

Score change-based school accountability rules and grade retention practices

(Latest version can be found [here](#))

Will Labadie

December 14, 2020

Abstract

Do accountability rules affect public school retention practices? Using a simple model of grade retention, I show that an administrator will retain students differently depending on the accountability ratings criteria he seeks to maximize. In particular, I find that an administrator will retain fewer students in order to maximize both school-wide average level and year-to-year change in scores than he will to maximize the average score alone, and that this is most pronounced in the final grade offered by his school. Using a novel dataset of school-grade level retention rates for 7 states in the U.S. and an event study design, I find that about 18% fewer students are retained on average each year when a state adds student growth to the accountability criteria by which schools are evaluated. This number roughly corresponds to around 100,497 fewer retained students each year nationwide, and \$1.4 billion saved in public school expenditures. I further find that administrators do retain significantly fewer students in the last grade offered by their schools - when retaining a student keeps him in a school's pool of test-takers for an additional year - implying that administrators use retention strategically. My results suggest that school administrators are willing to use retention as a tool for optimizing their schools' accountability ratings, and that the components of accountability systems do alter administrator behaviors.

1 Introduction

Each year, public school administrators endeavor to guide their schools to satisfactory ratings as defined by their state's adopted accountability system. It is well-established that administrators value maximizing their schools' ratings, and that they respond to the specific criteria by which their schools are rated (Rockoff and Turner, 2010; Chakrabarti, 2013). Accountability systems do not include targeted goals around student retention, but an administrator's retention practices could have significant impacts on her school's overall rating. To the extent that administrators decide whether or not to retain a student with an eye towards optimizing their schools' ratings, different ratings criteria could lead to different retention decisions.

In this study, I evaluate the impact of a common accountability measure, year-to-year changes in student standardized test scores, on public school retention policy. Retention has been shown to have profound effects on students in the short and long term, but the school-side determinants of retention are not well-understood (Schw-erdt, West and Winters, 2017; Eren, Lovenheim and Mocan, 2018). This paper sheds light on how administrator incentives could affect the retention decision.

I find that in the presence of score change-based accountability criteria, fewer public school students are held back, and that the effect grows over time. After 6 or more years under an accountability system that includes within-student score change, retention rates decrease by 18.3%. A back-of-the-envelope calculation using average enrollment counts from the Common Core of Data suggests that these estimates correspond to 100,497 fewer retained students each year nationwide. The reduction in years in the school system results in an average of \$1.4 billion less in student expenditures.¹ Back-of-the-envelope calculations based on other researchers' estimates of the effects of retention suggest that the reduction in retention also potentially avoids \$270 million in lost wages and \$91 million in crime-related costs incurred by retention of students, though it does potentially increase expenses associated with providing additional remedial math and reading courses by \$518 million.

I find that the effect of growth-based accountability criteria on retention is more pronounced in the final grade offered by a school. This result is consistent with the predictions of my statistical framework, and holds when the sample is restricted to elementary, middle, and high schools separately. This result suggests strategic retention on the part of school administrators, as the decision to retain or promote in the terminal grade either includes a student in the test-taking pool for an additional year, or removes him from the pool. That the differential effect of score change-based accountability criteria is negative suggests that administrators view retention as potentially harmful to year-over-year student score growth.

School administrators' incentives change with the accountability rules their schools are subject to. Through retention, an administrator has a degree of control over the level of test faced by a given student. If the administrator expects that a grade-4 student is unlikely to pass the grade 5 exam, but likely to pass the grade 4 exam, retaining that student could have a non-negligible effect on his school's accountability rating. Under a minimum-competency-based accountability system, strategic retention of this sort could have an especially clear effect on the school's accountability score. However, under a system that rates schools on both passing rate and within-student score change, the effect could be less clear. Because test scores are noisy measures of ability, the uncertainty of retaining a student for an additional year may outweigh the potential boon to the school passing rate since the value of passing one indicator is smaller than it would in a minimum-competency-based system. The uncertainty effect may be especially salient in the final grade offered by the school, since a promoted

¹Based on current per-student expenditures of \$13,847 (National Center of Education Statistics, 2019).

student will not contribute to the school's accountability score after her promotion, while a retained student will.

On the whole, my results suggest that fewer students are retained under score change-based systems. To the extent that minimum-competency-based systems encouraged retention of students that otherwise would not have been retained for maximizing passing rates, this is a good thing. Prior research on retention has found substantial negative effects in the long term (Eren, Lovenheim and Mocan, 2018; Brodaty, Gary-Bobo and Prieto, 2013). This is especially true for older retained students (Jacob and Lefgren, 2009). I find some evidence of higher retention among younger students, however. It is possible that score change-based accountability systems encourage more targeted retention, since administrators would be harmed by retaining students that would be developmentally harmed by retention. This also would be a positive result, since the weight of evidence seems to suggest that retention is best used as a highly individualized intervention (Jacob and Lefgren, 2004; Fruehwirth, Navarro and Takahashi, 2016).

To study the effects of score change-based accountability systems on retention practices, I develop a statistical framework in which students' test scores progress over time and administrators set score thresholds for each grade; students that score below the relevant score threshold are retained, while those that score above the threshold are promoted. This framework predicts different score threshold patterns for maximizing school-wide passing rate and school-wide score change, consistently predicting lower retention in the final grade offered by a school under a wide set of parameters. Because schools must satisfy both passing rate and score change-based criteria, the value of passing one is lower than it would be under a passing rate-only regime. The uncertainty of student ability translating to scores combined with the lower value of passing one criteria combines to suppress retention rates in the final grade.

To estimate the impact of score change-based accountability criteria on retention, I exploit the adoption of within-student score change into seven states' accountability criteria over time. I use a difference-in-differences strategy and a novel dataset assembled with the help of the states' education agencies to evaluate the effect of these accountability system changes. While the timing of states adopting the score change criteria is not random, it is exogenous to any school-level decision-making. Following Goodman-Bacon (2018), I employ an event-study design in my analysis. I study the differential effect in the final grade offered by a school with a triple-differences technique.

In this study, I show that administrators retain students differently depending on the accountability criteria they face. Almond, Lee and Schwartz (2016) show that public school administrators do exercise influence over the retention decision, though this study is the first to show that they do so strategically. A large literature exploring the unintended consequences of accountability systems has established already that administrators value satisfactory ratings (Rockoff and Turner, 2010). In this paper, I

establish retention as a tool that administrators are willing to use in pursuit of better ratings for their schools. This paper contributes to a large literature exploring the tools administrators use to maximize their schools' ratings, including highly individualized tools (Figlio and Winicki, 2005; Reback, 2008). I show that administrators alter the body of students that contribute to their schools' ratings. Cullen and Reback (2006) find evidence of schools altering their test-taking pool by strategically exempting certain students from taking state exams to improve their ratings.

