

**Effects of Large-Scale Early Math Interventions on Student Outcomes:
Evidence from Kentucky’s Math Achievement Fund**

Zeyu Xu
American Institutes for Research
zxu@air.org

Umut Özek
RAND
uozek@rand.org

Jesse Levin
American Institutes for Research
jlevin@air.org

Dong Hoon Lee
American Institutes for Research
dhlee@air.org

October 25, 2024

Acknowledgement: This work is supported by the National Science Foundation under Grant No. 2000483. The authors have no conflict of interest to declare. The authors would like to thank their research partners Karen Dodd, Erin Chavez, Aaron Butler, and Hannah Poquette from the Kentucky Department of Education; Barrett Ross from the Kentucky Center for Statistics; Kelly DeLong from the Kentucky Center for Mathematics; and Mary Lee Glore from Northern Kentucky University. The authors would also like to thank participants of the SREE conference for valuable feedback. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation, the Kentucky Department of Education, the Kentucky Center for Statistics, the Kentucky Center for Mathematics, or Northern Kentucky University.

Abstract: Addressing the educational needs of students in math early on is critical given that early gaps in math skills widen further over the course of schooling. This study examines the effectiveness and costs of Kentucky’s Math Achievement Fund – a unique state-level program that combines targeted interventions, peer-coaching, and close collaboration among teachers to improve math achievement in grades K-3. The program is found to improve not only math achievement, but also reading test scores and non-test outcomes including student attendance and disciplinary incidents. The benefits exist across students from various socioeconomic backgrounds, and they are slightly higher for racial minorities.

Author bios:

ZEYU XU, PhD, is a managing economist at the American Institutes for Research. His research focuses on college readiness, high school to college transition, and teacher labor market policies.

UMUT ÖZEK, is a senior economist at RAND. His current research focuses on immigrant students and English learners, implementation and consequences of educational accountability, K-12 remediation policies and their effects, and the design and effects of school choice programs.

JESSE LEVIN, PhD, is a principal research economist at American Institutes for Research. His research focuses on the economic evaluation of educational programs and investigating the adequacy and equity of school funding mechanisms

DONG HOON LEE, MPP, is a researcher at American Institutes for Research. His research focuses on inequalities in schools among immigrants and English learners, as well as policy solutions and their effects on these population groups.

Effects of Large-Scale Early Math Interventions on Student Outcomes: Evidence from Kentucky's Math Achievement Fund

Introduction

A long line of research demonstrates the economic value of mathematics skills to both individuals and the economy as a whole (e.g., Behrman, Ross, & Sabot, 2008; Glewwe, 1996; Hershbein & Kearney, 2014). Based on a review of evidence from the United States (e.g., Altonji & Pierret, 2001; Lazear, 2003; Murnane et al., 2000), Hanushek and Woessmann (2008) conclude that one standard deviation increase in mathematics performance at the end of high school translates into 12% higher annual earnings.

Adolescent mathematics performance, in turn, is strongly predicted by students' early-grade mathematics skills (Austin et al., 2021; Goldhaber et al., 2021; Watts, Duncan, Siegler, & Davis-Kean, 2014). Although the association between children's early skills and their later development is present in many domains, such association is twice as strong in mathematics as in reading after controlling for pre-school entry cognitive and socioemotional skills and family background (Duncan et al., 2007). More importantly, children who exhibit early delays in mathematical development also demonstrate lower growth rates in mathematics skills later (Morgan, Farkas & Wu, 2011; Baumert, Nagy, & Lehmann, 2012). Consequently, early gaps in mathematics skills widen further over the course of schooling (Geary et al., 2013; Loeb & Bassok, 2007; Magnuson & Duncan, 2006). How to diagnose and correct mathematics deficiencies early, therefore, is a challenge that our K–12 schools must confront *systematically*. That said, as pointed out in a National Academies report, attention to early mathematics learning

does not match its importance, and reading has been the focus of early childhood education (Cross, Woods & Schweingruber, 2009).¹

In this study, we examine the effects of Kentucky’s Math Achievement Fund (MAF) — a prime example of a state’s effort to intervene early to improve student mathematics achievement— on student outcomes using a difference-in-differences (DiD) design. To the best of our knowledge, MAF was the only fully implemented statewide K-3 mathematics intervention program at the time of this study. Kentucky’s early-grade mathematics performance mirrors the national trend, with roughly half of third-grade students proficient in mathematics in 2017. To help students in grades K–3 who struggle to meet grade-level expectations, MAF adopts a Response to Intervention (RtI) approach that consists of a universal screener, tiered interventions with increasing intensity, and regular progress monitoring. It also requires professional development for both the intervention teacher and classroom teachers in each grantee school. The goal of MAF is to identify mathematics deficiencies or the risks of mathematics difficulties early and provide remedial or preventative interventions before it becomes too late.

Our estimates suggest no statistically significant difference in mathematics achievement trends between MAF and non-MAF schools (i.e., schools that have never implemented MAF) in pre-treatment periods, yet the trends started to diverge gradually after MAF designation. For example, we find no effect of MAF designation on mathematics scores at the end of the first year of implementation, but student mathematics test scores in MAF schools rose by 0.05 standard deviation after 2 years relative to non-MAF schools, 0.08 standard deviation after 3 years, and 0.09 standard deviation after 4 years. Evidence also suggests positive spillover effects of the program on reading achievement, with students in MAF schools outperforming non-MAF school students by 0.06 standard deviation after 4 years. Additionally, MAF does not have any

unintended negative association with student behaviors that may result from labelling young children as not on track to meet grade expectations in mathematics. To the contrary, relative to non-MAF schools, disciplinary incidents and absences reduced significantly in MAF schools relative to non-MAF schools. Using detailed information obtained from administrative records and surveys administered to school MAF teams and coordinators, we find that the estimated four-year per-pupil cost of the MAF intervention is \$750 in 2018-19 dollars, which corresponds to an increase in math test scores by 0.09 standard deviation. As we discuss at the end of this study, closing gaps in mathematics skills early has important economic and equity implications. But even without considering those long-term benefits, these findings suggest that state-led early mathematics intervention programs like MAF could provide a cost-effective blueprint to raise mathematics achievement in early grades.

Policy Background

As provisioned in Kentucky statute KRS 158.844, MAF is a school-level program with two objectives: to provide teacher professional development and to address the needs of students in primary grades (Grades K-3) who are struggling to meet the grade level expectations for mathematics. The intervention is supposed to supplement, not replace, regular classroom instruction. In addition, MAF was also designed to help schools build capacity by overhauling mathematics instructional strategy and practice and strengthening collaboration among teachers. Through renewable, two-year local grants to approximately 90 schools per year, MAF awarded \$50,000 per school-year to support the hiring of one full-time mathematics intervention teacher (MIT) and training for the MIT and two classroom mathematics teachers (known as the Plus 2 teachers). The anticipated cost to operate a MAF program typically exceeded \$50,000, and grantee schools were required to provide additional funds (often from Title I funding) to cover the difference.

MAF requires a common structure of intervention with a set of prescribed components, but grantee schools also have discretion over how specific components are implemented, as described below.

Student Selection

All K-3 students in MAF schools are required to take a screener test to determine who are struggling to meet grade level or benchmark expectations for mathematics and thus need supplemental support. In the 2018–19 school year, MAF schools in Kentucky primarily used 3 screeners: Measures of Academic Progress (MAP) in 55 schools, iReady in 18 schools, and STAR in 15 schools. Screeners were typically administered three times a year. Students were added to and dismissed from intervention services throughout the year based on whether they have met grade-level expectations. Along with screener scores, several factors (e.g., teacher and parental input, scores on alternative tests, need for remediation in reading) were used in the selection of students for the intervention. In the first cohort of MAF schools (schools designated as MAF schools in the first grant cycle that started in 2015–16 school year), the share of students in grades K-3 chosen for MAF intervention ranged from 2 percent to 34 percent (with a median of 13 percent) in 2015–16 school year.²

In addition to universal screeners, easyCBM was used to monitor progress among intervention students. Kindergarten students took the Numbers and Operations Measurement, while students in Grades 1-3 took the Numbers and Operations Measurement and the Numbers and Operations in Algebra Measurement. easyCBM was administered three times a year (October, January, and May), and results were reported in a centralized database.

Requirements for Teachers

The MIT must be a highly trained, certified teacher with at least three years of teaching experience, with preference given to teachers with at least three years of primary teaching experience or training in mathematics intervention services for primary students. At least one-half of the MIT's time is required to be spent delivering mathematics interventions to lower-performing primary grade

students. Administrative records show that these interventions were typically provided in small groups of 4-6 students with each session lasting about 30-45 minutes. Typically, 5-6 sessions were offered every day for 5 days a week. Importantly, these sessions did not replace regular mathematics instruction, a key requirement of the MAF grant. A review of MIT weekly schedules demonstrates that MITs scheduled interventions around homeroom teachers' schedules so that students were not pulled during core instruction. In addition to direct instruction for targeted students, the MIT spent about 5.6 hours per week coteaching mathematics with regular classroom teachers and about 3.8 hours on lesson planning.

The MIT is also expected to engage in a significant number of professional development activities. Based on grantee performance reports, MITs typically participated in 8-12 events per year including webinars, multi-day training sessions, conferences, and courses offered by programs such as Add+VantageMR. The MIT is also expected to share instructional activities, books, and strategies obtained from these professional development opportunities with building staff during weekly team planning and professional learning communities (PLCs). MIT daily schedules suggest that the MITs spent about 2.3 hours per week on PLC activities.

In addition to the MIT, MAF requires grantee schools select two classroom teachers every year as Plus 2 teachers. Plus 2 teachers receive 10 days of professional development and attend, at the expense of the grantee schools, at least one state mathematics conference every year. Plus 2 teachers are required to be available for collaboration and co-teaching with the MIT throughout the school year. Along with the MIT, Plus 2 teachers are expected to lead professional learning with additional teachers to build capacity in the building.

Mathematics Intervention Packages

Grantee schools are required to use one of the four MAF-approved primary mathematics intervention packages: Add+VantageMR® (AVMR), Assessing Math Concepts®, Do the Math®, and Math Recovery Intervention Specialist® (MRIS). These packages are used to provide

interventions for targeted students, professional development for classroom teachers, and professional learning for intervention teachers. In the 2018–19 school year, out of the 97 MAF grantee schools, the overwhelming majority chose to use AVMR (76 schools) and MRIS (17 schools) to facilitate the training of intervention teachers and the delivery of intervention services.

Grantee Selection

Although we do not have details on the selection process such as the total number of applicants, application scoring rubrics, or any information on applicants who were not selected for funding, materials submitted by successful applicants shed light on the dimensions that the Kentucky Department of Education (KDE) valued most when selecting grantees. Specifically, each applicant addressed seven dimensions: identification of need, a description of the RtI framework, identification of which MAF-approved intervention package would be used, how students would be identified for targeted intervention, approaches to professional learning and sustainability, assessment and evaluation plan, and a budget.

Channels of Change

While improving the mathematics achievement of low-performing students is a key objective of MAF, program administrators and grantee performance reports indicate that the defining feature of MAF is its whole school approach. MAF's emphasis on coteaching, professional development for classroom teachers, and the expectation that the MIT and Plus 2 teachers share what they have learned from professional development activities with other teachers in the building suggest that the program could lead to positive spillover effects on other students in the school. To capture the full impact of MAF, this study examines the schoolwide effects of MAF instead of focusing narrowly on the outcomes of students who were referred to interventions.³

MAF could theoretically improve student achievement in mathematics through three channels. The first channel is through increased interaction between children and their teachers.

Because MAF requires that interventions not replace regular instruction, it results in a net gain in instruction time for intervention students. Empirical evidence shows that education interventions that add instruction time tend to improve student achievement significantly. For example, Taylor (2014) and Cortes, Goodman, and Nomi (2015) both used RD frameworks to study the effects of assigning students to an additional mathematics course and found positive effects on test scores in the short run. Similar evidence exists for reading, where Figlio, Holden, and Özek (2018) found significant benefits of additional instruction time using extended school days on student test scores.

The second channel is the delivery of more personalized support through MAF interventions in small group settings. Small groups allow the interventionist to better adapt instruction to individual student needs, and students are thought to learn better and faster if the content is within their zone of proximal development (Guthrie, 2015; Stipek, 2002; Vygotsky, 1978). Small groups are also thought to help establish emotional and peer support quickly, which is particularly important in the context of early childhood development and learning (Mashburn et al., 2008). As we review later, early childhood mathematics interventions delivered in small group settings have consistently produced positive improvements in student mathematics achievement (Pellegrini et al., 2018; Slavin, Lake, Davis, & Madden, 2011; Wanzek et al., 2016).

A third feature of MAF that may help improve student mathematics achievement is setting clear expectations for what the MIT should and should not do while encouraging close collaboration among teachers. This design feature protects the interventionist from distractions such as being asked to substitute for regular mathematics teachers in school. At the same time, the learning community approach to professional development has been demonstrated to

improve teacher performance (e.g., Buysse, Castro, & Peisner-Feinberg, 2010). Close collaboration between the MIT and classroom teachers also allows the interventionist to provide “just in time” support to children who need help with mathematics, and research has shown that working in a collaborative environment greatly improves teaching quality (Bransford, Brown, & Cocking, 1999).

Literature Review

There are hundreds of elementary grade mathematics interventions and more than ten thousand studies on their efficacy. Less than 1% of those studies are considered rigorous enough to be included in several best-evidence syntheses (Jacobse & Harskamp, 2011; Pellegrini et al., 2018; Savelsbergh et al., 2016; Slavin & Lake, 2008). However, the rigor of early mathematics intervention studies has improved in the last decade, and some patterns start to emerge. For example, two types of interventions are found to produce positive results consistently.

The first is tutoring programs that provide one-to-one or small group supplemental instruction, with average effect sizes about 0.30 standard deviations in student test scores.⁴ Tutoring programs typically require extensive professional development for the interventionists and the intervention tends to be intensive, meeting at least 3 times per week for at least 12 weeks, with each session lasting at least 30 minutes. The second type of promising interventions focuses on improving instructional process (such as cooperative learning, classroom management and motivation), with average effect sizes around 0.25 standard deviations in student test scores.⁵ All these programs pay close attention to student needs and behavior, and target instruction to student proficiency levels.⁶ Similar to tutoring programs, teacher professional development is an integral part of instructional process interventions. Peer coaching and teacher collaboration in instructional planning are often featured in teacher professional development. Taken together,

reviews of empirical evidence suggest that early mathematics intervention strategies that substantially affect children’s daily experiences by emphasizing personalization, engagement, and motivation are likely to produce beneficial results (Pellegrini et al., 2018).

Based on this literature, the MAF design—intensive small-group supplemental instruction by full-time MTs, personalized support, teacher peer-coaching, and collaboration between interventionists and regular classroom teachers—is consistent with the characteristics of successful early mathematics interventions. As such, our study complements the existing literature in this context in several important ways.

