




The Effects of Higher-Stakes Teacher Evaluation on Office Disciplinary Referrals

David D. Liebowitz, Lorna Porter & Dylan Bragg



To cite this article: David D. Liebowitz, Lorna Porter & Dylan Bragg (2022): The Effects of Higher-Stakes Teacher Evaluation on Office Disciplinary Referrals, Journal of Research on Educational Effectiveness, DOI: [10.1080/19345747.2021.2015496](https://doi.org/10.1080/19345747.2021.2015496)

To link to this article: <https://doi.org/10.1080/19345747.2021.2015496>


 [View supplementary material](#) 


 Published online: 26 Jan 2022.

 [Submit your article to this journal](#) 

 [View related articles](#) 


 [View Crossmark data](#) 

 This article has been awarded the Centre for Open Science 'Open Materials' badge.

 This article has been awarded the Centre for Open Science 'Preregistered' badge.



The Effects of Higher-Stakes Teacher Evaluation on Office Disciplinary Referrals

David D. Liebowitz^a , Lorna Porter^a, and Dylan Bragg^b

^aDepartment of Educational Methodology, Policy and Leadership, University of Oregon, Eugene, Oregon, USA; ^bEducational and Community Supports, College of Education, University of Oregon, Eugene, Oregon, USA

ABSTRACT

The effects of imposing accountability pressures on public school teachers are empirically indeterminate. In this paper, we study the effects of accountability in the context of teacher responses to student behavioral infractions in the aftermath of teacher evaluation reforms. We leverage cross-state variation in the timing of state policy implementation to estimate whether teachers change the rate at which they remove students from their classrooms. We find that higher-stakes teacher evaluation had no causal effect on the rates of disciplinary referrals, and we find no evidence of heterogeneous effects for grades subject to greater accountability pressures or in schools facing differing levels of disciplinary infractions. Our results are precisely estimated and robust to a battery of assumption and specification checks.

ARTICLE HISTORY

Received 15 July 2020
Revised 16 July 2021
Accepted 15 October 2021

KEYWORDS


Teacher evaluation; accountability; school discipline; difference-in-differences

Introduction

In response to financial incentives from the Obama administration's 2009 Race to the Top competition, 44 states implemented reforms to their teacher evaluation policies between 2011 and 2016 (NCTQ, 2016; Steinberg & Donaldson, 2016). Much of the public narrative and high-profile scholarship on teacher evaluation has been shaped by a focus on its use for high-stakes personnel decisions (e.g., Hanushek, 2009; Jackson et al., 2014). However, reforms to evaluation systems were also explicitly designed to promote the development of teachers' skills (Connally & Tooley, 2016; Donaldson, 2021). Former Secretary of Education Arne Duncan (2012) defended the value of teacher evaluation reforms by claiming that, "[b]etter evaluation systems improve classroom instruction" (p. 4). Thus, in principle, these reforms intended to improve teacher effectiveness through the combination of increased accountability pressures and added supports for skill development.

Among practicing teachers and their evaluators, it is common knowledge that the ability to create focused, orderly classroom learning environments has a dramatic

CONTACT David Liebowitz  daviddl@uoregon.edu  Department of Educational Methodology, Policy and Leadership, University of Oregon, 5267, Eugene, OR 97403, USA.

 Supplemental data for this article can be accessed online at <http://dx.doi.org/10.1080/19345747.2021.2015496>

© 2021 Taylor & Francis Group, LLC

influence on overall instructional effectiveness.¹ This knowledge is borne out in the emphasis that classroom observation rubrics place on behavioral management skills and in the weight that evaluators assign to these practices. In the median state, one-quarter of suggested evaluation rubric indicators are focused on teachers' skill in classroom management (Gilmour et al., 2019). Additionally, unruly classes are easily observable for teachers' evaluators; more so than, for example, alignment of instruction to grade-level standards (Cohen et al., 2020). In fact, principals report attending to student engagement over content-specific instructional practices during observations and feedback (Wieczorek et al., 2019). Thus, if teachers improved their classroom management practices in response to the past decade's evaluation reforms, they may have needed to remove fewer students from class by more effectively preventing and responding to misbehavior.

On the other hand, teachers might also have accomplished the goal of reducing disruptive behavior by imposing a lower floor of tolerance for misbehavior before sending a student out of class and to the school office.² Additional scrutiny and stress from higher-stakes evaluations may have created unintended incentives to create classrooms with orderly appearances during evaluative observations by removing disruptive students at higher rates. Further, to the extent that disruptive students create negative externalities on other students' learning (e.g., Carrell & Hoekstra, 2010), teachers can both increase their expected observation score and improve their average student's achievement by reducing incidences of active disruption in their classrooms.

In this paper, we test whether higher-accountability teacher evaluation policies affect how teachers respond to behavioral infractions in their classroom. Theory suggests that teacher evaluation reforms might have decreased the rate at which teachers removed students from class as their classroom management skills improved. Alternatively, these reforms might have increased the rate of removal as an unintended byproduct of added stress or a desire to present an image of an orderly classroom to external observers. We leverage Kraft et al.'s (2020) tally of the introduction of teacher evaluation policy reform in conjunction with disciplinary data from a large network of over 2,500 schools implementing a common behavior management system to estimate the impact of higher-stakes evaluation on Office Disciplinary Referrals (ODRs) in a difference-in-differences framework.

Our main findings are that higher-stakes teacher evaluation had no causal effect on the overall rate of Office Disciplinary Referrals. Our estimates are precise, and we can confidently rule out effects larger than a decrease of 0.21 referrals—or an increase of 0.04 referrals—that originate from teachers' classrooms in the average-sized school on an average day. We find no evidence of heterogeneous impacts across a variety of well-theorized dimensions that might influence how teacher evaluation influences disciplinary response. We subject our identification strategy and its assumptions to a host of robustness checks and consistently find that our results hold. In order to promote

¹Lazear's (2001) seminal work on the production of education lays out a theoretical relationship between instructional effectiveness and classroom behavior, but this phenomenon is dramatically understudied empirically.

²The typical mechanism by which teachers respond to student behavior that they have determined cannot be addressed in the classroom is to send students to a school administrator (e.g., principal, assistant principal, dean of students) in the school's office. Other approaches include waiting to speak to an administrator in the hallway. For the purposes of this paper, we describe all such events as Office Disciplinary Referrals (ODRs).

research transparency, we pre-registered our analytic plan in the Registry of Efficacy and Effectiveness Studies (REES #1748) prior to receipt of our data.

Our findings contribute to new literatures estimating the causal effects of accountability pressures on the black box of within-classroom behaviors by teachers and the effects of school policies on disciplinary processes. While there is a growing consensus that teacher evaluation policies increase voluntary exits, particularly among teachers rated poorly in observation- and value-added-based evaluation schemes (Cullen et al., 2019; Dee & Wyckoff, 2015; Loeb et al., 2015) and some promising evidence that teacher observation improves student outcomes (Burgess et al., 2021, Phipps, 2021; Taylor & Tyler, 2012), we know less about how high-stakes teacher evaluation policies change what teachers do in the classroom. Our failure to detect substantial changes in behavior management practices aligns with Phipps and Wiseman (2019) who find no evidence that teachers shift their focus to a particular instructional domain as the accountability pressure of an evaluator observation increases. Similarly, Garet et al. (2017) find no evidence of changes in instructional practices in the aftermath of receiving performance feedback.

We study phenomena closely related to those described in Holbein and Ladd (2017). Holbein and Ladd examine the effect of schoolwide accountability pressures, enacted through the labeling of a school as failing to make Adequate Yearly Progress under the No Child Left Behind Act, on the frequency of serious student misbehaviors leading to suspension. They find that increased schoolwide accountability pressured led to an increase in suspensions. While we also study accountability pressures, we study those that fall individually and exclusively on teachers. Additionally, we have available several outcome measures that allow us to assess the likelihood that our results are driven by changes in either student or educator behavior.

As a whole, we interpret our results as evidence that the introduction of higher-stakes accountability policies did not dramatically alter classroom disciplinary climates. Given the nature of our data, we are unable to rule out the possibility that our results are a product of heterogeneous responses to increased accountability by teachers' characteristics or skill. Nevertheless, we argue that our central contribution is to highlight the loose coupling of teacher accountability policy and classroom practices, as measured by teacher-initiated behavioral referrals. As implemented, teacher evaluation reforms were not sufficient to limit students' entry into the disciplinary pipeline. We reach this conclusion because we find no evidence of either positive (intended) or negative (unintended) consequences of teacher evaluation reforms on classroom disciplinary environments.

We begin by providing an overview of teacher evaluation and student discipline policies and processes. Next, we describe our disciplinary and policy data, and we present our empiricseal estimation framework. We then share our main results as well as evidence on the presence of heterogeneous effects. In our penultimate section, we present a host of assumption and robustness checks. Finally, we conclude with a discussion of how our findings provide insight into the complex relationship between evaluation policy and teachers' behaviors.

Teacher Evaluation and Student Discipline Policy

We briefly contextualize teacher evaluation policy within the nuanced literature on the broader topic of educational accountability. We then review the characteristics of U.S. teacher evaluation and student disciplinary policy reforms introduced in the past two decades. We then explore how these strands of educational policy might intersect in the classroom.

Accountability and Teacher Evaluation

External accountability pressures and incentives are central tools policy makers possess to improve teaching and learning conditions in schools. However, there is a complex relationship between incentive- and accountability-based policies and local actor behavior, particularly in the public sector (e.g., Dixit, 2002). A rich theoretical and empirical debate on the merits of accountability-driven policies in education exists, with much of the causal evidence finding either mixed or no effects.³ Even well-designed accountability and incentive policies can generate unintended responses, including educational triage (Ladd & Lauen, 2010; Neal & Schanzenbach, 2010; Reback, 2008), curriculum narrowing (Hamilton et al., 2005), and gaming (Figlio, 2006; Vogell, 2011).

Teacher evaluation, a particularly high-profile incarnation of educational accountability, seeks to leverage the individual application of incentives and sanctions to improve instructional practices and student outcomes. As with the broader literature on accountability, the evidence on the effectiveness of teacher evaluation is mixed.⁴ Of particular interest to our line of inquiry, teachers report that higher-stakes evaluation policies introduced in the last decade raised teacher anxiety as a result of worries about their self-efficacy (Ford et al., 2017). Thus, system leaders have a critical interest in understanding whether and in what ways teachers and principals respond to increased accountability in teacher evaluation.

As we note, 44 states implemented teacher evaluation reforms in response to the Race to the Top (RTTT) initiative. In fact, 35 states enacted reforms of teacher evaluation between 2009 and 2011 (NCTQ, 2011, 2017), but the exact timing of when these reforms came into effect spanned the subsequent six years.⁵ Our identification strategy exploits this exogenous federal shock to state policy, and what we will argue below are near-random cross-state differences in the timing of the subsequent implementation of teacher evaluation reforms.

While the particulars of each state's evaluation framework differ, as does district-level policy implementation, this variation is largely endogenous. Thus, we focus on the common accountability elements rather than on the intensity of accountability pressures across and within states and districts. Generally, state legislation and regulation defined

³See, among others, Brehm et al. (2017), Chakrabarti (2014), Chiang (2009), Deming et al. (2016), Eren (2019), Macartney (2016), Ozek (2012), and Reback et al. (2014). Deming and Figlio (2016) synthesize the literature on educational accountability.

⁴See, among others, Cullen et al. (2019), Kraft et al. (2020), Macartney et al. (2019), Pope (2019), Rothstein (2015), Steinberg and Sartain (2015), Stecher et al. (2018), Strunk et al. (2017) and Taylor and Tyler (2012). Liebowitz (2020) summarizes this nuanced literature.

⁵Alaska, Maine, Mississippi, New Jersey, North Dakota and Pennsylvania passed new teacher evaluation laws in 2012; Kentucky, South Carolina and Texas did so in 2013.

parameters for the evaluation process to which local districts' collective bargaining agreements were required to adhere. Thirty-five states promulgated model evaluation policies that many districts adopted as written, although others adopted them with modifications (NCTQ, 2017; Steinberg & Donaldson, 2016). All state reforms to teacher evaluation required that classroom observation of teaching practice be a part of a teacher's final rating, and in 59 percent of those states, these reforms established a minimum frequency of classroom observations for at least some teachers. In almost all cases, teacher evaluation reforms entailed adopting a common rubric for evaluating teachers' performance with multiple rating categories. In addition, 80 percent of states required some or all teachers to be evaluated based on student-learning gains (either through formal measurements of students learning, through teachers' contributions to students' progress toward locally determined learning objectives, or both). Some states additionally included measures of whole-school performance or parent-, student-, and peer-surveys of teacher competency. Eighty-three percent of new state teacher evaluation policies linked evaluation results to teachers' future professional development plans. Some required annual appraisals for all teachers, while in other states annual appraisal was limited to new or probationary teachers (Jacobs & Doherty, 2015; Steinberg & Donaldson, 2016; Winters & Cowen, 2013).

While summative teacher ratings have varying consequences, the policy reforms *as designed* represent substantial increases in accountability pressures on teachers. Over three-fifths (61 percent) of states instituted rules that led to the dismissal of teachers who were not rated Proficient, and almost half (48 percent) of states used evaluation results for tenure decisions (Steinberg & Donaldson, 2016). The Institute of Education Sciences (2014) judged the recent round of teacher evaluation reforms to have made strong shifts toward increased accountability. However, *as implemented*, most teachers across the country continued to receive ratings above the standard of proficiency, even after the introduction of clear observational rubrics and rating categories (Kraft & Gilmour, 2017). In many cases, principals felt unqualified, had insufficient time to conduct high-quality observations, or adapted the implementation of these policies to fit within their preexisting priorities (Donaldson & Woulfin, 2018; Kraft & Gilmour, 2016). Thus, the added accountability (our treatment) may have been too weak to result in substantial changes in teachers' practice; a consideration to which we return in the discussion of our findings.