The rest of the paper is structured as follows. Section 2 will elaborate on school accountability and the incentives of administrators. Section 3 lays out my statistical framework. Section 4 discusses the data I use to test the predictions of the framework. Section 5 details the empirical strategies I use and the results of the analyses. Section 6 evaluates the costs and benefits of the change in retention practices under score change-based accountability based on estimates from the literature on the effects of retention. Section 7 concludes.

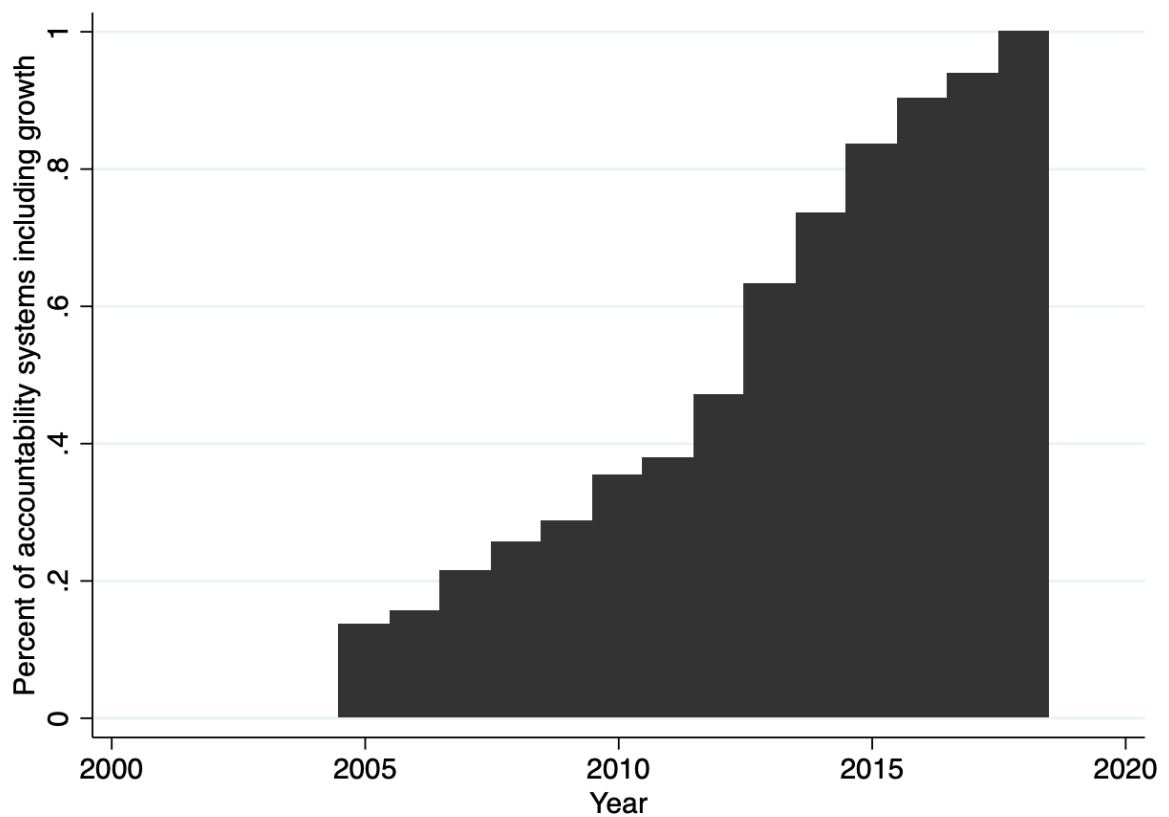
2 Background

Since the enactment of the No Child Left Behind (NCLB) Act of 2001, states have been required to evaluate public schools based on student performance on standardized math and English Language Arts (ELA) exams. Under NCLB, state education agencies were required to evaluate schools based on the rate at which students passed its standardized exam, and on the rate at which various subpopulations of students at the school passed the exam. However, NCLB allowed state education agencies a substantial degree of flexibility in designing their accountability systems; each agency could choose its own standardized exam to administer, and define what level of performance would be considered adequate, and accountability systems could and did evaluate schools on criteria beyond the minimum competency requirement. A common criteria, especially among states that had accountability systems in place prior to the enactment of NCLB, was year-to-year growth in individual students' scores.

The Every Student Succeeds Act (ESSA) of 2015 requires state education agencies to develop accountability systems that include multiple evaluation criteria, which must include "students' performance on the statewide assessment, high school graduation rates, and English language proficiency" as well as at least one additional measure of school quality that is left up to the individual agencies (Alliance for Excellent Education, 2016).

Figure 1 shows the percent of states whose accountability systems included a measure of student test score growth from the 2004-05 to 2017-18 school years. In the 2004-05 school year, 9 states evaluated schools partially based on the growth rate of students; by the 2017-18 school year, every state did. While states did not randomly choose whether or not to adopt growth into their accountability systems, the decision was exogenous to school-level decision-making, and affected the way in which schools were evaluated.

Figure 1: Presence of growth in accountability systems over time



Notes: Includes all 50 states and Washington, D.C.

Sources: Education Week Education Counts Research Center; individual states' education agencies.

Clearly, a school administrator's incentives are affected when the measures by which his school is evaluated change. Once student score growth is included in a school's rating, the administrator must attempt to satisfy the state-defined criterion along with the other indicators the state education agency considers in order to maximize his school's overall rating. A deep body of research has shown that authorities at schools use a number of tools to increase their school's overall rating. Reback (2008) finds that Texas students whose scores are relatively important to their school's overall rating perform better than expected, and finds evidence of finely targeted resource allocation and instruction to these "high-leverage" students; Figlio and Winicki (2005) find that schools in Virginia at risk of accountability-based sanctions increase the caloric content of school lunches on test days; Craig, Imberman and Perdue (2013) find that district administrators increase instructional budgets after facing a negative rating shock. While Cullen and Reback (2006) find that Texas schools manipulate the pool of test-takers via exemptions, no study has examined whether or not administrators hold students back strategically as well.

Under a minimum-competency based accountability system, a ratings-motivated

administrator might hold a student back for a number of reasons. If the student is in grade g and the administrator expects the student to pass the grade g exam in the coming school year, but not the grade $g + 1$ exam, he might retain the student to bolster the school's passing rate. This is particularly true if the student is from one of the subpopulations whose passing rates are explicitly valued under NCLB. The administrator might also retain the student if he believes that retention will positively impact the student's development, and have a lasting and positive impact on the student's future scores. Some research has shown that retention has positive short-run effects on student scores, but negative long-run effects on student development (Schwerdt, West and Winters, 2017); if a school is very likely to fail and its administrator estimates that the short-run effect of the student passing after being retained exceeds the later drag the student might have on the school's passing rate. On the other hand, if the administrator estimates that the long-run drag exceeds the short-run benefit, he would be more likely to promote the student.

