



# **HCEO WORKING PAPER SERIES**

Working Paper



HUMAN CAPITAL AND  
ECONOMIC OPPORTUNITY  
GLOBAL WORKING GROUP

The University of Chicago  
1126 E. 59th Street Box 107  
Chicago IL 60637

[www.hceconomics.org](http://www.hceconomics.org)

# Heterogeneity and Endogenous Compliance: Implications for Scaling Class Size Interventions\*

Karun Adusumilli      Francesco Agostinelli

Emilio Borghesan

April 2024

## Abstract

This paper examines the scalability of the results from the Tennessee Student-Teacher Achievement Ratio (STAR) Project, a prominent educational experiment. We explore how the misalignment between the experimental design and the econometric model affects researchers' ability to learn about the intervention's scalability. We document heterogeneity in compliance with class-size reduction that is more extensive than previously acknowledged and discuss its consequences for the evaluation of the experiment. Guided by this finding, we implement a new econometric framework incorporating heterogeneous treatment effects and endogenous class size determination. We find that the effect of class size on test scores differs considerably across schools, with only a small fraction of schools having significant benefits from reduced class sizes. We discuss the challenges this poses for the intervention's scalability and conclude by analyzing targeted class-size interventions.

Keywords: Scalability; Class Size Interventions; Treatment Effects Heterogeneity.

JEL Classification: C51, H52, I2, J13.

---

\*We thank Michael Dinerstein for useful comments and suggestions. We also thank the participants of the workshops and seminars where we presented our work. Adusumilli: Department of Economics, University of Pennsylvania, 133 S 36<sup>th</sup> St, Philadelphia, PA 19104, USA (email: akarun@sas.upenn.edu). Agostinelli: Department of Economics, University of Pennsylvania, 133 S 36<sup>th</sup> St, Philadelphia, PA 19104, USA (email: fagostin@sas.upenn.edu). Borghesan: Industrial Relations Section, Princeton University, 20 Washington Road, Princeton, NJ 08544, USA (email: emilio.borghesan@gmail.com).

# 1 Introduction

Whether reducing class size raises academic performance is a longstanding debate in empirical economics and policy discussions ([Mishel and Rothstein 2002](#)). While early observational studies pointed towards a small or null effect of class size, results from the Tennessee’s Student Teacher Achievement Ratio (STAR) experiment suggested otherwise. The experiment, conducted between 1985 and 1989, randomized kindergarten students at participating public schools into one of three class types: small, regular, and regular with a teacher’s aide. Multiple evaluations of STAR demonstrated significantly higher test scores for students attending the small class type, leading many researchers to conclude that class size reductions generated causal gains in student learning ([Folger and Breda 1989](#), [Finn and Achilles 1990](#), [Word et al. 1990](#), [Schanzenbach 2006](#)).

The findings from STAR inspired policymakers in various states to implement large-scale class size reductions. But these policies, such as Tennessee’s Project Challenge in 1989 and California’s Class Size Reduction law in 1996, saw disappointing results.<sup>1</sup> Prior research has showed how the null effects of these two policies can be explained by either a decline in teacher quality or by insufficient class size reductions in schools ([Hippel and Wagner 2018](#); [Jepsen and Rivkin 2009](#)). Our paper instead focuses on an overlooked aspect of the debate surrounding STAR and the scaling of experimental interventions. We ask whether and to what extent it is possible to learn about the scalability of class-size policies from the STAR experimental data itself.

We argue that the reduced-form evaluation and the two-stage least squares (2SLS) estimator commonly used in prior analyses are poorly-aligned with the STAR experimental design, thereby limiting the ability of researchers to learn about the intervention’s scalability from the experimental data. STAR randomized students into different class types. But, class types are not the same as class sizes, and school principals had some discretion in choosing target sizes for each class type. We document substantial heterogeneity in compliance, with different schools

---

<sup>1</sup>Project Challenge produced few benefits ([Achilles et al. 1995](#); [Hippel and Wagner 2018](#)), and initial analysis of California’s law found no effect of the statewide average reduction in class size of ten students ([Bohrnstedt and Stecher 1999, 2002](#); [Stecher and Bohrnstedt 2000](#)).

targeting both heterogeneous sizes (doses) and heterogeneous class-size reductions between treatment and control arms. In this context, reduced-form estimates conflate heterogeneous doses and impact effects. The 2SLS estimator instead identifies a weighted average of school-specific effects, but these weights depend on each school's endogenous compliance with the experimental design. Neither estimator identifies the policy-relevant relationship between class size and test scores, which is crucial for scaling class-size initiatives.

We therefore develop and apply a new econometric method that is better-aligned with the realities of STAR's implementation. We allow class size targets and the heterogeneous effects of class size on test scores to be jointly determined by a low-dimensional set of latent parameters. We classify schools into groups based on similar values of these latent parameters. The classification works so that class size targets and impact effects are independent within each group, but they may vary across groups, thereby allowing for grouped selection on unobserved gains. We use a Grouped Random Effects (GRE) methodology to simultaneously group schools and estimate parameters that govern the distribution of compliance behavior and impact effects (Adusumilli 2020). The number of groups is determined through an information criterion (BIC). BIC selects three groups as optimal.