Though we do not examine between-state differences in accountability intensity we do explore cross-grade differences. In particular, we hypothesize that grades 3-11 will be subject to greater levels of accountability as these are years in which mandated state testing occurs. Therefore, teachers of tested subjects in these grades experience accountability from both evaluative classroom observations and student test score outcomes.

Student Discipline

We explore the effects of teacher evaluation on student disciplinary outcomes because of the outsize impacts that entry into the school discipline system have on students' later life outcomes. Mounting evidence suggests credibly causal links between suspensions and involvement with the criminal justice system (Bacher-Hicks et al., 2019;

Sorensen et al., 2021). Office Disciplinary Referrals (ODRs), though less studied, are important precursors to these more severe disciplinary outcomes and provide important insights into teachers' responses to misbehavior. In this section, we describe recent policy developments related to student discipline, evidence on their implementation, and what happens when teachers remove students from their classroom.

In the late 1980s and 1990s, states and districts increasingly adopted sets of policies collectively known as "zero-tolerance discipline." Many of these policies originated in response to the federal Gun-Free School Act of 1994, but soon extended beyond firearm offenses (Curran, 2016). Under such policies, students committing one among a set of pre-specified disciplinary infractions were to be suspended or expelled from school, without discretion. In response to widespread concern about the long-term effects of exclusionary discipline and the disproportionate use of such disciplinary approaches for students of color, states and districts initiated reforms of many zero-tolerance laws in the 2000s and 2010s (Rafa, 2019; U.S. Department of Education & U.S. Department of Justice, 2014).

Some states have prohibited the use of suspension for less severe infractions (such as defiance or truancy) or the length of suspensions overall, while others prohibit the use of suspension in earlier grades except for extreme misbehavior (Steinberg & Lacoë, 2018; Anderson, 2020). Others require school districts to develop discipline plans that incorporate alternative discipline programs, such as Positive Behavioral Intervention and Supports (PBIS), Multi-Tiered Systems of Support (MTSS) and restorative justice practices (Rafa, 2019; Welsh & Little, 2018).

While discipline policy reforms have generally focused on suspension and expulsion, the process by which students are removed from class for lower-level infractions is less frequently a target of policy. To the student, Office Disciplinary Referrals (ODRs) are an entry point into the disciplinary system. ODRs cause students to miss instructional time, may lead to additional consequences ranging from lunch detention to suspension, and may represent a signal to students about how teachers perceive their behavior (Kennedy-Lewis & Murphy, 2016). Alternatively, they may provide students with opportunity to reflect, develop relationships with school administrators, and improve their future behavior.

The consequences of ODRs extend beyond the referred student to classmates and educators. ODRs have spillover effects as they influence the classroom composition during the time the student is out of class, which influences the structure of peer effects (Lazear, 2001). Further, ODRs requiring subjective interpretations of behaviors yield important signals about teachers' tolerance for student misbehavior and classroom management skills. Finally, ODRs impose significant time burdens on educators who receive students removed from the classroom. While no precise estimates of the amount of time that students spend outside of the classroom during each referral exist, some states mandate that students not be sent back to class sooner than 30 minutes, within the same class period, or before the principal has undertaken one of a set of prescribed disciplinary measures (e.g., Louisiana Revised Statute §17.416 A.(1)(b)(iii)). For our purposes, the two essential takeaways are: (1) state policy generally grants broad discretion to teachers to remove students from the classroom; and (2) these removals have significant, though imperfectly understood, implications for both students and educators.

The Relationship Between Teacher Evaluation and Student Discipline

Our analytic challenge is to understand the relationship between higher-stakes teacher evaluation and office disciplinary referrals. *Ex ante*, it is not obvious what the effects of greater accountability pressures might be on the rate of ODRs. On one hand, it is possible that greater accountability pressures may lead to an increase in the rate at which teachers send students out of class. In particular, if classroom observations are key contributors to teacher evaluation scores under high-stakes evaluation systems, teachers may be more likely to send students out of class for lesser infractions than under low-stakes evaluation conditions in the hopes that fewer disruptions occur during a supervisor's unanticipated visit. Further, if teacher evaluation scores are tied to assessment performance, teachers may utilize ODRs more frequently to send students out of class if they perceive doing so will result in a more effective learning environment for the majority of students.

On the other hand, it is possible that teacher evaluation reforms might lead to a decrease in the rate of ODRs. A decrease in referrals may be evidence of improved teacher skills resulting from accountability and professional supports introduced as part of the wave of teacher evaluation reforms. Classroom observation rubrics such as the Danielson Framework for Teaching (Danielson 1996), the CLASS (Mashburn et al., 2008) tool, and other such rubrics adopted by states and districts to conduct observations of teachers include measures of behavioral management skills, and the feedback teachers receive from those observations may serve as professional development. However, given other evidence that large-scale implementation of higher-stakes evaluation policies did not improve general instructional practices (Garet et al., 2017; Stecher et al., 2018), and that instructional guidance on classroom management practices from evaluation rubrics tends to be weak (Gilmour et al., 2019), our prior expectation is low that teachers' behavioral management skills will increase as a result of higher-stakes teacher evaluation policy implementation.

An alternative explanation for how the introduction of high-stakes evaluation policies might decrease ODRs is that evaluators only weakly observe and monitor teachers' behavior management skills. Given principals' difficulty in finding time to observe teaching practice (Kraft & Gilmour, 2016), they may use the frequency of teachers' referral of students to the office as a proxy for their skill in classroom management. This would create an incentive for teachers to reduce the frequency of ODRs in the aftermath of high-stakes evaluation reform but not be desirable from the policy maker's perspective. Still another potential explanation for negative effects of evaluation on ODRs is that increased administrator presence in classrooms might reduce misbehavior. Thus, positively signed coefficients of the effects of the introduction of new teacher evaluation systems would provide relatively straightforward interpretations of the effects of increased accountability on teachers' practice, while negatively signed ones would be ambiguous.

It is also possible that teachers respond differently to higher-stakes evaluation when the starting behavioral climate in their school differs or when effective disciplinary support strategies exist to respond to misbehavior. In the first case, schools with few initial behavioral incidents are unlikely to focus evaluative feedback on teachers' classroom management. In the second case, if evaluation reforms increased teachers' behavioral response skills but there were no schoolwide strategies to address misbehavior,

improvements in teachers' practice might be dominated by poor school culture systems and, therefore, be unobserved in the rates of ODRs. Alternatively, if evaluation pressures increase stress, teachers may be more prone to remove students from their rooms in schools with ill-defined criteria for when a student is to be sent out of class. Thus, exploring treatment heterogeneity by schoolwide climate and quality of support systems may reveal valuable insights into how the effects of evaluation differ across contexts.

Given the substantive importance, but empirical uncertainty, of the interaction between teacher evaluation and student discipline policies, we pose the following research questions:

1. How did the introduction of higher-stakes teacher evaluation policies affect the rates at which students were removed from the classroom as a consequence for misbehavior?
2. How, if at all, did the effect of the higher-stakes teacher evaluation on classroom disciplinary responses differ by school disciplinary climate and grade-level-specific accountability pressures?

Data and Sample

We present in the main text of the paper a brief description of our data and sample construction and share complete information in [Supplementary Appendix B](#).

Discipline Data

The primary data source for our analysis is the School-Wide Information System (SWIS) data system. This data system is used by schools that track behavioral data associated with the implementation of Positive Behavioral Implementation and Supports (PBIS) and maintained by the Education and Community Supports research and outreach unit at the University of Oregon. A particular strength of this dataset is that it records low-level disciplinary infractions rarely recorded in administrative data that reflect teacher (rather than administrator) behaviors and may represent critical entry points for students into an exclusionary discipline system. We draw on information from the 2006-07 through the 2017-18 school years. These data also include enrollment information that combine the best available information from schools' self-reports and their October 1 administrative count from the NCES Common Core Data. We supplement our main estimates with placebo tests using suspension data from the restricted-use Civil Rights Data Collection (CRDC).

Policy Implementation Data

To determine the timing of teacher evaluation reforms in 44 states, we draw on Kraft et al. (2020) who extend prior reviews by Steinberg and Donaldson (2016) and the National Council on Teacher Quality (NCTQ) (2016). [Supplementary Appendix Table A1](#) provides state-by-state policy implementation details and counts of schools by state in our data.

To address concerns that shifts in ODR activity reflect reforms in school discipline policy, we collect data on two categories of disciplinary policy reforms and test whether our results are robust to these policies. We compile information from the Compendium of School Discipline Laws and Regulations (Bezinque et al., 2018) on whether any state-level reforms to classroom and school discipline policies occurred during our analytic window.

Measures

Our primary predictor is a binary indicator for whether a state has implemented a higher-stakes teacher evaluation system. Following Kraft et al. (2020), we code *Implement Evaluation* as one in the first fall in which a state implements a new teacher evaluation policy and in all subsequent observations. It takes the value of zero prior to policy implementation and all observations in states that never implement evaluation reform.

We adopt two primary outcomes based on yearly counts of the total number of Office Disciplinary Referrals (ODRs), aggregated to the grade level. We scale these counts to reflect the number of ODRs that occur in a school of average size in our sample on a given day. Thus, our primary outcome reflects the number of referrals originating from a classroom per-500 students, per-day. We supplement this outcome with the rate of classroom-originating referrals that are a result of one of six behaviors defined by an expert panel (Greflund et al., 2014) as “subjective.”⁶ We anticipate that these subjective-classroom misbehaviors allow teachers more discretion in their disciplinary response and, thus, may be more sensitive to policy changes. We describe in the following section analytic approaches that take advantage of measures of ODRs that originate in locations other than the classroom (“Other”) and referrals originating from the classroom in response to one of 14 “objective” behavioral infractions.

We use a time-varying, binary indicator to measure schools’ successful implementation of PBIS. *Implement PBIS Well* takes the value of one in years in which schools successfully implement PBIS and zero otherwise.⁷

⁶We categorize subjective and objective behaviors data following Greflund et al. (2014): “Subjective behaviors were defined as behaviors that require not simply observing a discrete, objective event (e.g., a student smoking), but a significant value judgment regarding whether the intensity or quality of the behavior warrants an ODR (e.g., a student using inappropriate language). (...) The following behaviors were categorized as subjective: abusive language/inappropriate language/profanity, defiance/disrespect/insubordination/ non-compliance, harassment/bullying, disruption, dress code violation, and inappropriate display of affection. The following behaviors were categorized as less subjective: physical aggression/fighting, tardy, skipping, truancy, property damage/vandalism, forgery/theft, inappropriate location/out of bounds, use/possession of tobacco, alcohol, drugs, combustibles, weapons, bomb threat/false alarm, and arson. Three problem behaviors did not meet the inter-rater reliability criterion and were also not classified as subjective: lying/cheating, technology violation, and gang affiliation display” (pp. 220–221).

⁷To be classified as successfully implementing PBIS, schools had to meet one of the following thresholds: School-wide Evaluation Tool (SET): greater than or equal to 80 percent of expectations taught and overall implementation; Tiered Fidelity Inventory (TFI): Tier 1 ratio greater than or equal to 70 percent; Benchmark of Quality (BOQ: Total Ratio greater than or equal to 70 percent; Self-Assessment Survey (SAS): Implementation Average greater than or equal to 80; and Team Implementation Checklist (TIC): Implementation Average greater than or equal to 80 percent. We refer readers to McIntosh et al. (2013) and Mercer, McIntosh and Hoselton (2017) for details on the validation of these instruments. We do not use the continuous implementation scores as they represent substantially different scales and are not linked across instruments (Greflund et al., 2014; Mercer et al., 2017).

We include several other disciplinary measures related to our sensitivity tests. We draw on five waves of data from the CRDC, from the 2005 to 2006 school year to the 2015–2016 school year, that count the number of students suspended in a given year. We calculate the yearly proportion of students who are suspended in each school.⁸ Drawing on the Compendium of School Discipline Laws and Regulations, we code *Teacher authority to remove students from the classroom* and *Limitations, conditions, or exclusions for use of suspension and expulsion* as one in the first fall for all schools in states that enacted reforms to these policies between 2006 and 2018 ([Supplementary Appendix Table B2](#)).

Analytic Sample and Summary Statistics

Our analytic sample focuses on U.S. traditional public schools subject to state evaluation policies, for which we have outcome measures both before and after policy implementation. We restrict our sample to grade-year observations nested in schools which we observe at least four years before the adoption of high-stakes teacher evaluation and one year after the initial implementation year. We form our measures of Office Disciplinary Referrals from counts of referrals at the grade-school-year level. Thus, our main analytic sample includes 107,468 grade-school-year observations, nested in 20,137 school-year observations. These represent a total of 2,564 schools in 939 districts.

In [Table 1](#), we present summary statistics for the full sample, for schools and students located in states that never implemented high-stakes evaluation, and for schools and students in states that did. Our sample is not nationally representative. Schools in our sample are attempting to implement PBIS, have a standardized office referral process, use the SWIS data management system, and agree that their behavioral data can be used for research purposes. The strictest definition of the population to which our results generalize are these 2,564 schools. However, approximately 25,000 schools (or more than one-quarter of all U.S. public schools) were attempting to implement PBIS in some form in the 2016–17 school year, and about 11,000 of these schools used the SWIS data management system (Hoselton, 2018). We argue that we can cautiously generalize to this larger set of schools that differ from our sample only in the software application they use to track behavior or in their willingness to permit their data be used for research purposes.

Further, the demographic characteristics of our sample broadly match national racial and family income enrollment patterns; none differ by more than 4 percentage points. That said, schools attempting to implement PBIS might respond differently to higher-stakes teacher evaluation. Thus, [Table 1](#) is meant to situate our work within the larger U.S. traditional public school context, but not to suggest that our findings are directly generalizable to this context.