Under an accountability system that includes year-to-year growth, the incentives to retain change. On one hand, an administrator might be more likely to retain a student if he expects the short-run benefits to be high enough; it would improve his school's likelihood of passing both the state's minimum competency and growth standards. On the other hand, if retention has persistent negative effects, holding a student back will make it less likely that his school passes the state's growth standard.

3 Conceptual framework

Because the overall effect of retaining a student is unclear, I build and simulate a basic one-school model of grade retention. The model relies on a factor model formulation of student skill accumulation and skill measurement by standardized exams, following the basic structure of Cunha and Heckman (2008). In each period t , students at the school earn a score on the state's standardized exam. Assume that there is only one subject exam, and only one score counts towards the school's rating. Assume that student scores are represented by the following dynamic factor structure:

$$s_{it} = \mu_{it} + \alpha_{it}\theta_{it} + \epsilon_{it} \quad (1)$$

where i represents the individual student, and θ_{it} is a dynamic factor for each student. I normalize the factor loading $\alpha_{it} = 1$ for all t in this version of the paper. ϵ_{it} represents a random normal shock to the measured scores, representing the noisiness of test scores in measuring latent ability.

I assume the following form for students' skills production:

$$\theta_{it} = \gamma_0 + \gamma_1\theta_{it-1} + \eta_{it} \quad (2)$$

where η_{it} is independent across students and over time for the same students. This version of the model assumes that retention has no effect on the skills production function. There is much evidence that this is not the case; however, I choose not to

include it in this version of the model for simplicity. I assume that the initial distribution $\theta_{i0} \sim N(0,1)$. I assume that scores grow predictably by 0.5 points per year, and assign $\gamma_0 = 0.5$. Finally, I assume that measurement is relatively noisy while the noise in skill accumulation is relatively small: $\epsilon_{it} \sim N(0,0.5)$. and $\eta_{it} \sim N(0,0.25)$.

The administrator's objective is to allocate students across grades to maximize either the percent of students passing the state exam (passing rate) or both the passing rate and the percent of students exhibiting sufficient test score growth from one year to the next (sufficient growth rate). The school's passing rate in year t is given by

$$\Pi_t = \frac{\sum_{j=1}^G \sum_{i=1}^N \mathbb{1}(g_{it} = j) \times \mathbb{1}(s_{it} \geq \pi_j)}{N_t} \quad (3)$$

where g_{it} represents student i 's grade level in year t , G represents the final grade offered by the school, and N represents the total number of students. π_j represents the externally set passing score for the grade j exam. N_t represents the number of students whose scores count towards the school's rating in year t :

$$N_t = \sum_{j=1}^G \sum_{i=1}^N \mathbb{1}(g_{it} = j). \quad (4)$$

The school's sufficient growth rate is assumed to be measured as follows:

$$\Lambda_t = \frac{\sum_{j=1}^G \sum_{i=1}^N \mathbb{1}(g_{it} = j) \times \mathbb{1}(s_{it} - s_{it-1} \geq \lambda_j)}{N_t} \quad (5)$$

where λ_g represents the externally set target amount of score growth for a student in grade j , and all other objects are as defined above.

The school administrator is able to choose which passing standard a given student is subject to through retention. The administrator can choose to retain a grade g student, and her scores will be evaluated relative to π_g and λ_g , or to promote the student, and her scores will be evaluated relative to π_{g+1} and λ_{g+1} . I assume that administrators make the retention decision by setting promotion thresholds δ_g for each grade. Grade g students that score at or above δ_g are promoted to grade $g + 1$, while grade g students that score below δ_g are retained and repeat grade g . For a 3-grade school, administrators control students' grade level g based on the following:

$$g_{it} = \begin{cases} 0 & \text{if } a_{it} = 0 \\ 1 & \text{if } a_{it} = 1 \text{ or } s_{it-1} < \delta_1 \text{ \& } g_{it-1} = 1 \\ 2 & \text{if } s_{it-1} \geq \delta_1 \text{ \& } g_{it-1} = 1 \text{ or } s_{it-1} < \delta_2 \text{ \& } g_{it-1} = 2 \\ 3 & \text{if } s_{it-1} \geq \delta_2 \text{ \& } g_{it-1} = 2 \text{ or } s_{it-1} < \delta_3 \text{ \& } g_{it-1} = 3 \\ 4 & \text{if } s_{it-1} \geq \delta_3 \text{ \& } g_{it-1} = 3 \end{cases} \quad (6)$$

where a_{it} represents the number of years that student i has been in school. The scores of students in grades 0 and 4 do not count towards the school's rating.

An administrator of a school under a passing rate-based accountability system has the following optimization problem:

$$\max_{\delta_1, \delta_2, \delta_3} \Pi(\delta_1, \delta_2, \delta_3) \quad (7)$$

while an administrator of a school under a passing rate and score change-based accountability system has the following optimization problem:

$$\max_{\delta_1, \delta_2, \delta_3} \Pi(\delta_1, \delta_2, \delta_3) + \Lambda(\delta_1, \delta_2, \delta_3) \quad (8)$$

In both cases, I assume the administrator has two important constraints. First, he must promote students that pass the standardized exam:

$$\delta_g \leq \pi_g \quad (9)$$

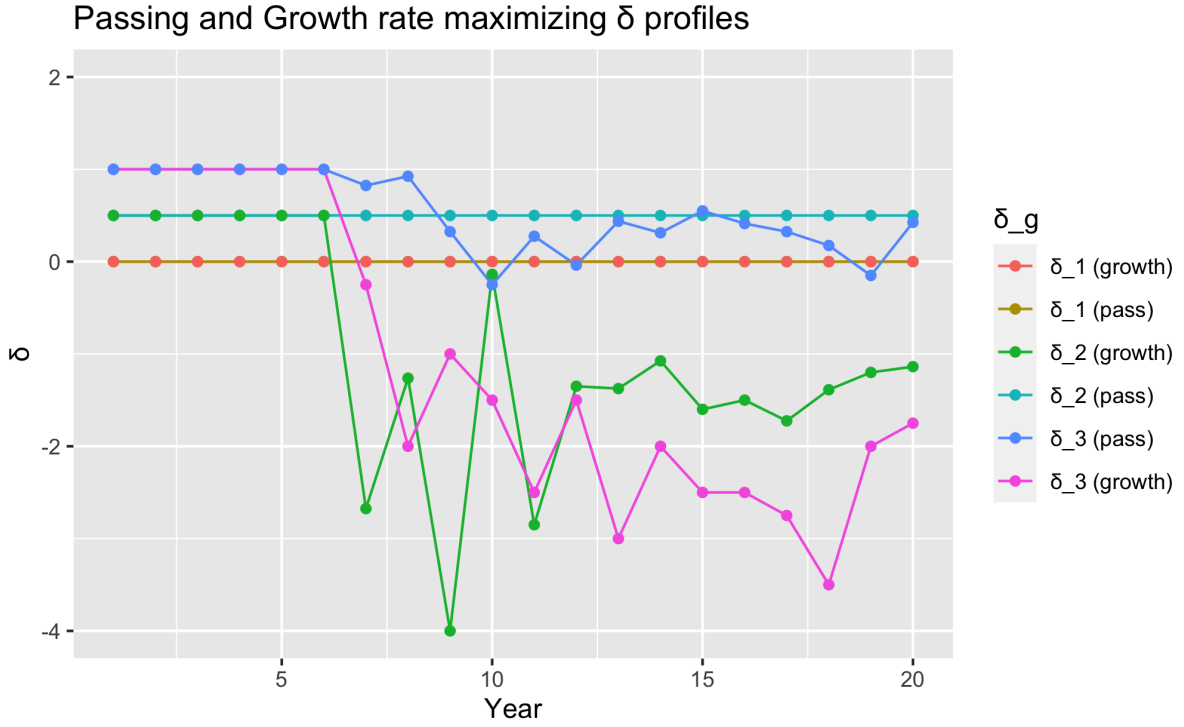
and he cannot retain a student in the same grade two years in a row

$$g_{it} \neq g_{it-2} \forall i. \quad (10)$$

These constraints are based on common rules adopted by school districts and state education agencies. The first has the additional feature of removing the option of retaining all passing students to maximize the school passing rate, which is appealing for the simulation of this model.

A Monte Carlo exercise, in which I perform 10,000 simulations given the parameter assumptions laid out above, yields some predictions about administrator behavior under the parameters assumed above. Figure 2 summarizes the results of the simulation. The results indicate that in the base case, an administrator attempting to maximize only passing rate should retain all failing students in grades 1 and 2, and should set the δ_3 threshold lower in grade 3, promoting students that passed and those that nearly passed. This strategy yields the highest average passing rate across simulations. An administrator interested in maximizing both passing and growth rate, on the other hand, should retain dramatically fewer students in grade 2 and especially in grade 3 relative to the passing-rate-only case.

Figure 2: Passing and growth rate-maximizing δ profiles



4 Data

I use a novel source of administrative data on retention rates at the school-grade level from the 2011-12 to the 2017-18 school year. I obtained this data from the education agencies of 7 states. It includes retention rates and school identifying information for every school in each state over this 7-year span. There are numerous types of schools in each state's public school system, which typically are either exempt from accountability rules, or the rules apply differently to them than they do to typical public schools. Some of the most common of these are charter schools, magnet schools, career and technical education (CTE) schools, and disciplinary alternative education schools. I omit these schools from my analysis, and focus on standard public schools, where accountability rules are likely to be most salient.

I combine this data with data on the components of states' accountability systems over the same period. I constructed this data set primarily using publicly available state statutes and administrative code regulations. I used data on state accountability system components from Education Week's Education Counts Research Center from the 2001-02 to 2011-12 school year to establish a baseline. This data was collected by Education Week via survey; state education agencies self-reported the components in their accountability systems each year, with documentation to support their responses. Their dataset includes indicators for whether or not a state "assigns ratings to all schools based on state-developed criteria", "uses measures of individual student

growth to rate schools”, and “uses measures of individual student growth for state ratings”. The data I collected from the state statutes and agencies fills in the same indicators for 2012-13 through 2017-18. Every state in my data had a criteria evaluating schools based on average student performance on standardized tests prior to adding a measure for individual student growth.

Different states define measures that rate schools on individual student growth differently. Most states that include a measure of individual student growth in school ratings measure the percent of students in each school whose test score increased by a sufficient amount from one year to the next (as defined by the state agency); some states define growth as the difference between a student’s score and a comparable group of students’ scores from the previous year, to attempt to average out idiosyncracies in test performance. However, all definitions emphasize year-to-year score increases.

I also rely on the Common Core of Data (CCD) in conjunction with this data for enrollment counts for various subpopulations in the school in each grade, the grades offered by each school, Title I status, charter status, and magnet status. The CCD variable defining a school’s grade span is crucial to my analyses on differential behavior by schools towards students in terminal grades. I use this data to construct a terminal grade indicator, which equals one if a given grade is the last offered by the school. Overall, my data covers around 12,400 public schools across 7 years in which states gradually adopt measures based on individual student score change in their accountability ratings systems.

I use this cross-state panel to analyze the effect that score change-based accountability systems have on public school retention practices. By using a cross-state panel, I am able to exploit the staggered adoption of score change-based accountability across states, comparing the difference in retention practices in states that adopt score change-based accountability to the difference in those that don’t, across the time period spanned by my data.

5 Empirical analysis

5.1 Measuring the Impact of Score change-based Accountability Systems on Retention Practices

I exploit the gradual adoption of score change-based accountability systems over time across states to study its effect on retention practices. The empirical strategy I use in this paper builds on a large literature in applied economics, which uses variation in timing of treatment to estimate a difference-in-differences effect.

My empirical strategy compares the difference in retention rates in states that adopt score change-based accountability systems before and after the switch to the difference for those that don’t over the same time period. Since all states had adopted some form

of score change-based accountability by the 2017-18 school year, a basic difference-in-difference design, in which some states are part of the “treatment” group and some are part of the “control” group, is not possible. Instead, I use an event-study design, which allows for analysis even with no untreated units (Abraham and Sun, n.d.). In addition, the event-study design allows for varying treatment effects over time, while the standard difference-in-differences framework is biased in the presence of treatment effects that vary over time.

My empirical approach relies on the exogeneity of the adoption of score change-based accountability systems. While the timing of a state’s adoption of a new accountability system and the content of that system are not likely to be random, it is likely to be exogenous to school administrators and their behavior. Accountability systems are adopted by state legislators, and retention practices are not mentioned in any of the accountability systems analyzed in this paper.

I use an event-study design of the following form:

$$\begin{aligned} r_{gcdst} = & \alpha + \beta_1 \mathbb{1}(t - T_s^* \leq -2)_{st} + \beta_2 \mathbb{1}(t - T_s^* \in [0, 2])_{st} \\ & + \beta_3 \mathbb{1}(t - T_s^* \in [3, 5])_{st} + \beta_4 \mathbb{1}(t - T_s^* \in [6, 8])_{st} \\ & + \gamma_c + \eta_g + \epsilon_{gcdst}. \end{aligned} \quad (11)$$

Here, r_{gcdst} represents the retention rate in grade g of campus c in district d of state s , year t , and T_s^* : year in which state s adopted student score change component in school ratings criteria. γ_c and η_g represent school and grade fixed effects respectively. In all specifications, I cluster my standard errors at the state-year level to allow for serial correlation, and because treatment is assigned at that level (Abadie et al., 2017).

