &= \frac{((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2,l}}{\hat{C}^{D,Z}}.
\end{aligned}$$

We note that $\hat{C}^{D,Z}$ is the sum of the numerator in each weight as one can see in (27). This implies that the weights sum to one:

$$\sum_{k \neq U} w_{IV,kU} + \sum_{k \neq U} \sum_{l > k} [w_{IV,kl}^k + w_{IV,kl}^l] = 1.$$

Completing the proof. □

B Proofs of the theorem and lemma in section 4.

In this section, we first prove Lemma 1 as a preparation.

B.1 Proof of Lemma 1

Proof. As one can see in the proof of Theorem 1, we have the following expression for the numerator of the TWFEIV estimand:

$$\hat{C}^{D,Z} = \sum_{k \neq U} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} \left[((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{D}_{kl}^{2 \times 2,k} + ((n_k + n_l)\bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2,l} \right].$$

We fix T and consider $N \rightarrow \infty$. We first derive the probability limit of $(n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2}$. By definition, we can rewrite $(n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2}$ as follows:

$$(n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2} = n_k n_u \bar{Z}_k (1 - \bar{Z}_k) \cdot \left[\left(\bar{D}_k^{POST(k)} - \bar{D}_k^{PRE(k)} \right) - \left(\bar{D}_U^{POST(k)} - \bar{D}_U^{PRE(k)} \right) \right].$$

By the law of large number (LLN), as $N \rightarrow \infty$, we obtain

$$\begin{aligned} & \left(\bar{D}_k^{POST(k)} - \bar{D}_k^{PRE(k)} \right) - \left(\bar{D}_U^{POST(k)} - \bar{D}_U^{PRE(k)} \right) \\ & \xrightarrow{p} \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[D_{it} | E_i = k] - E[D_{it} | E_i = U] \right] - \frac{1}{k - 1} \left[\sum_{t=1}^{k-1} E[D_{it} | E_i = k] - E[D_{it} | E_i = U] \right] \\ & = \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[D_{it}^k - D_{it}^\infty | E_i = k] \right] \\ & + \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[D_{it}^\infty | E_i = k] - \sum_{t=k}^T E[D_{it}^\infty | E_i = U] \right] \\ & - \frac{1}{k - 1} \left[\sum_{t=1}^{k-1} E[D_{it}^\infty | E_i = k] - E[D_{it}^\infty | E_i = U] \right] \\ & = \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[D_{it}^k - D_{it}^\infty | E_i = k] \right] \\ & = CAET_k^1(POST(k)). \end{aligned} \tag{28}$$

The second equality follows from the simple algebra and Assumption 5 (No anticipation for the first stage). The third equality follows from Assumption 6 (Parallel trend assumption in the treatment).

Next, we consider the probability limit of $n_k n_u \bar{Z}_k (1 - \bar{Z}_k)$. By the LLN and the Slutsky's theorem, we obtain

$$n_k n_u \bar{Z}_k (1 - \bar{Z}_k) \xrightarrow{p} Pr(E_i = k) Pr(E_i = U) \frac{T - (k - 1)}{T} \frac{k - 1}{T}. \tag{29}$$

Combining the result (28) with (29), by the Slutsky's theorem, we have

$$\begin{aligned} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2} & \xrightarrow{p} Pr(E_i = k) Pr(E_i = U) \frac{T - (k - 1)}{T} \frac{k - 1}{T} CAET_{k,t}^1(POST(k)) \\ & = w_{kU} CAET_k^1(POST(k)). \end{aligned} \tag{30}$$

where the weight w_{kU} is:

$$w_{kU} = Pr(E_i = k) Pr(E_i = U) \frac{T - (k - 1)}{T} \frac{k - 1}{T}.$$

By the completely same calculations, we also have

$$((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{D}_{kl}^{2 \times 2, k} \xrightarrow{p} w_{kl}^k CAET_k^1(MID(k, l)). \quad (31)$$

where the weight w_{kl}^k is:

$$w_{kl}^k = Pr(E_i = k)Pr(E_i = l) \frac{k-1}{T} \frac{l-k}{T}.$$

Next, we consider the probability limit of $((n_k + n_l)\bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2, l}$. By definition, we have

$$((n_k + n_l)\bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2, l} = n_k n_l (\bar{Z}_k - \bar{Z}_l) \bar{Z}_l \cdot \left[\left(\bar{D}_l^{POST(l)} - \bar{D}_l^{MID(k,l)} \right) - \left(\bar{D}_k^{POST(l)} - \bar{D}_k^{MID(k,l)} \right) \right].$$

By the law of large number (LLN), as $N \rightarrow \infty$, we have

$$\begin{aligned} & \left(\bar{D}_l^{POST(l)} - \bar{D}_l^{MID(k,l)} \right) - \left(\bar{D}_k^{POST(l)} - \bar{D}_k^{MID(k,l)} \right) \\ & \xrightarrow{p} \frac{1}{T - (l-1)} \left[\sum_{t=l}^T E[D_{it} | E_i = l] - E[D_{it} | E_i = k] \right] - \frac{1}{l-k} \left[\sum_{t=k}^{l-1} E[D_{it} | E_i = l] - E[D_{it} | E_i = k] \right] \\ & = \frac{1}{T - (l-1)} \left[\sum_{t=l}^T E[D_{it}^l - D_{it}^\infty | E_i = l] \right] - \frac{1}{T - (l-1)} \left[\sum_{t=l}^T E[D_{it}^k - D_{it}^\infty | E_i = k] \right] \\ & + \frac{1}{T - (l-1)} \left[\sum_{t=l}^T E[D_{it}^\infty | E_i = l] - \sum_{t=l}^T E[D_{it}^\infty | E_i = k] \right] \\ & + \frac{1}{l-k} \left[\sum_{t=k}^{l-1} E[D_{it}^k | E_i = k] - E[D_{it}^\infty | E_i = k] \right] \\ & + \frac{1}{l-k} \left[\sum_{t=k}^{l-1} E[D_{it}^\infty | E_i = k] - E[D_{it}^\infty | E_i = l] \right] \\ & = \frac{1}{T - (l-1)} \left[\sum_{t=l}^T E[D_{it}^l - D_{it}^\infty | E_i = l] \right] - \frac{1}{T - (l-1)} \left[\sum_{t=l}^T E[D_{it}^k - D_{it}^\infty | E_i = k] \right] \\ & + \frac{1}{l-k} \left[\sum_{t=k}^{l-1} E[D_{it}^k | E_i = k] - E[D_{it}^\infty | E_i = k] \right] \\ & = CAET_l^1(POST(l)) - CAET_k^1(POST(l)) + CAET_k^1(MID(k, l)) \end{aligned} \quad (32)$$

The second equality follows from the simple algebra and Assumption 5. The third equality follows from Assumption.