Our method uncovers sharp differences in the effectiveness of class size reductions across schools. Nearly all of the gains from reducing class size in STAR are driven by 29% of the schools in the sample. In this set of schools, a one-student reduction in class size causes a 0.09 standard deviation (sd) improvement in test scores. The relationship between class size and test scores is modest in the remaining 71% of schools. In fact, if the 29% of highly sensitive schools had been omitted from the experiment, 2SLS would have failed to detect any causal effect of class size on test scores.

We also find evidence of heterogeneous compliance. All schools under-complied with the intended experiment by creating class-size reductions that were too small. Additionally, schools with larger impact effects in absolute magnitude under-complied even further. We simulate counterfactual *intended* implementations of STAR – with treated class sizes ranging from 13 to 17 students, control

class sizes of between 22 and 25 students, and uniform compliance across schools – and find that it would have generated 18% larger test score gains.

Our analysis of the experimental data reveals both the challenges and opportunities in learning about scaling from the STAR intervention. Potential challenges include marked heterogeneity in impact effects across schools and the fact that schools with larger impact effects, whether positive or negative, exhibited lower compliance in reducing class sizes. We decompose the policy-relevant treatment effect, introduced in [Heckman and Vytlačil \(2001\)](#), into weighted sums of group-specific impact effects and group-specific compliance. A scaled-up version of the STAR experiment could therefore be undermined if a new population has a small proportion of schools with high impact effects (different weights) or low compliance. We show how researchers can forecast these weights in new populations by exploiting site-specific covariates to predict group membership as uncovered by the GRE model. Population heterogeneity can then turn into an opportunity rather than a challenge. If schools with high impact effects can be precisely identified, targeted interventions emerge as a cost-effective alternative to universal policies. To demonstrate this, we simulate a version of the STAR experiment that reduces class size at only the schools with the highest impact effects. Such an intervention would have generated test score gains that were 10.3% larger than the actual experiment while reducing the Black-white test score gap by 76.7% as much as the original experiment.

**Related Literature.** An extensive empirical literature uses observational and quasi-experimental designs to estimate the effect of class size on students' educational outcomes. The results are mixed, and there is debate over whether class size truly matters. For example, [Card and Krueger \(1992a\)](#) find that men have a higher return to schooling in states with higher quality schools (including lower pupil/teacher ratios).<sup>2</sup> On the other hand, [Heckman, Layne-Farrar, and Todd \(1995\)](#) argue that these effects vanish once the empirical model allows for nonlinear effects of school quality on students. [Hanushek \(1997\)](#) argues that the relationship between school quality and student achievements becomes insignif-

---

<sup>2</sup>[Card and Krueger \(1992b\)](#) find that the improvement in the quality of Black schools for the cohort of Southern-born men in 1960, 1970 and 1980 explained 20% of the narrowing of the Black-white earnings gap during the same period of time.

icant once researchers account for the heterogeneity in family inputs. [Angrist and Lavy \(1999\)](#) exploit the Maimonides' rule of classroom composition in Israeli public schools and find that class-size reductions have positive effects for fourth and fifth graders, but not for third graders. However, this finding is re-evaluated in [Angrist et al. \(2019\)](#) with more recent data and a larger sample, and the authors report that the effect of class-size reduction is now close to 0 for all grades. [Hoxby \(2000\)](#) exploits idiosyncratic variation in cohort population sizes to estimate the effect of class size on student learning and finds no effect. [Rivkin, Hanushek, and Kain \(2005\)](#) compare the relative importance of teachers and school quality and conclude that teacher quality matters more for children's academic achievement.

Analysis of the STAR experiment appears to rebut this conclusion. Multiple evaluations of STAR find positive short- and long-term effects on test scores and educational attainment for children who were randomly assigned to smaller classes. Because of the random assignment, results from STAR have been used to reinterpret previous research on the topic and to provide strong arguments in support of policy proposals to reduce pupil/teacher ratios (see [Krueger 1999](#), [Krueger and Whitmore 2001](#), [Schanzenbach 2006](#), [Chetty et al. 2011](#)). We argue, however, that these prior analyses of the STAR experiment do not take compliance heterogeneity into account, and this therefore limits the extent to which the conclusions drawn there are informative about a scaled-up version of the intervention.

Our work contributes to the growing literature that studies the ability of field experiments to inform policy decisions at scale. [List \(2022\)](#) tells the history of attempts to scale experimental interventions and highlights many potential pitfalls. [Heckman \(1992\)](#) discusses how endogenous selection into field experiments impacts the external validity of the estimated impacts. [Gechter et al. \(2023\)](#) consider the problem of a decision maker who wants to select sites for an RCT to maximize external validity. By treating STAR as a series of school-specific mini experiments, our GRE approach shows how the experimental data itself can inform both the threats and opportunities to scale their intervention in new contexts.

On the methodological side, our work is related to the large and growing literature on clustering methods and the EM algorithm. [Bonhomme and Manresa \(2015\)](#) and [Bonhomme, Lamadon, and Manresa \(2022\)](#) introduce Grouped Fixed

Effects (GFE) for panel data models. GFE categorizes units into groups so that each group has the same value of unobserved heterogeneity. Grouped Random Effects instead allows the distribution of unobserved heterogeneity to vary across groups. GFE is a special case of GRE, in which the random parameters are constrained to have no variance.

