

The Effect of State Teacher Evaluation Reforms on Achievement and Attainment

Joshua Bleiberg¹

Eric Brunner²

Erica Harbatkin³

Matthew A. Kraft¹

Matthew G. Springer⁴

April 2022

Abstract

Starting in 2009, most states began to adopt new high-stakes teacher evaluation systems spurred by Federal incentives and requirements. We examine the effects of these controversial and costly reforms on student achievement and attainment at a national scale by exploiting the staggered timing of implementation across states. We find precisely estimated null effects, on average, that rule out impacts as small as 1.5 percent of a standard deviation for achievement and 1 percentage point for high school graduation and college enrollment. We also find little evidence of heterogeneous effects across an index measuring system design rigor, specific design features, and district characteristics.

JEL Codes: I20, I28, J21

Keywords: Teacher evaluation, Student achievement, Education

Corresponding author, Joshua Bleiberg, can be reached at 164 Angell St., 2nd Floor, Providence, RI 02906 or joshua_bleiberg@brown.edu. Authors are listed in alphabetical order. The Spencer Foundation [Award#201700052] and the Institute for Education Sciences [Award # R305A170053] provided generous support to Matthew Kraft for this work. We are grateful for the feedback from Melissa Lyon, Danielle Edwards, Grace Falken, Alvin Christian, Alex Bolves, the participants in the Brown Half-Baked Research Series, the Annual Northeast Economics of Education Workshop, and the Society for Research on Educational Effectiveness annual conference.

1 Brown University, 2 University of Connecticut, 3 Michigan State University, and 4 University of North Carolina at Chapel Hill.

I. Introduction

The returns to improved performance evaluation systems have long been of interest to economists and employers. Evaluation systems have the potential to better align worker's effort with organizational goals as well as to inform employee skill development (Gibbons 1998; Prendergast 1999; Oyer and Schaefer 2011). We study efforts to strengthen performance evaluation in the K-12 public education system, which with more than 3.5 million teacher employees is one of the largest economic sectors in the U.S. Research demonstrates that teachers have large effects on a range of student outcomes, but that teacher effectiveness varies considerably (Chetty, Friedman, and Rockoff 2014; Petek and Pope 2016; Jackson 2018; Kraft 2019). Understanding the impacts of more rigorous and regular performance reviews for public school teachers is particularly important given the sizable potential gains from improving teacher productivity.

Between 2009 and 2017, 44 states and Washington, D.C. implemented major reforms to their teacher evaluation systems. Prior to the reforms, teacher evaluation was an infrequent and largely a perfunctory exercise that resulted in nearly all teachers receiving satisfactory ratings (Weisberg et al. 2009). Strong incentives by the federal government spurred states to reform evaluation systems by regularly evaluating teachers based on multiple measures (including student academic growth) and using performance ratings to inform professional development and personnel decisions. While most states made meaningful changes to their evaluation systems, the specific design features varied across each state and were implemented to differing degrees at the district level given the highly decentralized nature of the U.S. public education system (Kraft and Gilmour 2017).

In this paper, we examine how the statewide implementation of newly reformed teacher evaluation systems affected student achievement and educational attainment. We leverage variation in the timing of adoption of new teacher evaluation systems across states to identify the causal effects of these reforms in an event study and difference-in-differences (DiD) framework. We further explore potential heterogeneity in these effects given the substantial variation in the evaluation metrics and design features adopted by states. Our primary analyses combine data on the timing of state adoption of teacher evaluation reforms with comprehensive district-level student achievement data from 2009 to 2018 on standardized math and English Language Arts (ELA) exams from the Stanford Education Data Archive (SEDA). We augment this achievement data with data from the American Community Survey (ACS) Public Use Microdata Sample (PUMS) to examine the impact of teacher evaluation reforms on longer-run student attainment outcomes, namely high school graduation and college enrollment.

Understanding the *average* effect of these reform, as implemented in practice on a national scale, is critical for several reasons. Existing evidence on the effects of teacher evaluation reforms provides mixed evidence and is limited to a narrow set of districts and states.¹ Although both system design and implementation varied considerably, research suggests that many districts did engage in meaningful efforts to reform their teacher evaluation practices on the ground (Howell and Magazinnik 2017a; Howell and Magazinni 2020). A recent study using a similar identification strategy as ours found that the new state-level evaluation reforms raised the quality of new teachers but also decreased their job satisfaction and led to a substantial decline in

¹ See, for example, Taylor and Tyler 2012; Dee and Wyckoff 2015; Loeb, Miller, and Wyckoff 2015; Steinberg and Sartain 2015; Adnot et al. 2016; Sartain and Steinberg 2016; Stecher et al. 2018; Rodriguez, Swain, and Springer 2020; Cullen, Koedel, and Parsons 2021; Anderson, Cowen, and Strunk 2021; Dee, James, and Wyckoff 2021; Dotter, Chaplin, and Bartlett 2021; Sartain and Steinberg 2021).

the overall supply of newly licensed teacher candidates (Kraft et al. 2020). Thus, the net effect of evaluation reforms on student outcomes remains unclear.

These reforms were also highly controversial, leading to protests and lawsuits challenging their legitimacy in several states (McGuinn 2012; Government Accountability Office 2015a; Sawchuk 2015). Proponents argued that reforming teacher evaluation systems would allow districts to attract and retain more effective teachers by closely linking personnel decisions and compensation to rigorous, multi-measure evaluation ratings (Hanushek 2009). Opponents argued that the new high-stakes evaluation systems were based on invalid and unreliable metrics that would disincentivize cooperation and make the profession less attractive to prospective teachers (Murphy, Hallinger, and Heck 2013).

Efforts to implement teacher evaluation came with substantial financial and time costs as well. Chambers and colleagues (2013) estimate that the costs of implementing teacher evaluation systems in three large school districts was about four tenths of a percent of their total expenditures. Extrapolating from these findings, a back-of-the-envelope estimate suggests that public schools spend about \$2 billion each year on teacher evaluation systems.² New evaluation systems also created large demands on administrators' time to conduct frequent observations and complete considerable paperwork (Neumerski et al. 2018; Kraft and Christian in press), amounted to as much as 19 total days of work (Hess and Bell 2017).

We find that, on average, state teacher evaluation reforms had no discernable effect on student achievement in math or ELA. Estimates from event study models are small in magnitude

² See Exhibit C, which describes teacher evaluation expenditures as a percentage of total district expenditures. 0.4 percent is the average of costs in the 2011 and 2012 school year. Those two school years are the second and third years respectively that teacher evaluation reforms were in place. Total public school expenditures was \$604 billion in 2011 and \$601 billion in 2012 (National Center for Education Statistics 2019). 4 tenths of a percent of \$604 billion is approximately \$2.4 billion.

and statistically insignificant up to five years post-reform. Further, estimates from DiD specifications produce precisely estimated null effects on achievement; we can rule out positive effects of the reforms as small as 0.015 standard deviations in math and 0.009 standard deviations in ELA. Our estimates are precise enough to detect plausible effect sizes on achievement found in simulations of teacher deslection and dynamic models of evaluation reforms (Goldhaber and Hansen, 2010; Staiger and Rockoff, 2010; Rothstein, 2015; Leibowitz, 2021). We also find no evidence that teacher evaluation reforms impacted high school graduation or college enrollment rates and can rule out positive effects as small as 1 percentage point for both attainment measures.

We examine the robustness of these null results in several ways. First, we replicate our null findings using several newly developed two-way fixed effects (TWFE) estimators that address potential bias in the presence of heterogeneous treatment effects (Callaway and Sant’Anna 2020; Goodman-Bacon 2021; Sun and Abraham 2021). Second, we address the potential conflation of evaluation reforms with other related efforts to increase teacher accountability as well as a wide range of time-varying education reforms that occurred during our panel period. Our results are essentially unchanged when we control directly for these other state-level policy changes.

We then turn our focus to exploring whether these average estimates mask important treatment effect heterogeneity based on variation in evaluation system designs across states. To test this, we construct a state-level index of evaluation system design rigor based on 10 evaluation policy components commonly identified as key features of effective systems (Doherty and Jacobs 2015; Howell and Magazinnik 2017b).³ We also group system design elements into

³ See Appendix Table B1 for a full list of the components and their sources.

three broad categories motivated by the primary mechanisms through which proponents argued evaluation would benefit students. Overall, we find little evidence of heterogeneity based on either our index approach or using the broad categories of evaluation system design. Finally, we test for heterogeneous treatment effects across student body characteristics and find little evidence that teacher evaluation reforms impacted student achievement or attainment for any subgroup.

Our paper makes four primary contributions to the literature. First, and foremost, our nationally representative study provides the broadest and most generalizable evidence on the efficacy of teacher evaluation reforms in the U.S. Several studies of evaluation systems implemented in individual districts, such as Washington, D.C., and Chicago Public Schools, provide evidence that teacher evaluation reforms have the potential to improve teacher performance and student achievement, but the findings from those studies may lack generalizability (Taylor and Tyler 2012; Dee and Wyckoff 2015; Steinberg and Sartain 2015; Adnot et al. 2016; Sartain and Steinberg 2016; Dotter, Chaplin, and Bartlett 2021; Dee, James, and Wyckoff 2021). Other studies have found null or negligible effects of evaluation reforms on achievement (Stecher et al. 2018; Cullen, Koedel, and Parsons, 2021; Anderson, Cowen, and Strunk 2021). We illustrate how our results are consistent with this prior evidence by confirming that a small group of evaluation systems identified ex-post as exemplary did appear to raise student achievement. That these systems improved performance while most others that shared similar design features did not, suggests that other factors such as leadership quality, a sustained commitment to continuous improvement, school-level implementation details, and local labor market conditions may play a role in the success of evaluation reforms.

Second, we provide the first evidence on how teacher evaluation reforms affected students' longer-term outcomes. This is important given that prior research has documented how education interventions can affect longer-term outcomes even when effects on test scores are not present or fade out (Chetty et al. 2011; Bailey et al. 2020). Third, we add new empirical evidence to the large performance management literature in economics on the impact of evaluation systems on worker productivity (Baker 1992; Gibbons 1998; Bloom and Van Reenen 2007, 2011; Heinrich and Marschke 2010; Heinrich, Meyer, and Whitten 2010; Cappelli and Conyon 2018). Finally, we contribute broadly to the cross-disciplinary literature on the efficacy of scaling up promising programs and more specifically to the education literature that points to the pitfalls that prior reforms have faced when taken to scale across the decentralized U.S. public education system (Pressman and Wildavsky 1984; Coburn 2003; Honig 2006; Manna 2010; Gupta et al. 2021; Zhou et al. 2021).

The paper proceeds as follows. Section II describes the history and background of teacher evaluation reforms and reviews the related literature. Section III describes the data we assemble to examine the impact of teacher evaluation reforms on student achievement and educational attainment. Section IV outlines our empirical framework for isolating the causal effects of evaluation reforms on our outcomes of interest. We present our main findings in Section V and conclude in Section VI with a discussion of the implications of our results for policy and practice.

II. Background

Evaluation Reforms at the State and District Level

The widespread adoption of teacher evaluation reforms marked a shift from evaluation systems that relied primarily on teacher observation and typically had little, if any, connection

with teacher compensation or employment (Weisberg et al. 2009). The rapid uptake of teacher evaluation reforms came, in part, as a response to President Obama's \$4.35 billion federal Race to the Top (RTTT) program and its offer of large competitive grants to states that were struggling during the Great Recession addhere (Howell and Magazinnik 2017b). In particular, the application rubric for RTTT rewarded states for using student outcomes to evaluate teachers and inform personnel decisions with evaluation ratings. Additionally, the Obama administration required states to commit to teacher evaluation reforms in exchange for a waiver from the No Child Left Behind (NCLB) mandate to reach 100% proficiency by 2014.

Numerous studies confirm that the Federal government successfully leveraged the RTTT grant competition and NCLB waivers to spur widespread changes to state laws and regarding teacher evaluation (Wong 2015; NCTQ 2016; Howell and Magazinnik 2017b; Howell and Magazinni 2020). However, certain features of the new high-stakes evaluation systems promoted by the federal government were taken-up more readily than others. An early evaluation of progress implementing reforms across the 19 RTTT state shows that it induced the vast majority of winners to 1) requirement student achievement growth in evaluations, 2) adopt multicategory rating systems, 3) conduct annual evaluations, 4) require evaluations to be use for professional development, and 5) require evaluations to be used for dismissal decisions (Dragoset et al. 2016). During the time of the interviews with state department of education administrators in 2013, far fewer states required evaluations to be use for compensations and career advancement decisions.

The top-down federal push for evaluation reforms across states resulted in the rapid take-up of reforms, but sometimes failed to engage key stakeholders. Roughly a third of state RTTT winners reported that is was a challenge to maintain support from state legislatures and teachers' union for the reforms (Government Accountability Office 2015b). Many states also provided

school districts with some degree of autonomy in designing and implementing teacher evaluation, either by allowing local discretion within a state-designed system or permitting districts to develop their systems given a set of guidelines (Steinberg and Donaldson 2016). A salient question is whether state evaluation reforms provided districts with so much discretion over implementation that little changed in practice.

The available evidence suggests that state evaluation reforms did meaningfully impact evaluation practices on the ground, but that implementation fidelity varied considerably across districts. The Government Accountability Office surveyed a stratified random sample of 643 school districts across 19 RTTT states from November 2013 to April 2015 about their experiences implementing reforms (Government Accountability Office 2015b). Overall, 40 percent of district leaders indicated that educator effectiveness reforms were implemented with “High” or “Very High Quality,” with 36 percent describing their implementation as “Moderate Quality,” and only 7 percent indicating implementation was of “Low” or “Very Low Quality.” Fifty-one percent of district leaders indicated that capacity issues were “Not at all” or “Somewhat challenging” when implementing reforms related to teacher and principal effectiveness including evaluation reforms. Only 17 percent responded that it was “Very” or “Extremely challenging.” Financial capacity challenges stood out over organizational, human capital, and stateholder capacity as the largest challenge (29 percent responded that it was “Very challenging” or “Extremely challenging”) with 35 percent of respondents indicating their districts modified their plans because of these challenges and 5 percent deferring reforms entirely.

Data from the National Council for Teacher Quality (NCTQ) teacher contract database collected in 2019 also suggests that many large districts implemented key elements of new high-

stakes teacher evaluation systems (NCTQ 2022). The database of 148 school districts includes the 100 largest districts in the country, the largest district in each state, and member districts of the Council of Great City Schools. In 2019, three out of four district evaluation systems assigned at least some weight to student achievement. Eighty-six percent of districts required non-tenured teachers to be evaluated at least once every year, with almost half requiring annual evaluations for tenured teachers as well. The vast majority of large school districts (91 percent) required teachers to receive written feedback or to participate in a conference. Sixty-one percent of school districts used teacher evaluation as a criteria for dismissal and 42 percent of districts offered bonuses for strong evaluations.

