

Understanding Incentives in Teacher Subjective Evaluations

Andrew J Morgan *

September 4, 2025

Abstract

Employee evaluation is a central function of any firm, yet supervisor ratings that are compressed in the upper end of the distribution remain the norm, leading to questions about their informative value. I investigate ratings and behavior of principals and teachers in a novel setting where a district reformed compensation to be entirely based on performance. First, I document consistent increases in teacher subjective evaluations throughout the policy years, with mixed evidence on the relationship of this increase to teacher productivity. I then investigate how incentives specific to the performance pay policy influence ratings assignment. Focusing on a subgroup of teachers with the most to gain from inflated ratings, I find that teachers marginally close to a salary increase are assigned higher ratings unrelated to student achievement. I present suggestive evidence that an attempt by the district to penalize ratings that diverge from student achievement was ineffective in changing ratings behavior using a difference in differences approach. My findings suggest that in spite of inflation, ratings can inform overall teaching performance, but that the structure of the performance pay plan can influence misalignment with student achievement.

JEL Classification: I20, I28, H75, J33, M52

Keywords: Employee Evaluation, Teacher Salaries, Performance Pay, Teacher Productivity

*South Dakota State University Email: andrew.morgan@sdstate.edu.

1 Introduction

A central task of supervisors in and out of education is to assess the performance of and provide feedback to employees. However, use of subjective measures of performance, particularly when used in compensation, often comes with substantial biases. Supervisors are apt to evaluate employees more leniently than their true underlying performance, and evaluation distributions are frequently more compressed, generally coming from only the upper end of the ratings set (Weisberg et al. (2009), Frederiksen, Lange and Kriechel (2017)). In spite of this, subjective evaluations are used extensively across many industries. While there is ample evidence of the use of subjective evaluations of performance, less is known about the informative value of these evaluations, and ways in which firms can induce more accurate subjective ratings, particularly in firms where workers engage in complex, multi-task occupations.

I study how ratings evolve and respond in a relatively unique setting in a large, urban public school system in which principals and teachers face strong financial incentives for their performance, where a majority of teacher yearly compensation is determined by subjective evaluations. Principal pay is similarly determined by school performance, but principals also face constraints on the ratings they assign, which I explore in more detail.

I show how this compensation structure impacts reported ratings and the efficacy and reliability of these ratings in evaluating employees, benchmarking observation ratings to student achievement, using teacher value-added to compare against. While subjective ratings increase sharply over my sample period, ratings do not appear to become less correlated with impacts on student outcomes over time, with correlations similar to other, lower stakes settings.

However, it is less evident to what extent these ratings correspond to changes in employee productivity. I identify key potential mechanisms by which ratings may diverge from observed performance and document some evidence of this occurring. I present difference in discontinuity estimates that show that principals appear to assign ratings that diverge from student outcomes for teachers that face a higher financial incentive from increased ratings. I investigate the ratings assignment behavior from the outset of

this system as well as the response to a penalty in which supervisors receive a potential financial cost for assigning ratings that less closely align with an objective measure of employee quality. Using a difference in differences approach I show that supervisors do not significantly change their reporting behavior in response.

The evidence on the appropriate utilization of subjective employee evaluations is mixed. Use of subjective evaluations has been shown to be associated with higher levels of employee productivity (Gibbs et al. (2004)) and school accountability systems frequently employ subjective evaluations like classroom observations as a way to evaluate employees (NCTQ (2016)). At the same time, evaluations that are more lenient and compressed than the underlying distribution of employee quality has long been cited as a consistent problem in subjective evaluation systems (Prendergast and Topel (1996), Jawahar and Williams (1997), Golman and Bhatia (2012)). Inaccurate ratings can lead employees to not know their true productivity and blunt the incentive effects from performance pay systems, and indeed higher levels of compression have been shown to lead to lower productivity in the private sector (Kampkötter and Sliwka (2018)). These problems may be particularly pronounced in the public sector, where there may be weaker incentives on supervisors to maximize productivity.

Additionally, subjective evaluations have been shown to be both subject to bias from employer favoritism and to significant racial and gender biases, further calling into question the reliability of these types of performance ratings, particularly when ratings determine a significant portion of pay (Prendergast and Topel (1996), Jawahar and Williams (1997), Elvira and Town (2001), Moers (2005), Castilla (2012), Drake, Auletto and Cowen (2019)). When specifically used in evaluating teachers, performance ratings have also been shown to be subject to significant bias from classroom composition, with teachers assigned higher ability students showing much higher evaluation ratings (Whitehurst, Chingos and Lindquist (2014), Steinberg and Garrett (2016)).

Yet the vast majority of firms and state education agencies use some type of subjective performance in evaluating employees (Murphy and Cleveland (1991), NCTQ (2016)). In investigating the determinants and efficacy of subjective evaluations, most prior work

focuses on a narrowly-defined measure of productivity in a singular task or low-skilled environment (Bandiera, Barankay and Rasul (2009), Breuer, Nieken and Sliwka (2013)). Studies in most occupations do not allow for measures of individual productivity where output is determined multi-dimensionally. In contrast, investigating questions of efficacy and determinants of subjective evaluations in an education setting affords me a unique opportunity to construct well-defined measures of individual employee effectiveness – value-added to student test scores.

The proper way to determine effective teaching remains a consistent empirical question. Numerous studies have shown that there is substantial variation in teacher quality (Rockoff (2004), Rivkin, Hanushek and Kain (2005), Kane and Staiger (2008)). Additionally, prior evidence has found that principals are generally *able* to distinguish between low and high effectiveness teachers, but less is known about the decisions that go into how principals choose their ratings for teachers (Jacob and Lefgren (2008), Bacher-Hicks et al. (2019), Harris and Sass (2014)). While value-added has been shown to be an effective measure of teacher effectiveness, value-added measures are only available for a small subset of teachers (Chetty, Friedman and Rockoff (2014)). Thus the vast majority of teachers must be evaluated using alternative approaches. At the same time, value-added may only measure one aspect of student achievement – a student’s test score growth. While we would expect value-added and subjective evaluations to be related, the two measures inherently assess different aspects of productivity. My work fits into the broader literature on teacher evaluations and school accountability and investigates the reliability and potential improvement of the reliability of teacher evaluations.

At the same time, however, most prior work uses subjective evaluations that may not be applicable when evaluations are used in a high stakes context. Because ratings assigned when stakes are high could be more likely to be lenient (Jawahar and Williams (1997)), there could be substantial differences in the ratings assignment when applied to a setting where principal ratings determine compensation. More recent studies have shifted their attention to high-stakes evaluation systems, and the work that has been done generally shows that ratings remain modestly correlated with value-added measures

(Sartain et al. (2011), Kraft, Papay and Chi (2020)). Grissom and Loeb (2017) also study this question and show that evaluations made in low stakes environment differ slightly from evaluations that impact teacher renewal, but in both instances evaluations are quite high and compressed. My work follows more closely in this line of work, where ratings assignment affects both teacher and principal compensation.

The paper proceeds as follows. I detail the institutional background in Section 2. Section 3 then outlines a conceptual framework for supervisor rating assignment and identifies the potential trade-offs supervisors face when assigning ratings. I then proceed to describe the data I use for my empirical approaches in Section 4. I document ratings evolution in Section 5, present evidence on the potential mechanisms by which ratings can diverge from observed performance in Section 6, and conclude in Section 7.

2 Institutional Background

I explore subjective ratings in the broader context of a major compensation reform in a large, urban public school district in the Southwestern United States. Beginning in the 2012-2013 school year, the district undertook an extensive overhaul of the compensation and evaluation system for its teachers and principals. The district first implemented the principal reform and followed it up two years later in 2014-2015 with a parallel program for teachers. Both programs are a broad-based reform of the previous compensation systems with an intense focus on educator development. For both administrators and teachers, the prior compensation system determined primarily by education and experience was replaced with a system that is entirely determined by a formula that incorporated supervisor observations, surveys and student test scores.

The formula to determine compensation assigns points across these three main categories, with different weights for each of the three categories based on a teacher's subject, grade, and population taught. Total evaluation scores are calculated by a weighted combination of these factors. The district then assigns ranges, with fixed proportions of teachers in each range by year, for the two-year average of the total scores where each

range determines salary level. Compensation levels begin at \$47,000 for teachers in the Unsatisfactory level and increase to upwards of \$90,000 for teachers in the Master level. Each level increase corresponds to a \$3,000 and \$7,000 increase in salary. Additionally, once a teacher has reached a specific compensation level, their salary can be no lower for up to three years, and the district stipulates that if a teacher’s total evaluation score corresponds to a lower salary level for three consecutive years can their salary then be lowered.¹ Additionally, teachers receive the higher of their reformed salary or their salary in the 2014-2015 school year – their salary in the year prior to policy change – if available. I exploit this nonlinearity and stickiness in wage acquisition in a regression discontinuity design to identify one possible avenue of misaligned ratings.

The focus of my paper is on the more subjective nature of classroom observations, which constitute the majority of the total evaluation score for all teachers. Teachers in tested grades and subjects and that administer student surveys receive a 50% weighting on classroom observations while teachers without their own student test scores (i.e. grades and subjects not subject to state-standardized test scores) and without student survey scores receive the highest weight of 80%. Notably, observation scores make up a slightly larger portion of total evaluation ratings in the district than in systems with other similar policies (Putman et al. (2018)). This means that a substantial subset of teachers receive points from both a subjective classroom observations rating as well as from student achievement.² For this group of teachers, I estimate value-added to test scores as a way to generate a more objective measure of quality against which to compare more subjective classroom observations.

Principals are evaluated and compensated across a similar scoring system, where their total evaluation score is determined by their own outside observer’s subjective points assignment, as well as parent surveys and test scores of the school they oversee. Additionally, they are required to evaluate teachers by assigning points based on how well the

¹In practice, no teacher in my sample has moved to a lower compensation level in any year of my study.

²Student achievement points towards a teacher’s pay are determined by the maximum of three measures: student pass rate, test score growth compared to a student’s peer group, and a district-defined measure of value-added.

principal believes the teacher adheres to a set rubric, both on a year-long review and from in-classroom observations. These observations are intended to serve two main roles: first, to highlight improvement areas the principal believes the teacher should focus on and second to determine the classroom observations portion of teacher compensation score.

Figure 1 gives an example of the types of standards for which principals are asked to evaluate teachers. In total, principals evaluate teachers on 18 target categories across 4 domains, two of which are used exclusively as year-long reviews rather than strictly in-classroom metrics. For classroom observations, principals observe teachers between 5 and 9 times throughout the school year, decreasing with teacher evaluation level. Principals assign each category a score between 0 and 3 dependent on how well the principal perceives the teacher’s performance on that specific metric to be. Scores assigned to each of the categories are then averaged over all observation periods in the year, and average category scores are combined using a weighted average, where more weight is given to categories used in classroom observations. This process generates a metric for observation scores that ranges between 0 and 100 from which the district assigns points that partly determine evaluation level and compensation. Throughout the paper, I refer to this computed total observation score as the teacher’s “observation” score.

The district has also made it a focus to attempt to increase the accuracy of observations by requiring principals to receive annual training and certification in observation scoring. Additionally, recognizing that subjective measures of worker productivity are potentially prone to inaccuracy and lenient ratings, the district created a component of the principal’s own evaluation score that penalizes principals for having a higher average mismatch between scores they assign teachers in classroom observations and the scores teachers receive from the test score component of their own evaluation score.

While the majority of the principal score is determined in roughly the same way as teachers, there are two notable differences. The first is the penalty calculation mentioned above, and the second is a component determined by teacher development, where principals are rewarded for the average growth in their teachers’ total evaluation scores across years. Principals with the highest average teacher growth, relative to other principals in

the district in that year, receive higher points on this category. The mismatch penalty and the growth component constitute roughly 5% each towards a principal’s total pay.

With the introduction of the penalty metric, the district explicitly defines one way in which principals must adhere to a pre-defined accuracy of ratings. However, teachers without their own student test scores were not (could not be) included in this penalty calculation, allowing me to compare observation scores for two groups of teachers—those who were included and those who were excluded. Additionally, this penalty went into effect the second year of the program (2015-2016), giving me the opportunity to compare score differences for teachers with and without their own student test scores for one time period prior to implementation. This allows me to test the efficacy of this penalty by examining any divergence in the trajectory of subjective evaluations between these two groups in a difference in differences design.