Note that the LLN and the Slutsky's theorem implies

$$n_k n_l (\bar{Z}_k - \bar{Z}_l) \bar{Z}_l \xrightarrow{p} Pr(E_i = k)Pr(E_i = l) \frac{T - (l-1)}{T} \frac{l-k}{T}. \quad (33)$$

From the result of (32) with (33), we obtain

$$\begin{aligned} & ((n_k + n_l)\bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2, l} \\ & \xrightarrow{p} w_{kl}^l \left[CAET_l^1(POST(l)) - CAET_k^1(POST(l)) + CAET_k^1(MID(k, l)) \right]. \end{aligned} \quad (34)$$

where the weight w_{kl}^l is:

$$w_{kl}^l = Pr(E_i = k)Pr(E_i = l) \frac{T - (l - 1)l - k}{T} \frac{l - k}{T}.$$

To sum up, by combining (30),(31) with (34), we obtain

$$\begin{aligned} \hat{C}^{D,Z} &= \sum_{k \neq U} (n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} \left[((n_k + n_l)(1 - \bar{Z}_l))^2 \hat{V}_{kl}^{Z,k} \hat{D}_{kl}^{2 \times 2,k} + ((n_k + n_l)\bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2,l} \right] \\ &\xrightarrow{p} WCAET - \Delta CAET^1. \end{aligned}$$

where we define

$$\begin{aligned} WCAET &\equiv \sum_{k \neq U} w_{kU} CAET_k^1(POST(k)) + \sum_{k \neq U} \sum_{l > k} w_{kl}^k CAET_k^1(MID(k, l)) + w_{kl}^l CAET_l^1(POST(l)), \\ \Delta CAET^1 &\equiv w_{kl}^l [CAET_k^1(POST(l)) - CAET_k^1(MID(k, l))]. \end{aligned}$$

and the weights are:

$$w_{kU} = Pr(E_i = k)Pr(E_i = U) \frac{T - (k - 1)k - 1}{T} \frac{k - 1}{T}, \quad (35)$$

$$w_{kl}^k = Pr(E_i = k)Pr(E_i = l) \frac{k - 1}{T} \frac{l - k}{T}, \quad (36)$$

$$w_{kl}^l = Pr(E_i = k)Pr(E_i = l) \frac{T - (l - 1)l - k}{T} \frac{l - k}{T}. \quad (37)$$

Completing the proof. □

B.2 Preparation for the proof of Theorem 2.

Before we present the proof of Theorem 2, we show the following lemma.

Lemma 5. Suppose Assumptions 1-7 hold. If the treatment is binary, for all $k, l \in \{1, \dots, T\} (k \leq l)$, we have

$$\begin{aligned} \left[\sum_{t=k}^l E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] &= \sum_{t=k}^l E[D_{i,t}^k - D_{i,t}^\infty | E_i = k] \cdot E[Y_{i,t}(1) - Y_{i,t}(0) | E_i = k, CM_{k,t}] \\ &\equiv \sum_{t=k}^l CAET_{k,t}^1 \cdot CLATT_{k,t}. \end{aligned}$$

Proof. Because we assume a binary treatment, we have

$$\begin{aligned} \left[\sum_{t=k}^l E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] &= \left[\sum_{t=k}^l E[(D_{i,t}^k - D_{i,t}^\infty) \cdot (Y_{i,t}(1) - Y_{i,t}(0)) | E_i = k] \right] \\ &= \sum_{t=k}^l E[D_{i,t}^k - D_{i,t}^\infty | E_i = k] \cdot E[Y_{i,t}(1) - Y_{i,t}(0) | E_i = k, CM_{k,t}], \end{aligned}$$

where the first equality holds from Assumptions 1-3 and the second equality holds from Assumption 4. □

B.3 Proof of Theorem 2.

Proof. We fix T and consider $N \rightarrow \infty$. We first derive the probability limit of $\hat{w}_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2}$. We note that $\hat{w}_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2}$ is written as follows:

$$\hat{w}_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2} = \frac{(n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2}}{\hat{C}^{D,Z}} \cdot \frac{\hat{\beta}_{kU}^{2 \times 2}}{\hat{D}_{kU}^{2 \times 2}}.$$

We first consider the probability limit of $\frac{\hat{\beta}_{kU}^{2 \times 2}}{\hat{D}_{kU}^{2 \times 2}}$. Recall that we have already derived the probability limit of the denominator in the proof of Lemma 1:

$$\hat{D}_{kU}^{2 \times 2} \xrightarrow{p} \frac{1}{T - (k - 1)} \sum_{t=k}^T CAET_{k,t}^1. \quad (38)$$

We consider the numerator $\hat{\beta}_{kU}^{2 \times 2}$. By the law of large number (LLN), as $N \rightarrow \infty$, we obtain

$$\begin{aligned} & \left(\bar{Y}_k^{POST(k)} - \bar{Y}_k^{PRE(k)} \right) - \left(\bar{Y}_U^{POST(k)} - \bar{Y}_U^{PRE(k)} \right) \\ & \xrightarrow{p} \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[Y_{i,t} | E_i = k] - E[Y_{i,t} | E_i = U] \right] - \frac{1}{k - 1} \left[\sum_{t=1}^{k-1} E[Y_{i,t} | E_i = k] - E[Y_{i,t} | E_i = U] \right] \\ & = \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\ & + \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[Y_{i,t}(D_{i,t}^\infty) | E_i = k] - \sum_{t=k}^T E[Y_{i,t}(D_{i,t}^\infty) | E_i = U] \right] \\ & - \frac{1}{k - 1} \left[\sum_{t=1}^{k-1} E[Y_{i,t}(D_{i,t}^\infty) | E_i = k] - E[Y_{i,t}(D_{i,t}^\infty) | E_i = U] \right] \\ & = \frac{1}{T - (k - 1)} \left[\sum_{t=k}^T E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\ & = \frac{1}{T - (k - 1)} \sum_{t=k}^T CAET_{k,t}^1 \cdot CLATT_{k,t}. \quad (39) \end{aligned}$$

The first equality follows from the simple manipulation, Assumption 3 (Exclusion restriction in multiple time periods) and Assumption 5. The second equality follows from Assumption 7 (Parallel trend assumption in the outcome). The final equality follows from Lemma 5.

Combining the result (38) with (39), we obtain

$$\begin{aligned} & \frac{\hat{\beta}_{kU}^{2 \times 2}}{\hat{D}_{kU}^{2 \times 2}} \xrightarrow{p} \sum_{t=k}^T \frac{CAET_{k,t}^1}{\sum_{t=k}^T CAET_{k,t}^1} \cdot CLATT_{k,t} \\ & = CLATT_k^{CM}(POST(k)). \quad (40) \end{aligned}$$

Next, we consider the probability limit of $\frac{(n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2}}{\hat{C}^{D,Z}}$. By the LLN and the Slutsky's theorem, we obtain

$$\frac{(n_k + n_u)^2 \hat{V}_{kU}^Z \hat{D}_{kU}^{2 \times 2}}{\hat{C}^{D,Z}} \xrightarrow{p} \frac{Pr(E_i = k) Pr(E_i = U) T - (k - 1) k - 1}{C^{D,Z} T} CAET_k^1(POST(k)). \quad (41)$$

Here $C^{D,Z}$ is the probability limit of $\hat{C}^{D,Z}$ and its specific expression is already derived in Lemma 1.

