

NBER WORKING PAPER SERIES

SCREENING WITH MULTITASKING

Michael Dinerstein  
Isaac M. Opper

Working Paper 30310  
<http://www.nber.org/papers/w30310>

NATIONAL BUREAU OF ECONOMIC RESEARCH  
1050 Massachusetts Avenue  
Cambridge, MA 02138  
July 2022

We thank Ran Abramitzky, Timothy Bresnahan, Raj Chetty, Liran Einav, Caroline Hoxby, Susanna Loeb, and Derek Neal for their early comments on the paper. We also thank conference participants at 2017 APPAM and 2017 AEFPP and seminar participants at Brown University, the New York Federal Reserve, University of California - Irvine, and the University of Chicago Committee on Education for their helpful comments. We also thank Andrew McEachin, Christine Mulhern, and Lisa Abraham for their helpful feedback. We thank Terry Culpepper, Yiren Ding, Elena Istomina, Jora Li, and Jasper Snowden for research assistance. The research reported here was supported by the Institute of Education Sciences, U.S. Department of Education, through Grant R305A190148. The opinions expressed are those of the authors and do not represent views of the Institute or the U.S. Department of Education. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2022 by Michael Dinerstein and Isaac M. Opper. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Screening with Multitasking  
Michael Dinerstein and Isaac M. Opper  
NBER Working Paper No. 30310  
July 2022  
JEL No. D23,I21,J08,J24,J41

### **ABSTRACT**

What happens when employers would like to screen their employees but only observe a subset of output? We specify a model in which heterogeneous employees respond by producing more of the observed output at the expense of the unobserved output. Though this substitution distorts output in the short-term, we derive three sufficient conditions under which the heterogeneous response improves screening efficiency: 1) all employees place similar value on staying in their current role; 2) the employees' utility functions satisfy a variation of the traditional single-crossing condition; 3) employer and worker preferences over output are similar. We then assess these predictions empirically by studying a change to teacher tenure policy in New York City, which increased the role that a single measure – test score value-added – played in tenure decisions. We show that in response to the policy teachers increased test score value-added and decreased output that did not enter the tenure decision. The increase in test score value-added was largest for the teachers with more ability to improve students' untargeted outcomes, increasing their likelihood of getting tenure. We find that the endogenous response to the policy announcement reduced the screening efficiency gap – defined as the reduction of screening efficiency stemming from the partial observability of output – by 28%, effectively shifting some of the cost of partial observability from the post-tenure period to the pre-tenure period.

Michael Dinerstein  
Department of Economics  
University of Chicago  
1126 East 59th Street  
Chicago, IL 60637  
and NBER  
mdinerstein@uchicago.edu

Isaac M. Opper  
iopper@rand.org

## I Introduction

In many settings, the productivity of workers or institutions varies considerably.<sup>1</sup> Employers or policy-makers may therefore realize large gains from screening out low-performers after a probationary period, which has led to policy proposals for increased screening (Kraft et al., 2020). A key complication, though, is that worker or institutional output is multi-dimensional and rarely observed fully. In addition to adding noise to signals of employee type, such limited observability may distort signals by causing probationary workers to substitute effort toward observed output at the expense of unobserved output. During the probationary period, such substitution (or multitasking) leads to inefficient output and reduces the value of high-powered incentives (Holmstrom and Milgrom, 1991). In the screening context, however, there is an additional consideration: if the amount of the response varies across workers, then this substitution also affects which workers are retained – and hence output in the post-probationary period.

We investigate the implications of screening on a single output in a multi-dimensional output environment and, in particular, how the endogenous response to the screening policy affects its efficiency. Theoretically, we derive a set of sufficient conditions under which the endogenous response makes screening more efficient. Empirically, we study the teaching profession, where output is policy relevant and most incentive-based policy is tied to a teacher’s effect on test score value-added even though a recent literature has established that teachers affect a variety of student cognitive and behavioral outcomes (e.g., Jackson, 2018). We show that, relative to a policy in which the screening came as a surprise to the workers, the announcement of the screening policy leads to short-term losses in the untargeted measures due to multitasking but long-term gains due to improved screening. These gains have first-order policy effects: the endogenous response closes 28% of the reduction of screening efficiency stemming from the partial observability of output.

We start in Section II with a theoretical model that derives predictions about workers’ behavioral responses to announced screening policies. A set of heterogeneous workers vary in two dimensions – their returns to effort on two tasks. A worker is employed for up to two periods, and we consider a class of screening policies that make a worker’s probability of reaching the second period a function of first period output on the first task. We show any announced screening induces a multitasking problem, leading workers to spend too much effort in the first period on the task that they are screened on. This response, though, is heterogeneous and depends on worker type. Conditional on workers’ output in

---

<sup>1</sup>As examples, researchers have estimated large dispersion in quality across several types of agents: teachers (Chetty et al., 2014), doctors (Chan et al., 2022), and managers (Bertrand and Schoar, 2003). Estimated quality also varies across institutions like schools (Abdulkadiroğlu et al., 2020) and hospitals (Chandra et al., 2016). Quality may also differ across products, such as insurance plans (Abaluck et al., 2021).

the first dimension in the post-probationary period, we provide sufficient conditions such that workers higher on the second dimension have lower costs and higher preferences for substituting effort toward the first task: 1) all employees place similar value on staying in their current role; 2) the employees' utility functions satisfy a novel variation of the single-crossing condition; 3) employer and worker preferences over output are not too dissimilar. These teachers thus respond more strongly to incentives, which partially substitutes for the employer's inability to observe output on the second task. Hence, the endogenous response makes screening more efficient.

We then introduce our data in Section III, with a focus on how we measure teacher output in multiple dimensions. Here we build off of a recent literature that has explored teachers' effects on variables commonly collected by school districts but that capture different skills than test scores do (Gershenson, 2016; Jackson, 2018; Gilraine and Pope, 2020; Liu and Loeb, 2021). We extend beyond current test scores in two dimensions – timing and type of measure. For timing, we include measures of student outcomes in future years. For additional measures, we focus on student grades and attendance, which may capture non-cognitive or behavioral skills less related to what tests evaluate.<sup>2</sup> Consistent with the literature, we find that while teachers' effects are correlated across different output measures, these correlations are far from 1 and thus entail information not immediately revealed by test score value-added. We further anchor teachers' effects on different outputs to their effects on student graduation rates, which provides us with a one-dimensional measure of policy effectiveness.

In Section IV, we introduce the empirical context by describing teacher tenure policy in NYC and how it changes over time. Prior to the 2009-2010 school year, nearly all teachers received tenure following their third year in NYC, and the decision had no explicit consideration of a teacher's test score value-added. Then in 2009-2010, the district changed the process to incorporate test score value-added, which coincided with a sharp 30 percentage point decrease in tenure rates. We translate these policy details into an empirical strategy to estimate the causal effect of the screening policy's incentives on teacher output.<sup>3</sup> As the policy's implementation is sudden and only affects teachers yet to receive tenure, we have a variety of useful comparison groups that allow us to control for confounders. Already-tenured teachers do not see a change to their incentives and thus they let us control for yearly shocks to students and teachers in NYC. Pre-tenure teachers in years prior to the

---

<sup>2</sup>While the non-test score based measures are readily available to school districts, the future outcomes are not immediately available and thus district screening policy that requires immediate decisions on whether to give a teacher tenure does not necessarily have such measures at its disposal. As researchers analyzing data from prior years, we do not face that constraint and can thus characterize the policy's effect on multiple dimensions.

<sup>3</sup>Loeb et al. (2015) study how this tenure policy change affected teacher quality in the district by replacing marginal teachers at the end of their probationary period with new teachers.

policy change let us control for experience effects. And teachers for whom the policy change hit in the middle of their pre-tenure periods let us control for cohort or teacher effects that might reflect changes in entering cohorts' quality over time. We thus isolate causal effects by comparing how a cohort or teacher's outcomes change upon the introduction of the new policy, controlling for year and experience effects.

We use this empirical strategy in Section V to test how the policy change affected pre-tenure output. We focus on both targeted (test score value-added) output measures and untargeted output measures. We find that on the targeted measure the introduction of the screening policy increases output by  $0.015 - 0.033$  student standard deviations, depending on the estimator. This effect represents a meaningful fraction of the cross-sectional variation in teacher output in the unincentivized period. Based on the cross-sectional relationship between test score value-added and the other measures, we might expect the other measures to increase as well since the measures are positively correlated. Instead, we estimate that teacher output on the untargeted measures falls in response to the tenure policy. This negative effect on other outcomes is consistent with teachers substituting effort away from building student skills that persist (or are tested on future assessments) or non-tested skills and toward effort improving students' current test scores. Such substitution has immediate consequences. While the policy is focused on screening out low-performing teachers, in the pre-tenure period it has a positive effect on output in the targeted measure but negative effect on the untargeted measures.

Such changes could be temporary or persistent, depending on whether the change in effort alters teachers' future production functions. We find that a teacher's current test score value-added reverts to unincentivized levels (excluding the experience gradient) once she receives tenure and no longer faces the incentive. Thus, the policy had large, temporary incentive effects; the policy's total effect, however, also depends on how teachers are screened into and out of the district in the post-probationary period. How the policy affects screening in turn depends not only on the average incentive effect, but on how different teachers respond differently. Our model predicts that teachers with higher types (e.g., ability) in the untargeted task may substitute more toward the targeted task in the presence of incentives. In Section VI, we use forecasts from a multi-dimensional value-added model to estimate a teacher's type. We then test whether a teacher's response varies with her type and find that teachers with higher (unincentivized) output on the untargeted tasks respond more strongly to the implied incentive on the targeted task. Therefore, the behavioral response to the incentive selects positively on the untargeted dimension.

We put these pieces together in Section VII and estimate the behavioral responses' impact on the tenure policy's screening efficiency. Given that the policy is a threshold one (promote the top 67% based on test score value-added), the behavioral responses only

matter for marginal teachers; it is therefore unsurprising that we find that tenure outcomes change for just 4% of teachers. The teachers who only receive tenure due to the behavioral response, though, have considerably higher output on the untargeted dimension compared to the teachers who lose tenure (239% of the cross-sectional standard deviation). In contrast, these teachers have lower output on the targeted dimension, although these differences are minimal. The quality of the tenured teacher pool thus increases by 7.7% of a standard deviation on the untargeted dimension, while decreasing by 0.76% of a standard deviation on the targeted dimension. Together, these results suggest that the endogenous response to the policy reduced the screening efficiency gap – defined as the reduction of screening efficiency stemming from the partial observability of output – by 28%

In sum, the behavioral response changes the types of teachers screened in to tenure, while simultaneously distorting output in the pre-tenure period. In other words, the behavioral response shifted some the cost of partial observability of output from the post-tenure period to the pre-tenure period and therefore induces a trade-off between short-run and long-run output (Nichols and Zeckhauser, 1982). Because the changes in the two dimensions differ in the short-run versus long-run – targeted output is higher in the short-run and lower in the long-run while untargeted output is the reverse – whether the district benefits from the behavioral response depends on the length of the post-tenure period. We find that the district benefits from the behavioral response if teachers remain in the district between 3.5 and 45 years post-tenure.

## I.A **Related Literature**

This paper brings together the literatures on screening and multitasking. From the screening perspective, our focus is on a monopsonist employer’s on-the-job screening.<sup>4</sup> In many industries like education, this takes the form of up-or-out policies where the employee partially reveals her type during a probationary period and then the employer decides whether to keep the employee (often under favorable contract terms, like tenure).<sup>5</sup> The literature has focused on rationales that make these policies optimal: incentives to invest in firm-specific human capital when there is bilateral moral hazard (Kahn and Huberman, 1988; Prendergast, 1993); incentives to invest in general human capital (Waldman, 1990); maintenance of a stationary distribution of tenure within a firm (Rebitzer and Taylor, 2007); improving screening efficiency (O’Flaherty and Siow, 1992; Demougin and Siow, 1994; Chen and Lee,

---

<sup>4</sup>The broader screening literature considers policies where employees select into or out of a job (e.g., Bénabou and Tirole, 2016; Brown and Andrabi, 2021) or policies where the employer screens job applicants (e.g., Spence, 1973).

<sup>5</sup>In a more general set of contracts, Holmström (1999) argues that career concerns may lead to managerial moral hazard.

2009; Barlevy and Neal, 2019);<sup>6</sup> and (most related to our paper) reducing multitasking in the post-probationary period (Kou and Zhou, 2009). In contrast, we take the screening structure as given, but unlike these papers in which the employee’s type is one-dimensional, we consider a setting in which the employee’s type is multi-dimensional (Armstrong and Rochet, 1999).

Multidimensional screening is itself a vast literature, with Spence (1973) and Nichols and Zeckhauser (1982) offering models of multi-dimensional agents where a single-crossing condition means that the agent’s signal manipulation reveals her type and thus improves screening efficiency. A recent literature has extended the framework to include unobservable heterogeneity in manipulation costs. This addition lowers screening efficiency, especially at high stakes (Frankel and Kartik, 2019); in part motivated by this, Björkegren et al. (2020) develops an estimator of an optimal screening function that is robust to manipulation.<sup>7</sup> In contrast, we focus on the case where the manipulation cost stems from agents substituting from producing unobservable output to observable output. Unlike other forms of manipulation, this one affects the principal both directly, through distorted output in the probationary period, and indirectly, through a distorted signal.

Our focus on this form on manipulation costs means that we also speak to the larger literature on multitasking. The theoretical literature characterizes the optimal contracts when the principal cannot contract on all forms of output (Holmstrom and Milgrom, 1991; Baker, 1992). We take the contract as given and study its effects on output. The empirical literature usually studies the effect of targeted incentives on untargeted outcomes (e.g., Fryer and Holden, 2012; Hong et al., 2018), and in the education context this is often described as “teaching to the test.” These papers’ focus is on the effect of multitasking on output of a fixed set of workers. We also allow multitasking to affect screening efficiency by sending distorted signals, which means that we are particularly interested (both theoretically and empirically) in how the degree of multitasking differs across individuals in the population.

Our empirical context also places us in the economics of education literature. Efforts to improve teacher quality have generally split into two broad paths.<sup>8</sup> The first treats a teacher’s output as (mostly) fixed and quantifies the achievement gains from changing the composition of the teacher workforce through policies like removing low-performers (Hanushek, 2009; Chetty et al., 2014; Loeb et al., 2015).<sup>9</sup> The second treats a teacher’s

---

<sup>6</sup>Barlevy and Neal (2019)’s model rationalizes the combination of up-or-out policies and heavy workloads during the probationary period.

<sup>7</sup>A large empirical literature has studied the efficiency of different screening devices, such as hassle costs (e.g., Deshpande and Li, 2019).

<sup>8</sup>Both of these paths take the pool of potential teachers as fixed. Another strand of literature considers policies that affect the teaching pool more broadly (e.g., Nagler et al., 2015; Sutchter et al., 2016).

<sup>9</sup>There is also a vast literature on how to measure teacher effectiveness (e.g., Jacob and Lefgren, 2008; Staiger and Rockoff, 2010; Kane et al., 2013).

output as malleable and quantifies the achievement gains from incentive changes through policies like performance pay (e.g., Glewwe et al., 2010; Muralidharan and Sundararaman, 2011; Fryer, 2013; Dee and Wyckoff, 2015).<sup>10</sup> We argue that if teacher output is malleable, then screening and incentive policies are not separable, as high-stakes screening on output changes a teacher’s incentives. In other words, high-stakes screening necessarily alters teachers’ incentives and our empirical results speak to how important it is to account for teachers’ responses to these incentives when considering a teacher screening policy.

Performance pay for teachers has long been controversial due to the concern that it would cause teachers to “teach to the test.” Much of the empirical work on the topic looks at whether ranking teachers (or schools) based on unincentivized measures is similar to ranking based on incentivized measures (Papay, 2011; Corcoran et al., 2013; Cohodes, 2016), while other work compares the effects of performance pay on targeted and untargeted tests (Muralidharan and Sundararaman, 2011). More recently, an emerging literature has shown that teachers affect a range of student outcomes and that teacher effectiveness across these dimensions is only weakly correlated (Gershenson, 2016; Jackson, 2018; Kraft, 2019; Petek and Pope, 2016; Gilraine and Pope, 2020; Liu and Loeb, 2021). These correlations, though, are specific to the incentive environment (Neal, 2011). Thus, a change in incentives from a screening announcement might change the ordering of teachers in a given dimension and therefore which teachers receive tenure. Our empirical results show that teachers substitute across tasks that target different student outcomes in response to the change in tenure policy, which means the endogenous response to the policy along with the multi-dimensionality of teacher output create short-term losses (from multitasking) but long-term gains (from improved screening).

## II Model

We develop a two-period model of heterogeneous agents (teachers) exerting effort on two tasks. A principal (school district) decides which agents to retain for the second period. We present results in this section and provide proofs in Appendix A.

### II.A Model Intuition

Before diving into the details, we present the following simplified version of the model to build intuition. Consider a set of teachers who differ on their returns to effort on two tasks and suppose that each teacher’s utility is simply the sum of the two outputs. Now suppose that each day is split into three periods and all teachers are required to spend one period on

---

<sup>10</sup>Macartney et al. (2018) compares these types of policies directly in the same empirical context and estimates that performance pay has a higher return.



the first task and one period on the second task, but get to choose how to spend the third period; there is no cost of effort beyond the constraint that there are only three periods. If teachers allocate effort according to their comparative advantage, without the additional incentive of screening teachers who are better at task one than task two will choose to spend the third period working on task one and teachers who are better at task two than task one will choose to spend the third period working on task two. Crucially, this means that those individuals who are particularly good at task two spend more effort on task two – and less effort on task one – than their peers who are equally good at task one and less effective on task two. This allocation decision implies that, conditional on task one effectiveness, workers who are good on task two produce less of the first output and without any endogenous response to the policy would be less likely to be retained than individuals who are worse at task one.

Now consider what happens when all employees are told in advance that they will be evaluated on the first output and that all workers place a large value on staying in the profession is large. Although all employees will then have an incentive to shift their focus toward producing more of the first output, it is precisely those who are particularly effective on the second task who will be able to substitute, since those are the employees that are initially spending the third period working on the second task rather than the first. Thus, conditional on teachers' task one ability the size of response will be increasing in teachers' task two ability, making the observed output a better measure of the underlying worker quality.

This example made a number of simplifying assumptions, including that the size of the incentive was large enough that the teachers ignored any short-term costs of substituting effort in the initial period. In what follows, we develop a set of sufficient conditions which guarantee that the endogenous response improves screening efficiency, even when teachers account for the short-term cost of substituting effort in the initial period.

## II.B Model Set-Up

**Production of Student Outcomes and Teachers' Utility:** We assume each teacher has some fixed type  $\theta \in \Theta = [\theta_1^{min}, \theta_1^{max}] \times [\theta_2^{min}, \theta_2^{max}] \subset \mathbb{R}_+^2$  and in each period chooses effort level  $e \in \mathcal{E} \subset \mathbb{R}_+^2$ .<sup>11</sup> We assume that  $\theta_i^{min} > 0$  for  $i \in \{1, 2\}$  and that  $\mathcal{E}$  is a compact and convex set. The type and effort level combine to affect student outcomes  $x$  in the following way: if the  $k^{th}$  dimension of the individual's type is  $\theta_k$  and she exerts effort level

---

<sup>11</sup>Note that by holding  $\Theta$  fixed, we are assuming that the screening policy does not change the set of teachers in the population. If there is competition for employees' labor supply, the best employees may have the ability to switch employers rather than distort their effort (Bénabou and Tirole, 2016; Brown and Andrabi, 2021; Tincani, 2021).

$e_k$  on task  $k$ , then she improves student outcome  $k$  by  $e_k\theta_k$ , i.e.,  $x_k = e_k\theta_k$ . We will refer to the first outcome as the teacher’s “test score value-added.”

We assume that teachers get value from increasing students’ outcomes, but incur a cost of effort. This cost of effort, along with the constraint set, limits the amount of  $x$  that each teacher produces. Specifically, we assume teachers’ per-period utility is given by a function  $b(x) - c(e)$  for some concave benefit function  $b$  that is increasing in  $x$  and a strictly convex cost function  $c$  that is increasing in  $e$ . We assume that both  $b$  and  $c$  are twice differentiable, with  $\frac{\partial b^2}{\partial x_1 \partial x_2} \leq 0$  and  $\frac{\partial c^2}{\partial e_1 \partial e_2} > 0$ . These assumptions on the cross-derivatives, along with the assumption that  $\mathcal{E}$  is convex, capture the idea that teachers face a tradeoff when deciding how to allocate their effort. For technical reasons, we will also assume that  $b(x)$  is bounded – at least for  $\theta \in \Theta$  and  $e \in \mathcal{E}$  – and that  $c(e) < \infty$  for all  $e$  in the interior of  $\mathcal{E}$ .

Because we are not attempting to separate teachers’ benefits versus costs, we will instead generally write each teacher’s per-period utility as a function of  $e$  and  $\theta$ , i.e.  $u(e, \theta) = b(\theta_1 e_1, \theta_2 e_2) - c(e_1, e_2)$ . The assumptions on  $b$  and  $c$  imply the following assumptions on  $u$ : a)  $\frac{\partial^2 u}{\partial e_i \partial \theta_i} > 0$ ; b)  $\frac{\partial^2 u}{\partial e_1 \partial e_2} < 0$ , and c)  $\frac{\partial^2 u}{\partial e_i \partial \theta_j} \leq 0$  for  $i \neq j$ .<sup>12</sup> At times, we will write utility as a function of  $(x, \theta)$  instead of  $(e, \theta)$ .

We denote a teacher’s indirect utility – when she optimizes without dynamic considerations – as  $v(\theta) = \max_{e \in \mathcal{E}} u(e, \theta)$ . We denote the optimal choice of effort – again when she optimizes without dynamic considerations – as  $e^*(\theta) = \arg \max_{e \in \mathcal{E}} u(e, \theta)$ , with effort on task  $k$  denoted as  $e_k^*(\theta)$ . The effect on student outcome  $k$  under these optimal choices is then denoted as  $x_k^*(\theta)$ , which is equal to  $e_k^*(\theta) \cdot \theta_k$ .

An important implication of the assumptions on the teacher’s utility and choice set is that teachers choose to specialize on the task where they have comparative advantage, spending more effort on that task and less on the other task, rather than seeking similar levels of both outputs. Formally, this means that  $e_1^*(\theta)$  is increasing in  $\theta_1$  and decreasing in  $\theta_2$ , which implies that  $x_1^*(\theta)$  is likewise increasing in  $\theta_1$  and decreasing in  $\theta_2$ .

**Defining the Screening Policy:** We will assume that time can be split into discrete periods (e.g., school years) and that the pool of teachers is fixed within a period. We allow teachers to stay in the teaching profession for at most two periods, with the first period being the pre-tenure period and the second period being the post-tenure period. Whether an individual teacher is granted tenure, i.e., allowed to stay in the profession for the second period, is determined by a function  $p : X_1 \rightarrow [0, 1]$ . That is, the probability that a teacher with test score value-added in the first period of  $x_1$  is granted tenure is  $p(x_1)$ . For simplicity, we take  $p(x_1)$  to be exogenous and only assume that it is weakly increasing in  $x_1$ . Note that although a teacher affects multiple outcomes, whether she gets tenure only depends on

<sup>12</sup>The first comes from the fact that  $b$  is increasing in  $x_i$ ; the second comes directly from the assumptions on the cross derivatives of  $b$  and  $c$ ; the third comes from the assumption on the cross derivatives of  $b$ .

her output in the first dimension.

**Dynamic Optimization:** We assume the teachers place a weight equal to one on the first period and a weight of  $\lambda > 0$  on the second period.<sup>13</sup> We assume all individuals with the same type have the same (expected) outside option, which we denote as  $\tilde{v}(\theta)$ . We define  $\Delta v(\theta) = v(\theta) - \tilde{v}(\theta)$  and assume that  $\Delta v(\theta) > 0$  for all  $\theta$ , i.e., that all teachers want to stay in the profession if given the option. We also assume that  $\Delta v(\theta)$  is a continuously differentiable function of  $\theta$ .

It then follows that teachers choose effort in the first period to maximize:  $u(e, \theta) + \lambda \cdot p(e_1 \theta_1) \cdot \Delta v(\theta)$ . Even though the policy may solely seek to screen workers, it induces an incentive effect because whether the teacher stays depends on her test score value-added. We denote the first period optimal effort choices and effects on student outcomes for an individual of type  $\theta$  under some screening policy  $p$  as:

$$e^*(\theta|p) = \arg \max_{e \in \mathcal{E}} u(e, \theta) + \lambda \cdot p(e_1 \theta_1) \cdot \Delta v(\theta) \quad (1)$$

$$x_k^*(\theta|p) = e_k^*(\theta|p) \cdot \theta_k. \quad (2)$$

## II.C Incentive Effects of the Screening Policy

As is clear from Equation (1), the announcement of the screening policy changes pre-tenured teachers' incentives and therefore likely changes their output. In this subsection, we explore the impact of this change in incentives.

### II.C.1 How Does the Incentive Change Teachers' Output?

We start with the direction of the effect on teachers' output. Because the screening policy adds weight to the test score value-added, teachers exert more effort on the first task, which increases teachers' test score value-added in the pre-tenure period. If that makes effort on the other task more costly – or output on the other dimension less valued – then this extra effort spent increasing test score value-added leads to a decrease in other output. This is akin to the multitasking model of Holmstrom and Milgrom (1991) and reflects the concern that adding incentives to test score value-added measures would lead to more “teaching to the test.” Formally, we have the following theorem:

---

<sup>13</sup>This model is equivalent to a model with  $K$  pre-tenure periods and an infinite number of post-tenure periods in which the  $N^{th}$  period is discounted by a rate of  $\beta^N \in (0, 1)$ ; in this set-up  $\lambda = \frac{\beta^K}{1 - \beta^K}$ . Thus,  $\lambda$  may be less than one, due to teachers discounting the future, or greater than one, since the second period can be thought of as collapsing the post-tenure years of teaching into a single period.

**Theorem 1.** *For any weakly increasing screening function, we get that for every  $\theta$ :*

$$\begin{aligned} x_1^*(\theta|p) &\geq x_1^*(\theta) \\ x_2^*(\theta|p) &\leq x_2^*(\theta). \end{aligned}$$

### II.C.2 How Does the Response Vary across Workers?

Teachers increase their test score value-added in response to the incentive induced by the screening policy. But in the screening context, how the size of the response varies across teachers is particularly important as it will change who is screened out by the policy. In this section, we explore this issue.

We focus on how the responses differ among teachers who – absent the additional incentive – would have the same test score value-added, but different output on the other dimensions. That is, we focus on two teachers  $\theta$  and  $\theta'$  with  $x_1^*(\theta) = x_1^*(\theta') \equiv x_1^*$  and  $x_2^*(\theta) < x_2^*(\theta')$ .<sup>14</sup> We focus on this particular comparison for two reasons. First, it means that we can explore teachers' differing responses without specifying the particular shape of the screening function  $p(x_1)$ . Second, this comparison allows us to provide sufficient conditions on how the different responses to the incentive induced by the screening affect the efficiency of the screening itself.

**Different Size of Incentive:** We restricted our attention to teachers  $\theta$  and  $\theta'$  where any deviation away from their optimal output (absent the screening policy) has the same impact on the probability that they can stay for the post-tenure period – i.e.,  $p(x_1^*(\theta) + \Delta) = p(x_1^*(\theta') + \Delta)$  for any  $\Delta$  and any screening function  $p$ . As is clear from Equation (1), however, the size of the incentive depends not only on how changing  $x_1$  affects the likelihood that the teacher can stay into the post-tenure period, but also on how much she values staying in the profession relative to the outside option. Thus, part of what determines how different teachers respond is how the value of remaining in the teaching profession – relative to the outside option – depends on the underlying ability of the teacher.

**Different Responses to the Same Incentive:** The two teachers may also differ in the costs of deviating, which would cause them to respond differently even to the same incentive. To understand this mechanism, we consider how different individuals who initially produce the same  $x_1$  respond differently when presented with the same additional incentive to increase  $x_1$ . We formalize this identical additional incentive by adding a weakly increasing function of  $x_1$  –  $\lambda f(x_1)$  – to their original utility function. This can roughly be thought of as the screening function in our setting, but unlike in Equation (1) this function does not depend on  $\theta$ , once we condition on  $x_1$ . For technical reasons, we will assume that  $f$  is

<sup>14</sup>As we show in Appendix A, we can conclude that  $\theta'_k > \theta_k$  for  $k \in \{1, 2\}$ .

right-continuous, an assumption we will leave implicit in the lemma below and those in the Appendix.<sup>15</sup>

In our analysis it will help to define  $\tilde{u}(x_1, \theta)$  as the optimal utility individual  $\theta$  can get when constrained to produce  $x_1$ , i.e.,

$$\tilde{u}(x_1, \theta) \equiv \begin{cases} \max_{e \in \mathcal{E}} u(e, \theta) & \text{if } \exists e_1 \in \mathcal{E} \text{ s.t. } e_1 \theta_1 \geq x_1 \\ -\infty & \text{if } \forall e_1 \in \mathcal{E} \quad e_1 \theta_1 < x_1 \end{cases} \quad (3)$$

For ease of exposition, we will assume here that  $\tilde{u}(x_1, \theta)$  is differentiable in the feasible range and use  $\tilde{u}'(x_1, \theta)$  to denote  $\frac{\partial \tilde{u}(x_1, \theta)}{\partial x_1}$ .