First, studies that find positive effects on student mathematics learning tend to be small, typically with no more than a few hundred students carried out in a single district. In small-scale studies, developers and experimenters are closely involved, to an extent that is often not realistic if interventions were scaled up. Our study adds to the existing literature by studying whether promising intervention features will scale successfully in a statewide program. Second, large-scale evaluations of elementary mathematics interventions have produced reasonably solid evidence on what interventions do *not* work, but it remains unclear which intervention strategies may scale successfully. This is because large-scale studies, typically RCTs involving thousands of students across multiple districts or states, have been focusing on curricula, benchmark assessments and PD for mathematics content and pedagogy.⁷ To the best of our knowledge, there has been no evaluation of any statewide K-3 mathematics intervention programs. This is primarily because very few states have systematic, coherent early-grade mathematics policies in place. Third, existing mathematics intervention research generally focuses exclusively on mathematics achievement of students who were targeted for intervention, ignoring potential spillover effects on other students and outcomes or unintended adverse effects on student behaviors. Finally, because there are competing demands for attention to deficiencies in multiple domains like reading and socio-emotional competence

(Denton, Germino-Hausken, & West, 2000), our study estimates the cost-effectiveness of MAF to inform the best way to deploy educational resources.

Data

In our impact analysis, we use student-level administrative school records obtained from KDE. These data cover all students in grades K-5 between 2011-12 and 2018-19 school years and include reading and mathematics scores of all tested students (in grades 3 through 8) on the Kentucky Performance Rating for Educational Progress (K-Prep) test along with demographic information on students, such as race, gender, free- or reduced-price lunch eligibility, English language learner status, exceptional/special education status, student age, and schools attended. We also observe information about student disciplinary problems (number of disciplinary referrals, number of suspensions, and the total length of suspensions), and attendance (including days absent, days present, number of unexcused/excused tardy days, and the number of days possible).

Kentucky's administrative data also include detailed information on the interventions (including MAF) that students receive. The reporting of this information is at the student-intervention level, and it is mandatory. MAF intervention activities are populated by the mathematics team in each MAF school and checked three times a year by KDE to ensure high quality and timeliness of data entry. The intervention data include individual student-level information on start and end date of the intervention, intervention program (e.g., AVMR, MRIS), tier of instruction, duration of intervention (in minutes), frequency of intervention (times per week), intervention staff (the qualification level of the staff that is most directly providing the interventions services to the student), areas of student need (deficiency areas such as number

sense, mathematics reasoning, and measurement). We use this information to better understand the type, duration, and intensity of the intervention students receive.

There were two MAF grant cycles between the redesign in 2015-16 and the most recent school year we observe in our data (2018-19 school year). The top panel in Figure 1 presents the number of “ever-MAF” schools (i.e., elementary schools that were ever designated as MAF during this time frame) by MAF status. There were 107 elementary schools in Kentucky that were designated as MAF in the first grant cycle that started in 2015-16 school year. Of these schools, 51 remained as MAF in the second grant cycle that started in 2017-18 whereas 56 left the MAF designation. Thirty-nine schools were newly designated as MAF in the second grant cycle.

Table 1 examines the student characteristics in pre-MAF years (between 2011-12 and 2014-15 school years) for first-cohort MAF schools (schools that were first designated as MAF in 2015-16), second-cohort MAF schools (first time designation in 2017-18), and schools that were never designated as MAF between 2015-16 and 2018-19. The results suggest that the first- and second-cohort MAF schools were comparable along observed student characteristics: students in these schools had similar scores on math tests and were comparable along disciplinary incidents, attendance, race/ethnicity, socioeconomic status (as proxied by subsidized meal eligibility), English learner status, and special education status. That said, students in both first- and second-cohorts of MAF schools had lower math scores, were more likely to be eligible for subsidized meals, more likely to be involved in disciplinary incidents, had more absences, and were more (less) likely to be White (Black) compared to students in schools that were never designated as MAF between 2015-16 and 2018-19 school years.

In our cost analysis, we take a school perspective and therefore limit our data collection only to school staff who dedicated effort to the program. However, there is no formal program requirement or evidence of effort devoted on the part of volunteers or parents, leading us to believe that virtually all resources dedicated to the program are provided by schools. In turn, we use both administrative data and primary data collected through surveys administered to MITs and Plus 2 Teachers who were part of a MAF school intervention team in the 2018-19 school year in order to obtain information on the school resource effort devoted to the MAF program. The survey questionnaire was organized around the activities and personnel/nonpersonnel resources associated with the MAF program.

Survey administration took place in May and June of 2021 and were sent out to the designated MIT and principal, as well as 187 Plus2 teachers (P2Ts) from 97 MAF grantee schools that were active in the 2018-19 school year.⁸ Unfortunately, the survey data collection was conducted in the middle of the COVID-19 pandemic, which likely affected the survey response counts. Specifically, we received no responses from principals and responses from a total of 25 MITs and 8 P2Ts representing 30 schools. Among the 25 MITs who responded, 6 are dropped because they are from schools that are not included in the impact analysis.⁹ Ten of the remaining 19 MITs are from schools that are also first cohort MAF grant recipients, and 9 are from schools that received MAF grants in the second cohort only. Because the response rate among P2Ts is low, we decided to base the calculation of costs associated with these staff solely on extant administrative information pertaining to their required participation in professional development as part of the MAF program. Also, there was only a minimal required effort on the part of principals described in the program documentation.¹⁰ We acknowledge that by not fully capturing efforts associated with P2Ts and principals, the resulting cost estimates may be

considered a lower-bound of the overall program cost. However, the costs associated with these unaccounted efforts are likely to be small.

Empirical Framework

Determining Program Effectiveness

To estimate the effects of MAF on student outcomes, we primarily use a difference-in-differences (DiD) design. We estimate MAF effects for each cohort separately for two reasons. First, it is of policy interest to examine the extent to which program effects may have changed between cohorts. Program effects could change, for example, as the result of program maturing over time; they could also change as the program scales up and the initial enthusiasm wanes. Second, we observe four years of post-award data for the first cohort of MAF grantee schools. This allows us to examine the extent to which program effects, if they exist, persist for over time particularly in schools that have stopped receiving MAF support.

In our main analysis, we focus on MAF schools in the first cohort, comparing the differences in student outcomes in years after versus before the designation (2015-16 school year) with the same difference in never-MAF schools. Formally, using OLS we estimate the following two-way fixed-effects model:¹¹

$$Y_{igst} = \alpha + \beta MAF_s * \delta_t + \delta_t + \mu_s + \gamma_g + \varepsilon_{igst} \quad (1)$$

where Y_{igst} is the outcome of interest (test scores in math and reading standardized to zero mean and unit variance at the grade-by-year level, whether the student received a suspension, and % absent days) for student i in grade g and school s in year t , MAF_s is an indicator for MAF schools in the first cohort, δ_t is year fixed-effects (with 2014-15 school year serving as the omitted category), and μ_s and γ_g are school and grade fixed-effects respectively. Given that MAF is a schoolwide intervention, which means there is no variation within each school/time period in

treatment status among students, we cluster our standard errors at the school-level. In this setting, precisely estimated zero coefficients on the interaction between MAF_s and δ_t for years prior to 2014-15 would provide evidence that the parallel trends assumption (that is, treatment and comparison schools would follow similar trajectories in post-treatment years in the absence of the treatment) cannot be rejected, which is a critical identification assumption in the DiD design.

Panel (B) in Figure 1 checks the fidelity of implementation and examines student exposure to MAF in the first cohort MAF schools and never-MAF schools in the years before and after the MAF designation. In particular, this graph compares the share of students who ever received the MAF intervention over time between these two types of schools. The results indicate that by the end of the 4th year after the initial MAF designation, roughly 15 percent of all students in the first-cohort MAF schools ever received the MAF intervention compared to less than 1 percent in never-MAF schools (nearly all of these ever-MAF students in never-MAF schools had transferred from MAF schools).

Another concern in this context is that the first cohort MAF schools started implementing other interventions at the same time as MAF. For example, during the time frame we examine in this study, Kentucky also started implementing a K-3 reading intervention called Read to Achieve (RTA). If first-cohort MAF schools were more likely to implement RTA simultaneously compared to never-MAF schools, it would become harder to attribute the observed differences in student outcomes between the two types of schools in post-designation years to the causal effect of MAF. Panel (C) in Figure 1 repeats the same exercise in Panel (B), replacing MAF exposure with exposure to other interventions (including RTA). The share of students exposed to other interventions increased in both first-cohort MAF and never-MAF schools after 2014-15 school

year, but students in MAF schools were no more likely to receive another intervention compared to students in never-MAF schools.

Finally, recent developments in two-way fixed effects models suggest that traditional DiD estimators produce biased results when the timing of treatment implementation is staggered and the treatment effects are heterogeneous across groups or time (e.g., Callaway & Sant’Anna, 2021; de Chaisemartin & D’Haultfœuille, 2020). Separating the two cohorts of MAF grantee schools circumvents the complications related to staggered implementation. However, Sun and Abraham (2021) show that the estimated effect for year t could still be contaminated by both the weights and effects from periods other than year t . As robustness checks, we combine the two cohorts of grantees and implement the estimators developed by Callaway and Sant’Anna (2021) using both never-treated and not-yet-treated schools as the comparison. These details are presented after the main findings.

Determining Program Costs and Cost-Effectiveness

In order to gain a well-rounded understanding of whether the MAF program is a wise investment we cannot depend on estimates of impact alone. We must also consider the costs of implementing the program and assess these costs in relation to its effects through a cost-effectiveness analysis, which measures program cost per unit of outcome produced.

Our analysis employs the Ingredients Approach (Levin et al., 2018) to identify the costs of implementing the MAF program. The approach isolates the costs associated with the program by identifying the quantities of the personnel and nonpersonnel resources (ingredients) used in the implementation of the program and assigning corresponding prices to calculate their costs. Based on reviews of documentation and extant data, as well as discussion with MAF administrators, key ingredients involved in implementing of MAF consist of activities related to

intervention services, family engagement, assessment and monitoring, professional development, and administration (Table 2). School staff (MITs and P2Ts) surveys were used to quantify each ingredient. For personnel components, we asked staff to provide hours spent for a given program activity along with auxiliary costs related to activities in question, such as transportation and lodging/food costs for attending required formal professional development. We also collected data on staff years of experience and highest level of education attainment in order to apply appropriate compensation rates in the next step. Data for nonpersonnel resources such as software and equipment utilized during MAF-related activities, we asked respondents about the use of commonly used equipment and software.

The next step in the ingredients approach is to assign prices to the resources. The prices for personnel resources are compensation rates (salary plus benefits), which vary by years of experience and the highest level of education attainment. Importantly, the personnel compensation rates used reflect statewide averages, so that any variation in the subsequent costs across sites reflect differences in the qualifications of staff used and not the influence of local labor markets on the price of staff (see online Appendix A for more details). The unit prices of nonpersonnel resources involved in the MAF program were derived from information posted by major online retailers (e.g., BestBuy, Amazon, and HP).

The final step to calculate costs simply involves multiplying the quantity of each resource used by its corresponding unit price and summing across the resources. The final calculations of the overall program implementation cost for each school sums together the costs of personnel and annualized costs of nonpersonnel resources. As resource allocation data was only collected for the 2018-19 school year, the costs for this year are translated into present values for each of the previous three years and summed to provide a four-year implementation cost for each school

(in 2018-19 dollars) and expressed in per-pupil terms (More details can be found in online Appendix A). We calculate an enrollment-weighted average of the four-year program per-pupil costs across the cohort 1 grantee schools (using school K-5 enrollment as the weight), which serves as the key cost metric for the cost-effectiveness analysis.

This four-year weighted average per-student cost is then coupled with the estimated four-year MAF impacts on student outcomes to generate cost-effectiveness ratios. Each ratio shows the cost per unit of outcome produced by the intervention. For example, the cost-effectiveness of the MAF program with respect to outcome (o) is simply the average per-student cost ($\text{AverageCostPerStudent}_{MAF}$) divided by the impact (μ_o), as follows:

$$\text{Cost-Effectiveness}_o = \frac{\text{AverageCostPerStudent}_{MAF}}{\mu_o} \quad (2)$$

Results

Effectiveness of MAF

Figure 2 presents the event study estimates (along with their 95% confidence intervals) for our four main outcomes of interest: math (Panel A) and reading (Panel B) scores on K-Prep in grades 3 through 5; whether the student was involved in a disciplinary incident in Panel C and percent absent days in Panel D (both estimated separately for grades K-5 and K-3). Table 3 presents these coefficients, estimated without (in columns labeled as (I)) and with (columns labeled as (II)) student covariates (i.e., an indicator for subsidized meal receipt, race/ethnicity, gender, special education status, English learner status, foreign born indicator, and age).

The main takeaway from this analysis is that MAF had a significant positive effect on both test and non-test outcomes of students, especially beyond the first year after the designation. implementation. In particular, while we do not find any concerning evidence of differential pre-

treatment trends between first-cohort MAF and never-MAF schools in the outcomes of interest, we find significant differences in the years after MAF designation. For example, MAF increased student mathematics test scores by 0.05 standard deviation after 2 years, 0.08 standard deviation after 3 years, and 0.09 standard deviation after 4 years. Similar, yet slightly smaller effects emerge for reading scores, with second year effects of 0.03 standard deviation, third year effects of 0.04 standard deviation, and fourth year effects of 0.06 standard deviation. The gradually increasing effects of MAF is in line with (1) more students who received the MAF intervention in grades K-3 being tested in grades 3 through 5 and (2) the gradual implementation of effective practices regarding math instruction throughout the school.

We also find that MAF significantly reduced disciplinary incidents and student absences in participating schools. For example, MAF decreased the likelihood that students were involved in disciplinary incidents by 1 percentage points in the first year (by 25 percent of the dependent variable mean), 2 percentage points in the second year (by roughly 50 percent), and 3 percentage points in the third year (by 75 percent).¹² Similarly, MAF led to a decline in percent absent days of roughly 0.1 to 0.3 percentage points in the first four years (about 2 to 6 percent of the dependent variable mean).¹³ Panels (C) and (D) in Figure 2 also present the results for students in earlier grades (K-3) to see whether the effects were larger in the grades directly targeted by MAF interventions. The results suggest that the K-3 effects are mostly in line with the overall effects on non-test outcomes, which provides suggestive evidence that the effects of MAF designation goes above and beyond the students chosen for the intervention.

Our findings also indicate that the MAF effects persist even after schools leave the program, which provides further evidence that systemic, schoolwide shifts in MAF schools is an important driver behind the observed benefits. In particular, Figure 3 breaks down the analysis in

Figure 2 by whether the school remained as MAF in the second grant cycle (that started in 2017-18 school year) or left the MAF program. The estimated MAF effects are comparable for the two types of schools in the first two years of their designation (2015-16 and 2016-17 school years). What is more interesting is that we find significant benefits of the program on student outcomes in the third and fourth years after their designation when they no longer receive MAF support from the state. For example, for these schools, we find MAF effects of 0.1 standard deviation in math and significant reductions in disciplinary incidents and absences in the fourth year.