To summarize, a broad characterization of a prototypical large district that implemented new teacher evaluation reforms is one where teachers are evaluated annually on a multi-category scale based on administrators' ratings on an instructionally-aligned observation rubric and, in some cases, measures of student growth. Administrators typically provide some individualized performance feedback (often written) to teachers and use evaluation ratings to inform professional development and dismissal decisions. Few teachers are actually removed for poor performance, but teachers generally perceive dismissal as a threat and those rated below satisfactory leave at higher rates.

Mechanisms for Teacher Evaluation to Affect Student Outcomes

Theory predicts that performance evaluations are useful tools for improving worker output. Employers can use personnel evaluations to determine compensation and job responsibilities, as well as to provide feedback when objective measures are not available or cost-prohibitive (Baker 1992; Prendergast 1999). Data from performance evaluation systems can also provide information to leaders of public sector organizations to improve outcomes (Heinrich

2002). In the K-12 education sector, there are two potential mechanisms through which teacher evaluation may impact student achievement and attainment. First, evaluation reforms have the potential to change the composition of the teacher workforce by tying high-stakes personnel decisions such as dismissal and tenure decisions to performance ratings (Gordon, Kane, and Staiger 2006; Goldhaber and Hansen 2010; Staiger and Rockoff 2010; Liebowitz 2021; Sartain and Steinberg 2021). For example, several studies have found that new teacher evaluation systems increased voluntary turnover among lower-performing teachers (Loeb, Miller, and Wyckoff 2015; Steinberg and Sartain 2015; Rodriguez, Swain, and Springer 2020; Cullen, Koedel, and Parsons 2021). Similarly, evidence from a national study of teacher evaluation reforms found that these reforms increased the number of new teacher candidates who had attended more competitive undergraduate institutions but also decreased the overall supply of teacher candidates (Kraft et al. 2020).

Second, teacher evaluation may directly improve current teacher performance. Such improvements might reflect how the evaluation process promotes professional growth on the job and/or increased effort incentivized by dismissal threats or merit pay connected to evaluation scores (Firestone 2014; Donaldson and Papay 2015). The evaluation process itself may support ongoing improvements in teachers' practice if evaluators provide feedback and coaching, prompt teachers to reflect on their practices, or provide data that allow districts to match teachers with targeted professional development (Mintrop and Trujillo 2007; Springer 2010; Woulfin and Rigby 2017; Donaldson 2020; Donaldson and Firestone 2021; Galey-Horn and Woulfin 2021). Experimental studies of low-stakes observation and feedback by peers (Papay et al. 2020; Burgess, Rawal, and Taylor 2021) and administrators (Garet et al. 2018) have found some positive effects on achievement. However, field trials of training programs designed to improve

evaluator feedback in high-stakes settings found no improvements on feedback quality or student achievement (Mihaly et al. 2018; Kraft and Christian in press), while a recent quasi-experimental study found no evidence that teachers alter their professional improvement activities in response to evaluation ratings (Koedel et al. 2019).

Several quasi-experimental and experimental studies in large urban school districts point to the potential for evaluation systems to serve as engines for professional growth. Taylor and Tyler (2012) studied Cincinnati Public School's peer evaluation and feedback system. They found that being observed and evaluated by experienced, expert teachers and school principals improved teachers' ability to raise student achievement in math but not ELA. A similar study of France's national teacher evaluation system found that high-stakes observation and feedback by certified pedagogical inspectors improved teachers' contributions to student achievement (Briole and Maurin 2021).

Research on the District of Columbia Public Schools' high-stakes teacher evaluation system, DC IMPACT, has found positive and sustained effects on student achievement (Dee, James, and Wyckoff 2021). The DC IMPACT system is unique in that it uses master educators and administrators as observers, places substantial weight on test-based measures of teacher performance, offers large financial incentives tied to performance ratings, and has resulted in the dismissal of a non-trivial number of teachers rated as low performing. Studies provide evidence that multiple mechanisms improved performance on the job for teachers (Dee and Wyckoff 2015; Phipps and Wiseman 2021) and teacher quality overall via selective retention and replacement (Adnot et al. 2016).

Evidence from studies examining teacher evaluation systems that are more representative of those adopted at scale nationally in the U.S. is decidedly mixed. In an experimental study of a

pilot implementation of the new teacher evaluation system in Chicago Public Schools, Steinberg & Sartain (2015) found the pilot produced significant improvements in ELA achievement and positive but imprecisely estimated effects in math in the first year. However, the authors found no effect in either math or ELA among the cohort of schools that adopted the system in the second year, pointing to the challenges of sustaining effective evaluations at scale. An evaluation of the Gates Foundation’s Intensive Partnerships for Effective Teaching, which provided \$575 million to improve teacher evaluation across three large school districts and four charter management organizations, found that student achievement and graduation rates were largely unchanged after five years (Stecher et al. 2018). Finally, a recent evaluation of a suite of teacher labor market reforms in Michigan, including teacher evaluation, reduced tenure protections, and reduced collective bargaining power, found largely null effects on student achievement (Anderson, Cowen, and Strunk 2021).

III. Data

Treatment

We draw on data from Kraft et al. (2020) to define the treatment timing of teacher evaluation reforms. We consider a state to be treated in the first year when districts were required to enact the new evaluation system statewide. Figure 1, Panel A, shows the 44 states that reformed teacher evaluation systems throughout the country. California, Iowa, Montana, Nebraska, Vermont, and Wyoming did not reform their teacher evaluation systems. Washington, D.C. was the first to reform its evaluation system, in 2009, while states implemented reforms to their teacher evaluation systems between 2012 to 2017 (See Appendix Figure A1).⁴ The

⁴ Washington, D.C. does not contribute to the estimated effect of teacher evaluation on achievement because we do not observe pre-treatment math or ELA scores. We do observe pre-treatment attainment outcomes and leverage data from Washington, D.C. to identify those effects.

frequency of state reforms peaked in 2014 when 13 states reformed their teacher evaluation systems. The staggered timing in the rollout of reforms across states provides a unique opportunity to measure the effect of these evaluation systems on student outcomes.

<Insert Figure 1 Here>

We also collected data on 10 teacher evaluation policy components identified in the literature as key features of evaluation systems (NCTQ 2011; Doherty and Jacobs 2015; Howell and Magazinnik 2017b; NCTQ 2019). We then constructed an index equal to the number of teacher evaluation policy design components that states required districts to put in place (See Appendix Table B1). As illustrated in Figure 1, Panel B, there was substantial variation in the design rigor of new evaluation systems across states (See Appendix Table B2 for state-specific data).

In addition to examining counts of policy components, we group the 10 design components into three categories based on their policy rationales (See Appendix Table B3). Sixteen states adopted a collection of reforms focused on enhancing the reliability of teacher evaluation measures, 19 adopted either incentives or accountability systems, and 29 used evaluations to provide feedback or inform professional development.

Outcomes

We use district-by-grade-level data from the Stanford Education Data Archive (SEDA), which includes a nearly complete census of school districts, to capture student achievement on high-stakes standardized state tests (Reardon et al. 2021).⁵ The SEDA dataset links student performance across state-specific tests by norming scores relative to performance on the National

⁵ In a few cases entire state-years are excluded from SEDA (Reardon et al. 2021). For example, if fewer than 95 percent of students took the state test or if multiple tests were administered for the same content area in the same year, then the entire state-year is excluded from SEDA.

Assessment of Education Progress (NAEP). SEDA includes test score estimates for third through eighth grade in math and ELA from 2009 to 2018.⁶ Table 1 displays descriptive statistics for the full sample. We observe about 550,000 district-grade observations for both math and ELA.

<Insert Table 1 Here>

To measure educational attainment, we construct state-by-year level estimates of high school graduation rates and college enrollment from the American Community Survey (ACS) Public Use Microdata Sample (PUMS). To measure high school graduation for each year and state, we calculate the proportion of 18-year-olds born in a state who earned a high school diploma or equivalent certificate relative to the total number of 18-year-olds born in a state, and apply appropriate PUMS person weights. To measure college enrollment for each state and year, we divide the number of 22-year-old students born in a state and enrolled in college in each year by the total number of 22-year-olds in a state and year, again using PUMS person weights from 2008 to 2020.⁷ This procedure follows recent research on state education reforms that measures educational attainment based on the expected degree-earning age (Jackson, Wigger, and Xiong 2021; Rothstein and Schanzenbach 2021).⁸

Finally, we use the restricted NAEP student-level data on math and ELA achievement in fourth and eighth grades, available in odd-numbered years between 2003 to 2017, to replicate our core results. The NAEP assessment differs from the assessments used in SEDA in several relevant ways. First, the NAEP is not used for accountability purposes, removing any incentive for strategic behavior to increase scores. Second, the NAEP uses the same set of items for the

⁶ Test scores are aggregated in the SEDA up to the district-grade level and include all of the schools that fall within the borders of traditional public school districts.

⁷ Appendix Table A1 describes the number of treated states and observations across relative time. The analytic sample is “trimmed” to mitigate weak panel balance.

⁸ To avoid endogenous moves into states we use state of birth as a proxy for where a student attended school. Approximately 80 percent of students attend high school and college in their state of birth.

entirety of the study period, improving the validity of comparisons across time, and measuring a broad range of competencies. Finally, the NAEP is limited in that it is administered only every other year and each assessment wave only includes a sample of approximately 4,000 schools (Sikali 2019).

Controls

We supplement our main models with a parsimonious set of covariates. We add controls for the characteristics of schools and inputs to the educational production process related to student achievement or attainment. We measure all control variables prior to the first year of evaluation reforms and interact these baseline values with a time trend to control for potential differences in pre-treatment trends. This approach avoids including endogenous controls that may have been affected by the evaluation reforms themselves. In terms of school district characteristics, we include controls for district race and ethnicity (percent Black, percent Hispanic, percent Native American, and percent Asian), urbanicity, and total enrollment. Our education production process covariates include county level GDP, a poverty index, county unemployment rate, district-level student-teacher ratio, and district-level per-pupil expenditures.⁹ We also add covariates for baseline outcomes to control for pre-treatment differences in student achievement and attainment. Data for the covariates from the achievement outcome models are from the SEDA 2.1 and 4.0 covariate files (Reardon et al. 2021). We obtain county-level GDP from the U.S. Bureau of Economic Analysis (2021) and district-level student/teacher ratios and per-pupil instructional expenditures from the Common Core of Data (U.S. Department of Education 2021). In the models with attainment outcomes, we use a parallel set of covariates

⁹ Poverty index is estimated using socioeconomic status proxies. For more details, see Reardon et al. (2021).

measured at the state level from the NAEP, Common Core of Data, and Bureau of Economic Analysis.

IV. Method

We begin by fitting flexible event study models to test the parallel trends assumption and to explore the non-parametric evolution of any treatment effects:

$$Y_{sdgt} = \sum_{k=-5}^4 \tau_k 1(t = t_s^* + k) + \rho(\mathbf{X}'_{dt=2009} \times Year_t) + \alpha_d + \delta_g + \theta_t + \mu_{sdgt} \quad (1)$$

where Y_{sdgt} is a district-by-grade-by-year measure of mean achievement in grade g for district d in state s in year t (spring of school year). The term $1(t = t_s^* + k)$ represents a set of indicators for the years pre- and post-policy reform, with t_s^* denoting the year in which state s reformed its teacher evaluation system and $k \in [-5, 4]$. Tch_Eval_{st} equals 1 for states that reformed teacher evaluation systems and zero otherwise. \mathbf{X} is a vector of baseline covariates including the school district characteristics, education production process characteristics and baseline outcomes, discussed previously, all interacted with a linear time trend, $Year_t$. Each model also includes district fixed effects (α_d), grade fixed effects (δ_g), and year fixed effects (θ_t). The district fixed effects control for time-invariant district and state characteristics, including pre-treatment policies (e.g., standards-based reforms, teacher credentialing). The year and grade fixed effects control for year- and grade-specific shocks to achievement. μ is an idiosyncratic error term clustered at the state level. An alternative approach to estimating standard errors using the wild cluster bootstrap, which accounts for the small number of state clusters, produces very similar confidence intervals in our setting (Cameron, Gelbach, and Miller 2008; Roodman et al. 2019).

The coefficients of primary interest in Equation 1 are the τ_k 's, which represent the effect of teacher evaluation on our outcomes of interest k years before or after a reform. We measure

these effects relative to the year just prior to the reform ($k = -1$) so that τ_{-3} and τ_1 represent the average effect of reforms on our outcomes of interest three years prior to and one year after reform, respectively.

To examine the non-parametric effect of teacher evaluation on educational attainment, we adapt Equation 1 to focus on our state-by-year measures. The state-level attainment models follow the same specification as the district-level achievement models given by Equation 1, with a few differences. The baseline year in the attainment models is 2008 rather than 2009. The attainment models remove district and grade fixed effects, replacing them with state fixed effects. We also add baseline state-level controls (from 2008) for the percent of students eligible for free or reduced-price lunch (FRPL), percent Black, percent Hispanic, and average per-pupil expenditures, total student enrollment, NAEP scores, and the baseline outcome (either has a high school diploma or enrolled in college) all interacted with a linear time trend.

To improve precision, we complement our event studies with DiD specifications that take the following form:

$$Y_{sdgt} = \beta Tch_Eval_{st} + \rho(\mathbf{X}'_{at=2009} \times Year_t) + \alpha_d + \delta_g + \theta_t + \mu_{sdgt} \quad (2),$$

where Tch_Eval_{st} is an indicator that takes the value of unity if state s had enacted a teacher evaluation reform in year t and zero otherwise. All other variables are as defined in Equation 1. The coefficient of interest in Equation 2 is β , which is the DiD estimate of the effect of teacher evaluation averaged across the post-treatment years in our panel.

Our DiD framework relies on two key assumptions: 1) that comparison states provide a valid counterfactual for the trends in treated states in the absence of treatment; and 2) that there are no unobserved factors correlated with both our outcomes of interest and the timing of teacher evaluation reforms across states. We examine the first assumption visually and empirically using

the non-parametric event study. We also estimate a separate DiD model that includes state-specific linear time trends and examine the robustness of our results to the second assumption by fitting supplemental models that control for other education reforms that occurred within our panel window. The estimates from each approach are similar in sign and magnitude to those from our main DiD specification.¹⁰

Several recent studies have shown that estimates from standard event studies and DiD specifications relying on the staggered timing of treatment for identification may be biased in the presence of heterogeneous treatment effects (Callaway and Sant’Anna 2020; Goodman-Bacon 2021; Sun and Abraham 2021). Consequently, we also report results from alternative TWFE estimators robust to issues related to heterogeneous treatment effects (Cengiz et al. 2019; Baker, Larcker, and Wang 2021; Sun and Abraham 2021). As we report below, our results are very consistent across these alternative estimation approaches.