⁸We present additional information on the CRDC sample of schools in [Supplementary Appendix Table B1](#). In the sample of 343,015 school-year observations in the CRDC, the average school suspends 6 percent of its students per year ($SD = 0.09$). In auxiliary regressions, we find that districtwide patterns in ODRs are positively correlated to suspensions in the CRDC data, but imperfectly so. This aligns with our proposed use of the suspension data as a placebo test that measures a related, but different construct than classroom referrals.

Table 1. Descriptive statistics on school-wide information system (SWIS) data, 2006–2017.

	Full sample	NCES public (2006–2016)	Never evaluation	Ever Evaluation
Total schools	2,564	98,556	750	1,814
School-year observations	20,137		4,334	15,803
Grade-year observations	107,468		24,499	82,969
Total districts	939	13,647	313	626
IQR schools per district	2–10		1–8	2–12
School level				
Elem (K-6)	1,683		484	1,199
Middle (6–9)	405		97	308
High (9–12)	159		40	119
Multi-level grade span	317		128	189
School locale				
Rural	516	0.29	195	321
Town	502	0.14	187	315
Suburban	834	0.30	104	730
Urban	712	0.22	263	449
			Pre-2011 characteristics	
School characteristics				
Avg. school enrollment	510.2 (306.7)	521.2	437.6	509.2***
% Low income	0.52 (0.25)	0.49	0.49	0.52**
% American Indian/Native Alaskan	0.01 (0.05)	0.01	0.02	0.01***
% Asian/Pacific-Islander	0.05 (0.08)	0.05	0.05	0.03***
% Black	0.13 (0.19)	0.16	0.05	0.16***
% Hispanic	0.20 (0.23)	0.24	0.20	0.12***
% White Non-Hispanic	0.54 (0.31)	0.52	0.63	0.58***
% Schools by year implementing PBIS	0.71		0.60	0.62
Grade-level outcomes				
Daily referrals per 500 students – Classroom	2.03 (2.83) [1.69]		2.30	2.17**
Daily referrals per 500 students – other location	1.56 (1.77) [1.22]		1.97	1.84***
Daily referrals per 500 students – classroom subject	1.29 (1.95) [1.23]		1.61	1.43***
Daily referrals per 500 students – classroom object	0.48 (0.90) [0.40]		0.41	0.45***
% Schools under high-stakes evaluation				
2010–11	0.00			
2011–12	0.13			
2012–13	0.16			
2013–14	0.36			
2014–15	0.52			
2015–16	0.55			
2016–17	0.72			
2017–18	0.72			

Notes: * $p < .05$, ** $p < .01$, *** $p < .001$. Sample characteristics and outcomes weighted by school enrollment. Standard deviations, where applicable, in parentheses. Within-school standard deviations in brackets. 85 school-year observations have race/ethnicity data imputed to the same-school median. 40 school-year observations have race/ethnicity data imputed to the district-school year median. 490 school-year observations have low-income data imputed to the same-school median. 59 school-year observations have low-income data imputed to the district-school year value. 67 school-year observations have low income data imputed to the district median. Low-income and race enrollment capped at 100 percent of school enrollment. Outcomes above 99th percentile re-coded to value of 99th percentile. We use these imputations and outcome caps in all regression estimates. Schools implementing PBIS defined in text. National public school averages are from NCES Digest of Education Statistics between fall 2006 and 2016. NCES counts are averages of counts from 2006 to 2016, while NCES means are averages of yearly averages.

In our full sample, the average per-day number of referrals from classroom settings (Classroom) is 2.03 per 500 students and the average per-day number of referrals from all other locations within or beyond the school grounds (Other) is 1.56 referrals per 500 students. This implies that the average number of referrals across a 180-day school year for an average-sized school is around 650. Assuming each student spends 30 minutes out of class for each referral, this would suggest a total of 325 hours of lost instructional time each year and a considerable administrative staffing burden.⁹ The average per-day number of classroom-originating referrals that respond to behavior classified as “subjective” in nature was 1.29 referrals per 500 students. The average per-day number of referrals from classrooms that are “objective” behavioral infractions was 0.48 per 500 students.

There is considerable cross-school variability in the rate of referrals. The full-sample standard deviation for classroom-based referrals is 2.83 (1.77 for non-classroom referrals). In many schools, students miss substantial portions of the school year due to disciplinary events.

We have information on whether PBIS implementation was successful for slightly more than half of our sample. 61.5 percent of our grade-school-year observations are located in schools that conducted assessments of their implementation of PBIS practices. In 71 percent of these school-year observations, schools were assessed to be successfully implementing PBIS.

In [Figure 1](#), we map the differential timing, by state, of the implementation of higher-stakes evaluation policies as well as the six states that never enacted new policies. The majority of schools in our sample experience high stakes evaluation, though the overrepresentation of California schools in our data means that while 44 of 50 states experienced high-stakes teacher evaluation policies, only 72 percent of schools in our sample end up operating under a higher-stakes evaluation framework.

In [Figure 2](#), we plot the raw outcome data for schools in states that experienced high-stakes evaluation reform for the years before and after policy implementation. In Panel A, we present trends for ODRs that originate in the classroom, our first outcome of interest. We observe no evident discontinuity coinciding with the enactment of increased accountability measures for teachers or any change in the slope of these raw averages. We observe the same patterns in Panel B of [Figure 2](#) which displays classroom-originating referrals for infractions that are subjective in nature, our second outcome of interest. Thus, the descriptive evidence suggests that there may be limited effects of the implementation of higher-stakes evaluation on teachers’ disciplinary decisions; however, secular patterns may mask an underlying causal relationship. This motivates our identification strategy which we discuss in more detail next.

⁹While this number is a rough approximation given the lack of precise estimates, we believe it is a conservative one. Students are typically required to complete a reflection form and conference with a school administrator before returning to class. Students unprepared to return to class remain with the administrator for longer periods. If one administrator were responsible for all referrals in a 500 student school (a reality in many contexts), this would mean that 22.5 percent of her 8-hour work days over 180 school days would be devoted to these referrals.

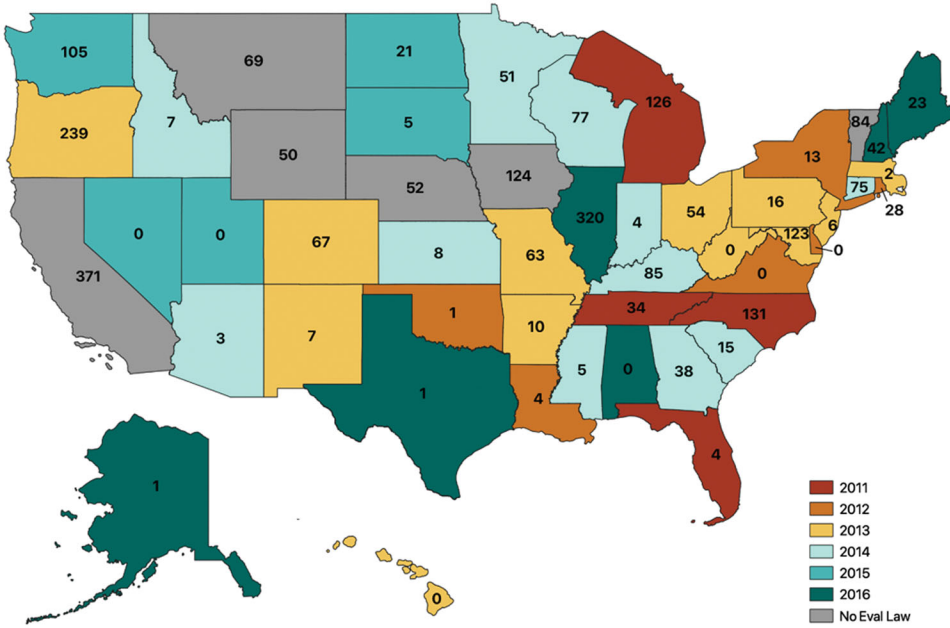


Figure 1. The timing of statewide teacher evaluation reforms and number of schools by state in analytic sample. *Notes:* Years represent the fall of the academic year in which new evaluation systems were fully implemented statewide. Numbers inside each state represent total schools in analytic sample ($n = 2,564$). Full list of states with schools in sample and timing of evaluation in [Supplementary Appendix Table A1](#).

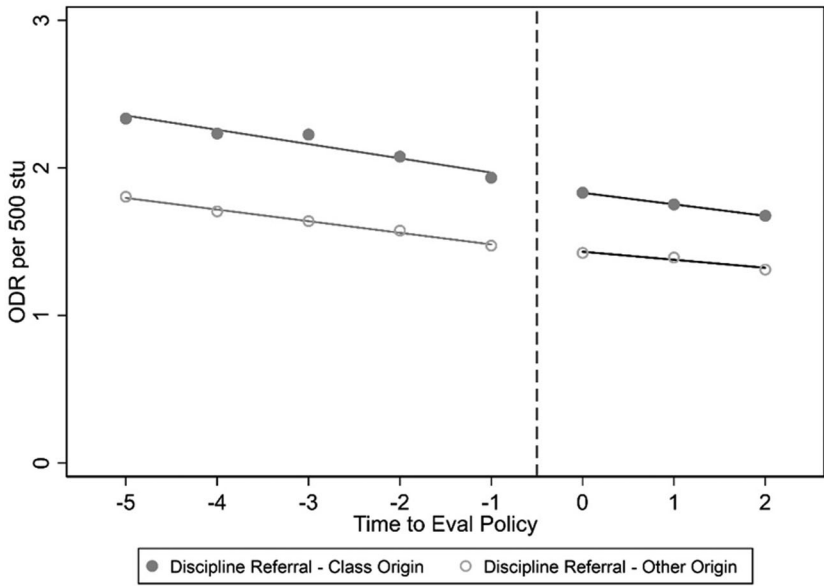
Analytic Plan

Our identification strategy relies on a two-way fixed effect difference-in-differences (DD) approach. Our first difference is the change in the rate of ODRs that may have been influenced by the change in evaluation policy for schools in states that experienced the teacher evaluation policy reform. Our second difference is the change in the rate of these ODRs for schools in states that had not yet experienced (or did not experience) the change. We present several complementary estimation strategies below, though each is a variant of our difference-in-differences research design framework.

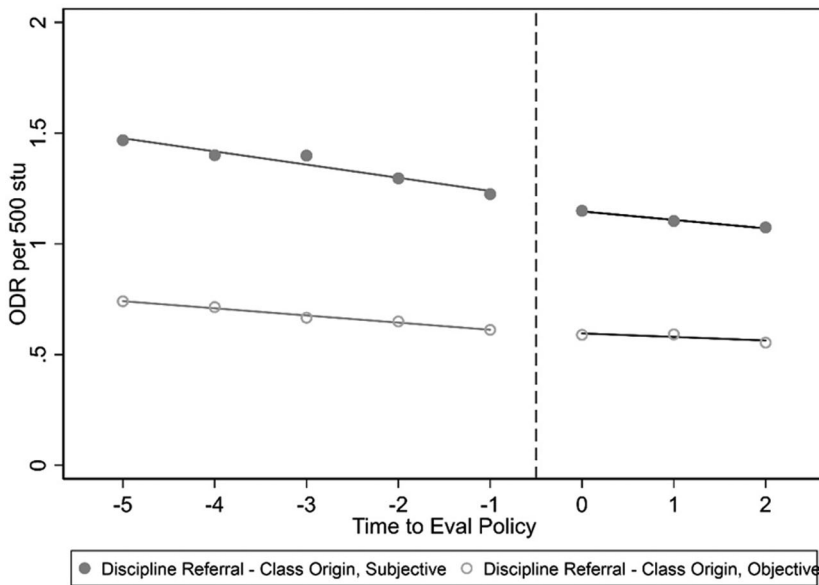
We begin by estimating a non-parametric event study. This approach allows us to flexibly estimate any pre-policy trends or time-varying treatment effects. We fit the following model:

$$ODR_{gjt} = \sum_{r=-6}^{3+} 1(t = t_s^* + r)\beta_r + (\mathbf{X}_{jt})\theta + \Delta_g + \Gamma_j + \Pi_t + \varepsilon_{gjt} \quad (1)$$

In simplified form, this represents the per-500-student per-day rate of Office Disciplinary Referrals (ODR_{gjt}) for each grade-year observation in grade g , school j , state s and time t , regressed on a series of indicators that take the value of 1 when the grade-year observation is a given number of years pre- and post-policy reform, with t_s^* indicating the year in which states implemented the higher-stakes teacher evaluation reform. This model includes grade- (Δ), school- (Γ) and year- (Π) fixed effects and a vector (\mathbf{X}) of school-level (j) background characteristics. We include a



Panel A. Location of ODR



Panel B. Objectivity of ODR

Figure 2. Average Office Disciplinary Referral rates before and after the introduction of teacher evaluation policies, by location of referral (Panel A) and objectivity of infraction (Panel B). *Notes:* Points represent weighted average of office disciplinary referrals for schools in states that implemented high-stakes evaluation. Line is best fit for those averages.

parsimonious set of plausibly exogenous school characteristic adjustments to capture school-specific characteristics and improve the precision of our estimates. As policy affects every school in a state, we do not anticipate that evaluation reforms would alter the demographic composition of a school.¹⁰ We include the following demographic characteristics measured at the school level: percent of students receiving free/reduced lunch, percent of students of various racial/ethnic backgrounds, and school enrollment.

We estimate our primary models using Weighted Least Squares in which we weight each observation by the grade-level enrollment. OLS estimates will return a heteroskedastic error term because estimates of ODRs will be known with more precision in grades (and schools and states) with a larger enrollment. Weighting our observations allows us to interpret our estimates as the effect of teacher evaluation on the rate of ODRs in the average-sized grade. The level of policy intervention is the state-year, but our errors are likely correlated across time within state. Thus, in our main estimates, we cluster standard errors at the state level, but we also test the robustness of our inferences to clustering at the state-year level and to estimating a two-way clustered standard error structure (Abadie et al., 2017; Bertrand et al., 2004; Cameron et al., 2011).