Combining the result (40) with (41), by the Slutsky's theorem, we have

$$\hat{w}_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2} \xrightarrow{p} w_{IV,kU} CLATT_k^{CM}(POST(k)). \quad (42)$$

where the weight $w_{IV,kU}$ is:

$$w_{IV,kU} = \frac{Pr(E_i = k)Pr(E_i = U)}{C^{D,Z}} \frac{T - (k - 1)}{T} \frac{k - 1}{T} \cdot CAET_k^1(POST(k)).$$

By the completely same argument, we also obtain

$$\hat{w}_{IV,kl}^k \hat{\beta}_{IV,kl}^{2 \times 2,k} \xrightarrow{p} w_{IV,kl}^k CLATT_k^{CM}(MID(k, l)). \quad (43)$$

where the weight $w_{IV,kl}^k$ is:

$$w_{IV,kl}^k = \frac{Pr(E_i = k)Pr(E_i = l)}{C^{D,Z}} \frac{k - 1}{T} \frac{l - k}{T} \cdot CAET_k^1(MID(k, l)).$$

Next, we derive the probability limit of $\hat{w}_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2,l}$. Recall that $\hat{w}_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2,l}$ is:

$$\hat{w}_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2,l} = \frac{((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2,l}}{\hat{C}^{D,Z}} \cdot \frac{\hat{\beta}_{kl}^{2 \times 2,l}}{\hat{D}_{kl}^{2 \times 2,l}}.$$

First note that in the proof of Lemma 1, we have already derived the probability limit of $\hat{D}_{kl}^{2 \times 2,l}$:

$$\begin{aligned} \hat{D}_{kl}^{2 \times 2,l} &\xrightarrow{p} CAET_l^1(POST(l)) - [CAET_k^1(POST(l)) - CAET_k^1(MID(k, l))] \\ &\equiv D_{kl}^{2 \times 2,l}. \end{aligned} \quad (44)$$

Here, to ease the notation, we define $D_{kl}^{2 \times 2,l}$ to be the probability limit of $\hat{D}_{kl}^{2 \times 2,l}$.

Next, we consider the probability limit of $\hat{\beta}_{kl}^{2 \times 2,l}$.

By the law of large number (LLN), as $N \rightarrow \infty$, we have

$$\begin{aligned}
\hat{\beta}_{kl}^{2 \times 2, l} &= \left(\bar{Y}_l^{POST(l)} - \bar{Y}_l^{MID(k, l)} \right) - \left(\bar{Y}_k^{POST(l)} - \bar{Y}_k^{MID(k, l)} \right) \\
&\xrightarrow{p} \frac{1}{T - (l - 1)} \left[\sum_{t=l}^T E[Y_{i,t} | E_i = l] - E[Y_{i,t} | E_i = k] \right] - \frac{1}{l - k} \left[\sum_{t=k}^{l-1} E[Y_{i,t} | E_i = l] - E[Y_{i,t} | E_i = k] \right] \\
&= \frac{1}{T - (l - 1)} \left[\sum_{t=l}^T E[Y_{i,t}(D_{i,t}^l) - Y_{i,t}(D_{i,t}^\infty) | E_i = l] \right] - \frac{1}{T - (l - 1)} \left[\sum_{t=l}^T E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\
&+ \frac{1}{T - (l - 1)} \left[\sum_{t=l}^T E[Y_{i,t}(D_{i,t}^\infty) | E_i = l] - \sum_{t=l}^T E[Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\
&+ \frac{1}{l - k} \left[\sum_{t=k}^{l-1} E[Y_{i,t}(D_{i,t}^k) | E_i = k] - E[Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\
&+ \frac{1}{l - k} \left[\sum_{t=k}^{l-1} E[Y_{i,t}(D_{i,t}^\infty) | E_i = k] - E[Y_{i,t}(D_{i,t}^\infty) | E_i = l] \right] \\
&= \frac{1}{T - (l - 1)} \left[\sum_{t=l}^T E[Y_{i,t}(D_{i,t}^l) - Y_{i,t}(D_{i,t}^\infty) | E_i = l] \right] - \frac{1}{T - (l - 1)} \left[\sum_{t=l}^T E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\
&+ \frac{1}{l - k} \left[\sum_{t=k}^{l-1} E[Y_{i,t}(D_{i,t}^k) | E_i = k] - E[Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\
&= \frac{1}{T - (l - 1)} \sum_{t=l}^T CAET_{l,t}^1 \cdot CLATT_{l,t} \\
&- \left[\frac{1}{T - (l - 1)} \sum_{t=l}^T CAET_{k,t}^1 \cdot CLATT_{k,t} - \frac{1}{l - k} \sum_{t=k}^{l-1} CAET_{k,t}^1 \cdot CLATT_{k,t} \right]. \tag{45}
\end{aligned}$$

The first equality follows from the simple algebra, Assumption 3, and Assumption 5. The second equality follows from Assumption 7. The final equality follows from Lemma 5.

Note that the LLN and the Slutsky's theorem yields

$$\frac{((n_k + n_l) \bar{Z}_k)^2 \hat{V}_{kl}^{Z,l} \hat{D}_{kl}^{2 \times 2, l}}{\hat{C}^{D,Z}} \xrightarrow{p} \frac{Pr(E_i = k) Pr(E_i = l) T - (l - 1) l - k}{C^{D,Z}} \frac{T - (l - 1) l - k}{T} \frac{l - k}{T} \cdot D_{kl}^{2 \times 2, l} \tag{46}$$

From the results of (44) and (45) with (46), we obtain

$$\begin{aligned}
\hat{w}_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2, l} &\xrightarrow{p} \sigma_{IV,kl}^l \cdot \left\{ \frac{1}{T - (l - 1)} \sum_{t=l}^T CAET_{l,t}^1 \cdot CLATT_{l,t} \right. \\
&- \left. \left[\frac{1}{T - (l - 1)} \sum_{t=l}^T CAET_{k,t}^1 \cdot CLATT_{k,t} - \frac{1}{l - k} \sum_{t=k}^{l-1} CAET_{k,t}^1 \cdot CLATT_{k,t} \right] \right\} \\
&= \sigma_{IV,kl}^l \cdot \left\{ CLATT_l(POST(l)) \right. \\
&- \left. [CLATT_k(POST(l)) - CLATT_k(MID(k, l))] \right\}. \tag{47}
\end{aligned}$$

where the weight $\sigma_{IV,kl}^l$ is:

$$\sigma_{IV,kl}^l = \frac{Pr(E_i = k)Pr(E_i = l) T - (l-1)l - k}{C^{D,Z}} \frac{T - (l-1)l - k}{T} \frac{l - k}{T}. \quad (48)$$

To sum up, by combining (42),(43) with (47), we obtain

$$\begin{aligned} \hat{\beta}_{IV} &= \left[\sum_{k \neq U} w_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^k \hat{\beta}_{IV,kl}^{2 \times 2,k} + w_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2,l} \right] \\ &\xrightarrow{p} WCLATT - \Delta CLATT. \end{aligned}$$

where we define

$$\begin{aligned} WCLATT &\equiv \sum_{k \neq U} w_{IV,kU} CLATT_k^{CM}(POST(k)) + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^k CLATT_k^{CM}(MID(k,l)) \\ &\quad + \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot CLATT_l(POST(l)), \\ \Delta CLATT &\equiv \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot [CLATT_k(POST(l)) - CLATT_k(MID(k,l))]. \end{aligned}$$

and the weights are:

$$w_{IV,kU} = \frac{Pr(E_i = k)Pr(E_i = U) T - (k-1)k - 1}{C^{D,Z}} \frac{T - (k-1)k - 1}{T} \cdot CAET_k^1(POST(k)), \quad (49)$$

$$w_{IV,kl}^k = \frac{Pr(E_i = k)Pr(E_i = l) k - 1 l - k}{C^{D,Z}} \frac{k - 1 l - k}{T} \cdot CAET_k^1(MID(k,l)), \quad (50)$$

$$\sigma_{IV,kl}^l = \frac{Pr(E_i = k)Pr(E_i = l) T - (l-1)l - k}{C^{D,Z}} \frac{T - (l-1)l - k}{T} \frac{l - k}{T}. \quad (51)$$

Completing the proof. \square

B.4 Proof of Lemma 2

Proof. We first simplify $CLATT_k^{CM}(W)$ and $D_{kl}^{2 \times 2,l}$ (defined in (44)) under Assumption 9. If we assume Assumption 9, $CLATT_k^{CM}(W)$ is:

$$\begin{aligned} CLATT_k^{CM}(W) &= \sum_{t \in W} \frac{Pr(CM_{k,t} | E_i = k)}{\sum_{t \in W} Pr(CM_{k,t} | E_i = k)} CLATT_{k,t} \\ &= \frac{1}{T_W} \sum_{t \in W} CLATT_{k,t} \\ &\equiv CLATT_k^{eq}(W). \end{aligned}$$

In addition, $D_{kl}^{2 \times 2,l}$ is:

$$\begin{aligned} D_{kl}^{2 \times 2,l} &= CAET_l^1(POST(l)) - [CAET_k^1(POST(l)) - CAET_k^1(MID(k,l))] \\ &= CAET_l^1. \end{aligned}$$

because we have $CAET_k^1(W) = CAET_k^1$.