V. Findings

Student Achievement

Event study estimates from models including baseline controls suggest that, on average, evaluation reforms did not affect students’ performance in math or ELA. As shown in Figure 2, Panel A, in the first year of treatment (i.e., year 0), we can rule out positive effects as small as 0.003 SD for math and 0.005 SD for ELA.¹¹ Estimated effects in subsequent years are less precise, but even five years after treatment, we can rule out positive effects as small as 0.04 SD in both math and ELA. Our event study estimates also provide strong evidence that differential

¹⁰ In auxiliary DiD models we add frequency weights for student enrollment. The weighted models yield similarly sized null effects.

¹¹ The event study estimates with and without controls are similar (see Appendix Table A2).

pre-trends do not drive our estimates: the pre-treatment estimates for all periods in math and ELA are individually and jointly indistinguishable from zero.

Our DiD estimates confirm these null effects and allow us to rule out small potential effects of teacher evaluation, averaged over all post-treatment years. Table 2, Panel A, includes the DiD estimates of the effect of teacher evaluation on student outcomes in math and ELA. The first column presents results without controls, and the second column includes baseline school, educational input, and achievement controls. After adding controls, we can rule out positive effects as small as 0.015 SD in math and about 0.009 SD in ELA.

Educational Attainment

Similar to our achievement findings, event study and DiD estimates suggest that teacher evaluation had little effect on educational attainment. Figure 2, Panel B, provides the estimated effect of teacher evaluation on high school graduation and college enrollment. The effect of teacher evaluation on high school graduation and the percent enrolled in college are both small in magnitude and indistinguishable from zero. Importantly, we once again we find no evidence of differential pre-treatment trends. Estimates from event study models are precise enough to rule out a 2.5 percentage point increase in high school graduation and college enrollment across all observed years post-reform (See Appendix Table A3).

DiD results also show a null effect of teacher evaluation on education attainment. Table 2, Panel B, presents the DiD estimates for educational attainment, pooling over all post-treatment years. Our most precise estimates from models with covariates allow us to rule out a 1 percentage point increase in high school graduation and college enrollment.

<Insert Table 2 Here>

<Insert Figure 2 Here>

Heterogeneity by Evaluation System Design

Our average estimates may mask important treatment effect heterogeneity due to variation in system design. We test for potential heterogeneous effects across states based on our index of the number of design features a state required school districts to put in place. In Table 3, we present models that interact the main treatment indicator with the continuous index of design rigor.¹² Overall, we find no evidence that high design rigor evaluation systems positively effected student achievement or attainment. The estimated coefficient on the interaction term between the treatment indicator and the design rigor index is statistically insignificant for three of our four outcomes. The one exception is ELA, where we find some evidence of negative differential effects. Specifically, states that required districts to put very few design components in place appear to have experienced small declines in student achievement post-reform. For example, our results in Table 3, Panel A, Column 4 suggest that in states which required districts to enact only two design features, the effect of teacher evaluation was -0.04 SD [95% CI: -0.07, -0.01].¹³

<Insert Table 3 Here>

The results in Table 3 suggest that, in general, the effect of teacher evaluation reform on student outcomes did not vary by the rigor of teacher evaluation system designs.¹⁴ We provide further evidence of these null effects by plotting event studies for states with a high number of

¹² In Table 3, the main effect of teacher evaluation is the effect of teacher evaluation for one state (i.e., Alabama) that implemented teacher evaluation, but did not choose a design that includes any of the components we observe in our index.

¹³ The effect of enacting two design features is equal to the main effect of teacher evaluation plus the index multiplied by 2 (i.e., Evaluation+(2 X Index)).

¹⁴ The results in Table 3 are similar when we use the first principal component from Principal Components Analysis.

design components compared with states with a low number of design components separately. Figure 3 shows the event studies where the blue estimates are the effect of teacher evaluation for strong design states (i.e., systems with seven or more design features) relative to comparison states that did not adopt any reforms. The black estimates are the effect of teacher evaluation for states that adopted weaker designs (i.e., between one and six teacher evaluation design components) relative to comparison states. The effect of teacher evaluation is null for states with both stronger and weaker designs in math. Consistent with the differential effects by design rigor from Table 3, the event studies show some evidence of small decreases in ELA scores one to two years after treatment for states with weak evaluation designs. As shown in Appendix Table A4, we find qualitatively similar results when estimating DiD models that pool across the post-treatment periods. The estimates with controls rule out positive effects as small as 0.039 SD in math, 0.040 SD in ELA, a 1 percentage point increase in high school students with diplomas, and a 2 percentage point increase in college enrollment for strong design states.¹⁵

<Insert Figure 3 Here>

Next, we use the 10 design components in our index to construct three non-mutually exclusive measures of specific policy rationales underlying teacher evaluation reforms: 1) reliable measurement; 2) incentives and accountability; and 3) professional development and feedback (see Appendix Table B2 for operationalizations of these dimensions and B3 for state counts). In Figure 4, Panel A, we plot event study estimates from states that adopted policy components to improve the reliability of teacher evaluation measures (e.g., use student test scores weighted at levels shown in research to yield reliable measures, at least two teaching

¹⁵ Appendix Figure A2 mirrors Figure 2 except we change the definition of high quality to states that implemented eight or more teacher evaluation components. We find similarly precise null effects for states that implemented reforms with eight or more teacher evaluation components.

observations, conduct student surveys). Figure 4, Panel B, plots estimates for states that tied incentives and accountability to teacher evaluation (e.g., bonuses, grant tenure). Figure 4, Panel C displays estimates for states that used teacher evaluation to inform professional development or provide feedback to teachers. The blue line shows the effect for evaluation systems with a specific policy rationale, and the black line traces the effect of evaluation systems without the specified policy rationale. Overall, the event study estimates depicted in Figure 4 show little evidence that the effect of teacher evaluation reform varied with specific design components; the estimated coefficients are generally small in magnitude and statistically insignificant.

<Insert Figure 4 Here>

Finally, in Figure 5, we attempt to reconcile our consistent null results with prior research documenting the positive effects of evaluation reforms in selected districts. In October of 2018, the National Council for Teacher Quality (NCTQ) released a report that profiled six district and state evaluation systems that were judged to have designed and implemented exemplary evaluation systems. These systems included Dallas Independent School District (DISD), Denver Public Schools (DPS), District of Columbia Public Schools (DCPS), Newark Public Schools (NPS), Tennessee (TN), and New Mexico (NM) (Putnam, Ross, and Walsh 2018). According to NCTQ, these exemplar systems successfully differentiated among teacher performance, retained higher-performing teachers and removed lower-performing teachers, and coincided with improvements in teacher evaluation ratings and student proficiency rates over time.

We test for differential effects among these exemplar systems by fitting models in which we disaggregate our indicator for treatment, *Tch_Eval*, into two mutually exclusive indicators identifying the implementation of new evaluation systems in 1) these exemplar districts and states and 2) all other states that adopted reforms (excluding the exemplary districts). Consistent

with prior evidence, we find medium-sized positive effects of the implementation of these exemplar evaluation systems on math and ELA achievement. Figure 5 illustrates both the null effects of evaluation among non-exemplary systems and the positive effects over time among exemplar systems rising to as high as 0.15 SD. In our pooled DiD model, we estimate a marginal significant positive effect of 0.09 SD in math and 0.07 SD effect in ELA (See Appendix Table A6).¹⁶

An important caveat to these analyses, is that these exemplar districts were selected ex-post by NCTQ based in part on their outcomes. Consequently, the results presented in Figure 5 and Appendix Table A6 should generally be viewed as descriptive rather than causal. Nevertheless, we view these results as providing evidence that is consistent both with the results of prior studies finding positive impacts of teacher evaluation reform in a small number of select districts, and the null effects found in studies that examine states and districts that implemented reforms that were more representative of those adopted at scale nationally.

<Insert Figure 5 Here>

VI. Robustness Checks

Treatment Timing

We employ two alternative approaches to our standard event study models to test their robustness to potential heterogeneity across states and over time. Our first approach utilizes a stacked DiD estimator that is robust to heterogeneous treatment effects in models with staggered timing of adoption (Cengiz et al. 2019; Baker, Larcker, and Wang 2021). Specifically, we create six datasets, one for each cohort of states that reformed teacher evaluation systems in the same year (i.e., 2012, 2013, 2014, 2015, 2016, 2017), including the states in each cohort and the six

¹⁶ DCPS does not contribute to the estimated effects because no pre-treatment data is observed in SEDA for DCPS. We run a parallel set of models using NAEP that do include Washington, DC and find similar results.

states that never reformed their evaluation systems. We append the six datasets and supplement the models described in equations 1 and 2 by adding district-by-cohort and year-by-cohort fixed effects. Our second approach estimates cohort-specific average treatment effect on the treated (CATT) developed by Sun and Abraham (2021) . This approach is novel in that it calculates weights to estimate the CATT to correct the potential for negative weights in DiD event study models with staggered timing of adoption. Both approaches avoid identifying effects from comparing late to early reformers.

The null effects of teacher evaluation on achievement and attainment are robust to both estimation strategies that account for heterogeneous treatment effects across cohorts. Figure 6 includes event studies for each of the achievement and attainment outcomes from both alternative estimation approaches along with our main estimates. Across outcomes, the magnitude and sign of the estimates in each of the three models are quite similar. The effect of teacher evaluation across relative time remains insignificant. Together, these results suggest that our estimated null effects of teacher evaluation are not biased by treatment effect heterogeneity by adoption cohort.

<Insert Figure 6 Here>

Parallel Trends

The null effect of teacher evaluation is robust to the inclusion of state-specific linear trends, which provides additional evidence that the parallel trends assumption is met. Appendix Table A7 includes results for the achievement and attainment outcomes with and without covariates augmented with state-specific linear trends.¹⁷ The achievement results are within 0.01

¹⁷ We present only effects without covariates for the attainment results because the state-level covariates interacted with the linear trends are collinear with the state-specific linear trends.

SD of the main results in Table 2. Similarly, the attainment results differ by less than 1 percentage point from the main results.

Contemporaneous Policies

Several other education policy reforms occurred contemporaneously during the period of adoption of teacher evaluation reforms. In particular, 17 states enacted reforms to teacher tenure between 2011 and 2014, with five eliminating tenure protections for new teachers and 12 increasing the number of probationary years for untenured teachers. Several states passed laws weakening collective bargaining for teachers between 2011 and 2016, with three restricting or eliminating mandatory collective bargaining and four eliminating mandatory union dues. Several states also enacted reforms to their school finance systems or adopted additional policies rewarded by RTTT (e.g., Common Core State Content Standards, school turnaround initiatives).¹⁸

Because these other reforms occurred in close temporal proximity to teacher evaluation reforms, they could bias our estimates of the impact of teacher evaluation reforms on student outcomes. To account for these potential confounding treatments, we specify models that add a vector of 19 time-varying education policies (Howell and Magazinnik 2017b; Kraft et al. 2020). As shown in Appendix Table A8, we find similarly precise null effects for achievement and attainment outcomes after adding state policy controls. We can rule out positive effects as small as 0.01 SD for achievement outcomes and 1 percentage point for attainment outcomes. The precisely estimated null effects suggest that unobserved education reforms do not bias the estimated effects of teacher evaluation.

¹⁸ See Kraft et al. (2020) for a complete listing of the education policy reforms that occurred contemporaneously during the sample timeframe.

Replicating Results in NAEP

We use the SEDA to measure student achievement in our preferred specification because it includes a near-census of school districts rather than a sampling of schools and is available every year rather than every other year. However, the state test scores used in the SEDA could reflect efforts to artificially raise scores due to the high-stakes attached to these tests (Booher-Jennings 2005; Neal and Schanzenbach 2010; Ballou and Springer 2017). To address this concern, we repeat our primary analyses using fourth- and eighth-grade math and ELA data from the low-stakes NAEP test. As shown in Appendix Table A9, consistent with our main results, we find null effects on achievement. We can rule out positive effects as small as 0.01 SD in math and 0.02 SD in ELA in models including controls.¹⁹ These results add further support for our primary analyses using the SEDA.

VII. Extensions

Academically Vulnerable Groups

Advocates framed teacher evaluation reforms as essential to closing racial and socioeconomic achievement gaps (Weisberg et al. 2009). Consequently, in Appendix Table A10, we extend our primary analyses based on SEDA test scores to test for heterogeneity across sub-populations of students from different racial and socioeconomic backgrounds. Specifically, in our primary DiD specifications, we add interactions between the main effect of teacher evaluation and the percent of students in a district-grade-year eligible for FRPL, percent Black, and percent Hispanic measured at baseline. To improve the interpretability of estimates, we standardize each variable to have a mean of zero and a standard deviation of one. We find little

¹⁹ These models control for the same baseline district characteristics in Equation 1 and add student covariates, including sex, race/ethnicity, free or reduced lunch eligibility, limited English proficiency, has individualized education plan, and modal age for grade. We also add controls for state baseline math and ELA scores in 2003, and an indicator for whether a school made Adequate Yearly Progress in 2003 (Reback et al. 2013).

evidence of heterogeneous effects. The estimated coefficient on the interaction between the treatment indicator and percent Hispanic for ELA and high school graduation is statistically significant and negative. This implies that, if anything, the reforms may have widened rather than closed achievement gaps between Hispanic and White students. However, the size of the effect is substantively small. The results in Appendix Table A10 suggest a 1 SD increase (20 percentage points) in the percent of Hispanic students leads to about a 0.03 SD decrease in ELA scores and a 0.3 percentage point decrease in high school graduation.

VIII. Conclusion

In this paper, we exploit the staggered timing of state teacher evaluation reforms to provide the first nationally representative evidence on how these reforms affected student achievement and educational attainment. We find that, on average, teacher evaluation reforms had no detectable effect on student achievement or attainment. We also find little evidence that the effect of teacher evaluation reforms varied depending on design rigor of the new evaluation systems states implemented or that teacher evaluation improved outcomes for the academically vulnerable groups it was intended to benefit. These null effects are robust to a wide range of specification checks, including alternative TWFE estimators, the inclusion of state-specific linear trends, and controlling for other contemporaneous education reforms.

As noted previously, several studies of evaluation systems implemented in individual districts, provide evidence that evaluation reform can improve student achievement (Taylor and Tyler 2012; Adnot et al. 2016; Dee and Wyckoff 2015; Dotter, Chaplin, and Bartlett 2021; James and Wyckoff 2020). Consistent with the results of those studies, we find positive effects on student achievement for a small set of states and districts with systems identified as exemplary

ex-post. This leads naturally to the question of why, at the national level, teacher evaluation reforms appear to have had little impact on student outcomes.

While we cannot provide a definitive answer to that question, we believe part of the answer is tied to the disconnect between the best practices for performance management systems and the actual design and implementation of new state systems. Despite the widespread adoption of teacher evaluation reforms, many states designed evaluation systems that only vaguely resembled the systems most reformers envisioned. The federal government used RTTT and NCLB waivers to influence the design of new teacher evaluation systems (Howell and Magazinnik 2017b), but this influence had its limits. For example, only 19 states adopted design features intended to link high-stakes accountability and incentives to performance ratings and only 16 established rigorous multi-measure evaluation systems.