Do certain student groups benefit more from MAF? Gender, for example, is widely documented to be associated with mathematics performance and STEM (Science, Technology, Engineering, and Math) choices (e.g., Card & Payne, 2021; Ellison & Swanson, 2010). Gender differences emerge in the early years of elementary school (Penner & Paret, 2008), and MAF may have differential impacts on boys and girls. In addition, research suggests that low-income and racially minority students tend to be more susceptible to being labelled as in need of interventions (e.g., Mokher et al., 2018; Ou, 2010; Papay et al., 2011). Finally, one possible unintended consequence of receiving MAF grants is that resource may be reallocated within school to focus more on grades K-3.

Figures 4A-4D examines this question and breaks down the analysis in Figure 2 by student grade (in tested grades) in Figure 4A, by whether the student received subsidized meals in Figure 4B, by student gender in Figure 4C, and by student race/ethnicity in Figure 4D. Table 4 presents the traditional, static DiD estimates (after MAF versus before MAF, in first-cohort MAF schools versus never MAF schools) obtained using the outcomes of interest in Figure 2 (with K-5 incidents and absence rates) overall and by student gender, subsidized meal receipt, and race/ethnicity. This table also presents the DiD results using an “ever-MAF” indicator that equals

1 if the student ever received the MAF intervention as the outcome to assess whether the MAF designation has a differential effect on the likelihood of receiving the intervention for different student groups.

The overarching conclusion from this analysis is that the benefits of MAF are widespread. In particular, we do not find any differences in MAF effects on outcomes of interest across grades, by student gender, or by student subsidized meal receipt although we find that MAF designation has a significantly larger effect on ever receiving the intervention for students who receive subsidized meals (an effect of 15 percentage points versus 8 for students who do not receive subsidized meals).

That said, the lack of overlap in confidence intervals suggest significant differences between White and non-White (mostly Black and Hispanic students). While MAF designation does not have a differential effect on the likelihood of receiving the intervention for these two student groups, the benefits on test scores and non-test outcomes are significantly larger for racial minorities. For example, we find MAF effects of 0.13 standard deviation in math and 0.10 standard deviation in reading for racial minorities compared to no statistically significant effects for White students. MAF led to declines in disciplinary incidents and absence rates among both racial minority and White students, but the decline in disciplinary incidents is significantly larger among racial minorities (6 percentage points, equivalent to 67 percent of the dependent variable mean for the student group) than among White students (0.7 percentage points, equivalent to 26 percent of the dependent variable mean). The MAF effects on absence rates are not statistically different between racial minorities and White students (a decline of 0.28 percentage points for racial minorities compared to 0.10 percentage points for White students).

What about differential effects by the share of students receiving the MAF intervention? For example, if the observed MAF effect is primarily driven by the effect of the intervention on treated students, then one would expect the effects to be larger in schools with a larger share of students receiving the intervention. The primary challenge in this analysis is that the number of students who directly receive intervention services under MAF in a school is endogenous and it depends partly on the effectiveness of the intervention in prior years (i.e., an effective intervention in the prior year would improve students' mathematics skills in the current year and hence reducing the share of students receiving intervention services). Therefore, in Figure 5 we break down the analysis in Figure 2 by the share of K-3 students who received direct interventions in first-cohort MAF schools in 2015-16 school year (first year of implementation). We find no significant differences in MAF effects for MAF schools with higher (above median) and lower (below median) share of K-3 students identified for MAF in the first year of implementation. This finding once again provides evidence that the policy impacts student outcomes in ways beyond its effect on students directly receiving intervention services.

An alternative explanation for the observed MAF effect is the possible effect of MAF designation on student composition. For example, if the designation leads to lower-performing students leaving these schools or higher-performing students entering, it could lead to improved outcomes in MAF schools on average compared to other schools. In online appendices, Figure 1 examines this possibility and checks the effects of MAF designation on observed student characteristics (kindergarten readiness for students in kindergarten, whether the student receives subsidized meals, race/ethnicity, special education status, English learner status, gender, and immigrant status). Overall, we do not find evidence of shifts in student composition large enough to explain the observed effects of MAF on student outcomes. This is consistent with the finding

presented in Table 3 that shows that the estimated effects of MAF remain virtually unchanged when we include student covariates in the regressions.

Effects for the Second MAF Cohort

How do these estimated compare with the effects on the second cohort of MAF schools? Online appendices Figure 2 repeats the analysis in Panels (B) and (C) in Figure 1 for the schools that were first designated as MAF in 2017-18 school year and online appendices Figure 3 presents the event study estimates obtained using these schools as the treatment group and never-MAF schools as the comparison group. Similar to the first cohort, MAF designation increases exposure to the MAF intervention, and we find no significant differences in exposure to other interventions between first-time second cohort MAF schools and never-MAF schools.

In terms of MAF effects on student outcomes, we find results that are somewhat in line with the findings using the first cohort MAF schools in the first two years after their designation. We find no significant effects on math scores and absence rates after 2 years although it is important to note that these coefficients are less precisely estimated as we have fewer treatment schools in this exercise and there is a slight upward trend in MAF effects on math scores in years after the designation. In contrast, we find significant benefits of MAF designation on reading scores and disciplinary incidents in the second year. In all cases, we find no concerning differences in pre-treatment trends between treatment and comparison groups.

Robustness Checks

In this section, we examine (1) the overall MAF effects across both grantee cohorts and (2) the robustness of the main findings to alternative comparison schools. Combining both MAF cohorts entails staggered treatment timing. Research demonstrates that traditional DiD estimates are biased (and often with little direct policy relevance) when treatment timing is staggered and

treatment effects vary across group and time periods (e.g., de Chaisemartin & D’Haultfœuille, 2020; Goodman-Bacon, 2021). Here, we follow the empirical strategy developed by Callaway and Sant’Anna (2021) that is robust to heterogeneous treatment effects in a staggered setting and estimate an event study regression of the following form:

$$Y_{igst} = \delta_t + \mu_s + \gamma_g + \sum_{l=-K, l \neq -1}^L \beta_l 1\{l_s = t - g_s\} + \varepsilon_{igst} \quad (3)$$

Similar to equation (1), Y_{igst} is the outcome of interest for student i in grade g and school s in year t . δ_t , μ_s and γ_g are year, school, and grade fixed-effects. Because implementation timing now varies by grantee cohort, the term $MAF_s * \delta_t$ in equation (1) is replaced by a series of relative time binary variables $1\{l_s = t - g_s\}$ that represent the l th year since school s first received MAF grants in year g_s . Thus, $l = 0$ for the year in which a school first received MAF. As is conventional and consistent with equation (1), the year immediately before MAF designation ($l = -1$) is the reference period and omitted. In other words, the treatment-control difference in outcomes for each period $l \neq -1$ is compared with the treatment-control difference in ($l = -1$). When $l \geq 0$, β_l estimates the cumulative effect $l + 1$ years after initial MAF designation. When $l < -1$, β_l estimates the placebo effect $|l|$ periods before MAF designation.

This event study model was estimated using two alternative sets of comparison schools: Never-MAF schools and not-yet-MAF schools. Our main analysis has so far used never-MAF schools as the comparison. However, those schools could be inherently different from schools that applied and received MAF grants in ways that are unobservable to researchers. As a result, not-yet-MAF schools were used as an alternative comparison group with the plausible assumption that these schools likely resembled schools that had received MAF grants earlier more closely than never-MAF schools. With only two cohorts of MAF schools, using not-yet-

MAF schools as the comparison in essence uses cohort 2 MAF schools as the control group for cohort 1 MAF schools.

These estimates are reported in Table 5 as Column III (never-MAF as the comparison) and IV (not-yet-MAF as the comparison). The primary DiD estimates for cohort 1 and cohort 2 are reproduced in Columns I and II for comparison purposes. Results obtained using Callaway and Sant’Anna estimators are highly consistent with the main findings. For math, the Callaway-Sant’Anna estimates are 0.01, 0.05, 0.06, and 0.08 standard deviation for the first four post-MAF years, respectively, when never-MAF schools are the comparison. These are nearly identical to the main DiD results reported for Cohort 1 MAF schools (that is, 0.00, 0.05, 0.08, and 0.09 for year t , $t+1$, $t+2$, and $t+3$, respectively). The coefficients are statistically significant at the end of the second, third, and fourth year since MAF was first implemented. When not-yet-MAF schools are the comparison, the estimate effects are limited to the first two post-MAF years, and they are 0.02 and 0.06 standard deviation. These are also in line with the main results, albeit statistically insignificant as the sample size is only 15% as large as the sample size used for the main results. Similarly, the estimated effects for other outcomes are robust to choices of samples, comparison groups, and estimators.

Cost-Effectiveness Estimates

The following section reports the results of both the cost and cost-effectiveness analyses. The findings provided for the main cost analysis are based on those first-cohort MAF schools that were two-time grant recipients and therefore implemented the program for four years from 2015-16 to 2018-19. The findings for the cost-effectiveness analysis are based on all schools in the main impact analysis sample, which includes both two- and four-year grant recipients that

received their initial MAF grants in 2015-16. As reported earlier, MAF impact estimates are very similar between grantees schools that received grants for two and four years (see Figure 3).

The overall four-year per-pupil cost of MAF is \$750 (Figure 6). As is the case with most education programs, most of the cost (\$726 or 96%) is attributable to personnel resources. The breakout of costs by program activities in Figure 7 shows that over half of the overall program cost (\$431 or 57%) is dedicated to providing direct intervention services to low-performing students, while just over a quarter (\$194 or 26%) is spent on professional development for program teachers. Smaller shares of the program cost are associated with assessing students and monitoring their progress (\$72 or 10%), program administration (\$25 or 3%), and engaging families (\$28 or 4%).

Figure 8 shows considerable variation in program cost across the grantee schools. We divide MAF grantee schools into terciles based on their per-pupil costs. It depicts the average per-pupil costs for low-, medium-, and high-cost MAF schools, as well as how these costs break out across the program activities.¹⁴ The average program cost per pupil for MAF schools in the top tercile is \$1,062, nearly twice as much as the per pupil cost for MAF schools in the bottom tercile (\$604). The cost breakout shows that the additional cost associated with the high-cost group can be attributed to greater spending across all activities, but especially on intervention services and administration.

Figure 9 illustrates how per-pupil program costs in 2018-19 vary with the scale at which MAF schools operate.¹⁵ It shows a clear pattern consistent with economies of scale where the per-student cost of schools operating with a smaller enrollment tends to be higher than those with larger enrollments. For instance, the schools with fewer than 300 K-5 students has annual per-

pupil costs around \$350, while the per-pupil costs of those with more than 600 students are close to \$150.

By combining the cost estimates and the estimated math achievement impact we find that the MAF program has a cost-effectiveness ratio of \$8,069 in terms the targeted outcome of interest, student achievement in mathematics. The cost-effectiveness ratio is an estimate, and therefore it is bounded by a confidence interval within which the true value lies with some degree of certainty. Using the Monte Carlo interval method documented in Dong et al., (2023), we find that the 95%-confidence interval of the estimated cost-effectiveness ratio to be between \$5,038 and \$19,735.¹⁶

We also searched the literature for comparable cost-effectiveness findings on interventions aimed at improving student math achievement outcomes. Table 6 reports cost-effectiveness ratios on math achievement scores from alternative interventions reported in Yeh (2010), Foster et al. (2013), and Barrett & VanDerHeyden (2020), which have been inflation adjusted to reflect 2018-19 dollars so that they can be compared to the MAF program cost-effectiveness ratio reported for math achievement. The comparison reveals that the MAF program CER estimate is less cost-effective than programs aimed at improving teacher quality such as the Appalachian Math and Science Partnership (Foster et al., 2013) and teachers obtaining a Master's degree (Yeh, 2010). While the point estimate of the MAF program suggests that it is more cost-effective than alternative programs aimed at improving classroom quality such as offering more rigorous math classes (Yeh, 2010) and class size reduction (Finn et al., 2001 and Nye et al., 2001), the ratios for these programs fall within the MAF 95%-confidence interval denoted by the shaded rows in Table 6 suggesting that their cost-effectiveness does not differ statistically from that of MAF. Note that without estimates of the confidence intervals for

the other interventions with lower cost-effectiveness ratios it is impossible to determine whether they differ statistically from that of the MAF program.

Concluding Remarks

Children starting their K–12 education today are more diverse in their knowledge and skills than before (Diamond et al., 2013), and K–12 schools are called upon to make up for deficiencies in foundational competencies that many entering children display due to their preschool experiences. Gaps in early-grade mathematics proficiency are stubbornly difficult to close, and early delays in mathematical development leads to lower growth rates in mathematics skills later. The MIT from one of the grantee schools wrote:

“The reason we wrote the grant originally was because of our fall to fall data. It always showed no growth, and actually seemed to increase in need as the students got older. We were doing a great job of working really hard to reduce the number of tier 3 students from fall to spring, but only to have an increase in the number of students needing math interventions steadily climb as students got older.”

Other MITs shared similar frustrations. How a state should design a coherent program that remedies and prevents mathematics development delays before they become intractable is an urgent question for state education policymakers to consider.

Findings from this study confirm that key elements of successful early mathematics interventions—small-group instruction, personalized support for students, peer-coaching, and close collaboration among teachers—can be scaled up moderately without losing effectiveness. To put the size of the MAF effect into perspective, the 4th year effect on mathematics achievement is equivalent to a class size reduction of 4-5 students (Angrist & Lavy, 1999) or replacing all novice teachers with teachers with 3-5 years of experience (Xu, Özek, & Hansen, 2015). But at a per-pupil cost of \$750 for four years of intervention, MAF could be considerably cheaper than these other policy alternatives.

Importantly, the MAF benefits were found across grade levels and sustained even after grantee schools stopped receiving MAF support. One possible explanation for the sustained impact is the combination of targeted support for K-3 students and attention to improving mathematics instruction in regular classrooms. Researchers have pointed out that mathematics interventions—even when successful—must be paired with classroom practices to maintain the gains in mathematics skills achieved during the interventions (Clements et al., 2013; Smith et al., 2013; and Watts, Duncan, Clements & Sarama, 2017). Because of the critical role that post-intervention learning environments play, it is essential to improve regular classroom teaching even if the sole objective of a program is to help students who are struggle with mathematics. In addition, all MAF grantees committed in their applications to continued support for at least some aspects of the program even after the funding expired. For example, many grantees proposed to use funds from Title I, extended school services, and the general fund to continue the provision of mathematics intervention packages and to support professional learning. A few grantees also planned to keep the role of MIT after MAF grant expiration.

MAF produced similar benefits for all students regardless of sex, subsidized meal status, or race/ethnicity. For some outcomes, the impact is larger for students who are eligible for subsidized meals and for racial minority students. These findings, if replicable in more schools and states, may hold the key to addressing inequalities in later life outcomes. This is because efforts to address achievement gaps in secondary schools tend to be too little, too late.¹⁷ Some of the proposed policy options are likely met with more challenges in practice than early-grade interventions.¹⁸

The extent to which a program like MAF can be further scaled up remains an open question. The success of MAF, in our view, is partially attributable to the management of the

program. With just over 90 grantee schools (fewer than 1/10 of all primary schools in Kentucky), MAF was closely managed by an experienced mathematics teacher at KDE. A centralized database was used to collect and monitor the number of intervention students and their progress multiple times a year. Individual MITs were contacted whenever discrepancy or data anomaly occurred. Before each school year, proposed daily schedules for all MITs were collected and approved. At the end of the school year, surveys were administered to collect feedback from MITs about their teaching and learning experiences. And at the end of each grant cycle, grantees that wished to renew were evaluated using a 14-point rubric. This level of involvement will be challenging when more schools are involved, and the quality of management will likely vary when multiple program managers are needed.