We use GRE over GFE for three reasons: First, it enables us to jointly model school-specific class size targets and treatment effects, which allows for correlation between compliance and impact effects across schools. Second, GRE estimates are more precise as they require fewer groups. GRE clusters observations so that the unobserved heterogeneity is approximately independent of the covariates within each group, which is weaker than requiring unobserved heterogeneity to be constant within each group as in GFE. Hence, it requires fewer groups, or, equivalently, produces smaller bias with the same number of groups. It also allows us to use posterior averaging to get better estimates of the marginal effects. Third, GRE simultaneously computes both the group assignments and estimates of group-specific parameters. By contrast, [Bonhomme, Lamadon, and Manresa \(2022\)](#) suggest a two-step method, which requires specifying moments for the first-stage group assignments and can be less accurate than one-step methods. [Bonhomme and Manresa \(2015\)](#) suggest a one-step method for linear panel data models, but our setting is nonlinear due to the need to model class sizes.

The GRE approach uses a modified version of the EM algorithm for computation, termed EAMP. This algorithm is closely related to recent developments in computer science on variational inference (see, e.g., [Blei, Kucukelbir, and McAuliffe 2017](#)), and the interpretation of EM as a variational optimization problem ([Neal and Hinton 1998](#), [Bergman et al. 2019](#)).

The Grouped Random Effects framework shares similarities with empirical Bayes estimation common in the literature on teacher value-added ([Chetty, Friedman, and Rockoff 2014](#)), neighborhood fixed effects ([Bergman et al. 2019](#)), and firm-level estimates of racial discrimination ([Kline and Walters 2021](#); [Kline, Rose, and Walters 2022](#)). In these papers, a Gaussian prior is used to shrink individual estimates in order to achieve lower MSE in the aggregate. Our GRE method employs Gaussian priors for the treatment effects as well, but we allow the prior param-

eters to vary across groups and we shrink school-specific treatment effects towards group-specific priors. This allows for greater flexibility in capturing complex patterns of unobserved heterogeneity. Additionally - and this is different from the empirical Bayes approach - we simultaneously shrink the estimates of both school-specific treatment effects and first stage compliance. Empirical Bayes estimation can thus be considered a special case of our approach.

## 2 Evaluation of STAR with Endogenous Class Size Reduction

### 2.1 The Tennessee STAR Experiment

Project STAR was a four-year longitudinal study of elementary school children in Tennessee between 1985 and 1989. In 1985 schoolchildren entering kindergarten at seventy-nine participating elementary schools were randomly allocated to one of three class types: small, regular, and regular with the addition of a teacher's aide. Teachers were randomly assigned to separate classes. The target size for the small classes was between thirteen and seventeen students, while the target size for the other two class types was between twenty-two and twenty-five students. Actual class size sometimes deviated from these targets due to schoolroom capacity constraints and attrition from participating schools.

The experimental design called for students to remain in the same class type through the third grade. However, a variety of reasons including sample attrition, behavioral issues, and parental complaints, led many students initially randomized into a small class to attend another class type in later grades or drop out of the experiment altogether. Of the 1,900 children initially randomized into a small class in kindergarten, only 857 (45%) attended a small class for all four years.<sup>3</sup> For this reason, we analyze the effects of class size on academic performance in kindergarten only, before any potentially endogenous reallocations of students to different class types may have taken place.<sup>4</sup>

---

<sup>3</sup>Further details about the experimental design are provided in [Boyd-Zaharias et al. \(2007\)](#).

<sup>4</sup>Prior research has found that assignment to the Regular + Aide group had no discernible effect on academic performance, so we omit this group from the analysis and concentrate instead on the small- and regular-sized class types ([Finn and Achilles 1990](#), [Word et al. 1990](#), [Folger and Breda 1989](#)). Appendix C shows that we find similar results when pooling together regular and regular + aide classrooms.



At the end of each year of the experiment, students were given a range of cognitive and noncognitive evaluations. We focus on three cognitive tests administered to the students at the end of kindergarten: the Stanford Achievement Tests (SAT) in math, reading, and word skills. The SAT uses item response theory to facilitate comparisons across students and across years for the same student. Our dependent variable is a simple average of scores on these three exams. For the 2.3% of students with scores on only one or two exams, we construct their mean score as the average of the available scores.