We then have the following lemma:

**Lemma 1.** *Consider two individuals  $\theta$  and  $\theta'$  with  $x_1^*(\theta) = x_1^*(\theta')$  and  $x_2^*(\theta) < x_2^*(\theta')$  and define:*

$$e^*(\theta|f) = \arg \max_{e \in \mathcal{E}} u(e, \theta) + \lambda f(x_1) \quad (4)$$

$$x_k^*(\theta|f) = e_k^*(\theta|f) \cdot \theta_k \quad (5)$$

Then  $x_1^*(\theta|f) \leq x_1^*(\theta'|f)$  for any weakly increasing function  $f(x_1)$  if either of the following are true:

- $\lambda$  is sufficiently large;
- $\tilde{u}'(x_1, \theta') - \tilde{u}'(x_1, \theta)$  is increasing in  $x_1$

While we leave the proof to Appendix A, we discuss the intuition here. Suppose that individual  $\theta$  moves from producing  $x_1^*(\theta)$  to producing  $x_1^*(\theta|f)$  in response to the incentive and consider the cost of individual  $\theta'$  increasing  $x_1$  by the same amount. As we show in Appendix A, from the assumption that  $x_1^*(\theta) = x_1^*(\theta')$  and  $x_2^*(\theta) < x_2^*(\theta')$  we can conclude that  $\theta_1 < \theta'_1$ . Thus, it is feasible for individual  $\theta'$  to also produce  $x_1^*(\theta|f)$  and individual  $\theta'$  has to increase their effect on the first task by less than individual  $\theta$  to achieve the same increase in output. To increase effort on the first task, however, they will have to decrease effort on the second task, either because of the constraint that  $e \in \mathcal{E}$  or to decrease the marginal cost of effort on the first task. This is more costly to individual  $\theta'$  than individual  $\theta$  since  $\theta_2 < \theta'_2$ , so without further assumptions it remains ambiguous which individual will

<sup>15</sup>Otherwise functions such as  $f(x_1) = \mathbf{1}(x_1 > \bar{x}_1)$  would leave some individuals wanting to produce the minimum amount of  $x_1$  that is *strictly* above  $\bar{x}_1$ , which leads to complications. By assuming that the function is right-continuous, we instead force the function to be  $\mathbf{1}(x_1 \geq \bar{x}_1)$ , so the individuals can settle on  $\bar{x}_1$ .

respond more.<sup>16</sup>

For large enough incentive, the reduction in  $e_2$  becomes unimportant relative to the maximizing the incentive and all that matters is how costly it is to produce  $x_1$ . Since  $\theta_1 < \theta'_1$ , it is cheaper for individual  $\theta'$  to produce  $x_1$  than individual  $\theta$  and so she will produce at least as much  $x_1$  with the incentive. One consequence is that if individuals have a large enough incentive to stay in the teaching profession, screening on  $x_1$  is roughly equivalent to screening on  $\theta_1$ . While subtle, this makes the screening more efficient since a policy that surprises employees by screening on one dimension of output is less efficient than one that screens on one dimension of ability; see Appendix B for more discussion. For smaller incentives, however, we cannot simply ignore the differential cost of reducing effort on the second task, which leads to the second sufficient condition. For some intuition on the economic meaning behind this condition, note that assuming that  $\tilde{u}'(x_1, \theta') - \tilde{u}'(x_1, \theta)$  is increasing in  $x_1$  is a way of formally assuming that  $\tilde{u}(x_1, \theta')$  is less concave than  $\tilde{u}(x_1, \theta)$ . This means that deviations away from  $x_1^*$  are less costly for  $\theta'$  than for  $\theta$  and so she responds more to a change in incentive. As we show in Appendix B, in the case of unconstrained optimization this condition can also be written as a condition on the third derivatives of  $u(x, \theta)$ ; of particular note, if  $u(x, \theta) = x_1 + x_2 - c(e)$  for some quadratic cost function  $c(e)$  – and hence the third derivatives are all zero – the condition will hold. We discuss this assumption in more detail, including how we can relax the assumption to one of single crossing condition on the marginal utility in Appendix B. We also show in there that the single crossing condition is not only sufficient but necessary if one wants to guarantee the result for all weakly increasing functions  $f(x_1)$ .

## II.D Impact of Incentive Effects on the Screening Efficiency:

Above, we discussed how individuals' response to the incentive induced by the screening policy varies both because of different incentives across individuals and because economic forces cause different individuals to respond differently even to the same incentive. From these results, it follows that we have a sufficient condition for when a higher ability teacher will respond more to the screening policy than a lower ability teacher who has the same test score value-added in the unincentivized world. We state this as a theorem:

**Theorem 2.** *Consider  $\theta < \theta'$  with  $x_1(\theta) = x_1(\theta')$ . Assume that  $\Delta v(\theta)$  is increasing in  $\theta$  and that either  $\tilde{u}'(x_1, \theta') - \tilde{u}'(x_1, \theta)$  is increasing in  $x_1$  or  $\lambda$  is sufficiently large. Then  $x_1^*(\theta|p) \leq x_1^*(\theta'|p)$  for any weakly increasing screening function  $p(x_1)$ .*

The theorem says that the individual with higher  $\theta$  responds more to the incentive,

---

<sup>16</sup>Specifically, this means that the same reduction in  $e_2$  leads to a larger reduction in  $x_2$  – and therefore  $u$  – for individual  $\theta'$  than for individual  $\theta$ .

which means that she is more likely to be retained.<sup>17</sup> Since we generally want to keep high  $\theta$  individuals in the profession, this seems to suggest that the incentive component of the screening policy actually makes the screening more efficient. Formalizing this result requires a bit of work, however, since for a fixed screening function the endogenous response also changes the fraction of teachers who are retained.

We start with the principal's value function – denoted  $V(\theta)$  – which determines how valuable she views keeping a teacher with skills  $\theta$  in the profession is. We make three assumptions about  $V(\theta)$ .

**Assumptions about  $V(\theta)$ .**

1.  $V(\theta)$  is increasing in  $\theta$ ;
2.  $\mathbb{E}[V(\theta)|x_1^*(\theta)]$  is increasing in  $x_1^*(\theta)$ ;
3.  $\frac{\partial V(\theta)}{\partial \theta_1} \leq \frac{\partial \Delta v(\theta)}{\partial \theta_1}$  and  $\frac{\partial V(\theta)}{\partial \theta_2} \geq \frac{\partial \Delta v(\theta)}{\partial \theta_2}$  for all  $\theta \in \Theta$ .

The first assumption says that the principal values the teacher's skills.<sup>18</sup> The second assumption does not follow from the first because  $x_1^*(\theta)$  is decreasing in  $\theta_2$ . In our empirical context, the observed measures  $x_1(\theta)$  and  $x_2(\theta)$  are positively correlated, which makes the assumption hold as long as the principal values both outputs.<sup>19</sup> The final assumption relates how the principal values a teacher's skills to how teachers value it themselves, as defined by how it affects their indirect utility of teaching relative to the outside option. It states that the principal values  $\theta_1$  by no more than the teachers themselves and values  $\theta_2$  by no less than the teachers themselves. While this is the strongest assumption, our motivation derives from the idea that the principal cares about both dimensions of skill but can only observe the first output. This assumption essentially states that she cares enough about the second dimension and, symmetrically, that she does not care too much about the first dimension.

Consider two screening functions  $p(x_1)$  and  $\tilde{p}(x_1)$ . We say that screening function  $\tilde{p}(x_1)$  is *more efficient* than function  $p(x_1)$  if both policies remove the same fraction of teachers and the average value of  $V(\theta)$  is higher for those teachers who remain after policy  $\tilde{p}(x_1)$

---

<sup>17</sup>This implies that for any increasing function of  $\theta$  – denoted  $V(\theta)$  – and any screening policy  $p$ , we have that the average value of  $V(\theta)$  is larger for those retained than those removed when comparing individuals who produce the same  $x_1$  in the unincentivized period, or that  $\mathbb{E}[V(\theta)p(x_1^*(\theta|p))|x_1^*(\theta)] \geq \mathbb{E}[V(\theta)|x_1^*(\theta)]\mathbb{E}[p(x_1^*(\theta|p))|x_1^*(\theta)]$  for any  $x_1^*(\theta)$ .

<sup>18</sup>One subtlety is that it rules out the case where the principal only cares about one of the outcomes. For example, if she defined value over the outcomes according to  $\tilde{V}(x) = x_1$ , then it would imply that  $V(\theta)$  is decreasing in  $\theta_2$  since  $x_1^*(\theta)$  is decreasing in  $\theta_2$ .

<sup>19</sup>If  $\mathbb{E}[V(\theta)|x_1^*(\theta)]$  is decreasing in  $x_1^*(\theta)$ , the ideal policy (absent incentives) would involve a screening function that is non-monotonic in  $x_1^*(\theta)$ . This gives the surprise screening policy an advantage over the announced screening policy, which has to account for the incentives change. We ignore this possibility because the screening rule in our empirical context is monotonic in test score value-added.

than after policy  $p(x_1)$ . In addition, as a matter of terminology, we will define an ex post (or surprise) screening policy as one that screens on  $x_1(\theta)$ , i.e., that screens on the equilibrium outcomes that occur absent the endogenous response.

We then conclude with the following theorem:

**Theorem 3.** *Assume the conditions on  $V(\theta)$  specified above and that the assumptions in Theorem 2 hold. Furthermore, assume that  $\theta$  is continuously distributed. Then for any ex post screening policy  $p(x_1)$ , there is a screening policy  $\tilde{p}(x_1)$  that is more efficient than  $p(x_1)$ :*

$$\begin{aligned}\mathbb{E}_{\Theta}[\tilde{p}(x_1^*(\theta|\tilde{p}))] &= \mathbb{E}_{\Theta}[p(x_1^*(\theta))] \\ \mathbb{E}_{\Theta}[V(\theta) \cdot \tilde{p}(x_1^*(\theta|\tilde{p}))] &\geq \mathbb{E}_{\Theta}[V(\theta) \cdot p(x_1^*(\theta))].\end{aligned}$$

The theorem states that under some assumptions the pre-tenure response to the changing incentives caused by the screening policy makes the screening more efficient. But it may still be the case that the pre-tenure distortions have a cost to the principal that is larger than the benefit generated from the improved screening. Our empirical work will consider effects on screening efficiency and how they compare to the size of the pre-tenure output distortion.

### III Data and Outcomes

#### III.A Data and Measuring Output

The model suggests that teachers respond to the incentives along multiple dimensions and that responses are likely to be heterogeneous across teachers. An ideal data set for empirical analysis thus has several features. First, it must include both measures of output that are targeted by incentives as well as measures of output that are untargeted and which capture skills distinct from the targeted measures. Relatedly, having a long-term outcome – and a long enough panel to observe the outcome – helps characterize the trade-off between gains in different short-run output measures based on how well they predict the long-term outcome. Second, the data set should track teachers over time and observe them with different levels of incentives. Such panel data, with within-teacher incentive variation, would allow the study of how different types of teachers respond differentially to incentives.

We meet these needs by turning to yearly administrative data from the New York City Department of Education. This data includes all NYC public schools from 2006-07 through 2014-15, which is a long enough panel to follow teachers over time and to observe longer-term outcomes for early cohorts. We infer a teacher’s output from the outcomes her students



achieve. For each student-year observation, we observe the student’s school and grade and end-of-year test scores, the outcome that personnel policy will target. We further observe other student outcomes – attendance rates and grades in tested and other subjects (for middle school students) – that we will use as untargeted outcomes, as we describe below, and students’ high school graduation status, a longer-run outcome. We also observe student demographic information and thus can construct a broad set of control variables that will help us isolate a teacher’s impact. Importantly, we observe a mapping between the student and the teacher she had in each subject-year. This mapping lets us link student outcomes to individual teachers and thus construct measures of teacher output.

On the teacher side, we can follow teachers over time and across schools within NYC. Our policy variation will depend on whether a teacher is tenured, and how the timing of her pre-tenure period lines up with policy changes. Our data include teachers’ experience each year and whether they are tenured in NYC.

We provide summary statistics for the student and teacher data in Table 1. Nearly 80% of our student sample is eligible for free or reduced price lunch (high-poverty), and just over 10% is an English language learner. We normalize the test scores to have mean 0 and standard deviation 1 for each grade-year. Students’ attendance rates are relatively high, averaging 94% with a standard deviation of 6%, while grades vary between 10 and 100 with means near 80% and standard deviations of about 10%. The mean teacher in our data has been in the district for nearly 9 years and at the same school for about 6 years. Just over one-fifth of the teacher-years correspond to the pre-tenure (probationary) period, which will be the focus of policy variation. NYC’s size leaves us with large sample sizes that enable us to have powerful empirical tests. We have nearly 29,000 teachers in the sample, and they contribute almost 100,000 teacher-years. At the subject-year level, we have an average of 32 students per teacher, such that the mean of a teacher’s students’ outcomes carries some signal of the teacher’s output.

### III.B **Constructing Teacher Outcomes**

We seek to estimate changes in teacher output along targeted and untargeted dimensions. In choosing student outcomes that may capture skill development separate from contemporaneous test score measures, we turn to a recent literature that estimates teacher effects on several student outcomes. Jackson (2018) shows that teachers’ effects on attendance and grades – student outcomes regularly collected by districts – capture information beyond a teacher’s effect on student test scores, while Petek and Pope (2016) and Kraft (2019) extend the analysis using psychological measurements. We follow this literature in choosing attendance (Gershenson, 2016; Liu and Loeb, 2021) and grades as behavioral outcomes that are not directly targeted by tenure policy.

In addition to heterogeneity in contemporaneous outcomes, the persistence of a teacher’s effect on a specific outcome may capture teacher heterogeneity relating to different modes of teaching. For instance, if a teacher focuses instruction on test-based memorization, the students may perform well on contemporaneous tests but not have the skills to build upon for future grades. We thus follow Carrell and West (2010) and Gilraine and Pope (2020) in measuring a teacher’s impact on the future realizations of test scores, grades, and attendance.

We will use these student outcomes for two purposes: measuring how teacher output changes in response to incentives and classifying teachers into heterogeneous types. As this introduces several identification and estimation challenges, we specify a statistical model of student outcomes. We provide an overview here, with more detail in Appendix C. Let  $i$  index students,  $j$  index teachers,  $c$  index classrooms,  $t$  index years, and  $k$  index outcomes. Student  $i$  has a vector of outcomes  $y_{it}$  where the  $k^{\text{th}}$  element of the vector is  $y_{i,t+\tau(k)}$ .  $\tau(\cdot)$  is a function that describes when an outcome is realized. For contemporaneous outcomes,  $\tau(k) = 0$ , while for outcomes realized in the future, like next year’s test scores,  $\tau(k) > 0$ . Without loss of generality, we place contemporaneous test scores as the first outcome ( $k = 1$ ). We standardize all outcomes to have mean 0 and standard deviation 1 for each grade-year. We model student outcomes as:

$$y_{i,t+\tau} = \Lambda X'_{it} + \sum_{e'} \rho_{e'} \mathbb{1}\{e_{jt} = e'\} + \mu_{jt} + \nu_{ct} + \phi_{c',t+1} \mathbb{1}(\tau \geq 1) + \phi_{c'',t+2} \mathbb{1}(\tau = 2) + \epsilon_{it} \quad (6)$$

where  $e_{jt}$  is a teacher’s experience level,  $X'_{it}$  is a vector of  $P$  student covariates that may include lagged outcomes, and  $\Lambda$  is a  $K \times P$  matrix of coefficients.  $\rho$  is a  $K \times 1$  vector of experience effects.  $\mu_{jt}$ ,  $\nu_{ct}$ , and  $\epsilon_{it}$  are  $K \times 1$  vectors of teacher-year effects, classroom effects, and student idiosyncratic variation, respectively.  $\phi_{c,t+\tau}$  are combined  $K \times 1$  classroom-teacher effects for students’ assignments in years after  $t$ . We start with an independence assumption:

**Assumption 1** (Independence of teacher, classroom, and student effects).

$$\begin{aligned} \mu_{jkt} &\perp (\nu_{clt}, \epsilon_{ilt}, \phi_{c',l,t+\tau}) | X_{it}, e_{jt} \quad \forall k, l, \tau \\ \nu_{ckt} &\perp (\mu_{jlt}, \epsilon_{ilt}, \phi_{c',l,t+\tau}) | X_{it}, e_{jt} \quad \forall k, l, c, c' \tau \\ \epsilon_{ikt} &\perp (\mu_{jlt}, \nu_{clt}, \phi_{c',l,t+\tau}) | X_{it}, e_{jt} \quad \forall k, l, \tau \\ \phi_{c',k,t+\tau} &\perp (\mu_{jlt}, \nu_{clt}, \epsilon_{ilt}) | X_{it}, e_{jt} \quad \forall k, l, \tau \end{aligned}$$

This assumption corresponds to selection of students to teachers based on observables ( $X_{it}$ ), but extended to a setting with multi-dimensional output. While we impose independence across effects that occur at different levels (e.g., teacher versus classroom), we allow

a given effect to be correlated across outcomes. Prior work has validated this assumption for most our measures with experimental variation (Kane and Staiger, 2008) or mover designs (Chetty et al., 2014; Delgado, 2021; Jackson, 2018; Gilraine and Pope, 2020). For our baseline analysis, we control for cubic functions of the  $t - 1$  outcome, with the exception of the subject-specific grade outcomes because often the lagged values are missing. In that case, we control for a cubic function of the  $t - 1$  test score.<sup>20</sup>

We estimate teacher  $j$ 's realized causal effect on outcome  $k$  in year  $t$ ,  $\hat{\mu}_{jkt}$ , by estimating Equation 6 with OLS for outcome  $k$ , constructing student-level residuals ( $y_{it} - \hat{\Lambda}X'_{it}$ ), and taking the teacher-year-subject mean over the residuals. This procedure yields our (noisy) measure of a teacher's annual realized output on each dimension. Teacher effects may be correlated across outcomes. A teacher who is effective at raising students' contemporaneous test scores may also be effective at raising students' future test scores. We do not yet specify a correlation structure for how a teacher's effectiveness varies over time. While pooling data across years would increase the precision of our teacher effect estimates, we seek to study teachers' responses to incentive changes. Thus, we do not want to impose structure on how a teacher's effectiveness varies over time and instead estimate teacher-year effects.

$\hat{\mu}_{jkt}$  will let us meet our first goal of measuring how a teacher's output changes in response to incentives. Our second goal is to classify teachers into heterogeneous types based on their unincentivized output. This introduces a few challenges. Teachers' output is not completely in their control, as classroom shocks or idiosyncratic student shocks mean that some years teachers may have higher or lower output than would be predicted by their type and effort. Thus, we want to develop a forecast of a teacher's mean output, where the forecast uses observed outcomes but adjusts for the presence of shocks. The precision of this forecast will be very limited if we only use data from a single year of a teacher's career. Hence, unlike the estimate of teacher's realized output, which was year-specific, we now want to pool data across a teacher's career to develop a forecast of her mean output in the unincentivized state. Pooling, though, imposes structure across years that might seem inconsistent with a setting where teacher output is a product of a teacher's type *and* the environment. We thus pool data only for years in the unincentivized state under the assumption that the constant environment lets us assume a teacher's type is fixed in expectation across years.

We now add an assumption about the structure of drift:

---

<sup>20</sup>This assumption is sufficient for all of our analysis. In fact, for our estimation of the causal effect of the change in incentives on teacher output, we can relax the assumption provided that any sorting of students (or classroom effects) to teachers is orthogonal to the policy variation. In this case, we do not need to control for  $X_{it}$  for identification purposes, but we may still do so to increase our estimates' precision. In Section V we show no evidence of sorting changes in response to the policy and show that our main results are robust to excluding controls.

**Assumption 2** (Joint stationarity of teacher effects without incentives). *The non-experience part of teacher value-added for each outcome follows a joint stationary process if there are no incentives. The covariances between the teacher’s value-added across years depend only on the number of years elapsed:*

$$\mathbb{E} [\mu_{jkt}|t] = \mathbb{E} [\mu_{jks}|t] = 0 \quad (7)$$

$$\text{Var} (\mu_{jkt}) = \sigma_{\mu_k}^2, \text{Cov} (\mu_{jkt}, \mu_{jk't}) = \sigma_{\mu_k \mu_{k'}} \quad (8)$$

$$\text{Cov} (\mu_{jkt}, \mu_{jk,t+s}) = \sigma_{\mu_k s} \quad (9)$$

$$\text{Cov} (\mu_{jkt}, \mu_{jk',t+s}) = \sigma_{\mu_k \mu_{k'} s} \quad (10)$$

for all  $k, k', t$  and  $s$ .

We follow Mulhern and Oppen (2021) in estimating the multi-year and multidimensional teacher value-added model; as Mulhern and Oppen (2021) discuss, if we further assume that the error terms are jointly normal, these estimates are also the empirical Bayes’ value-added measures and so we will generally refer to them as such. We estimate the model jointly across outcomes and using only data from unincentivized teacher-years, though the forecasts will apply to all years and provide the counterfactual of how the teacher would have affected her students’ outcomes in that if her incentives did not change. For each year, we construct a forecast of a teacher’s value-added using data from other (unincentivized) periods only. The forecast comes from a multidimensional empirical Bayes procedure (Mulhern and Oppen, 2021). We label the forecast  $\tilde{\mu}_{jt}$ .

These forecasts have several properties usually associated with current test score value-added. Within-teacher, the value-added on a given outcome is autocorrelated, which implies that the measures have some predictive power for adjacent years (Appendix Table A1). Second, for all measures, we see that teachers in their careers have rapid growth as they accrue experience (Appendix Figure A3). This experience profile is important because our empirical strategy will compare teachers at different experience points. We will thus need to control for changes in output we would expect from changes in experience, even in the absence of the policy. Further, the flattening out of the experience profile will allow us to pool teachers at high experience levels. Third, graduation outcomes are positively correlated with all of the measures, in a univariate sense (Appendix Table A2). We also see that value-added for each outcome is positively correlated with current test score value-added, though correlations are well below 1 (Appendix Table A1).

### III.C Combining Outcomes into an Index

Having multiple outcomes will allow us to estimate the effect of incentive changes on each outcome. But to interpret the results in the context of teachers’ targeting the development of specific skills, we seek to understand how distinct the outcomes are and their relative value in predicting a student’s long-term outcome. As the main distinction in the paper is whether an outcome is targeted by the tenure policy, we will keep current test score value-added as its own index and assess how the remaining outcomes co-vary.<sup>21</sup> In assessing these remaining outcomes, we face a missing data challenge where some teachers only have effects on a subset of outcomes. For instance, because we use grade 3-8 test scores, we do not have future test score value-added measures for 8th-grade teachers. We follow Mulhern and Opper (2021) in forecasting missing measures based on the covariances between the missing and non-missing measures.

We create the untargeted index by anchoring the measures to their relative predictiveness of whether a student graduates from high school.<sup>22</sup> Let  $Grad_{ijt}$  be whether student  $i$  graduated from high school, which we match to her teacher  $j$  in year  $t$ . We then regress graduation on measures of teacher  $j$ ’s value-added in year  $t$ :<sup>23</sup>

$$Grad_{ijt} = \omega' \tilde{\mu}_{jt} + v_{ij}, \quad (11)$$

where  $\omega$  is a vector of anchoring weights. We estimate using data from the unincentivized period and note that  $j$ ’s value-added in year  $t$  is estimated leaving out data from year  $t$ , so we avoid having the same classroom effects or student idiosyncratic shocks on both sides of the regression. We present the estimated weights in Appendix Table A3, which shows that most outcomes continue to predict graduation positively, with the largest coefficients on current attendance and future grades in untested subjects.<sup>24</sup> We use the estimated weights to construct two measures: targeted output ( $\tilde{\mu}_{jt}^T = \tilde{\mu}_{j1t}$ ) and an index of untargeted output

<sup>21</sup>Throughout, we will divide measures by whether they are ever targeted by the tenure policy (“targeted” vs. “untargeted”) and periods by whether they come with the stronger tenure incentives (“incentivized” vs. “unincentivized”).

<sup>22</sup>We also conduct a Principal Component Analysis. The first principal component loads heavily on grade measures while the second loads on measures of future outcomes (Appendix Table A4).

<sup>23</sup>Instead of an indicator for whether a student graduated, we use a graduation residual that residualizes the indicator with a cubic function of  $i$ ’s test scores in year  $t - 1$  and other student observable characteristics to improve precision.

<sup>24</sup>This anchoring regression uses cross-sectional teacher variation. Teachers, of course, may differ in other ways such that the cross-sectional relationship does not predict treatment effects when individual outcomes change. For our main results, we also present estimates that do not aggregate outcomes using the estimated anchoring weights. These measure-specific results also highlight that the results are not driven by the negative estimated weights on subject grades, and we get similar results with other weighting schemes that ensure all measures get positive weights.

$(\tilde{\mu}_{jt}^U = \frac{1}{\hat{\omega}_1} \sum_{k=2}^K \hat{\omega}_k \tilde{\mu}_{jkt})$ .<sup>25</sup> Because we divide the untargeted outcomes by  $\hat{\omega}_1$ , we measure both the targeted and index of untargeted output in units of (current) test score student standard deviations, which allows for comparisons in equivalent units that are common in the literature. We see that the index maintains the properties of the individual measures: high autocorrelation (Appendix Table A1), steep experience profile (Appendix Figure A3), and strong univariate predictor of graduation (Appendix Table A2).

Finally, because teacher output on the targeted measure tends to be quite correlated – in the cross-section, at least – with output on the untargeted measures, one might predict that when teachers respond to the increased test score value-added incentives, they produce commensurate gains in the untargeted measure. As a benchmark to compare the gains we estimate in Section V, we estimate a graduation anchoring function that does not condition on untargeted measures:

$$Grad_{ijt} = \gamma \tilde{\mu}_{j1t} + \nu_{ij}, \quad (12)$$

where  $\gamma$  is a scalar weight. For any estimated gains in the targeted measure, we can multiply them by  $\hat{\gamma}$  and compare that estimate to the estimated total gains  $\hat{\omega}_1 \hat{\mu}_{1t}$ . The difference between these is driven by changes in untargeted measures that deviate from the cross-sectional correlation.

## IV Context and Empirical Strategy

### IV.A Tenure Policy

In NYC, teachers become eligible for tenure after accumulating three years of teaching experience within the district. Once a teacher is eligible for tenure, the district may grant tenure, deny tenure, or extend the probationary (pre-tenure) period for further evaluation. Tenure denial makes the teacher ineligible to teach in the district while teachers who receive tenure are provided extra employment protections for the rest of their careers.

Before the 2009-2010 school year, nearly all eligible teachers in NYC received tenure. For example, during the 2007 and 2008 school years,<sup>26</sup> 94% of all eligible teachers received tenure (Loeb et al., 2015). The tenure process, however, changed dramatically starting in 2009. In November, 2009, Mayor Michael Bloomberg announced at a panel discussion at the Center for American Progress that not using student achievement scores to evaluate teachers up for tenure was “like saying to hospitals, ‘You can evaluate heart surgeons on any criteria you want - just not patient survival rates.’” He was therefore directing “our

<sup>25</sup>We also apply these weights to the unshrunk measures for further indices,  $\hat{\mu}_{jt}^T = \hat{\mu}_{jkt}$  and  $\hat{\mu}_{jt}^U = \frac{1}{\hat{\omega}_1} \sum_{k=2}^K \hat{\omega}_k \hat{\mu}_{jkt}$ .

<sup>26</sup>To simplify notation, we refer to each school year as the calendar year it ended; e.g., the 2009-2010 school year will be called the 2010 school year.

school Chancellor Joel Klein to ensure that principals actually use student achievement data to help evaluate teachers who are up for tenure this year.”<sup>27,28</sup> From that point forward, NYCDoE would begin to consider how effective teachers are at increasing their students’ test scores when deciding whether to give them tenure. This change, which mirrors many similar policies in other urban districts, proved controversial, as teachers argued decisions would incorporate unreliable measures and would affect the workplace environment negatively (McGuinn, 2012; Murphy et al., 2013; Bleiberg et al., 2021).

Perhaps due to the rapid implementation schedule and the contentious reaction among teachers, the details in how test score value-added affected tenure changed over time. For the 2010 school year, the district automatically coded a teacher as having “tenure in doubt” (“tenure likely”) if the 95% confidence interval of her value-added scores over the previous two years fell below (above) the median. Teachers whose confidence intervals included the median received no recommendation based on value-added. Although this coding only informed a final decision, the burden shifted onto the principal to argue why she was making a recommendation contrary to what the coding suggested.