Rather than estimate a separate coefficient for each post-adoption year, I group the event-study coefficients to gain precision. I cluster standard errors at the state-year level, to allow for serial correlation and correlation in the residuals of schools in the same state. The results of a placebo test show no pre-treatment effect and can be found in the first row of Table 1.

The administrator’s decision to retain is most impactful in the last grade offered at his school. If he expects a student to perform poorly if held back, promoting the student removes her from the pool of students whose test scores determine the school’s rating. To test the effect of score change-based accountability systems on retention practices in the last grade offered by a school, I estimate equation 2 separately for each grade, restricting attention to the first grade offered by a school, the middle grades offered, and the last grade offered. I also include a dummy variable identifying a grade as the last offered by a school as a third difference, estimating the following:

$$\begin{aligned} r_{gcdst} = & \alpha + \beta_1 \mathbb{1}(t - T_s^* \geq 0)_{st} \times \mathbb{1}(g = G_{ct}^T) + \beta_2 \mathbb{1}(t - T_s^* \geq 0)_{st} \\ & + \beta_3 \mathbb{1}(g = G_{ct}^T) + \gamma_c + \epsilon_{gcdst}, \end{aligned} \quad (12)$$

where G_{ct}^T represents the final grade offered at school c in year t . The results of a

placebo test show no statistically significant pre-treatment effect, and can be found in row 1 of Table 2.

5.2 Results

Table 1 presents the overall estimated effects of growth-based accountability on retention rates. Column (1) presents the results of estimating equation 2 without grade fixed effects, and column (2) includes them; column (2) represents my preferred specification. The results suggest modest negative effects on retention rates overall in the first few years after the adoption of growth-based accountability measures, though the estimated effects in these years are not statistically significant. After perhaps an adjustment period, retention rates fall by 0.620 percentage points - a 30% decrease from the pre-growth average.

Table 1: Short- and Long-run Effects of Score Change-based Accountability on Retention

	(1)	(2)
Years -5 to -2	0.097 (0.061)	0.136 (0.096)
Years 0 to 2	-0.060 (0.080)	-0.087 (0.085)
Years 3 to 5	-0.092 (0.087)	-0.122 (0.092)
Years 6 to 8	-0.589*** (0.094)	-0.620*** (0.098)
N	289820	289820
Grade FE		✓

Notes: Dependent variable is grade-level retention rate. Standard errors, clustered at the state-year level, are reported in parentheses. * denotes significance at 10%, ** denotes 5%, and *** denotes 1%. All regressions include school fixed effects.

Table 2 presents the main results on the effect of growth-based accountability on retention rates in the final grade offered by schools. The results suggest negative effects of a similar magnitude to the overall effect found in Table 1 in each post-treatment period; perhaps because the impact of promoting a student out of school is certain, the impact of the policy was more immediate in the terminal grades offered by schools.

Table 2: Differential Impact of Score Change-based Accountability in Terminal Grades

	(1)	(2)
Years -5 to -2	0.097 (0.061)	0.033 (0.063)
Years 0 to 2	-0.060 (0.080)	0.053 (0.079)
Years 3 to 5	-0.092 (0.087)	0.002 (0.087)
Years 6 to 8	-0.589*** (0.094)	-0.517*** (0.103)
$\mathbb{1}(g = G_c^T)$		-0.625*** (0.008)
Years -5 to -2 $\times \mathbb{1}(g = G_c^T)$		0.369*** (0.011)
Years 0 to 2 $\times \mathbb{1}(g = G_c^T)$		-0.687*** (0.059)
Years 3 to 5 $\times \mathbb{1}(g = G_c^T)$		-0.614*** (0.041)
Years 6 to 8 $\times \mathbb{1}(g = G_c^T)$		-0.524*** (0.089)
N	289820	289820

Notes: Dependent variable is grade-level retention rate. Standard errors, clustered at the state-year level, are reported in parentheses. * denotes significance at 10%, ** denotes 5%, and *** denotes 1%. All regressions include school fixed effects.

5.3 Effect heterogeneity

If a school administrator has certain beliefs about how different subpopulations might perform after being retained, they may be differentially impacted by the switch to growth-based accountability. In addition, since the performances of some subpopulations are heavily weighted in many accountability systems, the policy impact on retention practices may differ by subpopulation. To test for differential effects, I utilize two different approaches.

First, I consider only the case of Texas. The Texas Education Agency makes retention rates available at the subpopulation-grade-year level; as a result, I am able to test for the effect of Texas' switch to a growth-based accountability system on subgroup-specific retention rates. To do so, I use a simple difference-in-differences design. Assignment to treatment depends on the year of observation - Texas school ratings included a measure of student growth starting in the 2012-13 school year - and I use an indicator for whether or not a given grade is the last offered by a school as a measure of treatment intensity. Thus, I am comparing the retention rates of various subgroups in the last grade offered by a school to those of the same groups in all other grades offered by the school before and after exposure to a growth-based accountability system.

I estimate regressions of the following form:

$$r_{pgct} = \alpha + \beta_1 \mathbb{1}(t \geq 2013)_t \times \mathbb{1}(g = g_c^T)_{gct} + \beta_2 \mathbb{1}(t \geq 2013)_t + \beta_3 \mathbb{1}(g = g_c^T)_{gct} + n_{pct} + \eta_c + \epsilon_{gct} \quad (13)$$

where r_{pgct} represents the retention rate of subgroup p students in grade g at campus c in year t , g_c^T represents the last grade offered at campus c and n_{pct} represents the total number of students of subgroup p enrolled in school c in year t . I estimate this equation separately for elementary and middle schools.² The validity of this approach requires that retention rates in terminal and non-terminal grades exhibiting parallel trends in the pre-treatment period. Most subgroups of interest fail an informal parallel trend check; for this reason, I focus only on male, Black, Hispanic, and white students in this analysis. Figure 4 includes plots of the retention rate for each of these subgroups over time, from 2004-05 to 2017-18, in terminal and non-terminal grades.