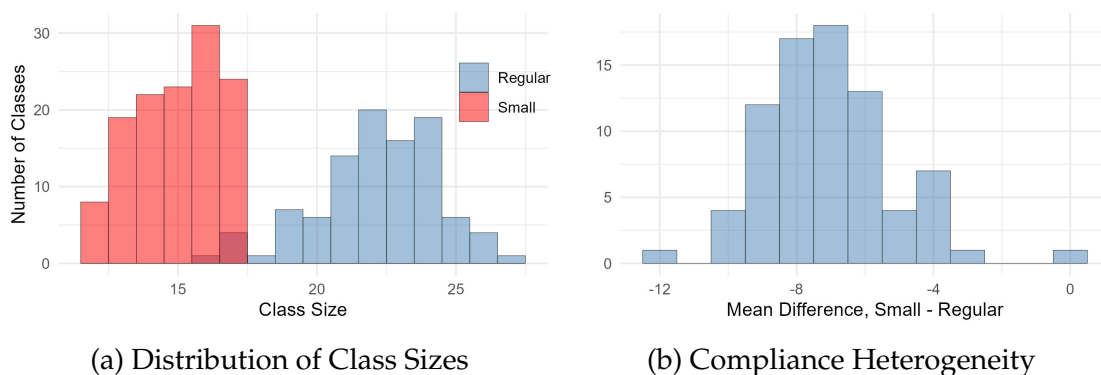
Project STAR collected additional demographic information on students and their teachers. These include information on student race, gender, whether students qualified for a free or reduced price lunch (a typical proxy for low family income), absences from school, as well as teacher demographic information and qualifications. [Krueger \(1999\)](#) show that the treatment and control groups in kindergarten do not differ by these observable characteristics, suggesting that the initial randomization of students into classroom types was not compromised.

Figure 1a illustrates the distribution of class sizes in kindergarten, where the smaller "treated" classes comprised twelve to seventeen students, and the regular "control" classes ranged from sixteen to twenty-seven students. Despite the binary classification created by the experiment, class size is a continuous treatment variable, and we will use the term dose throughout the text to stress this. Figure 1b illustrates the heterogeneity in dose reductions across schools. Although the average reduction in class size between control and treatment classrooms is seven students, the figure demonstrates that dose reductions varied considerably across schools, ranging between zero and twelve students. The next two sections address (i) how this heterogeneous compliance influences the interpretability of past estimates of class-size effects from the STAR experiment; and (ii) the degree to which heterogeneous compliance affects the STAR experiment's ability to inform the scalability of class-size reductions in new contexts.

## 2.2 The Implications of the Endogenous Class Size Reduction

Prior studies have claimed that the random assignment of students to different class types in the STAR experiment means that findings from STAR are more

**Figure 1: Distribution of Class Sizes and Class-Size Reductions**



Panel (a) plots the number of kindergarten classes of each size for regular and small classes in the Tennessee STAR experiment. Panel (b) plots the average difference in size between small and regular classrooms at each school.

credible than earlier observational studies which had yielded contrasting conclusions on the effects of class size (Schanzenbach 2016).<sup>5</sup> The first wave of studies evaluated the reduced form of the STAR randomization by estimating specifications like the following:

$$y_{is} = \beta_0 + \beta_1 \text{Small}_{is} + \beta_2 \text{Aide}_{is} + \eta_s + \varepsilon_{is} , \quad (1)$$

where  $\text{Small}_{is}$  is a binary variable that equals one if a student attends a small class, and  $\text{Aide}_{is}$  is a binary variable that equals one if a student attends a regular-sized class with a teacher and a teacher's aide. The omitted category is a regular-sized class with a teacher but no aide. This specification compares the average test scores of students who, due to the experiment, randomly attended different class types.

Since the experiment did not specify the exact size for each type, principals were free to choose from different class sizes within a targeted range for each treatment arm. This relatively free choice of class sizes plus some deviations from the targeted ranges generates the heterogeneity in dose reductions seen in Figure 1b.

Dose heterogeneity among students who attended the same class type makes

<sup>5</sup>See Hanushek 1986 for a review of the early observational studies of class size.

generalizing the reduced form findings difficult. An analogy to randomized controlled trials in medicine serves to illustrate this point. Consider an experiment that aims to test the effectiveness of a single-use medicine by randomly assigning patients to two groups. Patients in the treatment group can select the dose they would like to be administered. Patients in the control group are also able to select the dose as long as it was weakly less than the minimum dose received by anyone in the treatment group. It would then not be possible to determine the effectiveness of the treatment from the reduced form of this experiment because it confounds the effectiveness of the medicine with the differences in dose intensity chosen by individuals within and between treatment and control groups. The reduced form of the STAR experiment has the same limitation: it confounds class size effectiveness with heterogeneity in dose reductions across schools.

A second wave of studies on STAR aimed to produce more generalizable results by estimating education production functions using randomization into a small class as an instrument for the class size dose within a 2SLS framework (Krueger 1999). The within-school class-type randomization,  $Z_{is} \in \{0, 1\}$ , serves as an instrument for class size,  $n_{is}$ , in the following model:

$$\begin{aligned} n_{is} &= \mu_s + \delta_s Z_{is} + \nu_{is} \\ y_{is} &= \eta_s + \beta_{is} n_{is} + \epsilon_{is} , \end{aligned} \tag{2}$$

where we omit controls for ease of interpretation. Researchers using 2SLS assume that the within-school randomization is mean independent of unobserved determinants of test scores,  $\epsilon_{is}$ :

**Assumption 1 (No Selection on Unobservables).** *The STAR randomization between small and large class types  $Z_{is} \in \{0, 1\}$  generates random student compositions between different class types:  $E[\epsilon_{is}|Z_{is}, s] = E[\epsilon_{is}|s]$ .*

It is well known that the causal interpretation of 2SLS estimates with ordered multi-valued treatment doses is complicated by heterogeneous dose reductions across units, as discussed in Angrist and Imbens 1995, Rose and Shem-Tov 2021, and in the context of essential heterogeneity in Heckman, Urzua, and Vytlačil (2006). In the context of STAR, 2SLS fails to identify the policy-relevant marginal

effect of class size on test scores even under assumption 1. Proposition 1 shows that it instead identifies a weighted average of school-specific effects, where the weights depend on schools' endogenous dose reductions.