In subsequent years, a low value-added score no longer automatically coded a teacher as having her tenure in doubt. Yet teacher value-added scores continued to be a major focus of the tenure process. In 2011, teachers with low value-added scores were flagged as having an “Area of Concern” and those with high value-added scores were flagged as having “Notable Performance.” In 2013 and on, value-added scores remained part of the tenure evaluation process, but the district provided no explicit guidance on their use.

Mayor Bloomberg’s announcement signaled a major shift in the teacher tenure policy in NYC. These changes had two first-order effects: they lowered tenure rates and they made tenure decisions increasingly dependent on value-added scores. Indeed, we observe a significant decrease in the probability that a teacher received tenure following the reform. We plot the fraction of newly tenure-eligible (fourth year) teachers receiving tenure over time in Figure 1. We see that tenure rates fell precipitously from 97% in 2010 to 64% in 2012.<sup>29</sup> When teachers did not receive tenure, they could either have tenure denied or have their probationary (pre-tenure) period extended. Loeb et al. (2015), who have access to the specific tenure decisions, show that most non-tenured cases led to extensions

---

<sup>27</sup>Bloomberg’s full remarks from this talk are available at: <https://www.c-span.org/video/?290247-1/white-house-education-agenda-state-us-schools>

<sup>28</sup>Bloomberg made this announcement despite the fact that the New York State Legislature had banned the use of value-added in tenure decisions the year before, because in his words, “after a very close reading our lawyers tell use that the current law... only applies to teachers hired after July 1, 2008.” In addition, because the law expired by the time teachers hired after July 1, 2008 were up for tenure, it never had any affect on the tenure-granting process in NYCDoE.

<sup>29</sup>We do not have access to specific tenure decisions in our data; instead, for each year we observe whether the teacher has tenure.

of the probationary period. But these extensions were not merely delays leading to the same outcome, as most of the extended teachers left their schools or even the district.<sup>30</sup> Thus, while the policy did not lead to a large increase in tenure denial rates, the policy still dramatically decreased the fraction of teachers continuing in the district with tenure.

Tenure rates remained fairly flat through 2010, the first year following Bloomberg’s announcement, while the substantial decrease in tenure rates came a year later in 2011. This introduces an important question of how to measure the policy’s timing. While the prior analysis shows how tenure outcomes changed over time, the relevant policy timing impact is when it first affected teachers’ *incentives* to achieve higher value-added scores. The policy’s announcement – at the start of the 2010 school year – marked the point when teachers became aware that value-added scores would matter for future tenure decisions. We further provide evidence of the increased focus on tenure during the 2010 school year by examining the mentions of “tenure” on the teachers’ union (UFT) website; in Figure 2 we show a spike in mentions in 2010, the school year of the announcement. Thus, for our analysis we will treat the 2010 school year as the first under the new policy regime.

Beyond affecting aggregate tenure rates, the policy also tied tenure more closely to test score value-added measures. Using data and value-added measurement methods described in Section III, we assess the relationship between tenure and test score value-added before and after the policy change. In Figure 3, we bin teachers into ten deciles according to their value-added scores during their third year of experience. The y-axis is the fraction of teachers in each bin who have achieved tenure by the end of their fourth year. We plot the relationships separately for the periods before and after the change in the tenure process. Prior to the policy change, we see very little relationship between a teacher’s value-added score and her probability of receiving tenure. This is not surprising as the high tenure rates left little variation to explain. After the policy change, however, value-added scores become strong predictors of tenure outcomes.

We further explore the tenure rules and how they vary with teachers’ test score value-added and other measures, by estimating linear tenure rules separately for each policy regime. We estimate the following screening functions:

$$Tenure_{jt} = \pi_0^0 + \pi_1^0 \tilde{\mu}_{jt}^T + \pi_2^0 \tilde{\mu}_{jt}^U + v_{jt}^0 \quad (13)$$

$$Tenure_{jt} = \pi_0^1 + \pi_1^1 \tilde{\mu}_{jt}^T + \pi_2^1 \tilde{\mu}_{jt}^U + v_{jt}^1, \quad (14)$$

where superscript 0 is for pre-reform (2008-2009) tenure decisions and superscript 1 indicates post-reform (2011-2012) tenure decisions.  $Tenure_{jt}$  is an indicator for whether the teacher

---

<sup>30</sup>In Appendix Figure A1, we show that after the policy, teachers who would be entering their fourth year were much less likely to continue teaching in NYC.



receives tenure in the first year she is eligible to. We present the screening function coefficients in Table 2. In the pre-reform period, the mean tenure rate is high (97%), and neither test score value-added nor the index of other outcomes enters statistically significantly. The coefficients on these measures are also quite small, indicating precise zero estimates. In the second column, we show the estimated screening function in the post-reform period. As expected, the mean tenure rate falls (to 67%) while test score value-added is now a strong, and statistically significant, predictor of receiving on-time tenure. Interestingly, output on the untargeted measures does not enter the screening function significantly. This confirms that teachers' tenure-based incentives around non-test score value-added output did not change due to the policy.

## IV.B Concurrent Events

Although the change in tenure policy provides nice within-teacher variation in their incentives, the policy did not occur in a vacuum. For example, about two years before the change in tenure policy (in Fall 2007), student test score growth began to affect the grade each school received (Rockoff and Turner, 2010). Around the same time, NYC established a pilot program that distributed information about teacher value-added to 112 principals (Rockoff et al., 2012). Within the next year, in Fall of 2008, these value-added measures were distributed to every teacher in NYC, although they were told that “they won’t be used in tenure determinations or the annual rating process.”<sup>31</sup> As discussed above, this decision was reversed in the following year by Mayor Bloomberg. It was not until 2012, however, that the teacher data reports were made public.

At the same time as the policy was being announced and implemented, NYC was feeling the beginnings of the Great Recession, which likely affected the types of teachers who were entering the teaching market (Nagler et al., 2015). Furthermore, the Great Recession caused a large financial shortfall for the NYC government, leading the Chancellor to institute a hiring freeze of new teachers. While there were some exceptions, the fraction of first-year teachers in NYC plunged from around 10% of total teachers to a mere 2% (Appendix Figure A2).

## IV.C Empirical Strategy

Our first goal is to estimate the behavioral response of pre-tenured teachers to the added incentive. While posing potential threats to identification, the concurrent events described above also suggest the types of comparisons that might isolate a causal effect. Unlike the

---

<sup>31</sup>This quote comes from the New York Times article on the policy titled, “Teachers to Be Measured Based on Students’ Standardized Test Scores,” which was written on October 1, 2008 and cites the quote as coming from a memo written by Chancellor Klein.

school report card and value-added information dissemination policies, the tenure reform only affected a portion of the teachers: the untenured teachers. In short, the change in the tenure process caused a sudden increase in untenured teachers' incentives to increase their value-added scores relative to their tenured colleagues, which suggests the use of a difference-in-difference estimator. Consider teacher  $j$  with outcome  $y_{jt}$  in school year  $t$ . A simple difference-in-difference specification would then be:

$$y_{jt} = \tau \text{UntenuredPost}_{jt} + \nu \text{Untenured}_{jt} + \eta_t + \epsilon_{jt}, \quad (15)$$

where  $\text{UntenuredPost}_{jt}$  is an indicator for  $j$  being untenured in year  $t$  and the tenure policy change being in place. This specification thus compares how untenured and tenured teachers' outcomes differentially changed once the tenure policy was announced.

While straightforward, this specification is not sufficient since the composition of these two groups is changing over time for several reasons. First, reduced hiring of new teachers during the Great Recession may have led recent new teacher cohorts to differ from prior cohorts. If more selective hiring led to more productive incoming teachers, then we might infer that the tenure policy improved untenured teachers' output when it was simply selection. Second, the Great Recession or the tenure incentives themselves may have affected which NYC public school teachers choose to remain teaching in the district. Third, since the tenure reform directly affects teacher tenure rates tenure status is endogenous with respect to the policy. In the extreme, if low-output teachers have delayed tenure decisions but keep teaching in NYC, then the composition of the group of untenured teachers would decline in quality over time.<sup>32</sup>

We deal with these challenges by making three adjustments to the base difference-in-difference specification. First, we address the endogeneous selection into receiving tenure by making a sample restriction that will apply to all of our specifications. Namely, we exclude teachers once they have become tenure eligible – that is, they have accumulated at least 3 prior years of experience – under the new tenure policy. These excluded observations correspond to periods when some of these teachers may have tenure while others may still be in their probationary period. If we had left these teachers in the analysis, then comparisons across (current) tenure status would likely involve a large degree of selection.<sup>33</sup> We will

<sup>32</sup>While complicated the identification Great Recession does have one advantageous factor for our analysis; a weak labor market plausibly increases the incentive NYC teachers have to get tenure, and so the Great Recession indirectly increased the strength of our policy instrument.

<sup>33</sup>Alternatively, we could have left these observations in the analysis and instrumented for tenure incentive with whether the teacher was in the *standard* probationary period – i.e., fewer than 3 years of prior experience. We choose to restrict the sample because we worry about monotonicity of the proposed instrument. When a cohort of teachers advances from 2 to 3 years of prior experience and the instrument switches values, those teachers who did receive tenure see a large decrease in their incentives while those teachers whose probationary periods were extended have even higher incentives than before.

therefore only use tenured teachers who received tenure prior to the policy change, when tenure rates were nearly 100%, such that both the untenured and tenured teachers in our analysis will have been subject to minimal involuntary attrition.

Second, we assign each teacher  $j$  to a cohort  $m(j)$  based on the year  $j$  started teaching in NYC and include cohort fixed effects in the empirical specification. If the Great Recession changes selection among new teachers based on productivity levels, then the cohort fixed effects should control for any cross-cohort differences upon entry. Thus, instead of looking across-teachers, we will use the fact that a few cohorts of teachers were teaching in NYC (in the probationary period) when the policy was announced. They therefore spent their initial year(s) teaching without test score value-added incentives, then were suddenly told that their tenure would depend on their test score value-added. We will thus explore how they responded to this sudden change, on both the targeted and untargeted tasks.

Adding cohort fixed effects, though, complicates the untenured versus tenured comparison, as within-cohort changes may differ between these groups for reasons unrelated to the policy itself. In particular, a large literature on the shape of the teaching experience profile has documented the possibility that early-career teachers have a steeper profile (Rockoff, 2004; Rice, 2013) and it is precisely the early-career teachers who are untenured.<sup>34</sup> Hence, as we add cohort fixed effects, we also introduce a vector of fixed effects for each level of prior experience in NYC; we combine teachers with six or more years of prior experience. This second adjustment implies that we will identify the policy effects based on how output growth rates change differentially between untenured and tenured teachers based on the policy implementation, beyond what we would expect from generic experience effects.

We summarize how these research design choices affect identification in Table 3. We show teacher cohorts - based on when they entered NYC - in the rows and the academic years in the columns. In the intersection of a row and column, we show the cohort's incentive status, which is whether the cohort's teachers are in the probationary period *and* NYC has already incorporated value-added into the tenure decision process. Because we include cohort fixed effects, our identification will come from cohorts (rows) for whom the incentive status varies across columns (i.e., the 2008 and 2009 cohorts).

We translate this variation into the following empirical specification:

$$y_{jt} = \tau Incentive_{jt} + \lambda_e + \nu_m + \eta_t + \epsilon_{jt}, \quad (16)$$

where  $e = e(j, t)$  is  $j$ 's level of prior experience in year  $t$  and  $Incentive_{jt}$  is an indicator for  $j$  being in the first 3 years of teaching and the tenure policy change being in place. Our

<sup>34</sup>Wiswall (2013) and Papay and Kraft (2015) show that using different identifying restrictions can generate a linear experience profile. Our framework could accommodate imposing a linear experience profile for later-career teachers.

parameter of interest is  $\tau$ , which is the coefficient of our covariate that indicates whether teacher  $j$  had incentives in time  $t$ . We cluster our standard errors by teacher and find similar standard errors when we cluster by school. We will also show permutation tests where we estimate a distribution of placebo effects by permuting the exposed cohorts (Abadie et al., 2010; Idoux, 2021).

This strategy does not deal with the other compositional concerns though. If the tenure policy or Great Recession changes attrition patterns, our comparisons will still be confounded. We will therefore run analyses that include teacher fixed effects instead of cohort fixed effects:

$$y_{jt} = \tau Incentive_{jt} + \lambda_e + \nu_j + \eta_t + \epsilon_{jt}. \quad (17)$$

Any selective within-cohort attrition will then be controlled for, provided it is based on time-invariant differences across teachers (rather than, say, growth rates).

To summarize, our empirical strategy relies on the assumption that after accounting for yearly shocks and selection into and out of teaching based on time-invariant unobservables, the returns to experience for teachers in the probationary period would not have changed post-2009 in the absence of the policy change. We will attribute any systematic deviation from the “unincentivized” experience profile to the policy change.

#### IV.D Selective Attrition

Before we turn to the effects on output, we return to the possibility of selective voluntary attrition. Indeed, attrition may be a response to the policy itself, as teachers unlikely to get tenure may attrit early, once the policy is announced. While the specification with teacher fixed-effects in theory controls for this attrition, we still explore attrition and its relationship with our instrument.

In Figure 4, we plot survival curves for cohorts grouped by their exposure to the treatment. The solid blue line, for instance, shows the survival curve for teachers in never exposed cohorts – i.e., those who had already completed the standard probationary period before the policy change. At the other extreme, the dashed orange line shows cohorts always exposed – i.e., those who entered the district after the policy’s announcement. We show unincentivized years with dots and incentivized years with diamonds. The survival curves largely lie on top of each other, indicating no differential attrition based on policy exposure.

Even with similar cohort-level attrition rates, we could still have selective attrition if the composition of the attriters changes in response to the policy. We again split cohorts based on policy exposure and plot the difference in mean test score value-added between stayers (teachers who stay in the district) and leavers (Figure 5). Again, we see no discernible relationship. For instance, in the third year of teaching, the blue (never exposed) and green

(exposed for two years) lines lie on top of each other, despite the green line representing a cohort that had the incentive treatment in their third year of teaching. This lack of a systematic pattern of attrition thus makes us less worried that selective attrition explains our results.

## V Estimated Effects on Teacher Output

### V.A Targeted Measure

We start by showing the effect of the tenure policy on the targeted measure: current test score value-added.

We show the empirical strategy visually in Figure 6, where we group teaching cohorts based on their policy exposure (i.e., never exposed, exposed for one year, etc.). We then plot each composite cohort’s mean test score value-added (after netting out year and experience effects) during the first three years of their probationary period. We see that in cohorts’ first years in NYC, they differ significantly in terms of their output. The always exposed cohorts, represented by the orange line, have higher initial test score value-added than the other composite cohorts. This could be a consequence of the policy, as these are the only cohorts whose first years teaching occur under the new tenure policy. But as highlighted above, the differential selection into teaching may also be driven by other factors like the Great Recession. Hence, we are hesitant to argue that the reform is the sole driver of this difference.

Instead, we will look at how output evolves over time, *within-cohort*. We see that for one year of prior teaching experience, the green line (which represents the cohorts exposed for 2 years) transitions from not being treated with tenure incentives to being treated. This is also the cohort that has the largest gain between years 0 and 1. Then for output after two years of prior teaching experience, we see that the cohort newly treated in this year (the red line) jumps up to “join” the other treated lines. The fact that the deviations from the unincentivized pattern occur in the specific years the cohorts receive the incentive treatment increases our confidence that our empirical strategy is picking up response to the policy.

We complement the graphical evidence with estimation of Equations 16 and 17, where we use current test score value-added as the outcome. We present the estimates in Table 4. In column (1), we show that that – consistent with Figure 6 – the specification that includes cohort fixed effects suggests that teachers responded in meaningful ways to the change in incentive. Specifically, it suggests that the tenure incentive increased value-added by 0.033 student standard deviations ( $\sigma$ ) and is statistically significant at the 1% level. While this specification includes cohort fixed effects to control for changes in selection into teaching over time, the estimated effect may still reflect some compositional changes from attrition.

For example, if inexperienced teachers who are unsure whether they would like to make a career out of teaching are more likely to attrit during the Great Recession, and these teachers tend to have lower value-added, then we might incorrectly attribute the output gains to response to the policy’s incentive changes. The visual results of Section IV.D suggest this is not a huge issue, but we confirm this by replacing cohort fixed effects with teacher fixed effects in our preferred specification. As shown in column (2) of Table 4 when we do so, we find that the effect of the policy is slightly smaller but still large:  $0.015\sigma$ . It remains statistically significant at the 10% level. In Appendix Figure A4, we conduct a permutation test and show our estimates are more extreme than all but 1 placebo estimate (column 1) or than all placebo estimates (column 2).

At the bottom of the table, we present the implied effect on graduation rates based on estimates of the cross-sectional relationship between test score value-added and graduation ( $\hat{\gamma}$  from estimating Equation 12). The implied impact on (residualized) graduation rates is 0.0028 or 0.0013 for the specifications with cohort and teacher fixed effects, respectively. But because this extrapolation exercise projects graduation onto test score value-added without any other controls, these predicted gains assume the cross-sectional relationship between test score value-added and untargeted output holds fixed. As an alternate exercise, we can extrapolate the test score value-added effects assuming other skills are unchanged by estimating Equation 11 and multiplying the test score value-added policy effect by  $\hat{\omega}_1$ . As shown at the bottom of Table 4, the implied effect on graduation rates is smaller. Thus, in evaluating the effect of the behavioral response on long-term output during the probationary period, the policymakers need to know whether the structure of skill production remains constant or whether it changes with the policy (e.g., due to multitasking).

These estimates are economically significant. The response to incentives ( $0.015 - 0.033\sigma$ ) is equivalent to 11-23% of the cross-sectional standard deviation of forecasted teacher value-added (0.142).

We explore heterogeneity in treatment effects based on subject tested and grade level in Appendix Table A5. We find that the effect is slightly larger in elementary grades and similar across subjects. But we lack precision to either reject equal effects by level or subject or to rule out meaningful differences.

Thus, the policy had large incentive effects, which increased teachers’ immediate output in the targeted measure. But because teachers produce output in multiple dimensions, multitasking may have led teachers to substitute out of effort developing untargeted skills.

## V.B *Untargeted Measures*

We test for the policy’s effects on untargeted measures with the same specifications (Equations 16 and 17), but varying the type of teacher output. We start by using our index of

untargeted measures, which we specified in Section III, and present the results in the last two columns of Table 4. Unlike output in the targeted measure, the index of untargeted measures does not increase in response to the tenure policy change. Instead, as shown in columns (3) and (4), we find statistically significant decreases that are larger in magnitude than the increases in the targeted measure. In Appendix Figure A4, we conduct a permutation test and show our estimates are more extreme than all placebo estimates. The negative effect is surprising if one naively extrapolates from the positive cross-sectional relationship between measures and the positive effect on test score value-added to predict the policy response on the other measures. But as we developed in our model, this substitution away from untargeted measures is consistent with teachers multitasking and distorting effort toward tasks that affect targeted measures.

As with the targeted measure, we can extrapolate, based on the cross-sectional relationship, what this policy effect implies for graduation rates. We find negative changes in implied graduation rates due to the reduction in output in the untargeted measures. These negative changes more than cancel out the positive effects from the targeted measures.

We further decompose this effect to show how each of our untargeted measures changes in response to the policy. We present the teacher fixed effects results in Table 5.<sup>35</sup> We examine the effects on future test scores (in  $t + 1$  and  $t + 2$ ), current and future attendance, current and future grades, and current and future grades in untested subjects. Consistent with the effect on the index, we find statistically significant decreases in  $t + 1$  test scores and  $t + 1$  grades in both tested and untested subjects (and a marginally significant decrease on  $t + 1$  attendance). These effects are reasonably large. In standard deviation units, these estimates dominate the effect on current test scores. We find a statistically significant increase in current grades in the tested subject, arguably the measure most closely related to the targeted one. For the other outcomes, we fail to reject no change. While null effects may not necessarily indicate evidence of multitasking – say, the teachers maintain effort on untargeted measures and effort on the targeted measure has no spillover effect – negative effects likely indicate that teachers are reducing their effort on untargeted measures (unless spillovers are negative, which we consider unlikely).

Though our untargeted measures do not necessarily map neatly into lower-dimensional representations of skills, we might expect that future test scores are more likely to capture development of cognitive skills while attendance and grades are more likely to capture development of behavioral or non-cognitive skills (Jackson, 2018). We see similar effects across measures in these two groups. Instead, we see a clearer division between the effects on current versus future measures. Thus, teachers may be substituting toward tasks with

<sup>35</sup>Appendix Table A6 shows the estimates from a specification with cohort fixed effects.

short-run payoffs but that fail to build (cognitive or non-cognitive) skills that persist.<sup>36</sup>

One worry about the untargeted outcomes is that, unlike the targeted outcome, lagged outcomes may be unavailable or poor predictors of current or future outcomes. If lagged outcomes are necessary controls for estimating causal effects, then we could be picking up effects that are not causal. As we mentioned in Section IV, our identification strategy (for this section) does not require causal estimates of teacher effects but rather that students do not change sorting to teachers systematically in response to the policy. We test this directly by replacing the outcome in Equation 17 with a predetermined student characteristic. We show the estimates in Table 6. We find precise zero effects of the incentive on student sorting to teachers based on several observable characteristics. Given no change in student sorting, we can estimate the effects of the incentive on the (unresidualized) untargeted outcomes. We present the results in Appendix Table A7 and find very similar patterns to our results using residualized outcomes.

A second worry is that future outcomes depend on the following teacher’s actions. For our analysis, we might incorrectly attribute changes in, say, value-added on  $t + 1$  test scores, as reflecting the current teacher’s response to incentives rather than the subsequent teacher’s response. In Appendix Tables A8 and A9, we show that controlling for the treatment status of the subsequent teachers hardly changes our estimates.

## V.C Persistence

In Section VII, we will look at the effects of these behavioral responses on predicted output in the post-probationary period. Such effects could operate through two channels: changes to screening and changes to teacher output (even once tenure incentives no longer matter). We will focus on the former in Section VII, but now we examine the latter. Changes to teacher output could persist if teachers respond to the incentives by making investments that change their future production functions. For instance, a teacher might develop new lesson plans and continue to use them after receiving tenure. We test whether the responses to incentives persist, even once the incentives disappear, with the following specification (for teacher fixed effects):

$$y_{jt} = \tau Incentive_{jt} + \phi PostIncentive_{jt} + \lambda_e + \nu_j + \eta_t + \epsilon_{jt}. \quad (18)$$

$PostIncentive_{jt}$  is an indicator for whether the teacher has tenure but faced the new tenure policy at some point during her probationary period. We present the estimates, for models

---

<sup>36</sup>We find the strongest effects for future grades. The Principal Component Analysis shows these measures entering both the first and second components with large weights, which suggests substitution is not along the lines that principal components capture.



with cohort or teacher fixed effects, in Table 7. In this analysis, we drop our sample restriction to allow teachers to be in the analysis for both the (incentivized) probationary period and the post-probationary period. For the targeted measure, we estimate that teachers revert back to their pre-incentive levels (excluding the experience profile). For the untargeted measures we estimate a coefficient in the same direction as the incentive effect, though we fail to reject complete reversion. Thus, we do not find strong evidence that teachers' responses to incentives had persistent effects by changing their future teaching output.

## VI Heterogeneous Responses

We now extend our empirical model to allow for heterogeneous responses across teachers based on differences in their output in the unincentivized regime. This is a relevant source of heterogeneity for two reasons. First, because tenure is absorbing, once teachers have been screened for tenure, their output incentives disappear. As we saw in Section V, teachers' post-tenure output reverts to their unincentivized output (once we have removed the experience gradient). Thus, the teachers' unincentivized output arguably provides the best measure of how valuable a teacher will be to the district after she receives tenure. If the teachers with the highest unincentivized output are the ones most likely to respond to the incentive, then the behavioral response will make the screening component of the policy even more effective.

Second, such heterogeneity maps nearly into our theory model developed in Section II. While our model's main source of heterogeneity is the teacher's unobserved type  $\theta$ , the model results are based on output in the unincentivized state, denoted in the model as  $x_1^*$  for the targeted output and  $x_2^*$  for the untargeted output.<sup>37</sup> Specifically, the model provides conditions which imply that, conditional on  $x_1^*$ , teachers with a higher  $x_2^*$  will increase their targeted output in the incentivized state more than those with lower  $x_2^*$ . Hence, we can use what we observe without the tenure incentives to classify teachers and test whether teachers responded to the tenure incentives in a similar way.

To classify teachers, we use the fact that the multi-year value-added model from Section III yields forecasts for each teacher-year-measure. Specifically,  $\tilde{\mu}_{jt}^T$  is the forecast for teacher  $j$ 's test score value-added (absent the common experience profile) in year  $t$  and  $\tilde{\mu}_{jt}^U$  is the forecast for an index of the untargeted measures. We start by examining heterogeneous responses based on  $\tilde{\mu}_{jt}^T$  or  $\tilde{\mu}_{jt}^U$  separately before classifying heterogeneous responses based on their joint distribution.

---

<sup>37</sup>Our model does not have a stochastic component to output whereas observed output does. The predictions therefore pertain to forecasted output – i.e., the predictable part of output that is attributable to the teacher.

As in Section V, we first show our result graphically, in Figure 7. Here, we fix a cohort and year and regress a teacher’s unshrunk test score value-added (i.e.,  $\hat{\mu}_{jt}^T$ ) on  $\tilde{\mu}_{jt}^U$  and plot the coefficients. In the first year of teaching, this coefficient is about 0.045 for all three cohorts. Note that for all of the cohorts plotted, the first year teaching comes prior to the tenure reform, so this coefficient captures the positive cross-sectional relationship between a teacher’s effect on current test scores and other output that exists where there are no additional incentives. We then show how the cross-sectional relationship changes as experience accrues and as incentives change due to the tenure reform. Focusing on the cohort never exposed to the reform, we see the cross-sectional relationship holds steady. For the other cohorts, in contrast, we see immediate (and persistent) increases in the regression coefficient once the reform is implemented. The sudden change in the joint distribution indicates that the response to the reform’s incentives was stronger among teachers with higher output on the untargeted measures (in the unincentivized period).

To show the results in regression form, we modify our main estimating equations (Equations 16-17) to allow for heterogeneous impacts (in the incentivized period) based on output forecasts from the unincentivized period. The modified teacher fixed effects specification is:

$$\begin{aligned} \hat{\mu}_{jt}^T &= \tau Incentive_{jt} + \xi_1 \tilde{\mu}_{jt}^T Incentive_{jt} + \xi_2 \tilde{\mu}_{jt}^U Incentive_{jt} \\ &+ \sum_{e'} \pi_1^{e'} \mathbb{1}\{e_{jt} = e'\} \tilde{\mu}_{jt}^T + \sum_{e'} \pi_2^{e'} \mathbb{1}\{e_{jt} = e'\} \tilde{\mu}_{jt}^U \\ &+ \lambda_e + \nu_j + \eta_t + \epsilon_{jt}. \end{aligned} \tag{19}$$

We control for the predictiveness of  $\tilde{\mu}_{jt}^T$  and  $\tilde{\mu}_{jt}^U$ , without incentives, flexibly by letting it vary by experience level. Because we are interested in how teachers respond to the incentive by shifting into the targeted measure, our outcome is the teacher’s mean (current) test score residual ( $\hat{\mu}_{jt}^T$ ).

We present the estimates in Table 8. In columns (1) and (2), which differ based on the level of fixed effects, we focus only on heterogeneity based on  $\tilde{\mu}_{jt}^T$  (i.e., imposing  $\xi_2 = 0$  and  $\pi_2^{e'} = 0 \forall e'$ ). We see some evidence that teachers with higher forecasted value-added in the targeted measure respond more strongly to the tenure policy incentives, though we lack statistical precision to make confident statements. In columns (3) and (4), we compare teachers’ responses based on  $\tilde{\mu}_{jt}^U$  and find that higher forecasted value-added in the index of untargeted measures predicts a larger increase in current test scores. This is consistent with Figure 7.

Columns (5) and (6) include both forms of heterogeneity jointly. As in the one-dimensional heterogeneity analysis, we see that teachers with higher forecasts of current test score value-added may respond more than teachers with lower forecasts, but we cannot statistically

reject no differential response. We also see that the response to the incentive is stronger for teachers with higher forecasts of the untargeted measures ( $\tilde{\mu}_{jt}^U$ ), and we can reject equal response at the 5% level. In Appendix Figure A5, we conduct a permutation test and show this estimate is more extreme than all placebo estimates. While the estimated coefficient on  $\tilde{\mu}_{jt}^U$  is smaller than the estimated coefficient on  $\tilde{\mu}_{jt}^T$ , cross-sectional differences in  $\tilde{\mu}_{jt}^U$  are over 3 times as large as cross-sectional differences in  $\tilde{\mu}_{jt}^T$ . Thus, a one standard deviation difference in  $\tilde{\mu}_{jt}^T$  translates to a similar estimated response as a one standard deviation difference in  $\tilde{\mu}_{jt}^U$ .