Even at the current scale, grantee schools reported several implementation challenges. These include difficulty in recruiting qualified MITs, turnovers of MIT and Plus 2 teachers, insufficient number of training slots, and difficulty in creating an intervention schedule that did not conflict with any core classes. These issues are likely to become more acute if more schools were to be included in programs like MAF.

In addition to how well programs like MAF could be further scaled up, future research should also investigate how accurately existing K-3 mathematics screeners reflect and predict mathematics skills development. Some MITs reported concerns that students identified for mathematics intervention might have deficiency in language instead of mathematic skills. Although we find no research evidence for early-grade mathematics screeners, a study on interim assessments used for reading screening in North Carolina concluded that its K-2 screeners did not adequately identify students at risk of scoring below proficient on the state reading assessment at the end of grade 3 (Koon et al., 2020). There is a need for more research so that

districts can make informed choices of early-grade mathematics screeners that are based on their reliability, validity, and classification accuracy (Petscher et al., 2019).

Notes

1. Although nearly half of all states have specific requirements for early literacy assessments and interventions, only three states at the time of writing had policies that mandate additional mathematics support for K–3 students based on screening test results (Georgia Kindergarten Inventory of Developing Skills, Maine’s pilot Numeracy4ME, and Kentucky’s Mathematics Achievement Fund which is the focus of the current study). A handful of other states require numeracy screening before Grade 3, but these states do not stipulate interventions based on test results.

2. MAF excludes students with Individualized Educational Plans (IEPs) related to math goals. Students with IEPs receive specialized instruction from the special education teachers.

3. Ideally, we would also want to estimate the extent to which MAF improves the outcomes of targeted students. The use of universal screeners in grantee schools lends support to a potential regression discontinuity design. Unfortunately, an examination of program implementation indicates that screener scores were one of the many inputs used in intervention referrals. An empirical analysis of screener scores confirms that the designated cutoffs on universal screeners are not binding, and there is no significant discontinuity in treatment probability at the cutoffs.

4. Examples of tutoring programs include *Mathematics Recovery* (Smith, Cobb, Farran, Cordray & Munter, 2013), *Galaxy Mathematics* (Fuchs et al., 2013), *Pirate Mathematics* (Fuchs et al., 2010), *Number Rockets* (Gersten et al., 2015), and *FocusMATH* (Styers & Baird-Wilkerson, 2011).

5. Examples of instructional process programs include *Team Assisted Individualization (TAI)* (Stevens & Slavin, 1995; Slavin & Karweit, 1985), *PAX Good Behavior Game* (Weis, Osborne, & Dean, 2015), and *Individualizing Student Instruction* (Connor et al., 2018).

6. By contrast, interventions that focus on mathematics curricula (textbooks), benchmark assessment and teacher professional development for mathematics knowledge or pedagogy have generally produced little effect on early grade mathematics learning (Pellegrini et al., 2018). For example, *Mathematics Solutions* and *Cognitively Guided Instructions*, two widely used mathematics professional development programs that focus on improving teachers' mathematics knowledge and pedagogy, are found to have no discernible effect on student mathematics learning (Jacob, Hill, & Corey, 2017; Schoen, LaVenita, & Tazar, 2018). Mathematics curricula are found to have no detectable relationship with 4th- and 5th-grade students' mathematics achievement growth (Blazar et al., 2019).

7. These studies generally show null program effects (e.g., Garet et al., 2016; Randel et al., 2016; Newman et al., 2012).

8. Not all the second cohort schools were eligible for the analysis. Out of the 97 schools, 4 schools joined in the second year of the grant cycle, 4 schools did not serve students in grades 3-5 and had no outcome data.

9. Three of the six schools only serve grades K through 2 and therefore have no state standardized test scores, two schools received off-cycle grants (i.e., they joined MAF in the second year of the second MAF grant cycle), and one school does not have complete historical student performance data that are needed to estimate the MAF impact.

10. Specifically, program documentation does not detail any effort on the part of principals outside of a single 6-hour training session on mathematics and cognitive coaching that occurs in the second year of the grant cycle.

11. Examples of studies that use repeated cross-sectional data at the individual level to estimate group-level effects include Bertrand & Mullainathan (2004) and Rouse, et al. (2013).

12. It is important to note that this decline could be driven by two factors: (1) change in student behavior and/or (2) increased leniency among teachers and school administrators for student misbehavior after MAF designation.

13. There may be concerns that multiple comparisons will increase the likelihood of false positives. In online appendices Tables 2 and 3, we present the p -values that corresponds to the comparisons presented in tables 3 and 4, respectively, along with p -values adjusted for comparing multiple outcomes using two methods: sharpened false discover rate q -values and Bonferroni correction. The q -values are the preferred adjustment because Bonferroni is known to increase the probability of false negatives (e.g., Anderson, 2008; Menyhart et al., 2021).

14. Specifically, we calculated the K-5 enrollment-weighted averages of per-pupil program costs for the seven least costly, six most costly schools, and six medium-cost MAF schools from which resource allocation data was collected.

15. The 19 schools are composed of 10 first-cohort MAF schools and 9 second-cohort MAF schools. The 9 second-cohort schools were plotted in Figure 9 because we included them in the imputation model for the 89 first-cohort schools with missing cost data. However, the 9 second-cohort MAF schools were not included in the impact analysis sample which consists only of the first-cohort schools (2015-16 school year).

16. Estimation of the 95%-confidence interval bounding the point estimate of the MAF program mathematics achievement cost-effectiveness was facilitated by the tool developed as part of the study by Dong, et al. (2023).

17. Transition interventions, for example, are designed to help secondary school students get ready for college. Even though 39 states have offered transition intervention programs as of 2017, empirical evidence suggests that they had no effect on helping students succeed in college.

18. For example, academic tracking (e.g., only advanced mathematics students can take algebra 2 in 8th grade) is shown to be related to within-school segregation of racial/ethnic groups (e.g., Clotfelter et al., 2021). But tracking, especially in mathematics, reflects differential skills accumulation over many years. Policies that call for opening access to advanced secondary mathematics courses to all students raise concerns about instructional challenges, diminished rigor, and extra workload for underprepared students to simultaneously learn advanced topics and close preexisting gaps in mathematics skills. Addressing gaps in mathematics skills early appears to be a more sensible policy choice.

References

- Altonji, J. G., & Pierret, C. R. (2001). Employer Learning and Statistical Discrimination. *Quarterly Journal of Economics*, 116(1): 313–50.
- Anderson, M. L. (2008). Multiple inference and gender differences in the effects of early intervention: A reevaluation of the Abecedarian, Perry Preschool, and Early Training Projects. *Journal of the American statistical Association*, 103(484), 1481-1495.
- Austin, W., Figlio, D., Goldhaber, D., Hanushek, E., Kilbride, T., Koedel, C., ... & Strunk, K.O. (2021). *Academic Mobility in U.S. Public Schools: Evidence from Nearly 3 Million Students*. CALDER Working Paper No. 227-0821-2.
- Authors. (2015). Blinded per journal guidelines.
- Authors. (2018). Blinded per journal guidelines.
- Barrett, C. & VanDerHeyden, A. (2020). A cost-effectiveness analysis of classwide math intervention, *Journal of School Psychology*, 80, 54-65.
- Baumert, J., Nagy, G., & Lehmann, R. (2012). Cumulative advantages and the emergence of social and ethnic inequality: Matthew effects in reading and mathematics development within elementary schools?. *Child Development*, 83(4), 1347-1367.
- Behrman, J. R., Ross, D. R., & Sabot, R. (2008). Improving Quality versus Increasing the Quantity of Schooling: Estimates of Rates of Return from Rural Pakistan. *Journal of Development Economics*, 85(1–2): 94–104.
- Bertrand, M., Duflo, E., & Mullainathan, S. (2004). How much should we trust differences-in-differences estimates?. *The Quarterly Journal of Economics*, 119(1), 249-275.
- Blazar, D., Heller, B., Kane, T., Polikoff, M., Staiger, D., Carrell, S.,...& Kurlaender, M. (2019). *Learning by the Book: Comparing math achievement growth by textbook in six Common*

- Core states*. Research Report. Cambridge, MA: Center for Education Policy Research, Harvard University.
- Bransford, J. D., Brown, A. L., & Cocking, R. R. (1999). The design of learning environments. *How people learn: Brain, mind, experience, and school*, 119-142.
- Buyse, V., Castro, D.C., and Peisner-Feinberg, E. (2010). Effects of a Professional Development Program on Classroom Practices and Outcomes for Latino Dual Language Learners. *Early Childhood Research Quarterly*, 25(1), 94-206.
- Callaway, B., & Sant’Anna, P. H. (2021). Difference-in-differences with multiple time periods. *Journal of Econometrics*, 225(2), 200–230.
- Card, D., & Payne, A. A. (2021). High school choices and the gender gap in STEM. *Economic Inquiry*, 59(1), 9-28.
- Chambers, J. (1999). Measuring Resources in Education: From Accounting to the Resource Cost Model Approach. Working Paper No. 1999-16. U.S. Department of Education. National Center for Education Statistics.
- Chambers, J., & Parrish, T. (1994). Developing a resource cost database. In W.S. In Barnett (Ed.), *Cost analysis for education decisions: Methods and examples* (Vol. 4, pp. 7-21). Greenwich, CT: JAI.
- Clements, D. H., Sarama, J., Wolfe, C. B., & Spitler, M. E. (2013). Longitudinal evaluation of a scale-up model for teaching mathematics with trajectories and technologies: Persistence of effects in the third year. *American Educational Research Journal*, 50(4), 812–850.
- Clotfelter, C. T., Ladd, H. F., Clifton, C. R., & Turaeva, M. R. (2021). School segregation at the classroom level in a southern ‘New destination’ state. *Race and Social Problems*, 13(2), 131–160.

- Cortes, K. E., Goodman, J. S., & Nomi, T. (2015). Intensive math instruction and educational attainment long-run impacts of double-dose algebra. *Journal of Human Resources*, 50(1), 108-158.
- Cross, C. T., Woods, T. A., & Schweingruber, H. E. (2009). *Mathematics learning in early childhood: Paths toward excellence and equity*. National Academies Press.
- de Chaisemartin, C., & d'Haultfoeuille, X. (2020). Two-way fixed effects estimators with heterogeneous treatment effects. *American Economic Review*, 110(9), 2964–2696.
- Diamond, K.E., Justice, L.M., Siegler, R.S., & Snyder, P.A. (2013). *Synthesis of IES Research on Early Intervention and Early Childhood Education*. (NCSE 2013-3001). Washington, DC: National Center for Special Education Research, Institute of Education Sciences, U.S. Department of Education.
- Dong, N., Maynard, R., Kelcey, B., Spybrook, J., Li, W., Bowden, A. B., & Pham, D. (2023). Advantages of Monte Carlo Confidence Intervals for Incremental Cost-Effectiveness Ratios: A Comparison of Five Methods. Manuscript submitted to *Journal of Research on Educational Effectiveness*.
- Duncan, G. J., Dowsett, C. J., Claessens, A., Magnuson, K., Huston, A. C., Klebanov, P., ... & Sexton, H. (2007). School readiness and later achievement. *Developmental psychology*, 43(6), 1428.
- Ellison, G., & Swanson, A. (2010). The gender gap in secondary school mathematics at high achievement levels: Evidence from the American Mathematics Competitions. *Journal of Economic Perspectives*, 24(2), 109-128.

- Figlio, D., Holden, K. L., & Ozek, U. (2018). Do students benefit from longer school days? Regression discontinuity evidence from Florida's additional hour of literacy instruction. *Economics of Education Review*, 67, 171-183.
- Finn, J. D., S. B. Gerber, C. M. Achilles, & J. Boyd-Zaharias. 2001. "The enduring effects of small classes." *Teachers College Record* 103: 145-183.
- Foster, J., Toma, E. & Troske, S. (2013). Does Teacher Professional Development Improve Math and Science Outcomes and Is It Cost Effective? *Journal of Education Finance*, 38(3), 255-275.
- Geary, D.C., Hoard, M. K., Nugent, L., & Bailey, H. D. (2013). Adolescents' functional numeracy is predicted by their school entry number system knowledge. *PLoS ONE*, 8(1): e54651.
- Glewwe, P. (1996). The relevance of standard estimates of rates of return to schooling for education policy: A critical assessment. *Journal of Development economics*, 51(2), 267-290.
- Goldhaber, D., Wolff, M., & Daly, T. (2021). Assessing the Accuracy of Elementary School Test Scores as Predictors of Students' High School Outcomes. CALDER Working Paper No. 235-0821-2
- Guthrie, G. (2015). The formalistic education paradigm in Papua New Guinea. *Contemporary PNG Studies*, 22, 33.
- Hanushek, E., & Woessmann, L. (2008). The role of cognitive skills in economic development. *Journal of Economic Literature*, 46(3), 607-668.
- Hershbein, B., Harris, B., & Kearney, M. (2014). *Major decisions: Graduates' earnings growth and debt repayment*. Washington, DC: Hamilton Project, Brookings Institution.

http://www.hamiltonproject.org/assets/files/major_decisions_graduates_earnings_growth_debt_repayment.pdf.

Jacobse, A. E., & Harskamp, E. G. (2011). *A meta-analysis of the effects of instructional interventions on students' mathematics achievement*. Groningen: GION, Gronings Instituut voor Onderzoek van Onderwijs, Opvoeding en Ontwikkeling, Rijksuniversiteit Groningen.

Koon, S., Foorman, B., & Galloway, T. (2020). *Identifying North Carolina students at risk of scoring below proficient in reading at the end of grade 3* (REL 2020–030). U.S. Department of Education, Institute of Education Sciences, National Center for Education Evaluation and Regional Assistance, Regional Educational Laboratory Southeast.
<https://files.eric.ed.gov/fulltext/ED607256.pdf>.

Lazear, Edward P. 2003. "Teacher Incentives." *Swedish Economic Policy Review*, 10(3): 179–214.

Levin, H., McEwan, P., Belfield, C., Bowden, B., Shand, R., (2018). Economic Evaluation in Education: Cost-Effectiveness and Benefit-Cost Analysis. Thousand Oaks, CA: Sage Publications.

Loeb, S., & Bassok, D. (2007). Early childhood and the achievement gap. *Handbook of research in education finance and policy*, 517-534.

Mashburn, A.J., Pianta, R.C., Hamre, B.K., Downer, J.T., Barbarin, O.A., Bryant, D., Burchinal, M. Early, D.M., and Howes, C. (2008). Measures of Classroom Quality in Pre-Kindergarten and Children's Development of Academic, Language, and Social Skills. *Child Development*, 79(3), 732-749.