5.4 Results

Tables 3, 4, 5, and 6 present the effects of growth-based accountability on the retention rates of male, Black, Hispanic, and white students respectively in Texas. I estimate equation 4 separately for elementary, middle, and high schools, as retention rates vary considerably across the three school types. The results presented in Table 3 show that the switch to growth-based accountability had a relatively small but significant effect on the retention rate of elementary school boys. The retention rate among boys in terminal grades decreases by .137 percentage points after the switch to growth-based accountability - a 4% decrease relative to the pre-growth average retention rate. This is largely fueled by a decrease in the retention rate of high school boys in terminal grades (i.e. high school graduates), and partly by a decrease in the retention rate of elementary school boys. The rate at which elementary school boys are held back decreases by .056 percentage points after the switch to growth-based accountability - a 1.6% decrease relative to the pre-growth average retention rate. The results of Tables 4, 5, and 6 show similar patterns - modest decreases in terminal-grade retention rates after the shift to growth-based accountability. Notably, the decreases for Black and Hispanic students seem to be driven primarily by decreases in the terminal-grade retention rate in middle school rather than in elementary or high school. Overall, I do not find compelling evidence that any of these subgroups are more affected by the policy change than others.

²High schoolers that repeat a course are counted as retained by TEA; for this reason, I choose to exclude them from my analysis.

Table 3: Effect of Score Change-based Accountability on Retention Rate Among Male Students: Texas

	Full sample	Elementary school	Middle school	High school
Post-growth	-0.014 (0.016)	-0.009 (0.009)	-0.003 (0.045)	-0.130 (0.121)
$\mathbb{1}(g = G_c^T)$	0.065** (0.024)	0.102*** (0.019)	0.075* (0.031)	-0.597*** (0.142)
Post-growth $\times \mathbb{1}(g = G_c^T)$	-0.137*** (0.037)	-0.056* (0.023)	-0.075 (0.051)	-0.660* (0.286)
N	245702	181322	41247	23133

Notes: Dependent variable is grade-level retention rate. Standard errors, clustered at the district-year level, are reported in parentheses. * denotes significance at 10%, ** denotes 5%, and *** denotes 1%. All regressions include school fixed effects and control for the total enrolled male students at the school.

Table 4: Effect of Score Change-based Accountability on Retention Rate Among Black Students: Texas

	Full sample	Elementary school	Middle school	High school
Post-growth	0.002 (0.002)	0.002 (0.001)	0.004 (0.003)	-0.001 (0.011)
$\mathbb{1}(g = G_c^T)$	0.025*** (0.007)	0.008* (0.004)	0.018 (0.011)	0.039 (0.023)
Post-growth $\times \mathbb{1}(g = G_c^T)$	-0.020** (0.007)	-0.005 (0.004)	-0.011 (0.010)	-0.040 (0.026)
N	309871	232045	52542	25284

Notes: Dependent variable is grade-level retention rate. Standard errors, clustered at the district-year level, are reported in parentheses. * denotes significance at 10%, ** denotes 5%, and *** denotes 1%. All regressions include school fixed effects and control for the total enrolled Black students at the school.

Table 5: Effect of Score Change-based Accountability on Retention Rate Among Hispanic Students: Texas

	Full sample	Elementary school	Middle school	High school
Post-growth	-0.005 (0.006)	0.009* (0.005)	0.018 (0.021)	-0.056 (0.043)
$\mathbb{1}(g = G_c^T)$	0.025* (0.011)	0.035** (0.013)	0.056** (0.018)	0.082 (0.055)
Post-growth $\times \mathbb{1}(g = G_c^T)$	-0.028* (0.012)	-0.020 (0.011)	-0.060* (0.027)	-0.056 (0.074)
N	261830	197226	45884	18720

Notes: Dependent variable is grade-level retention rate. Standard errors, clustered at the district-year level, are reported in parentheses. * denotes significance at 10%, ** denotes 5%, and *** denotes 1%. All regressions include school fixed effects and control for the total enrolled Hispanic students at the school.

Table 6: Effect of Score Change-based Accountability on Retention Rate Among White Students: Texas

	Full sample	Elementary school	Middle school	High school
Post-growth	0.003 (0.001)	-0.001 (0.001)	0.006 (0.004)	0.018 (0.010)
$\mathbb{1}(g = G_c^T)$	0.021*** (0.006)	0.015* (0.007)	0.012 (0.008)	0.083* (0.034)
Post-growth $\times \mathbb{1}(g = G_c^T)$	-0.016** (0.006)	-0.010 (0.006)	-0.004 (0.010)	-0.094** (0.034)
N	309790	237318	51782	20690

Notes: Dependent variable is grade-level retention rate. Standard errors, clustered at the district-year level, are reported in parentheses. * denotes significance at 10%, ** denotes 5%, and *** denotes 1%. All regressions include school fixed effects and control for the total enrolled white students at the school.

6 Costs and benefits of score change-based accountability systems

In this section, I use estimates of the effects of retention from the literature to calculate some back-of-the-envelope costs and benefits associated with the decrease in retention rates caused by score change-based accountability systems. The results of my analyses suggest that retention rates drop on average by around 18% after a state implements a score change-based accountability system, and that the change is more immediate in terminal grades. The total enrollment across grades 3-8 in the U.S. was 22,890,943 for the 2018-19 school year. The average retention rate in a school under a score change-based system was 2.020 in my data. Using these numbers, I calculate that around 457,819 students are retained each year, and that score change-based accountability systems lead to around 100,497 fewer retained students each year. National Center of Education Statistics (2019) reports that 2019 average per-student expenditures in the US were \$13,847. Assuming that a retained student spends an additional year in school relative to the case in which he is not retained, these numbers taken together suggest that the reduction in retentions caused by score change-based accountability lead to around \$1.392 billion saved in student expenditures.

I combine estimates from the literature with various sources to calculate per-student costs and benefits of retention. These calculations and sources can be found in Table 7. Prior research has shown that retention has substantial negative long-term effects on students, especially those retained in later grades.

Based on estimates from Brodaty, Gary-Bobo and Prieto (2013), I calculate that re-

tention decreases the beginning-of-career wage of a student by \$2,682 (in 2019 dollars) due to delayed entry into the job market. If retained students accept lower-paying starting jobs in the future, this could affect their future wages, as well. Deveraux (2002) estimates that about 60% of the wage differential between two individuals that started jobs at the same time could be explained by the difference in the starting wage. Through their effect on retention practices, score change-based accountability systems stop the retention of 100,500 students, and thus around \$270 million in lost beginning-of-career wages, and potentially even more in future wages.

Based on estimates from Eren, Lovenheim and Mocan (2018), I calculate that a student retained in the 8th grade will cause \$338 more in violent crime costs in expectation, and \$564 more in drug offense costs. In total, score change-based accountability systems avoid around 58,000 violent crimes based on figures from Federal Bureau of Investigation (2019) and \$91 million in crime-related costs. Beyond these costs, violent crimes cause long-term socio-emotional problems for their victims, reduce physical activity of neighbors, and lower house values (Langton, 2014; Janke et al., 2016; Taylor, 1995).