3 Conceptual Framework

Principals face various and potentially competing incentives in assigning ratings. The main task involves observing the performance of their employees and assigning ratings based on the productivity they observe. Principals may prefer to accurately assess performance, but at the same time may care about the well-being and retention of their employees. Interpersonal factors based on the long-term nature of a supervisor-employee relationship and costs associated with replacing teachers may lead to principals assigning ratings that are more lenient. For instance if a principal would in isolation assign an accurate, but low, rating to a teacher, they may choose instead to assign a rating that is more lenient and higher if they face interpersonal pressure from a teacher, especially if they believe that they would incur high costs associated with replacing that teacher.³

Following the prior literature, I model supervisors as having a distaste for ratings that less closely align with observed performance (defined as “inaccurate”), coming from two

³This section narratively follows the models described in prior literature on supervisor ratings. See Kampkötter and Sliwka (2018), Golman and Bhatia (2012), and Jawahar and Williams (1997) for more in depth discussion of modeling supervisor accuracy behavior.

primary factors. First, inaccuracy likely affects outcomes for the school and district. If teachers are not rated accurately, this impedes their ability to effectively develop skills and weakens the alignment between performance incentives and productivity. If pay is determined by school outcomes (as it is in the principal’s case), the supervisor may suffer from lower productivity stemming from less well-developed workers. Second, to the extent that accuracy is verifiable by the district, principals may face reprimand if their reporting deviates too far from true performance.

At the same time, supervisors also likely care about the well-being of their employees, especially so when ratings are directly tied to pay, and may face large costs in replacing employees if low ratings induce exit.⁴

I explore this interpersonal component of this framework in 6.1. I identify a subset of teachers that should have the highest likelihood to receive potential special treatment – those that just missed hitting a higher compensation level the year prior. If principals take expected pay into account, we might observe a higher observation score for the teachers who just missed out on hitting a higher compensation level relative to the teachers who did not. I then compare this to the change in student achievement to determine if ratings changes correspond with changes in student outcomes, or solely higher subjective ratings.

A rather unique aspect of the reform I study is the district-wide implementation of a penalty component for inaccurate ratings from principals. Naturally, this penalty should serve to incentivize principals to more accurately rate teachers. Given prior evidence of high ratings, I would expect to see that any change in accuracy would likely come from a decrease in ratings. In my analysis in 6.2, I exploit the fact that the penalty calculation was introduced later in the program and that some teachers were not included in the penalty calculation to explore the effect of this penalty on rating behavior.

⁴See Prendergast and Topel (1996), Prendergast (2002), Sliwka (2007), and Giebe and Gürtler (2012) for a discussion on employee favoritism and interpersonal relationships in the firm.

4 Data

The panel nature of my administrative data allows me to generate precise measures of individual subjective and objective performance. Because I also know the date, score, supervisor, and employee for each observation period, I am also able to construct differences in observation scores by the number of times a supervisor observes an employee. I use files from the district's administrative systems to construct teacher and principal panel data sets for the school years 2012-2013 to 2018-2019. These data include detailed measures of each of the components and sub-components of the evaluation systems as well as district-calculated total scores from each of these measures. The data also include the raw data needed to calculate these measures, including data from each observation period for each teacher and an identification to link each observation period to an observing principal.

I merge these data with data from the state administrative system containing employment and demographic information for all staff in the district. State data detail where and in what capacity each district employee serves in each year, with demographic information on each employee's years of professional experience, sex, and race and ethnicity. Using these data, I am also able to link each teacher to the students that they teach in each year.

These data also contain a rich panel of demographics and test scores for each student in the district in each year. They also list classroom enrollment for each of the roughly 150,000 students in the district in each year of my sample period. I merge these to student demographic information containing sex, race and ethnicity, free and reduced price lunch status, limited English proficiency, as well as special education status. Using these data, I am able to link students to teachers in order to construct value-added to student test scores for each teacher to generate a traditionally more objective measure of teacher evaluations.

Student test scores come from the state standardized tests, which were administered to every student in grades 3-8 reading and math. I standardize all student test scores to the state averages of each grade and year, so that each score is normalized around the state mean of 0 and standard deviation 1 within grade and year.

Virtually all (93%) of the approximately 10,000 teachers in the district in each year are observed in classroom and assigned a total evaluation score in each year. Teachers that are not required to be evaluated include certain types of special education teachers and teachers classified as guest or substitute teachers. Each teacher that is evaluated is required to be observed between 5 and 9 times throughout the year, with teachers with lower evaluation levels required to receive more observations.

Table 1 presents information on teacher evaluation components as well as demographic information for two groups of teachers. Because of data limitations, I can only identify the total number of observations for the years 2014-2015 to 2017-2018, and so this table excludes statistics for teachers in 2018-2019. Columns 1 and 2 break down these summary statistics by the two analysis groups I use in this paper – all teachers with evaluation scores and a subgroup that contains both evaluations and for which I can construct test score value-added measures.⁵ The second group is comprised of grades 3-8 reading and math subject teachers. Across all years, the average total observation score is roughly 73 out of 100, coming from an average of 8 total observations per year for both groups of teachers. The value-added sample is slightly less experienced with an average of around 9 years of experience compared to 10 years for the full sample. Approximately 70-80% of teachers are female, and a little under a third hold an advanced degree. A little over one-third of teachers in both samples identify as Black or African American, and a little under a third each identify as white or as Hispanic.

5 Documenting Ratings Evolution

Given the unique policy setting in the district, I first document the evolution of the observation scores over time and compare to prior work. I find a pattern of consistently increasing observation scores, indicative of either increasingly more lenient ratings, more productive employees, or both. To determine the degree to which observation scores correspond to changes in student outcomes, I then document the relationship between

⁵Roughly 97% of teachers for whom I can construct value-added have evaluation scores in each year—for roughly 80 teachers out of approximately 2500 in each year.

observation scores and value-added throughout my sample period, finding a modest relationship that persists across years but with potential differences for teachers with different levels of experience.

Tying pay to evaluations can create strong incentives for leniency and at the same time could also improve student outcomes if teachers are incentivized to develop more skills and/or exert more effort. Indeed, taking the district as a whole, Hanushek et al. (2023) use a synthetic control approach to show that student achievement increases relative to comparable districts in the state. Because an increase in a teacher’s ratings may reflect improvement to teaching or be reflective of lenient and inaccurate ratings, it is important to explore how changes in observation scores correspond to changes in student outcomes at the teacher level to determine how indicative ratings are for teacher effectiveness.

I first document the time trends in observation scores I observe. Figure 2 shows that the broad distributions of observation scores continually shift towards the right since the start of the teacher reform in the 2014-2015 school year. By the last year for which I have data, the most common value of total observation scores hovers around 95 points (out of 100). Table 2 shows the mean observation ratings is increasing across all years of my data. While the overall standard deviation of observation scores does not substantially change over my sample period, ratings are still rather compressed, evident in less willingness on the part of principals to rate employees using the lower ends of the ratings scale. Across years, I observe little usage of the two lowest ratings. Table 2 shows that roughly 56% of teachers in the first year receive evaluations corresponding to an average score within the two highest levels (66 points), meaning that on average very few teachers receive a “progressing” or “unsatisfactory” average rating across all metrics. This number increases to over 80% of teachers in the last year of my study period. Additionally, an increasing number of teachers show virtually no avenue for improvement, receiving an evaluation score of 95 points or higher. In the latter years, roughly 20% of teachers receive scores at or above 95 points.

We can also see that this change has happened slowly and progressively throughout the sample period. The mean growth in a teacher’s observation score year-on-year is similarly

high. Table 3 shows the changing distribution of the growth of observation scores over time. The mean teacher receives a modest increase in score each year, ranging from around 6 points to a little over 3 points in the last year of my sample. The change in observation scores remains right skewed throughout the years, with roughly 25% of teachers receiving double-digit increases in observation score in any particular year. However, there does exist a sizable fraction of teachers with any decrease in score from the year prior. Roughly 30% of teachers in any year receive a lower observation score than in the year prior. Decreases remain small relative to increases, however, leading to an overall continuous increase in scores.

To document the degree to which observation scores may diverge from student outcomes, I estimate a yearly value-added to each teacher’s students’ math or reading test scores to compare against. I construct value-added estimates for each teacher in the following way:

$$A_{igjt} = f(A_{ijgt-1}) + X\beta + \delta_{gt} + \gamma_{jt} + \varepsilon_{igjt} \quad (1)$$

A_{igjt} measures each student’s i math or reading test score in grade g at time t with teacher j . I regress this on a cubic in past student achievement and control for student demographics, X , including student sex, race and ethnicity, free and reduced price lunch status, special education status and limited English proficiency. X also includes school-by-year means of each of these student demographic variables. I additionally control for grade-by-year fixed effects, δ_{gt} . γ_{jt} captures the teacher-by-year fixed effect and represents the teacher value-added to student achievement in each year.

The use of value-added also allows me to have a reliable measure of individual objective performance, and I directly measuring how related subjective measures are to objective measures of individual performance in this high stakes context. Table 4 shows that overall correlations between observations and value-added are positive and modest in magnitude. Correlations between both reading and math value-added fluctuate throughout the years, hovering around 0.2 to 0.3, with math being slightly more positively correlated than reading—roughly in line with what prior work has found (Jacob and Lefgren (2008),

Kraft, Papay and Chi (2020)). Notably, for both subsets of teachers, correlations do not appear to exhibit an increasing or decreasing pattern, suggesting that there is no rank order shifting across the years – that principals remain fairly consistent in rating the most effective teachers more highly.

A concern with subjective evaluations in a system where there may be a strong incentive towards lenient ratings is that scores may be inflated and thus be less useful in determining a teacher’s future productivity. Given the trend towards higher ratings I observe, I follow a similar procedure to that outlined in Bacher-Hicks et al. (2019) and Rockoff and Speroni (2010) in identifying the predictive validity of subjective observations on future student outcomes. If increases in ratings are becoming less predictive of student achievement, we would expect to see the relationship between predicted observation scores and student achievement fall.

To determine if this is the case, I regress outcomes of students matched to their teacher in a particular year on predicted yearly-standardized observation scores. To generate predicted observations scores and to account for measurement error and fluctuations in measurements, I follow the procedure outlined in the following three regressions. I first regress the two-year prior year observation score $Score_{jt}$ on the score from the one-year prior, $Score_{jt-1}$ (Equation 2). I capture this coefficient, $\hat{\beta}$, and multiply it to the observation score from t to generate a predicted observation score for teacher j in the following year, \widehat{Score}_{jt+1} (Equation 3). Finally, I match students to the teachers they have in each year, and regress student outcomes in year $t+1$ on predicted observation score, controlling for student demographics, prior achievement scores and prior absences and discipline in each regression as well.

$$Score_{jt} = \beta Score_{jt-1} + \varepsilon_{jt-1} \quad (2)$$

$$\widehat{Score}_{jt+1} = \hat{\beta} * Score_{jt} \quad (3)$$

$$Outcome_{ijt+1} = \delta \widehat{Score}_{jt+1} + X\Omega + \varepsilon_{ijt+1} \quad (4)$$

Table 5 shows the results from this process. If classroom observations less accurately

predict student outcomes over time, we should expect to see the coefficient on predicted observations decreasing over time. However, this does not appear to be the case. The coefficients on math and reading test scores in the top and bottom panels, respectively, show that the coefficient for predicted score in 2016-2017 (Column 1) is no smaller for math and reading test scores than the coefficient in 2018-2019 (Column 3). Having a teacher that is one standard deviation higher in predicted observation score increases math test scores by around 0.11 standard deviations and around 0.073 standard deviations for reading teachers. Observation scores overall appear to be moderately correlated with student achievement, with stable predicted estimates over time.