- Menyhart, O., Weltz, B., & Györfy, B. (2021). MultipleTesting. com: A tool for life science researchers for multiple hypothesis testing correction. *PLoS One*, 16(6), e0245824.
- Mokher, C. G., Leeds, D. M., & Harris, J. C. (2018). Adding it up: How the Florida College and Career Readiness Initiative impacted developmental education. *Educational Evaluation and Policy Analysis*, 40(2), 219–242.
- Morgan, P. L., Farkas, G., & Wu, Q. (2011). Kindergarten children's growth trajectories in reading and mathematics: Who falls increasingly behind?. *Journal of learning disabilities*, 44(5), 472-488.
- Murnane, R. J., Willett, J. B., Duhaldeborde, Y. & Tyler, J. (2000). How Important Are the Cognitive Skills of Teenagers in Predicting Subsequent Earnings. *Journal of Policy Analysis and Management*, 19(4): 547–68.
- Nye, B., L. V. Hedges, and S. Konstantopoulos. 2001. "Are effects of small classes cumulative? Evidence from a Tennessee experiment." *Journal of Educational Research*, 94, 336-345.
- Ou, D. (2010). To leave or not to leave? A regression discontinuity analysis of the impact of failing the high school exit exam. *Economics of Education Review*, 29(2), 171–186.
- Papay, J. P., Murnane, R. J., & Willett, J. B. (2011). How performance information affects human-capital investment decisions: The impact of test-score labels on educational outcomes (NBER Working Paper No. 17120). National Bureau of Economic Research.
- Pellegrini, M., Lake, C., Inns, A., & Slavin, R. E. (2018, October). Effective programs in elementary mathematics: A best-evidence synthesis. In *annual meeting of the Society for Research on Educational Effectiveness*, Washington, DC.
- Penner, A. M., & Paret, M. (2008). Gender differences in mathematics achievement: Exploring the early grades and the extremes. *Social Science Research*, 37(1), 239-253.

- Petscher, Y., Fien, H., Stanley, C., Gearin, B., Gaab, N., Fletcher, J.M., & Johnson, E. (2019). *Screening for Dyslexia*. Washington, DC: U.S. Department of Education, Office of Elementary and Secondary Education, Office of Special Education Programs, National Center on Improving Literacy.
- Rouse, C. E., Hannaway, J., Goldhaber, D., & Figlio, D. (2013). Feeling the Florida heat? How low-performing schools respond to voucher and accountability pressure. *American Economic Journal: Economic Policy*, 5(2), 251-281.
- Savelsbergh, E. R., Prins, G. T., Rietbergen, C., Fechner, S., Vaessen, B. E., Draijer, J. M., & Bakker, A. (2016). Effects of innovative science and mathematics teaching on student attitudes and achievement: A meta-analytic study. *Educational Research Review*, 19, 158–172.
- Slavin, R. E., & Lake, C. (2008). Effective programs in elementary mathematics: A best-evidence synthesis. *Review of Educational Research*, 78(3), 427–515.
- Slavin, R. E., Lake, C., Davis, S., & Madden, N. A. (2011). Effective programs for struggling readers: A best-evidence synthesis. *Educational Research Review*, 6(1), 1–26.
- Smith, T. M., Cobb, P., Farran, D. C., Cordray, D. S., & Munter, C. (2013). Evaluating math recovery: Assessing the causal impact of a diagnostic tutoring program on student achievement. *American Educational Research Journal*, 50(2), 397–428.
- Stipek, D. (2002). At what age should children enter kindergarten? A question for policy makers and parents. *Society for Research in Child Development Policy Report*, 16, 13–16.
- Sun, L., & Abraham, S. (2021). Estimating dynamic treatment effects in event studies with heterogeneous treatment effects. *Journal of Econometrics*, 225(2), 175–199.

- Taylor, E. (2014). Spending more of the school day in math class: Evidence from a regression discontinuity in middle school. *Journal of Public Economics*, 117, 162–181.
- Wanzek, J., Vaughn, S., Scammacca, N., Gatlin, B., Walker, M. A., & Capin, P. (2016). Meta-analyses of the effects of tier 2 type reading interventions in grades K-3. *Educational Psychology Review*, 28(3), 551–576.
- Watts, T. W., Clements, D. H., Sarama, J., Wolfe, C. B., Spitler, M. E., & Bailey, D. H. (2017). Does Early Mathematics Intervention Change the Processes Underlying Children's Learning?. *Journal of Research on Educational Effectiveness*, 10(1), 96–115.
- Watts, T. W., Duncan, G. J., Siegler, R. S., & Davis-Kean, P. E. (2014). What's past is prologue: Relations between early mathematics knowledge and high school achievement. *Educational Researcher*, 43(7), 352-360
- Xu, Z., Özek, U., & Hansen, M. (2015). Teacher performance trajectories in high-and lower-poverty schools. *Educational Evaluation and Policy Analysis*, 37(4), 458-477.
- Yeh, S. S. (2010). The cost effectiveness of 22 approaches for raising student achievement. *Journal of Education Finance*, 36(1), 38-75.

Tables and Figures

Table 1. Student Characteristics by School MAF Status, Elementary Schools between 2011-12 and 2014-15 School Years

	School in First MAF Cohort	School in Second MAF Cohort	Never MAF School
Math scores (Grades 3-5)	-0.019 (0.974)	-0.019 (0.960)	0.008 (1.005)
KG ready	0.351 (0.477)	0.339 (0.473)	0.356 (0.479)
% absent days	5.052 (4.741)	4.682** (4.383)	4.641** (4.566)
Disciplinary incident	0.034 (0.181)	0.037 (0.189)	0.029 (0.167)
Received subsidized meals	0.661 (0.473)	0.655 (0.475)	0.602** (0.490)
Special education	0.113 (0.316)	0.109 (0.311)	0.103** (0.304)
English learner	0.028 (0.165)	0.026 (0.159)	0.035 (0.183)
White	0.882 (0.323)	0.887 (0.316)	0.828** (0.378)
Black	0.070 (0.255)	0.067 (0.250)	0.116*** (0.320)
Hispanic	0.058 (0.234)	0.059 (0.235)	0.064 (0.245)
Asian	0.008 (0.090)	0.008 (0.091)	0.015*** (0.121)
Foreign born	0.001 (0.034)	0.001 (0.038)	0.003*** (0.055)
Age	8.771 (1.890)	8.808 (2.026)	8.882** (1.892)
Number of unique students	80,947	32,591	436,972
Number of unique schools	107	39	720

Notes: Standard deviations are given in parentheses. Math scores are standardized at the grade-year level to zero mean and unit variance. *, **, and *** represent that the means for the corresponding group are statistically different than the first column (students in the first MAF cohort of schools) at 10, 5, and 1 percent, respectively.

Table 2. Key MAF Program Activities and their Definitions Accounted for in the Cost Analysis.

Program Activity	Definition
Intervention Services	Delivery of MAF program instruction to students including preparing lessons, providing direct instruction, and collaborating on instruction with other teachers.
Family Engagement	Engagement with parents through conferences and events.
Assessment and Monitoring	Administration of assessments, identification of MAF students, and monitoring of student progress.
Professional Development	Participation in professional development opportunities related to the MAF program, including formal training provided by Kentucky Center for Mathematics (KCM) and provision of professional development and coaching by Mathematics Intervention Teacher to other staff within schools. ¹
Administration	Planning and budgeting, scheduling service delivery, and entering data on service provision and student progress.

¹ See description of training on KCM website: <https://www.kentuckymathematics.org/maf.php>.

Table 3. Effects of MAF Designation on Student Outcomes: Event Study Estimates

Treatment: first cohort MAF schools, Comparison: never-MAF								
	Math scores		Reading scores		Disciplinary incidents		% absent days	
	(I)	(II)	(I)	(II)	(I)	(II)	(I)	(II)
2011-12 SY	0.035 (0.026)	0.030 (0.025)	-0.003 (0.017)	-0.008 (0.017)	-0.006 (0.005)	-0.006 (0.005)	-0.105* (0.062)	-0.102* (0.062)
2012-13 SY	0.017 (0.024)	0.013 (0.023)	-0.007 (0.017)	-0.011 (0.016)	-0.000 (0.007)	-0.000 (0.007)	-0.043 (0.059)	-0.036 (0.058)
2013-14 SY	0.024 (0.017)	0.023 (0.017)	0.016 (0.012)	0.015 (0.012)	-0.006 (0.004)	-0.005 (0.004)	-0.049 (0.053)	-0.043 (0.051)
2015-16 SY	0.007 (0.017)	0.004 (0.017)	-0.005 (0.013)	-0.008 (0.013)	-0.009*** (0.003)	-0.009*** (0.003)	-0.148*** (0.057)	-0.137** (0.058)
2016-17 SY	0.053** (0.024)	0.053** (0.024)	0.030* (0.017)	0.030* (0.017)	-0.019*** (0.004)	-0.019*** (0.004)	-0.099* (0.055)	-0.093* (0.055)
2017-18 SY	0.075*** (0.026)	0.075*** (0.025)	0.043** (0.019)	0.041** (0.019)	-0.027*** (0.005)	-0.027*** (0.005)	-0.216*** (0.069)	-0.215*** (0.068)
2018-19 SY	0.099*** (0.028)	0.093*** (0.028)	0.065*** (0.021)	0.057*** (0.021)	-0.022*** (0.005)	-0.022*** (0.005)	-0.271*** (0.072)	-0.273*** (0.072)
Mean of Y	0.005	0.005	0.006	0.007	0.039	0.040	4.755	4.753
N	1,123,838	1,120,885	1,123,838	1,120,884	2,383,336	2,368,845	2,351,178	2,337,473
School FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Grade FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Student covariates	No	Yes	No	Yes	No	Yes	No	Yes

Notes: Robust standard errors, clustered at the school level, are given in parentheses. The estimates represent the estimated coefficient on the interaction term given with 2014-15 school year (the year before the policy took effect) serving as the baseline category. Student covariates include an indicator for subsidized meal receipt, race/ethnicity, gender, special education status, English learner status, foreign born indicator, and age. *, **, and *** represent statistical significance at 10, 5, and 1 percent, respectively.

Table 4. Effects of MAF Designation on Receiving MAF Services and Student Outcomes, DiD Estimates, Overall and by Subgroup

Treatment: first cohort MAF schools, Comparison: never-MAF							
	All	Female	Male	Subsidized meals	No subsidized meals	White	Non-White
Ever received MAF services	0.125*** (0.005)	0.134*** (0.005)	0.117*** (0.005)	0.147*** (0.005)	0.078*** (0.005)	0.125*** (0.005)	0.128*** (0.008)
90% CI	[0.118, 0.133]	[0.125, 0.143]	[0.109, 0.125]	[0.138, 0.155]	[0.071, 0.086]	[0.117, 0.133]	[0.115, 0.142]
Mean of Y	0.012	0.013	0.011	0.016	0.006	0.013	0.009
Math scores	0.040** (0.020)	0.036* (0.020)	0.044** (0.021)	0.043** (0.020)	0.025 (0.026)	0.025 (0.020)	0.128*** (0.034)
90% CI	[0.007, 0.072]	[0.003, 0.068]	[0.009, 0.078]	[0.010, 0.077]	[-0.017, 0.068]	[-0.008, 0.059]	[0.079, 0.185]
Mean of Y	0.005	0.006	0.005	-0.242	0.414	0.068	-0.312
Reading scores	0.031* (0.016)	0.022 (0.017)	0.037** (0.017)	0.024 (0.018)	0.032* (0.019)	0.019 (0.016)	0.104*** (0.035)
90% CI	[0.005, 0.057]	[-0.005, 0.050]	[0.009, 0.066]	[-0.005, 0.053]	[0.001, 0.064]	[-0.007, 0.046]	[0.046, 0.162]
Mean of Y	0.007	0.082	-0.065	-0.231	0.401	0.081	-0.363
Disciplinary incident	-0.016*** (0.003)	-0.009*** (0.002)	-0.023*** (0.004)	-0.023*** (0.004)	-0.005*** (0.002)	-0.007*** (0.002)	-0.058*** (0.009)
90% CI	[-0.021, -0.011]	[-0.012, -0.006]	[-0.029, -0.016]	[-0.029, -0.017]	[-0.008, -0.002]	[-0.011, -0.004]	[-0.074, -0.043]
Mean of Y	0.040	0.017	0.060	0.054	0.015	0.031	0.087
% absent days	-0.133*** (0.049)	-0.161*** (0.050)	-0.105* (0.054)	-0.158*** (0.060)	-0.104** (0.045)	-0.098** (0.050)	-0.275*** (0.085)
90% CI	[-0.214, -0.052]	[-0.244, -0.078]	[-0.194, -0.017]	[-0.257, -0.060]	[-0.179, -0.029]	[-0.181, -0.016]	[-0.415, -0.135]
Mean of Y	4.753	4.759	4.747	5.494	3.474	4.822	4.392
N (tested grades)	1,120,885	549,092	571,789	698,038	422,847	935,134	185,751
N (all grades)	2,368,845	1,149,404	1,219,418	1,486,883	881,962	1,986,014	382,831

Notes: Robust standard errors, clustered at the school level, are given in parentheses. The estimates represent the estimated coefficient on the interaction term given with 2014-15 school year (the year before the policy took effect) serving as the baseline category. All regressions control for school, year, and grade fixed-effects, and the student covariates including an indicator for subsidized meal receipt, race/ethnicity, gender, special education status, English learner status, foreign born indicator, and age. *, **, and *** represent statistical significance at 10, 5, and 1 percent, respectively.