Based on estimates from Manacorda (2012) and estimates on the returns to education from Kolesnikova (2010), I calculate that an individual retained in middle school will earn \$1,000-\$4,000 less in yearly wages than he would in the absence of the retention. Other researchers have found a host of negative effects associated with decreased wages, including increased chances of divorce and decrease chances of marrying, decreased access to childcare, and increased levels of obesity and hypertension (Fremstad and Boteach, 2015; Census Bureau, 2011; Leigh, 2013). Through their effect on retention practices, score change-based accountability systems avoid \$101-404 million in lost annual wages, as well as additional costs associated with low wages.

On the other hand, some research has found retention to have academic benefits, particularly in the short-term. Schwerdt, West and Winters (2017) find that retained students enrolled in fewer remedial courses later in life; based on their findings, I calculate that retained students spend around \$4,034 less on remedial reading courses and \$1,121 less on remedial math courses because they were retained. Jacob and Lefgren (2004) find evidence of gains to standardized tests scores among retained 3rd graders, and Schwerdt, West and Winters (2017) find that retention increases high school GPA's. While both of these benefits are real and non-negligible, the financial implications of them are less clear-cut, and I was unable to calculate a concrete number based on these estimates. Through their effect on retention practices, score change-based accountability systems increase spending on remedial reading and math courses by around \$518 million, and may incur additional costs related to reduced GPA and foregone early test score bumps for students.

By suppressing retention, score change-based accountability systems avoid many of the long-term negative consequences of retention that may have otherwise been experienced by around 2.7 million students every year. This effect reduces student expenses by the public school system by around \$1.4 billion, increases annual wages of the would-be-retained by \$100-400 million in total, and avoids \$91 million in crime-related costs. Additional negative externalities may exist beyond these. The change does incur \$518 million in remedial course enrollment fees for students not retained that otherwise would have been.

7 Conclusion

This paper provides evidence that public school administrators use retention as a tool in optimizing their school's rating. They respond to changes in the way the school is evaluated - particularly a switch to school ratings that incorporate year-to-year individual student score change - by changing the rate at which they retain students. My estimates, based on a 7-state sample, suggest that schools retain 18% fewer students after operating under a growth-based system for at least 6 years, and that the effect is more pronounced in the final grade offered by a school, where promotion of a student evicts them from the school's pool of test-takers. I find some suggestive evidence that Black and Hispanic students are more likely to be promoted out of the last grade offered by their middle school after a switch to score change-based accountability, and that boys are more likely to be promoted out of the last grade offered by their elementary school after the switch. These results suggest that administrator expectations about the effects of retention may differ by gender, age, and race. Overall, administrators choose to retain less when student score change matters to their school ratings, particularly when they are able to remove the student from the test-taking pool via promotion, suggesting that administrators expect retention to be harmful to a student's score change in the short term.

Table 7: Estimated effects of retention

Mechanism	Per-student yearly cost/benefit of retention	Source for effect of retention
Decreased beginning-of-career wages	\$2,682 decrease in beginning-of- career wage	Brodaty, Gary- Bobo and Prieto (2013) ¹
Increased probability of committing violent crime as an adult	\$338 more in expected violent crime costs	Eren, Loven- heim and Mo- can (2018) ²
Increased probability of drug conviction	\$564 more in ex- pected drug of- fense costs	Eren, Loven- heim and Mo- can (2018) ³
Decreased educational attainment	\$1,000-\$4,000 decrease in earnings	Manacorda (2012) ⁴
Lower remedial reading course enrollment	\$4,034 less in remedial course expenses	Schwerdt, West and Winters (2017) ⁵
Lower remedial math course enrollment	\$1,121 less in remedial course expenses	Schwerdt, West and Winters (2017) ⁶

¹ Brodaty, Gary-Bobo and Prieto (2013) calculate that a one-year delay into the job market caused by retention decreases beginning-of-career wages by 9%. The number given combines this estimate with the average 2019 entry-level salary of \$32,592 (ZipRecruiter, 2020).

² Eren, Lovenheim and Mocan (2018) estimate a 58.44% increase in probability of committing violent crime as an adult when retained in eighth grade. Federal Bureau of Investigation (2019) report 1,203,808 violent crimes in the U.S. in 2019, and Miller, Cohen and Wiersema (1996) estimate that violent crime imposed annual costs of \$426 million in the US (\$694.49 million in 2019 dollars). These numbers suggest that each violent crime costs \$577 in 2019 dollars. The number given in row two combines this per-crime cost with the estimate given in Eren, Lovenheim and Mocan (2018).

³ Eren, Lovenheim and Mocan (2018) estimate a 10.02% increase in probability of a drug conviction when retained in eighth grade. Olson and Stout (1991) estimate that the cost of investigating and arresting a drug offender in 1989 was \$2,711 (\$5,635 in 2019 dollars). The number given in row three combines this per-offense cost with the estimate given in Eren, Lovenheim and Mocan (2018).

⁴ Kolesnikova (2010) suggests a 10% return to a year of education for practical purposes. Social Security Administration (2020) report the average wage of an American in 2019 was \$51,916, suggesting a rough return of \$5,000 per year of education. The number given in row 4 combines this with the estimate given in Manacorda (2012) of 0.2-0.8-year decrease in educational attainment caused by retention.

⁵ Douglas-Gabriel (2016) reports that the average cost of a remedial course at a four-year institution was around \$3,000 in 2016 (\$3,201.57 in 2019 dollars). The number given in row 5 combines this with the estimate of a 1.26-course decrease in remedial reading course enrollment due to retention given in Schwerdt, West and Winters (2017).

⁶ Douglas-Gabriel (2016) reports that the average cost of a remedial course at a four-year institution was around \$3,000 in 2016 (\$3,201.57 in 2019 dollars). The number given in row 6 combines this with the estimate of a 0.35-course decrease in remedial math course enrollment given in Schwerdt, West and Winters (2017).