I also document a general increase in ratings that is coming from across all experience levels of teachers. Figure 3 shows the average observation scores for all teachers in the district by experience level for a given year. In each year, ratings are higher for teachers with higher levels of experience up to roughly 5 years of experience and then level off, similar to what is observed in experience profiles to value-added where, in the early years, there exists an increase in the evaluation metric and then a leveling off after 5-10 years. Note, however, that these plots do not plot the returns to experience per se, but rather the broad means of observation scores for teachers with a given level of experience in each year.⁶ Increases appear to be occurring in all experience levels, with teachers with any given number of years of experience having higher average observation scores than similarly experienced teachers in any prior year.

There is also some evidence that teachers at higher levels of experience have higher ratings increases than teachers at lower levels of experience. Figure 4 plots the difference between the 2018-2019 school year ratings and the 2014-2015 ratings for teachers with a given level of experience in each year. While score increase is positive for all experience levels, the difference is increasing nearly throughout the experience distribution. Estimates for teachers past roughly 15 years of experience fluctuate, but are generally higher than teachers with 10 years of experience, and the difference for teachers with 10 years of experience is generally higher than teachers with 5 or fewer years. While observation

⁶I cap experience at 25 years, but the pattern holds when including teachers with up to 40 years of experience. This constitutes around 95% of all teachers.

scores are not expected to be perfectly aligned with productivity measured by student test scores, this is somewhat contrary to what we would expect of teacher professional development from the literature on value-added.

To the extent that observations and student achievement are related, if increases in teacher evaluations by experience level correspond to increases in student achievement, we should expect to see the difference in value-added between the two years for more experienced teachers be higher than the difference in value-added for less experienced teachers as well. This should especially be the case given that the district has placed a strong focus on improvement to student achievement. I present comparable experience figures to Figure 4 for yearly math and reading value-added to determine if the difference in value-added is changing at a commensurate level to observation scores. However, Figures 5 and 6 show differences in value-added by experience with no similar pattern. The change in value-added for teachers with higher levels of experience is no different than the change in value-added for teachers with lower levels of experience. Together, these suggest that observation scores are increasing for more highly experienced teachers with no noticeable increase in value-added for these same teachers. Thus, it appears to be the case that more experienced teachers may be receiving ratings that do not accurately reflect improved student achievement, relative to less experienced teachers.

Potential explanations of this relationship with experience could include an increased pressure principals face from more experienced teachers holding higher influence in the school or from higher professional development. Kalogrides, Loeb and Beteille (2013) find that higher teacher experience is associated with higher prior student achievement within a school, suggesting that more experienced teachers may be better able to influence a principal when it comes to classroom assignment. At the same time, Kraft, Papay and Chi (2020) document a similar increase in observation scores in a returns to experience model and attribute the increase to professional development not measured by test scores, but can only examine this question for teachers with 10 or fewer years of experience due to data limitations.

Taken together, use of subjective evaluations in this high stakes context does not

appear to be driving a divergence between classroom observations and student outcomes overall, though there is some suggestive evidence that this differs by teacher experience. Observation scores are moderately related to value-added, and they appear to not be losing predictive power over time; however magnitudes are relatively modest in each case. Given the steady growth in scores over time, it is important to determine the ways in which principals may be rating teachers that may or may not reflect true productivity growth. In the next section, I turn to directly investigating specific ways in which ratings may be more likely to be applied inaccurately.

6 Divergence between and Ratings and Performance

Given the incentive structures faced by both principals and teachers, we might expect leniency to differ for different groups of teachers. Principals may face a financial penalty for inaccurately rating teachers included in the penalty, but face no such trade-off for the teachers excluded from the calculation. Teachers who just missed out on hitting a compensation level in the prior year may have strong incentives surrounding the receipt of ratings in the following year, relative to those who achieved a higher compensation level. Investigating the magnitude of these factors is necessary in determining the significance of the incentives principals and teachers face in ratings behavior and assignment.

I explore two main empirical approaches to identify the extent and determinants of observation score bias in teaching. First, to investigate the effect of strong financial incentives on ratings, I present regression discontinuity estimates that compare subjective and objective measures of teacher performance for teachers who fall just to the left or right of a compensation cutoff the year prior. If principals are assigning more lenient ratings due to the discontinuous effects on teachers' salary, we would expect to see an increase in ratings for those teachers just to the left of the cutoff relative to those just to the right in the year prior, with no comparable increase in student achievement. Second, to estimate if principals respond to the penalty by assigning more accurate ratings, I construct difference in differences estimates, comparing observations ratings of teachers

who were and were not included in a principal’s penalty calculation before and after the introduction of the penalty. If the penalty was effective in reducing leniency, we should see this manifest as lower observation scores relative to teachers who were not included, after the introduction of the penalty

6.1 Teacher Financial Incentives Can Alter Ratings

I present evidence that suggests that principals may be more inclined to assign more lenient ratings if they believe this could help increase a teacher’s pay a meaningful amount. Since teachers are only subject to a pay increase if their score puts them in a higher compensation level, it may be the case that teachers and principals have the strongest incentive to be more lenient for teachers who just missed out on attaining a higher compensation level in the year prior. Teachers in this group should have the highest incentive to lobby for more lenient scores, put in higher effort to pass the threshold through higher student achievement and teacher development, or both. Additionally, because teachers receive a salary protection for at least three years after achieving a compensation level, teachers just to the right of a compensation cutoff should have a discontinuously lower incentive than teachers just to the left.

To test this, I estimate regression discontinuity estimates for teachers just around a compensation threshold for their observation scores and math and reading value-added in the next year. While there is some incentive no matter where a teacher is in the distribution of ratings for continuous improvement, the incentive for lenient ratings should be much less strong for teachers who are “far” enough away from a compensation cutoff. I present estimates showing that while value added does not appear to be discontinuously changing for teachers around a compensation cutoff, observation scores do differ by a modest amount.

At the same time, other factors may also be changing around a threshold. Teachers that do not achieve a higher level in the year prior may lose motivation if they believe their low rating does not reflect their true ability. Principals may also consciously or subconsciously rate based on a teacher’s evaluation level, such that principals rate teachers

who are just above a threshold substantially higher than those just below solely because they have a higher rating level.⁷ If this is the case, principals may be rating teachers more highly simply based on a recognition of higher evaluation level. Because of this, we might expect to see a higher rating for teachers just above the cutoff relative to those below, in the absence of other factors.

To account for these factors, I present difference in discontinuity estimates using teachers who did and did not have grandfathered pay, following Grembi, Nannicini and Troiano (2016), to isolate the effects of salary changes around the thresholds. Teachers without a financial incentive would likely still be subject to demotivating factors, and principals' ratings could still reflect asymmetric information for this group as well, so that taking the difference between these groups differences out these effects.

Specifically, I compare discontinuities for two groups of teachers, those with and without grandfathered salary from prior to the implementation of the reform. I estimate the following regression

$$\begin{aligned}
C_{jt+1} = & \beta_0 \textit{above}_{jt} + \beta_1 \textit{point}_{jt} + \beta_2 \textit{above} * \textit{point}_{jt} \\
& + \beta_3 \textit{fincentive}_{jt} + \beta_4 \textit{fincentive} * \textit{point}_{jt} \\
& + \beta_5 \textit{fincentive} * \textit{above}_{jt} + \beta_6 \textit{fincentive} * \textit{above} * \textit{point}_{jt} + \varepsilon_{jt+1}
\end{aligned} \tag{5}$$

where C represents either the observation score or math or reading value added of teacher j in year $t + 1$. I regress this on whether or not a teacher's score in the prior year placed them in a higher compensation level, *above*, the centered distance from a threshold value in the year prior, *point*, and the interaction between these variables. The variable *fincentive* _{jt} is an indicator for whether or not a teacher has a financial incentive for hitting a compensation level in the next year, i.e that they did *not* have grandfathered pay that was higher than the next compensation level they face. β_5 represents the difference in the

⁷Cullen et al. (2016) show some evidence of this occurring for educators in a similar system, where principals just rated "unacceptable" face disproportionately higher labor market penalties than otherwise equal principals.

discontinuity between these two groups, isolating the financial incentive faced by teachers and eliminating potential other factors changing around this cutoff. I estimate this using a bandwidth of 4 points on either side of the cutoffs.

Results of the difference in discontinuities regressions reveal a significant and substantial drop in observation scores just after the threshold for teachers with a financial incentive relative to the discontinuity for those without an incentive. Column 1 of Table 6 shows that there is a 5 point decrease in observation scores for teachers who just hit a compensation level in the year prior, with no corresponding significant change in value-added (columns 2 and 3). This disconnect between reported performance and observed student outcomes suggests that teachers receive more lenient ratings as a result of being “close” to receiving more pay, and that student achievement growth is no different on either side of the threshold.

I also present results estimating the discontinuities around thresholds separately between the two groups of teachers who just missed out on attaining a compensation level in the prior year. Specifically, I estimate two regression discontinuity equations in the following way

$$C_{jt+1} = \beta_0 above_{jt} + \beta_1 point_{jt} + \beta_2 above * point_{jt} + \varepsilon_{jt+1} \quad (6)$$

where dependent and independent variables are as defined above. β_0 represents the linear approximation of the discontinuity for just achieving a higher compensation level on a teacher’s evaluation components in the next year. Main estimates use a bandwidth of 4 points and I present estimates for multiple bandwidths as well.

Figure 7 presents regression discontinuity plots of classroom observation scores for teachers who were not subject to grandfathered pay around the threshold and thus had a financial incentive (Figure 7a) and teachers who did not have a financial incentive (Figure 7b). We see some evidence in Figure 7a that teachers within 4 points of a cutoff on the left of a threshold have slightly higher observation scores in the next year after missing the cutoff than teachers just to the right of the threshold. Table 7, Columns 1 and 2 of Panel A show that the discontinuity with and without controlling for teacher

experience, degree status and sex results is an insignificant roughly 1.5 point decrease after the threshold. In contrast, teachers on the right side of the threshold and without a financial incentive in fact have a much higher observation score than teachers who missed the cutoff, potentially indicating that principals do assign higher ratings based on evaluation levels. Table 7, Columns 1 and 2 of Panel B shows these teachers receive a substantial but insignificant increase of roughly 3.5 points.

The comparable estimates for value-added do not show similar changes. When turning to Figures 8 and 9 we see small, insignificant and inconsistent results for math and reading value-added. These results are not subject to information asymmetry from evaluation levels across the threshold, so we should not expect to see differences between teachers with and without a financial incentive. Teachers just to the left of the cutoff, both with and without financial incentive appear to have slightly higher but insignificant math value-added, with estimates from the full model in Column 4 of Table 7 showing a 0.008 and -0.174 change in math value-added for teachers with and without a financial incentive, respectively. The case is reversed for reading. Figure 9 and Column 6 of Table 7 show an estimated discontinuity of around 0.04 and 0.07 for reading value-added for the two groups. While there may be some slight evidence of increased observation scores, no such pattern holds for math or reading value-added.

Table 8 shows that these results hold for varying bandwidths, but remain noisy. Columns 1-3 show results for observation scores, math and reading value-added for teachers who were within 5 points of the threshold in the year prior, and Columns 4-6 show for these for a bandwidth of 3. Both sets of results suggest a slight, insignificant fall in subjective evaluation scores for teachers just after the threshold, with no clear change in value-added for teachers with a financial incentive, with a slight, insignificant increase in observation scores and no change in value-added for teachers without a financial incentive.

A concern with any regression discontinuity design is that individuals can perfectly manipulate the side of the threshold onto which they fall, leading to biased estimates from incomparable individuals on either side of the cutoff. In this instance, while I make the case that individual teachers and principals may try to manipulate scores around a

cutoff, they do not appear able to perfectly influence which side of the threshold they fall. Given that thresholds for each compensation level varies from year to year and is assigned by the district, teachers (and principals) cannot know ahead of time *precisely* how much they must increase their scores from the prior year. The empirical evidence also bears this out. Figure 10 shows no discontinuity in the number of teachers on either side of the threshold for both groups of teachers and Table 9 shows virtually no change in observable characteristics for these teachers as well. Neither experience nor advanced degree status changes meaningfully at the cutoff for either group of teachers, but advanced degree status does have a marginally significant discontinuity around the cutoff for teachers with a financial incentive.

Taken together, the difference in discontinuity estimates suggest that teachers with a financial incentive do receive meaningfully higher subjective ratings, which does not transfer to student achievement. Separate regression discontinuity estimates provide some evidence that principals may be more inclined to assign ratings differently solely based on evaluation levels, with again little difference in student achievement between groups of teachers.