**Proposition 1.** *Suppose Assumption 1 holds, and that school-specific treatment effects are linear. Moreover, suppose that the original STAR randomization is relevant, i.e., it generates significant differences in class size within each school. The 2SLS estimator of model (2) identifies a weighted average effect of marginal effects, where the weights depend on the endogenous compliance behavior of schools with respect to the class size reduction in the experimental setting:*

$$\beta_1^{2sls} = \sum_{s \in \mathcal{S}} \beta_s \frac{\phi_s \bar{Z}_s (1 - \bar{Z}_s) (\mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 1] - \mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 0])}{\sum_{s \in \mathcal{S}} \phi_s \bar{Z}_s (1 - \bar{Z}_s) (\mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 1] - \mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 0])}. \quad (3)$$

In the above equation,  $\beta_s = \mathbb{E}[\beta_{is}|s]$ ,  $\phi_s$  is the fraction of students attending school  $s$ ,  $\bar{Z}_s$  is the fraction of individuals in school  $s$  who are in small (treated) classrooms, and

$$\begin{aligned} \mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 1] &= \frac{1}{n_{s,treat}} \sum_{c \in s \cap treat} n_{cs} (n_{cs} - \bar{n}_s), \\ \mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 0] &= \frac{1}{n_{s,control}} \sum_{c \in s \cap control} n_{cs} (n_{cs} - \bar{n}_s), \\ \bar{n}_s &= \frac{\sum_{c \in s} n_{cs}^2}{\sum_{c \in s} n_{cs}}. \end{aligned}$$

The quantities  $\mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 1]$  and  $\mathbb{E}[\Delta n_{is}|s, Z_{i,s} = 0]$  represent the deviations between treated and control class sizes from the average class size at school  $s$ , denoted by  $\bar{n}_s$ .

The average marginal effect of class size on test scores,  $\mathbb{E}[\beta_{is}]$ , is a main policy-relevant parameter of interest, but proposition 1 shows that 2SLS will only identify  $\mathbb{E}[\beta_{is}]$  in two special cases: (i) when impact effects are homogeneous across schools, namely  $\beta_s = \beta$  for all  $s \in \mathcal{S}$ ; or (ii), when schools have equal levels of compliance regarding dose reductions:  $\mathbb{E}[\Delta n_{is}|s, Z_{i,s}] = \mathbb{E}[\Delta n_{is}|Z_{i,s}]$  and  $\bar{Z}_s = \bar{Z}$ . Both cases impose strong restrictions on either impact heterogeneity

or on schools' objective functions. In the first case, an extensive literature has demonstrated the existence of heterogeneous treatment effects in educational settings (Carneiro, Heckman, and Vytlačil 2011, Walters 2018, Borghesan and Vasey 2024). The second case will be violated if, for example, schools choose a particular reduction in class size between small and regular class types, thereby self-selecting into a particular intensity of treatment on the basis of their unobserved gains from the experiment.

Apart from these special cases, the 2SLS estimand represents a weighted average of school-specific marginal effects of class size on test scores. The weights in (3) depend on three factors, each of which may be correlated with the effect of class size on test scores at that school. They are highest in schools that enroll many students, where  $\phi_s$  is large, they are higher in schools with an equal number of treated and control students, which maximizes  $\bar{Z}_s(1 - \bar{Z}_s)$ , and they are higher in schools that create a large difference between the sizes of treatment and control classrooms.<sup>6</sup> The last two factors are endogenous in the STAR experiment.

Table 1 provides an example with three schools to demonstrate how endogenous choices may affect the weights in the 2SLS estimand. School A has 45 students and creates a treatment class of 15 and a control class of 30. School B has 45 students and creates a treatment class of 20 and a control class of 25. School C has 60 students and creates two treatment classes of 15 students and one control class of 30 students. School A has a low population, a low treated fraction of 1/3, but a large dose reduction between treated and control classrooms. School B has a higher fraction of treated students than school A but has a lower dose reduction. Finally, School C is associated with a high value for all components that affect the 2SLS weights. School C is then overweighted in the 2SLS estimand relative to its size, while school B is significantly underweighted. In this example, two of the three factors affecting the weights - the fraction of students receiving treatment and the dose reduction - are under the school's control in the STAR experiment.

---

<sup>6</sup>The weights will be nonnegative at any school where the average treated classroom has weakly fewer students than the average control classroom. This was the case with all schools in the STAR experiment.

**Table 1:** Example of How IV Weights Are Determined

	$N$	Treated	Control	$\phi_s$	$\bar{Z}_s(1 - \bar{Z}_s)$	Dose Reduction	Weight
A	45	15	30	0.3	0.22	-15	0.35
B	45	20	25	0.3	0.25	-5	0.13
C	60	15, 15	30	0.4	0.25	-15	0.52

The table provides an example of the control and treatment class sizes created by three hypothetical schools.  $N$  refers to the enrollment at the school, the second and third columns refer to the treated and control classroom sizes created, and the remaining columns correspond to elements of the IV weights in (3).