References

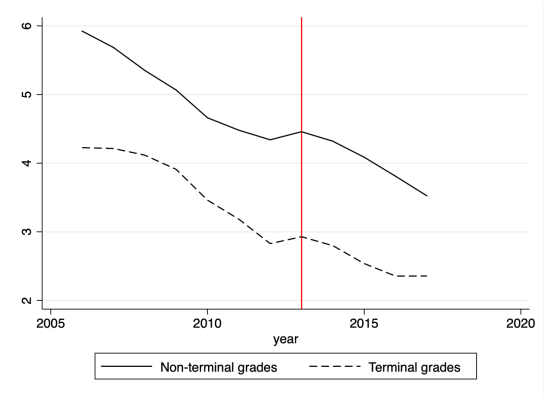
- Abadie, Alberto, Susan Athey, Guido Imbens, and Jeffrey Wooldridge.** 2017. "When Should You Adjust Standard Errors for Clustering?"
- Abraham, Sarah, and Liyang Sun.** n.d.. "Estimating Dynamic Treatment Effects in Event Studies With Heterogeneous Treatment Effects." *Journal of Econometrics*.
- Alliance for Excellent Education.** 2016. "Every Student Succeeds Act Primer: Personalized Learning." Alliance for Excellent Education April, Washington, D.C.
- Almond, Douglas, Ajin Lee, and Amy Ellen Schwartz.** 2016. "Retention Heterogeneity in New York City Schools."
- Brodaty, Thomas O, Robert J Gary-Bobo, and Ana Prieto.** 2013. "Does Speed Signal Ability? The Impact of Grade Retention on Wages."
- Chakrabarti, Rajashri.** 2013. "Vouchers, public school response, and the role of incentives: Evidence from florida." *Economic Inquiry*, 51(1): 500–526.
- Craig, Steven G., Scott A. Imberman, and Adam Perdue.** 2013. "Does it pay to get an A? School resource allocations in response to accountability ratings." *Journal of Urban Economics*, 73(1): 30–42.
- Cullen, Julie Berry, and Randall Reback.** 2006. "Tinkering Toward Accolades: School Gaming under a Performance Accountability System."
- Cunha, Flavio, and James J Heckman.** 2008. "Formulating , Identifying and Estimating the Technology of Cognitive and Noncognitive Skill Formation." *The Journal of Human Resources*, 43(4): 738–782.
- Douglas-Gabriel, Danielle.** 2016. "GradePoint: Remedial classes have become a hidden cost of college." *Washington Post*, , (11/27/2017).
- Eren, Ozkan, Michael F. Lovenheim, and Naci H. Mocan.** 2018. "The effect of grade retention on adult crime: Evidence from a test-based promotion policy."
- Federal Bureau of Investigation.** 2019. "Violent Crime."
- Figlio, David N., and Joshua Winicki.** 2005. "Food for thought: The effects of school accountability plans on school nutrition." *Journal of Public Economics*, 89(2-3): 381–394.
- Fruehwirth, Jane Cooley, Salvador Navarro, and Yuya Takahashi.** 2016. "How the Timing of Grade Retention Affects Outcomes : Identification and Estimation of Time-Varying Treatment Effects." *Jornal of Labor Economics*, 34(4): 979–1021.
- Goodman-Bacon, Andrew.** 2018. "Difference-in-Differences With Variation in Treatment Timing."

- Jacob, Brian A., and Lars Lefgren.** 2004. "Remedial education and student achievement: A regression-discontinuity analysis." *Review of Economics and Statistics*, 86(1): 226–244.
- Jacob, Brian A., and Lars Lefgren.** 2009. "The Effect of Grade Retention on High School Completion." *American Economic Journal: Applied Economics*, 1(3): 33–58.
- Kolesnikova, Natalia.** 2010. "The return to education isn't calculated easily." *Regional Economist*, , (January): 1–12.
- Manacorda, Marco.** 2012. "The cost of grade retention." *Review of Economics and Statistics*, 94(2): 596–606.
- Miller, Ted R., Mark A. Cohen, and Brian Wiersema.** 1996. "Victim Costs and Consequences: A New Look."
- National Center of Education Statistics.** 2019. "Public School Expenditures."
- Oison, D.E., and L.S. Stout.** 1991. "Cost of Processing a Drug Offender Through the Criminal Justice System." Chicago.
- Reback, Randall.** 2008. "Teaching to the rating : School accountability and the distribution of student achievement." *Journal of Public Economics*, 92: 1394–1415.
- Rockoff, Jonah, and Lesley J. Turner.** 2010. "Short-run impacts of accountability on school quality." *American Economic Journal: Economic Policy*, 2(4): 119–147.
- Schwerdt, Guido, Martin R. West, and Marcus A. Winters.** 2017. "The effects of test-based retention on student outcomes over time: Regression discontinuity evidence from Florida." *Journal of Public Economics*, 152: 154–169.
- Social Security Administration.** 2020. "Measures of Central Tendency for Wage Data."
- ZipRecruiter.** 2020. "What is the Average Entry Level Salary by State."

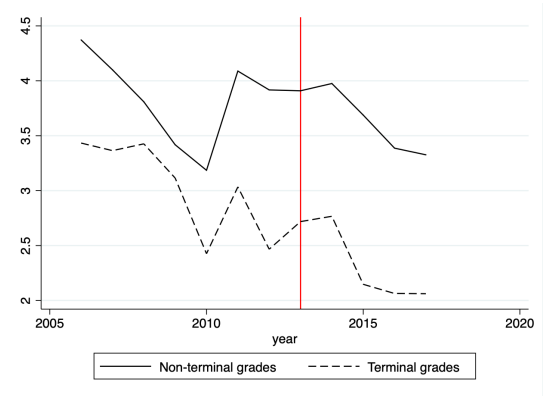
8 Tables & Figures

A Parallel trends figures for effect heterogeneity analysis

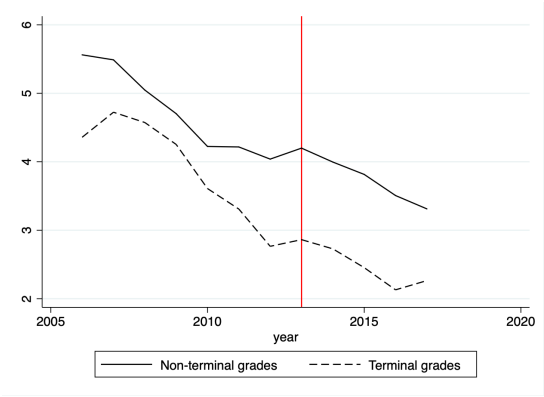
Figure A.1: Parallel trend checks



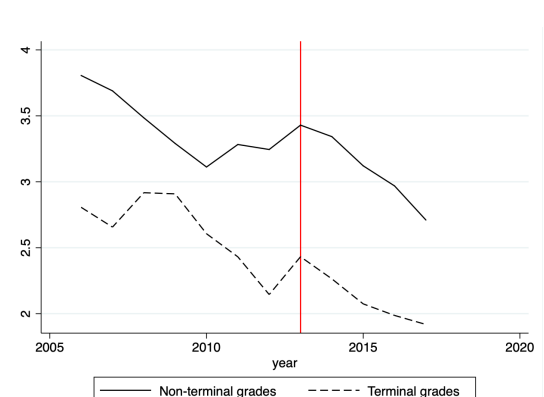
(a) Retention rate among male students



(b) Retention rate among Black students



(c) Retention rate among Hispanic students



(d) Retention rate among white students