Even when states adopted more rigorous design features, these features were not always sustained over time or implemented in ways that resembled the high-stakes systems shown to have positive effects in prior research (NCTQ 2019). Such systems appear to have been organizationally, economically, and politically challenging to scale across a diverse and decentralized U.S. public education system. For example, nationally, less than one percent of teachers were rated as unsatisfactory under the new evaluation systems, with performance-based dismissals being exceedingly rare (Kraft and Gilmour 2017). Similarly, states that did link evaluation to compensation often offered small bonuses of only a few hundred to a thousand dollars and set the bar so low that most teachers qualified for the bonuses (NCTQ 2019). As a result, the accountability components of teacher evaluation systems were often designed and implemented in ways that rendered them low-stakes (Aldeman and Chuong 2014).

Evidence also suggests that evaluation reforms were sometimes implemented in ways that resulted in unintended consequences, a further possible explanation for our null results and the small negative effects we find in some contexts. Prior research documents that teacher evaluation reforms decreased job satisfaction and perceived autonomy among new teachers (Kraft et al. 2020). New evaluation systems created large demands on administrators' time to conduct frequent observations and complete considerable paperwork, displacing other more potentially productive activities (Neumerski et al. 2018). Many districts also placed unrealistic expectations on administrators to provide critical feedback to teachers, narrowing the scope, depth, and quality of feedback teachers received (Hunter and Springer in press; Kraft and Christian in press).

Firms in the private sector often fail to implement best management practices and performance evaluation systems because of imperfectly competitive markets and the costs of implementing such policies and practices (Bloom and Van Reenen 2007). These same factors are likely to have influenced the design and implementation of teacher evaluation reforms. Unlike firms in a perfectly competitive market with incentives to implement management and evaluation systems that increase productivity, school districts and states face less competitive pressure to innovate. Similarly, adopting evaluation systems like the one implemented in Washington D.C. requires a significant investment of time, money, and political capital. Many states and districts may have believed that the costs of fully adopting high-stakes evaluations outweighed the benefits, and subsequently evaluation reforms failed to improve student outcomes.

Reference List

- Adnot, Melinda, Dee, Thomas, Katz, Veronica, and Wyckoff, James, "Teacher turnover, teacher quality, and student achievement in DCPS," *Educational Evaluation and Policy Analysis*, 39 (2016), 54–76.
- Aldeman, Chad, and Chuong, Carolyn, "Teacher Evaluations in an Era of Rapid Change: From 'Unsatisfactory' to 'Needs Improvement'." (Washington, D.C., Bellwether Education Partners, 2014).
- Anderson, Kaitlin P., Cowen, Joshua M., and Strunk, Katharine O., "The impact of teacher labor market reforms on student achievement: Evidence from Michigan," *Education Finance and Policy*, 1 (2021), 1–43.
- Bailey, Drew H., Duncan, Greg J., Cunha, Flávio, Foorman, Barbara R., and Yeager, David S., "Persistence and fade-out of educational-intervention effects: Mechanisms and potential solutions," *Psychological Science in the Public Interest*, 21 (2020), 55–97 (Sage Publications Sage CA: Los Angeles, CA).
- Baker, Andrew, Larcker, David F., and Wang, Charles C.Y., "How Much Should We Trust Staggered Difference-In-Differences Estimates?," SSRN (2021).
- Baker, George, "Incentive contracts and performance measurement," *Journal of Political Economy*, 100 (1992), 598–614.
- Ballou, Dale, and Springer, Matthew G., "Has NCLB encouraged educational triage? Accountability and the distribution of achievement gains," *Education Finance and Policy*, 12 (2017), 77–106.
- Bleiberg, Joshua, and Harbatkin, Erica, "Teacher Evaluation Reform: A Convergence of Federal and Local Forces," *Educational Policy*, 34 (2020), 918–952.
- Bloom, Nicholas, and Van Reenen, John, "Measuring and explaining management practices across firms and countries," *The Quarterly Journal of Economics*, 122 (2007), 1351–1408.
- , "Human resource management and productivity," in *Handbook of Labor Economics*, David Card and Orley Ashenfelter, eds. (Amsterdam, Elsevier, 2011).
- Booher-Jennings, Jennifer, "Below the bubble: 'Educational triage' and the Texas accountability system," *American Educational Research Journal*, 42 (2005), 231–268.
- Briole, Simon, and Maurin, Éric, "There's always room for improvement: the persistent benefits of repeated teacher evaluations," (Padua, Italy, 2021).
- Burgess, Simon, Rawal, Shenila, and Taylor, Eric S., "Teacher Peer Observation and Student Test Scores: Evidence from a Field Experiment in English Secondary Schools," *Journal of Labor Economics*, 39 (2021), 1155–1186.
- Callaway, Brantly, and Sant'Anna, Pedro HC, "Difference-in-differences with multiple time periods," *Journal of Econometrics*, (2020).
- Cameron, A. Colin, Gelbach, Jonah B., and Miller, Douglas L., "Bootstrap-based improvements for inference with clustered errors," *The Review of Economics and Statistics*, 90 (2008), 414–427 (The MIT Press).
- Cappelli, Peter, and Conyon, Martin J., "What do performance appraisals do?," *Industrial and Labor Relations Review*, 71 (2018), 88–116.
- Cengiz, Doruk, Dube, Arindrajit, Lindner, Attila, and Zipperer, Ben, "The effect of minimum wages on low-wage jobs," *The Quarterly Journal of Economics*, 134 (2019), 1405–1454.
- Chambers, J., Brodzia de los Reyes, I., and O'Neil, C., "How much are districts spending to implement teacher evaluation systems," (Washington, D.C., RAND, 2013).

- Chetty, Raj, Friedman, John N., Hilger, Nathaniel, Saez, Emmanuel, Schanzenbach, Diane Whitmore, and Yagan, Danny, "How does your kindergarten classroom affect your earnings? Evidence from Project STAR," *The Quarterly journal of economics*, 126 (2011), 1593–1660 (MIT Press).
- Chetty, Raj, Friedman, John N., and Rockoff, Jonah E., "Measuring the impacts of teachers I: Evaluating bias in teacher value-added estimates," *The American Economic Review*, 104 (2014), 2593–2632.
- Coburn, Cynthia E., "Rethinking scale: Moving beyond numbers to deep and lasting change," *Educational Researcher*, 32 (2003), 3–12.
- Cullen, Julie Berry, Koedel, Cory, and Parsons, Eric, "The compositional effect of rigorous teacher evaluation on workforce quality," *Education Finance and Policy*, 16 (2021), 7–41.
- Dee, Thomas S., James, Jessalynn, and Wyckoff, Jim, "Is Effective Teacher Evaluation Sustainable? Evidence from District of Columbia Public Schools," *Education Finance and Policy*, 16 (2021), 313–346.
- Dee, Thomas S., and Wyckoff, James, "Incentives, selection, and teacher performance: Evidence from IMPACT," *Journal of Policy Analysis and Management*, 34 (2015), 267–297.
- Doherty, Kathryn, and Jacobs, Sandi, "State of the States 2015: Evaluating Teaching, Leading and Learning," (Washington, D.C., National Council on Teacher Quality, 2015).
- Donaldson, Morgaen, "Multidisciplinary Perspectives on Teacher Evaluation: Understanding the Research and Theory," (New York, Routledge, 2020).
- Donaldson, Morgaen, and Firestone, William, "Rethinking teacher evaluation using human, social, and material capital," *Journal of Educational Change*, 22 (2021), 1–34.
- Donaldson, Morgaen, and Papay, John, "Teacher Evaluation for Accountability and Development," in *Handbook of Research in Education Finance and Policy* (Routledge, 2015).
- Dotter, Dallas, Chaplin, Duncan, and Bartlett, Maria, "Impacts of School Reforms in Washington, DC on Student Achievement," (Mathematica Policy Research, 2021).
- Dragoset, Lisa, Thomas, Jaime, Herrmann, Mariesa, Deke, John, James-Burdumy, Susanne, Graczewski, Cheryl, Boyle, Andrea, Tanenbaum, Courtney, Giffin, Jessica, and Upton, Rachel, "Race to the Top: Implementation and Relationship to Student Outcomes," *National Center for Education Evaluation and Regional Assistance*, (2016) (ERIC).
- Firestone, William A., "Teacher evaluation policy and conflicting theories of motivation," *Educational researcher*, 43 (2014), 100–107.
- Galey-Horn, Sarah, and Woulfin, Sarah I, "Muddy Waters: The Micropolitics of Instructional Coaches' Work in Evaluation," *American Journal of Education*, 127 (2021), 441–470.
- Garet, Michael S., Wayne, Andrew J., Brown, Seth, Rickles, Jordan, Song, Mengli, and Manzeske, David, "The Impact of Providing Performance Feedback to Teachers and Principals," (Washington, D.C., National Center for Education Evaluation and Regional Assistance, 2018).
- Gibbons, Robert, "Incentives in organizations," *Journal of Economic Perspectives*, 12 (1998), 115–132.
- Goldhaber, Dan, and Hansen, Michael, "Using performance on the job to inform teacher tenure decisions," *American Economic Review*, 100 (2010), 250–55.
- Goodman-Bacon, Andrew, "Difference-in-differences with variation in treatment timing," *Journal of Econometrics*, (2021).

- Gordon, Robert James, Kane, Thomas J., and Staiger, Douglas, “Identifying Effective Teachers Using Performance on the Job,” (Washington, D.C., The Hamilton Project, 2006).
- Government Accountability Office, “Race to the Top: Survey of State Education Agencies’ Capacity to Implement Reform,” (Washington, D.C., 2015a).
- , “Race to the Top: Survey of School Districts’ Capacity to Implement Reform (GAO-15-317SP, April 2015), an E-supplement to GAO-15-295,” (2015b).
- Gupta, Snigdha, Supplee, Lauren H., Suskind, Dana, and List, John A., “Failed to Scale: Embracing the Challenge of Scaling in Early Childhood,” in *The Scale-Up Effect in Early Childhood and Public Policy* (Amsterdam, Routledge, 2021).
- Hanushek, Eric A., “Teacher deselection,” in *Creating a new teaching profession*, Dan Goldhaber and Jane Hanaway, eds. (Urban Institute, 2009).
- Heinrich, Carolyn J., “Outcomes-based performance management in the public sector: implications for government accountability and effectiveness,” *Public Administration Review*, 62 (2002), 712–725.
- Heinrich, Carolyn J., and Marschke, Gerald, “Incentives and their dynamics in public sector performance management systems,” *Journal of Policy Analysis and Management*, 29 (2010), 183–208.
- Heinrich, Carolyn J., Meyer, Robert H., and Whitten, Greg, “Supplemental education services under No Child Left Behind: Who signs up, and what do they gain?,” *Educational Evaluation and Policy Analysis*, 32 (2010), 273–298.
- Honig, Meredith I., “New directions in education policy implementation: Confronting complexity,” (Suny Press, 2006).
- Howell, William G., and Magazinni, Asya, “Financial Incentives in Vertical Diffusion: The Variable Effects of Obama’s Race to the Top Initiative on State Policy Making,” *State Politics & Policy Quarterly*, 20 (2020), 185–212 (Cambridge University Press).
- Howell, William G., and Magazinnik, Asya, “Presidential Prescriptions for State Policy: Obama’s Race to the Top Initiative,” *Journal of Policy Analysis and Management*, 36 (2017a), 502–531.
- , “Presidential Prescriptions for State Policy: Obama’s Race to the Top Initiative,” *Journal of Policy Analysis and Management*, 36 (2017b), 502–531.
- Hunter, Seth B, and Springer, Matthew G., “Critical Feedback Characteristics, Teacher Human Capital, and Early-Career Teacher Performance: A Mixed-Methods Analysis Using Written Feedback from Formal Evaluation Conferences,” *Educational Evaluation and Policy Analysis*, (in press).
- Jackson, C. Kirabo, “What do test scores miss? The importance of teacher effects on non-test score outcomes,” *Journal of Political Economy*, 126 (2018), 2072–2107.
- Jackson, C. Kirabo, Wigger, Cora, and Xiong, Heyu, “Do school spending cuts matter? Evidence from the great recession,” *American Economic Journal: Economic Policy*, 13 (2021), 304–35.
- Koedel, Cory, Li, Jiayi, Springer, Matthew G., and Tan, Li, “Teacher performance ratings and professional improvement,” *Journal of Research on Educational Effectiveness*, 12 (2019), 90–115.
- Kraft, Matthew, “Teacher effects on complex cognitive skills and social-emotional competencies,” *Journal of Human Resources*, 54 (2019), 1–36.

- Kraft, Matthew, Brunner, Eric J., Dougherty, Shaun M., and Schwegman, David J., "Teacher accountability reforms and the supply and quality of new teachers," *Journal of Public Economics*, 188 (2020), 104212.
- Kraft, Matthew, and Christian, Alvin, "Can teacher evaluation systems produce high-quality feedback? An administrator training field experiment," *American Educational Research Journal*, (in press).
- Kraft, Matthew, and Gilmour, Allison, "Revisiting the widget effect: Teacher evaluation reforms and the distribution of teacher effectiveness," *Educational Researcher*, 46 (2017), 234–249.
- Liebowitz, David D., "Teacher evaluation for accountability and growth: Should policy treat them as complements or substitutes?," *Labour Economics*, (2021), 102024.
- Loeb, Susanna, Miller, Luke C., and Wyckoff, James, "Performance screens for school improvement: The case of teacher tenure reform in New York City," *Educational Researcher*, 44 (2015), 199–212.
- Manna, Paul, "Collision course: Federal education policy meets state and local realities," (Washington, D.C., CQ Press, 2010).
- McGuinn, Patrick, "Stimulating reform: Race to the Top, competitive grants and the Obama education agenda," *Educational Policy*, 26 (2012), 136–159.
- Mihaly, Kata, Schwartz, Heather L., Oppen, Isaac M., Grimm, Geoffrey, Rodriguez, Luis, and Mariano, Louis T., "Impact of a Checklist on Principal-Teacher Feedback Conferences Following Classroom Observations," (Austin, TX, Regional Educational Laboratory Southwest, 2018).
- Mintrop, Heinrich, and Trujillo, Tina, "The practical relevance of accountability systems for school improvement: A descriptive analysis of California schools," *Educational Evaluation and Policy Analysis*, 29 (2007), 319–352.
- Murphy, Joseph, Hallinger, Philip, and Heck, Ronald H., "Leading via teacher evaluation: The case of the missing clothes?," *Educational Researcher*, 42 (2013), 349–354.
- National Center for Education Statistics, "Table 236.10. Summary of expenditures for public elementary and secondary education and other related programs, by purpose: Selected years, 1919-20 through 2017-18," (National Center for Education Statistics, 2019).
- NCTQ, "State of the States 2011: Trends and Early Lessons on Teacher Evaluation and Effectiveness Policies," *National Council on Teacher Quality* (2011).
- , "State-by-state evaluation timeline briefs," (2016).
- , "Teacher & Principal Evaluation Policy. State of the States 2019," *National Council on Teacher Quality* (2019).
- , "Teacher Contract Database," <https://www.nctq.org/contract-database/category/Evaluation> (2022).
- Neal, Derek, and Schanzenbach, Diane Whitmore, "Left behind by design: Proficiency counts and test-based accountability," *The Review of Economics and Statistics*, 92 (2010), 263–283.
- Neumerski, Christine M., Grissom, Jason A., Goldring, Ellen, Rubin, Mollie, Cannata, Marisa, Schuermann, Patrick, and Drake, Timothy A., "Restructuring instructional leadership: How multiple-measure teacher evaluation systems are redefining the role of the school principal," *The Elementary School Journal*, 119 (2018), 270–297.
- Oyer, P., and Schaefer, S., "Personnel economics: Hiring and incentives," in *Handbook of Labor Economics*, David Card and Orley Ashenfelter, eds. (Amsterdam, Elsevier, 2011).