### 2.3 Lessons from STAR for Class-Size Reduction Policies at Scale

Is the compliance heterogeneity of STAR a flaw or a virtue? We believe that the STAR experiment serves as a natural laboratory to learn about the challenges of scaling class-size reduction interventions. The endogenous compliance behavior exhibited by schools within the STAR experiment, though posing a potential challenge for specific econometric evaluation methodologies, also serves as a valuable resource. It provides an opportunity to gain insights into the factors influencing the adoption of policies at scale, as well as potential bottlenecks (DellaVigna, Linos, and Kim 2023). In STAR, each school represents a distinct experiment, which allows the researcher to study how impact heterogeneity and dose reductions vary systematically across schools. This approach aligns with the recommendation of Al-Ubaydli, List, and Suskind (2020) that researchers leverage multi-site trials to investigate the variability in program impacts across diverse populations and situational contexts. They argue that “the design of multi-site trials can provide empirical content into why effects might not scale and give empirical hints about where more research is necessary before scaling.”

One way to learn about impact heterogeneity in the STAR experiment would be to estimate the relationship between class size and test scores separately by school. Even if impact heterogeneity and dose reductions were correlated *across schools*, the within-school random assignment of students and teachers to classes

of varying sizes would produce unbiased estimates of school-specific impact effects. However, limited variation in class size within each school would then lead to noisy estimates, amplifying the scalability threat posed by low-powered studies (Al-Ubaydli, List, and Suskind 2020).<sup>7</sup> Such a strategy would be suboptimal for researchers and policymakers who intend to extrapolate from the STAR experiment to potential class-size reductions in new contexts. We therefore propose a new econometric framework that aims to balance the twin attractions of statistical efficiency and allowing for unobserved heterogeneity in impact effects and compliance. Our approach will thus shed insight on how these two forces affect the scalability of the STAR experiment.

### 3 From STAR Randomization to Policy-Relevant Effects

In this section, we introduce a model aimed at uncovering heterogeneous effects of class size on test scores in the presence of potentially endogenous compliance. Our method clusters schools into a fixed number of groups,  $k = 1, \dots, K$ , that share a common set of underlying parameters that collectively govern (1) the distribution of class sizes within each school and, (2) the marginal effects of class size on test scores.

#### 3.1 Test Score Model

Test scores are generated by the following equation for student  $i$  in school  $s$  belonging to group  $k$ :

$$\begin{aligned}
 y_{isk} &= \eta_s + \beta_{isk}n_{isk} + x_{isk}'\theta + \epsilon_{isk} , \\
 \beta_{isk} &\sim N(\mu_k, \Sigma_k) , \\
 \epsilon_{isk} &\sim N(0, \sigma_{\epsilon,k}^2),
 \end{aligned} \tag{4}$$

---

<sup>7</sup>The median school in the experiment features two distinct classes: one small and one regular.

where  $n_{isk}$  is individual  $i$ 's class size and  $x_{isk}$  is a vector of observable student characteristics.<sup>8</sup> In the model, impact effects are heterogeneous and are assumed to follow a Normal distribution with mean,  $\mu_k$ , and variance,  $\Sigma_k$ , that are allowed to differ by group. The term  $\mu_k$  describes the expected effect of a 1-unit increase in class size for a randomly selected student attending a school in group  $k$ . Student demographics,  $x_{isk}$ , instead have non-random effects on test scores. Test scores additionally depend on a full set of school fixed effects,  $\eta_s$ , and a student-specific shock,  $\epsilon_{isk}$ , with mean 0 and a variance  $\sigma_{\epsilon,k}^2$  that is allowed to vary by group. The model therefore incorporates group-specific heteroskedasticity.

The random assignment of teachers and students to classrooms is implicitly imposed via two modeling choices. First, the distribution of treatment effects is school- but not classroom-specific, so that what matters for the marginal effect of class size on test scores is the match between the student and the school. Second, assignment to the treatment or control group only affects test scores through its impact on classroom size,  $n_{isk}$ .

### 3.2 Class Size Model

Each school is associated with a vector of class sizes,  $\mathbf{n}_{sk}^{(t)}$  and  $\mathbf{n}_{sk}^{(c)}$ , for the treated and control groups.<sup>9</sup> We posit that these class sizes are generated according to the following multinomial model:

$$\begin{aligned} \mathbf{n}_{sk}^{(c)} &\sim \text{multinomial}(p_{sk}^{(c)}) \text{ for all } k \text{ in the control group,} \\ \mathbf{n}_{sk}^{(t)} &\sim \text{multinomial}(p_{sk}^{(t)}) \text{ for all } k \text{ in the treatment group.} \end{aligned} \quad (5)$$

Here,  $p_{sk}^{(c)} \equiv \{p_{sk1}^{(c)}, \dots, p_{skL}^{(c)}\}$ ,  $p_{sk}^{(t)} \equiv \{p_{sk1}^{(t)}, \dots, p_{skM}^{(t)}\}$  are school-specific probability distributions over the size of each class. The support for these distributions is the

---

<sup>8</sup>We depart from [Krueger \(1999\)](#) in using raw scores rather than percentile ranks as the outcome variable. Causal models of the ranks of a dependent variable do not admit a traditional interpretation of coefficients as marginal effects. In addition, our model assumes that outcomes are independently distributed conditional on group membership, while the use of ranks would induce dependence across groups. Finally, a model of test score ranks is inappropriate for many forms of counterfactual analysis, as a reduction in class size for every student may have beneficial effects on academic outcomes without causing any discernible effect on ranks.