The coefficients of interest are the seven β_r ($-5 \leq r \leq 1$) which represent the effect of evaluation reforms on rates of ODRs r years before and after the policy introduction.¹¹ We measure all effects compared to the year prior to the reform ($r = -1$), and we assign all non-treated schools to the same pre-policy year ($r = -1$). In practice, this means that all observations of non-treated schools contribute to the reference category against which we measure treated schools' pre- and post-policy implementation ODR rates.

Next, we extend Equation 1 into the pre/post difference-in-differences framework where we pool estimates across years to increase precision and test the global effects of the policy:

$$ODR_{gkst} = \beta_1 EVAL_{st} + (\mathbf{X}_{jt})\theta + \Delta_g + \Gamma_j + \Pi_t + \mu_{gkst} \tag{2}$$

The indicator $EVAL_{st}$ takes the value of 1 if the school is in a state that is in a year with a high-stakes evaluation system. β_1 is the causal parameter of interest. All other terms are defined as above.

We also relax the assumption of the standard difference-in-differences model of time-invariant treatment effects in Equation (2) by adding a linear time trend:

$$ODR_{gkst} = \beta_1 EVAL_{st} + \beta_2 EVAL \times YEAR_{st} + \beta_3 YEAR_{st} + (\mathbf{X}_{jt})\theta + \Delta_g + \Gamma_j + \Pi_t + \nu_{gkst} \tag{3}$$

¹⁰In fact, we regress the seven school demographic characteristics on our evaluation indicator and reject the null in only one instance. Evaluation implementation predicts a small decrease in the FRPL composition of a school (Beta: 1.25 p.p., SE: 0.52). Given the multiple hypotheses we test and the small magnitude of the coefficient, we take these results as consistent with our claim that school demographic characteristics are exogenous to policy implementation, though we present estimates without these adjustments in all cases to address this concern.

¹¹In our event-study results, we estimate coefficients for all available data (including binned categories for years 6+ pre-, 2 years post, and 3+ years post evaluation) but only interpret years 5 through 1 to ensure that we only compare units that are observable for all treatment timing years. Including estimates outside the 5 to +1 bandwidth mixes relative-time effects with compositional shifts to schools that we are and are not able to observe for these years. In our main difference-in-differences estimates, however, we pool pre- and post-treatment periods to take advantage of the full range of data which allows us to include schools that implement new evaluation policies for up to seven years.

where $YEAR_{st}$ is a linear time trend for states, centered around the year the state implemented the high-stakes teacher evaluation policy.¹² The interaction term $EVAlxYEAR_{st}$ allows for the relative time trends among schools in treated states to differ post-reform. The coefficient on the main effect of treatment (β_1) identifies the immediate response of the introduction of high-stakes evaluation of ODRs and the coefficient on the interaction term (β_2) captures linear deviations from the average effect. The time-trend coefficient (β_3) tests for any differential trends in the pre-reform period among states that introduced high-stakes evaluation systems.¹³

To better understand the effects of accountability pressures, we examine differences between grade levels under greater and lesser accountability and differences across schools with varying rates of referrals prior to the higher-stakes evaluation era. Specifically, we hypothesize that grades 3–11 will be subject to greater levels of accountability as these are years in which high-stakes testing occurs in schools. Therefore, teachers of tested subjects in these grades experience accountability from both evaluative classroom observations and student test score outcomes. We theorize that ODRs from grades K–2 and 12, will be less sensitive to the introduction of high-stakes teacher evaluation. We note here the pre-registration (Registry of Efficacy and Effectiveness Studies in Education 1748.2) of our analytic plan in which we propose to explore the higher accountability applied to grades 3–11.¹⁴ Additionally, we explore whether teachers respond differently to increased accountability when the starting behavioral climate in

¹²While this approach as well as the event study models are similar in spirit to the Comparative Interrupted Time Series (CITS) design, there are several important distinctions between the CITS and the difference-in-differences with multiple time points approaches. The CITS approach models different pre-trends for treatment and counterfactual groups. The DD approach assumes (and then tests for) parallel trends between treatment and counterfactual groups. This stricter assumption about pre-trends means that the DD approach need not rely on linear (or higher-order) extrapolations from pre-trends to estimate intercept- and slope-shifts. Given that our setting meets the stricter parallel trends assumptions of the differences-in-differences approach, we fit all models using this approach.

¹³We may be concerned that the estimates from Equations (2) and (3) will be biased as a result of unobserved state-level factors that, contemporaneous with the introduction of high-stakes teacher evaluation, also affect ODRs. Triple difference (DDD) estimates that leverage alternative, potentially unaffected, outcomes help us address these sources of bias. We model these as follows:

$$ODR_{gijst} \beta_1 EVAlxAFFECT_{st} + \beta_2 EVAl_{st} + \beta_3 AFFECT_{st} + (AFFECT_{st} \cdot \Gamma_j) \phi + (AFFECT_{st} \cdot \Pi_t) \delta + (X_{jt}) \theta + \Delta_g + \Gamma_j + \Pi_t + u_{gijst} \cdot AFFECT_{st}$$

is an indicator variable that takes the value of one if the observation is one in which we would anticipate the introduction of high-stakes evaluation policies will affect the rate of ODRs or affect the rate more intensively. We contrast locations in which ODRs occur, specifically comparing classroom-originating ODRs, which we anticipate would be influenced by changes in the teacher evaluation policies, and non-classroom-originating ODRs, which we anticipate would not be affected by the policy changes. Alternatively, we contrast the type of infraction (subjective or objective) resulting in the ODR. β_j represents the effect of the introduction of the high-stakes evaluation policy on anticipated affected outcomes, compared to unaffected outcomes in states that had not yet or never adopted the evaluation policy. We adjust for unexplained within-school and within-year heterogeneity in affected outcomes by interacting our $AFFECT$ indicator with year – and school-indicators. As we show below, we find null effects for all of our double difference models, and so we do not feature our triple difference framework prominently. We do present these results in Tables A13 and A14 and, as expected, they also return precise zeros.

¹⁴We do not present K–2 and 12 as a counterfactual of *no* accountability increases, but we would anticipate that any effects in these grades would be less intense than tested grades. We recognize that not all teachers in grades 3–11 teach a tested subject, but in the presence of an effect of greater test-score-based accountability pressures in these grades, we would nevertheless expect to see an average treatment effect in these grade bands. While all schools in states in our sample require high-stakes assessments in grades 3–8, high-school assessment requirements vary. All states require students to test at some point in grades 9–11. Some states require testing only in one of these grades, other states require testing across multiple years, still others allow student discretion on the grade in which students take tests. Our estimates are even closer to zero when we restrict our definition of higher-accountability grades to 3–8 (class: 0.026 (0.067); subjective-class: 0.022 (0.067)).

their school differs by interacting ODR-rates in the year prior to the implementation of evaluation policies with our policy indicator. We thank several readers of early drafts of our paper for highlighting this important possible source of heterogeneous response, and we acknowledge that this analysis was not part of our pre-registration plan so should be considered as exploratory.

We also examine the extent to which the implementation of effective disciplinary support strategies serves to moderate the effects of greater accountability. Specifically, what are the effects on ODRs when schools develop better systems of behavioral supports in the context of higher-stakes evaluation? The fixed-effects structure of our analysis means that our estimates are of the effect of within-school improvements in the implementation of PBIS and the interaction of these improvements with the introduction of high-stakes evaluation. However, we note explicitly the exploratory nature of this analysis as the successful implementation of PBIS is clearly an endogenous characteristic of the school. A negative coefficient on the interaction term would suggest that the combination of teacher evaluation reforms and well-implemented behavioral support systems produce improvements in teacher skills. A null effect on the interaction with a larger-signed estimate on the main *Implement Evaluation* term would imply that schools with robust student support systems are less responsive to added external accountability pressures. Finally, a positively signed coefficient on the *Implement Evaluation* term with an opposite signed coefficient on the interaction would suggest that, in settings with ill-defined criteria for office referrals, increased accountability pressures caused teachers to remove more students from class.

For β_1 to be an unbiased estimator in the above models, we make three assumptions about our research design: (1) schools and grades in untreated states (and not-yet-treated states) provide a valid counterfactual for schools and grades in treated (or already-treated) states; (2) there are no unobserved simultaneous shocks correlated with our outcomes and the introduction of higher-stakes teacher evaluation reforms; and (3) the estimates for each grade and year are appropriately pooled to create the full sample Average Treatment Effect on the Treated (ATT) estimate.

In our full sample, our first assumption depends on the fact that schools in states that did not or had-not-yet adopted evaluation reform were equal in expectation to schools that had. We present information in [Table 1](#) on the differences in schools and students which did and did not experience high-stakes evaluation. The most notable difference is the larger school size for schools that ever experienced high-stakes evaluation. While there are some baseline differences in the characteristics of schools and students in states that did and did not experience high-stakes evaluation, we account for this in our difference-in-differences estimation framework. As we compare the difference in values before and after the policy change with differences in values over the same time period in locales that did not or had not yet experienced the policy change, starting differences between treated and untreated locales does not threaten the validity of our design. Nevertheless, for external validity purposes, it is reassuring that the outcome values prior to the start of the era of evaluation reform are quite close in schools located in states that did and did not experience high-stakes teacher evaluation, all within 0.2 referrals per-500 students per day. We more formally test for violations of the first assumption by examining whether schools in states that implemented evaluation were

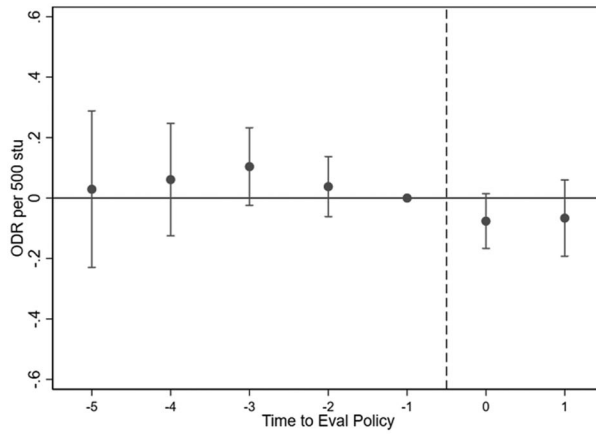
on parallel trends prior to evaluation implementation as those that did not (or had not yet) implemented evaluation.

In some models, we restrict our sample to only those observations in states that ever enacted teacher evaluation. In these specifications, our estimates rely only on variation in the timing of the implementation of teacher evaluation reforms, and we need only assume that *when* (but not *whether*) states implemented evaluation was as-good-as random. State legislatures enacted policy reforms in response to a common, exogenous federal shock (RTTT), but the date of mandatory, statewide adoption of the new evaluation systems differed. Cross-state differences in the length of time between passing and implementing evaluation reforms result from plausibly exogenous factors such as state legislative negotiating processes, the length of state legislative sessions and the duration of pilot phases (NCTQ, 2017). In [Supplementary Appendix Table A2](#), we present evidence that highlights that, on both observable demographic characteristics and office referral rates, there are no evident patterns distinguishing schools that implemented evaluation earlier as opposed to later in our analytic timeframe. We cannot, however, rule out all potential endogenous reasons for variation in the timing of implementation.

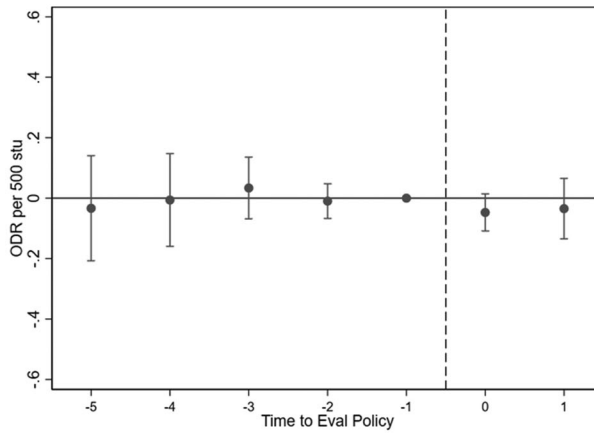
If other policy reforms are contemporaneous with evaluation policy reforms, our estimates of the effects of evaluation reform could be biased. Thus, we test our second assumption with a set of placebo models in which we estimate the effects of teacher evaluation on alternate outcomes that measure disciplinary events from outside the classroom (Other) or that are in response to behavioral events so severe that we do not anticipate they would be affected by evaluation policy reform alone (Objective). [Figure 2](#) provides suggestive evidence that the trends in these placebo outcomes are unresponsive to policy implementation, and we test this formally below. While these are not perfectly clean falsification tests as it is possible that evaluation reforms could affect overall school climate and therefore these outcomes as well, we expect they will be less sensitive to the reforms than our primary outcomes. In addition to using placebo outcomes, we also create fictional dates for evaluation policy implementation that precede the actual years of adoption and test whether these false evaluation years predict changes in referral rates. As a final test of the second assumption, we fit a series of models in which we include other teacher accountability and discipline policy reforms to assess whether they predict changes in ODR rates or moderate the main effect of teacher accountability pressures.¹⁵

There has been a recent explosion in the econometric literature documenting frequent failures of the third assumption, particularly in the context of differential timing as an identification strategy (e.g., Athey & Imbens, 2018; Borusyak & Jaravel, 2017; Callaway & Sant'Anna, 2020; de Chaisemartin & D'Haultfoeuille, 2020; Ferman & Pinto, 2019; Gibbons et al., 2018; Goodman-Bacon, 2021; Imai & Kim, 2020). We test our fixed effects for the presence of time-varying weights and negative weighting, and we replicate

¹⁵Kraft et al. (2020) seek to rule out threats to their identification strategy from contemporaneous teacher and accountability policy reforms, such as the implementation of Common Core Standards or licensure tests. These are less relevant to our identification strategy as we are ultimately interested in whether and how increased accountability shifts teachers' classroom practices. Our results are robust to the inclusion of policy indicators for the reform of tenure laws and weakening of collective bargaining (see [Supplementary Appendix Tables A15 and A16](#)). To the extent that our estimates of teacher evaluation reforms are influenced by other accountability-related policy reforms, this would imply that our results are evidence of overall accountability pressures on teacher practice, rather than specific to teacher evaluation.