This result is important for policy and consistent with our model. The behavioral response to the screening-induced incentives leads teachers to substitute into the targeted measure. We see higher degrees of substitution among teachers who are better (without incentives) on the untargeted measure. The screening policy thus leads to a multitasking problem that lowers output on untargeted measures in the short-run but increases output on untargeted measures in the long-run by selecting different teachers.<sup>38</sup> We quantify this long-term effect in the next section.

## VII Estimated Effects on Screening Efficiency

In the previous section, we showed that the multitasking response to the screening policy was heterogenous across teachers. The behavioral response thus affects who receives tenure and therefore may impact the screening efficiency of the policy. In this section, we quantify these effects.

For each teacher  $j$ , we require 4 objects for our calculation: (a) forecasted targeted output in the probationary period, with incentives,  $x_{j1}^w$ ; (b) forecasted targeted output in the probationary period, without incentives,  $x_{j1}^{wo}$ ; (c) forecasted output in the post-probationary period,  $x_{j1}^p$  and  $x_{j2}^p$ . We then define a tenure screening policy that keeps the top  $p\%$  of teachers according to their output on the targeted dimension. Letting  $r(\cdot)$  be a function that converts teacher's output on the targeted measure to a percentile ranking, we calculate the expected post-tenure output (for output dimension  $d$ ) with behavioral responses as:

$$\mathbb{E}(x_d|w) = \mathbb{E}(x_{jd}^p | r(x_{j1}^w) > p) \quad (20)$$

---

<sup>38</sup>Our model offers two different sufficient assumptions for producing this result – lower cost to substituting and higher desire to substitute. While we do not have proxies for the  $\tilde{u}$  function, we can examine whether teachers with higher  $x_2$ , conditional on  $x_1$ , value the post-tenure period more. Specifically, we compare voluntary attrition rates for tenured teachers based on their (unincentivized) output on the targeted and untargeted dimensions. In Appendix Table A10, we find that both dimensions negatively predict voluntary attrition, which is consistent with the model's assumption.

and the expected post-tenure output without behavioral responses as:

$$\mathbb{E}(x_d|wo) = \mathbb{E}(x_{jd}^p|r(x_{j1}^{wo}) > p). \quad (21)$$

The difference in these expectations is the impact of the behavioral response on the composition of tenured teachers. Note that the incentives only matter through the selection margin because if teachers are granted tenure, the amount they produce does not depend on how much they responded to the probationary period incentives.

We use our prior analysis to estimate the 4 necessary objects. We use our multi-year value-added model from Section III to forecast  $x_{j1}^{wo}$ :

$$x_{j1}^{wo} = \tilde{\mu}_{jt}^T. \quad (22)$$

This forecast relies on data only from the unincentivized periods. We further impose that post-probationary output also matches this forecast (for all dimensions):

$$x_{j1}^p = \tilde{\mu}_{jt}^T \quad \text{and} \quad x_{j2}^p = \tilde{\mu}_{jt}^U. \quad (23)$$

This assumption is motivated by our analysis showing that once they receive tenure, teachers revert to their pre-incentives level of output.<sup>39</sup>

Finally, we use the estimated coefficients from Equation 19, to relate targeted output with and without a behavioral response:

$$x_{j1}^w = x_{j1}^{wo} + \hat{\tau} + \hat{\xi}_1 \tilde{\mu}_{jt}^T + \hat{\xi}_2 \tilde{\mu}_{jt}^U \quad (24)$$

The key assumption is that we have summarized the heterogeneity in treatment effects with our specification.<sup>40</sup> The behavioral response will only change the composition of tenured teachers if  $\hat{\xi}_2 \neq 0$ . Otherwise, any behavioral response might increase short-run output but does not change the ordering of teachers.<sup>41</sup>

We focus our analysis on the cohort of teachers that entered NYC in 2007 (and were

---

<sup>39</sup>For the analysis, we consider the screening for a fixed cohort, based on forecasted output in a specific year,  $t$ .

<sup>40</sup>Our economic object of interest is how the treatment effect heterogeneity relates to a teacher's (estimated) type. While for some analyses, any other forms of treatment effect heterogeneity that are orthogonal to a teacher's (estimated) type would not affect the conclusions, here we plug these predictions into a non-linear screening function.

<sup>41</sup>The requirement that  $\hat{\xi}_2 \neq 0$  to change selection is specific to the tenure policy that keeps a fixed fraction of teachers. Compositional changes would still occur under counterfactual policies that apply an absolute threshold for receiving tenure, even with  $\hat{\xi}_2 = 0$ . Furthermore, we focus on a threshold policy with a deterministic rule. If, instead, the policy had some stochastic component, then  $\hat{\xi}_1 > 0$  would imply that increased dispersion in test score value-added, which in turn would increase the signal districts have in its tenure decisions. This could be an additional screening benefit to explore in future work.

first eligible for tenure in 2010) and apply a 67% tenure rate, which is similar to the tenure rate in NYC post-tenure reform.<sup>42</sup> In Figure 8, we show the joint distribution of  $x_{j1}^{wo}$  and  $x_{j1}^w$ . In Section VI we showed that conditional on  $\tilde{\mu}_{jt}^U$  we saw a limited differential response along the targeted dimension. But because  $\tilde{\mu}_{jt}^U$  covaries with  $\tilde{\mu}_{jt}^T$ , the joint distribution of targeted output with and without the behavioral response can rotate. We see indeed see a slight rotation of the relationship between  $x_{j1}^{wo}$  and  $x_{j1}^w$ . The distribution of test score value-added becomes more dispersed without much change in teachers' relative positions.

To more directly understand how the behavioral response affects screening along the two dimensions, in Figure 9 we order teachers based on their forecasted two-dimensional output in the post-probationary period along the x- and y-axes. We divide each teacher into one of four categories and label the teacher's point on the figures accordingly. The categories indicate both whether the teacher receives tenure when there is no behavioral response and whether the teacher receives tenure when there is a behavioral response. The solid and dashed lines show the tenure cutoffs in the two regimes, under a hypothetical policy in which 67% of teachers receive tenure. In the top panel, we zoom in on teachers with test score value-added between  $-0.1\sigma$  and  $0.1\sigma$  – i.e., those close to the tenure cutoffs – while in the bottom panel we zoom out and show all teachers.

When there is no behavioral response, tenure depends solely on one's unincentivized test score value-added, so we see a flat line at the cutoff. When there is a behavioral response, in contrast, we see that the tenure cutoff line rotates clockwise due to the fact that those with higher unincentivized value-added on the untargeted measures (the x-axis) respond more to the policy. Thus, teachers with high levels of unincentivized value-added on the untargeted measures may still receive tenure despite low levels of unincentivized test score value-added. The triangles show the teachers who receive tenure only when there is no behavioral response and the circles show the teachers who receive tenure only when there is a behavioral response. These groups are relatively small, comprising of only 4.3% of all teachers. We only see a limited number of changed tenure decisions because the tenure policy employs a threshold rule such that many behavioral responses are among teachers who are ex post inframarginal. But despite a very concentrated change in the composition of tenured teachers, these teachers vary dramatically in their  $\tilde{\mu}_{jt}^U$ . The *mean* difference in  $\tilde{\mu}_{jt}^U$  across teachers shifted into tenure versus out of tenure is  $1.73\sigma$ , or 239% of the cross-sectional standard deviation in  $\tilde{\mu}_{jt}^U$ . The differences in mean  $\tilde{\mu}_{jt}^T$  are much smaller:  $0.024\sigma$ , or 18% of the cross-sectional standard deviation.

We quantify the changes in screening from the behavioral response ( $\mathbb{E}(x_d|w) - \mathbb{E}(x_d|wo)$ ) and present the results in Table 9. We see that mean  $x_{j1}^p$  falls by just  $0.001\sigma$  (0.76% of the cross-sectional standard deviation) while mean  $x_{j2}^p$  increases by  $0.056\sigma$  (7.7% of the

<sup>42</sup>Appendix Table A11 shows the results for the 2006 cohort.

cross-sectional standard deviation).

Teachers' behavioral responses to the policy thus change the composition of the tenured employees. Does the district prefer the new composition? This depends on how the district values output in the targeted measure versus output in the untargeted measures. In our theoretical model, we showed that as long as principals value the untargeted output at least as much as teachers do, then the district prefers the new composition. Without imposing the theoretical assumptions, we can make quantitative statements based on our estimates. In particular, because we anchored the measures to their predictiveness of graduation rates, they are measured on comparable units. Hence, we see that the district's screening has become more efficient in graduation units, as the improved selection on the untargeted measures exceeds the worsened selection on the targeted measures. If the district has reasons to value the outputs beyond their predictiveness of graduation rates, we estimate that the district prefers the new composition unless it values output in the targeted measure to the untargeted measure at a rate more than 56 times their predictiveness of graduation.

If we return to the graduation units, we can add the two measures up as "mean total output," which we show for the different policies in the last column of Table 9. We see that the behavioral response causes the tenured teachers' predicted total output to increase by  $0.055\sigma$ . We can also compare our results to the gains the district would achieve if it could (infeasibly) observe all forms of output and assign tenure according to the sum. We show the associated output of the tenured teachers under this policy in the second-to-last row of Table 9. We find, predictably, that total output is highest in this infeasible policy. Thus, the feasible policies do not achieve the first-best. If we define the "screening efficiency gap" as the difference in mean total output between an infeasible policy that screens on both dimensions and the ex post screening policy that screens on test score value-added without the behavioral response, we estimate that this screening efficiency gap is  $0.198\sigma$ . We find that the behavioral response to the policy closes 28% of the screening efficiency gap.

This improvement in the screening efficiency comes at the cost of distorting short-term output. A final cost-benefit analysis thus compares the endogenous response's short-run effects during the probationary period with the long-run effects on screening efficiency. Clearly, comparisons between the short-run and long-run depend on the district's discount rate and how many years tenured teaches are likely to work. Note that on the targeted measure there is large increase in short-term output and relatively small decrease in long-term output, while on the untargeted measures there is a large increase in the long-term output and a relatively small decrease in the short-term output. In other words, on both measures the benefits are larger than the costs, which suggests that the endogenous response improved overall efficiency. We can quantify this by noting that with no discounting, teachers would need to stay in the profession for 45 years for the (negative) long-run effects to

overtake the (positive) short-run effects on the targeted output, but the (positive) long-run effect overtakes the (negative) short-run effect after 3.5 years for the untargeted output.<sup>43</sup> Thus, provided tenured teachers remain in the district for between 3.5 and 45 years, the endogenous response leads to gains in both targeted and untargeted output.

## VIII Conclusion

It is a longstanding concern that evaluating individuals or institutions on a single measure will cause that measure to lose meaning, a worry social sciences often refer to as Campbell's Law.<sup>44</sup> Nowhere has this been more discussed than in the education setting and, in particular, using test scores to evaluate students, teachers, or schools. In developing Campbell's law, for example, Donald Campbell used test scores as an example, writing that: "Achievement tests may well be valuable indicators of general school achievement under conditions of normal teaching... but when test scores become the goal of the teaching process they lose their value of indicators of educational status."<sup>45</sup> In this paper, we illustrate both theoretically and empirically that the opposite can also be true – that evaluating individuals or institutions on a single measures may instead make the measure more informative, rather than less, of the individual or institution's underlying ability.

---

<sup>43</sup>We assume that the short-run effects last the 3 years of the probationary period.

<sup>44</sup>A similar idea is referred to as Goodhart's Law, which states that "when a measure becomes a target, it ceases to be a good measure." Goodhart's Law, like the Lucas Critique, was initially referred to macroeconomic models and monetary policy, but have since been applied to a number of different contexts.

<sup>45</sup>Holmstrom and Milgrom (1991) also used test scores were also used as the motivating example of multitasking.

## References

- Abadie, Alberto, Alexis Diamond, and Jens Hainmueller**, “Synthetic Control Methods for Comparative Case Studies: Estimating the Effect of California’s Tobacco Control Program,” *Journal of the American Statistical Association*, 2010, *105* (490), 493–505.
- Abaluck, Jason, Mauricio Caceres Bravo, Peter Hull, and Amanda Starc**, “Mortality Effects and Choice across Private Health Insurance Plans,” *The Quarterly Journal of Economics*, 2021, *136* (3), 1557–1610.
- Abdulkadiroğlu, Atila, Parag A Pathak, Jonathan Schellenberg, and Christopher R Walters**, “Do Parents Value School Effectiveness?,” *American Economic Review*, 2020, *110* (5), 1502–39.
- Armstrong, Mark and Jean-Charles Rochet**, “Multi-Dimensional Screening:: A User’s Guide,” *European Economic Review*, 1999, *43* (4-6), 959–979.
- Baker, George P**, “Incentive Contracts and Performance Measurement,” *Journal of Political Economy*, 1992, *100* (3), 598–614.
- Barlevy, Gadi and Derek Neal**, “Allocating Effort and Talent in Professional Labor Markets,” *Journal of Labor Economics*, 2019, *37* (1), 187–246.
- Bénabou, Roland and Jean Tirole**, “Bonus Culture: Competitive Pay, Screening, and Multitasking,” *Journal of Political Economy*, 2016, *124* (2), 305–370.
- Bertrand, Marianne and Antoinette Schoar**, “Managing with Style: The Effect of Managers on Firm Policies,” *The Quarterly journal of economics*, 2003, *118* (4), 1169–1208.
- Björkegren, Daniel, Joshua E Blumenstock, and Samsun Knight**, “Manipulation-proof machine learning,” *arXiv preprint arXiv:2004.03865*, 2020.
- Bleiberg, Joshua, Eric Brunner, Erica Harbatkin, Matthew A. Kraft, and Matthew G. Springer**, “The Effect of Teacher Evaluation on Achievement and Attainment: Evidence from Statewide Reforms,” Technical Report, Working Paper 2021.
- Brown, Christina and Tahir Andrabi**, “Inducing Positive Sorting through Performance Pay: Experimental Evidence from Pakistani Schools,” *University of California at Berkeley Working Paper*, 2021.
- Carrell, Scott E and James E West**, “Does Professor Quality Matter? Evidence from Random Assignment of Students to Professors,” *Journal of Political Economy*, 2010, *118* (3), 409–432.



- Chan, David C, Matthew Gentzkow, and Chuan Yu**, “Selection with Variation in Diagnostic Skill: Evidence from Radiologists,” *The Quarterly Journal of Economics*, 2022, *137* (2), 729–783.
- Chandra, Amitabh, Amy Finkelstein, Adam Sacarny, and Chad Syverson**, “Health Care Exceptionalism? Performance and Allocation in the US Health Care Sector,” *American Economic Review*, 2016, *106* (8), 2110–44.
- Chen, Zhao and Sang-Ho Lee**, “Incentives in Academic Tenure under Asymmetric Information,” *Economic Modelling*, 2009, *26* (2), 300–308.
- Chetty, Raj, John N. Friedman, and Jonah E. Rockoff**, “Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-Added Estimates,” *American Economic Review*, 2014, *104* (9), 2593–2632.
- Cohodes, Sarah R.**, “Teaching to the Student: Charter School Effectiveness in Spite of Perverse Incentives,” *Education Finance and Policy*, 2016, *11* (1), 1–42.
- Corcoran, Sean, Jennifer L. Jennings, and Andrew A. Beveridge**, “Teacher Effectiveness on High- and Low-Stakes Tests,” 2013.
- Dee, Thomas and James Wyckoff**, “Incentives, Selection, and Teacher Performance: Evidence from IMPACT,” *Journal of Policy Analysis and Management*, Spring 2015, *34* (2), 267–297.
- Delgado, William**, “Heterogeneous Teacher Effects, Comparative Advantage, and Match Quality,” 2021.
- Demougin, Dominique and Aloysius Siow**, “Careers in Ongoing Hierarchies,” *The American Economic Review*, 1994, pp. 1261–1277.
- Deshpande, Manasi and Yue Li**, “Who Is Screened out? Application Costs and the Targeting of Disability Programs,” *American Economic Journal: Economic Policy*, 2019, *11* (4), 213–48.
- Frankel, Alex and Navin Kartik**, “Muddled Information,” *Journal of Political Economy*, 2019, *127* (4), 1739–1776.
- Fryer, Roland G.**, “Teacher Incentives and Student Achievement: Evidence from New York City Public Schools,” *Journal of Labor Economics*, 2013, *31* (2), 373–427.
- Fryer, Roland G and Richard T Holden**, “Multitasking, Learning, and Incentives: A Cautionary Tale,” 2012.

- Gershenson, Seth**, “Linking Teacher Quality, Student Attendance, and Student Achievement,” *Education Finance and Policy*, 2016.
- Gilraine, Michael and Nolan G. Pope**, “Making Teaching Last: Long- and Short-Run Value-Added,” *Working Paper*, 2020.
- Glewwe, Paul, Nauman Ilias, and Michael Kremer**, “Teacher Incentives,” *American Economic Journal: Applied Economics*, 2010, 2 (3), 205–27.
- Hanushek, Eric A.**, “Teacher Deselection,” *Creating a New Teaching Profession*, 2009, 168, 172–173.
- Holmström, Bengt**, “Managerial Incentive Problems: A Dynamic Perspective,” *The Review of Economic Studies*, 1999, 66 (1), 169–182.
- Holmstrom, Bengt and Paul Milgrom**, “Multitask Principal-Agent Analyses: Incentive Contracts, Asset Ownership and Job Design,” *Journal of Law, Economics and Organization*, 1991, 7, 24–52.
- Hong, Fuhai, Tanjim Hossain, John A List, and Migiwa Tanaka**, “Testing the Theory of Multitasking: Evidence from a Natural Field Experiment in Chinese Factories,” *International Economic Review*, 2018, 59 (2), 511–536.
- Idoux, Clemence**, “Integrating New York City Schools: The Role of Admission Criteria and Family Preferences,” 2021.
- Jackson, C. Kirabo**, “What Do Test Scores Miss? The Importance of Teacher Effects on Non-Test Score Outcomes,” *Journal of Political Economy*, 2018, 126 (5), 2072–2107.
- Jacob, Brian A. and Lars Lefgren**, “Can Principals Identify Effective Teachers? Evidence on Subjective Performance Evaluations in Education.,” *Journal of Labor Economics*, 2008, 26 (1), 101–136.
- Kahn, Charles and Gur Huberman**, “Two-Sided Uncertainty and” Up-or-Out” Contracts,” *Journal of Labor Economics*, 1988, 6 (4), 423–444.
- Kane, Thomas J. and Douglas O. Staiger**, “Estimating Teacher Impacts on Student Achievement: An Experimental Evaluation,” *NBER*, 2008, (14607).
- , **Daniel F. McCaffrey, Trey Miller, and Douglas O. Staiger**, *Have We Identified Effective Teachers? Validating Measures of Effective Teaching Using Random Assignment*, Seattle, WA: Bill and Melinda Gates Foundation, 2013.

- Kou, Zonglai and Min Zhou**, “Multi-tasking vs. Screening: A Model of Academic Tenure,” *CCES, Fudan University Working Paper*, 2009.
- Kraft, Matthew A**, “Teacher Effects on Complex Cognitive Skills and Social-Emotional Competencies,” *Journal of Human Resources*, 2019, *54* (1), 1–36.
- , **Eric J Brunner, Shaun M Dougherty, and David J Schwegman**, “Teacher Accountability Reforms and the Supply and Quality of New Teachers,” *Journal of Public Economics*, 2020, *188*, 104212.
- Liu, Jing and Susanna Loeb**, “Engaging Teachers Measuring the Impact of Teachers on Student Attendance in Secondary School,” *Journal of Human Resources*, 2021, *56* (2), 343–379.
- Loeb, Susanna, Luke C. Miller, and James Wyckoff**, “Performance Screens for School Improvement: The Case of Teacher Tenure Reform in New York City,” *Educational Researcher*, 2015, *44* (4).
- Macartney, Hugh, Robert McMillan, and Uros Petronijevic**, “Teacher Value-Added and Economic Agency,” Technical Report, National Bureau of Economic Research 2018.
- McGuinn, Patrick**, “Stimulating Reform: Race to the Top, Competitive Grants and the Obama Education Agenda,” *Educational Policy*, 2012, *26* (1), 136–159.
- Mulhern, Christine and Isaac Opper**, “Measuring and Summarizing the Multiple Dimensions of Teacher Effectiveness,” 2021.
- Muralidharan, Karthik and Venkatesh Sundararaman**, “Teacher Performance Pay: Experimental Evidence from India,” *Journal of Political Economy*, 2011, *119* (1), 39–77.
- Murphy, Joseph, Philip Hallinger, and Ronald H Heck**, “Leading via Teacher Evaluation: The Case of the Missing Clothes?,” *Educational Researcher*, 2013, *42* (6), 349–354.
- Nagler, Markus, Marc Piopiunik, and Martin R. West**, “Weak Markets, Strong Teachers: Recession at Career Start and Teacher Effectiveness,” *NBER Working Paper*, 2015.
- Neal, Derek**, “The Design of Performance Pay in Education,” in “Handbook of the Economics of Education,” Vol. 4, Elsevier, 2011, pp. 495–550.
- Nichols, Albert L and Richard J Zeckhauser**, “Targeting Transfers through Restrictions on Recipients,” *The American Economic Review*, 1982, *72* (2), 372–377.

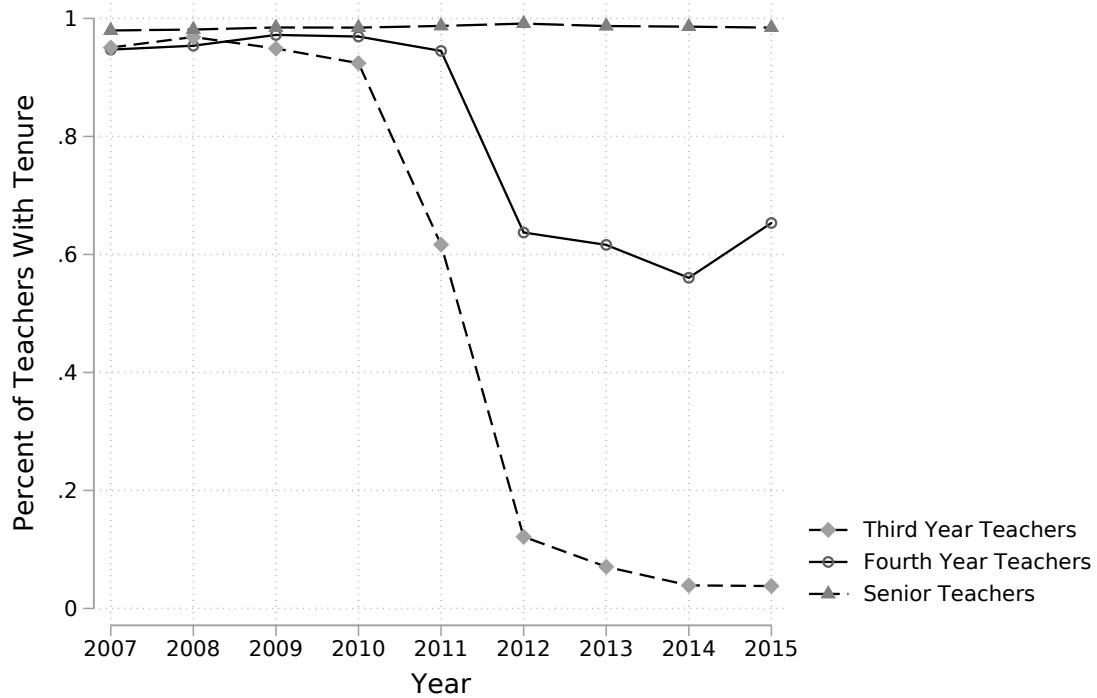
- O’Flaherty, Brendan and Aloysius Siow**, “On the Job Screening, Up or Out Rules, and Firm Growth,” *Canadian Journal of Economics*, 1992, pp. 346–368.
- Papay, John P.**, “Different Tests, Different Answers: The Stability of Teacher Value-Added Estimates Across Outcome Measures,” *American Education Research Journal*, 2011, *48*, 163–193.
- Papay, John P and Matthew A Kraft**, “Productivity Returns to Experience in the Teacher Labor Market: Methodological Challenges and New Evidence on Long-Term Career Improvement,” *Journal of Public Economics*, 2015, *130*, 105–119.
- Petek, Nathan and Nolan G. Pope**, “The Multidimensional Impact of Teachers on Students,” *Working Paper*, 2016.
- Prendergast, Canice**, “The Role of Promotion in Inducing Specific Human Capital Acquisition,” *The Quarterly Journal of Economics*, 1993, *108* (2), 523–534.
- Rebitzer, James B and Lowell J Taylor**, “When Knowledge Is an Asset: Explaining the Organizational Structure of Large Law Firms,” *Journal of Labor Economics*, 2007, *25* (2), 201–229.
- Rice, Jennifer King**, “Learning from Experience? Evidence on the Impact and Distribution of Teacher Experience and the Implications for Teacher Policy,” *Education Finance and Policy*, 2013, *8* (3), 332–348.
- Rockoff, Jonah and Lesley J Turner**, “Short-run Impacts of Accountability on School Quality,” *American Economic Journal: Economic Policy*, 2010, *2* (4), 119–47.
- Rockoff, Jonah E.**, “The Impact of Individual Teachers on Student Achievement: Evidence from Panel Data,” *The American Economic Review*, 2004, *94* (2), 247–252.
- Rockoff, Jonah E., Douglas O. Staiger, Thomas J. Kane, and Eric S. Taylor**, “Information and Employee Evaluation: Evidence from a Randomized Intervention in Public Schools,” *American Economic Review*, 2012, *102* (7), 3184–3213.
- Spence, Michael**, “Job Market Signaling,” *The Quarterly Journal of Economics*, 1973, *87* (3), 355–374.
- Staiger, Douglas O. and Jonah E. Rockoff**, “Searching for Effective Teachers with Imperfect Information,” *Journal of Economic Perspectives*, Summer 2010, *24* (3), 97–118.

- Sutcher, Leib, Linda Darling-Hammond, and Desiree Carver-Thomas**, “A Coming Crisis in Teaching? Teacher Supply, Demand, and Shortages in the U.S.,” Technical Report, Learning Policy Institute 2016.
- Tincani, Michela M**, “Teacher Labor Markets, School Vouchers, and Student Cognitive Achievement: Evidence from Chile,” *Quantitative Economics*, 2021, *12* (1), 173–216.
- Waldman, Michael**, “Up-or-Out Contracts: A Signaling Perspective,” *Journal of Labor Economics*, 1990, *8* (2), 230–250.
- Wiswall, Matthew**, “The Dynamics of Teacher Quality,” *Journal of Public Economics*, 2013, *100*, 61–78.