We then rewrite the probability limit of $\hat{w}_{IV,kU}\hat{\beta}_{IV,kU}^{2 \times 2}$, $\hat{w}_{IV,kl}^k\hat{\beta}_{IV,kl}^{2 \times 2,k}$ and $\hat{w}_{IV,kl}^l\hat{\beta}_{IV,kl}^{2 \times 2,l}$ respectively. First, the probability limit of $\hat{w}_{IV,kU}\hat{\beta}_{IV,kU}^{2 \times 2}$ and $\hat{w}_{IV,kl}^k\hat{\beta}_{IV,kl}^{2 \times 2,k}$ is simplified to:

$$\begin{aligned}\hat{w}_{IV,kU}\hat{\beta}_{IV,kU}^{2 \times 2} &\xrightarrow{p} w_{IV,kU}CLATT_k^{CM}(POST(k)) \\ &= w_{IV,kU}CLATT_k^{eq}(POST(k)).\end{aligned}\quad (52)$$

$$\begin{aligned}\hat{w}_{IV,kl}^k\hat{\beta}_{IV,kl}^{2 \times 2,k} &\xrightarrow{p} w_{IV,kl}^kCLATT_k^{CM}(MID(k, l)) \\ &= w_{IV,kl}^kCLATT_k^{eq}(MID(k, l)).\end{aligned}\quad (53)$$

where the weights $w_{IV,kU}$, $w_{IV,kl}^k$ are:

$$w_{IV,kU} = \frac{Pr(E_i = k)Pr(E_i = U)}{C^{D,Z}} \frac{T - (k - 1)}{T} \frac{k - 1}{T} \cdot CAET_k^1, \quad (54)$$

$$w_{IV,kl}^k = \frac{Pr(E_i = k)Pr(E_i = l)}{C^{D,Z}} \frac{k - 1}{T} \frac{l - k}{T} \cdot CAET_k^1. \quad (55)$$

Next, we reconsider the probability limit of $\hat{w}_{IV,kl}^l\hat{\beta}_{IV,kl}^{2 \times 2,l}$.

First, we note that the probability limit of $\hat{w}_{IV,kl}^l$ is simplified to:

$$\begin{aligned}\hat{w}_{IV,kl}^l &= \frac{((n_k + n_l)\bar{Z}_k)^2\hat{V}_{kl}^{Z,l}\hat{D}_{kl}^{2 \times 2,l}}{\hat{C}^{D,Z}} \xrightarrow{p} \frac{Pr(E_i = k)Pr(E_i = l)}{C^{D,Z}} \frac{T - (l - 1)}{T} \frac{l - k}{T} \cdot D_{kl}^{2 \times 2,l} \\ &= \frac{Pr(E_i = k)Pr(E_i = l)}{C^{D,Z}} \frac{T - (l - 1)}{T} \frac{l - k}{T} \cdot CAET_l^1.\end{aligned}\quad (56)$$

Here the second equality follows from $D_{kl}^{2 \times 2,l} = CAET_l^1$.

Second, the probability limit of $\hat{\beta}_{IV,kl}^{2 \times 2,l}$ simply reduces to:

$$\begin{aligned}\hat{\beta}_{IV,kl}^{2 \times 2,l} &= \frac{\hat{\beta}_{kl}^{2 \times 2,l}}{\hat{D}_{kl}^{2 \times 2,l}} \\ &\xrightarrow{p} \frac{1}{CAET_l^1} \cdot \frac{1}{T - (l - 1)} \sum_{t=l}^T CAET_l^1 \cdot CLATT_{l,t} \\ &\quad - \frac{1}{CAET_l^1} \cdot \left[\frac{1}{T - (l - 1)} \sum_{t=l}^T CAET_k^1 \cdot CLATT_{k,t} - \frac{1}{l - k} \sum_{t=k}^{l-1} CAET_k^1 \cdot CLATT_{k,t} \right] \\ &= CLATT_l^{eq}(POST(l)) \\ &\quad - \frac{1}{CAET_l^1} \cdot [CLATT_k(POST(l)) - CLATT_k(MID(k, l))].\end{aligned}\quad (57)$$

where we use (45) and $D_{kl}^{2 \times 2,l} = CAET_l^1$.

Combining the result (57) with (56), by the Slutsky's theorem, we have

$$\hat{w}_{IV,kl}^l\hat{\beta}_{IV,kl}^{2 \times 2,l} \xrightarrow{p} w_{IV,kl}^lCLATT_l^{eq}(POST(l)) - \sigma_{IV,kl}^l \cdot [CLATT_k(POST(l)) - CLATT_k(MID(k, l))]. \quad (58)$$

where the weight $w_{IV,kl}^l$ is:

$$w_{IV,kl}^l = \frac{Pr(E_i = k)Pr(E_i = l)}{C^{D,Z}} \frac{T - (l - 1)}{T} \frac{l - k}{T} \cdot CAET_l^1, \quad (59)$$

and $\sigma_{IV,kl}^l$ is already defined in (48).

Finally, from the result (52) and (53) with (58), we obtain

$$\begin{aligned}\hat{\beta}_{IV} &= \left[\sum_{k \neq U} \hat{w}_{IV,kU} \hat{\beta}_{IV,kU}^{2 \times 2} + \sum_{k \neq U} \sum_{l > k} \hat{w}_{IV,kl}^k \hat{\beta}_{IV,kl}^{2 \times 2,k} + \hat{w}_{IV,kl}^l \hat{\beta}_{IV,kl}^{2 \times 2,l} \right] \\ &\xrightarrow{p} WCLATT - \Delta CLATT.\end{aligned}$$

where we define:

$$\begin{aligned}WCLATT &\equiv \sum_{k \neq U} w_{IV,kU} CLATT_k^{eq}(POST(k)) + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^k CLATT_k^{eq}(MID(k, l)) \\ &\quad + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^l CLATT_l^{eq}(POST(l)), \\ \Delta CLATT &\equiv \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot [CLATT_k(POST(l)) - CLATT_k(MID(k, l))].\end{aligned}$$

Completing the proof. □

B.5 Proof of Lemma 3

Proof. First, we show $\Delta CLATT = 0$. Under Assumption 9 and Assumption 11, we have $CLATT_{e,t} = CLATT_e$. This implies:

$$\begin{aligned}\Delta CLATT &= \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot [CLATT_k(POST(l)) - CLATT_k(MID(k, l))] \\ &= 0.\end{aligned}$$

because we have $CLATT_k(POST(l)) - CLATT_k(MID(k, l)) = 0$.