Panel A. Classroom ODRs



Panel B. Subjective Classroom ODRs

Figure 3. Non-parametric event study displaying effect of high-stakes teacher evaluation reforms on rate of per-500-student, per-day Office Disciplinary Referrals (ODRs), by location and subjectivity. *Notes:* Point estimates for years pre- and post-evaluation reforms and corresponding 95 percent confidence intervals derived from event study model describe in Equation 1 that is weighted by grade enrollment, includes grade, school and year fixed effects and time-varying school characteristics, with standard errors clustered at state level. Full coefficients reported in Columns IIa and IIc of Supplementary Appendix Table A3.

our standard two-way fixed effect models using de Chaisemartin and D’Haultfoeuille’s (2018) time-corrected Wald (Wald-TC) estimator.

Results

Event-Study Estimates

We find no evidence that rates of Office Disciplinary Referrals (ODRs) changed in the aftermath of the introduction of higher-stakes teacher evaluation policies. In Figure 3, we present results from Equation 1 for both referrals originating in the classroom (Panel A) and referrals originating in the classroom that are subjective in nature (Panel

B). While there is some visual evidence that ODRs decline somewhat after the introduction of high-stakes evaluation policies, all estimates fall within the 95 percent confidence interval and are small in magnitude.

Independently, the results in [Figure 3](#) provide little evidence that there were trends in the rate of classroom-based or subjective referrals prior to the implementation of teacher evaluation policies. This suggests the first assumption required of our DD estimates holds and we have no need to instrument with leads in our event study (Freyaldenhoven et al., 2019).

Difference-in-Differences Estimates

Results from our main difference-in-differences estimates confirm that we find no causal effect of the implementation of high-stakes evaluation on rates of classroom- or classroom-subjective ODRs. In [Table 2](#), we present the results of [Equations \(2\) and \(3\)](#). While our estimates are consistently signed and of nearly identical magnitude, in all cases we fail to reject the null. Due to their parsimony and consistency with models that include additional controls, we select Models I and IV as our preferred estimates for classroom- and classroom-subjective referrals, respectively because the absence of additional controls requires us to make fewer assumptions about parallel pre-trends across the covariates' sub-strata (e.g., Callaway & Sant'Anna, 2021).¹⁶

We estimate these effects with precise zeros. In our preferred estimates, we can confidently rule out ranges of effects greater than a decrease of 0.21 referrals or an increase of 0.04 referrals per-500 students, per day for classroom referrals and a decrease of 0.14 or an increase of 0.06 referrals per-500 students, per day for subjective-classroom referrals. These confidence intervals correspond to a 0.07 standard deviations (*SD*) decrease and a 0.01 *SDs* increase or a 0.07 *SDs* decrease and a 0.03 *SDs* increase for classroom and subjective referrals, respectively.¹⁷ We find no evidence of differential post-evaluation policy implementation trends in Models III or VI.

Heterogeneity of Effects

Strong Systems of Behavioral Supports

We find no evidence that improvements in schools' implementation of Positive Behavioral Interventions and Supports (PBIS) practices serves to moderate the effects of accountability. We present results in [Table 3](#) of a series of estimates in the 61.5 percent of grade-school-year observations for which we have measures of PBIS implementation. In Models I and V we first re-estimate the results from [Table 2](#) and examine the main effect of evaluation implementation in this sub-sample of observations. For these schools, the effects are even closer to zero. We then introduce the time-varying effect of

¹⁶This represents a departure from our pre-registered preferred estimates which included demographic covariates in our model specification to improve precision and adjust for any remaining observable bias in our models. The new methodological insights post-date our pre-registration. As we present the pre-registered results (Models II and V) alongside these preferred results and the coefficients differ by only 0.005 and 0.001 of a referral per-500 students per day, we believe adopting these as our preferred results is fully consistent with our pre-registration plan.

¹⁷We scale the precision of these null effects to the standard deviation of our outcomes across the full analytic sample. When we scale our outcome to the *within-school* standard deviation of our outcomes, our 95 percent confidence intervals are 0.12 to +0.02 *SD* and 0.11 to +0.04 *SD* units for the main effects of evaluation on classroom and subjective referrals, respectively.

Table 2. The effect of teacher evaluation reforms on office disciplinary referrals, by location and subjectivity.

	A. Class			B. Subjective		
	I	II	III	IV	V	VI
Implement evaluation	-0.084 (0.063)	-0.089 (0.063)	-0.083 (0.072)	-0.041 (0.049)	-0.042 (0.050)	-0.054 (0.043)
Implement evaluation * Trend			0.046 (0.043)			0.007 (0.024)
Time trend			-0.017 (0.032)			0.004 (0.022)
School composition controls		X	X		X	X
Impact estimate (SDs)	-0.030	-0.031	-0.029	-0.021	-0.022	-0.028
[95% C.I. (SDs)]	[-0.073, 0.014]	[-0.075, 0.012]	[-0.079, 0.021]	[-0.070, 0.028]	[-0.072, 0.029]	[-0.071, 0.016]
Grade-year observations (N)	107,468	107,468	107,468	107,468	107,468	107,468
School-year observations	20,137	20,137	20,137	20,137	20,137	20,137
R-squared	0.559	0.559	0.559	0.550	0.550	0.550
Sample outcome mean [SD]		2.03 [2.83]			1.29 [1.95]	

Notes: * $p < .05$, ** $p < .01$, *** $p < .001$. Cells report estimates and associated standard errors clustered at the state level in parentheses. School controls include time-varying enrollment, proportion low-income and race/ethnicity. All models include grade, school and year fixed-effects and are weighted by grade enrollment.

Table 3. The moderating effect of positive behavioral interventions and supports (PBIS) on the effect of teacher evaluation reforms on office disciplinary referrals, by location and subjectivity.

	A. Classroom				B. Subjective			
	I	II	III	IV	V	VI	VII	VIII
Implement evaluation	-0.045 (0.071)	-0.083 (0.108)	-0.089 (0.109)	-0.244 (0.197)	-0.054 (0.067)	-0.075 (0.083)	-0.078 (0.084)	-0.193 (0.108)
Implement PBIS well		-0.116 (0.063)	-0.116 (0.063)	-0.105 (0.063)		-0.086 (0.045)	-0.086 (0.045)	-0.081 (0.045)
Implement evaluation * PBIS		0.036 (0.097)	0.037 (0.097)	0.171 (0.188)		0.018 (0.059)	0.019 (0.059)	0.116 (0.100)
Implement evaluation * Trend				0.139 (0.090)				0.064 (0.049)
Implement evaluation * Trend * PBIS				-0.087 (0.084)				-0.064 (0.044)
Time trend				0.002 (0.035)				0.016 (0.033)
School composition controls			X	X			X	X
Grade-year observations (<i>N</i>)	66,075	66,075	66,075	66,075	66,075	66,075	66,075	66,076
School-year observations	12,309	12,309	12,309	12,309	12,309	12,309	12,309	12,309
<i>R</i> -squared	0.602	0.602	0.602	0.602	0.584	0.584	0.584	0.584

Notes: * $p < .05$, ** $p < .01$, *** $p < .001$. Cells report estimates and associated standard errors clustered at the state level in parentheses. School controls include time-varying enrollment, proportion low-income and race/ethnicity. All models include grade, school and year fixed-effects and are weighted by grade enrollment. Models I and V replicate results from main DD estimate on PBIS implementation sub-sample. Fewer observations reflect subset of grade-school-year observations (61.5 percent) reporting PBIS implementation information. For context, the sample means (standard deviations) of the outcomes are Classroom ODRs: 2.03 (2.83) and Classroom-Subjective ODRs: 1.29 (1.95).

PBIS implementation and its interaction with teacher evaluation. In our most-parsimonious, preferred specifications (Models II and VI), we can confidently rule out ranges of moderating effects greater than a decrease of 0.15 referrals or an increase of 0.23 referrals per-500 students, per day for classroom referrals and a decrease of 0.10 or an increase of 0.13 referrals per-500 students, per day for subjective-classroom referrals. These correspond to 95 percent confidence intervals of -0.05 SDs to $+0.08$ SDs for classroom referrals and -0.05 SDs to $+0.07$ SDs for subjective referrals.¹⁸ We observe no post-evaluation implementation time trends.

Pre-Policy Referral Rates

We do not find any evidence of heterogeneity of effects by the rates of disciplinary referrals in the year prior to evaluation policy implementation. In Table 4, we present results in Models 1 and V in which we re-estimate our primary models on the sub-sample of grade-school-year observations in states that ever experienced evaluation and that are not observed in the year immediately prior to policy implementation ($t=-1$). We exclude this year so as not to interact our policy predictor with a value that is on both the left- and right-hand sides of our estimating equations. These results are consistent with our main estimates. In Models II–IV and VI–VIII, both the main effect of evaluation implementation and its interaction with the pre-policy rate of referral are indistinguishable from zero.¹⁹

¹⁸When we scale our outcome to the within-school standard deviation, the 95 percent confidence interval on the moderating effects of PBIS are 0.09 to $+0.13$ and -0.08 to $+0.11$ SD units for classroom and subjective referrals.

¹⁹We similarly find no effects on quadratic terms for pre-policy referral rates (class: 0.000 (0.004); subjective: 0.000 (0.006)) and in models where we average referral rates from the two years prior to policy implementation and then leave these two years out (class: 0.011 (0.044); subjective: 0.019 (0.038)).

Table 4. The effect of teacher evaluation reforms on Office Disciplinary Referrals, by pre-evaluation implementation referral rates.

	A. Classroom			B. Subjective		
	I	II	III	IV	V	VI
Implement evaluation	-0.100 (0.082)	-0.120 (0.088)	-0.123 (0.088)	-0.064 (0.050)	-0.101 (0.054)	-0.102 (0.053)
Implement evaluation * ODR _r = -1		0.011 (0.041)	0.012 (0.041)		0.031 (0.032)	0.032 (0.032)
School composition controls			X			X
Grade-year observations (N)	74,452	74,452	74,452	74,452	74,452	74,452
School-year observations	14,175	14,175	14,175	14,175	14,175	14,175
R-squared	0.546	0.546	0.547	0.523	0.523	0.523

Notes: * $p < .05$, ** $p < .01$, *** $p < .001$. Cells report estimates and associated standard errors clustered at the state level in parentheses. School controls include time-varying enrollment, proportion low-income and race/ethnicity. All models include grade, school and year fixed-effects, are weighted by grade enrollment, and include only states ever-under higher-stakes evaluation law. Models I and V re-estimate main effects of evaluation from Table 2 on sample with all observations in year prior to evaluation implementation omitted. Models II–III and V–VI estimate effects of evaluation interacted with rate in year prior to evaluation implementation. For context, the sample means (standard deviations) of the outcomes are Classroom ODRs: 2.03 (2.83) and Classroom-Subjective ODRs: 1.29 (1.95).

Table 5. The effect of teacher evaluation reforms on Office Disciplinary Referrals, by grade-level accountability pressures, location and subjectivity.

	A. Class (3–11 only)			B. Subjective (3–11 only)		
	I	II	III	IV	V	VI
Implement evaluation	-0.092 (0.066)	-0.098 (0.067)	-0.096 (0.076)	-0.052 (0.056)	-0.054 (0.058)	-0.075 (0.045)
Implement evaluation * Trend			0.054 (0.047)			0.008 (0.027)
Time trend			-0.018 (0.033)			0.008 (0.025)
School composition controls		X	X		X	X
Grade-year observations (N)	64,437	64,437	64,437	64,437	64,437	64,437
School-year observations	19,630	19,630	19,630	19,630	19,630	19,630
R-squared	0.586	0.586	0.586	0.573	0.573	0.573

Notes: * $p < .05$, ** $p < .01$, *** $p < .001$. Cells report estimates and associated standard errors clustered at the state level in parentheses. School controls include time-varying enrollment, proportion low-income and race/ethnicity. All models include grade, school and year fixed-effects and are weighted by grade enrollment. For context, the sample means (standard deviations) of the outcomes are Classroom ODRs: 2.03 (2.83) and Classroom-Subjective ODRs: 1.29 (1.95).

Intensity of Accountability Pressures

We find no evidence of heterogeneous effects for grades that should be subject to more intensive accountability pressures. In Table 5, we present results in which we restrict our grade-school-year observations to those that represent grades 3 through 11. We again fail to reject the null hypothesis and can confidently rule out small effects, both using effect size metrics and substantive interpretations. We present the corresponding event study estimates for high-accountability grades in Supplementary Appendix Figure A2 and Table A4. We present analogous results for the moderating effects of successful PBIS implementation in higher-accountability grades (3–11) in Supplementary Appendix Table A5.