## IX Tables and Figures

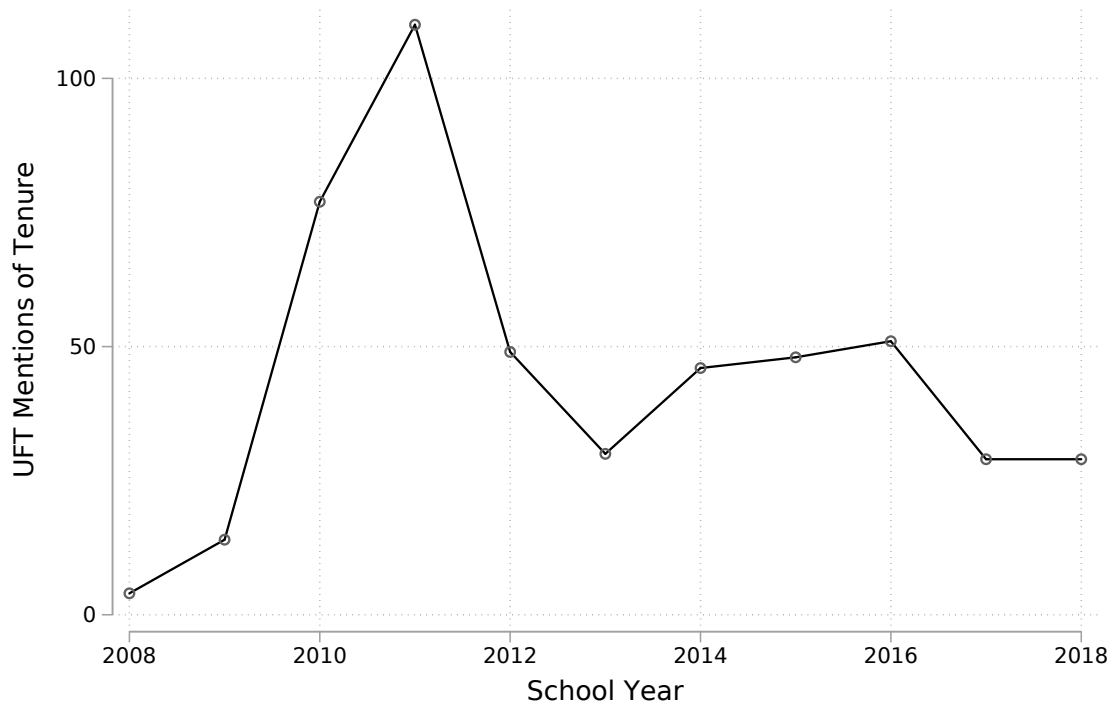
### IX.A Figures

Figure 1: Changes in Tenure Rates



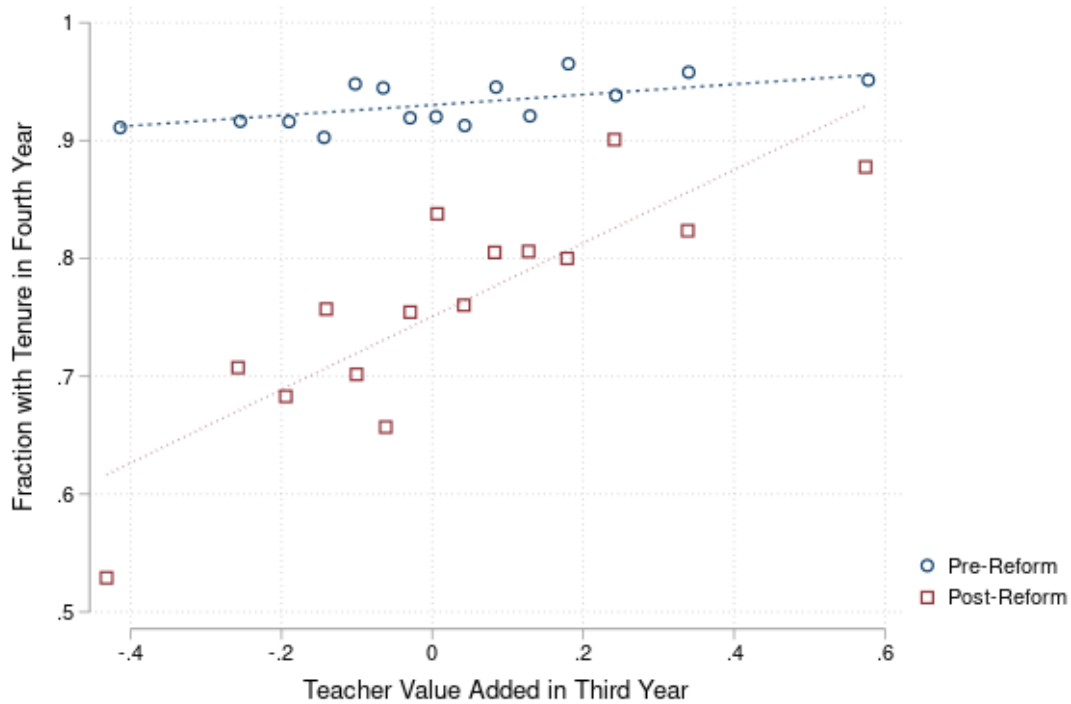
Note: The figure shows the fraction of teachers with tenure over time, split by years of experience. Teachers with fewer than three prior years of experience are typically not yet tenure eligible while teachers in their fourth year or later are tenure eligible. “Senior Teachers” have six or more prior years of experience.

Figure 2: Union Website Mentions of Tenure



Note: The figure shows the number of mentions of “tenure” on the teachers union (UFT) website, over time. The policy was announced in the 2010 school year.

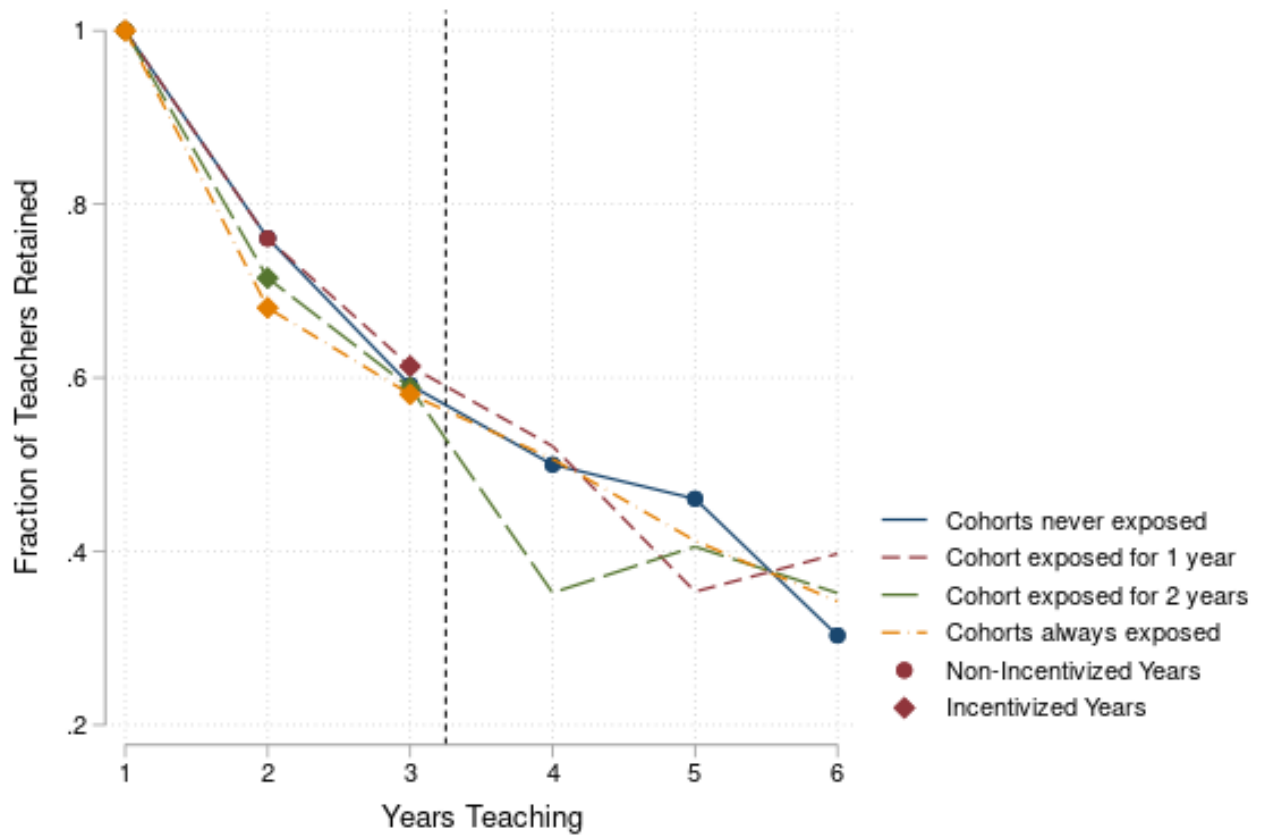
Figure 3: Tenure Probability by Value-Added



Note: This figure is a binscatter that groups teachers into twenty ventiles according to their (unshunken) test score value-added scores during their third year of experience. The y-axis is the fraction of teachers in each bin who have received tenure by the end of that year (entering their fourth year). We plot the relationships separately for the periods before and after the changes in the tenure process.

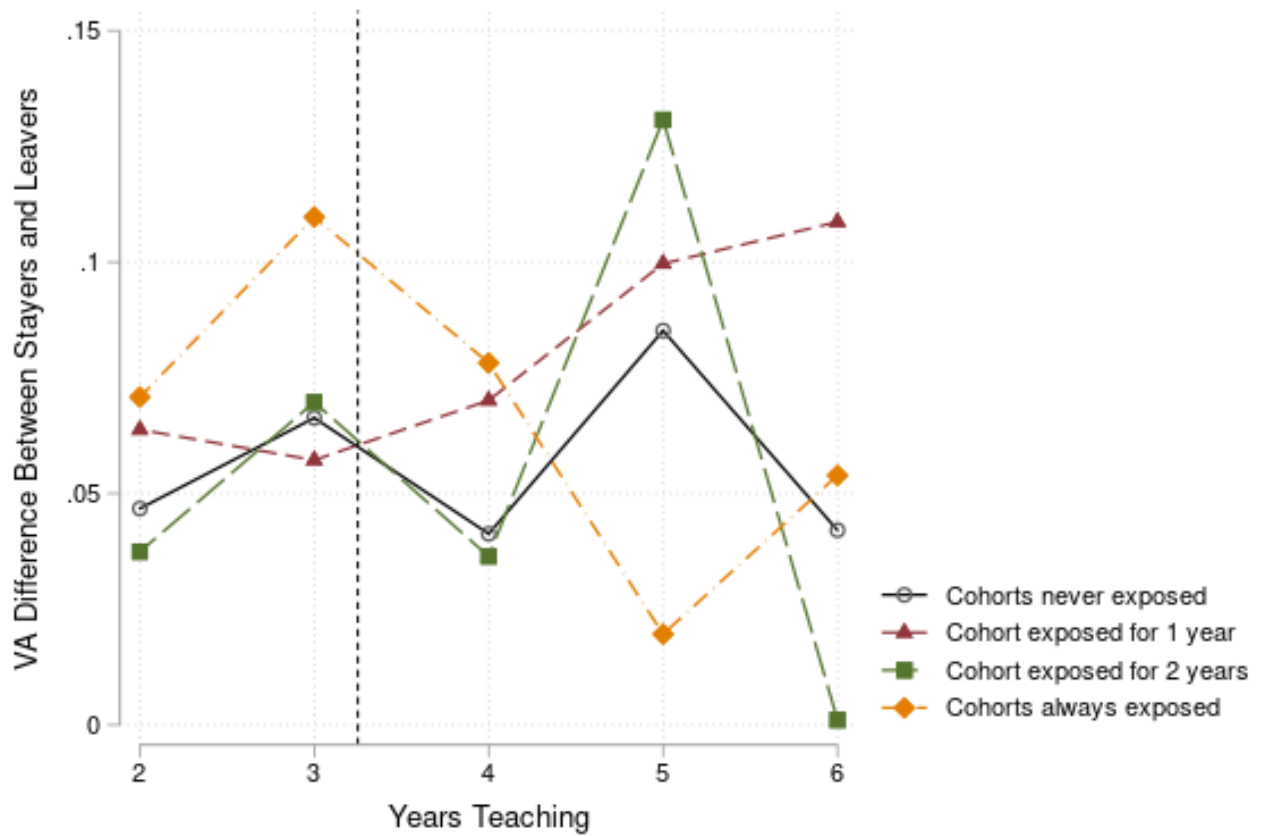


Figure 4: Survival Curves



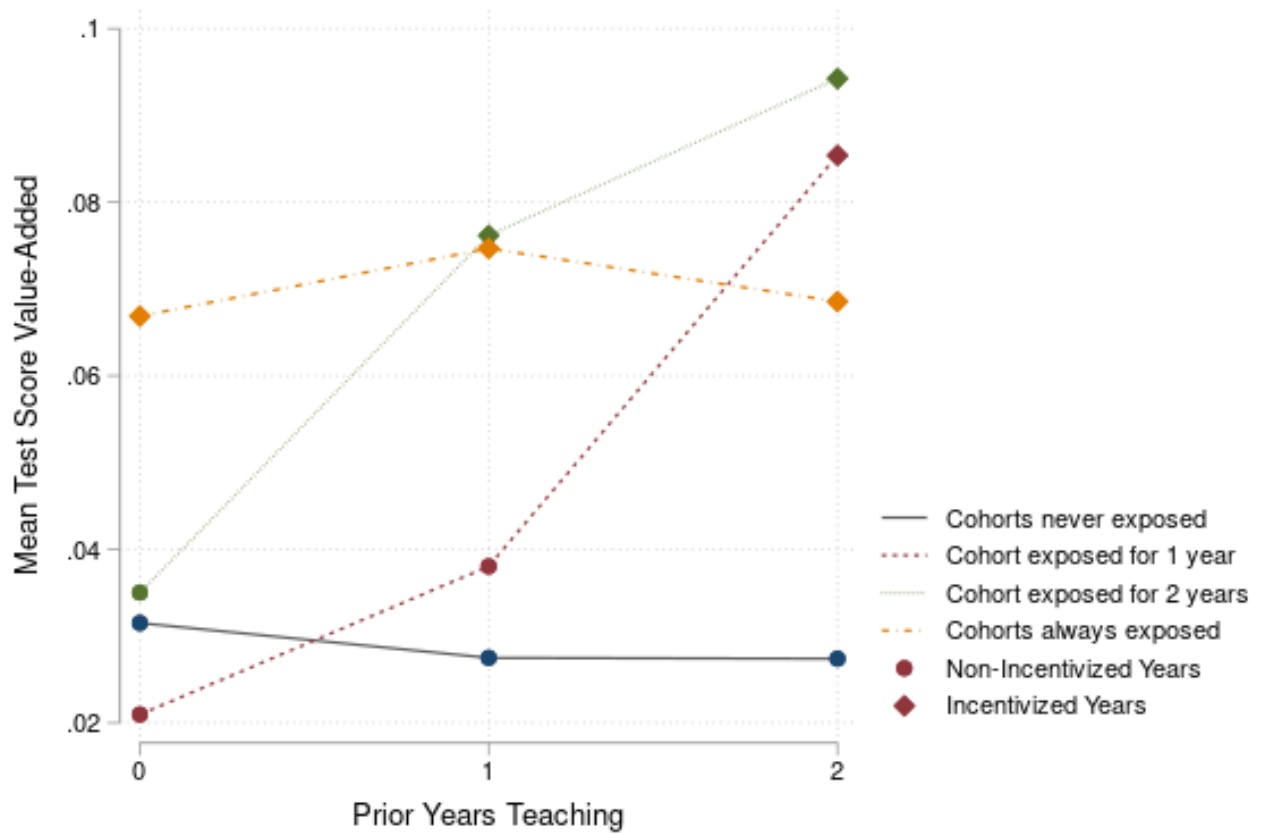
Note: This figure shows the fraction of teachers who stay teaching in the district, by years of experience. For instance, about 60% of teachers reach at least 3 years of experience. We plot separate lines for teacher cohorts based on the number of years they were exposed to the new tenure policy. Diamonds designate years in the probationary period under the new tenure policy.

Figure 5: Attrition Related to Test Score Value-Added



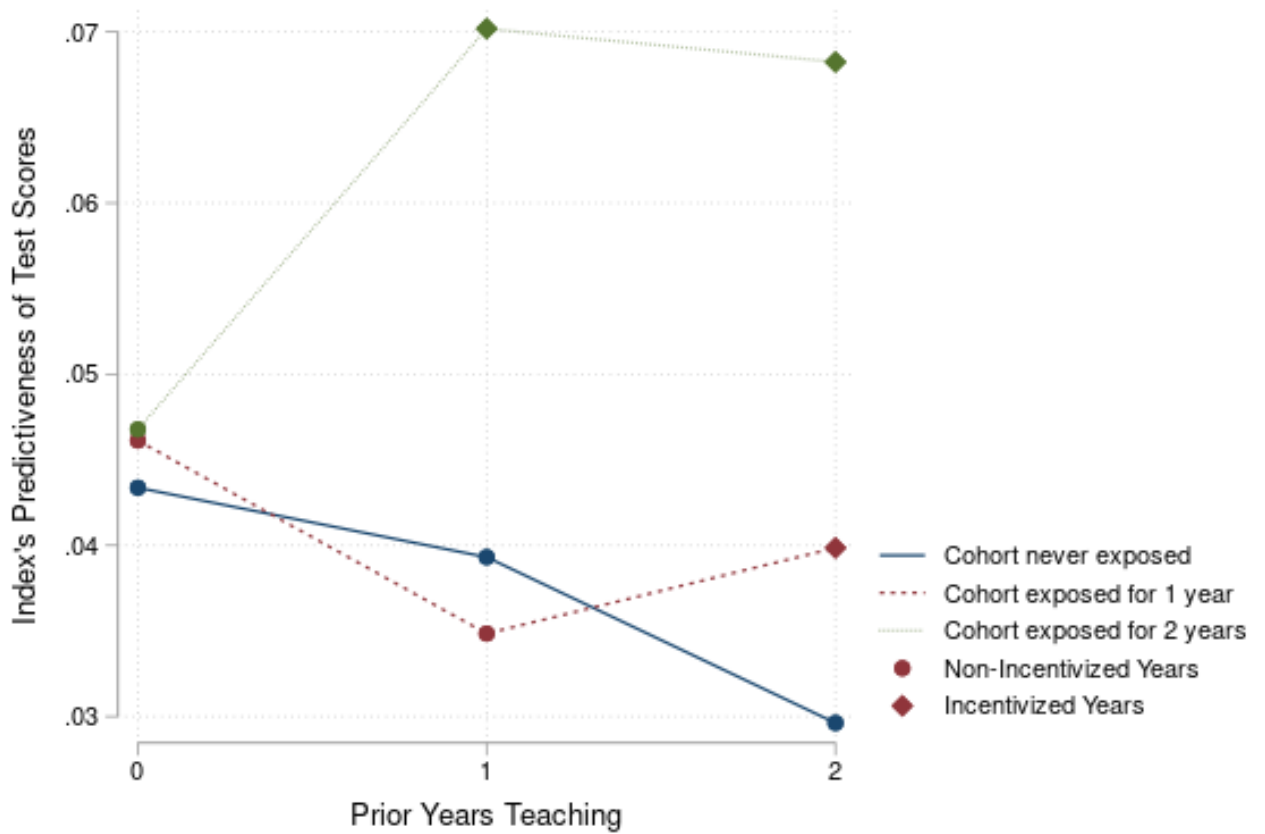
Note: This figure shows composition differences between teachers who stay for the following year and teachers who leave the district, split by years of experience. Each point represents the mean test score value-added difference between stayers and leavers at a specific experience level. The vertical line shows the end of the standard probationary period. We plot separate lines for teacher cohorts based on the number of years they were exposed to the new tenure policy.

Figure 6: Change in Test Score Value-Added



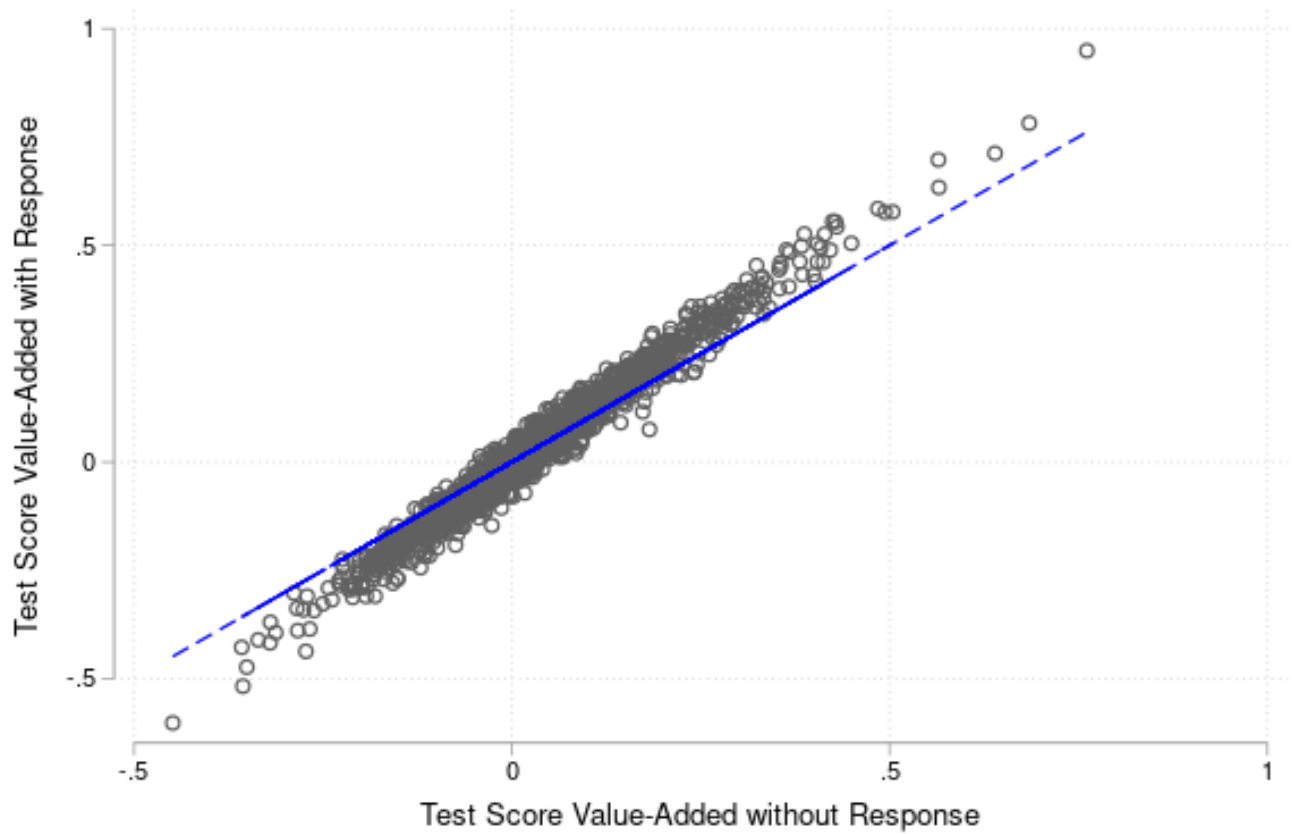
Note: This figure plots mean test score value-added for each cohort of teachers in the standard probationary period and at each level of prior experience. We plot separate lines for teacher cohorts based on the number of years they were exposed to the new tenure policy. Diamonds designate years in the probationary period under the new tenure policy.

Figure 7: Heterogeneous Response



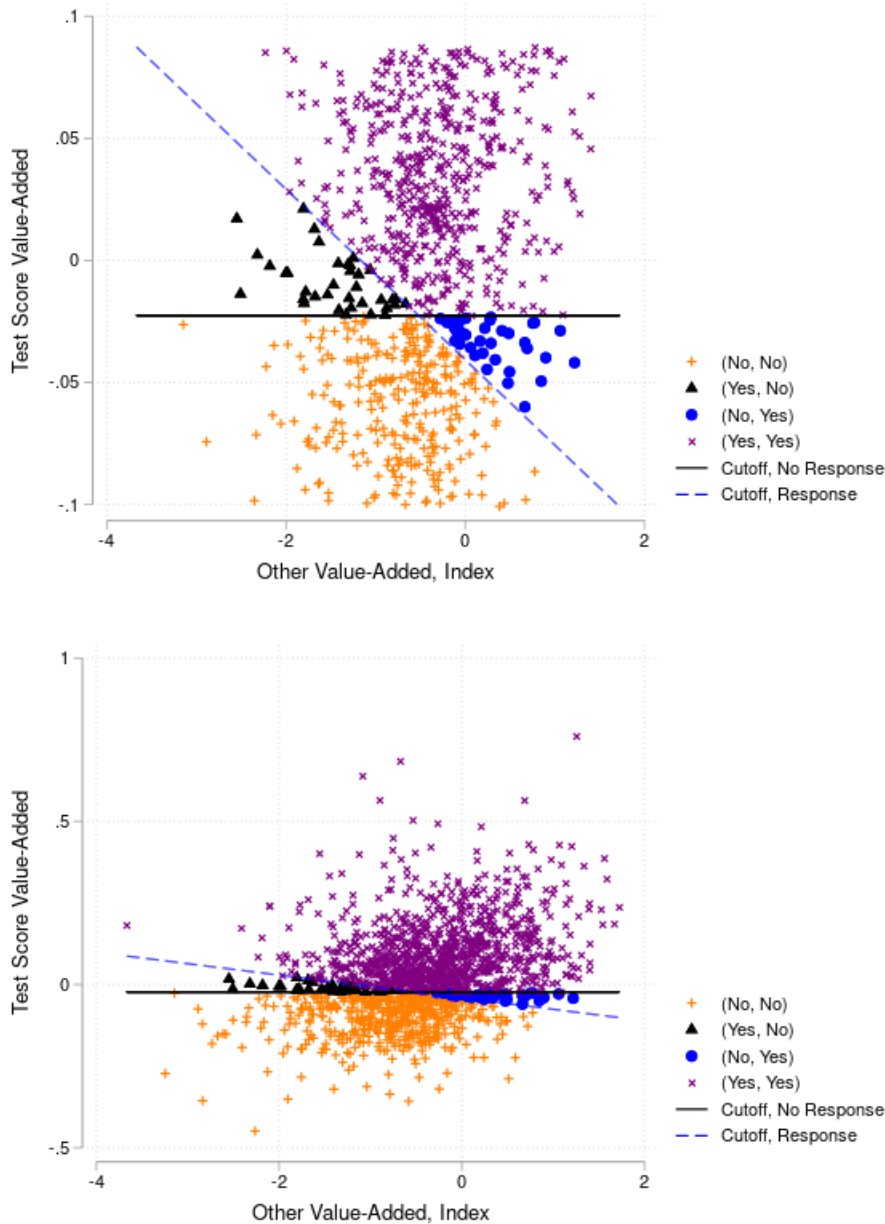
Note: This figure plots the predictiveness of the untargeted index forecast on test scores. We estimate this predictiveness by running a cross-sectional regression of a teacher's mean test score residuals on her forecasted value-added in the index of untargeted measures. We run a separate regression for each year and plot the coefficient on the forecasted value-added in the index. We plot separate lines for teacher cohorts based on the number of years they were exposed to the new tenure policy. Diamonds designate years in the probationary period under the new tenure policy. The "always exposed" cohorts are not included because this figure relies on forecasted test score value-added, as estimated in unincentivized periods. Because they are "always exposed," these cohorts do not have data to form such a forecast.

Figure 8: Response in Test Score Value Added



Note: This figure shows the joint distribution of teachers' test score value-added with and without the behavioral response to the policy. We construct the response based on our estimates. The blue line is the 45-degree line.

Figure 9: Tenure Predictions, with and without Behavioral Responses



Note: This figure shows the joint distribution of teacher forecasted value-added in the untargeted measure index (x-axis) and in test scores (y-axis). Each teacher in the 2007 cohort is plotted based on her forecasted value-added in unincentivized periods. The different symbols and colors reflect whether, according to our estimates, the teacher would receive tenure if (a) no teachers' behavior responded to the policy change and (b) if all teachers' behavior responded. For instance, "(No, Yes)" designates teachers who would only receive tenure when there are behavioral responses. The lines show the tenure cutoffs without behavioral responses (solid black) and with behavioral responses (dashed blue). The top panel zooms in on the test score value-added range between  $-0.1\sigma$  and  $0.1\sigma$  while the bottom panel zooms out and shows all teachers.

## IX.B Tables

Table 1: Summary Statistics

	Obs.	Mean	Std. Dev.	Min	Max
<b><i>Student Data</i></b>					
Male	4,602,185	0.49	0.50	0.00	1.00
Asian	4,602,185	0.17	0.38	0.00	1.00
Black	4,602,185	0.27	0.45	0.00	1.00
Hispanic	4,602,185	0.39	0.49	0.00	1.00
White	4,602,185	0.16	0.36	0.00	1.00
High-Poverty	4,602,185	0.80	0.40	0.00	1.00
English Language Learner	4,602,185	0.11	0.31	0.00	1.00
Middle Schooler	4,602,185	0.57	0.49	0.00	1.00
Math Score	2,325,544	0.00	1.00	-7.29	4.13
ELA Score	2,276,641	-0.00	1.00	-12.59	8.46
Attendance Rate	4,595,853	0.94	0.06	0.00	1.00
Grade in Math	1,032,550	80.36	11.99	10.00	100.00
Grade in ELA	815,248	79.19	11.30	10.00	100.00
Grade in Untested Subjects	2,573,954	81.81	9.35	10.00	100.00
<b><i>Teacher Data</i></b>					
Years Teaching at Current School	97,687	6.12	5.59	0.00	49.08
Years Teaching in District	97,687	8.34	6.66	0.00	49.17
In Probationary Period	98,520	0.21	0.41	0.00	1.00
<b><i>Counts</i></b>					
Number of Students	899,291				
Number of Student-Years	2,478,028				
Number of Student-Year-Subjects	4,602,185				
Number of Teachers	28,946				
Number of Teacher-Years	98,520				
Number of Teacher-Year-Subjects	145,021				

This table shows summary statistics for our student and teacher estimation samples. “High-Poverty” indicates eligibility for free or reduced price lunch. Math and ELA scores are normalized to have mean 0 and standard deviation 1. The probationary period is the pre-tenure period.

Table 2: Relationship between Tenure Decisions and Output

	On-Time Tenure	On-Time Tenure
Targeted Output	-0.00974 (0.0130)	0.260*** (0.0582)
Untargeted Output	-0.000118 (0.00152)	0.0127 (0.0130)
Constant	0.973*** (0.00356)	0.658*** (0.0177)
Sample	2008-2009	2011-2012
Mean DV	0.972	0.670
N	2137	700

This table shows the relationship between teacher output and whether she receives tenure by the beginning of her fourth year (on-time). Targeted output is the (unshrunk) mean test score residual. Untargeted output is the (unshrunk) index of other measures. An observation is a teacher in her fourth year of experience. The columns cover samples before and after the reform, respectively.

Table 3: Summary of Empirical Strategy

Entry Cohort	2007	2008	2009	2010	2011	2012
2007	0	0	0	-	-	-
2008		0	0	1	-	-
2009			0	1	1	-
2010				1	1	1

This table summarizes the empirical strategy by showing how we use within-cohort variation. Rows are the year the teacher entered the district and the columns are (spring) academic years. A blank entry means the teacher is not yet in the district. A dash means that the teacher is out of the standard probationary period. A blue zero indicates the teacher is in the probationary period, but *without* test score value-added incentives. A red one indicates the teacher is in the probationary period, *with* test score value-added incentives.