Next, we consider $WCLATT$. Since we have $CLATT_k^{eq}(W) = CLATT_k$, we obtain:

$$\begin{aligned}WCLATT &= \sum_{k \neq U} w_{IV,kU} CLATT_k + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^k CLATT_k + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^l CLATT_l \\ &= \sum_{k \neq U} CLATT_k \left[w_{IV,kU} + \sum_{j=1}^{k-1} w_{IV,jk}^k + \sum_{j=k+1}^K w_{IV,kj}^k \right].\end{aligned}$$

Completing the proof. □

C Extensions in section 5

C.1 Non-binary, ordered treatment

This subsection considers a non binary, ordered treatment. We show Lemma 6 below that is analogous to Lemma 5 in a binary treatment. If we use Lemma 6 instead of Lemma 5 in the proof of Theorem 2 and Lemmas 2-3, we obtain the theorem and the lemmas which replace $CLATT_{e,k}$ with $CACRT_{e,k}$.

Lemma 6. Suppose Assumptions 1-7 hold. If treatment is a non-binary, ordered, for all $k, l \in \{1, \dots, T\} (k \leq l)$, we have

$$\begin{aligned} & \left[\sum_{t=k}^l E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\ &= \sum_{t=k}^l E[D_{i,t}^k - D_{i,t}^\infty | E_i = k] \cdot \sum_{j=1}^J w_{t,j}^k \cdot E[Y_{i,t}(j) - Y_{i,t}(j-1) | E_i = k, D_{i,t}^k \geq j > D_{i,t}^\infty] \\ &\equiv \sum_{t=k}^l CATT_{k,t}^1 \cdot CACRT_{k,t}. \end{aligned}$$

where the weight $w_{t,j}^k$ is:

$$w_{t,j}^k = \frac{\Pr(D_{i,t}^k \geq j > D_{i,t}^\infty | E_i = k)}{\sum_{j=1}^J \Pr(D_{i,t}^k \geq j > D_{i,t}^\infty | E_i = k)}.$$

Proof. By the similar argument in the proof of lemma 5, one can show that

$$\begin{aligned} & \left[\sum_{t=k}^l E[Y_{i,t}(D_{i,t}^k) - Y_{i,t}(D_{i,t}^\infty) | E_i = k] \right] \\ &= \left[\sum_{t=k}^l \sum_{j=1}^J \Pr(D_{i,t}^k \geq j > D_{i,t}^\infty | E_i = k) \cdot E[Y_{i,t}(j) - Y_{i,t}(j-1) | E_i = k, D_{i,t}^k \geq j > D_{i,t}^\infty] \right] \quad (60) \end{aligned}$$

$$\begin{aligned} & E[D_{i,t}^k - D_{i,t}^\infty | E_i = k] \\ &= \sum_{j=1}^J \Pr(D_{i,t}^k \geq j > D_{i,t}^\infty | E_i = k). \quad (61) \end{aligned}$$

Combining the result (60) with (61), we obtain the desired result. \square

C.2 Unbalanced panel case

In this section, we consider an unbalanced setting. We use the notation for a panel data setting, but the discussions and the results are the same if we consider an unbalanced repeated cross section setting.

Proof of Theorem 3

Let $N_{e,t}$ be the sample size for cohort e at time t and $N = \sum_e \sum_t N_{e,t}$ be the total number of observations. We consider the following two way fixed effects instrumental variable regression:

$$\begin{aligned} Y_{i,t} &= \mu_i + \delta_t + \alpha Z_{i,t} + \epsilon_{i,t}, \\ D_{i,t} &= \gamma_i + \zeta_t + \pi Z_{i,t} + \eta_{i,t}. \end{aligned}$$

We define $\hat{Z}_{i,t}$ to be the residuals from regression $Z_{i,t}$ on the time and individual fixed effects.

From the FWL theorem, the TWFEIV estimator $\hat{\beta}_{IV}$ is:

$$\begin{aligned}\hat{\beta}_{IV} &= \frac{\sum_i \sum_t \hat{Z}_{i,t} Y_{i,t}}{\sum_i \sum_t \hat{Z}_{i,t} D_{i,t}} \\ &= \frac{\sum_e \sum_t N_{e,t} \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} \hat{Z}_{e(i),t} Y_{e(i),t}}{\sum_e \sum_t N_{e,t} \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} \hat{Z}_{e(i),t} D_{e(i),t}} \\ &= \frac{\sum_e \sum_t N_{e,t} \hat{Z}_{e,t} \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} Y_{e(i),t}}{\sum_e \sum_t N_{e,t} \hat{Z}_{e,t} \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} D_{e(i),t}},\end{aligned}$$

where the third equality follows from the fact that $\hat{Z}_{i,t}$ only varies across cohort and time level.

We note that by the definition of $\hat{Z}_{e,t}$, we have

$$\sum_t N_{e,t} \hat{Z}_{e,t} = 0 \quad \text{for all } e \in \mathcal{S}(E_i), \quad (62)$$

$$\sum_e N_{e,t} \hat{Z}_{e,t} = 0 \quad \text{for all } t \in \{1, \dots, T\}. \quad (63)$$

To ease the notation, we define the sample mean for a random variable $R_{i,t}$ in cohort e at time t as follows:

$$R_{e,t} \equiv \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} R_{e(i),t}.$$

Here, we note that we can express $Y_{e,t}$ in the following:

$$\begin{aligned}Y_{e,t} &= \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} Y_{e(i),t} \\ &= \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} [Y_{e(i),t}(D_{i,t}^\infty) + Z_{e,t} \cdot (Y_{e(i),t}(D_{i,t}^e) - Y_{e(i),t}(D_{i,t}^\infty))] \\ &= Y_{e,t}(D_{i,t}^\infty) + Z_{e,t} \cdot (Y_{e,t}(D_{i,t}^e) - Y_{e,t}(D_{i,t}^\infty)).\end{aligned} \quad (64)$$

where the second equality follows from Assumptions 1-3 and Assumption 5.

First, we consider the probability limit of the numerator in the TWFEIV estimator. By using (62) and (63), we obtain

$$\begin{aligned}\sum_e \sum_t N_{e,t} \hat{Z}_{e,t} \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} Y_{e(i),t} &= \sum_e \sum_t N_{e,t} \hat{Z}_{e,t} Y_{e,t} \\ &= \sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [Y_{e,t} - Y_{e,1} - (Y_{1,t} - Y_{1,1})].\end{aligned} \quad (65)$$

To further develop the expression, we use (64):

$$\begin{aligned}Y_{e,t} - Y_{e,1} - (Y_{1,t} - Y_{1,1}) &= Y_{e,t}(D_{i,t}^\infty) - Y_{e,1}(D_{i,1}^\infty) - [Y_{1,t}(D_{i,t}^\infty) - Y_{1,1}(D_{i,1}^\infty)] \\ &\quad + Z_{e,t} \cdot (Y_{e,t}(D_{i,t}^e) - Y_{e,t}(D_{i,t}^\infty)) - Z_{e,1} \cdot (Y_{e,1}(D_{i,1}^e) - Y_{e,1}(D_{i,1}^\infty)) \\ &\quad - [Z_{1,t} \cdot (Y_{1,t}(D_{i,t}^1) - Y_{1,t}(D_{i,t}^\infty)) - Z_{1,1} \cdot (Y_{1,1}(D_{i,1}^e) - Y_{1,1}(D_{i,1}^\infty))] \quad (66)\end{aligned}$$

Substituting (66) into (65), we obtain:

$$\begin{aligned}
& \sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [Y_{e,t} - Y_{e,1} - (Y_{1,t} - Y_{1,1})] \\
&= \sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [Y_{e,t}(D_{i,t}^\infty) - Y_{e,1}(D_{i,1}^\infty) - [Y_{1,t}(D_{i,t}^\infty) - Y_{1,1}(D_{i,1}^\infty)]] \\
&+ \sum_e \sum_t N_{e,t} \hat{Z}_{e,t} Z_{e,t} \cdot (Y_{e,t}(D_{i,t}^e) - Y_{e,t}(D_{i,t}^\infty)), \tag{67}
\end{aligned}$$

where the second equality holds from (62) and (63).