- Papay, John P., Taylor, Eric S., Tyler, John H., and Laski, Mary E., "Learning job skills from colleagues at work: Evidence from a field experiment using teacher performance data," *American Economic Journal: Economic Policy*, 12 (2020), 359–88.
- Petek, Nathan, and Pope, Nolan, "The multidimensional impact of teachers on students," (University of Chicago Working Paper, 2016).
- Phipps, Aaron R., and Wiseman, Emily A., "Enacting the rubric: Teacher improvements in windows of high-stakes observation," *Education Finance and Policy*, 16 (2021), 283–312.
- Prendergast, Canice, "The provision of incentives in firms," *Journal of Economic Literature*, 37 (1999), 7–63.
- Pressman, Jeffrey L., and Wildavsky, Aaron, "Implementation: How great expectations in Washington are dashed in Oakland; Or, why it's amazing that federal programs work at all, this being a saga of the Economic Development Administration as told by two sympathetic observers who seek to build morals on a foundation," (Berkley, CA, University of California Press, 1984).
- Putnam, Hannah, Ross, Elizabeth, and Walsh, Kate, "Making a Difference: Six Places Where Teacher Evaluation Systems Are Getting Results.," (Washington, D.C., National Council on Teacher Quality, 2018).
- Reardon, Sean, Ho, Andrew, Shear, Benjamin, Fahle, Erin, Kalogrides, Demetra, and Chavez, Belen, "Stanford Education Data Archive (Version 4.0)," (2021).
- Rodriguez, Luis A., Swain, Walker A., and Springer, Matthew G., "Sorting through performance evaluations: The influence of performance evaluation reform on teacher attrition and mobility," *American Educational Research Journal*, 57 (2020), 2339–2377.
- Roodman, David, Nielsen, Morten Ørregaard, MacKinnon, James G., and Webb, Matthew D., "Fast and wild: Bootstrap inference in Stata using boottest," *The Stata Journal*, 19 (2019), 4–60 (SAGE Publications).
- Rothstein, Jesse, "Teacher quality policy when supply matters," *American Economic Review*, 105 (2015), 100–130.
- Rothstein, Jesse, and Schanzenbach, Diane Whitmore, "Does Money Still Matter? Attainment and Earnings Effects of Post-1990 School Finance Reforms," (National Bureau of Economic Research, 2021).
- Sartain, Lauren, and Steinberg, Matthew P., "Teachers' Labor Market Responses to Performance Evaluation Reform: Experimental Evidence from Chicago Public Schools," *Journal of Human Resources*, 51 (2016), 615–655.
- , "Can Personnel Policy Improve Teacher Quality? The Role of Evaluation and the Impact of Exiting Low-Performing Teachers," *EdWorkingPapers.com* (Annenberg Institute at Brown University, 2021).
- Sawchuk, Stephen, "Teacher Evaluation Heads to the Courts," *Education Week* (2015) (Oct. 5, 2021).
- Sikali, Emmanuel, "NAEP 2017 National and State Mathematics and Reading, and Puerto Rico Mathematics (Grades 4 & 8) Restricted-Use Data Files," (NCES, 2019).
- Springer, Matthew G., "Performance Incentives: Their Growing Impact on American K-12 Education," (Brookings Institution Press, 2010).
- Staiger, Douglas O., and Rockoff, Jonah E., "Searching for effective teachers with imperfect information," *Journal of Economic perspectives*, 24 (2010), 97–118.

- Stecher, Brian M., Holtzman, Deborah J., Garet, Michael S., Hamilton, Laura S., Engberg, John, Steiner, Elizabeth D., Robyn, Abby, Baird, Matthew D., Gutierrez, Italo A., Peet, Evan D., Brodziak de los Reyes, Iliana, Fronberg, Kaitlin, Weinberger, Gabriel, Hunter, Gerald P., and Chambers, Jay, "Improving Teaching Effectiveness: Final Report: The Intensive Partnerships for Effective Teaching Through 2015–2016," (Washington, D.C., RAND, 2018).
- Steinberg, Matthew P., and Donaldson, Morgaen, "The new educational accountability: Understanding the landscape of teacher evaluation in the post-NCLB era," *Education Finance and Policy*, 11 (2016), 340–359.
- Steinberg, Matthew P., and Sartain, Lauren, "Does teacher evaluation improve school performance? Experimental evidence from Chicago's Excellence in Teaching project," *Education Finance and Policy*, 10 (2015), 535–572.
- Sun, Liyang, and Abraham, Sarah, "Estimating dynamic treatment effects in event studies with heterogeneous treatment effects," *Journal of Econometrics*, 225 (2021), 175–199 (Elsevier).
- Taylor, Eric S., and Tyler, John H., "The effect of evaluation on teacher performance," *The American Economic Review*, 102 (2012), 3628–3651.
- U.S. Bureau of Economic Analysis, "CAGDP2 Gross domestic product (GDP) by county and metropolitan area," (2021).
- U.S. Department of Education, "Common Core of Data," (2021).
- Weisberg, Daniel, Sexton, Susan, Mulhern, Jennifer, and Keeling, David, "The Widget Effect: Our National Failure to Acknowledge and Act on Differences in Teacher Effectiveness. Executive Summary," *New Teacher Project*, (2009).
- Wong, Kenneth K., "Federal ESEA waivers as reform leverage: Politics and variation in state implementation," *Publius: The Journal of Federalism*, 45 (2015), 405–426 (Oxford University Press).
- Woulfin, Sarah L., and Rigby, Jessica G., "Coaching for coherence: How instructional coaches lead change in the evaluation era," *Educational Researcher*, 46 (2017), 323–328.
- Zhou, Jin, Baulos, Alison, Heckman, James J., and Liu, Bei, "The Economics of Investing in Early Childhood," *The Scale-Up Effect in Early Childhood and Public Policy: Why Interventions Lose Impact at Scale and What We Can Do About It*, (2021).
- Zhou, Jin, Alison Baulos, James J. Heckman, and Bei Liu. 2021. "The Economics of Investing in Early Childhood: Importance of Understanding the Science of Scaling." In *The Scale-Up Effect in Early Childhood and Public Policy* (pp. 76-97). Routledge.

TABLE I

Analytic Sample Descriptive Characteristics

Characteristic	Mean	SD	N	Source
ELA Score	0.041	0.38	491,944	SEDA
Math Score	0.041	0.41	460,401	SEDA
High School Graduation	61.3	6.48	520	ACS
College Enrollment	62.6	6.07	520	ACS
Percent White	0.74	0.27	460,287	SEDA
Percent Black	0.08	0.17	460,287	SEDA
Percent Hispanic/Latinx	0.13	0.20	460,287	SEDA
Percent Native American	0.03	0.10	460,287	SEDA
Percent Asian	0.02	0.05	460,287	SEDA
Total Enrollment (Ks)	327.17	979.73	460,287	SEDA
Urban/City	0.07	0.26	460,287	SEDA
GDP Chained \$s (100Ks)	23.55	68.24	460,401	BEA
Poverty Index	0.13	0.07	460,287	SEDA
Unemployment Rate	0.07	0.03	460,287	SEDA
Student Teacher Ratio	15.12	4.16	447,509	CCD
Per-Pupil Expenditures in Ks	6.768	3.60	459,906	CCD

Note: SEDA=Stanford Education Data Archive; ACS=American Community Survey; BEA=Bureaus of Economic Analysis; CCD=Common Core of Data. Table 1 includes descriptive statistics for units included in the analytic sample from the regressions for each outcome. Covariate descriptive restricted are estimated using the district data from the SEDA math sample.

TABLE II
Effect of Teacher Evaluation: Difference-in-Differences Models

	(1)	(2)	(3)	(4)
Panel A. Achievement				
Outcome	Math	Math	ELA	ELA
Teacher Evaluation	-0.0184 (0.0126)	-0.0080 (0.0115)	-0.0220 (0.0120)	-0.0098 (0.0096)
District FE	X	X	X	X
Grade FE	X	X	X	X
Year FE	X	X	X	X
District Ed Controls		X		X
Local SES Controls		X		X
Achievement Controls		X		X
n	460,401	460,401	491,944	491,944
Panel B. Attainment				
Outcome	HS Grad	HS Grad	College Enroll	College Enroll
Teacher Evaluation	0.3996 (0.6677)	-0.0718 (0.6363)	0.0885 (0.6785)	0.0860 (0.7002)
State FE	X	X	X	X
Year FE	X	X	X	X
State Ed Controls		X		X
State SES Controls		X		X
Attainment Controls		X		X
n	520	520	520	520

Note: Models with achievement outcomes include district fixed effects, grade fixed effects, year fixed effects, and baseline covariates measured in 2009 interacted with a linear year trend: Percent Black, Percent Hispanic, Percent Native American, Percent Asian, Total Enrollment, Urban/City, GDP, Poverty Index, Unemployment Rate, Student Teacher Ratio, Per-Pupil Expenditures, ELA Score, and Math Score. Models with attainment outcomes include state fixed effects, year fixed effects, and baseline covariates measured in 2009 interacted with a linear year trend: Percent Black, Percent Hispanic, Percent Native American, Percent Asian, Total Enrollment, Urban/City, GDP, Percent FRPL, Unemployment Rate, Student Teacher Ratio, Per-Pupil Expenditures, either baseline High School Graduation or Baseline College Enrollment. Standard errors are clustered by state. *p < 0.05, ** p < 0.01, ***p < 0.001.

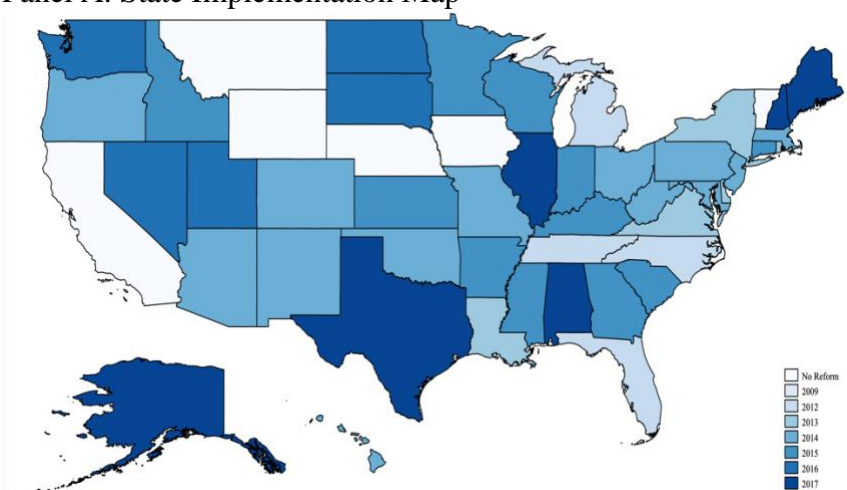
TABLE III
Regressing Continuous Teacher Quality Index on Outcomes

	(1)	(2)	(3)	(4)
Panel A. Achievement				
Outcome	Math	Math	ELA	ELA
Teacher Evaluation	-0.0474 (0.0275)	-0.0318 (0.0248)	-0.0686** (0.0214)	-0.0590** (0.0217)
Teacher Evaluation X Index	0.0052 (0.0048)	0.0043 (0.0044)	0.0085* (0.0037)	0.0090* (0.0039)
District FE	X	X	X	X
Year FE	X	X	X	X
Grade FE	X	X	X	X
District Ed Controls		X		X
Local SES Controls		X		X
Achievement Controls		X		X
n	460,401	460,401	491,944	491,944
Panel B. Attainment				
Outcome	HS Grad	HS Grad	College Enroll	College Enroll
Teacher Evaluation	1.5430 (0.9753)	1.3745 (0.8877)	-0.3799 (1.0460)	-0.8878 (0.9822)
Teacher Evaluation X Index	-0.2035 (0.1508)	-0.2520 (0.1406)	0.0834 (0.1721)	0.1689 (0.1677)
State FE	X	X	X	X
Year FE	X	X	X	X
State Ed Controls		X		X
State SES Controls		X		X
Attainment Controls		X		X
n	520	520	520	520

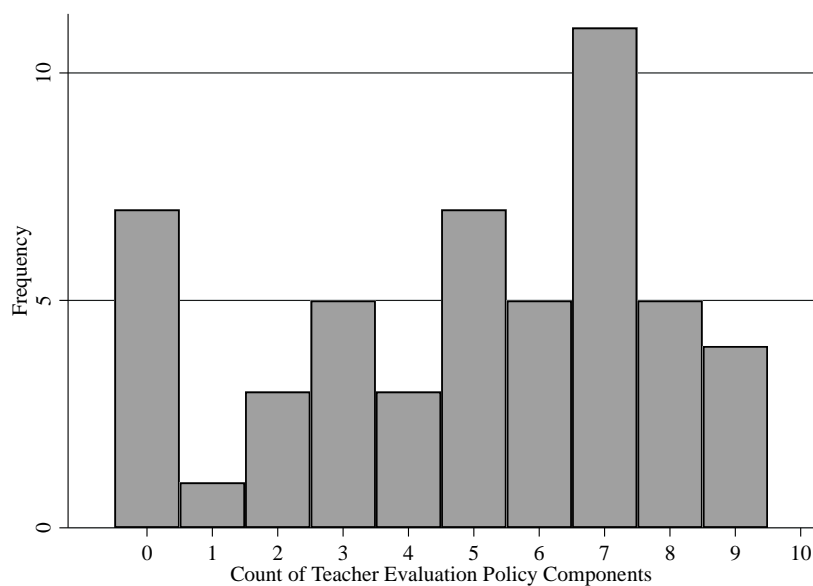
Note: See notes in Table 2 for a full list of covariates. Standard errors are clustered by state. The main effect of teacher evaluation is the effect of teacher evaluation for one state (i.e., Alabama) that implemented teacher evaluation, but did not choose a design that includes any of the components we observe in our index. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Teacher Evaluation Implementation