Table 4: Effect of Policy Change on Probationary Period Output

	Test Score	Test Score	Untargeted Index	Untargeted Index
Incentive	0.0330*** (0.00718)	0.0149* (0.00791)	-0.0562* (0.0295)	-0.0644** (0.0307)
Fixed Effects	Cohort	Teacher	Cohort	Teacher
Implied Grad (Uni)	.0028	.0013	-.0023	-.0026
Implied Grad (Multi)	.0008	.0004	-.0022	-.0025
N Teachers	27296	22488	19225	16724
Mean DV	0.0910	0.0952	0.0259	0.0244
N	132486	127678	102906	100405

This table shows the causal effect of the tenure policy change on targeted and untargeted output in the probationary period. The columns switch between cohort and teacher fixed effects. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. The teachers with more than 3 years of experience are only included if they finished the standard probationary period before the tenure policy change. “Implied Grad” is the implied change in graduation rate from our anchoring regressions, where “(Uni)” comes from a univariate regression with just the targeted (or untargeted) output and “(Multi)” comes from a multivariate regression with both sets of measures. All outcome units are test score student standard deviations.

Table 5: Effect of Policy Change on Specific Untargeted Outcomes

	Score 1	Score 2	Attend	Attend 1	Grades	Other Grades	Grades 1	Other Grades 1
Incentive	-0.0263** (0.0115)	-0.00648 (0.0143)	0.00282 (0.00245)	-0.0187* (0.0113)	0.0372** (0.0179)	-0.00310 (0.0169)	-0.0630*** (0.0203)	-0.0405** (0.0197)
Teacher FEs	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Mean DV	0.0482	0.0433	0.00386	0.0277	0.0246	-0.00840	0.0362	0.0380
N	105861	84359	127444	106376	32984	43418	56362	70194

This table shows the causal effect of the tenure policy change on individual untargeted (residualized) outcomes in the probationary period. The “1” or “2” in the column headers indicate the measure’s number of years into the future. “Grades” are in the tested subject while “Other” grades are in untested subjects. All variables are standardized at the grade-year level to have mean 0 and standard deviation 1 in the full population. All regressions include teacher fixed effects. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. The teachers with more than 3 years of experience are only included if they finished the standard probationary period before the tenure policy change.

Table 6: Effect of Policy Change on Predetermined Student Characteristics

	Lag Score	Lag Other Score	Male	Black	Hispanic	Asian	ELL	Poverty
Incentive	0.000390 (0.0141)	-0.00311 (0.0142)	0.00121 (0.00407)	0.00180 (0.00467)	0.00111 (0.00498)	-0.000643 (0.00345)	-0.0000865 (0.00506)	-0.00301 (0.00671)
Teacher FEs	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Mean DV	-0.0175	-0.0144	0.488	0.290	0.397	0.158	0.119	0.813
N	127678	127678	127678	127678	127678	127678	127678	127678

This table shows the causal effect of the tenure policy change on predetermined student characteristics in the probationary period. “Score” is the test score in the same subject while “Other Score” is in the opposite subject. “ELL” are English language learners. “Poverty” indicates students eligible for free or reduced price lunch. All regressions include teacher fixed effects. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. The teachers with more than 3 years of experience are only included if they finished the standard probationary period before the tenure policy change.

Table 7: Effect of Policy Change on Post-Probationary Period Output

	Test Score	Untargeted Index
Incentive	0.0185*** (0.00656)	-0.0421* (0.0223)
Post-Incentive	-0.0114 (0.00960)	-0.0366 (0.0439)
Fixed Effects	Teacher	Teacher
N Teachers	1252	17126
Mean DV	0.0972	0.0286
N	135361	107261

This table shows the causal effect of the tenure policy change on targeted and untargeted output. We estimate separate coefficients for the standard probationary period (first 3 years) and for after the teacher has tenure (“Post”). Each of these indicators is interacted with whether the new tenure policy is in place. All regressions include teacher fixed effects. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. All outcome units are test score student standard deviations.

Table 8: Heterogeneous Responses to the Incentive

	Score	Score	Score	Score	Score	Score
Incentive	0.0142** (0.00591)	0.00125 (0.00844)	0.0241*** (0.00749)	0.0185*** (0.00636)	0.0148** (0.00591)	0.0122 (0.00844)
Incentive * Targeted VA	0.0703 (0.0567)	0.217* (0.115)			0.0546 (0.0565)	0.164 (0.102)
Incentive * Untargeted VA			0.651* (0.324)	1.128** (0.484)	0.0158* (0.00888)	0.0407** (0.0172)
FEs	Cohort	Teacher	Cohort	Teacher	Cohort	Teacher
SD VA1	.142	.142			.142	.142
SD VA2			.482	.482	.482	.482
Mean DV	0.00594	0.00594	0.00594	0.00594	0.00594	0.00594
N	84399	84399	84399	84399	84399	84399

This table shows the causal effect of the tenure policy change on targeted output, split by a teacher’s forecasted targeted and untargeted value-added. The forecasts are based on a value-added model estimated with data from the unincentivized period. Regressions include either cohort or teacher fixed effects. Each value-added measure is interacted with experience indicators. “SD VA1” and “SD VA2” are the cross-sectional standard deviations of targeted and untargeted forecasted value-added. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. The teachers with more than 3 years of experience are only included if they finished the standard probationary period before the tenure policy change. Outcome units are test score student standard deviations.

Table 9: Output under Different Screening Regimes

	Obs.	Mean Targeted Output	Mean Untargeted Output	Mean Total Output
<i>Teachers' Tenure under Different Responses</i>				
Never Tenured	510	-0.104	-0.786	-0.890
Only Tenured w/o Behavioral Response	36	-0.009	-1.444	-1.453
Only Tenured w/ Behavioral Response	36	-0.033	0.286	0.253
Always Tenured	1,077	0.100	-0.237	-0.136
<i>Tenured Teachers under Different Policies</i>				
Screening w/o Behavioral Response	1,113	0.097	-0.276	-0.179
Screening w/ Behavioral Response	1,113	0.096	-0.220	-0.124
(Infeasible) Screening on Both Dimensions	1,113	0.064	-0.046	0.019
<i>Gains Relative to Infeasible First-Best</i>				
Gains (Fraction)				0.279

This table shows the mean output for different groups of teachers and under different policies. The sample is teachers who started in the district in 2007. The top panel splits teachers into four groups based on whether they would receive tenure in a regime without a behavioral response and whether they would receive tenure in a regime with a behavioral response. The middle panel shows the mean output associated with the set of teachers receiving tenure under different policies. The first two policies are screening on the targeted measure, without and with a behavioral response. The last policy is an infeasible policy screening on the sum of output across both dimensions. The final panel shows the fraction of gains the behavioral response achieves relative to the distance between the screening without behavioral response regime and the infeasible policy screening on the sum of output. “Mean Targeted” output is the mean forecasted test score value-added. “Mean Untargeted Output” is the mean forecasted value-added on the untargeted index. “Mean Total Output” is the sum of the mean forecasted value-added across the targeted and untargeted measures. All outcome units are test score student standard deviations.

## A Proofs

### A.1 Main Proofs

**Theorem 1.** *For any weakly increasing screening function, we get that for every  $\theta$ :*

$$\begin{aligned} x_1^*(\theta|p) &\geq x_1^*(\theta) \\ x_2^*(\theta|p) &\leq x_2^*(\theta). \end{aligned}$$

*Proof.* As we do throughout this paper, we define  $\tilde{u}(x_1, \theta)$  as the individual  $\theta$ 's utility when producing  $x_1$  and then optimizing over  $x_2$ . Formally, we get that:

$$\tilde{u}(x_1, \theta) \equiv \begin{cases} \max_{e \in \mathcal{E}} u(e, \theta) & \text{if } \exists e_1 \in \mathcal{E} \text{ s.t. } e_1 \theta_1 \geq x_1 \\ -\infty & \text{if } \forall e_1 \in \mathcal{E} \quad e_1 \theta_1 < x_1 \end{cases} \quad (25)$$

With this definition and that  $x_1^*(\theta)$  is an optimum, we get that  $\tilde{u}(x_1^*(\theta), \theta) \geq \tilde{u}(x_1, \theta)$  for all  $x_1 < x_1^*(\theta)$ . It is then also clear that  $\tilde{u}(x_1^*(\theta), \theta) + \lambda p(x_1^*(\theta)) \Delta v(\theta) \geq \tilde{u}(x_1, \theta) + \lambda p(x_1) \Delta v(\theta)$  for all  $x_1 < x_1^*(\theta)$  and so  $x_1^*(\theta|p) \geq x_1^*(\theta)$ .

Next, we define  $x_2^*(x_1, \theta)$  to be  $\theta$ 's optimal choice  $x_2$  when producing  $x_1$ . Importantly, this function is the same regardless of whether the utility includes the term  $\lambda p(x_1) \Delta v(\theta)$  or not, since that expression does not depend on  $x_2$ . In addition, our assumptions on the cross derivatives and the assumption that  $\mathcal{E}$  is convex gives us that  $x_2^*(x_1, \theta)$  is decreasing in  $x_1$ . Since  $x_1^*(\theta|p) \geq x_1^*(\theta)$  it therefore follows that  $x_2^*(\theta|p) \leq x_2^*(\theta)$ . □

**Lemma 1.** *Consider two individuals  $\theta$  and  $\theta'$  with  $x_1^*(\theta) = x_1^*(\theta')$  and  $x_2^*(\theta) < x_2^*(\theta')$  and define:*

$$e^*(\theta|f) = \arg \max_{e \in \mathcal{E}} u(e, \theta) + \lambda f(x_1) \quad (4)$$

$$x_k^*(\theta|f) = e_k^*(\theta|f) \cdot \theta_k \quad (5)$$

*Then  $x_1^*(\theta|f) \leq x_1^*(\theta'|f)$  for any weakly increasing function  $f(x_1)$  if either of the following are true:*

- $\lambda$  is sufficiently large;
- $\tilde{u}'(x_1, \theta') - \tilde{u}'(x_1, \theta)$  is increasing in  $x_1$

*Proof.* To show that the theorem holds for sufficiently large  $\lambda$ , let  $e_1^{max}$  be the maximum possible effort level on task one, while still satisfying the constraint that  $e \in \mathcal{E}$ . Then

there exists  $\bar{\lambda}$  such that  $\forall \lambda \geq \bar{\lambda}$ , we have  $\tilde{u}(\theta e_1^{max}, \theta') + \lambda f(\theta e_1^{max}) > \tilde{u}(x_1, \theta') + \lambda f(x_1)$  for every  $x_1$  such that  $f(x_1) < f(\theta e_1^{max})$ . Note that this stems from the fact that  $\theta'_1 > \theta_1$  and so it will require less effort (on task one) to produce  $x_1 = \theta e_1^{max}$  for individual  $\theta'$  than individual  $\theta$ . This, along with our assumption that  $b(x)$  is bounded and that  $c(e) < \infty$  for every  $e \in \mathcal{E}$  is what allows us to conclude that the inequality holds. Of course, the inequality implies that individual  $\theta'$  will prefer  $\theta e_1^{max}$  to any  $x_1$  with a lower value of  $f$  and we also know that  $x_1^*(\theta|f) \leq \theta e_1^{max}$  from the constraint set  $\mathcal{E}$ . That result is sufficient if  $f$  is strictly increasing, but since it is only weakly increasing we need to add a few more technical details. Specifically, let  $\tilde{x}_1 = \min\{x_1 | f(x_1) = f(\theta e_1^{max})\}$ , which is well-defined if  $f$  is right-continuous. From the concavity of  $u$  and the reasons discussed above, we can then conclude that  $x_1^*(\theta|f) \leq \max\{\tilde{x}_1, x_1^*(\theta)\} \leq x_1^*(\theta'|f)$ .

See Appendix B for a proof that the theorem holds if  $\tilde{u}(x_1, \theta') - \tilde{u}(x_1, \theta)$  is increasing in  $x_1$  and a more detailed discussion of that condition.  $\square$

**Theorem 2.** *Consider  $\theta < \theta'$  with  $x_1(\theta) = x_1(\theta')$ . Assume that  $\Delta v(\theta)$  is increasing in  $\theta$  and that either  $\tilde{u}'(x_1, \theta') - \tilde{u}'(x_1, \theta)$  is increasing in  $x_1$  or  $\lambda$  is sufficiently large. Then  $x_1^*(\theta|p) \leq x_1^*(\theta'|p)$  for any weakly increasing screening function  $p(x_1)$ .*

*Proof.* Define  $\tilde{x}_1^*(\theta'|p) = \arg \max u(x, \theta') + \lambda p(x_1) \cdot \Delta v(\theta)$ , i.e., the optimal choice of individual  $\theta'$  if her value of staying in the profession relative to the outside option were  $\Delta v(\theta)$  instead of  $\Delta v(\theta')$ . From the fact that  $\Delta v(\theta)$  is increasing in  $\theta$  and  $\theta < \theta'$ , it follows that  $x_1^*(\theta'|p) \geq \tilde{x}_1^*(\theta'|p)$ . But from the previous theorem and the assumption that  $\tilde{u}'(x_1, \theta') - \tilde{u}'(x_1, \theta)$  is increasing in  $x_1$ , we get that  $\tilde{x}_1^*(\theta'|p) \geq x_1^*(\theta|p)$ . Thus,  $x_1^*(\theta|p) \leq x_1^*(\theta'|p)$ .  $\square$

**Theorem 3.** *Assume the conditions on  $V(\theta)$  specified above and that the assumptions in Theorem 2 hold. Furthermore, assume that  $\theta$  is continuously distributed. Then for any ex post screening policy  $p(x_1)$ , there is a screening policy  $\tilde{p}(x_1)$  that is more efficient than  $p(x_1)$ :*

$$\begin{aligned} \mathbb{E}_\Theta \left[ \tilde{p}(x_1^*(\theta|\tilde{p})) \right] &= \mathbb{E}_\Theta \left[ p(x_1^*(\theta)) \right] \\ \mathbb{E}_\Theta \left[ V(\theta) \cdot \tilde{p}(x_1^*(\theta|\tilde{p})) \right] &\geq \mathbb{E}_\Theta \left[ V(\theta) \cdot p(x_1^*(\theta)) \right]. \end{aligned}$$

*Proof.* Start with any ex post screening policy. Since  $\mathbb{E}[V(\theta)|x_1(\theta)]$  is increasing in  $x_1(\theta)$ , it follows that an ex post screening policy with a threshold screening function, i.e.,  $p(x_1) = \mathbf{1}(x_1 \geq \bar{x}_1)$  for some  $\bar{x}_1$ , is more efficient at screening than the initial ex post screening policy. Next, choose  $\tilde{x}_1$  such that the announced screening policy with a threshold screening function  $\tilde{p}(x_1) = \mathbf{1}(x_1 \geq \tilde{x}_1)$  retains the same fraction of teachers, i.e.  $\mathbb{E}_\Theta \left[ \tilde{p}(x_1^*(\theta|\tilde{p})) \right] = \mathbb{E}_\Theta \left[ p(x_1^*(\theta)) \right]$ . We show in Lemma A.4 that under the assumption that

$x_1^*(\theta)$  is continuously distributed such an  $\tilde{x}$  exists. Using Lemma A.2, we can conclude that  $\mathbb{E}_\Theta \left[ V(\theta) \cdot \tilde{p}(x_1^*(\theta|\tilde{p})) \right] \geq \mathbb{E}_\Theta \left[ V(\theta) \cdot p(x_1^*(\theta)) \right]$ , which by definition means that the announced screening policy is more efficient at screening than the ex post screening policy. Since we did not specify the initial ex post screening policy the result that there is always an announced screening policy that is more efficient at screening is true for all ex post screening policy policies.  $\square$

## A.2 Supporting Lemmas

**Lemma A.1.** *Suppose that  $x_1^*(\theta) = x_1^*(\theta')$  and  $x_2^*(\theta') > x_2^*(\theta)$ . Then  $\theta'_k > \theta_k$  for  $k \in \{1, 2\}$ .*

*Proof.* First, we note that  $x_1^*(\theta)$  is increasing in  $\theta_1$  and decreasing in  $\theta_2$ . Thus, if  $\theta'_1 \geq \theta_1$  and  $\theta'_2 \leq \theta_2$ , then by our assumptions we would have that  $x_1^*(\theta') > x_1^*(\theta)$ . (This assumes that the comparative statics are strict.). We can similarly rule out the fact that  $\theta'_1 \leq \theta_1$  and  $\theta'_2 \geq \theta_2$ . We can also rule out the fact that  $\theta'_2 < \theta_2$ . (Otherwise, it would be then be cheaper for  $\theta$  than  $\theta'$  to move from  $x_2^*(\theta)$  to  $x_2^*(\theta')$ , which would contradict the assumption of optimization.) This rules out any other case than  $\theta'_k > \theta_k$  for  $k \in \{1, 2\}$ .  $\square$

**Lemma A.2.** *Assume the conditions in Theorem 2 and the assumptions on  $V(\theta)$  outlined in the paper. Then consider a ex post screening policy with a threshold screening function  $p(x_1) = \mathbf{1}(x_1 \geq \bar{x}_1)$  and an announced screening policy with a threshold screening function  $p(x_1) = \mathbf{1}(x_1 \geq \tilde{x}_1)$ . Define:*

$$A = \{\theta | x_1(\theta) < \bar{x}_1 \ \& \ x_1(\theta|p) \geq \tilde{x}_1\}$$

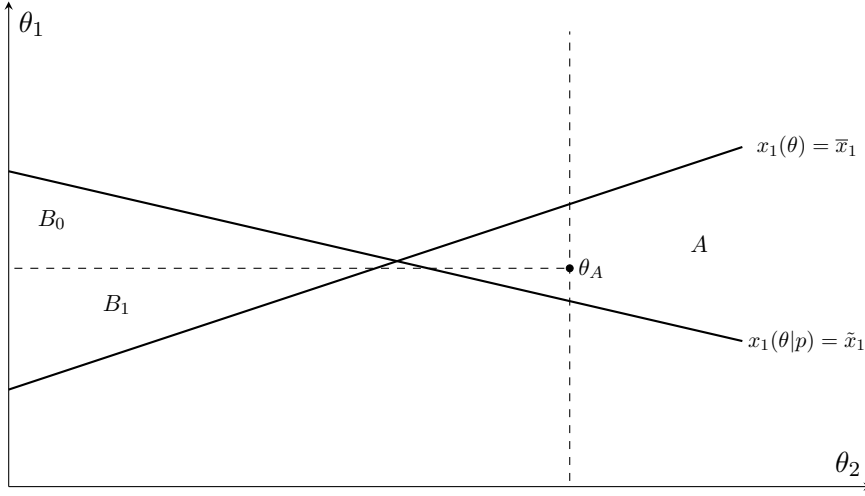
$$B = \{\theta | x_1(\theta) \geq \bar{x}_1 \ \& \ x_1(\theta|p) < \tilde{x}_1\}$$

*Then  $\forall \theta^A \in A, \theta^B \in B$ , we have that  $V(\theta^A) \geq V(\theta^B)$ .*

*Proof Sketch.* To illustrate the intuition, we plot the parameter space  $(\theta_1, \theta_2)$  below. On this space we draw a line through all of the points in which  $x_1(\theta) = \bar{x}_1$ . From our assumptions, every  $\theta$  point below this line has  $x_1(\theta) < \bar{x}_1$  and so is not retained under the ex post screening policy, while every  $\theta$  above the line has  $x_1(\theta) \geq \bar{x}_1$  and so is retained. While we could similarly label the line  $x_1(\theta|p) \geq \tilde{x}_1$  in the announced policy, there will likely be bunching there and so the line should more precisely be defined as partitioning the space into the  $\theta$ 's that are retained under the announced screening policy and those that are not.

From the results of Theorem 2, we get that the slope of the line where  $x_1(\theta|p) = \tilde{x}_1$  is less positive than the line where  $x_1(\theta) = \bar{x}_1$ , as they are illustrated below. From this, we can conclude that  $\forall \theta^A \in A, \theta^B \in B$ , we have that  $\theta_2^A > \theta_2^B$ . If  $\theta_B \in B_1$  in the figure below, i.e., that  $\theta_1^A > \theta_1^B$ , then  $\theta_A > \theta_B$  and so it follows directly that  $V(\theta_A) > V(\theta_B)$ . The challenging case is therefore when  $\theta_1^B > \theta_1^A$ , i.e., when  $\theta_B \in B_0$ .

Under the assumption that  $\theta_1^B > \theta_1^A$  and  $\theta_2^A > \theta_2^B$ , however, it follows that moving from  $x_1(\theta)$  to  $\tilde{x}_1$  is “cheaper” for  $\theta_B$  than for  $\theta_A$  in the sense that  $\tilde{u}(\tilde{x}_1, \theta_B) - u(x_1(\theta_B), \theta_B) \geq \tilde{u}(\tilde{x}_1, \theta_A) - u(x_1(\theta_A), \theta_A)$ . This, along with the fact that  $\theta_A$  does increase to  $\tilde{x}_1$  under the announced screening policy but  $\theta_B$  does not, allows us to infer that  $\Delta v(\theta_A) > \Delta v(\theta_B)$ . Combining this result with the fact that  $\theta_1^B > \theta_1^A$  and  $\theta_2^A > \theta_2^B$  and the assumptions regarding  $V$ , we get that  $V(\theta_A) > V(\theta_B)$ .



□

*Proof.* Throughout, we will assume that  $A, B \neq \emptyset$ , since otherwise the statement is vacuous. We fix  $\theta^A \in A$  and  $\theta^B \in B$ . We start by ruling out the possibility that  $\theta^A \leq \theta^B$ , i.e., that  $\theta_k^A \leq \theta_k^B$  for  $k \in \{1, 2\}$ . If it were, then there would exist a  $\tilde{\theta}$  such that  $\tilde{\theta}_2 = \theta_2^B$ ,  $\tilde{\theta}_1 < \theta_1^B$ , and  $x_1(\theta^A) = x_1(\tilde{\theta})$  (see Lemma A.3 for the proof). From Lemma 1, we then get that  $x_1(\tilde{\theta}|p) \geq x_1(\theta^A|p)$ . However, since  $\tilde{\theta}_2 = \theta_2^B$  and  $\tilde{\theta}_1 < \theta_1^B$ , we also get that  $x_1(\theta^B|p) \geq x_1(\tilde{\theta}|p)$  since  $x_1(\theta|p)$  is increasing in  $\theta_1$ . Therefore,  $x_1(\theta^B|p) \geq x_1(\theta^A|p)$ , which is a contradiction to the fact that  $\theta^A \in A$  and  $\theta^B \in B$ .

Similarly, it also cannot be the case that  $\theta_2^B \geq \theta_2^A$  and  $\theta_1^B \leq \theta_1^A$ . This follows from the fact that  $x_1(\theta)$  is increasing in  $\theta_1$  and decreasing in  $\theta_2$  and  $x_1(\theta^B) > x_1(\theta^A)$ . Thus, we can infer that  $\theta_2^B \leq \theta_2^A$ .

If  $\theta_1^B \leq \theta_1^A$  then  $\theta^A \geq \theta^B$  and so by assumption  $V(\theta^A) \geq V(\theta^B)$ . Thus, in what follows we will assume that  $\theta_1^B \geq \theta_1^A$  and  $\theta_2^B \leq \theta_2^A$  and show that in this case  $V(\theta^A) \geq V(\theta^B)$ . To do so, we define  $\tilde{u}(x_1, \theta)$  as in the paper, i.e.,

$$\tilde{u}(x_1, \theta) \equiv \begin{cases} \max_{e \in \mathcal{E}} u(e, \theta) & \text{if } \exists e_1 \in \mathcal{E} \text{ s.t. } e_1 \theta_1 \geq x_1 \\ -\infty & \text{if } \forall e_1 \in \mathcal{E} \quad e_1 \theta_1 < x_1 \end{cases} \quad (26)$$

Since  $\theta^A \in A$  and  $\theta^B \in B$ , we know that  $\theta^A$  finds it worth producing  $\tilde{x}_1$  in the pre-tenure period under the announced screening policy, while  $\theta^B$  does not. From this, we can conclude that:

$$v(\theta^A) \leq \tilde{u}(\tilde{x}_1, \theta^A) + \lambda \cdot \Delta v(\theta^A) \quad (27)$$

$$v(\theta^B) \geq \tilde{u}(\tilde{x}_1, \theta^B) + \lambda \cdot \Delta v(\theta^B) \quad (28)$$

Rearranging, we get that:

$$\lambda \cdot [\Delta v(\theta^A) - \Delta v(\theta^B)] \geq [(\tilde{u}(\tilde{x}_1, \theta^B) - v(\theta^B)) - (\tilde{u}(\tilde{x}_1, \theta^A) - v(\theta^A))] \quad (29)$$

Finally, from the envelope theorem we get that that  $\tilde{u}(x, \theta) - v(\theta)$  is increasing in  $\theta_1$  and decreasing in  $\theta_2$  for every  $x$ . Thus,  $[(\tilde{u}(\tilde{x}_1, \theta^B) - v(\theta^B)) - (\tilde{u}(\tilde{x}_1, \theta^A) - v(\theta^A))] \geq 0$  and so  $\Delta v(\theta^A) \geq \Delta v(\theta^B)$ . This, combined with the assumptions about  $V$  and the fact that  $\theta_2^B \leq \theta_2^A$  and  $\theta_1^B \geq \theta_1^A$  implies that  $V(\theta^A) \geq V(\theta^B)$ .  $\square$

**Lemma A.3.** *Suppose that  $\theta^A \in A$  and  $\theta^B \in B$  and  $\theta^A \leq \theta^B$ . Then exists a  $\tilde{\theta}$  such that  $\tilde{\theta}_2 = \theta_2^B$ ,  $\tilde{\theta}_1 \leq \theta_1^B$ , and  $x_1^*(\theta^A) = x_1^*(\tilde{\theta})$ .*

*Proof.* Since  $x_1^*(\theta)$  is decreasing in  $\theta_2$  and by assumption  $x_1^*(\theta^A) < x_1^*(\theta^B)$  and  $\theta_2^B \leq \theta_2^A$ , it follows that  $x_1^*(\theta_1^A, \theta_2^B) < x_1^*(\theta^A)$ . From Berge's maximization theorem,  $x_1^*(\theta)$  is upper hemicontinuous and so from the intermediate value theorem as we increase  $\theta_1$  from  $(\theta_1^A, \theta_2^B)$  to  $(\theta_1^B, \theta_2^B)$  there must be a  $\tilde{\theta}_1 \in (\theta_1^A, \theta_1^B)$  such that  $x_1^*(\theta^A) = x_1^*(\tilde{\theta}_1, \theta_2^B)$ .  $\square$

**Lemma A.4.** *Suppose that  $\theta$  is continuously distributed. Then for every  $p \in (0, 1)$  there exists a  $\tilde{x}$  such that under the policy  $\tilde{p}(x_1) = \mathbf{1}(x_1 \geq \tilde{x}_1)$ , we get that  $\mathbb{E}_\Theta [\tilde{p}(x_1^*(\theta|\tilde{p}))] = p$ .*

*Proof.* Define  $g(\tilde{x}_1) \equiv \mathbb{E}_\Theta [\tilde{p}(x_1^*(\theta|\tilde{p}))]$ . If  $g$  is a continuous function of  $\tilde{x}_1$ , the intermediate value theorem implies the result.

To show that  $g$  is continuous, we consider a sequence that converges to  $\tilde{x}_1$ , i.e.,  $(\tilde{x}_{1,n}) \rightarrow \tilde{x}_1$ . If  $g(\tilde{x}_{1,n}) \rightarrow g(\tilde{x}_1)$ , then it follows that  $g$  is continuous.