6.2 Principal Penalties Don't Alter Behavior

The district recognized that ratings are likely to be influenced by factors other than observed performance, as indicated by ratings that are compressed in the upper end of the distribution in other settings, and so instituted a penalty mechanism intending to tie principals ratings to more objective and observable measures of teacher performance. To investigate the effect of the penalty I exploit the fact that the penalty only came into effect in the second year of the program and that only a subset of teachers were included in the penalty calculation, though all teachers receive subjective ratings. I calculate difference in differences estimates for each year, comparing teachers who were and were not included in the calculation before and after the introduction of the penalty. This roughly corresponds to exploring the difference in each year between teachers who do and do not teach in tested grades and subjects, before and after the penalty was introduced.

Because these two groups of teachers face slightly different pay calculations and work in different classroom settings, their baseline level of observations may differ. To the extent this matters, this should be captured by the level differences in the period before the penalty was introduced. If the penalty was effective in preventing higher and more lenient ratings, we might expect to see a divergence between the two groups of teachers, with slower growth for teachers who are included in the penalty calculation. The incentive structure for all teachers is such that principals may be inclined to assign lenient ratings to both groups, but because of the penalty, principals are constrained in some way in assigning ratings that are lower for teachers with student test scores.

Specifically, I estimate

$$score_{jt} = \beta tested_group_j + \theta_t + \sum_{t=2015}^{2019} \theta_t * tested_group_j + \varepsilon_{jt} \quad (7)$$

where *tested_group* equals 1 for teachers *j* who were included in a principal’s penalty calculation (taught in tested grades and subjects in year *t*), and 0 otherwise. θ represents year fixed effects and the interactions between each θ and the *tested_group* dummy variable represents the effect of being in the group contributing to the penalty in each year. Some estimates also include teacher level demographic controls for years of professional experience and having an advanced degree.

Examining average observation scores between these two groups of teachers suggest that principals do not appear to be differentially altering their rating behavior for the teachers that were and were not included in the penalty calculation. Figure 11 plots the average observation scores for teachers with (blue line) and without their own student test scores (red line).⁸ The differences between the two lines after the introduction of the penalty in the 2015-2016 school year do not appear to be meaningfully different from the year in which the penalty was not in place. The difference between the observation scores for teachers that were and were not included in the penalty calculation is highest in the

⁸Note that this sample differs slightly from the value-added sample used in the other portions of this paper. The district calculates the penalty based on points gained from all teachers with their own student achievement—not just those for whom value-added can be calculated – and thus constitutes a broader group of teachers.

2015-2016 school year, where teachers that were not included in the penalty calculation have an approximately 1 point higher average score, but all differences are small in every other year.

Difference in differences estimates in Table 10 show the average observation score difference for each year relative to the year before the penalty was introduced. The estimates imply that observation scores for each year are not significantly changed between the two groups after the penalty. The only significant effect comes in 2017, the sign of which implies, if anything, that teachers with their own student test scores increase in average observation score relative to the pre-period, counter to the incentive structure of the penalty. Changing teacher demographics does not explain this result, as after accounting for teacher demographic controls in Column 2 the estimates remain small and largely insignificant. This evidence suggests that principals are not meaningfully taking the penalty into account when rating teachers.

Principals could be rating the teacher groups differently for reasons other than the penalty metric. By differencing out the ratings difference between the two groups before the penalty was introduced, the effect of this should be mitigated. However, because observation data does not exist prior to the 2014-2015 school year, I am not able to fully test the existence of such trends between the groups prior to the penalty introduction, and rely on only one pre-year as a sufficient estimate of the difference in the absence of the penalty. However, the fact that there is minimal difference between the groups prior to the penalty suggests that principals treat these groups of teachers similarly with or without an incentive to do so.

The evidence presented in this section suggests that careful attention should be paid to designing performance pay systems. The ratings principals assign appear to be more misaligned with performance in the presence of strong financial incentives their employees face and principals do not appear responsive to a modest financial penalty influencing their ratings behavior. It is important for policy-makers to account for the different incentives faced by both principals and teachers, and to consider the different avenues by which observations can deviate from objective measures of performance.

7 Conclusion

Given the increased focus on teacher evaluations and school accountability, I show that careful consideration should be taken when designing subjective evaluation systems to account for the incentives faced by both supervisors and employees. Understanding the ways in which supervisors assign ratings is key to implementing a system that aligns incentives between the district, supervisors, and teachers while still allowing for ratings that provide meaningful feedback.

I add to the growing literature on the use and determinants of subjective evaluation in the education sector, speaking to the concerns in these systems of lenient ratings that may less accurately reflect true productivity. I find that implementing a system of subjective evaluations tied to performance compensation in a large, urban public school system produced evaluations that grew sharply over time, though similarly to ratings in lower-stakes settings. At the same time, student achievement grew considerably during this time period, though evidence is more mixed on if this increase fully corresponds to increased employee productivity. While subjective evaluations remain predictive of student achievement over time, I present evidence that suggests that this may differ for more experienced teachers.

Supervisors may also face strong incentives towards ratings that less accurately reflect employee productivity in certain instances. Ratings are significantly higher when teachers face strong financial incentives from their ratings, with no corresponding increase in other measures of performance. Teachers who just miss out on hitting a higher compensation level in the year prior have significantly higher observation scores in the next year compared to teachers who did attain a higher level, with no significant change in student achievement.

When ratings do diverge, altering the behavior for supervisors may also be difficult given the other incentives principals face to retain teachers and maintain school performance. Modest financial penalties appear to have no impact in changing the way ratings are assigned. If the focus of the district is to identify effective educators and maintain student achievement, designing a system that penalizes principals for ratings that diverge

from district-defined measures may not be effective when principals and teachers also face strong financial incentives from these systems. One potential solution is to observe teachers solely using outside observers, though that may come with other concerns, such as outside observers potentially being less knowledgeable and able to offer meaningful feedback.

References

- Bacher-Hicks, Andrew, Mark J Chin, Thomas J Kane, and Douglas O Staiger.** 2019. “An Experimental Evaluation of Three Teacher Quality Measures: Value-added, Classroom Observations, and Student Surveys.” *Economics of Education Review*, 73: 1–15.
- Bandiera, Oriana, Iwan Barankay, and Imran Rasul.** 2009. “Social Connections and Incentives in the Workplace: Evidence From Personnel Data.” *Econometrica*, 77(4): 1047–1094.
- Breuer, Kathrin, Petra Nieken, and Dirk Sliwka.** 2013. “Social Ties and Subjective Performance Evaluations: An Empirical Investigation.” *Review of Managerial Science*, 7(2): 141–157.
- Castilla, Emilio J.** 2012. “Gender, Race, and the New (Merit-Based) Employment Relationship.” *Industrial Relations (Berkeley)*, 51: 528–562.
- Chetty, Raj, John N Friedman, and Jonah E Rockoff.** 2014. “Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-Added Estimates.” *The American Economic Review*, 104(9): 2593–2632.
- Cullen, Julie Berry, Gregory Phelan, Eric A Hanushek, and Steven G Rivkin.** 2016. “Performance Information and Personnel Decisions in the Public Sector: The Case of School Principals.”
- Drake, Steven, Amy Auletto, and Joshua M Cowen.** 2019. “Grading Teachers: Race and Gender Differences in Low Evaluation Ratings and Teacher Employment Outcomes.” *American Educational Research Journal*, 56(5): 1800–1833.
- Elvira, Marta, and Robert Town.** 2001. “The Effects of Race and Worker Productivity on Performance Evaluations.” *Industrial Relations (Berkeley)*, 40(4): 571–590.

- Frederiksen, Anders, Fabian Lange, and Ben Kriechel.** 2017. "Subjective Performance Evaluations and Employee Careers." *Journal of Economic Behavior and Organization*, 134: 408–429.
- Gibbs, Michael, Kenneth A Merchant, Wim A. Van der Stede, and Mark E Vargus.** 2004. "Determinants and Effects of Subjectivity in Incentives." *The Accounting Review*, 79(2): 409–436.
- Giebe, Thomas, and Oliver Gürtler.** 2012. "Optimal Contracts for Lenient Supervisors." *Journal of Economic Behavior and Organization*, 81(2): 403–420.
- Golman, Russell, and Sudeep Bhatia.** 2012. "Performance Evaluation Inflation and Compression." *Accounting, Organizations and Society*, 37(8): 534–543.
- Grembi, Veronica, Tommaso Nannicini, and Ugo Troiano.** 2016. "Do Fiscal Rules Matter?" *American Economic Journal. Applied Economics*, 8(3): 1–30.
- Grissom, Jason A., and Susanna Loeb.** 2017. "Assessing Principals' Assessments: Subjective Evaluations of Teacher Effectiveness in Low- and High-Stakes Environments." *Education Finance and Policy*, 12(3): 369–395.
- Hanushek, Eric A, Jin Luo, Andrew J Morgan, Minh Nguyen, Ben Ost, Steven G Rivkin, and Ayman Shakeel.** 2023. "The Effects of Comprehensive Educator Evaluation and Pay Reform on Achievement." National Bureau of Economic Research Working Paper 31073.
- Harris, Douglas N, and Tim R Sass.** 2014. "Skills, Productivity and the Evaluation of Teacher Performance." *Economics of Education Review*, 40: 183–204.
- Jacob, Brian A, and Lars Lefgren.** 2008. "Can Principals Identify Effective Teachers? Evidence on Subjective Performance Evaluation in Education." *Journal of Labor Economics*, 26(1): 101–136.

- Jawahar, I. M, and Charles R Williams.** 1997. "Where All the Children Are Above Average: The Performance Appraisal Purpose Effect." *Personnel Psychology*, 50(4): 905–925.
- Kalogrides, Demetra, Susanna Loeb, and Tara Beteille.** 2013. "Systematic Sorting: Teacher Characteristics and Class Assignments." *Sociology of Education*, 86(2): 103–123.
- Kampkötter, Patrick, and Dirk Sliwka.** 2018. "More Dispersion, Higher Bonuses? On Differentiation in Subjective Performance Evaluations." *Journal of Labor Economics*, 36(2): 511–549.
- Kane, Thomas J, and Douglas O Staiger.** 2008. "Estimating Teacher Impacts on Student Achievement: An Experimental Evaluation."
- Kraft, Matthew A, John P Papay, and Olivia L Chi.** 2020. "Teacher Skill Development: Evidence from Performance Ratings by Principals." *Journal of Policy Analysis and Management*, 39(2): 315–347.
- Moers, Frank.** 2005. "Discretion and Bias in Performance Evaluation: The Impact of Diversity and Subjectivity." *Accounting, Organizations and Society*, 30(1): 67–80.
- Murphy, Kevin R, and Jeanette N Cleveland.** 1991. "Performance Appraisal: An Organizational Perspective." Allyn and Bacon.
- NCTQ.** 2016. "State Evaluation Briefs." National Council on Teacher Quality.
- Prendergast, Canice.** 2002. "Uncertainty and Incentives." *Journal of Labor Economics*, 20(S2): S115–S137.
- Prendergast, Canice, and Robert H. Topel.** 1996. "Favoritism in Organizations." *The Journal of Political Economy*, 104(5): 958–978.
- Putman, Hannah, Elizabeth Ross, Kate Walsh, Kelli Lakis, and Kency Nittler.** 2018. "Making a Difference: Six Places where Teacher Evaluation Systems Are Getting Results." National Council on Teacher Quality.