From (67), as $N \rightarrow \infty$, we obtain

$$\begin{aligned}
& \sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [Y_{e,t} - Y_{e,1} - (Y_{1,t} - Y_{1,1})] \\
&\xrightarrow{p} \sum_e \sum_t E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot \left\{ E[Y_{e,t}(D_{i,t}^\infty)|E_i = e] - E[Y_{e,1}(D_{i,1}^\infty)|E_i = e] \right. \\
&\quad \left. - (E[Y_{1,t}(D_{i,t}^\infty)|E_i = 1] - E[Y_{1,1}(D_{i,1}^\infty)|E_i = 1]) \right\} \\
&+ \sum_e \sum_t E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot E[Z_{e,t} \cdot (Y_{e,t}(D_{i,t}^e) - Y_{e,t}(D_{i,t}^\infty))|E_i = e] \\
&= \sum_e \sum_{t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot E[(Y_{e,t}(D_{i,t}^e) - Y_{e,t}(D_{i,t}^\infty))|E_i = e] \\
&= \sum_e \sum_{t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CATT_{e,t}^1 \cdot CLATT_{e,t}, \tag{68}
\end{aligned}$$

where $n_{e,t}$ is population share and $E[\hat{Z}_{i,t}|E_i = e]$ in cohort e at time t . The first equality follows from Assumption 1 and Assumption 7. The second equality follows from Assumption 5.

Next, we consider the probability limit of the numerator. We note that the structure in the numerator is same as the one in the denominator. Therefore, by the same argument, we have:

$$\begin{aligned}
& \sum_e \sum_t N_{e,t} \hat{Z}_{e,t} \frac{1}{N_{e,t}} \sum_i^{N_{e,t}} D_{e(i),t} \\
&\xrightarrow{p} \sum_e \sum_{t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CATT_{e,t}^1. \tag{69}
\end{aligned}$$

Combining the result (69) with (68), we obtain

$$\begin{aligned}
& \hat{\beta}_{IV} \xrightarrow{p} \beta_{IV} \\
&= \frac{\sum_e \sum_{t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CATT_{e,t}^1 \cdot CLATT_{e,t}}{\sum_e \sum_{t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CATT_{e,t}^1} \\
&= \sum_e \sum_{t \geq e} w_{e,t} \cdot CLATT_{e,t}.
\end{aligned}$$

where the weight $w_{e,t}$ is:

$$w_{e,t} = \frac{E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CATT_{e,t}^1}{\sum_e \sum_{t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CATT_{e,t}^1}.$$

Completing the proof.

Supplementary of Theorem 3

We provide another representation of Theorem 3. We assume that there does not exist a never exposed cohort, that is, we have $\infty \notin \mathcal{S}(E_i)$, and define $l = \max\{E_i\}$ to be the last exposed cohort.

Lemma 7. Suppose Assumptions 1-7 hold. Assume a binary treatment and an unbalanced panel setting. If there does not exist a never exposed cohort, i.e., we have $\infty \notin \mathcal{S}(E_i)$, the population regression coefficient β_{IV} is:

$$\beta_{IV} = \sum_e \sum_{l-1 \geq t \geq e} w_{e,t}^1 \cdot CLATT_{e,t} + \sum_e \sum_{t \geq l} w_{e,t}^2 \cdot \Delta_{e,t}, \quad (70)$$

where $\Delta_{e,t}$ is:

$$\frac{CAET_{e,t}^1 \cdot CLATT_{e,t} - CAET_{l,t}^1 \cdot CLATT_{l,t}}{CAET_{e,t}^1 - CAET_{l,t}^1},$$

and the weights $w_{e,t}^1$ and $w_{e,t}^2$ are:

$$w_{e,t}^1 = \frac{E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CAET_{e,t}^1}{\sum_e \left(\sum_{l-1 \geq t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CAET_{e,t}^1 + \sum_{t \geq l} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot (CAET_{e,t}^1 - CAET_{l,t}^1) \right)}, \quad (71)$$

$$w_{e,t}^2 = \frac{E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot (CAET_{e,t}^1 - CAET_{l,t}^1)}{\sum_e \left(\sum_{l-1 \geq t \geq e} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot CAET_{e,t}^1 + \sum_{t \geq l} E[\hat{Z}_{i,t}|E_i = e] \cdot n_{e,t} \cdot (CAET_{e,t}^1 - CAET_{l,t}^1) \right)}. \quad (72)$$

We note that when there is no never exposed cohort, we can only identify each $CLATT_{e,t}$ before the time period $l = \max\{E_i\}$ for cohort $e \neq l$, exploiting the time trends of the unexposed treatment and outcome for cohort l . This implies that in equation (70), each $\Delta_{e,t}$ is the bias term occurring from the bad comparisons performed by TWFEIV regressions. In a given application, we can estimate $CLATT_{e,t}$, $CLATT_{l,t}$, and the associated weights $w_{e,t}^1$, $w_{e,t}^2$ by constructing the consistent estimators, using (73) and (74) below.

Proof. We consider the case where there is no never exposed cohort, i.e., we have $\infty \notin \mathcal{S}(E_i)$. In this case, by using the last exposed cohort $l = \max\{E_i\}$, we obtain

$$\begin{aligned} \hat{\beta}_{IV} &= \frac{\sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [Y_{e,t} - Y_{e,1} - (Y_{l,t} - Y_{l,1})]}{\sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [D_{e,t} - D_{e,1} - (D_{l,t} - D_{l,1})]} \\ &= \frac{\sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [D_{e,t} - D_{e,1} - (D_{l,t} - D_{l,1})] \cdot \widehat{WDID}_{e,t}}{\sum_e \sum_t N_{e,t} \hat{Z}_{e,t} [D_{e,t} - D_{e,1} - (D_{l,t} - D_{l,1})]}. \end{aligned}$$

where we define

$$\widehat{WDID}_{e,t} \equiv \frac{[Y_{e,t} - Y_{e,1} - (Y_{l,t} - Y_{l,1})]}{[D_{e,t} - D_{e,1} - (D_{l,t} - D_{l,1})]}.$$

From the Law of Large Numbers and the same argument in the proof of Theorem 2, we have

$$\widehat{WDID}_{e,t} \xrightarrow{p} \begin{cases} 0 & (t < e) \\ CLATT_{e,t} & (l-1 \geq t \geq e) \\ \frac{CAET_{e,t}^1 \cdot CLATT_{e,t} - CAET_{l,t}^1 \cdot CLATT_{l,t}}{CAET_{e,t}^1 - CAET_{l,t}^1} & (t \geq l) \end{cases} \quad (73)$$