Table 5. Effects of MAF Designation on Student Outcomes: Robustness Checks

	Math scores				Reading scores			
	(I)	(II)	(III)	(IV)	(I)	(II)	(III)	(IV)
t-6		0.046 (0.050)	-0.034* (0.020)			0.021 (0.030)	-0.019 (0.019)	
t-5		0.019 (0.048)	-0.003 (0.025)			0.012 (0.031)	-0.006 (0.024)	
t-4	0.030 (0.025)	0.010 (0.036)	-0.009 (0.012)	0.016 (0.023)	-0.008 (0.017)	-0.006 (0.027)	-0.002 (0.010)	0.011 (0.020)
t-3	0.013 (0.023)	0.015 (0.036)	0.005 (0.015)	0.018 (0.029)	-0.011 (0.016)	0.009 (0.023)	0.018 (0.013)	0.037 (0.027)
t-2	0.023 (0.017)	-0.001 (0.025)	-0.015 (0.014)	-0.041* (0.024)	0.015 (0.012)	-0.002 (0.016)	-0.011 (0.010)	-0.036 (0.022)
t	0.004 (0.017)	0.031 (0.026)	0.005 (0.014)	0.017 (0.025)	-0.008 (0.013)	0.021 (0.019)	-0.004 (0.011)	0.006 (0.021)
t+1	0.053** (0.024)	0.046 (0.035)	0.045** (0.020)	0.064 (0.041)	0.030* (0.017)	0.044** (0.022)	0.027** (0.013)	0.040 (0.027)
t+2	0.075*** (0.025)		0.064** (0.025)		0.041** (0.019)		0.033* (0.018)	
t+3	0.093*** (0.028)		0.082*** (0.028)		0.057*** (0.021)		0.047** (0.020)	
N	1,120,885	1,021,085	1,179,373	162,253	1,120,884	1,021,083	1,179,371	162,252
Pre-trend test (χ^2)			7.814	3.527			9.026	4.424
p-value			0.452	0.317			0.340	0.219

	Disciplinary incidents				% absent days			
	(I)	(II)	(III)	(IV)	(I)	(II)	(III)	(IV)
t-6		0.002 (0.007)	0.011 (0.009)			-0.038 (0.104)	0.054 (0.073)	
t-5		0.013 (0.008)	0.000 (0.008)			0.023 (0.104)	0.060 (0.077)	
t-4	-0.006 (0.005)	0.013* (0.007)	0.002 (0.004)	-0.005 (0.010)	-0.102* (0.062)	0.095 (0.096)	0.023 (0.046)	-0.027 (0.084)
t-3	-0.000 (0.007)	0.006 (0.006)	-0.002 (0.004)	-0.005 (0.010)	-0.036 (0.058)	-0.037 (0.084)	-0.007 (0.045)	-0.045 (0.080)
t-2	-0.005 (0.004)	0.011*** (0.004)	0.001 (0.003)	0.011 (0.008)	-0.043 (0.051)	-0.045 (0.080)	0.112 (0.072)	0.161* (0.087)
t	-0.009*** (0.003)	-0.009* (0.005)	-0.007*** (0.003)	-0.014*** (0.005)	-0.137** (0.058)	-0.037 (0.084)	-0.156** (0.077)	-0.135 (0.083)

t+1	-0.019*** (0.004)	-0.010* (0.006)	-0.014*** (0.003)	-0.014** (0.006)	-0.093* (0.055)	-0.097 (0.090)	-0.165** (0.075)	-0.144 (0.097)
t+2	-0.027*** (0.005)		-0.022*** (0.005)		-0.215*** (0.068)		-0.300*** (0.105)	
t+3	-0.022*** (0.005)		-0.018*** (0.006)		-0.273*** (0.072)		-0.421*** (0.114)	
N	2,368,845	2,153,398	2,501,438	364,741	2,337,473	2,122,664	2,469,065	362,101
Pre-trend test (χ^2)			13.218	2.548			6.462	3.551
p-value			0.105	0.467			0.596	0.314

Notes: Robust standard errors, clustered at the school level, are given in parentheses. The year immediately before the policy took effect (t-1) serves as the baseline category. All models include school, year, and grade fixed effects and student covariates including an indicator for subsidized meal receipt, race/ethnicity, gender, special education status, English learner status, foreign born indicator, and age. Column (I) includes cohort 1 only. Column (II) includes cohort 2 only. Columns (III) and (IV) include both cohorts, with never-MAF schools and not-yet-MAF schools as the comparison group, respectively. *, **, and *** represent statistical significance at 10, 5, and 1 percent, respectively.

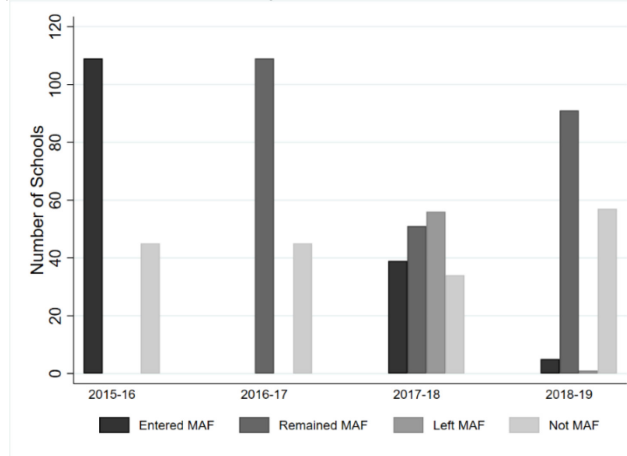
Table 6. Relative Cost-Effectiveness of Interventions on Math Achievement

Interventions	Effect Sizes (standard deviation)	Cost Per Student (in 2019)	CER (cost per 1 standard deviation)
Classwide math intervention ^a	0.18	\$255	\$1,415
Appalachian Math and Science Partnership (AMSP) training ^b	0.03	\$57	\$1,889
Obtaining a master's degree for teachers ^c	0.22	\$898	\$4,082
MAF – Monte Carlo lower bound _{95%-CI}	NA	NA	\$5,038
MAF – Point estimate	0.09	\$750	\$8,069
Students taking more rigorous math courses ^c	0.20	\$2,443	\$12,215
Class size reduction, TN STAR Project 1 ^d	0.13	\$1,763	\$13,667
Class size reduction, TN STAR Project 2 ^e	0.09	\$1,763	\$19,589
MAF – Monte Carlo upper bound _{95%-CI}	NA	NA	\$19,735

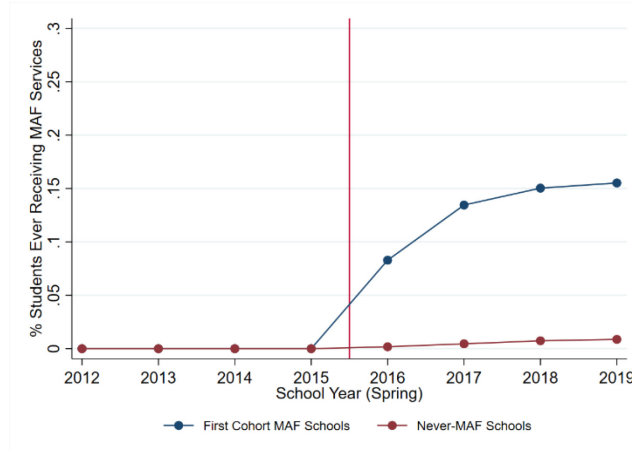
Notes: ^a Barrett & VanDerHayden (2020), ^b Foster et al. (2013), ^c Yeh (2010), ^d Finn et al. (2001), ^e Nye et al. (2001). Shaded rows represent upper and lower bounds of 95%-confidence interval around MAF program cost-effectiveness ratio point estimate.

Figure 1. Exposure to MAF at the School and Student-Level

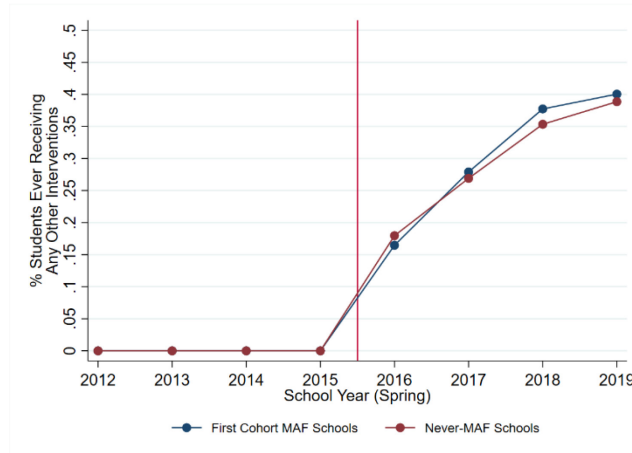
(A) Number of schools by MAF Status: Ever-MAF Schools



(B) % of ever-MAF students: First-cohort MAF versus never-MAF schools

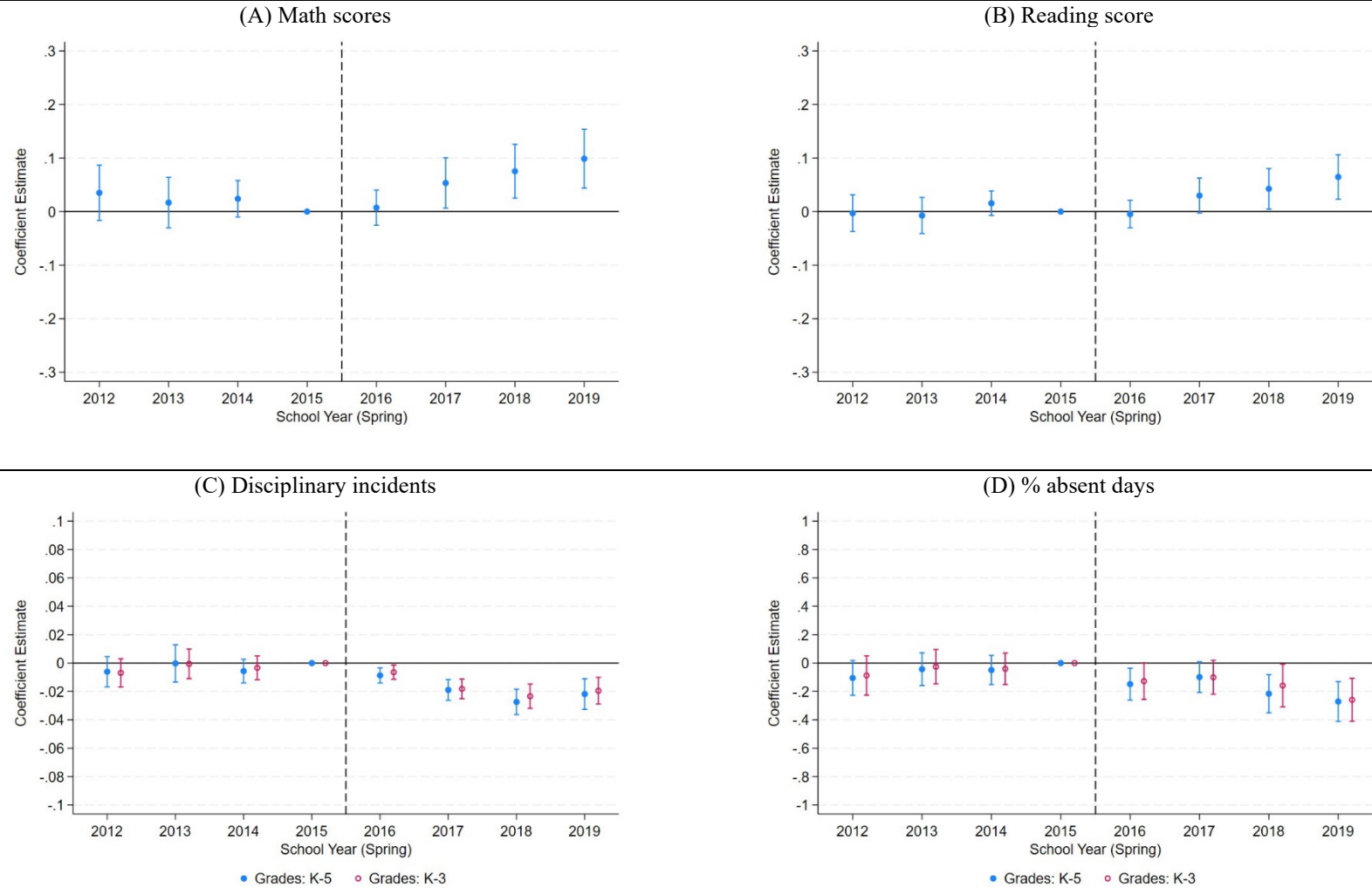


(C) % of students ever-identified for other services: First-cohort MAF versus never-MAF schools



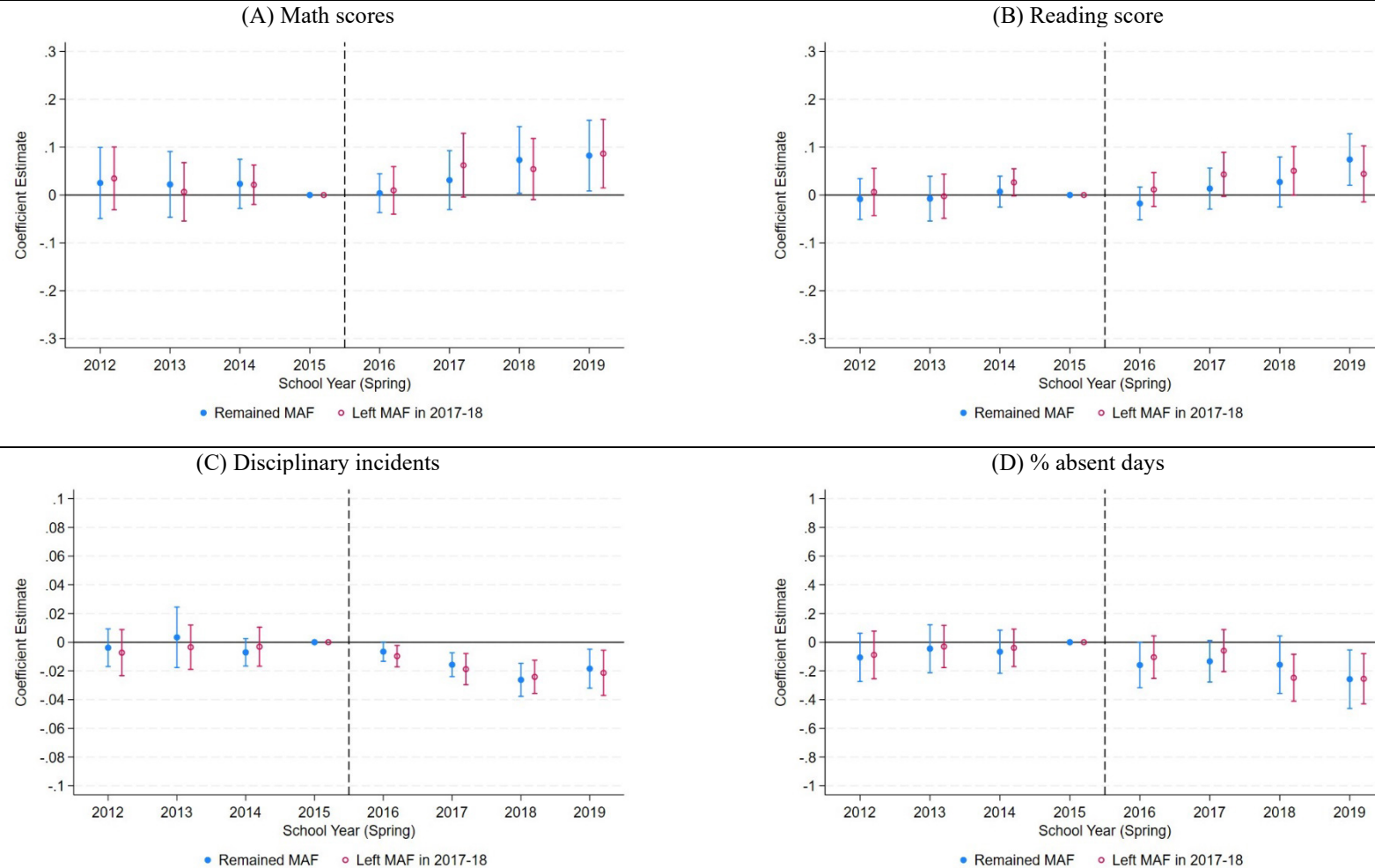
Notes: The top panel presents the number of elementary schools that were ever identified as MAF after 2015-16 by MAF status (entered, remained, left, and not MAF) by school year between 2015-16 and 2018-19 school years whereas the bottom two panels presents the percentage of students in first-cohort MAF schools (elementary schools that were first designated as MAF in 2015-16 school year) and never-MAF schools who have ever received MAF services (panel B) or other services (panel C) up to (and including) that school year.