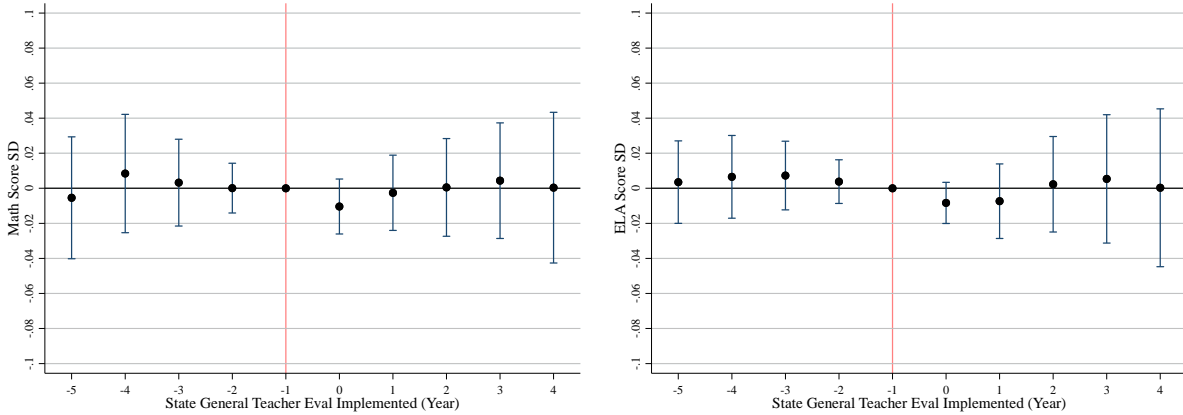
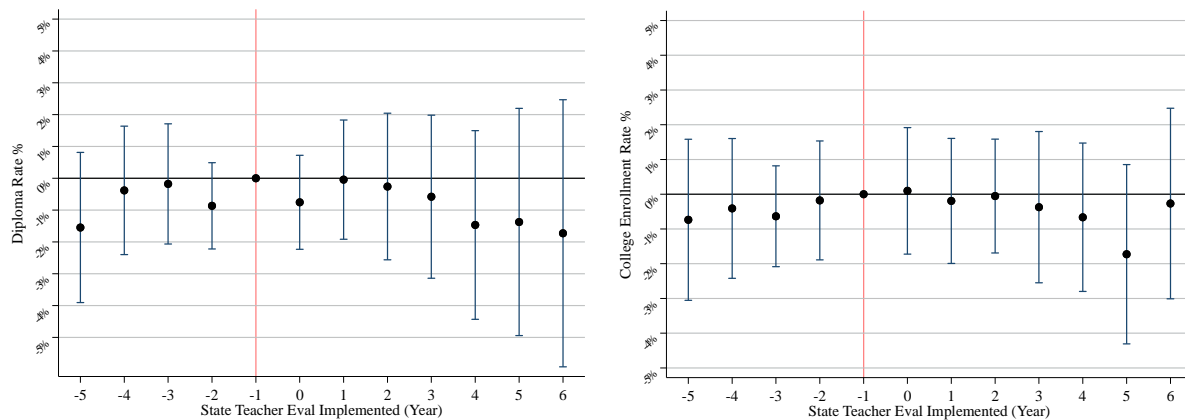
Panel A. State Implementation Map



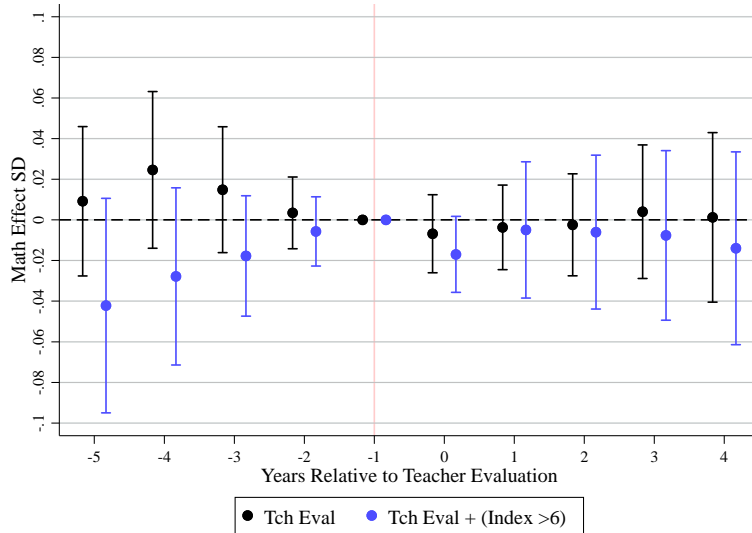
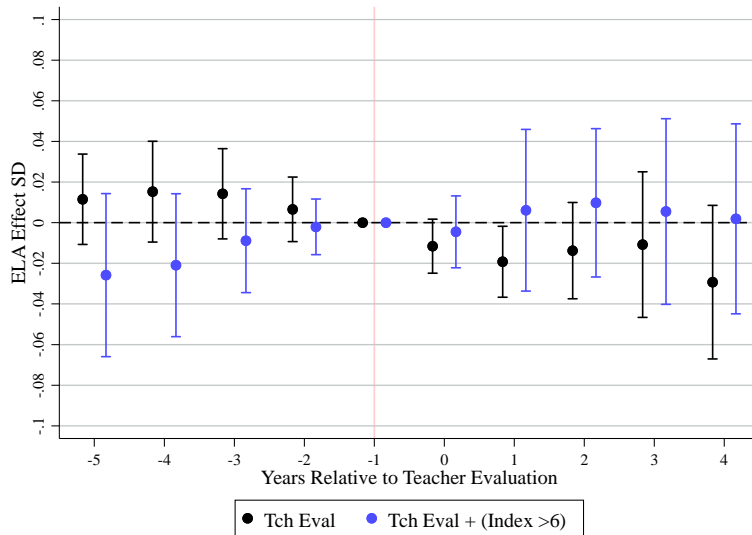
Panel B. Histogram of Teacher Evaluation Reform Quality Index



Note: The index for comparison states is zero even if they implemented a component of teacher evaluation reform. See Appendix B for details on the components of the index. All years are the spring of the school year.

FIGURE II**Event Study: Effects on Achievement and Attainment****Panel A. Achievement****Panel B. Attainment**

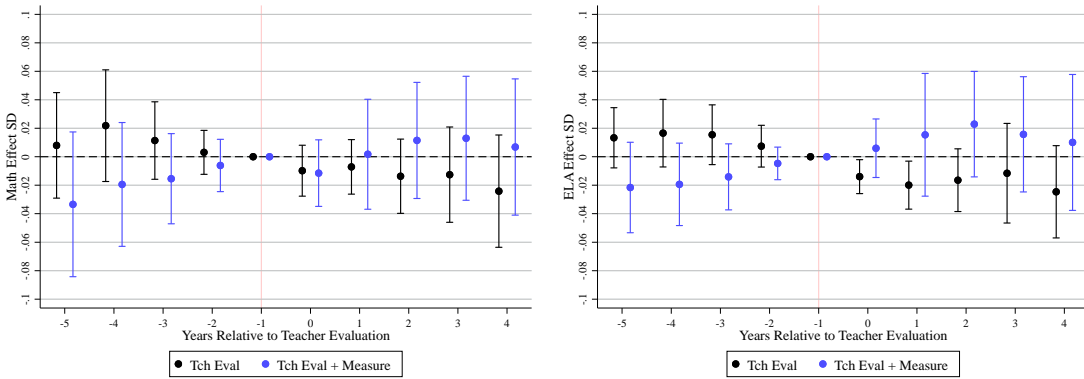
Note: Models with achievement outcomes include district fixed effects, grade fixed effects, year fixed effects, and baseline covariates measured in 2009 interacted with a linear year trend: percent Black, percent Hispanic, percent Native American, percent Asian, total enrollment, urban/city, GDP, poverty index, unemployment rate, student teacher ratio, per-pupil expenditures, ELA score, and math score. Models with attainment outcomes include state fixed effects, year fixed effects, and baseline covariates measured in 2009 interacted with a linear year trend: percent Black, percent Hispanic, percent Native American, percent Asian, total enrollment, urban/city, GDP, percent FRPL, unemployment rate, student teacher ratio, per-pupil expenditures, and either baseline high school graduation or baseline college enrollment. Standard errors are clustered by state.

FIGURE III**Event Study: Heterogeneity by Index****Panel A. Math****Panel B. ELA**

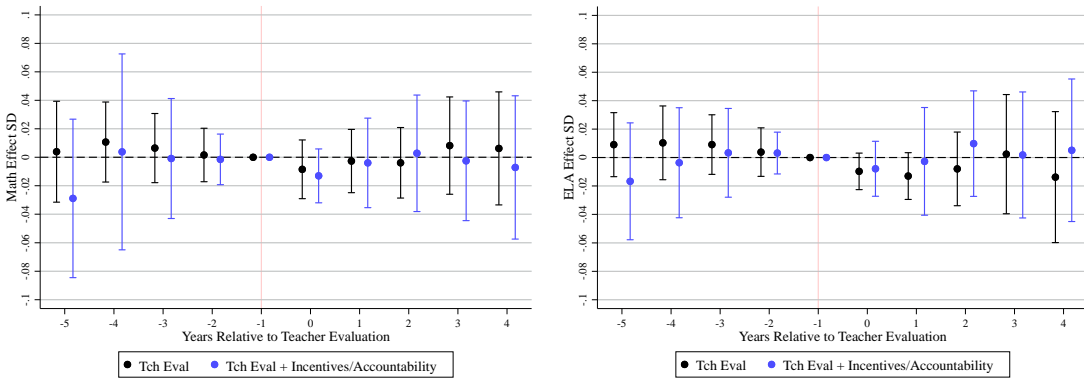
Note: Models include two estimates for each relative time period: the main event study dummies and a set of event study dummies interacted with a time-invariant indicator equal to one for states that had an index from 7 to 10. The black estimates “Tch Eval” are the main event study dummies. The blue estimates “Tch Eval + (Index >6)” are the linear combination of the estimates for the “high group” estimates and main event study estimate from the same relative time period. 20 states have an index from 7 to 10. Model specification found in notes for Figure 2. Standard errors are clustered by state.

FIGURE IV
Event Study: Heterogeneity by System Design

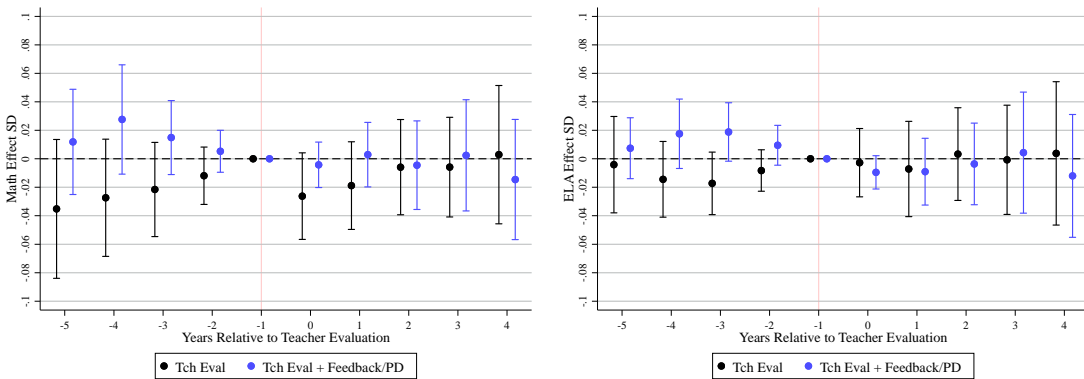
Panel A. Measurement



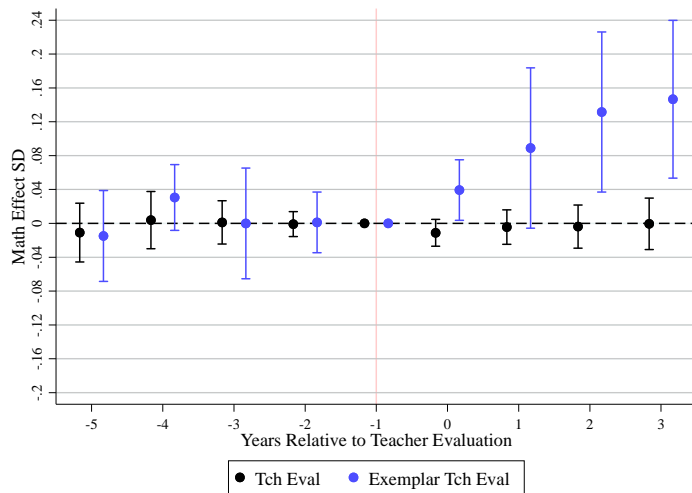
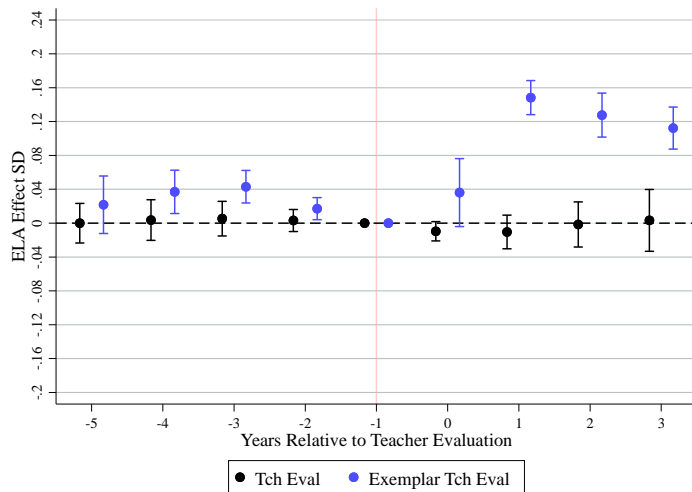
Panel B. Accountability and Incentives