By assumption  $\theta$  is continuously distributed, and we will denote its probability density function as  $f(\theta)$ . Then we can write  $g(\tilde{x}_1)$  as:

$$g(\tilde{x}_1) = \int \mathbf{1}(\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) \geq 0) f(\theta) d\theta. \quad (30)$$

Define  $h(\theta) \equiv \mathbf{1}(\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) \geq 0)$  and  $h_n(\theta) \equiv \mathbf{1}(\tilde{u}(\tilde{x}_{1,n}, \theta) + \lambda \Delta v(\theta) - v(\theta) \geq 0)$ . Also define  $e_{1,max}$  as  $\sup\{e_1 | e \in \mathcal{E}\}$ . From Lemma A.5,  $h_n(\theta) \rightarrow h(\theta)$  at every  $\theta$  such that  $\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) \neq 0$  and  $\theta_1 e_{1,max} \neq \tilde{x}_1$ . Furthermore, from



our assumption that  $V(\theta)$  is increasing in  $\theta_1$  and that  $\frac{\partial V(\theta)}{\partial \theta_1} \leq \frac{\partial \Delta v(\theta)}{\partial \theta_1}$  and the envelope condition, we get that  $\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta)$  is strictly increasing in  $\theta_1$  at every point where  $\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) = 0$ . This means there is at most one  $\theta_1$  for every  $\theta_2$  such that  $\tilde{u}(\tilde{x}_{1,n}, \theta) + \lambda \Delta v(\theta) - v(\theta) = 0$ . Under the assumption that  $\theta$  is continuously distributed, this implies that the set  $\{\theta | \tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) = 0 \cup \theta_1 e_{1,max} = \tilde{x}_1\}$  has zero measure. Together with the result from Lemma A.5, this implies that  $h_n \rightarrow h$  pointwise almost everywhere. From the bounded convergence theorem, it follows that  $g_n \rightarrow g$  and so  $g$  is continuous.  $\square$

**Lemma A.5.** *Let  $e_{1,max} = \sup\{e_1 | e \in \mathcal{E}\}$ . Then  $h_n(\theta) \rightarrow h(\theta)$  at every  $\theta$  such that  $\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) \neq 0$  and  $\tilde{x}_1 \neq \theta_1 e_{1,max}$ .*

*Proof.* Consider some  $\theta$  such that  $\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) \equiv \epsilon \neq 0$ . We first consider the case in which  $\tilde{x}_1 < e_{1,max} \theta_1$ . We then know that  $\tilde{u}(\tilde{x}_1, \theta)$  is continuous in  $\tilde{x}_1$  and so there exists a  $\delta > 0$  such that  $|\tilde{u}(\tilde{x}_1, \theta) - \tilde{u}(x_1, \theta)| < \frac{|\epsilon|}{2}$  for every  $x_1$  such that  $|x_1 - \tilde{x}_1| < \delta$ . This implies that  $\mathbf{1}(\tilde{u}(x_1, \theta) + \lambda \Delta v(\theta) - v(\theta) \geq 0) = \mathbf{1}(\tilde{u}(\tilde{x}_1, \theta) + \lambda \Delta v(\theta) - v(\theta) \geq 0)$  for every  $x_1$  such that  $|x_1 - \tilde{x}_1| < \delta$ . It thus follows that there exists an  $N$  such that  $h_n(\theta) = h(\theta)$  for all  $n > N$ .<sup>46</sup>

If  $\tilde{x}_1 > e_{1,max} \theta_1$ , then  $\tilde{u}(\tilde{x}_1, \theta) = -\infty$ . Again, there exists an  $N$  such that  $h_n(\theta) = h(\theta)$  for all  $n > N$  since for all  $n > N$  we get that  $\tilde{x}_{1,n} > e_{1,max} \theta_1$ .  $\square$

## B Model Appendix

### B.1 The Inefficiency of Screening on Output

Much of the intuition for our theoretical results can be seen in a comparison between two infeasible screening regimes. One – which we refer to as “ex post screening” – screens on the teachers’ output without the additional incentive of the screening policy in place, i.e., on  $x_1^*(\theta)$  and the other – which we refer to as “ability screening” – screens on the first dimension of ability, i.e., on  $\theta_1$ . Since the principals need not worry about incentives in either regime, it is easy to see that if  $u$  is differentiable and the constraint that  $e \in \mathcal{E}$  does not bind, then the two regimes would be equivalent if the principal observed both dimensions of ability and both dimensions of output. There is an underlying inefficiency with screening on output, however, that appears when only one of the two outputs is observed.

This inefficiency can be summarized with the following lemma.

---

<sup>46</sup>This is because  $\tilde{x}_{1,n} \rightarrow \tilde{x}_1$  means that there exists some  $N$  such that  $|\tilde{x}_{1,n} - \tilde{x}_1| < \delta$  for all  $n > N$ .

**Lemma B.1.** *In any ex post screening policy, the probability that individual  $\theta$  is retained is weakly decreasing in  $\theta_2$  for every  $\theta_1$ .*

*Proof.* This follows directly from the definition of an ex post screening policy and the fact that  $x_1^*(\theta)$  is decreasing in  $\theta_2$ .  $\square$

In some sense, we can think of this lemma as being roughly akin to a screening version of the traditional multitasking problem. The traditional multitasking problem highlights an inherent inefficiency when the principal can only add an incentive to a single output: it reduces effort related to the production of the other outcomes. The screening version outlined in the lemma above instead highlights an inherent inefficiency when the principal can only screen on a single output: an individual who is effective at increasing the other output is less likely to be retained than one who is ineffective at increasing the other output and equally effective at increasing the output screened on. This implies the following theorem:

**Thm B.1.** *The ability screening regime is more efficient than the ex post screening regime.*

*Proof.* Consider any ex post screening policy  $p(x_1)$  and define an ability screening regime as:

$$\tilde{p}(\theta_1) = \mathbb{E}[p(x_1^*(\theta))|\theta_1] \quad (31)$$

Because  $p$  and  $\tilde{p}$  retain the same fraction of teachers of each  $\theta_1$ , they retain the same fraction of all teachers. Furthermore, since  $p(x_1^*(\theta))$  is decreasing in  $\theta_2$ , for any increasing function  $V(\theta)$ , we have that:  $\mathbb{E}[V(\theta)\tilde{p}(\theta_1)|\theta_1] \geq \mathbb{E}[V(\theta)p(x_1^*(\theta))|\theta_1]$  for every  $\theta_1$ . Since that inequality holds for every  $\theta_1$ , we get that:  $\mathbb{E}[V(\theta)\tilde{p}(\theta_1)] \geq \mathbb{E}[V(\theta)p(x_1^*(\theta))]$ .  $\square$

Again, the intuition is straightforward. Principals would like to keep teachers with both high  $\theta_1$  and  $\theta_2$  and Lemma B.1 shows that an ex post screening policy that retains individuals with high  $x_1^*(\theta)$  biases one toward removing teachers with high  $\theta_2$ . While not actively retaining individuals with high  $\theta_2$  the way an optimal policy would, screening on ability rather than output does at least remove this bias and is thus more efficient.

Finally, note that theorem does not say that ex post screening on  $x_1^*(\theta)$  would necessarily reduce the average of  $x_2^*(\theta)$  in the population. If  $\theta_1$  and  $\theta_2$  are highly correlated, ex post screening on  $x_1^*(\theta)$  would improve the average  $x_2^*(\theta)$  in the population; however, Theorem B.1 says that  $x_2^*(\theta)$  could be increased even more if the policy maker were able to screen directly on  $\theta_1$  rather than on  $x_1^*(\theta)$ .

## B.2 Different Responses to the Same Incentive

Suppose there is a manager who cares about two dimensions of employee output, but who can only observe one of those dimensions. In this principal-agent problem, often referred

to as the multitasking problem (Holmstrom and Milgrom, 1991), the manager may shy away from implementing as large an incentive on the output she observes if doing so causes employees to shift away from time spent producing other dimension of output. But suppose that she also must evaluate the employees based on the single observed dimension at the end of the period; does that change her decision about the optimal size of incentive? Stated differently, does the substitution pattern induced by the added incentive cause the first dimension to be a more or less noisy signal of the employees' underlying ability?

As we discuss in the paper, the answer depends on how different individuals respond differently to the incentive. Here, we expand on the results in the paper to show that a very similar condition on the utility function discussed in Lemma 1 is not only sufficient, but – in a limited sense – necessary to ensure that higher ability individuals respond more to the incentive than lower ability individuals who produce the same  $x_1$  without the incentive. We then discuss in a bit more detail the economic meaning behind the condition and provide a few examples of specific utility functions that do (and do not) meet the criteria.

### B.2.1 Formal Model and Definitions:

We use the same model of employee output as described in Section II.B and consider the case where each employee gets paid – or has to pay – an additional  $f(x_1)$  when producing  $x_1$ . For notational simplicity, we drop the  $\lambda$  term in front of  $f$  in this section; doing so in Section II.B allowed us to formally consider the case where the incentive got large, but that is not the focus here and so we will stop the over-parameterization of the  $f$  function. If  $f(x_1)$  is increasing, then this serves as an additional positive incentive; if it is decreasing it serves as an additional negative incentive. A key point here is that all individuals receive the same incentive and value it equivalently in utility terms.

As in the paper, it helps to define  $\tilde{u}(x_1, \theta)$  as the optimal utility individual  $\theta$  can get when constrained to produce at least  $x_1$ , i.e.,

$$\tilde{u}(x_1, \theta) \equiv \begin{cases} \max_{e \in \mathcal{E}} u(e, \theta) & \text{if } \exists e_1 \in \mathcal{E} \text{ s.t. } e_1 \theta_1 \geq x_1 \\ -\infty & \text{if } \forall e_1 \in \mathcal{E} \quad e_1 \theta_1 < x_1 \end{cases} \quad (32)$$

We will use  $\tilde{u}'(x_1, \theta)$  to denote  $\frac{\partial \tilde{u}(x_1, \theta)}{\partial x_1}$  when  $\tilde{u}$  is differentiable and while differentiability is not an necessary assumption in our results, it makes the notation and intuition more straightforward. Throughout this discussion, we will also focus on cases where  $f$  is small enough that we do not have to worry about  $\tilde{u}(x_1, \theta)$  equalling  $-\infty$  for either of the two individuals we consider, which simplifies the exposition.

We next turn to conditions on  $\tilde{u}$  that allow us to make conclusions about which individuals respond more or less to the added incentive. The condition, which is both sufficient

and – in a limited sense we discuss below – necessary, is a single crossing condition on the derivative of  $\tilde{u}$ , defined formally below.

**Definition 1.** We say that  $\tilde{u}(x_1, \theta)$  has the single crossing condition on  $\tilde{u}'$  if every  $\theta' > \theta$  with  $x_1^*(\theta) = x_1^*(\theta') \equiv x_1^*$ ,  $\tilde{u}(x_1, \theta') - \tilde{u}(x_1, \theta)$  is strictly decreasing for all  $x_1 < x_1^*$  and  $\tilde{u}(x_1, \theta') - \tilde{u}(x_1, \theta)$  is strictly increasing for all  $x_1 > x_1^*$ .

We call this a single crossing condition because – in the case where  $\tilde{u}$  is differentiable – it implies that  $\tilde{u}'(x_1, \theta)$  and  $\tilde{u}'(x_1, \theta')$  cross a single time, at  $x_1^*$ . We state this formally in the following remark:

**Remark 1.** When  $\tilde{u}(x_1, \theta)$  is differentiable, the two definitions are equivalent:

1.  $\tilde{u}(x_1, \theta)$  has the single crossing condition on  $\tilde{u}'$  as defined above;
2. For every  $\theta' > \theta$  with  $x_1^*(\theta) = x_1^*(\theta') \equiv x_1^*$ , we have that  $\tilde{u}'(x_1, \theta') > \tilde{u}'(x_1, \theta)$  for every  $x_1 > x_1^*$  and  $\tilde{u}'(x_1, \theta') < \tilde{u}'(x_1, \theta)$  for every  $x_1 < x_1^*$ .

We deliberately choose to call this a single crossing condition on  $\tilde{u}'$  to highlight the similarities and differences between our results and the traditional results on comparative statics. While traditional comparative statics aims at understanding how the optimal choice of  $x$  varies according to characteristics  $\theta$ , we consider a slightly different question and aim to understand how the *change* in the optimal choice of  $x$  vary according to characteristics  $\theta$  when the individuals are presented with an identical *change* in incentives. Interestingly, the result mirrors the result from comparative statics, although it now hinges on increasing differences in the *marginal*, i.e., change in, utility function rather than the utility function itself. Just as in traditional comparative statics, we can also relax our condition slightly from a single-crossing condition to an increasing differences condition, which is no longer necessary but is sufficient; again, this increasing differences condition is on the marginal utility rather than the utility function itself and can be stated as a condition on the second derivatives of  $\tilde{u}$ . We state this formally in the remark below.

**Remark 2.** If  $\tilde{u}''(x_1, \theta') > \tilde{u}''(x_1, \theta)$  for every  $\theta' > \theta$  with  $x_1^*(\theta) = x_1^*(\theta') \equiv x_1^*$ , then  $\tilde{u}'(x_1, \theta)$  has the single crossing condition on  $\tilde{u}'$  as defined above.

A key difference between our results and traditional comparative statics results is that ours are much more limited, in that we do not compare the change in optimal choices of any two individuals but only individuals who choose the same output absent the additional incentive. This also handles the challenges inherent to settings with multidimensional types. By narrowing our comparison to individuals who choose the same output absent the additional incentive, we essentially reduce the types to a single dimension. Because we condition on a choice ( $x_1(\theta)$ ) rather than a type, the analysis does not rely on standard single-dimensional screening results.

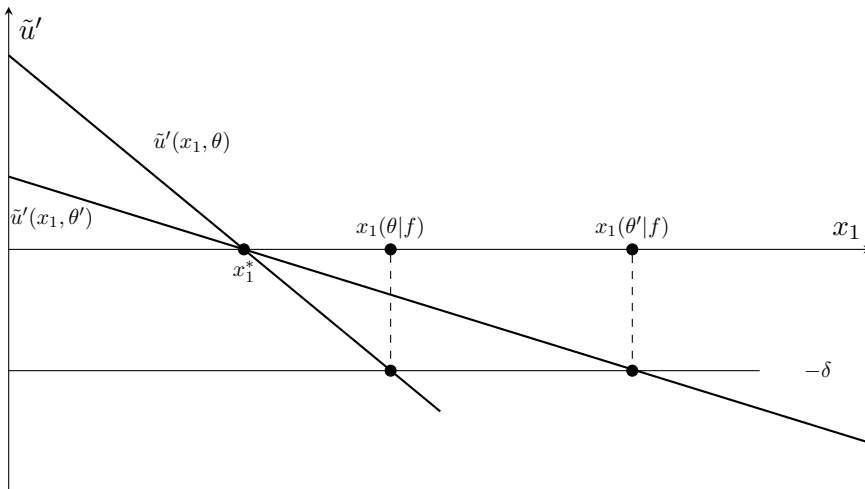
**B.2.2 Theory Results and Proofs:**

Given the definition of the single crossing condition on  $\tilde{u}'$  from above, we can finally turn to the results. First, we present the result that the single crossing condition is sufficient to ensure that  $\theta'$  responds more to the incentive than  $\theta$ .

**Thm B.2.** *Assume that  $\tilde{u}(x_1, \theta)$  has the single crossing condition on  $\tilde{u}'$ , and consider any  $\theta' > \theta$  with  $x_1^*(\theta) = x_1^*(\theta')$ . Then:*

- *For every weakly increasing function  $f(x_1)$ , we have  $x_1^*(\theta'|f) \geq x_1^*(\theta|f)$ . Similarly, for every weakly decreasing function  $x_1^*(\theta'|f) \leq x_1^*(\theta|f)$ .*
- *Further assume that  $\tilde{u}$  is differentiable. Then for every strictly increasing differentiable function  $f(x_1)$ , we have  $x_1^*(\theta'|f) > x_1^*(\theta|f)$ . Similarly, for every strictly decreasing continuous function  $x_1^*(\theta'|f) < x_1^*(\theta|f)$ .*

*Proof Sketch.* Suppose for simplicity that  $f(x_1) = \delta \cdot x_1$  for some  $\delta > 0$  and consider the figure below, where the curves  $\tilde{u}'(x_1, \theta')$  and  $\tilde{u}'(x_1, \theta)$  are both plotted on a graph with  $\tilde{u}'$  on the y-axis and  $x_1$  on the x-axis. From the assumptions on  $u$ , we get that both  $\tilde{u}'(x_1, \theta')$  and  $\tilde{u}'(x_1, \theta)$  are downward sloping and if we assume  $\tilde{u}$  is convex, the optimal choice of  $x_1^*(\theta|f)$  is the point where  $\tilde{u}'(x_1, \theta) = -\delta$ . Clearly  $x_1^*(\theta|f) > x_1^*$  and from the single crossing condition on  $\tilde{u}'$ , we get that  $\tilde{u}'(x_1, \theta') > \tilde{u}'(x_1, \theta)$  for all  $x_1 > x_1^*$ , so it follows that  $\tilde{u}'(x_1^*(\theta|f), \theta') > -\delta$ . Thus,  $x_1^*(\theta'|f) > x_1^*(\theta|f)$ .



□

*Proof.* Consider any weakly increasing function  $f(x_1)$ . Clearly,  $x_1^*(\theta|f) \geq x_1^*$  and  $x_1^*(\theta'|f) \geq x_1^*$ . If  $x_1^*(\theta|f) = x_1^*$ , then we are done, so we will assume in what follows that  $x_1^*(\theta|f) > x_1^*$ .

We then consider any  $x_1 \in (x_1^*, x_1^*(\theta|f))$ . Since  $x_1^*(\theta|f)$  is an optimizer, we get that  $\tilde{u}(x_1, \theta) + f(x_1) \leq \tilde{u}(x_1^*(\theta|f), \theta) + f(x_1^*(\theta|f))$  or that  $\tilde{u}(x_1^*(\theta|f), \theta) - \tilde{u}(x_1, \theta) \geq f(x_1) - f(x_1^*(\theta|f))$ . Further, since  $x_1^*(\theta|f) > x_1 \geq x_1^*$  we get from the single crossing condition on  $\tilde{u}'$  that:  $\tilde{u}(x_1^*(\theta|f), \theta') - \tilde{u}(x_1^*(\theta|f), \theta) > \tilde{u}(x_1, \theta') - \tilde{u}(x_1, \theta)$ . Rearranging, it follows that:

$$\tilde{u}(x_1^*(\theta|f), \theta') - \tilde{u}(x_1, \theta') > \tilde{u}(x_1^*(\theta|f), \theta) - \tilde{u}(x_1, \theta) \geq f(x_1) - f(x_1^*(\theta|f)).$$

Thus, individual  $\theta'$  would choose  $x_1^*(\theta|f)$  over  $x_1$  for all  $x_1 \in (x_1^*, x_1^*(\theta|f))$ , which along with the fact that  $x_1^*(\theta'|f) \geq x_1^*$ , proves that  $x_1^*(\theta'|f) \geq x_1^*(\theta|f)$ . The proof that  $x_1^*(\theta'|f) \leq x_1^*(\theta|f)$  for any weakly decreasing function  $f$  is identical.

To prove the second bullet point, that if  $\tilde{u}$  is differentiable and  $f(x_1)$  is a strictly increasing differentiable function we get that  $x_1^*(\theta'|f) > x_1^*(\theta|f)$ , we note that the single crossing condition implies:  $\tilde{u}'(x_1^*(\theta|f), \theta') > \tilde{u}'(x_1^*(\theta|f), \theta)$ . Furthermore, since  $x_1^*(\theta|f)$  is an optimum, we get that  $\tilde{u}'(x_1^*(\theta|f), \theta) = -f'(x_1^*(\theta|f))$  for interior  $x_1^*(\theta|f)$ . Together, this implies that  $\tilde{u}'(x_1^*(\theta|f), \theta') + f'(x_1^*(\theta|f)) > 0$ . Thus, for small enough  $\epsilon > 0$ , we get that  $\tilde{u}(x_1^*(\theta|f) + \epsilon, \theta') + f(x_1^*(\theta|f) + \epsilon) > \tilde{u}(x_1^*(\theta|f), \theta') + f(x_1^*(\theta|f))$  and so  $\theta'$  would choose  $x_1^*(\theta|f) + \epsilon$  over  $x_1^*(\theta|f)$ . Together with the previous result that  $x_1^*(\theta'|f) \geq x_1^*(\theta|f)$ , we conclude that  $x_1^*(\theta'|f) > x_1^*(\theta|f)$ . Again, the proof is identical to show that  $x_1^*(\theta'|f) < x_1^*(\theta|f)$  if  $f$  is strictly decreasing. □

The above results imply that the single-crossing condition on  $\tilde{u}'$  is sufficient, in that it ensures that  $\theta'$  responds more to the change in incentives than  $\theta$ . We next show that it is also a necessary condition for a general  $f$  function. Formally, we have the following theorem:

**Thm B.3.** *Suppose  $\tilde{u}(x_1, \theta)$  does not have the single crossing condition on  $\tilde{u}'$ . Then there exists a strictly increasing function such that  $x_1^*(\theta'|f) \leq x_1^*(\theta|f)$  or there exists a strictly decreasing function such that  $x_1^*(\theta'|f) \geq x_1^*(\theta|f)$ .*

*Proof.* We will initially assume that the failure of the single crossing condition on  $\tilde{u}'$  occurs by there being some  $\tilde{x}_1 > x_1^*$  such that  $\tilde{u}'(\tilde{x}_1, \theta') = \tilde{u}'(\tilde{x}_1, \theta)$  for some  $\theta' > \theta$  with  $x_1^*(\theta) = x_1^*(\theta') \equiv x_1^*$ . We will then show that there exists a strictly increasing differentiable  $f(x_1)$  such that  $x_1^*(\theta|f) = x_1^*(\theta'|f) = \tilde{x}_1$ .

Specifically, for  $f(x_1) = \Delta x_1$  with  $\Delta = -\tilde{u}'(\tilde{x}_1, \theta)$ , we get that  $\tilde{u}'(\tilde{x}_1, \theta) + f'(\tilde{x}_1) = \tilde{u}'(\tilde{x}_1, \theta') + f'(\tilde{x}_1) = 0$ . From the assumption that  $u(x, \theta)$  is concave in  $x$  and  $\mathcal{E}$  is convex,  $\tilde{x}_1$  is the optimal choice of  $x_1$  for both  $\theta'$  and  $\theta$  under the added incentive  $f(x_1)$  and so  $x_1^*(\theta|f) = x_1^*(\theta'|f) = \tilde{x}_1$ . □

## C Value-Added Estimation

In this appendix, we describe the different forms of value-added we use in the paper and how we estimate them. Our estimation procedure follows Mulhern and Opper (2021), although Mulhern and Opper (2021) does not control for experience in their estimates.

### C.1 Residualizing Outcomes

Let  $i$  index students,  $j$  index teachers,  $c$  index classrooms, and  $t$  index years. Let  $\tau()$  be a function that describes when an outcome is realized. For contemporaneous outcomes,  $\tau(k) = 0$ , while for outcomes realized in the future, like next year’s test scores,  $\tau(k) > 0$ . Our statistical model of outcomes, for a specific subject-level, is:

$$y_{i,t+\tau} = \Lambda X'_{it} + \sum_{e'} \rho_{e'} \mathbb{1}\{e_{jt} = e'\} + \mu_{jt} + \nu_{ct} + \phi_{c',t+1} \mathbb{1}(\tau \geq 1) + \phi_{c',t+2} \mathbb{1}(\tau = 2) + \epsilon_{it} \quad (33)$$

where  $e_{jt}$  is a teacher’s experience level (with all teachers with six or more years of prior experience grouped into one level).

We have 4 types of outcomes:

1. **Targeted outcome:** test scores in year  $t$
2. **Untargeted outcomes:** test scores in year  $t + 1$ , test scores in year  $t + 2$ , attendance rate in year  $t$ , attendance rate in year  $t + 1$ , grades in tested subject in year  $t$ , grades in tested subject in year  $t + 1$ , grades in untested subjects in year  $t$ , grades in untested subjects in year  $t + 1$
3. **Index of untargeted outcome:** an index of the above outcomes (constructed below)
4. **Long-term outcome:** whether the student graduates high school on-time

For ease of exposition, label the four outcomes (at the student-level)  $y_{it}^1, \bar{y}_{it}^2, y_{it}^3, y_{it}^4$ .

In a first step, we standardize each outcome in  $y_{it}^1$  and  $\bar{y}_{it}^2$  to have mean 0 and standard deviation 1 for each grade-year in NYC.

We then residualize outcomes in  $y_{it}^1, \bar{y}_{it}^2$ , and  $y_{it}^4$  by regressing them on a set of observable characteristics and teacher fixed effects:

$$y_{i,t+\tau} = \Lambda X'_{it} + \sum_{e'} \rho_{e'} \mathbb{1}\{e_{jt} = e'\} + \mu_j + v_{it}. \quad (34)$$

where  $v_{it} = \mu_{jt} - \mu_j + \nu_{ct} + \phi_{c',t+1} \mathbb{1}(\tau \geq 1) + \phi_{c',t+2} \mathbb{1}(\tau = 2) + \epsilon_{it}$ . We run separate regressions for each outcome, subject (math or ELA), and level (elementary or middle) combination.

We let the set of controls,  $X'_{it}$ , vary by outcome. For all outcomes, we include year dummy variables and indicators for whether the student receives free or reduced price lunch and whether the student is an English language learner, male, Black, Hispanic, and Asian. For lagged outcomes we use:

- Cubic polynomials in  $t - 1$  test scores for each subject – used for test scores in  $t, t + 1, t + 2$ , subject grades in  $t$  and  $t + 1$ , untested subject grades in  $t$  and  $t + 1$ , graduation
- Cubic polynomial in  $t - 1$  attendance rate – used for attendance rate in  $t$  and  $t + 1$ .

For each student, we construct two residuals:

1.  $\hat{v}_{it}^1 = y_{i,t+\tau} - \hat{\Lambda}X'_{it} - \sum_{e'} \hat{\rho}_{e'} \mathbb{1}\{e_{jt} = e'\}$
2.  $\hat{v}_{it}^2 = y_{i,t+\tau} - \hat{\Lambda}X'_{it}$

The residuals differ in whether the teacher’s experience effects are included.

## C.2 Constructing Measures of Teacher Output in Each Year

We then construct two (noisy) measures of a teacher’s output in year  $t$  (for each subject-level) by taking the mean of the two student residuals over each teacher-year combination:

1.  $\hat{\mu}_{jkt}^1 = \frac{1}{N_{jkt}} \sum_{i \in \mathcal{I}_{jkt}} \hat{v}_{ikt}^1$
2.  $\hat{\mu}_{jkt}^2 = \frac{1}{N_{jkt}} \sum_{i \in \mathcal{I}_{jkt}} \hat{v}_{ikt}^2$

where  $\mathcal{I}_{jkt}$  is the set of  $N_{jkt}$  students with outcome  $k$  who are taught by  $j$  in year  $t$ . We construct these measures for each outcome in  $y_{it}^1$ ,  $\bar{y}_{it}^2$ , and  $y_{it}^4$ .

For analysis in Section V we use the measure that includes experience effects ( $\hat{\mu}_{jkt}^2$ ) when it is the outcome variable. The exception is Figure 6, where we use the version without experience effects ( $\hat{\mu}_{jkt}^1$ ) to show the flatness of the curve for unexposed cohorts. For analysis in Sections VI and VII, we use the version without experience effects ( $\hat{\mu}_{jkt}^1$ ) as the outcome because we project it onto shrunken value-added measures that exclude the experience profile.

## C.3 Constructing Forecasts of Teacher Output

The prior measures are noisy estimates of a teacher’s realized output in a given year. For classifying teachers according to their unincentivized output, we construct forecasts that incorporate data from multiple years. We construct forecasts for each outcome in  $y_{it}^1$  and  $\bar{y}_{it}^2$ .

We follow Mulhern and Opper (2021) and refer there for the details. The key estimation points are:



- The estimates are from a joint Empirical Bayes procedure where the estimates are shrunk jointly.
- We estimate using data from unincentivized periods only. We produce estimates for all years as if they were unincentivized (even if they were actually incentivized).
- We estimate using the non-experience residuals ( $\hat{\mu}_{jkt}^1$ ).
- We allow for the non-experience component of a teacher’s effect to drift over time. We let drift rates vary depending on the difference in years between measures, where we estimate a constant drift rate for year differences at least 3.
- We keep teacher-subject-levels with at least ten students with an outcome.
- In forecasting a teacher’s output in year  $t$ , we use data from all years except year  $t$  (i.e., a leave-out estimator).
- For each  $jt$ , some of the output measures may be missing. In these cases, we predict the missing measures with forecasts based on the non-missing measures. Specifically, we estimate a separate joint Empirical Bayes procedure for each combination of non-missing measures and use it to forecast the missing measures. The identifying assumption is that the missingness is random conditional on the forecast of the non-missing measures.
- Our inclusion of  $\phi_{c',t+1}\mathbb{1}(\tau \geq 1) + \phi_{c',t+2}\mathbb{1}(\tau = 2)$  in the model means that the correlation structure of a teacher’s students’ residuals varies with the overlap in their classes in  $t + 1$  and  $t + 2$ . We incorporate this in constructing the forecasts.

We denote these forecasts as  $\tilde{\mu}_{jkt}$  and use them in Sections VI and VII. The exception is the first four columns in Table 8, where we show how treatment effects vary with heterogeneity in *only* the targeted or untargeted forecasts, rather than jointly. For these columns, we construct forecasts from Empirical Bayes procedures that only include the relevant outcomes ( $y_{it}^1$  or  $\bar{y}_{it}^2$ ).