<sup>9</sup> $\mathbf{n}_{sk}^{(t)}$  and  $\mathbf{n}_{sk}^{(c)}$  represent vectors of counts of the number of classes in school  $s$  of each size. The class size experienced by each student in equation (4),  $n_{isk}$ , is a scalar integer-valued variable.



same as the support for the class size vectors,  $\mathbf{n}_{sk}^{(t)}$  and  $\mathbf{n}_{sk}^{(c)}$ , and is equal to their observed support in the data:  $\{12, \dots, 17\}$  and  $\{16, \dots, 27\}$  respectively (so that  $L = 6$  and  $M = 12$ ).

The model assumes that class sizes are drawn from school-specific multinomial distributions, with multiple classes in the same school of a particular type (either treatment or control) representing independent draws from the same distribution. One can intuitively think of  $p_{sk}^{(c)}$  and  $p_{sk}^{(t)}$  as denoting the target proportions of control and treatment classroom sizes for each school. The observed class sizes are then random deviations from this target. The assumption that  $p_{sk}^{(c)}$  and  $p_{sk}^{(t)}$  are school-specific rather than classroom-specific is consistent with the random assignment of students and teachers to classrooms. We do not model the assignment of treatment and control status to the classroom, since these were determined randomly in the experiment.

Importantly, we model class size levels and not simply the difference in class size between treatment and control classrooms. It is quite plausible that schools with different control class sizes are associated with different impact effects, causing selection bias. Sorting on gains could result if schools with large classrooms prior to the experiment simultaneously have large impacts and choose a larger (or smaller) class-size reduction. By modeling the class size levels for both treatment and control class types, we are able to account for both selection bias and sorting on gains.

Note that we do not take the number of students in each school as given. It is endogenous in our model, determined by the sizes of each class in the school. Endogenizing school size in this manner is useful if school-level factors determining cohort size are correlated with the impact effects,  $\beta_{isk}$ .

### 3.3 Endogeneity and Grouped Random Effects

We adopt a Grouped Random Effects approach, following [Adusumilli \(2020\)](#), to simultaneously estimate (4) and (5) and to group schools that have a common set of underlying parameters that jointly determine class size and impact effects. These common parameters,  $\rho_k$ , are constant within each group, and they serve as priors over each group's random parameters,  $(\beta_{isk}, p_{sk}^{(c)}, p_{sk}^{(t)})$ . By then recovering

the parameters governing the distribution of  $(p_{sk}^{(c)}, p_{sk}^{(t)})$ , we are able to isolate the intensity of the class size reduction between treatment and control classes within each school.

We use GRE as it enables us to jointly model both treatment effects and class sizes. While the Grouped Fixed Effects approach in [Bonhomme and Manresa \(2015\)](#) allows for nonlinear models, the algorithm they develop could only be applied to the linear model for test scores and not to the class size model in (5). Such an approach would only group schools based on heterogeneity in  $\beta_{isk}$  but would fail to capture the correlation between  $\beta_{isk}$  and  $n_{sk}^{(c)}$  and  $n_{sk}^{(t)}$  that can cause selection bias and sorting on gains.

Let  $k \in 1, \dots, K$  denote the set of groups, and  $w_s(k)$  the group assignment, with  $w_s(k) = 1$  if school  $s$  is in group  $k$ . For the prior on the treatment effect coefficients,  $\beta \equiv \{\beta_{isk} : s = 1, \dots, S; s(i) = s; w_s(k) = 1, k = 1, \dots, K\}$ , we specify

$$\pi(\beta|\gamma) := \prod_s \prod_{i:s(i)=s} \prod_k N(\beta_{isk} | \mu_k, \Sigma_k)^{w_s(k)}, \quad (6)$$

so that the student-specific marginal effects,  $\beta_{isk}$ , are random draws from a group-specific normal distribution with mean  $\mu_k$  and variance  $\Sigma_k$ . This assumes that the distribution of unobserved student and school characteristics affecting the relationship between class size and test scores is similar across all schools within the same group.