Similarly, we obtain

$$[D_{e,t} - D_{e,1} - (D_{l,t} - D_{l,1})] \xrightarrow{p} \begin{cases} 0 & (t < e) \\ CAET_{e,t}^1 & (l-1 \geq t \geq e) \\ CAET_{e,t}^1 - CAET_{l,t}^1 & (t \geq l) \end{cases} \quad (74)$$

Combining the result (73) with (74) and by the Slutsky's theorem, we obtain

$$\beta_{IV} = \sum_e \sum_{l-1 \geq t \geq e} w_{e,t}^1 \cdot CLATT_{e,t} + \sum_e \left[\sum_{t \geq l} w_{e,t}^2 \cdot \frac{CAET_{e,t}^1 \cdot CLATT_{e,t} - CAET_{l,t}^1 \cdot CLATT_{l,t}}{CAET_{e,t}^1 - CAET_{l,t}^1} \right].$$

Completing the proof. \square

C.3 Random assignment of the instrument adoption date

First, we set up the additional notations. We define $LATE_k^{CM}(W)$ and $LATE_k(W)$ analogous to $CLATT_k^{CM}(W)$ and $CLATT_k(W)$ in section 4:

$$LATE_k^{CM}(W) \equiv \sum_{t \in W} \frac{AE_{k,t}^1}{\sum_{t \in W} AE_{k,t}^1} LATE_{k,t},$$

$$LATE_k(W) \equiv \frac{1}{T_W} \sum_{t \in W} AE_{k,t}^1 LATE_{k,t},$$

where we replace $CAET_{k,t}^1$ and $CLATT_{k,t}$ with $AE_{k,t}^1$ and $LATE_{k,t}$ respectively in $CLATT_k^{CM}(W)$ and $CLATT_k(W)$.

Theorem 4 below presents the TWFEIV estimand under Assumptions 1 - 5 and Assumption 14 (Random assignment assumption of adoption date E_i).

Theorem 4. Suppose Assumptions 1-5 and 14 holds. Then, the population regression coefficient β_{IV} consists of two terms:

$$\beta_{IV} = WLATE - \Delta LATE.$$

where we define:

$$WLATE \equiv \sum_{k \neq U} w_{IV,kU} LATE_k^{CM}(POST(k)) + \sum_{k \neq U} \sum_{l > k} w_{IV,kl}^k LATE_k^{CM}(MID(k, l))$$

$$+ \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot LATE_l(POST(l)),$$

$$\Delta LATE \equiv \sum_{k \neq U} \sum_{l > k} \sigma_{IV,kl}^l \cdot [LATE_k(POST(l)) - LATE_k(MID(k, l))].$$

The weights $w_{IV,kU}$, $w_{IV,kl}^k$ and $\sigma_{IV,kl}^l$ are the same in the proof of Theorem 2.

Theorem 4 is analogous to Theorem 2, but $CAET_{e,l}^1$ and $CLATT_{e,l}$ are replaced by $AE_{e,l}^1$ and $LATE_{e,l}$ respectively because we assume a random assignment of adoption date. If we consider the restrictions on the effects of the instrument on the treatment and outcome as in section 4.2, the similar arguments hold as in Theorem 2, Lemma 2 and Lemma 3.

Proof. First, we note that Assumption 14 implies Assumption 6 and Assumption 7. Therefore, we obtain the result in Theorem 2 under Assumptions 1-5 and 14.

By noticing that we have $CATT_{e,l}^1 = AE_{e,l}^1$ and $CLATT_{e,l} = LATE_{e,l}$ under Assumption 14, we obtain:

$$\begin{aligned} CLATT_k^{CM}(W) &= LATE_k^{CM}(W), \\ CLATT_k(W) &= LATE_k(W). \end{aligned}$$

By replacing $CLATT_k^{CM}(W)$ and $CLATT_k(W)$ with $LATE_k^{CM}(W)$ and $LATE_k(W)$ in Theorem 2, we obtain the desired result. \square

D Proofs and discussions in section 7

In this appendix, we first derive equation (22). We then discuss the causal interpretation of the covariate-adjusted TWFEIV estimand under staggered DID-IV designs, imposing the additional assumptions.

D.1 Decomposing the between IV coefficient

Let $\hat{C}_b^{D,\tilde{z}}$ denote the covariance between $D_{i,t}$ and $\tilde{z}_{k,t}$, the between term of $\tilde{z}_{i,t}$. The between IV coefficient $\hat{\beta}_{b,IV}^z$ is:

$$\hat{\beta}_{b,IV}^z = \frac{\hat{C}(Y_{i,t}, \tilde{z}_{k,t})}{\hat{C}(D_{i,t}, \tilde{z}_{k,t})} = \frac{\hat{C}(Y_{i,t}, \tilde{z}_{k,t})}{\hat{C}_b^{D,\tilde{z}}}. \quad (75)$$

To derive equation (22), we decompose the covariance between $Y_{i,t}$ and $\tilde{z}_{k,t}$. To do so, we first split the between term $\tilde{z}_{k,t}$ into the between term of $Z_{i,t}$ and the between term of $\tilde{p}_{i,t}$:

$$\begin{aligned} \tilde{z}_{k,t} &= [(\bar{Z}_{k,t} - \bar{Z}_k) - (\bar{Z}_t - \bar{Z})] - [(\bar{p}_{k,t} - \bar{p}_k) - (\bar{p}_t - \bar{p})] \\ &\equiv \tilde{Z}_{k,t} - \tilde{p}_{k,t}. \end{aligned}$$

Then, we have

$$\begin{aligned} \hat{C}(Y_{i,t}, \tilde{z}_{k,t}) &= \frac{1}{NT} \sum_i \sum_t Y_{i,t} [(\bar{z}_{k,t} - \bar{z}_k) - (\bar{z}_t - \bar{z})] \\ &= \frac{1}{T} \sum_k n_k \sum_t \bar{Y}_{k,t} [(\bar{z}_{k,t} - \bar{z}_k) - (\bar{z}_t - \bar{z})] \\ &= \sum_k \sum_l n_k n_l \frac{1}{T} \sum_t (\bar{Y}_{k,t} - \bar{Y}_{l,t}) \tilde{Z}_{k,t} - \sum_k \sum_l n_k n_l \frac{1}{T} \sum_t (\bar{Y}_{k,t} - \bar{Y}_{l,t}) \tilde{p}_{k,t} \\ &= \sum_k \sum_{l>k} (n_k + n_l) [\hat{C}_{kl}^{D,Z} \hat{\beta}_{IV,kl}^{2 \times 2} - \hat{C}_{b,kl}^p \hat{\beta}_{b,IV,kl}^p]. \end{aligned} \quad (76)$$

$\hat{\beta}_{IV,kl}^{2 \times 2}$ is an estimator obtained from an IV regression of $Y_{i,t}$ on $D_{i,t}$ with $\tilde{Z}_{k,t}$ as the excluded instrument in (k, l) cell subsample. $\hat{C}_{kl}^{D,Z}$ is the covariance between $D_{i,t}$ and $\tilde{Z}_{k,t}$ in (k, l) cell

subsample. $\hat{\beta}_{b,IV,kl}^p$ is an estimator obtained from an IV regression of $Y_{i,t}$ on $D_{i,t}$ with $\tilde{p}_{k,t}$ as the excluded instrument in (k, l) cell subsample. $\hat{C}_{b,kl}^p$ is the covariance between $D_{i,t}$ and $\tilde{p}_{k,t}$ in (k, l) cell subsample.