- Rivkin, Steven G, Eric A Hanushek, and John F Kain.** 2005. "Teachers, Schools, and Academic Achievement." *Econometrica*, 73(2): 417–458.
- Rockoff, Jonah E.** 2004. "The Impact of Individual Teachers on Student Achievement: Evidence from Panel Data." *The American Economic Review*, 94(2): 247–252.
- Rockoff, Jonah E, and Cecilia Speroni.** 2010. "Subjective and Objective Evaluations of Teacher Effectiveness." *The American Economic Review*, 100(2): 261–266.
- Sartain, Lauren, Sara Ray Stoelinga, Eric R Brown, Stuart Luppescu, Kavita Matsko, Frances Miller, Claire Durwood, Jennie Jiang, and Danielle Glazer.** 2011. "Rethinking Teacher Evaluation in Chicago: Lessons Learned from Classroom Observations, Principal-Teacher Conferences, and District Implementation." Consortium on Chicago School Research.
- Sliwka, Dirk.** 2007. "Trust as a Signal of a Social Norm and the Hidden Costs of Incentive Schemes." *The American Economic Review*, 97(3): 999–1012.
- Steinberg, Matthew P., and Rachel Garrett.** 2016. "Classroom Composition and Measured Teacher Performance: What Do Teacher Observation Scores Really Measure?" *Educational Evaluation and Policy Analysis*, 38(2): 293–317.
- Weisberg, Daniel, Susan Sexton, Jennifer Mulhern, and David Keeling.** 2009. "The Widget Effect." *The Education Digest*, 75(2): 31–.
- Whitehurst, Grover J, Matthew M Chingos, and Katharine M Lindquist.** 2014. "Evaluating Teachers with Classroom Observations: Lessons Learned in Four Districts." Brookings Institution.

Figures

2.3 DELIVERY: Facilitates clear, cohesive, and purposeful learning experiences				
ESSENTIAL TEACHER SKILLS & ACTIONS	Exemplary (3 Points)	Proficient (2 Points)	Progressing (1 Point)	Unsatisfactory (0 Points)
<ul style="list-style-type: none"> Supports objectives, prior learning, and all student populations based on their subject, grade, and level with appropriate instructional strategies Delivers content clearly, accurately, and coherently Incorporates appropriately varied digital and/or print and/or hands-on instructional resources Emphasizes the value and connection of content to overall learning and prior knowledge Combines differentiated and relevant instructional strategies and questioning techniques to maintain appropriate pace and engagement 	<p>Consistently and effectively presents the content:</p> <ul style="list-style-type: none"> Logically, coherently, and in a grammatically correct fashion Supporting the learning of the posted objective(s) Building on content previously mastered Supporting all student populations based on their subject, grade, and level Supporting cross-curricular learning Allowing for student input <p>Consistently and appropriately uses multiple, differentiated:</p> <ul style="list-style-type: none"> Strategies/materials Questioning techniques Academic language Technologies <p>to engage and emphasize key concepts and their value with no irrelevant information</p> <p>Instructions, procedures, and material usage for participating in activities are clear to all or nearly all students</p>	<p>Consistently presents the content:</p> <ul style="list-style-type: none"> Logically, coherently, and in a grammatically correct fashion Supporting the learning of the posted objective(s) Building on content previously mastered Supporting all student populations based on their subject, grade, and level <p>Consistently uses multiple, differentiated:</p> <ul style="list-style-type: none"> Strategies/materials Questioning techniques Academic language Technologies <p>to engage and emphasize key concepts and their value with little to no irrelevant information</p> <p>Instructions, procedures, and material usage for participating in activities are clear to most students</p>	<p>Generally presents content logically and coherent fashion, but:</p> <ul style="list-style-type: none"> Some parts are unclear, grammatically inaccurate, or developmentally inappropriate May not effectively support the learning of the posted objective(s) May not build on content previously mastered May not support all student populations based on their subject, grade, and level <p>Uses limited:</p> <ul style="list-style-type: none"> Verbal and nonverbal techniques to convey concepts and their value Academic language with some irrelevant information <p>Instructions, procedures, and material usage for participating in activities are clear to some students</p>	<p>Presents content and purpose:</p> <ul style="list-style-type: none"> In a confusing way, using unclear, grammatically incorrect, and/or incoherent language With little to no evidence of instruction in support of the posted objective(s) Does not build on content previously mastered Does not support all student populations based on their subject, grade, and level <p>Rarely uses:</p> <ul style="list-style-type: none"> Verbal and nonverbal techniques to convey concepts and their value Academic language with some irrelevant or inaccurate information <p>Instructions, procedures, and material usage for participating in activities are clear to very few students</p>

May 2019 Revision

Figure 1: An example of one of the 18 metrics on which principals are asked to evaluate teachers in the classroom. Principals assign points on each metric, from 0 to 3 according to how well they believe a teacher follows this rubric.

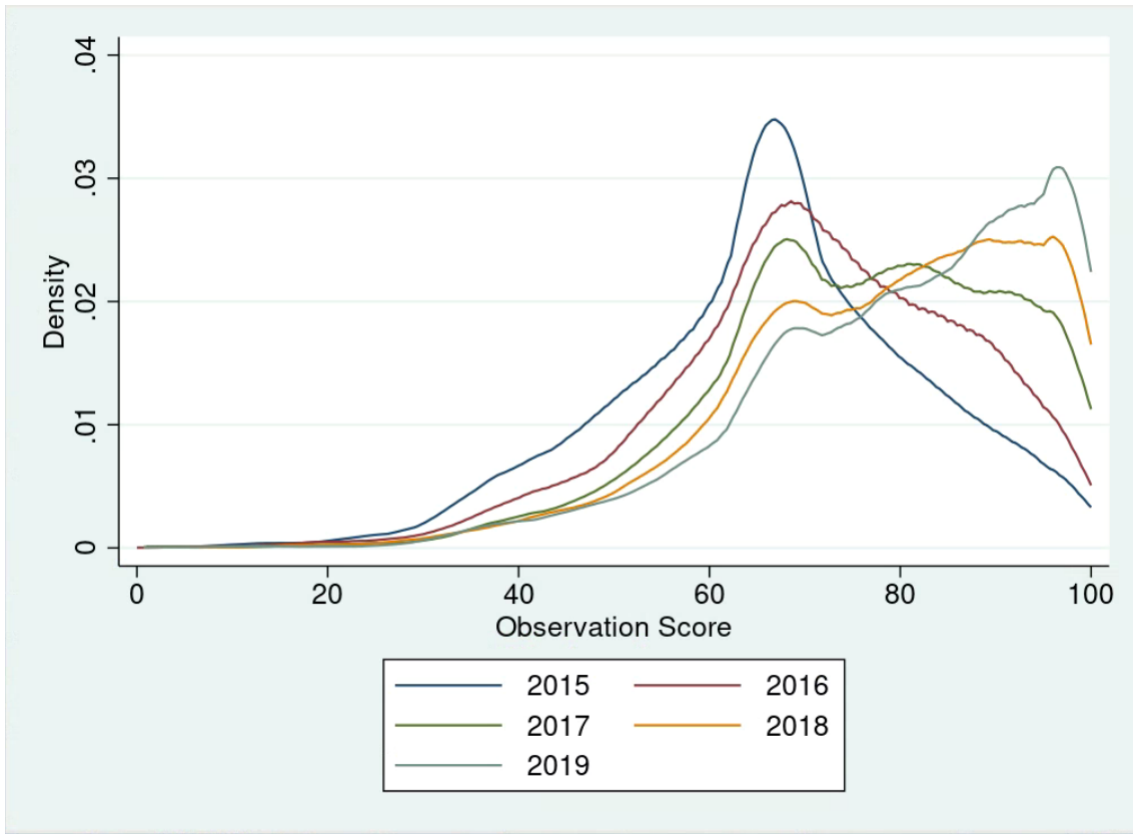


Figure 2: Kernel density distributions of the total observation scores for all teachers in each school year. School years are denoted by the Spring calendar year of each school year.

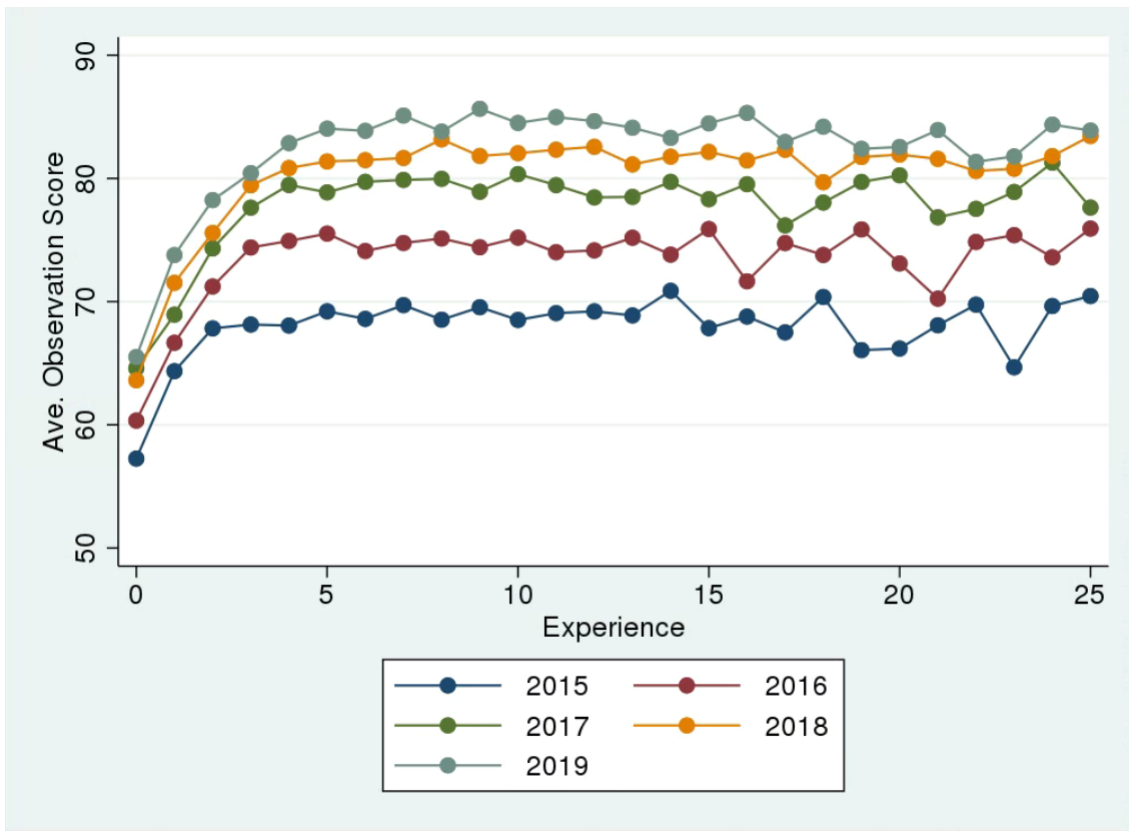


Figure 3: Average observation scores by teacher experience. Each dot represents the average observation score for teachers with that level of experience in each year. Each line represents the pattern of average observation scores for one school year.

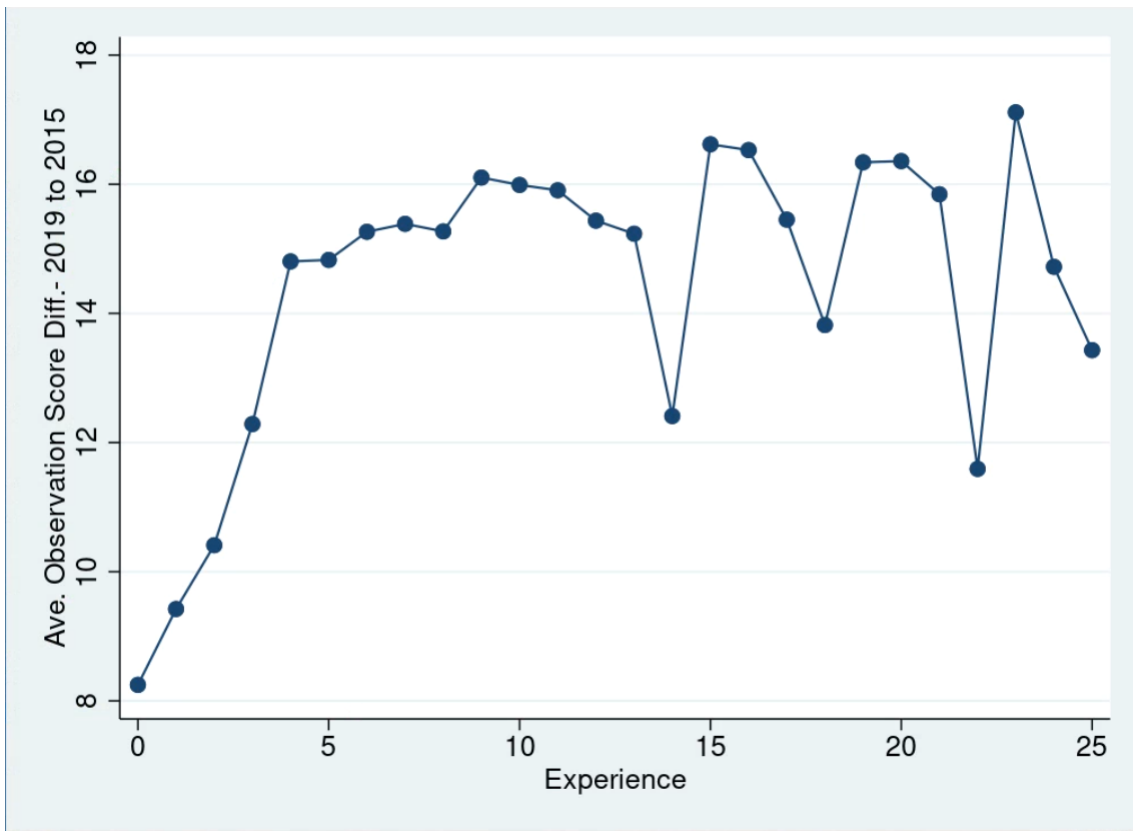


Figure 4: The difference between average observation scores in 2018-2019 and 2014-2015 for teachers in each year with a given level of experience. Each dot represents the difference in the average classroom observation score for teachers with that level of experience between the two years.