Figure 2. Effects of MAF Designation on Student Outcomes: Event Study Estimates



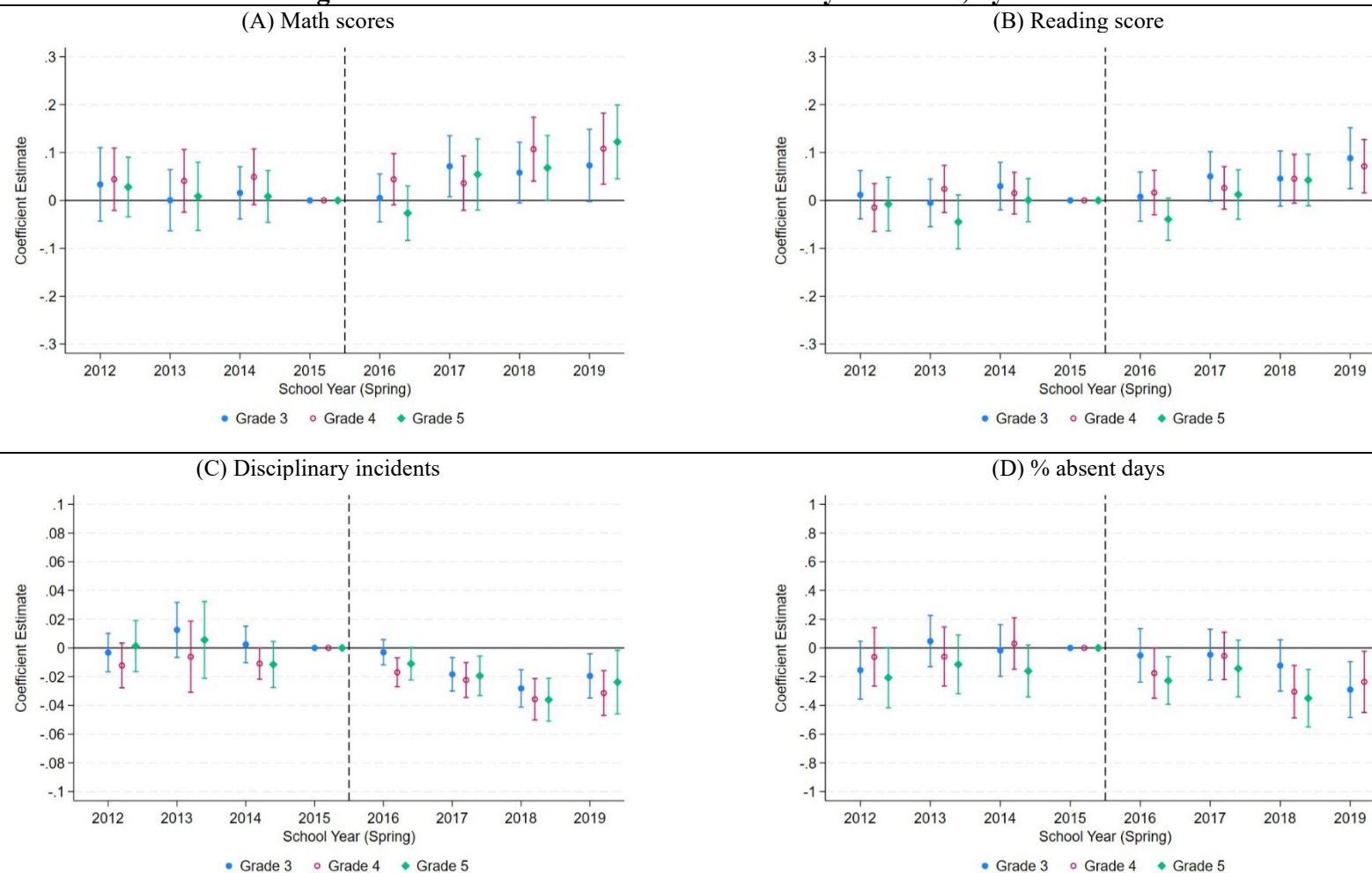
Notes: Each dot presents the estimated coefficient on the interaction term between the MAF indicator and the indicator for the corresponding school year, with the spikes providing the 95% confidence intervals. The estimates are obtained using schools in the first cohort of MAF schools as the treatment and never MAF schools as the comparison group with school, grade, and year fixed effects. Standard errors are clustered at the school level.

Figure 3. Effects of MAF Designation on Student Outcomes: Event Study Estimates, by School Duration in MAF



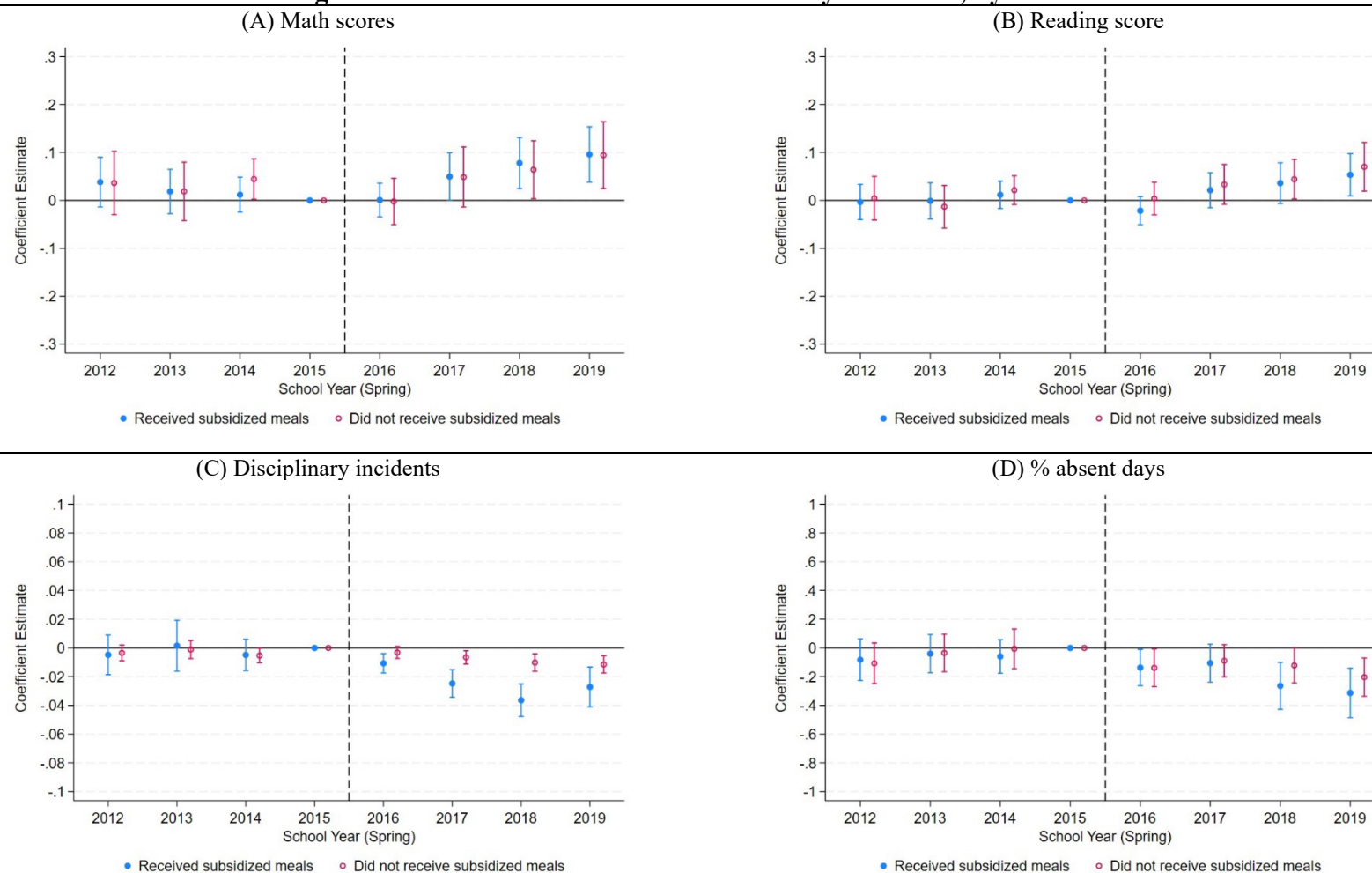
Notes: Each dot presents the estimated coefficient on the interaction term between the MAF indicator and the indicator for the corresponding school year, with the spikes providing the 95% confidence intervals. The estimates are obtained using schools in the first cohort of MAF schools that left MAF at the end of 2017-18 and those that remained as MAF until 2018-19 as the treatment groups and never-MAF schools as the comparison group with school, grade, and year fixed-effects. Standard errors are clustered at the school level.

Figure 4A. Effects of MAF Designation on Student Outcomes: Event Study Estimates, by Grade



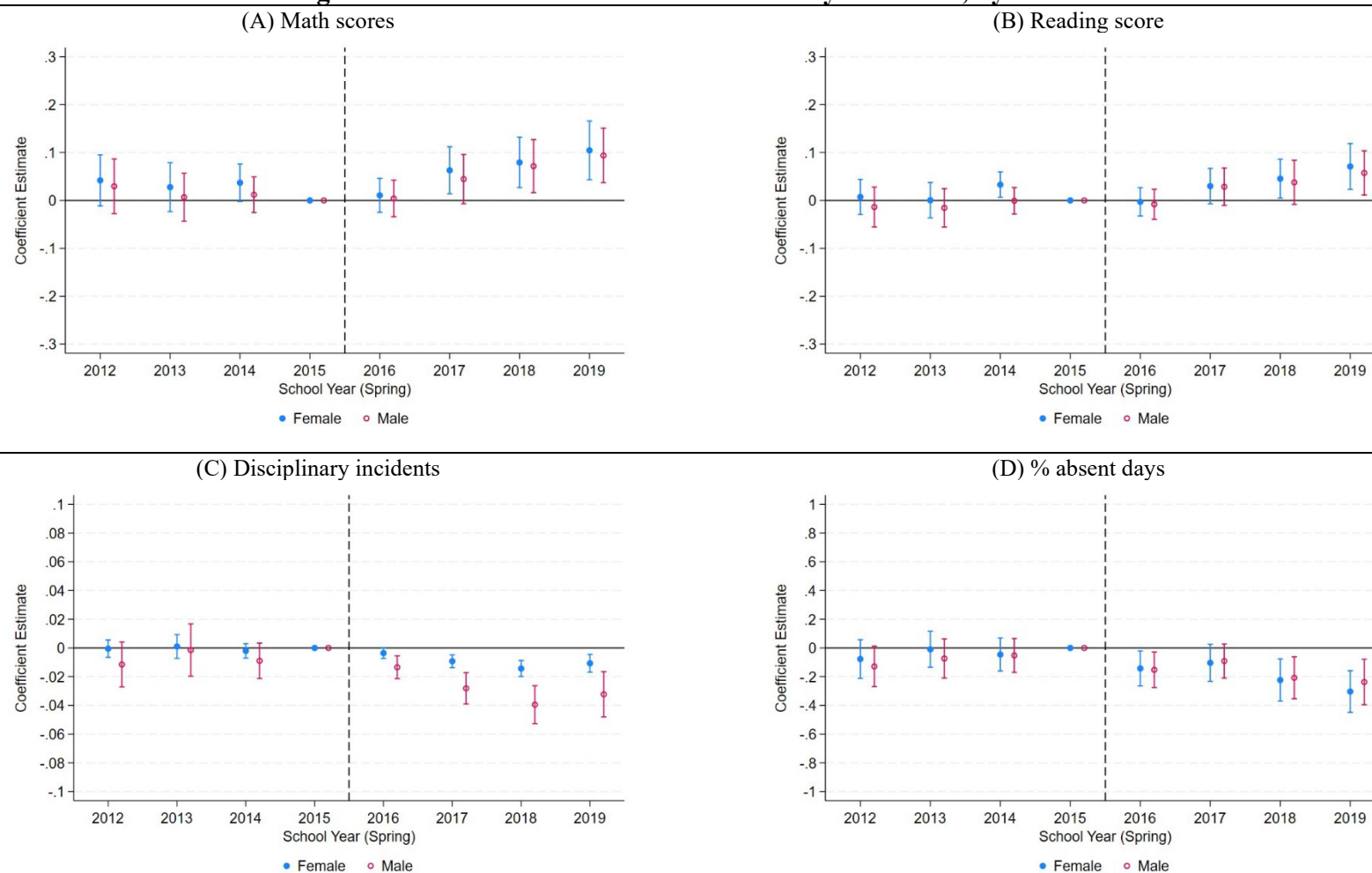
Notes: Each dot presents the estimated coefficient on the interaction term between the MAF indicator and the indicator for the corresponding school year, with the spikes providing the 95% confidence intervals. The estimates are obtained using schools in the first cohort of MAF schools as the treatment and never MAF schools as the comparison group with school, grade, and year fixed-effects. Standard errors are clustered at the school level.

Figure 4B. Effects of MAF Designation on Student Outcomes: Event Study Estimates, by Student Subsidized Meal Receipt



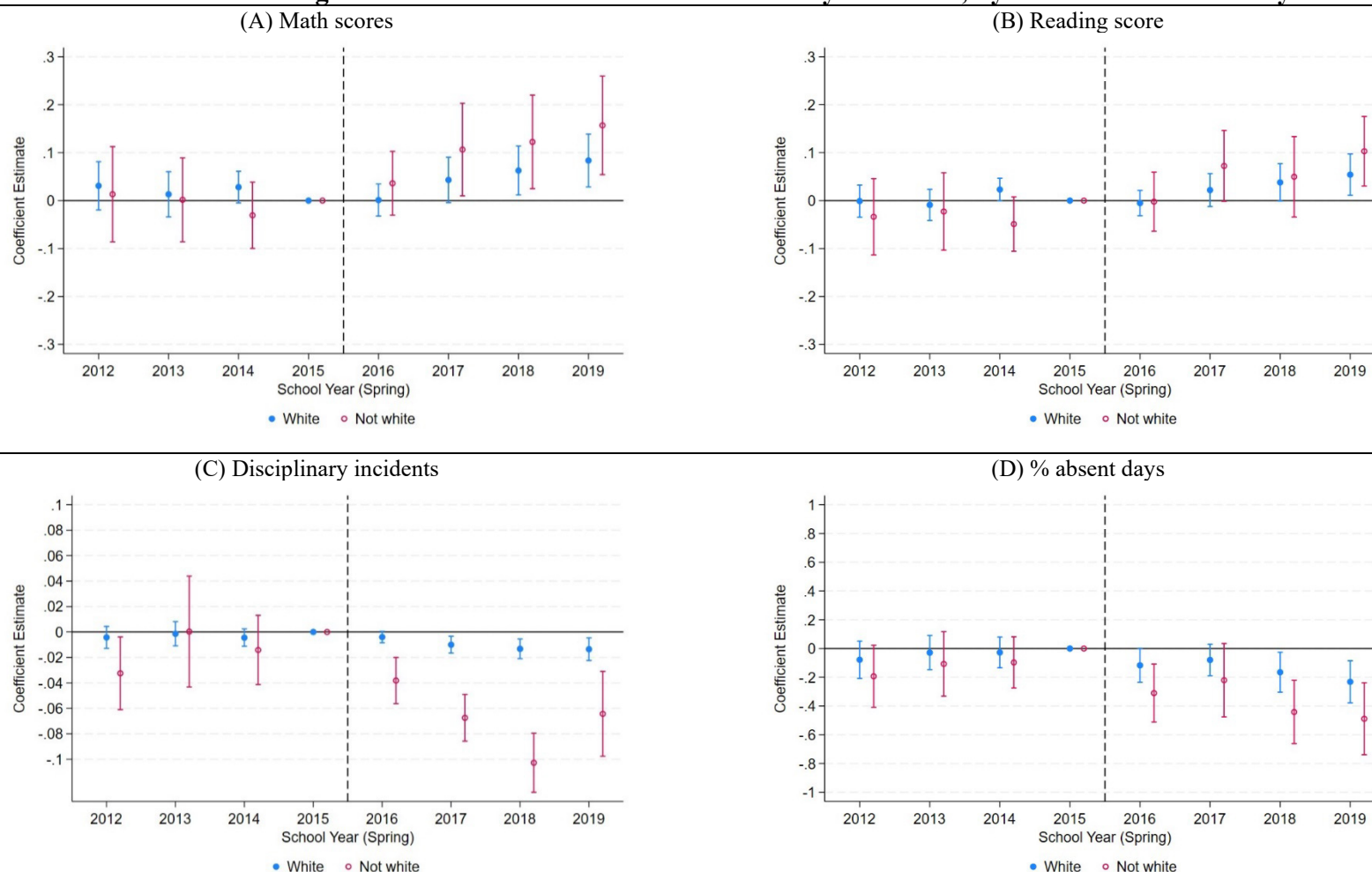
Notes: Each dot presents the estimated coefficient on the interaction term between the MAF indicator and the indicator for the corresponding school year, with the spikes providing the 95% confidence intervals. The estimates are obtained using schools in the first cohort of MAF schools as the treatment and never MAF schools as the comparison group with school, grade, and year fixed-effects. Standard errors are clustered at the school level.