By combining (75) with (76), we obtain equation (22).

D.2 Causal interpretation of the covariate-adjusted TWFEIV estimand

This section considers the causal interpretation of the covariate-adjusted TWFEIV estimand β_{IV}^X . To simplify the analysis, we first make the following assumptions. Goodman-Bacon (2021) also make similar assumptions to investigate the causal interpretation of the covariate-adjusted TWFEIV estimand in Appendix B.

- (i) Time-varying covariates $X_{i,t}$ are not affected by instrument (policy shock).
- (ii) Time-varying covariates $X_{i,t}$ do not vary within cohorts.
- (iii) The coefficients obtained from regressing $\tilde{Z}_{i,t}$ on $\tilde{X}_{i,t}$ in (k, l) cell subsample are the same regardless of the pair (k, l) .

Because Assumption (ii) implies that the within term is equal to zero, the covariate-adjusted TWFEIV estimator $\hat{\beta}_{IV}^X$ simplifies to

$$\hat{\beta}_{IV}^X = \sum_k \sum_{l>k} s_{b,kl} \hat{\beta}_{b,IV,kl}^z.$$

Assumption (iii) guarantees that $\hat{\beta}_{b,IV,kl}^z$ is equal to the between coefficient obtained from estimating equation (14) in (k, l) subsample, which we denote $\hat{\beta}_{b,IV,kl}^{z,X}$ hereafter. To see this formally, let $\tilde{p}_{i,t}^{kl} \equiv \hat{\Gamma}_{k,l} \tilde{X}_{i,t}$ denote the linear projection obtained from regressing $\tilde{Z}_{i,t}$ on $\tilde{X}_{i,t}$ in (k, l) subsample and let $\tilde{p}_{j,t}^{kl}$ denote the between term of $\tilde{p}_{i,t}^{kl}$ in cohort j . We note that $\tilde{p}_{j,t}^{kl} \neq \tilde{p}_{j,t}$ (the between term of $\tilde{p}_{i,t}$) holds in general because $\tilde{p}_{i,t} = \hat{\Gamma} \tilde{X}_{i,t}$ is estimated using the whole sample. Then, we have

$$\begin{aligned} \hat{\beta}_{b,IV,kl}^z &= \frac{\hat{C}(Y_{i,t}, \tilde{z}_{j,t})}{\hat{C}(D_{i,t}, \tilde{z}_{j,t})} = \frac{\hat{C}(Y_{i,t}, \tilde{Z}_{j,t} - \tilde{p}_{j,t}^{kl}) + \hat{C}(Y_{i,t}, \tilde{p}_{j,t}^{kl} - \tilde{p}_{j,t})}{\hat{C}(D_{i,t}, \tilde{Z}_{j,t} - \tilde{p}_{j,t}^{kl}) + \hat{C}(D_{i,t}, \tilde{p}_{j,t}^{kl} - \tilde{p}_{j,t})} \\ &= \frac{\hat{C}(D_{i,t}, \tilde{Z}_{j,t} - \tilde{p}_{j,t}^{kl}) \hat{\beta}_{b,IV,kl}^{z,X} + \hat{C}(D_{i,t}, \tilde{p}_{j,t}^{kl} - \tilde{p}_{j,t}) \hat{\beta}_{b,IV,kl}^{dif}}{\hat{C}(D_{i,t}, \tilde{Z}_{j,t} - \tilde{p}_{j,t}^{kl}) + \hat{C}(D_{i,t}, \tilde{p}_{j,t}^{kl} - \tilde{p}_{j,t})}, \quad j = k, l, \end{aligned}$$

where $\hat{\beta}_{b,IV,kl}^{dif}$ is an estimator obtained from an IV regression of $Y_{i,t}$ on $D_{i,t}$ with the difference $\tilde{p}_{j,t}^{kl} - \tilde{p}_{j,t}$ as the excluded instrument. Because Assumption (iii) ($\tilde{p}_{j,t}^{kl} = \tilde{p}_{j,t}$) implies $\hat{C}(D_{i,t}, \tilde{p}_{j,t}^{kl} - \tilde{p}_{j,t}) = 0$, we obtain $\hat{\beta}_{b,IV,kl}^z = \hat{\beta}_{b,IV,kl}^{z,X}$.

Hereafter, we assume the identifying assumptions in staggered DID-IV designs and Assumption (i)-(iii). We focus on the between coefficient $\hat{\beta}_{b,IV,kl}^z = \hat{\beta}_{b,IV,kl}^{z,X}$ as it clarifies how covariates affect the interpretation of the TWFEIV estimand:

$$\hat{\beta}_{b,IV,kl}^z = \frac{\hat{C}(Y_{i,t}, \tilde{Z}_{j,t}) - \hat{C}(Y_{i,t}, \tilde{p}_{j,t}^{kl})}{\hat{C}(D_{i,t}, \tilde{Z}_{j,t}) - \hat{C}(D_{i,t}, \tilde{p}_{j,t}^{kl})}, \quad j = k, U.$$

Then, by the similar calculations in the proof of Theorem 2, we obtain

$$\begin{aligned}\hat{C}(Y_{i,t}, \tilde{Z}_{j,t}) &= \hat{V}_{kU}^z \cdot \hat{D}_{kU}^{2 \times 2} \beta_{kU,IV}^{2 \times 2} \\ &\xrightarrow{p} V_{kU}^z \cdot CAET_k^1(POST(k)) \cdot CLATT_k^{CM}(POST(k)),\end{aligned}\quad (77)$$

and

$$\begin{aligned}\hat{C}(Y_{i,t}, \tilde{p}_{j,t}^{kl}) &= \frac{n_{kU}(1 - n_{kU})}{T} \sum_t (\bar{Y}_{kt} - \bar{Y}_{Ut}) \cdot [(\bar{p}_{k,t}^{kU} - \bar{p}_k^{kU}) - (\bar{p}_{U,t}^{kU} - \bar{p}_U^{kU})] \\ &\xrightarrow{p} \frac{N_{kU}(1 - N_{kU})}{T} \sum_t \{E[Y_{i,t}(D_{i,t}^\infty) | E_i = k] - E[Y_{i,t}(D_{i,t}^\infty) | E_i = U]\} [(p_{k,t}^{kU} - p_k^{kU}) - (p_{U,t}^{kU} - p_U^{kU})] \\ &\quad + \frac{N_{kU}(1 - N_{kU})}{T - (k - 1)} \sum_{t \geq k} \underbrace{CAET_k \cdot CLATT_{k,t}}_{CAIET_{k,t}} \cdot [(p_{k,t}^{kU} - p_k^{kU}) - (p_{U,t}^{kU} - p_U^{kU})],\end{aligned}\quad (78)$$

where N_{kU} and $[(p_{k,t}^{kU} - p_k^{kU}) - (p_{U,t}^{kU} - p_U^{kU})]$ are the probability limits of n_{kU} and $[(\bar{p}_{k,t}^{kU} - \bar{p}_k^{kU}) - (\bar{p}_{U,t}^{kU} - \bar{p}_U^{kU})]$, respectively. Equations (77) and (78) indicate that covariates affects the causal interpretation of $\hat{\beta}_{b,IV,kU}^z$ in two ways. First, it additionally introduce the covariance between the difference in unexposed outcomes and the difference in the variation of the linear projection for cohorts k and U (the first term in equation (77)). Second, it additionally introduce the covariance between the $CAIET_{k,t}$ and the difference in the variation of the linear projection for cohorts k and U (the second term in equation (78)).