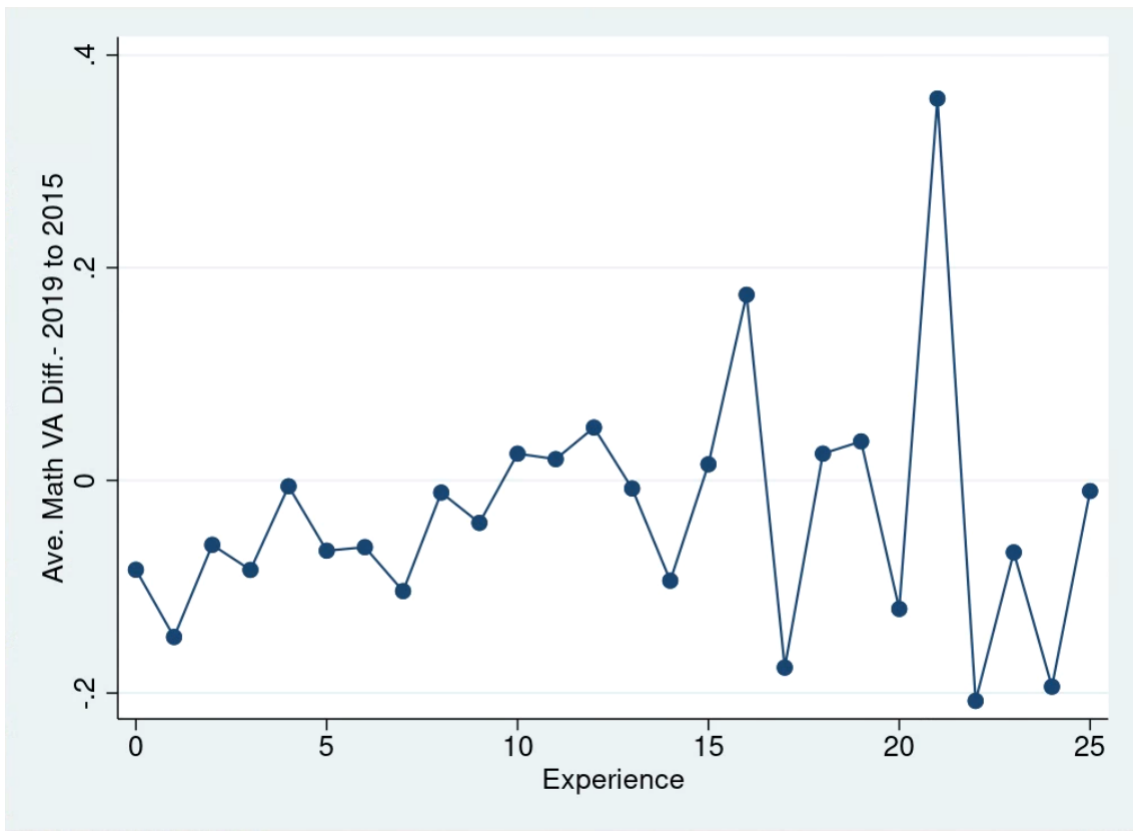


Figure 5: The difference between the average math value-added in 2018-2019 and 2014-2015 for teachers in each year with a given level of experience. Each dot represents the difference in the average math value-added score for teachers with that level of experience between the two years.

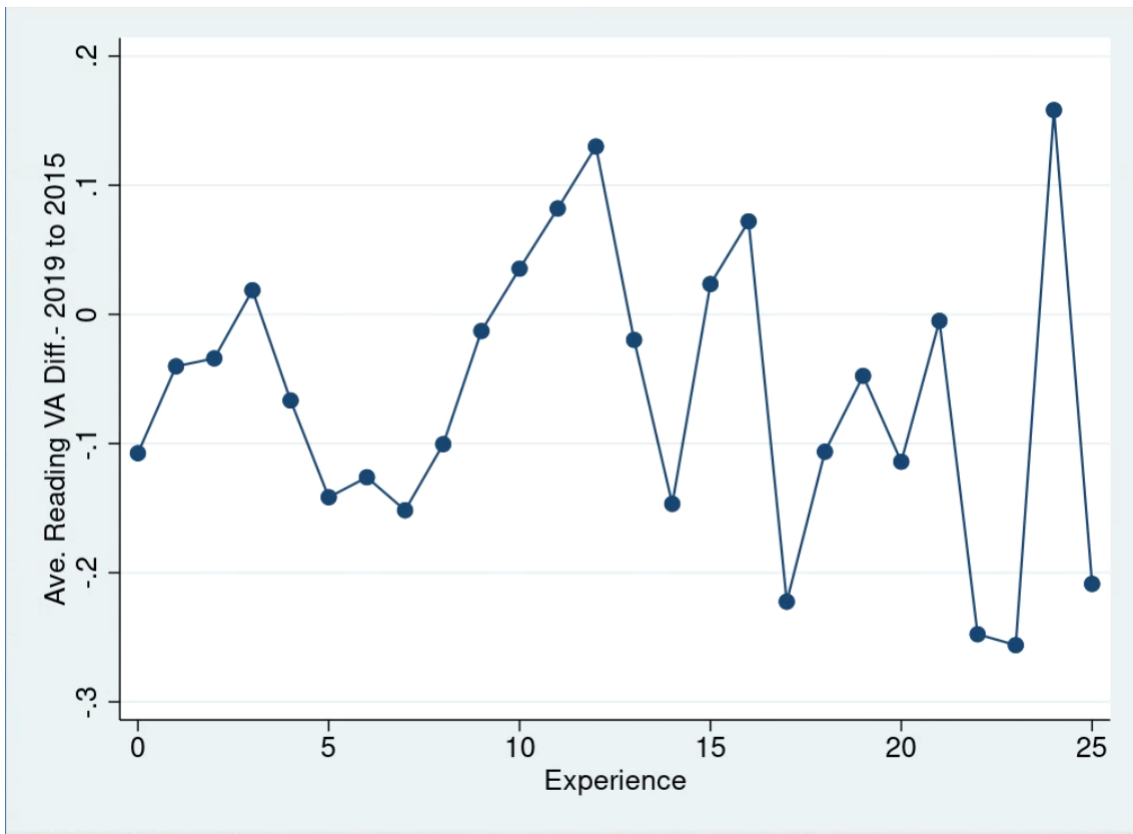
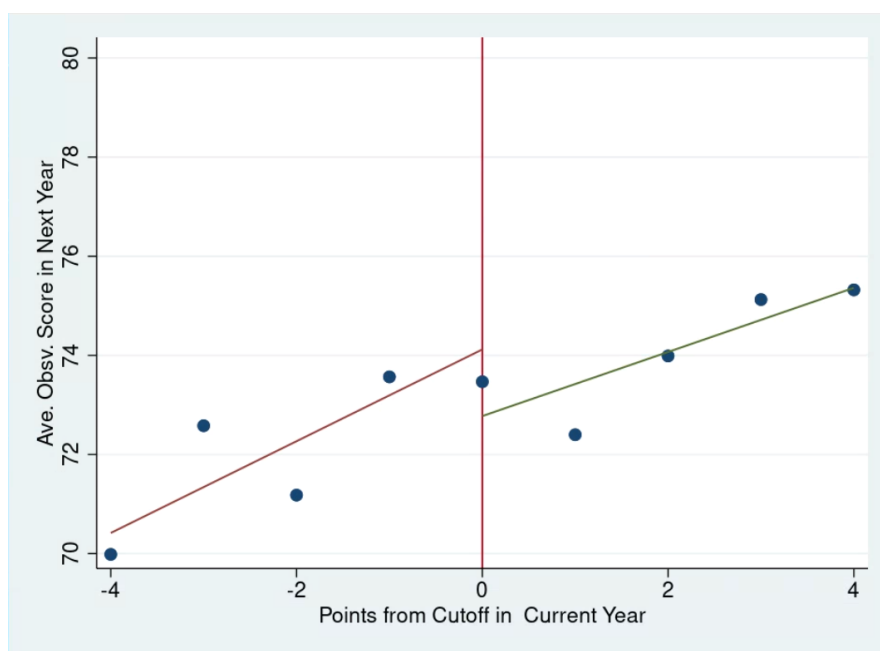
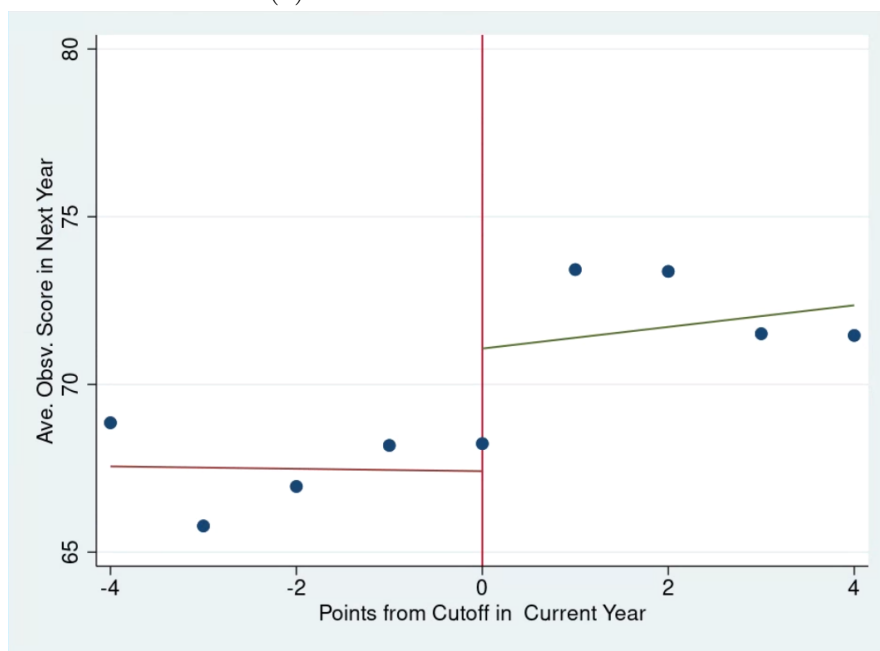


Figure 6: The difference between the average reading value-added in 2018-2019 and 2014-2015 for teachers in each year with a given level of experience. Each dot represents the difference in the average reading value-added score for teachers with that level of experience between the two years.