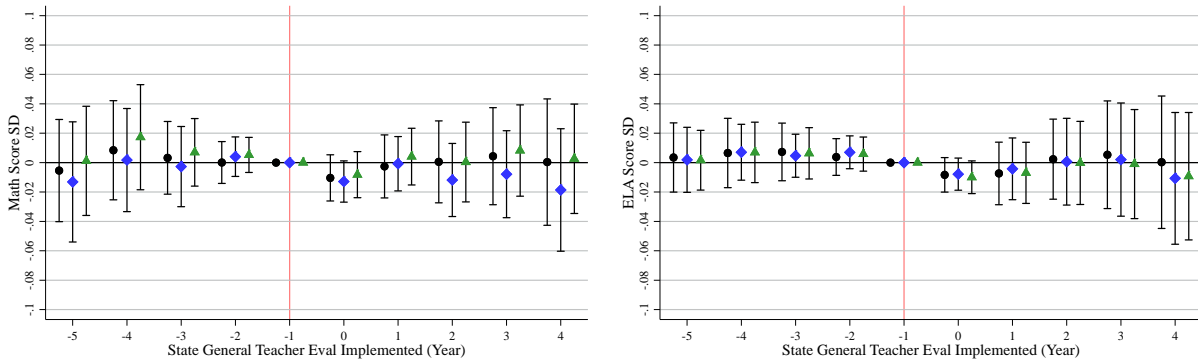
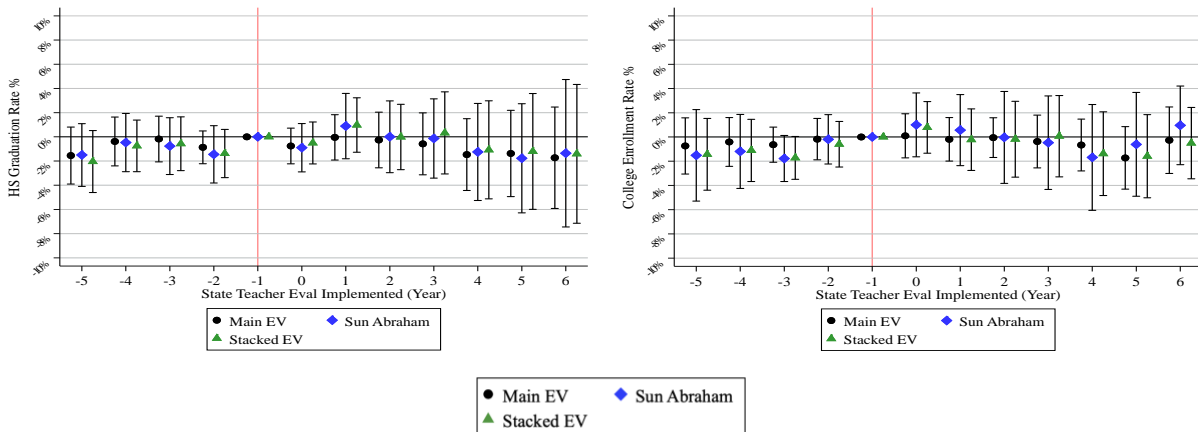
Panel B. Feedback and Professional Development



Note: Models include the main event study dummies and a set of event study dummies interacted with a time-invariant indicator equal to one for a specified system design. The black estimates are the main event study dummies. The blue estimates are the linear combination of the estimates for the “high group” estimates and main event study estimate from the same relative time period. See Table B3 for state system design details. Model specification found in notes for Figure 2. Standard errors are clustered by state.

FIGURE V**Event Study: Heterogeneity by Exemplar Evaluation Systems****Panel A. Math****Panel B. ELA**

Note: Exemplar teacher evaluation systems include: DISD, DCPS, DPS, NPS, TN, NM. Models include two estimates for each relative time period: the main event study dummies and a set of event study dummies interacted with a time-invariant indicator equal to one for states that were exemplar districts. We present effects up to 4 years after adoption of evaluation systems because Tennessee is the only exemplar system we observe outcomes for 5 years after treatment. The black estimates “Tch Eval” are the main event study dummies. The blue estimates “Exemplar Tch Eval” are the effect of teacher evaluation for the exemplar districts and states. Model specification found in notes for Figure 2. Standard errors are clustered by state.

FIGURE VI**Event Study and Estimates Robust to Heterogenous Effects Across Cohorts****Panel A. Achievement****Panel B. Attainment**

Note: Model specification found in notes for Figure 2. Main event study duplicates the results from Figures 2. Diamonds indicate CATT estimates and triangle are stacked event study estimates. Each model includes six stacks with a cohort of treated states and six never treated states. Standard errors are clustered by state by stack.

APPENDIX TABLE A.1

Observations and States Across Relative Time			
Relative Time	Treated States	N	Trimmed
Panel A. Achievement			
Pre -8	6	22,896	X
Pre -7	11	28,133	X
Pre -6	22	54,788	X
Pre -5	34	88,631	
Pre -4	39	99,255	
Pre -3	43	108,143	
Pre -2	41	102,721	
Pre -1	38	97,970	
Post 0	41	94,711	
Post 1	41	88,263	
Post 2	38	72,867	
Post 3	33	70,386	
Post 4	19	41,662	
Post 5	9	11,308	X
Post 6	5	9,361	X
Panel B. Attainment			
Pre -8	6	6	X
Pre -7	11	11	X
Pre -6	22	22	X
Pre -5	35	35	
Pre -4	40	40	
Pre -3	44	44	
Pre -2	44	44	
Pre -1	44	44	
Post 0	45	45	
Post 1	44	44	
Post 2	44	44	
Post 3	44	44	
Post 4	38	38	
Post 5	33	33	
Post 6	22	22	X
Post 7	9	9	X
Post 8	4	4	X

Note: Treated states indicates the number of treated states observable for a specified relative time period. For achievement outcomes, N indicates the number of district-grade observations pooled across subject for a specified relative time period. For attainment outcomes the unit of analysis is state so the unique number of states and number of observations is identical.

APPENDIX TABLE A.2

Event Study Achievement Effects		
	(1)	(2)
Panel A. Math		
-5 Pre	0.0215 (0.0165)	-0.0054 (0.0173)
-4 Pre	0.0256 (0.0162)	0.0084 (0.0168)
-3 Pre	0.0150 (0.0123)	0.0032 (0.0123)
-2 Pre	0.0054 (0.0073)	0.0001 (0.0071)
0 Post	-0.0157 (0.0081)	-0.0104 (0.0078)
1 Post	-0.0137 (0.0111)	-0.0026 (0.0107)
2 Post	-0.0187 (0.0145)	0.0005 (0.0139)
3 Post	-0.0205 (0.0181)	0.0044 (0.0164)
4 Post	-0.0248 (0.0233)	0.0004 (0.0214)
District FE	X	X
Grade FE	X	X
Year	X	X
District Ed Controls		X
Local SES Controls		X
Achievement Controls		X
n	460,401	460,401
Panel B. ELA		
-5 Pre	0.0235 (0.0146)	0.0035 (0.0117)
-4 Pre	0.0222 (0.0136)	0.0065 (0.0118)
-3 Pre	0.0183 (0.0107)	0.0073 (0.0098)
-2 Pre	0.0093 (0.0067)	0.0038 (0.0062)
0 Post	-0.0144* (0.0062)	-0.0083 (0.0058)
1 Post	-0.0193 (0.0113)	-0.0073 (0.0106)
2 Post	-0.0161 (0.0152)	0.0023 (0.0136)
3 Post	-0.0203 (0.0209)	0.0054 (0.0182)
4 Post	-0.0315 (0.0261)	0.0003 (0.0224)
n	491,944	491,944

Note: See notes in Table 2 for a full list of covariates. Model 1 includes state and year fixed effects. Model 2 adds district education, SES, and achievement controls. Standard errors are clustered by state. *p < 0.05, ** p < 0.01, ***p < 0.001.

APPENDIX TABLE A.3

Event Study Attainment Effects		
	(1)	(2)
Panel A. HS Graduation		
-5 Pre	-1.9309 (1.1552)	-1.1564 (1.0726)
-4 Pre	-0.7242 (1.0252)	-0.1103 (0.9090)
-3 Pre	-0.4345 (0.9235)	-0.0142 (0.9061)
-2 Pre	-0.9788 (0.6393)	-0.7651 (0.6396)
0 Post	-0.5336 (0.7072)	-0.7897 (0.7150)
1 Post	0.4394 (0.8247)	-0.0916 (0.8138)
2 Post	0.4157 (0.9060)	-0.4042 (1.0090)
3 Post	-0.0831 (0.9239)	-1.2201 (1.0187)
4 Post	0.2264 (1.0375)	-1.2667 (1.2675)
5 Post	-0.6339 (1.0944)	-2.4965 (1.5312)
n	520	520
Panel B. College Enrollment		
-5 Pre	-0.5615 (1.0750)	-0.2984 (1.2211)
-4 Pre	-0.3048 (0.9649)	-0.0749 (1.0438)
-3 Pre	-0.5730 (0.6965)	-0.4004 (0.7653)
-2 Pre	-0.1636 (0.8210)	-0.0582 (0.8535)
0 Post	0.0430 (0.9159)	-0.0185 (0.9316)
1 Post	-0.2863 (0.9525)	-0.4124 (0.9685)
2 Post	-0.1475 (0.9568)	-0.3558 (1.0330)
3 Post	-1.0473 (1.3214)	-1.3349 (1.4345)
4 Post	-0.8688 (1.3485)	-1.2801 (1.4766)
5 Post	-1.6587 (1.5795)	-2.2053 (1.8105)
n	520	520

Note: See notes in Table 2 for a full list of covariates. Model 1 includes state and year fixed effects. Model 2 adds state covariates and attainment controls. Standard errors are clustered by state. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

APPENDIX TABLE A.4

Effect of Rigorously Designed Teacher Evaluation				
	(1)	(2)	(3)	(4)
Panel A. Achievement				
Outcome	Math	Math	ELA	ELA
Teacher Evaluation	-0.0273 (0.0148)	-0.0161 (0.0135)	-0.0361** (0.0133)	-0.0248* (0.0109)
Eval X High Quality	0.0221 (0.0208)	0.0203 (0.0201)	0.0357* (0.0177)	0.0385* (0.0169)
District FE	X	X	X	X
Year FE	X	X	X	X
Grade FE	X	X	X	X
District Ed Controls		X		X
Local SES Controls		X		X
Achievement Controls		X		X
n	460,401	460,401	491,944	491,944
Panel B. Attainment				
Outcome	HS Grad	HS Grad	College Enroll	College Enroll
Teacher Evaluation	0.4409 (0.7144)	0.0914 (0.6527)	0.1002 (0.6657)	0.0601 (0.6558)
Eval X High Quality	-0.0924 (0.7450)	-0.3787 (0.7039)	-0.0261 (0.7286)	0.0587 (0.7211)
State FE	X	X	X	X
Year FE	X	X	X	X
State Ed Controls		X		X
State SES Controls		X		X
Attainment Controls		X		X
n	520	520	520	520

Note: See notes in Table 2 for a full list of covariates. High quality indicates an index value of 7, 8, or 9. For full list of covariates see Table 1. Standard errors are clustered by state. *p < 0.05, **p < 0.01, ***p < 0.001.

APPENDIX TABLE A.5

Moderation Analysis with Theoretical Constructs						
	(1)	(2)	(3)	(4)	(5)	(6)
Panel A. Achievement						
Outcome	Math	Math	Math	ELA	ELA	ELA
Eval	-0.0206 (0.0127)	-0.0109 (0.0137)	0.0069 (0.0187)	-0.0261* (0.0099)	-0.0177 (0.0101)	0.0073 (0.0165)
Eval X Measurement	0.0396* (0.0184)			0.0533** (0.0170)		
Eval X Incent/Account		0.0064 (0.0206)			0.0180 (0.0181)	
Eval X Feedback/PD			-0.0219 (0.0198)			-0.0254 (0.0188)
District FE	X	X	X	X	X	X
Year FE	X	X	X	X	X	X
Grade FE	X	X	X	X	X	X
District Ed Controls	X	X	X	X	X	X
Local SES Controls	X	X	X	X	X	X
Achievement Controls	X	X	X	X	X	X
n	460,401	460,401	460,401	491,944	491,944	491,944
Panel B. Attainment						
Outcome	HS Grad	HS Grad	HS Grad	College Enroll	College Enroll	College Enroll
Teacher Evaluation	-0.0141 (0.2112)	0.0380 (0.2189)	0.0881 (0.3036)	0.2515 (0.5912)	0.2240 (0.7161)	0.7133 (0.7188)
Eval X Measurement	0.0443 (0.2273)			-0.8687 (0.8514)		
Eval X Incent/Account		-0.0828 (0.2195)			-0.7212 (0.6925)	
Eval X Feedback/PD			-0.1234 (0.2594)			-1.1436 (0.7160)
State FE	X	X	X	X	X	X
Year FE	X	X	X	X	X	X
State Ed Controls	X	X	X	X	X	X
State SES Controls	X	X	X	X	X	X
Attainment Controls	X	X	X	X	X	X
n	520	520	520	520	520	520

Note: See notes in Table 2 for a full list of covariates. Each model includes all fixed effects and controls. See Appendix Table B3 for a full list of states that belong to each construct. *p < 0.05, ** p < 0.01, ***p < 0.001.

APPENDIX TABLE A.6

Effects of Exemplar Evaluation Systems				
	(1)	(2)	(3)	(4)
Outcome	Math	Math	ELA	ELA
Teacher Evaluation	-0.0220 (0.0125)	-0.0105 (0.0114)	-0.0247* (0.0118)	-0.0120 (0.0094)
Exemplar	0.0855 (0.0549)	0.0925 (0.0527)	0.0544* (0.0256)	0.0702* (0.0296)
District FE	X	X	X	X
Year FE	X	X	X	X
Grade FE	X	X	X	X
District Ed Controls		X		X
Local SES Controls		X		X
Achievement Controls		X		X
n	440,565	440,565	471,797	471,797

Note: Exemplar teacher evaluation systems include: DISD, DCPS, DPS, NPS, TN, NM. Models include two estimates for each relative time period: the main event study dummies and a set of event study dummies interacted with a time-invariant indicator equal to one for states that were exemplar districts. We present effects up to 4 years after adoption of evaluation systems because Tennessee is the only exemplar system we observe outcomes for 5 years after treatment. The black estimates “Tch Eval” are the main event study dummies. The blue estimates “Exemplar Tch Eval” are the effect of teacher evaluation for the exemplar districts and states. Model specification found in notes for Figure 2. Standard errors are clustered by state. *p < 0.05, ** p < 0.01, ***p < 0.001.

APPENDIX TABLE A.7

Controlling for State-Specific Linear Trends		
	(1)	(2)
Panel A. Achievement		
Outcome	Math	Math
Teacher Evaluation	-0.0044 (0.0125)	-0.0044 (0.0125)
District FE	X	X
Grade FE	X	X
Year FE	X	X
District Ed Controls		X
Local SES Controls		X
Achievement Controls		X
State-Specific Trends	X	X
n	460,401	460,401
Outcome	ELA	ELA
Teacher Evaluation	-0.0044 (0.0080)	-0.0043 (0.0080)
District FE	X	X
Grade FE	X	X
Year FE	X	X
District Ed Controls		X
Local SES Controls		X
Achievement Controls		X
State-Specific Trends	X	X
n	491,944	491,944
Panel B. Attainment		
Outcome	HS Grad	College Enroll
Teacher Evaluation	-0.1086 (0.8974)	0.2126 (1.0088)
State FE	X	X
Year FE	X	X
State-Specific Trends	X	X
n	520	520

Note: See notes in Table 2 for a full list of covariates. Standard errors are clustered by state. Covariates in the models with attainment outcomes are interacted with linear time trends and are perfectly collinear with the state-specific trends. *p < 0.05, ** p < 0.01, ***p < 0.001.