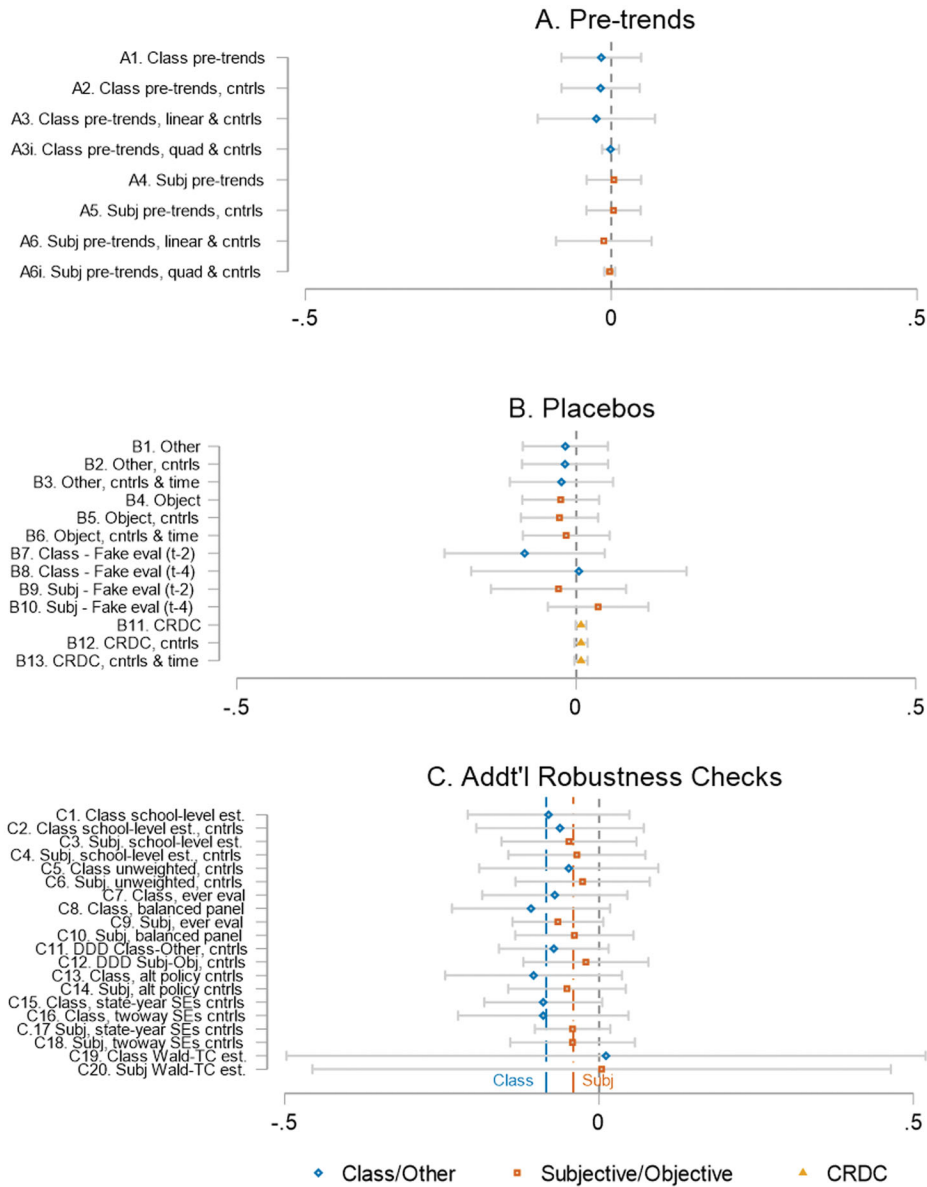


Figure 4. Point estimates and 95 percent confidence intervals of tests of assumptions on main difference-in-differences analysis (grades K-12). Notes: To meet difference-in-differences assumptions, all 95 confidence intervals should overlap with zero in Panels A and B. Estimates in Panel C should overlap with preferred DD model point estimates for classroom (dash: -0.084) and classroom, subjective ODRs (dash-point: -0.041). As these are indistinguishable from zero in the population, all estimates in Panel C should also overlap with zero. Data for estimates B11–B13 from Civil Rights Data Collection. Full set of point estimates available in [Supplementary Appendix Tables A6, A8, A9, A12, A13, A15, A17 and A18](#).

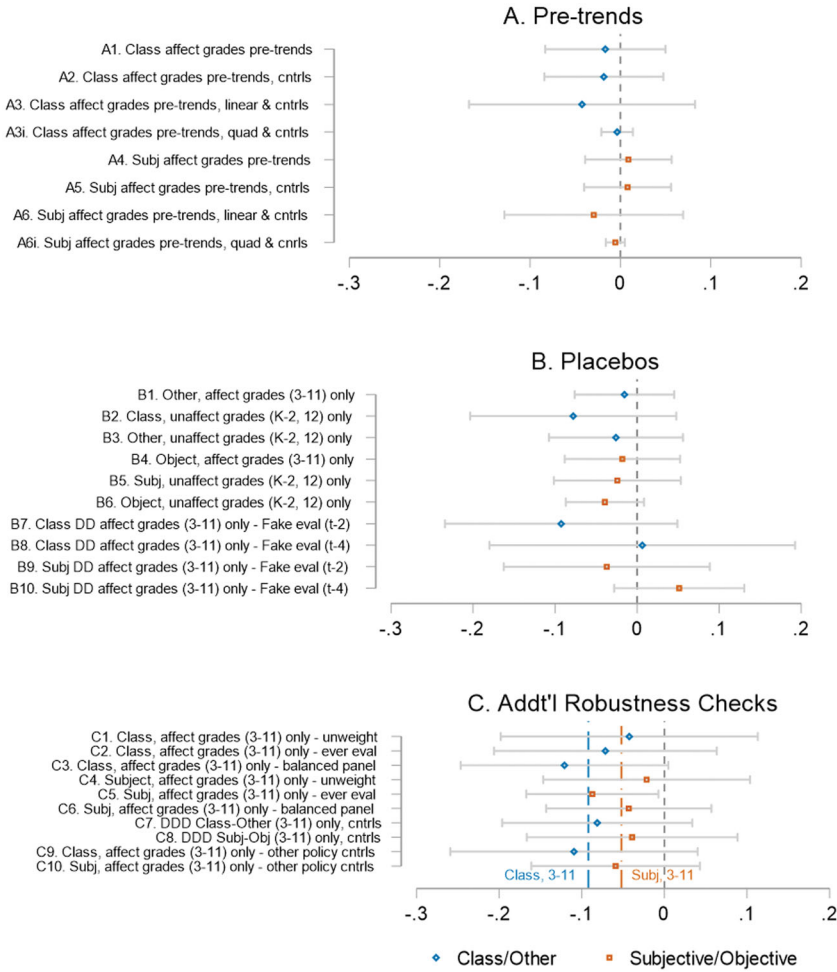


Figure 5. Point estimates and 95 percent confidence intervals of tests of assumptions on high-accountability grade difference-in-differences analysis (grades 3–11). *Notes:* To meet difference-in-differences assumptions, all 95 confidence intervals should overlap with zero in Panels A and B. Estimates in Panel C should overlap with preferred DD model point estimates for classroom (dash: -0.092) and classroom, subjective ODRs (dash-point: -0.052). As these are indistinguishable from zero in the population, all estimates in Panel C should also overlap with zero. Full set of point estimates available in [Supplementary Appendix Tables A7, A10, A11, A14, A16, and A19](#).

Assumption Checks and Sensitivity Analyses

The three central assumptions of our identification strategy hold across multiple tests. Given the extensive robustness checks we conduct, for the purpose of parsimony we display the relevant results from these checks in [Figures 4 and 5](#) in the main text of the paper and display the full set of coefficients and statistics in [Supplementary Appendix A](#).

Pre-Trends

Our tests of the parallel trends assumption reinforce the graphical evidence from [Figures 2 and 3](#) that schools in states that did not implement higher-stakes evaluation (or had not yet) provide valid counterfactuals. We present results of our tests of parallel pre-trends in Panel A of [Figures 4 and 5](#) (corresponding to [Supplementary Appendix Tables A6 and A7](#)). If the assumption holds, these coefficients should be indistinguishable from zero, which in all cases is true. Note that the pre-trend coefficients without covariate adjustments are also indistinguishable from zero (Estimates A1 and A4). As the unconditional parallel trends assumption is met, our models are robust to concerns raised by Callaway and Sant'Anna (2021) about inaccurate treatment effect estimates in the presence of heterogeneous treatment effects and covariate-specific time trends. Estimates A3i and A6i present coefficients on quadratic time trends and are also indistinguishable from zero.

Placebo Tests

We find no evidence that the introduction of higher-stakes evaluation affects outcomes that we do not anticipate these reforms would influence. Further, when we adjust the date of policy implementation to create falsification tests, we find no evidence that these artificial policy implementation dates influenced referral rates. We present these results in Panel B of [Figures 4 and 5](#) (corresponding to [Supplementary Appendix Tables A8–A11](#)). Estimates of the effect of evaluation policy on unaffected outcomes should be indistinguishable from zero. Similarly, estimates of the effect of placebo evaluation dates in years before the policy was actually implemented should also be zero (or at least substantially attenuated in pooled pre- and post- tests).

We present Estimates B1–B6 in [Figure 4](#) as evidence that evaluation policy implementation had no effect on ODRs from locations other than the classroom and on ODRs for behavioral infractions that involved objective reasons to send students to the office. In estimates B7–B10 we demonstrate that using a date of evaluation implementation two or four years before the actual implementation date is not predictive of changes in ODR rates.

In Estimates B1–B6 in [Figure 5](#), we present the corresponding placebo outcome tests in our analysis of potential heterogeneous effects in higher-accountability grades. Here, we use both our main outcomes (classroom and objective ODRs) in the theoretically “unaffected” grades (K-2, 12) as well as the secondary outcomes (other location and subjective ODRs) in the theoretically “affected” grades (3–11). In Estimates B7–B10, we present analogous falsification tests where we use a date two or four years prior to the actual policy implementation and examine the effects on just the high-accountability grades. Again, all estimates are indistinguishable from zero.

In Estimates B11–B13 of [Figure 4](#), we present results in which we test the effect of evaluation policy implementation on rates of suspension from the Civil Rights Data Collection (CRDC) sample. In this national sample of schools, we note that the introduction of evaluation policies yields a small and non-significant positive coefficient of 0.4 to 0.7 percentage points in the proportion of suspended students. Thus, we interpret the CRDC findings as consistent with our other falsification tests that there were no endogenous shifts in state discipline policy.

Alternate Sample, Policy, Specification and Weighting Approaches

In our last set of robustness checks, we present further evidence that the schools in untreated states provide valid counterfactuals, that our results are not driven by concurrent policies, and that our results are robust to the method of weighting individual fixed effect ATEs into a pooled estimate. In Panel C of [Figures 4 and 5](#) (corresponding to [Supplementary Appendix Tables A12–A19](#)), we present these results. In addition to the zero line, we also include, for reference, the point estimates from our preferred models (Models I and IV in [Tables 2 and 5](#)). Results of these robustness checks should overlap with the main results. Given that our main models find that teacher evaluation policy reform does not change the rate of ODRs, they should also overlap with zero.

We first test the robustness of our results to a slight expansion of our main analytic sample such that we can generalize our findings to a larger set of schools. Our primary sample is comprised of 107,468 grade-school-year observations that represent 20,137 school-year observations. However, our data includes outcomes reported only at the school level for an additional 268 school-year observations. We present results from re-estimating [Equations \(2\) and \(3\)](#) using data aggregated at the school level in Estimates C1–C4 in [Figure 4](#). The estimates are essentially identical to our grade-level models.

Using alternate weighting approaches, our estimates also generalize to the average school (rather than the average-sized school) in our sample. In Estimates C5 and C6 in [Figure 4](#) and C1 and C4 of [Figure 5](#), we present results from models with standard OLS estimators and again fail to reject the null with respect to the main results or zero.

We find no evidence that our results are driven by endogenous differences between schools in states that do and do not adopt high-stakes evaluation. Miller et al. (2019) recently added to the literature finding that fixed-effects models which rely on selection into identification often return biased results due to endogenous differences in those units which select the treatment. Thus, in Estimates C7 and C9 of [Figure 4](#) and C2 and C5 of [Figure 5](#) we present results in which we restrict our sample to only those grade-year observations nested in states which ever implemented evaluation. These results, therefore, identify causal effects only off of differential timing of *when*, and not *whether*, states enacted evaluation policy. We again fail to reject the null.

We may also be concerned that our difference-in-difference results are driven by events substantially removed from policy enactment, particularly when we observe these time periods for only some units. Given the start and end periods of our data, the maximal years pre- and post-teacher evaluation reform that we can see for all observations is 5 years pre- and 1 year after the initial policy implementation. In Estimates C8 and C10 of [Figure 4](#) and C3 and C6 of [Figure 5](#), we restrict our sample to grade-year observations during this frame. In all of these estimates but one, we reject the null. When we estimate the effect of evaluation reform for high-accountability grades (3–11) on subjective ODRs in the sample of states that ever implemented evaluation, we find that it modestly reduced the rate of ODRs. Given the multiple hypothesis tests we conduct and the small magnitude of the estimate, we choose to interpret this estimate as consistent with our main findings which are modestly negative in magnitude but indistinguishable from zero.

As expected, our triple-difference estimates (Estimates C11 and C12 of [Figure 4](#) and C7 and 8 of [Figure 5](#)) that difference out the change in non-classroom or objective-rationale referrals from the change in our primary outcomes return estimates even closer to zero.

We find no evidence that alternate teacher accountability or school discipline policy reforms either predict any changes in the rates of disciplinary referrals or that they moderate the effects of teacher evaluation. In Estimates C13 and C14 of [Figure 4](#) and C9 and C10 of [Figure 5](#), we present the results of adjusting the main effect of teacher evaluation implementation for the adoption of these other policy reforms. The results are indistinguishable from our main estimates and zero. [Supplementary Appendix Tables A15 and A16](#) further demonstrate that none of the policy reforms individually predicts changes in ODR rates. We also estimate the effects of bundles of accountability policies separately from the effects of discipline reforms. All results are consistent.

Alternative approaches to clustering standard errors do not change our statistical inferences. In Estimates C15 and C17 of [Figure 4](#), we cluster standard errors at the level of intervention: the state-year. In estimates C16 and C18, we implement Cameron et al.'s (2011) two-way clustering approach at the state and year level. Neither of these approaches changes our inference.