Figure 4C. Effects of MAF Designation on Student Outcomes: Event Study Estimates, by Student Gender



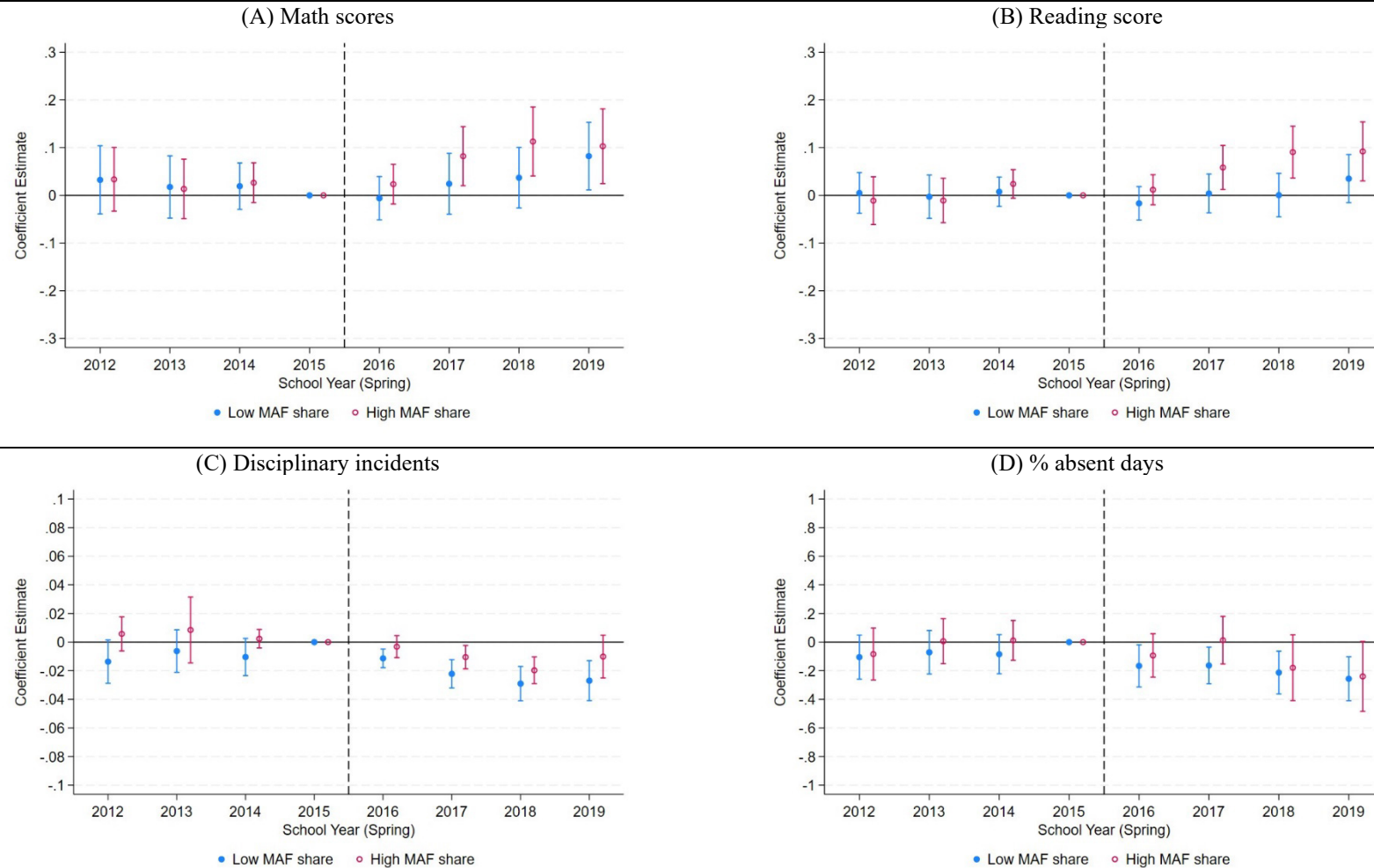
Notes: Each dot presents the estimated coefficient on the interaction term between the MAF indicator and the indicator for the corresponding school year, with the spikes providing the 95% confidence intervals. The estimates are obtained using schools in the first cohort of MAF schools as the treatment and never MAF schools as the comparison group with school, grade, and year fixed-effects. Standard errors are clustered at the school level.

Figure 4D. Effects of MAF Designation on Student Outcomes: Event Study Estimates, by Student Race/Ethnicity



Notes: Each dot presents the estimated coefficient on the interaction term between the MAF indicator and the indicator for the corresponding school year, with the spikes providing the 95% confidence intervals. The estimates are obtained using schools in the first cohort of MAF schools as the treatment and never MAF schools as the comparison group with school, grade, and year fixed-effects. Standard errors are clustered at the school level.

Figure 5. Effects of MAF Designation on Student Outcomes: Event Study Estimates, by the Share of K-3 Students Receiving MAF Intervention in 2015-16 in First Cohort MAF Schools



Notes: Each dot presents the estimated coefficient on the interaction term between the MAF indicator and the indicator for the corresponding school year, with the spikes providing the 95% confidence intervals. The estimates are obtained using schools in the first cohort of MAF schools as the treatment and never MAF schools as the comparison group with school, grade, and year fixed-effects. Standard errors are clustered at the school level.

Figure 6. Average MAF Cost Per K-5 Pupil by Personnel and Nonpersonnel for MAF Implementation Schools (in 2018-19 Dollars)

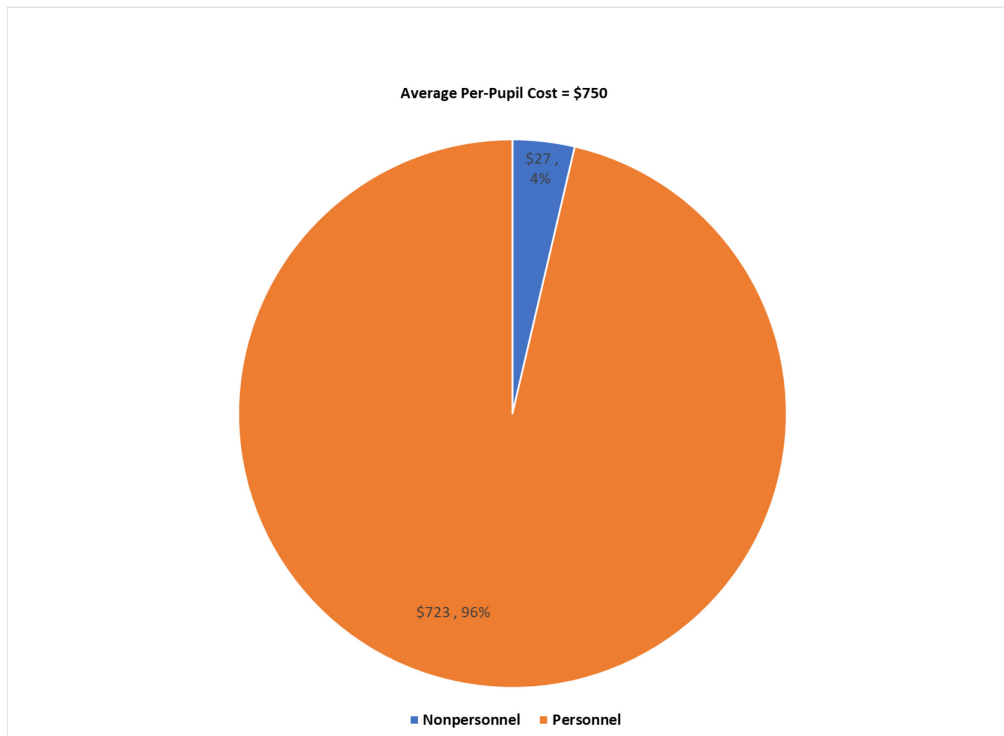


Figure 7. Average MAF Cost Per K-5 Pupil by Program Activity for MAF Implementation Schools (in 2018-19 Dollars)

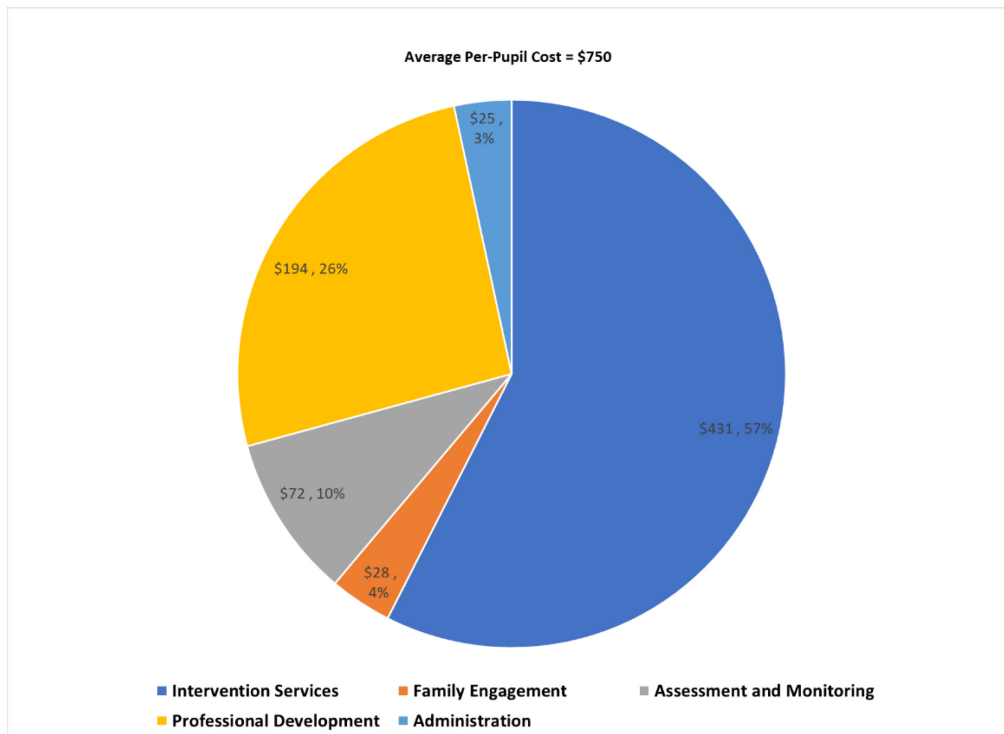


Figure 8. Average MAF Program Cost Per Pupil by Activity in High, Medium and Low Cost MAF Implementation Schools (in 2018-19 Dollars)

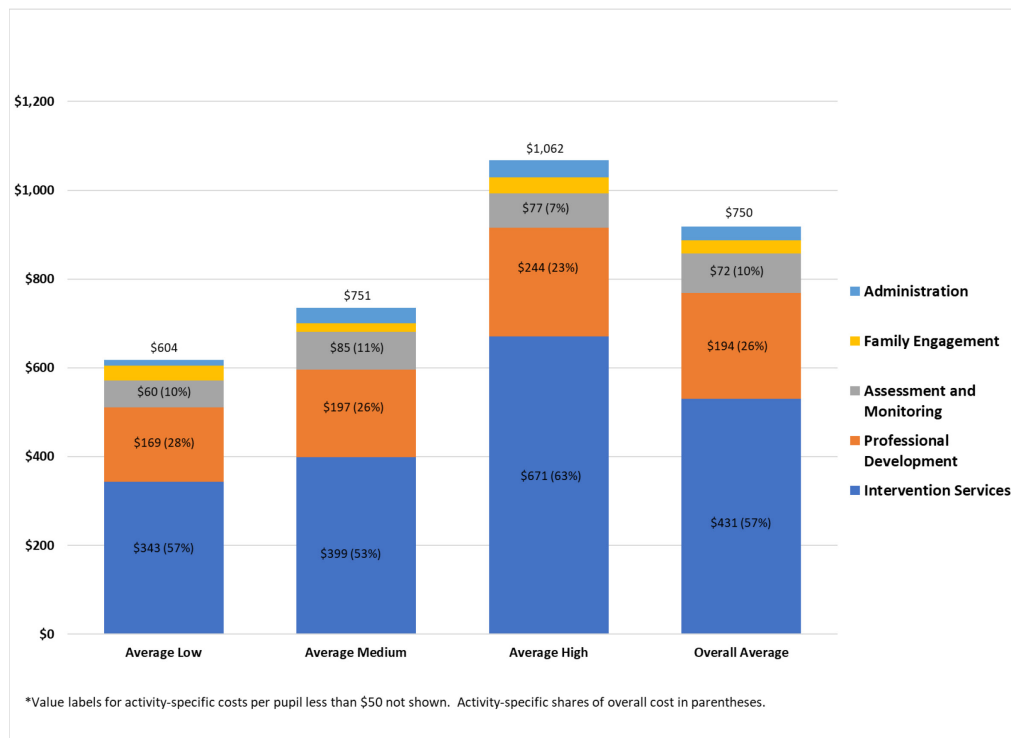


Figure 9. Relationship Between Annual Program Cost Per K-5 Pupil and K-5 Enrollment for MAF Implementation Schools (in 2018-19 Dollars)