APPENDIX TABLE A.8

Controlling for Time Varying State Policies				
	(1)	(2)	(3)	(4)
Panel A. Achievement				
Outcome	Math	Math	ELA	ELA
Teacher Evaluation	-0.0123 (0.0136)	-0.0033 (0.0106)	-0.0206 (0.0138)	-0.0077 (0.0099)
District FE	X	X	X	X
Grade FE	X	X	X	X
Year FE	X	X	X	X
District Ed Controls		X		X
Local SES Controls		X		X
Achievement Controls		X		X
State Policies	X	X	X	X
n	455,388	455,388	486,663	486,663
Panel B. Attainment				
Outcome	HS Grad	HS Grad	College Enroll	College Enroll
Teacher Evaluation	0.5227 (0.6655)	0.1217 (0.6573)	-0.5513 (0.6321)	-0.5141 (0.6298)
State FE	X	X	X	X
Year FE	X	X	X	X
State Ed Controls		X		X
State SES Controls		X		X
Attainment Controls		X		X
State Policies	X	X	X	X
n	513	513	513	513

Note: See notes in Table 2 for a full list of covariates. Standard errors are clustered by state. Policy covariates from Kraft et al (2020) and Howell & Magazinnik (2017b) include eliminate tenure, increase probationary period, weaken collective bargaining, eliminate mandatory union dues, won Race to the Top, implement Common Core, basic skills licensure tests, content area licensure tests, pedagogical knowledge licensure tests, Common Core assessment, charter authorizer, charter building funds, charter cap, school turnaround, alternative teacher certification, vouchers, high school exit exams, summative testing, and school finance reform interacted with state quartiles of median household income (2000). *p < 0.05, ** p < 0.01, ***p < 0.001.

APPENDIX TABLE A.9

Replicating results using the Low-Stakes NAEP Assessment		
	(1)	(2)
Panel A. Math		
Teacher Evaluation	-0.063*	-0.024
	(0.026)	(0.013)
District FE	X	X
Grade FE	X	X
Year FE	X	X
Student Controls		X
District Ed Controls		X
Achievement Controls		X
n	1,480,590	1,480,590
Panel B. ELA		
Teacher Evaluation	-0.058	0.002
	(0.044)	(0.011)
District FE	X	X
Grade FE	X	X
Year FE	X	X
Student Controls		X
District Ed Controls		X
Achievement Controls		X
n	1,397,020	1,397,020

Note: Student covariates include sex, race/ethnicity, Free and Reduced Price Lunch Eligibility, Limited English Proficiency, has Individualized Education Plan, and modal age for grade.

District covariates includes all the district characteristics included in Table 1. NAEP samples sizes rounded in accordance with NCES restricted use rules. Achievement characteristics include state baseline math and ELA scores in 2003 and a school level indicator of whether a school made Adequate Yearly Progress. NAEP results use student-level inverse probability weights. Standard errors are clustered by state. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

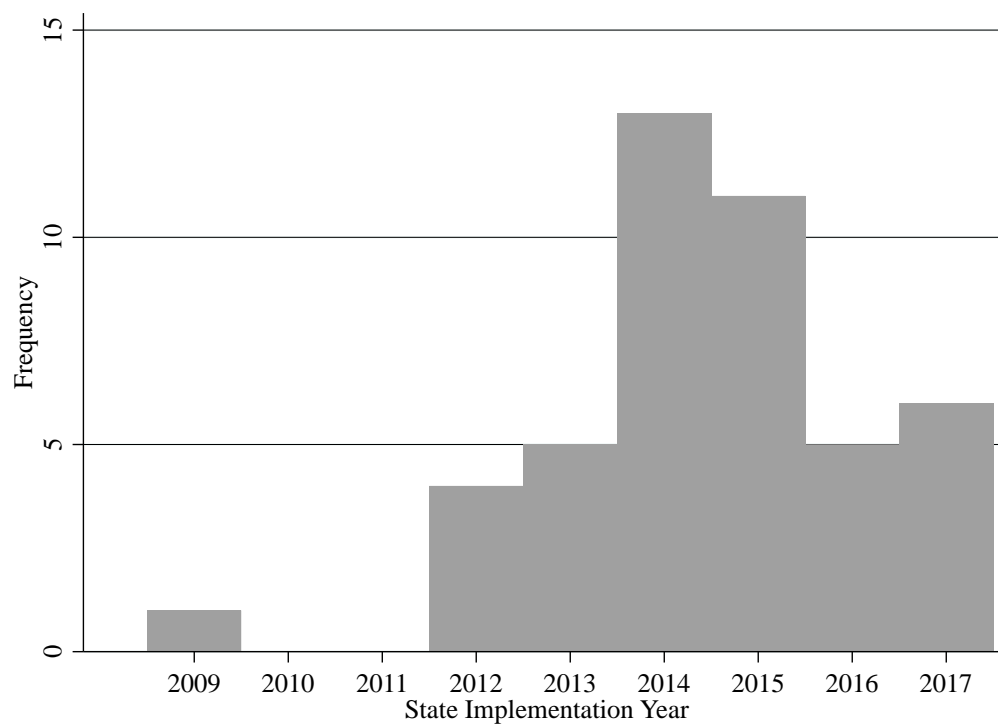
APPENDIX TABLE A.10

Differential Effects for Sub-Groups				
	(1)	(2)	(3)	
Panel A. Math				
Teacher Evaluation	-0.0125 (0.0120)	-0.0103 (0.0117)	-0.0140 (0.0123)	
Teacher Evaluation X Percent FRPL	-0.0138 (0.0092)			
Teacher Evaluation X Percent Black		-0.0007 (0.0052)		
Teacher Evaluation X Percent Hispanic			-0.0108 (0.0056)	
n	450,163	450,163	450,163	
Panel B. ELA				
Teacher Eval	-0.0210 (0.0129)	-0.0201 (0.0127)	-0.0152 (0.0090)	
Teacher Evaluation X Percent FRPL	-0.0069 (0.0068)			
Teacher Evaluation X Percent Black		-0.0019 (0.0052)		
Teacher Evaluation X Percent Hispanic			-0.0179*** (0.0051)	
n	480,801	480,801	480,801	
Panel C. High School Graduation				
Teacher Eval	-0.1558 (0.6235)	-0.1179 (0.6235)	-0.0908 (0.6279)	
Teacher Evaluation X Percent FRPL	-1.0503 (0.6475)			
Teacher Evaluation X Percent Black		-0.4382 (0.7577)		
Teacher Evaluation X Percent Hispanic			-0.7385 (0.5881)	
n	520	520	520	
Panel D. College Enrollment				
Teacher Eval	0.1392 (0.7012)	0.1937 (0.7116)	0.0897 (0.6985)	
Teacher Evaluation X Percent FRPL	0.7934 (0.8557)			
Teacher Evaluation X Percent Black		1.0920 (0.7692)		
Teacher Evaluation X Percent Hispanic			0.5846 (0.6704)	
n	520	520	520	

Note: Models with achievement outcomes includes district, year, and grade fixed effects, district education controls, local SES controls, and achievement controls. Models with attainment outcomes include state, year fixed effects, state education, state SES controls, and attainment controls. See notes in Table 2 for a full list of covariates. Standard errors are clustered by state. Poverty rate, percent Black, and percent Hispanic are all measured at baseline (2009) and standardized. *p < 0.05, ** p < 0.01, ***p < 0.001.

APPENDIX FIGURE A1

Implementation of Evaluation Reforms by Year

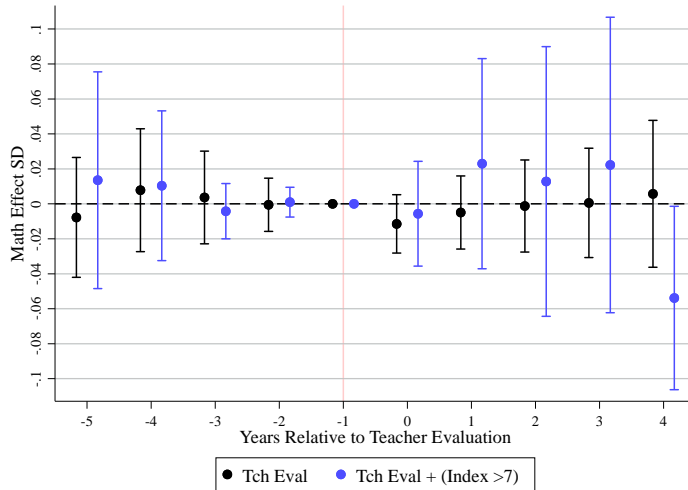


Note: All years are the spring of the school year.

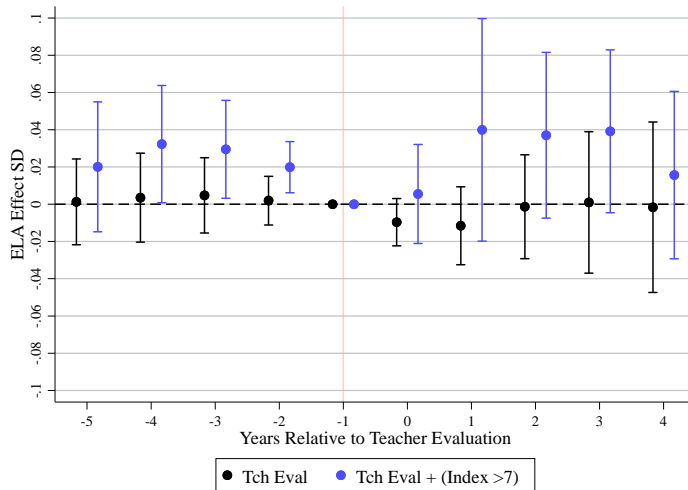
APPENDIX FIGURE A2

Event Study Rigorous Design (Index 8 to 10)

Panel A. Math



Panel B. ELA



Note: Models include two estimates for each relative time period: the main event study dummies and a set of event study dummies interacted with a time-invariant indicator equal to one for states that had an index from 8 to 10. The black estimates “Tch Eval” are the main event study dummies. The blue estimates “Tch Eval + (Index > 7)” are the linear combination of the estimates for the “high group” estimates and main event study estimate from the same relative time period. 9 states have an index from 8 to 11: CT, DC, DE, GA, LA, NJ, RI, TN, and UT. Model specification found in notes for Figure 2. Standard errors are clustered by state.

APPENDIX TABLE B.1**Teacher Evaluation Reform Components**

Category	Variable	Descriptions	Source	State #
Accountability/ Incentive	Fire Teachers	Tenured and untenured teachers rated “ineffective” may be removed from their position.	Howell & Magazinnik (2017)	28
Accountability/ Incentive	Grant Tenure	Teacher evaluation ratings used to grant tenure and/or full certification.	Howell & Magazinnik (2017)	29
Accountability/ Incentive	Bonus	Providing additional compensation to teachers rated “highly effective”.	Howell & Magazinnik (2017)	20
Accountability/ Incentive	Career Ladder	Providing additional responsibilities to teachers rated “highly effective”.	Howell & Magazinnik (2017)	11
Measurement	Multiple Categories	Evaluations have three or more rating categories.	Howell & Magazinnik (2017)	38
Measurement	Observations Required	Observations are a required component of teacher evaluations.	Doherty & Jacobs (2015)	27
Measurement	Student Survey	Student surveys are a required component of teacher evaluations.	Doherty & Jacobs (2015)	7
Measurement	Student data	Student test scores (e.g., growth scores, value-added) with a weight of 20-50 percent are a required component of teacher evaluations.	Bleiberg & Harbatkin (202); Doherty & Jacobs (2015)	21
Feedback/PD	Feedback Required	Teachers receive feedback based on their evaluations.	Doherty & Jacobs (2015)	35
Feedback/PD	Inform PD	Teacher evaluations inform coaching, induction support, and/or professional development.	Howell & Magazinnik (2017)	36

Note: Howell & Magazinnik (2017) do not include data for DC. Design features for DC were determined using the NCTQ State of the State reports from three years were used were used (NCTQ 2011, 2019; Doherty and Jacobs 2015).

APPENDIX TABLE B.2

Teacher Evaluation Categorical Constructs and Quality Measures

Category	Descriptions	State #
Measurement	Teacher evaluation systems include at least three of the following components: (1) Student test scores weighted 20 to 50 percent; (2) observations [at least two explicitly required]; (3) student surveys; (4) Evaluations have three or more rating categories.	16
Accountability/ Incentive	Teacher evaluation systems include at least three of the following components: (1) Evaluation used to either grant tenure or (2) remove teachers from their position and evaluations used for either (3) promotions or (4) bonuses.	19
Feedback/PD	Teachers must receive feedback based on their evaluation; have their evaluation inform coaching, induction support and/or professional development.	29
Low Quality	State index value is 0 to 3.	10
Medium Quality	State index value is 4 to 6.	15
High Quality	State index value is 7 to 9.	20

APPENDIX TABLE B.3**State Teacher Evaluation Component Measures by State**

State	Ever Adopted	Measurement	Accountability/ Incentive	Feedback/PD	Index
AK	1	0	0	0	4
AL	1	0	0	0	0
AR	1	0	1	1	7
AZ	1	0	0	1	5
CA	0	0	0	0	0
CO	1	0	1	1	7
CT	1	1	1	1	9
DC	1	1	1	0	8
DE	1	0	1	1	7
FL	1	0	1	1	7
GA	1	1	1	1	9
HI	1	1	0	1	8
IA	0	0	0	0	0
ID	1	0	0	0	3
IL	1	0	0	1	6
IN	1	1	1	0	7
KS	1	0	0	0	4
KY	1	1	0	1	6
LA	1	1	1	1	8
MA	1	0	1	1	8
MD	1	0	0	0	3
ME	1	1	0	1	7
MI	1	0	1	1	7
MN	1	0	0	0	3
MO	1	0	0	1	3
MS	1	0	0	0	1
MT	0	0	0	0	0
NC	1	0	0	0	5
ND	1	0	0	0	2
NE	0	0	0	0	0
NH	1	0	1	1	6
NJ	1	1	0	1	7
NM	1	1	0	1	5
NV	1	0	1	1	7
NY	1	0	1	1	6
OH	1	1	1	0	7
OK	1	1	1	0	7
OR	1	0	0	0	3
PA	1	1	0	0	5
RI	1	1	1	1	9
SC	1	0	0	0	2
SD	1	0	0	1	5
TN	1	1	1	1	8
TX	1	0	0	1	2
UT	1	1	1	1	9
VA	1	0	0	0	4
VT	0	0	0	0	0
WA	1	0	0	1	5
WI	1	0	0	1	6
WV	1	0	0	1	5
WY	0	0	0	0	0