We find no evidence that our results are driven or biased by greater treatment weights imposed on units treated in the middle of the policy window due to greater conditional variance in treatment (Goodman-Bacon, 2021), that there are any negative weights on our individual unit-year observations, or that our results are sensitive to alternate mechanisms for weighting the unit-year-specific average treatment effects (ATEs). In Panel A of [Supplementary Appendix Figure A3](#), we present the distribution of weights on each school-year fixed effect by year that the unit was treated. We find no evidence of systematic variation in fixed effect weight by year of implementation of teacher evaluation reform. In Panel B of [Figure A3](#), we plot the school-year-level ATE against its weight. Notably, there are few outlying values; thus, we are relatively unconcerned with the recent concerns raised about fixed effects estimates in our sample.

We formally compare our results with de Chaisemartin and D'Haultfoeuille's (2018) Wald-TC estimator in Estimates C19 and C20 of [Figure 4](#). Though the exact coefficients on these estimates are slightly different than our main models, they are nevertheless extremely small in magnitude and statistically indistinguishable from either zero or our main model coefficients. We present the graphical event study using the Wald-TC estimator in [Supplementary Appendix Figure A4](#). Note, however, that the results using this estimator are less precise, and we are unable to rule out relatively large effects given this approach.²⁰

Finally, we note that our finding that successfully implementing PBIS has no moderating effect on the implementation of high-stakes evaluation is robust to alternate sample construction. In [Supplementary Appendix Tables A20 and A21](#), we present results that restrict the sample to grade-year observations in states that ever implemented teacher evaluation as well as ones that restrict the sample to 5-years-pre- and 1-year post-evaluation implementation. Results across all models are equivalent to our main models.

²⁰We implement de Chaisemartin and D'Haultfoeuille's (2018) Wald-TC estimator using the May 2019 version of the `did_multipleGT` Stata package. Current versions of `did_multipleGT` implement the DID_M estimator (de Chaisemartin & D'Haultfoeuille, 2020); these return essentially identical results. Replications of the Wald-TC results will return slightly different values due to the use of bootstrapping for obtaining standard errors (we use 50 replications). To reduce computing resource demands, we estimate these results using our school-year sample, though in practice this does not meaningfully affect our standard errors as we cluster them at the state level.

Conclusion

Policy makers and school system leaders have a critical interest in understanding whether, and if so how, educators respond to external accountability pressures. Designers of accountability-based policies must carefully weight the purported benefits of the policy against its potential harms. In this paper we find that, in the context of pressures from higher-stakes teacher evaluation policies, teachers do not, on average, alter their responses to students' classroom misbehavior. Across a variety of specifications and robustness checks, we find no evidence that the rates of removing students from class changed in the aftermath of these policy reforms. Furthermore, we find no evidence that when schools improve at developing systems of behavioral supports that these serve to moderate any effects of evaluation implementation.

We find that higher-stakes teacher evaluation reforms did not alter the rates of ODRs, and our confidence intervals rule out all but substantively meaningless effects. We contextualize our findings in the framework Jacob et al. (2019) provide for assessing null results. Our estimated impact of higher-stakes evaluation is both substantively small and has a tight confidence interval. The most extreme value within the confidence intervals of our main results represents a change of 0.2 referrals in the daily classroom referral rate in a school of 500 students. With classes of 25 students, this would represent a change of roughly 0.01 referrals per teacher in a day. That said, the teacher evaluation reforms at the start of the 2010s were not directly targeted at reducing ODRs, nor by some measures was higher-stakes teacher evaluation implemented intensively. However, the time, human resource and political costs of teacher evaluation reforms were quite high. We believe our null findings are important in adding to the existing literature on the limits of widely adopted higher-stakes evaluation policy reforms in producing changes in classroom practices or student achievement (Bleiberg et al., 2021; Garet et al., 2017; Stecher et al., 2018). Our findings do not speak to more intensively implemented versions of similar policies, nor do they speak to other policies with a more explicit focus on improving classroom behavioral climate.

We introduce a novel outcome to the causal literature that is rarely present in administrative data but is a critical pre-cursor to exclusionary school discipline. Teachers' classroom management practices and disciplinary responses are key mechanisms for student engagement and are frequently students' first points of entry into school disciplinary systems. Growing evidence indicates that suspensions harm both near-term academic performance and the future school attendance of suspended students (Anderson, 2020; Lacoé & Steinberg, 2018). This is particularly salient given recent evidence on the causal effects of exclusionary discipline on students' future college enrollment and involvement in the criminal justice system (Bacher-Hicks et al., 2019). Evidence indicates that zero tolerance policies increase suspensions (Curran, 2016) and alternatives such as PBIS decrease them when implemented well (Horner et al., 2009). However, the general equilibrium effects, including on non-suspended students, of alternative policy approaches are indeterminate (Steinberg & Lacoé, 2018), which may be a function of heterogeneous effects and uneven policy implementation (Skiba, 2015; Welsh & Little, 2018). Our results suggest that teacher accountability policies, either on their own or coupled with well-implemented systems to promote positive behavior, are not sufficient to limit students' entry into the disciplinary pipeline.

We make considerable efforts to test that we have selected appropriate counterfactuals to ensure that our findings do not reflect countervailing forces, either secular trends or unobserved shocks, that mask the true effect of increased accountability. These threats are difficult to fully disprove. Furthermore, data limitations prevent us from fully modeling teachers' classroom-based disciplinary responses. To better understand the data generating process, we would benefit from records of each instance of student misbehavior, including those that do not result in an office referral. That said, this is a challenge that plagues nearly all research that employs disciplinary measures as its outcome. Ideally, we would like to observe heterogeneity in teachers' responses to accountability pressures by demographic characteristics, preparation pathways, professional experiences, assessed skill and more. Additionally, we are unable to distinguish whether our findings reflect low-level intensity in the implementation of teacher evaluation (Kraft & Gilmour, 2017) or limits in the ability of accountability pressures to influence teacher practice. We also do not observe the date of students' disciplinary infractions, and so are unable to test Figlio's (2006) finding that impending state accountability assessments increase the severity of disciplinary responses. As we note above, our broadest population of inference is the 25,000 schools attempting to implement PBIS in the United States, which may not generalize to the full population of public schools; though, the fact that successful PBIS implementation does not moderate our main finding may assuage this concern somewhat. Various sources of classical measurement error may limit the precision of our estimates; however, the magnitudes of all our point estimates are so small that we do not consider this a critical concern. Finally, in some states we observe only two years of referral data under higher-stakes teacher evaluation; and thus may be unable to detect the increasing effects of treatment over time for these later adopters. These outstanding issues present promising opportunities for future analysis of the effects of accountability pressures on classroom choices by teachers, including those beyond pedagogy.

Despite these limitations, our findings contribute to the limited understanding of the effects of accountability policy inside the black-box of classroom practice. Our results capture only one dimension of teacher response to higher-stakes evaluation policies and should be understood in the context of the broader literature. However, for those hoping for dramatic improvements in teaching practice as well as for those concerned about serious unintended consequences of high-stakes evaluation policy, our findings present another reminder of the loose-coupling between accountability policy, teacher behavior and classroom practices.

Acknowledgment

We are grateful to the Education and Community Supports (ECS) research unit at the University of Oregon for access to the confidential School-Wide Information System data. We thank Kent McIntosh and Angus Kittelman for answering various data-related questions and providing substantive feedback. We thank Kaitlin Anderson, Chris Curran, Glen Waddell, Anwasha Guha, several anonymous referees, participants at the Association of Public Policy and Management Fall Conference, the University of Oregon Applied Micro-Econometrics seminar, and the Education Policy Collaborative Annual Meeting for their feedback. All errors are our own.

ORCID

David D. Liebowitz  <http://orcid.org/0000-0001-7375-6034>

Data availability statement

The data can be obtained by filing a request directly with ECS: <https://ecs.uoregon.edu/research-projects/>. Replication materials are available at: 10.17605/OSF.IO/9X8PU.

Open Scholarship



This article has earned the [Center for Open Science](#) badges for Open Materials and Preregistered through Open Practices Disclosure. The materials are openly accessible at [OSF.IO/9X8PU](https://osf.io/9X8PU) and <https://sreereg.icpsr.umich.edu/sreereg/search/search> and enter: Liebowitz for search term OR direct download via: <https://sreereg.icpsr.umich.edu/sreereg/subEntry/2506/pdf?section=all&action=download> <https://sreereg.icpsr.umich.edu/sreereg/subEntry/2507/pdf?section=all&action=download>. To obtain the author's disclosure form, please contact the Editor.

References

- Abadie, A., Athey, S., Imbens, G., & Wooldridge, J. (2017). *When should you adjust standard errors for clustering?* NBER Working Paper Series No. 24003, Cambridge, MA.
- Anderson, K. P. (2020). Academic, attendance, and behavioral outcomes of a suspension reduction policy: Lessons for school leaders and policy makers. *Educational Administration Quarterly*, 56(3), 435–471. <https://doi.org/10.1177/0013161X19861138>
- Athey, S., & Imbens, G. (2018). *Design-based analysis in difference-in-differences settings with staggered adoption*. NBER Working Paper Series No. 24963, Cambridge, MA. <https://doi.org/10.3386/w24963>
- Bacher-Hicks, A., Billings, S. B., & Deming, D. J. (2019). *The school to prison pipeline: Long-run impacts of school suspensions on adult crime*. NBER Working Paper Series No. 26257, Cambridge, MA.
- Bertrand, M., Duflo, E., & Mullainathan, S. (2004). How much should we trust differences-in-differences estimates? *The Quarterly Journal of Economics*, 119(1), 249–275. <https://doi.org/10.1162/003355304772839588>
- Bezinque, A., Garcia, K., Darling, K., & Stuart-Cassel, V. (2018). *Compendium of school discipline laws and regulations for the 50 states, Washington, D.C. and the U.S. territories*. Washington, DC. <http://safesupportivelearning.ed.gov/school-discipline-compendium>
- Borusyak, K., & Jaravel, X. (2017). *Revisiting event study designs* (SSRN Working Paper). SSRN Working Papers. <https://doi.org/10.2139/ssrn.2826228>
- Bragg, D. (2019). *School-wide information system: Dataset D0098*. University of Oregon.
- Brehm, M., Imberman, S. A., & Lovenheim, M. F. (2017). Achievement effects of individual performance incentives in a teacher merit pay tournament. *Labour Economics*, 44, 133–150. <https://doi.org/10.1016/j.labeco.2016.12.008>
- Burgess, S., Rawal, S., & Taylor, E. S. (2021). Teacher peer observation and student test scores: Evidence from a field experiment in English secondary schools. *Journal of Labor Economics*, 39(4), 1155–1186. <https://doi.org/10.1086/712997>
- Callaway, B., & Sant'Anna, P. H. C. (2021). Difference-in-differences with multiple time periods. *Journal of Econometrics*, 225(2), 200–230. <https://doi.org/10.1016/j.jeconom.2020.12.001>

- Cameron, A., Gelbach, J. B., & Miller, D. L. (2011). Robust inference with multiway clustering. *Journal of Business & Economic Statistics*, 29(2), 238–249. <https://doi.org/10.1198/jbes.2010.07136>
- Carrell, S. E., & Hoekstra, M. L. (2010). Externalities in the classroom: How children exposed to domestic violence affect everyone's kids. *American Economic Journal: Applied Economics*, 2(1), 211–228. <https://doi.org/10.1257/app.2.1.211>
- Chakrabarti, R. (2014). Incentives and responses under no child left behind: Credible threats and the role of competition. *Journal of Public Economics*, 110, 124–146. <https://doi.org/10.1016/J.JPUBECO.2013.08.005>
- Chiang, H. (2009). How accountability pressure on failing schools affects student achievement. *Journal of Public Economics*, 93(9–10), 1045–1057. <https://doi.org/10.1016/J.JPUBECO.2009.06.002>
- Cohen, J., Hutt, E., Berlin, R., & Wiseman, E. (2020). The change we cannot see: Instructional quality and classroom observation in the era of common core. *Educational Policy*, 1–27. <https://doi.org/10.1177/0895904820951114>
- Connally, K., & Tooley, M. (2016). *Beyond ratings: Re-envisioning state teacher evaluation systems as tools for professional growth*.
- Cullen, J. B., Koedel, C., & Parsons, E. (2019). The compositional effect of rigorous teacher evaluation on workforce quality. *Education Finance and Policy*, 1–85. https://doi.org/10.1162/edfp_a_00292
- Curran, F. C. (2016). Estimating the effect of state zero tolerance laws on exclusionary discipline, racial discipline gaps, and student behavior. *Educational Evaluation and Policy Analysis*, 38(4), 647–668. <https://doi.org/10.3102/0162373716652728>
- Danielson, C. (1996). *Enhancing professional practice: A framework for teaching*. ASCD.
- de Chaisemartin, C., & D'Haultfoeuille, X. (2018). Fuzzy difference-in-differences. *The Review of Economic Studies*, 85(2), 999–1028. <https://doi.org/10.1093/restud/rdx049>
- de Chaisemartin, C., & D'Haultfoeuille, X. (2020). Two-way fixed effects estimators with heterogeneous treatment effects. *American Economic Review*, 110(9), 2964–2996. <https://doi.org/10.1257/aer.20181169>
- Dee, T. S., & Wyckoff, J. (2015). Incentives, selection, and teacher performance: Evidence from IMPACT. *Journal of Policy Analysis and Management*, 34(2), 267–297. <https://doi.org/10.1002/pam.21818>
- Deming, D. J., Cohodes, S., Jennings, J., & Jencks, C. (2016). School accountability, postsecondary attainment, and earnings. *Review of Economics and Statistics*, 98(5), 848–862. https://doi.org/10.1162/REST_a_00598
- Deming, D. J., & Figlio, D. (2016). Accountability in US education: Applying lessons from K–12 experience to higher education. *Journal of Economic Perspectives*, 30(3), 33–56. <https://doi.org/10.1257/jep.30.3.33>
- Dixit, A. (2002). Incentives and organizations in the public sector: An interpretive review. *The Journal of Human Resources*, 37(4), 696–727. <https://doi.org/10.2307/3069614>
- Donaldson, M. L. (2021). *Multidisciplinary perspectives on teacher evaluation: Understanding the research and theory*. Routledge.
- Donaldson, M. L., & Woulfin, S. (2018). From tinkering to going “rogue”: How principals use agency when enacting new teacher evaluation systems. *Educational Evaluation and Policy Analysis*, 40(4), 531–556. <https://doi.org/10.3102/0162373718784205>
- Duncan, A. (2012, July 23). The Tennessee Story. *The Huffington Post*. https://www.huffpost.com/entry/the-tennessee-story_b_1695467
- Eren, O. (2019). Teacher incentives and student achievement: Evidence from an advancement program. *Journal of Policy Analysis and Management*, 38(4), 867–890. <https://doi.org/10.1002/pam.22146>
- Ferman, B., & Pinto, C. (2019). Inference in differences-in-differences with few treated groups and heteroskedasticity. *The Review of Economics and Statistics*, 101(3), 452–467. https://doi.org/10.1162/rest_a_00759