## C.4 Constructing the Index of Untargeted Output

We create the untargeted index by anchoring the measures to their relative predictiveness of whether a student graduates from high school:

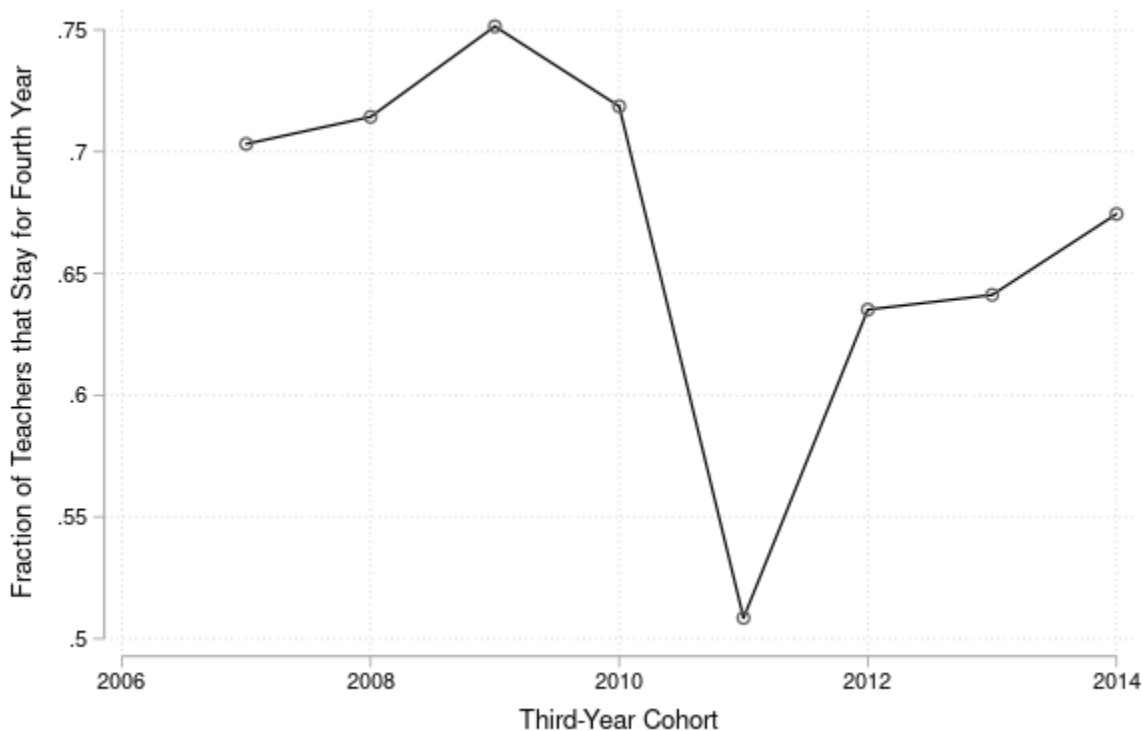
$$\hat{\mu}_{jlt}^1 = \omega' \tilde{\mu}_{jlt} + v_{ij}, \quad (35)$$

where  $l$  corresponds to the graduation outcome and  $\omega$  is a vector of anchoring weights. We estimate using data from the unincentivized period.

We use the estimated weights to construct two measures: targeted output ( $\tilde{\mu}_{jt}^T = \tilde{\mu}_{j1t}$ ) and an index of untargeted output ( $\tilde{\mu}_{jt}^U = \frac{1}{\hat{\omega}_1} \sum_{k=2}^K \hat{\omega}_k \tilde{\mu}_{jkt}$ ). We also apply these weights to the unshrunk measures for further indices,  $\hat{\mu}_{jt}^T = \hat{\mu}_{jkt}$  and  $\hat{\mu}_{jt}^U = \frac{1}{\hat{\omega}_1} \sum_{k=2}^K \hat{\omega}_k \hat{\mu}_{jkt}$ .

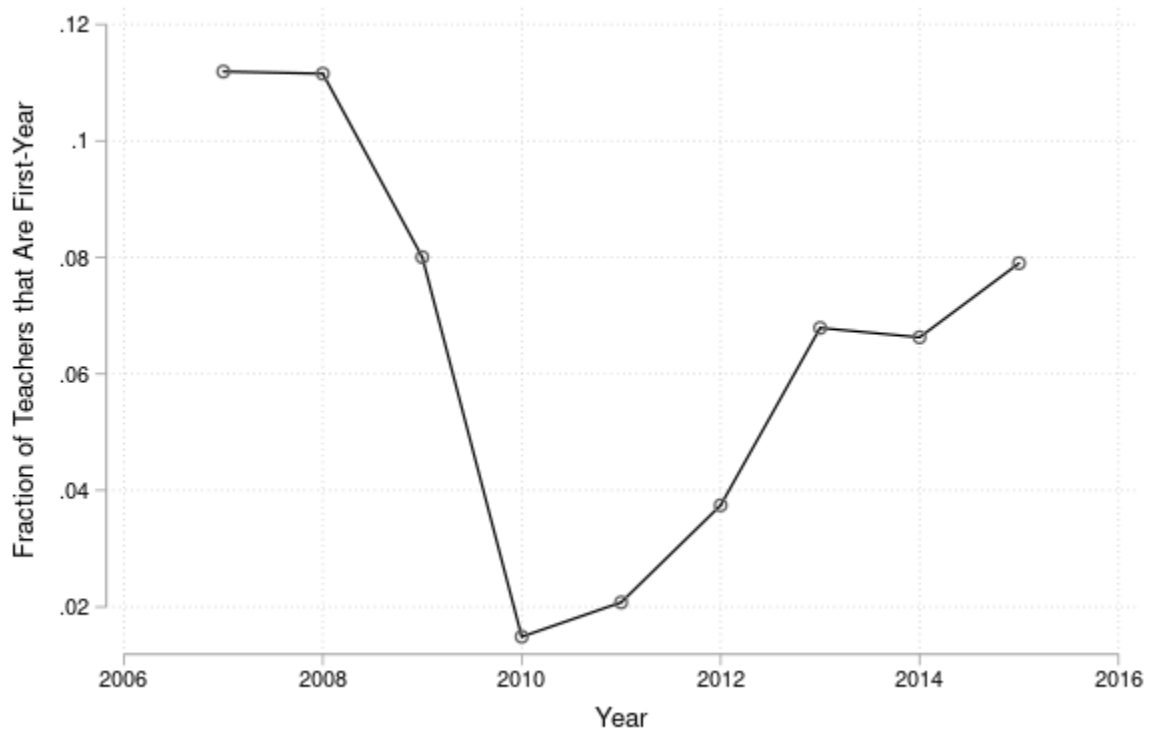
## D Appendix Figures

Figure A1: Teachers' Persistence to Fourth Year of Teaching



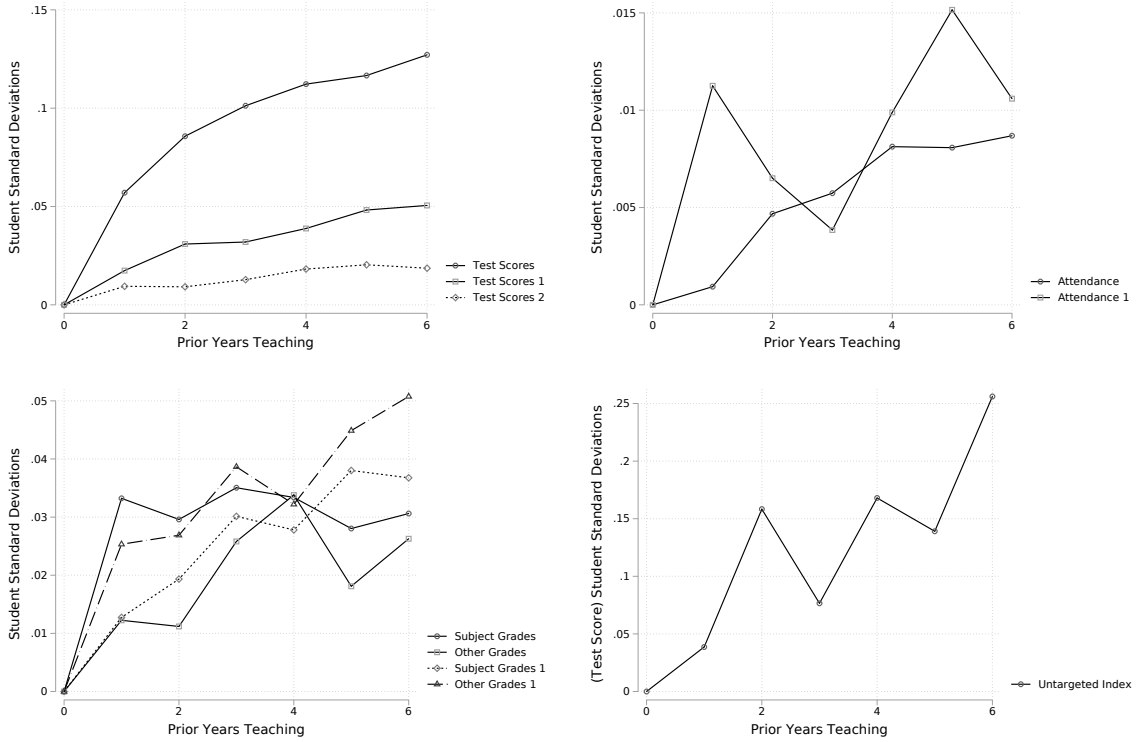
Note: This figure shows the fraction of teachers who were in the district in their third year of teaching who remain in the district their fourth year. The x-axis classifies cohorts based on the academic year when their third year of experience occurred.

Figure A2: New Teachers' Fraction of Workforce



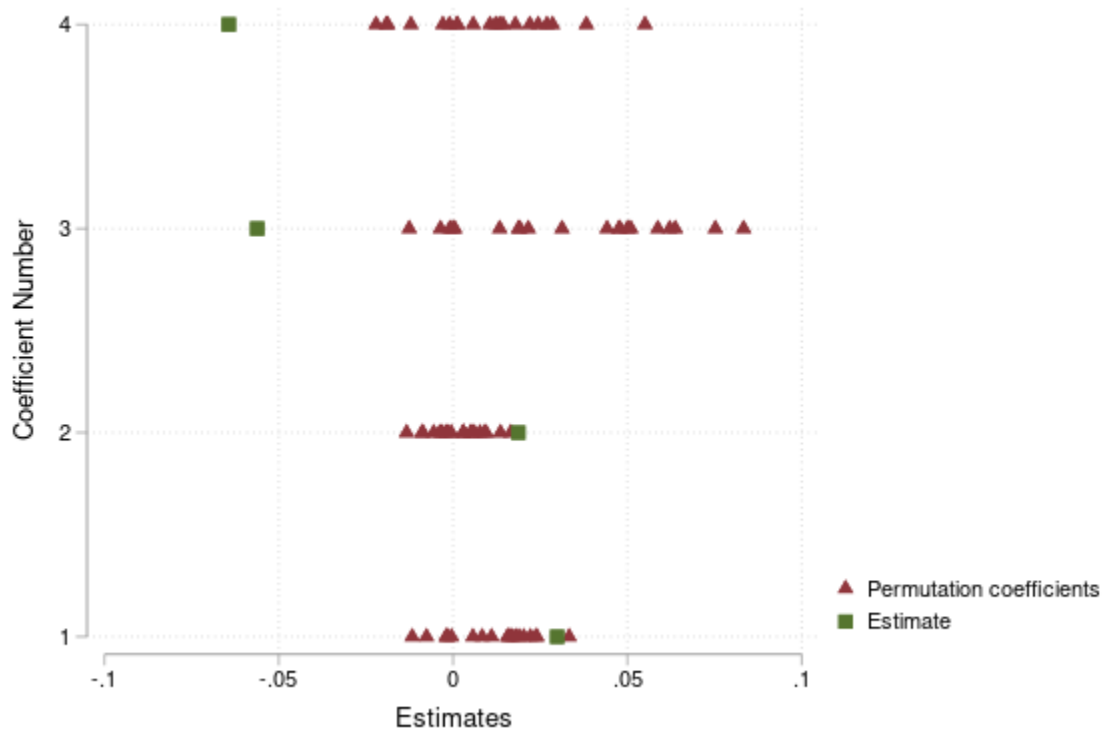
Note: This figure shows the fraction of each year's teacher workforce that first-year teachers comprise in NYC.

Figure A3: Experience profiles



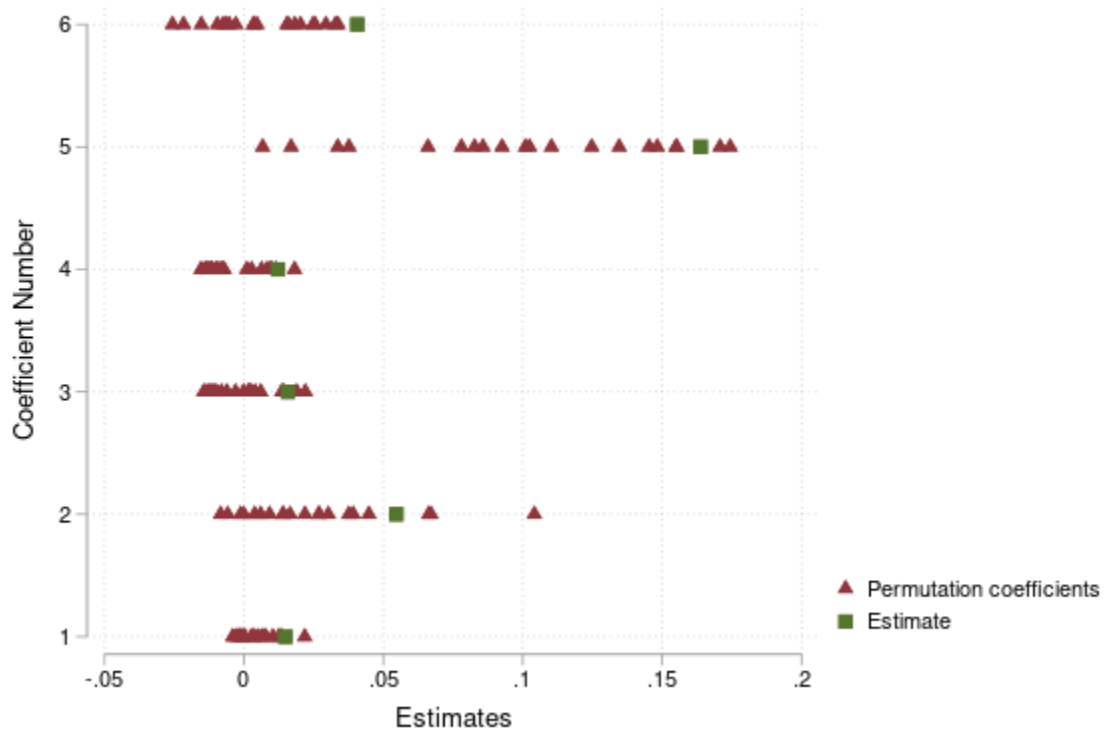
Note: This figure shows the estimated experience profile for our measures, where output in the first year of teaching is normalized to 0 and teachers with six or more years of experience are grouped into one category. The top left panel shows the score measures (current and future), the top right panel shows the attendance measures, the bottom left panel shows the grades measures, and the bottom right panel shows the untargeted index. In the first three panels, units are student standard deviations on each measure. In the bottom right panel, the units are student standard deviations on test scores.

Figure A4: Permutation Tests for Table 4 Estimates



Note: This figure shows permutation tests for the coefficients on “Incentive” in Table 4. The labeled rows correspond to the columns in Table 4. For each test, we assign the policy change to a different set of cohorts. In the correctly specified timing, the 2010 policy affects the cohorts with 0-2 years of prior experience. In the placebo timings, we let the policy affect cohorts with 3-5, 4-6, 5-7, etc. years of prior experience (up to 22-24). We maintain the structure of the policy change and the sample restrictions we impose in our main analysis. The correctly specified regression is labeled with a green square while the placebo estimates are red triangles.

Figure A5: Permutation Tests for Table 8 Estimates



Note: This figure shows permutation tests for the coefficients in columns (5) and (6) of Table 8. Labeled rows (1)-(3) correspond to the estimate in column (5) in Table 8 (from top to bottom). Labeled rows (4)-(6) correspond to the estimate in column (6) in Table 8 (from top to bottom). For each test, we assign the policy change to a different set of cohorts. In the correctly specified timing, the 2010 policy affects the cohorts with 0-2 years of prior experience. In the placebo timings, we let the policy affect cohorts with 3-5, 4-6, 5-7, etc. years of prior experience (up to 22-24). We maintain the structure of the policy change and the sample restrictions we impose in our main analysis. The correctly specified regression is labeled with a green square while the placebo estimates are red triangles.

## E Appendix Tables

Table A1: Correlation between Value-Added in  $t$  and  $t - 1$

	Corr b/t VA $t$ and VA $t-1$	Corr b/t VA $t$ and Test Score VA $t-1$	Corr b/t VA $t$ and Index VA $t-1$
Test Score	0.436	0.436	0.161
Untargeted Index	0.560	0.160	0.560
Test Score $t+1$	0.357	0.196	0.201
Test Score $t+2$	0.446	0.118	0.249
Attendance	0.553	0.024	-0.121
Attendance $t+1$	0.265	0.028	0.052
Subject Grades	0.535	0.046	0.121
Other Grades	0.565	0.102	0.344
Subject Grades $t+1$	0.290	0.045	0.178
Other Grades $t+1$	0.453	0.074	0.357

This table shows correlations between a teacher's (shrunk) value-added measure in  $t$  and various (unshrunk) value-added measures in  $t - 1$ . The columns show correlations with lagged value-added in (1) the same outcome, (2) the targeted measure ("Test Score"), and (3) the index of untargeted measures. We include the targeted measure ("Test Score"), the index of untargeted measures, and each untargeted measure separately.

Table A2: Outcomes' Univariate Relationship to Graduation

	Grad	Grad	Grad	Grad	Grad	Grad	Grad	Grad	Grad
Test Score VA	0.0879*** (0.00727)								
Test Score 1 VA		0.0635*** (0.00631)							
Test Score 2 VA			0.0716*** (0.00798)						
Attendance VA				0.0868*** (0.0119)					
Attendance 1 VA					0.0358*** (0.00611)				
Subject Grade VA						0.0618*** (0.00435)			
Other Grade VA							0.0697*** (0.00339)		
Subject Grade 1 VA								0.0923*** (0.00617)	
Other Grade 1 VA									0.106*** (0.00481)
N Teachers	11689	11694	11689	11695	11695	11692	11670	11678	11676
Mean DV	-0.0154	-0.0154	-0.0153	-0.0154	-0.0155	-0.0153	-0.0153	-0.0153	-0.0153
N	32066	32089	32044	32157	32149	32078	31937	31997	31991

This table shows the univariate regression of (the residual of) whether a student graduated from high school on the targeted forecasted measure or on each untargeted forecasted measure. The forecasts come from our multi-year multi-dimensional value-added model that is estimated on unincentivized periods only. Forecasts are constructed for all periods and leave out data from that year. The regression includes only observations from the unincentivized period. In the variable labels, the number indicates the number of years in the future the outcome is realized. All variables are in (separate) standard deviation units.



Table A3: Anchoring Outcomes to Graduation

	Grad
Test Score VA	0.0334*** (0.0113)
Test Score 1 VA	0.0300* (0.0173)
Test Score 2 VA	0.0246* (0.0129)
Attendance VA	0.231*** (0.0211)
Attendance 1 VA	0.0134 (0.0156)
Subject Grade VA	-0.0149** (0.00693)
Other Grade VA	0.0352*** (0.00549)
Subject Grade 1 VA	-0.0295** (0.0142)
Other Grade 1 VA	0.111*** (0.0108)
N Teachers	11670
Mean DV	-0.0153
N	31937

This table shows the regression of (the residual of) whether a student graduated from high school on the targeted and untargeted forecasted measures. The forecasts come from our multi-year multi-dimensional value-added model that is estimated on unincentivized periods only. Forecasts are constructed for all periods and leave out data from that year. The regression includes only observations from the unincentivized period. In the variable labels, the number indicates the number of years in the future the outcome is realized. All variables are in (separate) standard deviation units.

Table A4: Principal Component Analysis

	First Component	Second Component
Test Score 1 VA	0.156	0.376
Test Score 2 VA	0.151	0.346
Attendance VA	0.015	0.124
Attendance 1 VA	0.051	0.282
Subject Grades VA	0.422	-0.208
Other Grades VA	0.804	-0.350
Subject Grades 1 VA	0.186	0.450
Other Grades 1 VA	0.302	0.525

This table shows the first two components of a Principal Component Analysis on the value-added on each untargeted measure.

Table A5: Effect of Policy Change on Probationary Period Output – By Subject and Level

	Score	Score	Score	Score	Index	Index	Index	Index
Incentive	0.0148 (0.0103)	0.0160* (0.00957)	0.0226* (0.0130)	0.00872 (0.00906)	-0.0657* (0.0386)	-0.0578 (0.0387)	-0.0653* (0.0372)	-0.0669 (0.0513)
Fixed Effects	Teacher	Teacher	Teacher	Teacher	Teacher	Teacher	Teacher	Teacher
Sample	Math	ELA	Elem	Middle	Math	ELA	Elem	Middle
N Teachers	13184	13868	13389	9909	11961	12494	9036	8030
Mean DV	0.123	0.0810	0.121	0.0542	0.0290	0.0255	-0.0152	0.0915
N	57796	59586	84107	42934	48665	49973	67931	32268

This table shows the causal effect of the tenure policy change on targeted and untargeted output in the probationary period. The columns show the effects in different subsamples, split by the tested subject (Math or ELA) and the level of school (elementary or middle). All regressions include teacher fixed effects. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. The teachers with more than 3 years of experience are only included if they finished the standard probationary period before the tenure policy change. All outcome units are test score student standard deviations.

Table A6: Effect of Policy Change on Specific Untargeted Outcomes, Cohort Fixed Effects

	Score 1	Score 2	Attend	Attend 1	Grades	Other Grades	Grades 1	Other Grades 1
Incentive	-0.00572 (0.00981)	0.0209 (0.0137)	0.00491** (0.00212)	0.000419 (0.00900)	0.0139 (0.0186)	0.0218 (0.0159)	-0.0351* (0.0181)	-0.0188 (0.0175)
Teacher FEs	No	No	No	No	No	No	No	No
Mean DV	0.0496	0.0548	0.00417	0.0314	0.0251	-0.00873	0.0398	0.0409
N	110038	87893	132253	110572	36847	48078	60916	74374

This table shows the causal effect of the tenure policy change on individual untargeted (residualized) outcomes in the probationary period. The “1” or “2” in the column headers indicate the measure’s number of years into the future. “Grades” are in the tested subject while “Other” grades are in untested subjects. All variables are standardized at the grade-year level to have mean 0 and standard deviation 1 in the full population. All regressions include cohort fixed effects. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. The teachers with more than 3 years of experience are only included if they finished the standard probationary period before the tenure policy change.

Table A7: Effect of Policy Change on Specific Untargeted Raw Outcomes, Teacher Fixed Effects

	Score	Score 1	Score 2	Attend	Attend 1	Grades	Other Grades	Grades 1	Other Grades 1
Incentive	0.0135 (0.0133)	-0.0252 (0.0156)	-0.0151 (0.0192)	0.00435 (0.0109)	-0.0242* (0.0134)	0.0337 (0.0214)	0.000386 (0.0199)	-0.0652*** (0.0230)	-0.0429* (0.0219)
Teacher FEs	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Mean DV	-0.0533	-0.0483	-0.0432	-0.0445	-0.0429	-0.0741	-0.102	-0.0613	-0.0719
N	127678	106647	85244	127480	106756	33217	43512	57508	71371

This table shows the causal effect of the tenure policy change on individual untargeted (non-residualized) outcomes in the probationary period. The “1” or “2” in the column headers indicate the measure’s number of years into the future. “Grades” are in the tested subject while “Other” grades are in untested subjects. All variables are standardized at the grade-year level to have mean 0 and standard deviation 1 in the full population. All regressions include teacher fixed effects. An observation is a teacher-subject-year. Standard errors are clustered by teacher. The sample covers years 2006 to 2014. The teachers with more than 3 years of experience are only included if they finished the standard probationary period before the tenure policy change.

Table A8: Effect of Policy Change on Outcomes, Controlling for Next Two Years’ Teachers

	Untargeted Index	Untargeted Index	Untargeted Index	Untargeted Index	Score 2	Score 2
Incentive	-0.0749* (0.0389)	-0.0794** (0.0389)	-0.0837** (0.0364)	-0.0837** (0.0364)	-0.00476 (0.0135)	-0.00582 (0.0135)
Incentive 1		-0.00913 (0.0253)		-0.0102 (0.0198)		-0.0267*** (0.00814)
Incentive 2		-0.170*** (0.0239)		0.00720 (0.0168)		-0.0163** (0.00740)
Fixed Effects	Cohort	Cohort	Teacher	Teacher	Teacher	Teacher
N Teachers	11324	11324	11324	11324	11202	11202
Mean DV	0.0487	0.0487	0.0487	0.0487	0.0277	0.0277
N	63796	63796	63796	63796	63361	63361

This table shows how our estimates of the effect of the policy change on teachers’ untargeted outcomes varies depending on whether we control for the treatment status of the teachers in the two subsequent years. Adjacent columns show the regression with and without these controls, where we restrict the samples to the teachers for whom we can classify the two subsequent teachers. “Score 2” is the test score in year  $t + 2$  and is in test score  $t + 2$  student standard deviation units while the untargeted index includes all of the untargeted outcomes and is in test score  $t$  student standard deviation units. “Incentive 1” is the fraction of teacher  $j$ ’s students in year  $t$  that have an incentivized  $t + 1$  teacher, and “Incentive 2” is the fraction that have an incentivized  $t + 2$  teacher.

Table A9: Effect of Policy Change on Outcomes, Controlling for Next Year’s Teachers

	Score 1	Score 1	Attend 1	Attend 1	Grades 1	Grades 1	Other Grades 1	Other Grades 1
Incentive	-0.0262** (0.0116)	-0.0264** (0.0116)	-0.0348*** (0.0114)	-0.0346*** (0.0114)	-0.0523** (0.0235)	-0.0508** (0.0235)	-0.0428** (0.0201)	-0.0424** (0.0201)
Incentive 1		-0.0116* (0.00672)		0.00868 (0.00625)		0.0522*** (0.0154)		0.0169** (0.00848)
Fixed Effects	Teacher	Teacher	Teacher	Teacher	Teacher	Teacher	Teacher	Teacher
N Teachers	11299	11299	11310	11310	7407	7407	8231	8231
Mean DV	0.0266	0.0266	0.00370	0.00370	0.0234	0.0234	0.0276	0.0276
N	63719	63719	63726	63726	31770	31770	36882	36882

This table shows how our estimates of the effect of the policy change on teachers’ untargeted outcomes varies depending on whether we control for the treatment status of the teacher in the subsequent year. Adjacent columns show the regression with and without these controls, where we restrict the samples to the teachers for whom we can classify the subsequent teacher. All outcomes are realized in year  $t + 1$ . Units are student standard deviations for the respective outcome. “Incentive 1” is the fraction of teacher  $j$ ’s students in year  $t$  that have an incentivized  $t + 1$  teacher.

Table A10: Voluntary Attrition by Targeted and Untargeted Value-Added

Voluntary Attrition	
Targeted VA	-0.278*** (0.0430)
Untargeted VA	-0.0210** (0.00914)
N Teachers	5834
Mean DV	0.394
N	8693

This table shows how teachers’ voluntary attrition rates vary with their targeted and untargeted forecasted value-added. We consider tenured teachers in their seventh year of teaching in the district and determine whether the teachers left the sample before the end of our data. If so, we label them as voluntary attrition. Targeted and untargeted forecasted value-added are estimated using data from unincentivized periods only and are both in student test score standard deviation units.

Table A11: Output under Different Screening Regimes – 2006 Cohort

	Obs.	Mean Targeted Output	Mean Untargeted Output	Mean Total Output
<i>Teachers' Tenure under Different Responses</i>				
Never Tenured	477	-0.119	-0.771	-0.890
Only Tenured w/o Behavioral Response	50	-0.009	-1.262	-1.271
Only Tenured w/ Behavioral Response	50	-0.041	0.239	0.198
Always Tenured	1,024	0.111	-0.222	-0.111
<i>Tenured Teachers under Different Policies</i>				
Screening w/o Behavioral Response	1,074	0.106	-0.271	-0.165
Screening w/ Behavioral Response	1,074	0.104	-0.201	-0.097
(Infeasible) Screening on Both Dimensions	1,074	0.065	0.010	0.074
<i>Gains Relative to Infeasible First-Best</i>				
Gains (Fraction)				0.285

This table shows the mean output for different groups of teachers and under different policies. The sample is teachers who started in the district in 2006. The top panel splits teachers into four groups based on whether they would receive tenure in a regime without a behavioral response and whether they would receive tenure in a regime with a behavioral response. The middle panel shows the mean output associated with the set of teachers receiving tenure under different policies. The first two policies are screening on the targeted measure, without and with a behavioral response. The last policy is an infeasible policy screening on the sum of output across both dimensions. The final panel shows the fraction of gains the behavioral response achieves relative to the distance between the screening without behavioral response regime and the infeasible policy screening on the sum of output. “Mean Targeted” output is the mean forecasted test score value-added. “Mean Untargeted Output” is the mean forecasted value-added on the untargeted index. “Mean Total Output” is the sum of the mean forecasted value-added across the targeted and untargeted measures. All outcome units are test score student standard deviations.