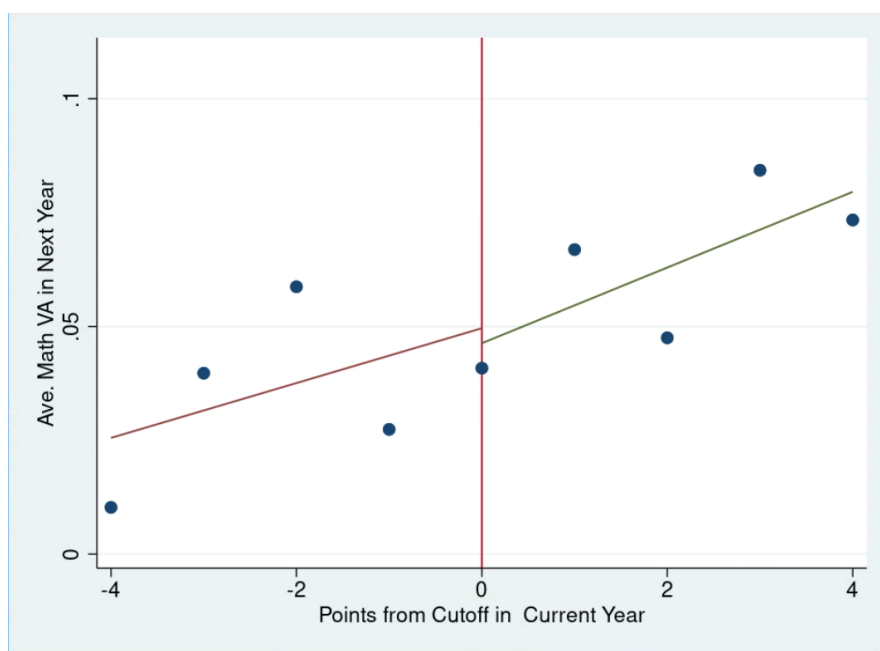


(a) With Financial Incentive

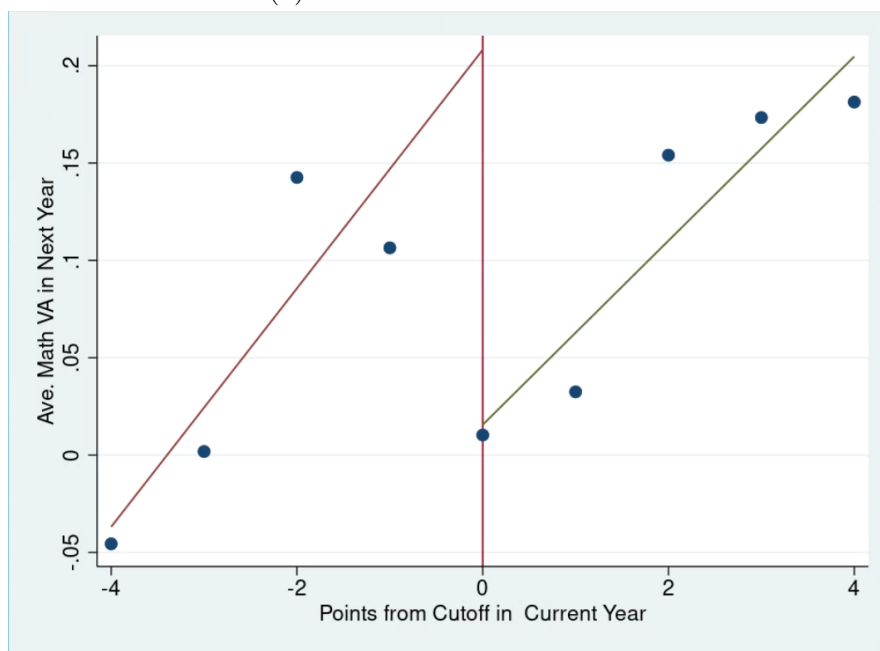


(b) Without Financial Incentive

Figure 7: Discontinuity plots of classroom observations scores in the following year around a compensation cutoff in the year prior for (a) teachers with a financial incentive to achieve the higher compensation level and (b) teachers without such an incentive. Estimates are pooled for every year of my sample period. Each dot represents the average observation score for the teachers at that points-distance from the cutoff in the year prior.

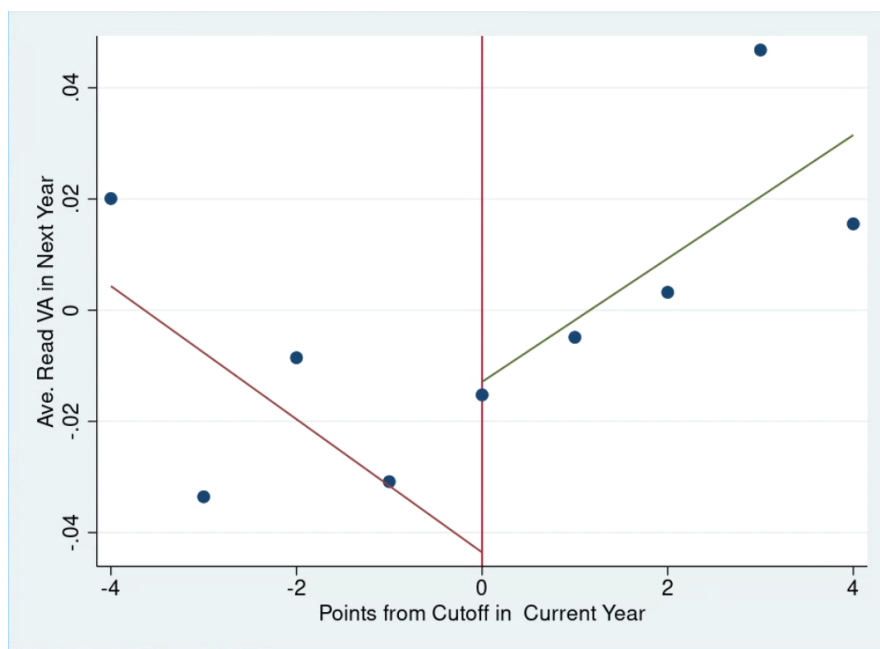


(a) With Financial Incentive

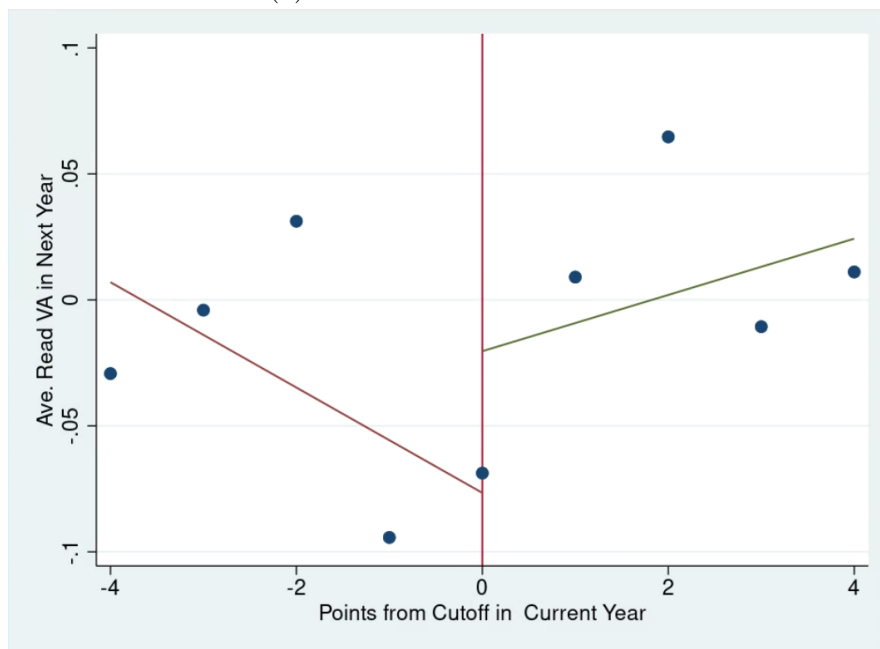


(b) Without Financial Incentive

Figure 8: Discontinuity plots of math value-added in the following year around a compensation cutoff in the year prior for (a) teachers with a financial incentive to achieve the higher compensation level and (b) teachers without such an incentive. Estimates are pooled for every year of my sample period. Each dot represents the average math value-added for the teachers at that points-distance from the cutoff in the year prior.

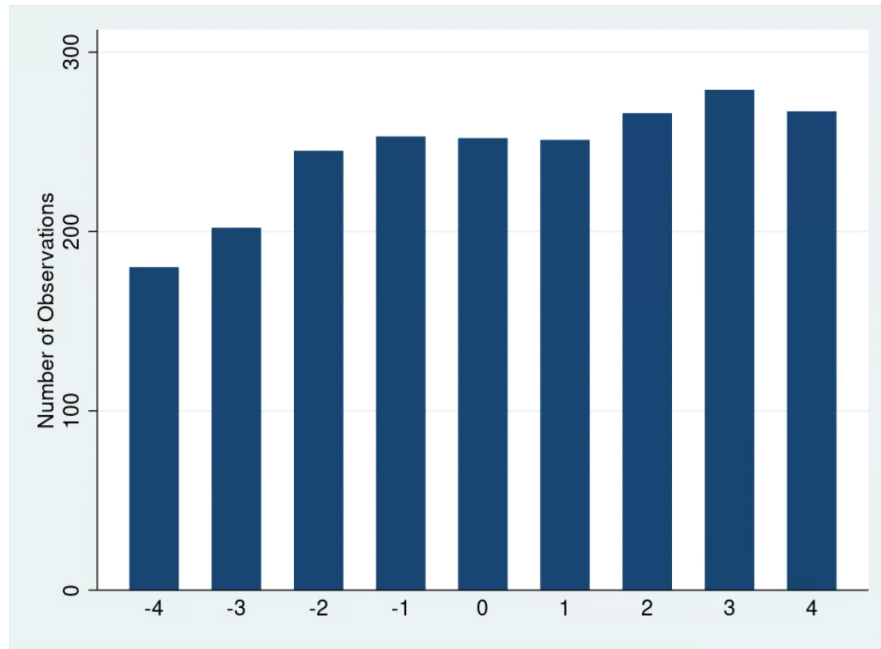


(a) With Financial Incentive

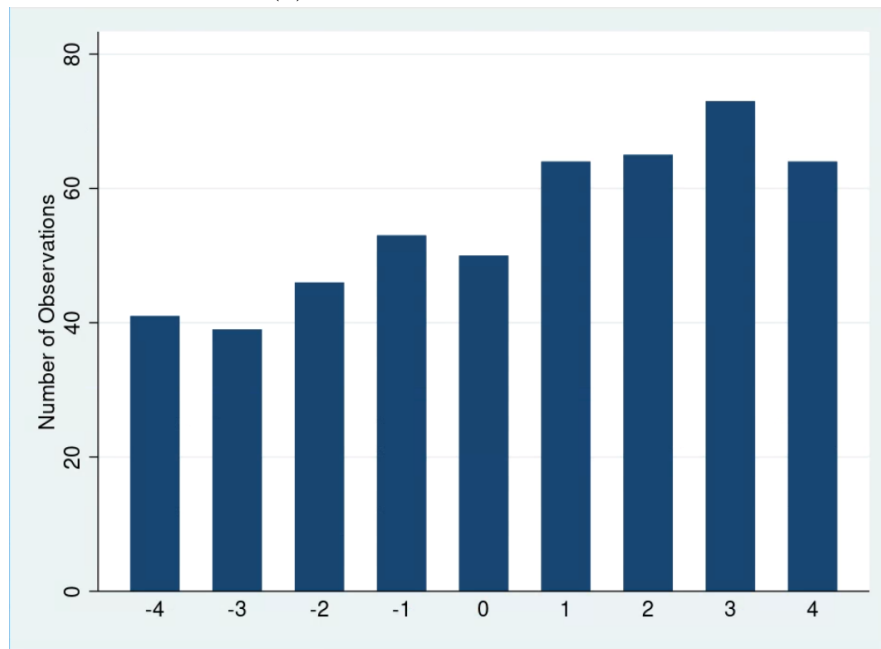


(b) Without Financial Incentive

Figure 9: Discontinuity plots of reading value-added in the following year around a compensation cutoff in the year prior for (a) teachers with a financial incentive to achieve the higher compensation level and (b) teachers without such an incentive. Estimates are pooled for every year of my sample period. Each dot represents the average reading value-added for the teachers at that points-distance from the cutoff in the year prior.