- Figlio, D. N. (2006). Testing, crime and punishment. *Journal of Public Economics*, 90(4–5), 837–851. <https://doi.org/10.1016/J.JPUBECO.2005.01.003>
- Ford, T. G., Van Sickle, M. E., Clark, L. V., Fazio-Brunson, M., & Schween, D. C. (2017). Teacher self-efficacy, professional commitment, and high-stakes teacher evaluation policy in Louisiana. *Educational Policy*, 31(2), 202–248. <https://doi.org/10.1177/0895904815586855>
- Freyaldenhoven, S., Hansen, C., & Shapiro, J. M. (2019). Pre-event trends in the panel event-study design. *American Economic Review*, 109(9), 3307–3338. <https://doi.org/10.1257/aer.20180609>
- Garet, M. S., Wayne, A. J., Brown, S., Rickles, J., Song, M., & Manzeske, D. (2017). *The impact of providing performance feedback to teachers and principals* (NCESS 2018-4001).
- Gibbons, C. E., Suárez Serrato, J. C., & Urbancic, M. B. (2018). Broken or fixed effects? *Journal of Econometric Methods*, 8(1). <https://doi.org/10.1515/jem-2017-0002>.
- Gilmour, A. F., Majeika, C. E., Sheaffer, A. W., & Wehby, J. H. (2019). The coverage of classroom management in teacher evaluation rubrics. *Teacher Education and Special Education*, 42(2), 161–174. <https://doi.org/10.1177/0888406418781918>
- Goodman-Bacon, A. (2021). Difference-in-differences with variation in treatment timing. *Journal of Econometrics*, 225(2), 254–277. <https://doi.org/10.1016/j.jeconom.2021.03.014>
- Greflund, S., McIntosh, K., Mercer, S. H., & May, S. L. (2014). Examining disproportionality in school discipline for aboriginal students in schools implementing PBIS. *Canadian Journal of School Psychology*, 29(3), 213–235. <https://doi.org/10.1177/0829573514542214>
- Hamilton, L. S., Berends, M., & Stecher, B. M. (2005). *Teachers' responses to standards-based accountability* (Rand Working Papers No. WR-259-EDU), Santa Monica, CA. https://www.rand.org/pubs/working_papers/WR259.html
- Holbein, J. B., & Ladd, H. F. (2017). Accountability pressure: Regression discontinuity estimates of how No Child Left Behind influenced student behavior. *Economics of Education Review*, 58, 55–67. <https://doi.org/10.1016/J.ECONEDUREV.2017.03.005>
- Horner, R. H., Sugai, G., Smolkowski, K., Eber, L., Nakasato, J., Todd, A. W., & Esperanza, J. (2009). A randomized, wait-list controlled effectiveness trial assessing school-wide positive behavior support in elementary schools. *Journal of Positive Behavior Interventions*, 11(3), 133–144. <https://doi.org/10.1177/1098300709332067>
- Hoselton, R. (2018). SWIS 2017-18 summary report. Eugene, OR.
- IES. (2014). State requirements for teacher evaluation policies promoted by race to the top, Washington, DC. <https://ies.ed.gov/ncee/pubs/20144016/pdf/20144016.pdf>
- Imai, K., & Kim, I. S. (2020). On the use of two-way fixed effects regression models for causal inference with panel data. *Political Analysis*, 1–11. <https://doi.org/10.1017/pan.2020.33>
- Jacob, R. T., Doolittle, F., Kemple, J., & Somers, M.-A. (2019). A framework for learning from null results. *Educational Researcher*, 48(9), 580–589. <https://doi.org/10.3102/0013189X19891955>
- Jacobs, S., & Doherty, K. (2015). State of the states 2015: Evaluating teaching, leading and learning.
- Kennedy-Lewis, B. L., & Murphy, A. S. (2016). Listening to “frequent flyers”: What persistently disciplined students have to say about being labeled as “bad. *Teachers College Record: The Voice of Scholarship in Education*, 118(1), 1–40. <https://doi.org/10.1177/016146811611800106>
- Kraft, M. A., Brunner, E. J., Dougherty, S. M., & Schwegman, D. J. (2020). Teacher accountability reforms and the supply and quality of new teachers. *Journal of Public Economics*, 188, 104212. <https://doi.org/10.1016/j.jpubeco.2020.104212>
- Kraft, M. A., & Gilmour, A. F. (2016). Can principals promote teacher development as evaluators? A case study of principals' views and experiences. *Educational Administration Quarterly*, 52(5), 711–753. <https://doi.org/10.1177/0013161X16653445>
- Kraft, M. A., & Gilmour, A. F. (2017). Revisiting the widget effect: Teacher evaluation reforms and the distribution of teacher effectiveness. *Educational Researcher*, 46(5), 234–249. <https://doi.org/10.3102/0013189X17718797>
- Lacoe, J., & Steinberg, M. P. (2018). Do suspensions affect student outcomes? *Educational Evaluation and Policy Analysis*, 0(0). <https://doi.org/10.3102/0162373718794897>

- Ladd, H. F., & Lauen, D. L. (2010). Status versus growth: The distributional effects of school accountability policies. *Journal of Policy Analysis and Management*, 29(3), 426–450. <https://doi.org/10.1002/pam.20504>
- Lazear, E. (2001). Educational production. *The Quarterly Journal of Economics*, 116(3), 777–803. <https://doi.org/10.1162/00335530152466232>
- Liebowitz, D. D. (2020). *Teacher evaluation for growth and accountability: Under what conditions does it improve student outcomes?* (Unpublished Working Paper). Eugene, OR. https://scholar.harvard.edu/files/dliebowitz/files/teacher_eval_review_oct_2020.pdf
- Loeb, S., Miller, L. C., & Wyckoff, J. (2015). Performance screens for school improvement. *Educational Researcher*, 44(4), 199–212. <https://doi.org/10.3102/0013189X15584773>
- Macartney, H. (2016). The dynamic effects of educational accountability. *Journal of Labor Economics*, 34(1), 1–28. <https://doi.org/10.1086/682333>
- Macartney, H., McMillan, R., & Petronijevic, U. (2019). *Teacher value-added and economic agency* (NBER Working Paper Series No. 24747). Cambridge, MA. <https://doi.org/10.3386/w24747>
- Mashburn, A. J., Pianta, R. C., Hamre, B. K., Downer, J. T., Barbarin, O. A., Bryant, D., Burchinal, M., Early, D. M., & Howes, C. (2008). Measures of classroom quality in prekindergarten and children's development of academic, language, and social skills. *Child Development*, 79(3), 732–749. <https://doi.org/10.1111/j.1467-8624.2008.01154.x>
- McIntosh, K., Mercer, S., Hume, A., Frank, J. L., Turri, M., & Mathews, S. (2013). Factors related to sustained implementation of schoolwide positive behavior support. *Exceptional Children*, 79(3), 293–311.
- Mercer, S. H., McIntosh, K., & Hoselton, R. (2017). Comparability of fidelity measures for assessing tier 1 school-wide positive behavioral interventions and supports. *Journal of Positive Behavior Interventions*, 19(4), 195–204. <https://doi.org/10.1177/1098300717693384>
- Miller, D., Shenhav, N., & Grosz, M. (2019). *Selection into identification in fixed effects models, with application to head start* (NBER Working Paper Series No. 26174). Cambridge, MA. <https://doi.org/10.3386/w26174>
- NCTQ. (2011). *State of the states: Trends and early lessons on teacher evaluation and effectiveness policies*.
- NCTQ. (2016). *State-by-state evaluation timeline briefs*.
- NCTQ. (2017). *State teacher policy database*. <https://www.nctq.org/yearbook>
- Neal, D., & Schanzenbach, D. (2010). Left behind by design: Proficiency counts and test-based accountability. *Review of Economics and Statistics*, 92(2), 263–283. <https://doi.org/10.1162/rest.2010.12318>
- Ozek, U. (2012). *One day too late? Mobile students in the era of accountability* (CALDER Working Paper Series No. 82). caldercenter.org/sites/default/files/WP_82_Final.pdf
- Phipps, A., & Wiseman, E. (2019). Enacting the rubric: Teacher improvements in windows of high-stakes observation. *Education Finance and Policy*, 1–51. https://doi.org/10.1162/edfp_a_00295
- Phipps, A. R. (2021). *Unintended consequences of teacher performance pay: A theory on incentives and evidence from Washington, D.C.* (unpublished Working Paper). <https://sites.google.com/view/aaronhippsecon/research>
- Pope, N. G. (2019). The effect of teacher ratings on teacher performance. *Journal of Public Economics*, 172, 84–110. <https://doi.org/10.1016/J.JPUBECO.2019.01.001>
- Rafa, A. (2019). *The status of school discipline in state policy*. Denver, CO. www.ecs.org/wp-content/uploads/The-Status-of-School-Discipline-in-State-Policy.pdf
- Reback, R. (2008). Teaching to the rating: School accountability and the distribution of student achievement. *Journal of Public Economics*, 92(5–6), 1394–1415. <https://doi.org/10.1016/J.JPUBECO.2007.05.003>
- Reback, R., Rockoff, J., & Schwartz, H. L. (2014). Under pressure: Job security, resource allocation, and productivity in schools under No Child Left Behind. *American Economic Journal: Economic Policy*, 6(3), 207–241. <https://doi.org/10.1257/pol.6.3.207>
- Rothstein, J. (2015). Teacher quality policy when supply matters. *American Economic Review*, 105(1), 100–130. <https://doi.org/10.1257/aer.20121242>

- Skiba, R. J. (2015). Interventions to address racial/ethnic disparities in school discipline: Can systems reform be race-neutral? In *Race and social problems* (pp. 107–124). New York: Springer.
- Sorensen, L. C., Bushway, S. D., & Gifford, E. J. (2021). Getting tough? The effects of discretionary principal discipline on student outcomes. *Education Finance and Policy*, 1–74. https://doi.org/10.1162/edfp_a_00341
- Stecher, B., Holtzman, D., Garet, M., Hamilton, L., Engberg, J., Steiner, E., ... Chambers, J. (2018). *Improving teaching effectiveness: Final report: The intensive partnerships for effective teaching through 2015–2016*. RAND Corporation.
- Steinberg, M. P., & Donaldson, M. L. (2016). The new educational accountability: Understanding the landscape of teacher evaluation in the post-NCLB era. *Education Finance and Policy*, 11(3), 340–359. https://doi.org/10.1162/EDFP_a_00186
- Steinberg, M. P., & Lacoë, J. (2018). Reforming school discipline: School-level policy implementation and the consequences for suspended students and their peers. *American Journal of Education*, 125(1), 29–77. <https://doi.org/10.1086/699811>
- Steinberg, M. P., & Sartain, L. (2015). Does teacher evaluation improve school performance? Experimental evidence from Chicago's excellence in teaching project. *Education Finance and Policy*, 10(4), 535–572. https://doi.org/10.1162/EDFP_a_00173
- Strunk, K. O., Barrett, N., & Lincove, J. A. (2017). *When tenure ends: The short-run effects of the elimination of Louisiana's teacher employment protections on teacher exit and retirement*. <https://educationresearchalliancenola.org/files/publications/041217-Strunk-Barrett-Lincove-When-Tenure-Ends.pdf>
- Taylor, E. S., & Tyler, J. H. (2012). The effect of evaluation on teacher performance. *American Economic Review*, 102(7), 3628–3651. <https://doi.org/10.1257/aer.102.7.3628>
- U.S. Department of Education, & U.S. Department of Justice. (2014). *Dear colleague letter on the nondiscriminatory administration of school discipline*.
- Vogell, H. (2011, July 26). Investigation into APS cheating finds unethical behavior across every level. *Atlanta Journal-Constitution*, p. 1. <https://www.ajc.com/news/local/investigation-into-aps-cheating-finds-unethical-behavior-across-every-level>
- Welsh, R. O., & Little, S. (2018). The school discipline dilemma: A comprehensive review of disparities and alternative approaches. *Review of Educational Research*, 88(5), 752–794. <https://doi.org/10.3102/0034654318791582>
- Wieczorek, D., Clark, B., & Theoharis, G. (2019). Principals' instructional feedback practices during race to the top. *Leadership and Policy in Schools*, 18(3), 357–381. <https://doi.org/10.1080/15700763.2017.1398336>
- Winters, M. A., & Cowen, J. M. (2013). Who would stay, who would be dismissed? An empirical consideration of value-added teacher retention policies. *Educational Researcher*, 42(6), 330–337. <https://doi.org/10.3102/0013189X13496145>