(a) With Financial Incentive



(b) Without Financial Incentive

Figure 10: Plots of the number of teacher observations around a compensation cutoff for (a) teachers with a financial incentive to achieve the higher compensation bin and (b) teachers without such an incentive. Observations are pooled over all years of the sample period. The height of each bar represents the number of teachers at that distance from a cutoff.

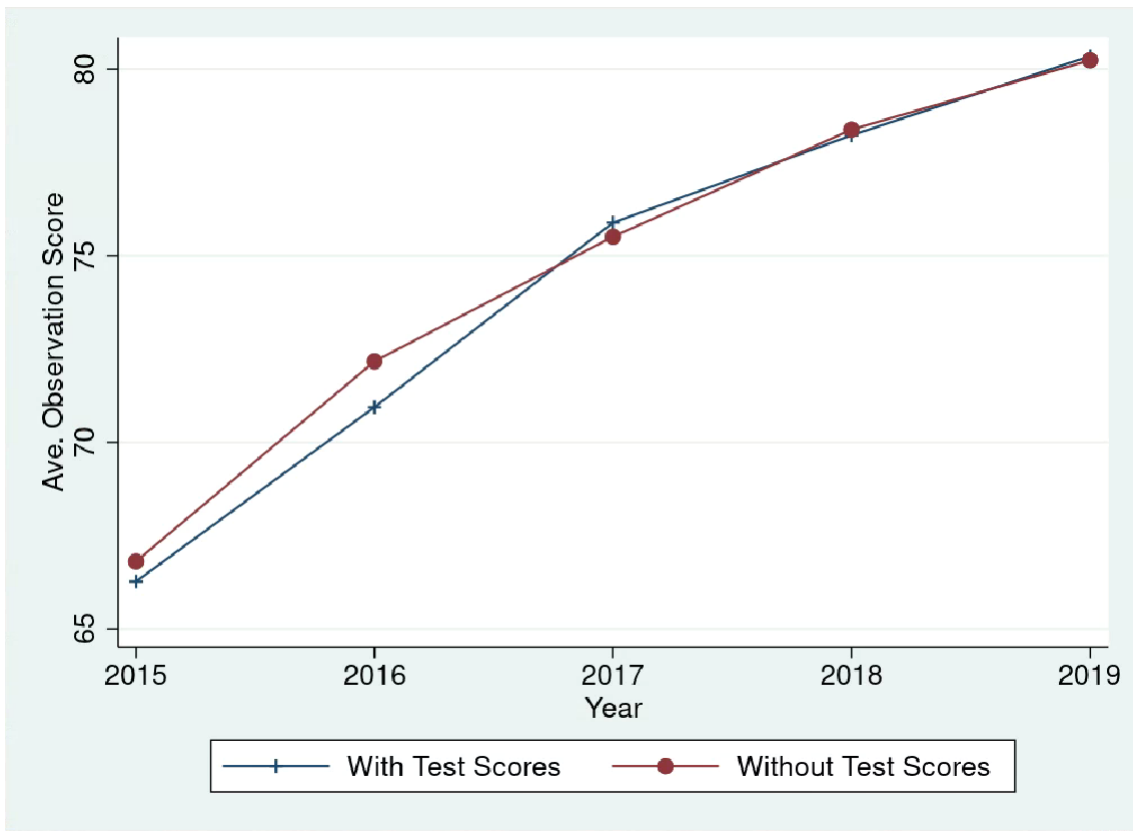


Figure 11: The trends in observation score means for teachers with and without student test scores. The red line represents teachers without student test scores and thus were not included in the principal penalty calculation. The blue line represents the trend for teachers with their own student test scores and were included.

Tables

Table 1: Summary Statistics for Teachers

	Full Sample	Value-Added Sample
Observation Score	73.18 (16.12)	73.47 (16.39)
# of Observations	8.033 (2.591)	8.082 (2.624)
Experience	9.795 (9.768)	8.894 (9.184)
Female	0.713 (0.453)	0.782 (0.413)
Advanced Degree	0.292 (0.455)	0.281 (0.449)
White	0.308 (0.462)	0.273 (0.446)
Black	0.371 (0.483)	0.398 (0.490)
Hispanic	0.277 (0.448)	0.289 (0.453)
Teacher-Year Observations	39142	9611

Summary statistics represented here are for the school years 2014-2015 to 2017-2018. Standard deviations in parentheses. # of Observations represents the average number of classroom observations for a teacher in each year. Experience represents average years of professional experience. Advanced Degree represents the fraction holding a Master's degree or higher.

Table 2: Mean and Variance of Observation Scores by Year

Year	Mean	Std. Dev.	Fraction 66+	Fraction 95+
2015	66.42	15.89	0.56	0.04
2016	71.31	15.58	0.68	0.06
2017	75.75	15.41	0.75	0.10
2018	78.29	15.72	0.79	0.15
2019	80.30	15.68	0.82	0.20

Each row describes classroom observations for one school year. Fraction 66+ signifies the fraction of teachers that received an average score of 66 or higher, representing the fraction of teacher who received an average metric score in the two highest categories. Fraction 95+ represents the fraction of teachers with scores of at least 95 points out of 100.

Table 3: Distribution of Observation Score Growth by Year

Year	Mean	10%	25%	50%	75%	90%	Fract. Neg. Growth
2016	5.94	-9	-1	6	13	21	0.24
2017	5.25	-8	-1	5	12	20	0.28
2018	3.43	-9	-2	3	10	17	0.32
2019	3.20	-9	-2	2	9	17	0.32

Each row presents summary statistics for the change in teachers' observation scores from the prior school year. "Fract. Neg. Growth" represents the fraction of teachers who received any decrease in observation score from the year prior.

Table 4: Correlations Between Classroom Observation Scores and Value-Added

Year	Math	Reading
2015	0.216	0.174
2016	0.270	0.203
2017	0.192	0.132
2018	0.252	0.185
2019	0.300	0.213

Each row represents the Pearson's correlation coefficient between classroom observations and either math or reading value-added in each school year.

Table 5: Predicted Evaluations on Student Achievement

Math Test Score			
	(1)	(2)	(3)
	2017	2018	2019
Pred. Obsv.	0.115*** (0.0160)	0.120*** (0.0183)	0.109*** (0.0153)
Observations	41660	40421	43559
Reading Test Score			
	(1)	(2)	(3)
	2017	2018	2019
Pred. Obsv.	0.0703*** (0.00984)	0.0803*** (0.00944)	0.0758*** (0.00931)
Observations	47050	46340	47919

Standard errors in parentheses.

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Each regression is run at the student level. Standard errors are clustered at the teacher level. Each regression estimates the math or reading test score for students in that year based on the predicted classroom observation score for that student's teacher, based on observations from the two prior years. Each regression also controls for student demographics including student sex, race and ethnicity, free and reduced price lunch status, special education status and limited English proficiency, as well as prior year test score, absences and disciplinary infractions.

Table 6: Difference in Discontinuity Estimates

	(1)	(2)	(3)
	Obsv. Score	Math VA	Reading VA
Above Cutoff W/ Incentive	-5.011** (2.041)	0.0625 (0.100)	0.0137 (0.0832)
Teacher Controls	Y	Y	Y
Observations	2672	1689	2050

Standard errors in parentheses

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Each regression is run at the teacher-year level, using robust standard errors. Difference in discontinuity estimates show the difference in the discontinuities between teachers with and without a financial incentive around a threshold, on the average classroom observation scores, math or reading value-added for teachers who just achieved a higher compensation level in the prior year. Teacher controls include teacher years of experience and whether or not a teacher has a Master's degree or higher.

Table 7: Regression Discontinuity Estimates

Panel A: Teachers With Financial Incentive						
	(1)	(2)	(3)	(4)	(5)	(6)
	Obsv. Score	Obsv. Score	Math VA	Math VA	Reading VA	Reading VA
Above Cutoff	-1.663 (1.220)	-1.425 (1.204)	-0.00193 (0.0523)	0.00820 (0.0519)	0.0352 (0.0394)	0.0347 (0.0393)
Teacher Controls	N	Y	N	Y	N	Y
Observations	2172	2172	1369	1369	1661	1661
Panel B: Teachers Without Financial Incentive						
	(1)	(2)	(3)	(4)	(5)	(6)
	Obsv. Score	Obsv. Score	Math VA	Math VA	Reading VA	Reading VA
Above Cutoff	3.656 (2.663)	3.556 (2.664)	-0.193 (0.142)	-0.174 (0.143)	0.0562 (0.123)	0.0725 (0.118)
Teacher Controls	N	Y	N	Y	N	Y
Observations	495	495	325	325	382	382

Standard errors in parentheses

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Each regression is run at the teacher-year level using robust standard errors. Regression discontinuity estimates show the difference in average classroom observation scores, math or reading value-added for teachers who just achieved a higher compensation level in the prior year. Panel A presents these estimates for teachers who faced a financial incentive around the cutoff, while Panel B presents these estimates for teachers who had grandfathered pay protection and so did not face a financial incentive. Teacher controls include teacher years of experience and whether or not a teacher has a Master's degree or higher.

Table 8: Regression Discontinuity Estimates-Variied Bandwidths

Panel A: Teachers With Financial Incentive						
	Bandwidth 5			Bandwidth 3		
	(1)	(2)	(3)	(4)	(5)	(6)
	Obsv. Score	Math VA	Reading VA	Obsv. Score	Math VA	Reading VA
Above Cutoff	-1.455 (1.075)	-0.0386 (0.0450)	0.0294 (0.0347)	-0.772 (1.403)	0.0219 (0.0594)	0.00992 (0.0436)
Teacher Controls	Y	Y	Y	Y	Y	Y
Observations	2610	1624	1937	1731	1089	1297
Panel B: Teachers Without Financial Incentive						
	Bandwidth 5			Bandwidth 3		
	(1)	(2)	(3)	(4)	(5)	(6)
	Obsv. Score	Math VA	Reading VA	Obsv. Score	Math VA	Reading VA
Above Cutoff	3.967 (2.416)	-0.0464 (0.131)	0.0527 (0.106)	1.354 (3.027)	-0.185 (0.176)	0.120 (0.145)
Teacher Controls	Y	Y	Y	Y	Y	Y
Observations	599	387	444	390	251	302

Standard errors in parentheses

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Each regression is run at the teacher-year level using robust standard errors. Regression discontinuity estimates show the difference in average classroom observation scores, math or reading value-added for teachers who just achieved a higher compensation level in the prior year. Panel A presents these estimates for teachers who faced a financial incentive around the cutoff, while Panel B presents these estimates for teachers who had grandfathered pay protection and so did not face a financial incentive. Teacher controls include teacher years of experience and whether or not a teacher has a Master's degree or higher. Columns 1-3 show estimates for a distance of 5 around a threshold, and Columns 4-6 shows estimates for a distance of 3.

Table 9: Regression Discontinuity Estimates-Falsification Tests

Panel A: Teachers With Financial Incentive		
	(1)	(2)
	Experience	Advanced Degree
Above Cutoff	0.169 (0.442)	-0.0631* (0.0348)
Observations	2172	2172
Panel B: Teachers Without Financial Incentive		
	(1)	(2)
	Experience	Advanced Degree
Above Cutoff	0.396 (1.560)	0.134 (0.100)
Observations	495	495

Standard errors in parentheses

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Each regression is run at the teacher-year level using robust standard errors. Regression discontinuity estimates show the difference in professional years of experience and holding a Master's Degree or higher for teachers who just achieved a higher compensation level in the prior year.

Table 10: Difference in Differences of Observation Score by Penalty Inclusion

	(1)	(2)
	Ave. Obsv. Score	Ave. Obsv. Score
2016	-0.773 (0.479)	-0.738 (0.474)
2017	0.879* (0.467)	0.791* (0.461)
2018	0.296 (0.471)	0.244 (0.465)
2019	0.285 (0.472)	0.269 (0.466)
Teacher Controls	N	Y
Observations	48564	48564

Standard errors in parentheses

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Each regression is run at the teacher-year level using robust standard errors. Difference in differences estimates show the difference in average classroom observation scores between teachers who were and were not included in the principal penalty calculation. Teacher controls in the second column include teacher years of experience and whether or not a teacher has a Master's degree or higher.