The school-specific multinomial probabilities for class size,  $\mathbf{p}^{(c)} \equiv \{p_{sk}^{(c)}\}_{s=1}^S$ ,  $\mathbf{p}^{(t)} \equiv \{p_{sk}^{(t)}\}_{s=1}^S$ , are assumed to have the following group-specific Dirichlet priors:

$$\begin{aligned} \pi(\mathbf{p}^{(c)}|\boldsymbol{\eta}^{(c)}) &= \prod_s \prod_k \text{Dirichlet}(p_{sk}^{(c)} | \eta_k^{(c)})^{w_s(k)}, \\ \pi(\mathbf{p}^{(t)}|\boldsymbol{\eta}^{(t)}) &= \prod_s \prod_k \text{Dirichlet}(p_{sk}^{(t)} | \eta_k^{(t)})^{w_s(k)}, \end{aligned} \quad (7)$$

where  $\boldsymbol{\eta}^{(c)} \equiv \{\eta_k^{(c)}\}_{k=1}^K$ ,  $\boldsymbol{\eta}^{(t)} \equiv \{\eta_k^{(t)}\}_{k=1}^K$  are group-specific, and  $\text{Dirichlet}(p|\eta)$  denotes the pdf of the Dirichlet distribution with parameter  $\eta$  evaluated at  $p$ . The Dirichlet distribution is the conjugate prior for the multinomial distribution. It is chosen both for computational tractability, and because it is a distribution

over probability distributions on the unit simplex. It is, therefore, an appropriate choice for generating probability vectors that govern the sizes of control and treatment classes at each school. Overall, the set of parameters to be estimated are then  $\rho_k := \{\mu_k, \Sigma_k, \eta_k^{(c)}, \eta_k^{(t)}\}_k$ .

### 3.4 Probability model

Let  $\mathbf{y} \equiv \{y_{isk}\}_{i,s,k}$  and  $\mathbf{x} \equiv \{x_{isk}\}_{i,s,k}$  denote the vectors of test scores and covariates, and group the variance of the test score disturbances as  $\sigma_\epsilon^2 \equiv \{\sigma_{\epsilon,k}^2\}_k$ . Given the treatment effect heterogeneity  $\beta$ , and the class size vectors  $\mathbf{n}^{(c)} := \{n_{isk}^{(c)}\}_{i,s,k}$  and  $\mathbf{n}^{(t)} := \{n_{isk}^{(t)}\}_{i,s,k}$ , equation (4) implies that we can model test scores using the log-likelihood

$$\begin{aligned} \ln p(\mathbf{y}|\beta, \mathbf{n}^{(c)}, \mathbf{n}^{(t)}, \mathbf{x}, \sigma_\epsilon^2, \theta) &:= \sum_s \sum_{i:s(i)=s} \ln p(y_{isk}|\beta_{isk}, n_{isk}, x_{isk}, \sigma_{\epsilon,k}^2, \theta, s), \\ &:= \sum_s \sum_{i:s(i)=s} -\frac{1}{2} \ln \sigma_{\epsilon,k(s)}^2 - \frac{(y_{isk} - \beta_{isk} n_{isk} - x'_{isk} \theta - \eta_s)^2}{2\sigma_{\epsilon,k(s)}^2}. \end{aligned}$$

Similarly, the multinomial class size model in (5) implies that the distribution of the sizes of treatment and control classrooms is given by

$$\ln p(\mathbf{n}^{(c)}|\mathbf{p}^{(c)}) := \sum_{s=1}^S \sum_{g=1}^{N^{(c)}(s)} \sum_{k=1}^K w_s(k) \left\{ \sum_{l=0}^{L-1} \mathbb{I}\{n_{skg}^{(c)} \equiv 12 + l\} p_{skl}^{(c)} \right\}, \quad (8)$$

$$\ln p(\mathbf{n}^{(t)}|\mathbf{p}^{(t)}) := \sum_{s=1}^S \sum_{g=1}^{N^{(t)}(s)} \sum_{k=1}^K w_s(k) \left\{ \sum_{m=0}^{M-1} \mathbb{I}\{n_{skg}^{(t)} \equiv 16 + m\} p_{skm}^{(t)} \right\}, \quad (9)$$

where  $N^{(c)}(s)$  and  $N^{(t)}(s)$  are the observed number of control and treatment classes, respectively, in school  $s$ . The expressions in (8) and (9) show that the log-likelihood of obtaining the vector of control and treated class sizes,  $\mathbf{n}^{(c)}$  and  $\mathbf{n}^{(t)}$ , is equal to the sum, over all groups, of an indicator for group membership times the school-group-specific multinomial probabilities for the observed class sizes.

Conditional on the unobserved heterogeneity,  $(\beta, \mathbf{p}^{(c)}, \mathbf{p}^{(t)})$ , the joint likelihood of

the observations  $(\mathbf{y}, \mathbf{n}^{(c)}, \mathbf{n}^{(t)})$  is given by

$$p(\mathbf{y}, \mathbf{n}^{(c)}, \mathbf{n}^{(t)} | \boldsymbol{\beta}, \mathbf{p}^{(c)}, \mathbf{p}^{(t)}, \mathbf{x}, \boldsymbol{\sigma}_\epsilon^2, \theta) := p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{n}^{(c)}, \mathbf{n}^{(t)}, \mathbf{x}, \boldsymbol{\sigma}_\epsilon^2, \theta) \cdot p(\mathbf{n}^{(c)} | \mathbf{p}^{(c)}) \cdot p(\mathbf{n}^{(t)} | \mathbf{p}^{(t)}).$$

Additionally, in view of (6) and (7), the prior distribution of unobserved heterogeneity is

$$\pi(\boldsymbol{\beta}, \mathbf{p}^{(c)}, \mathbf{p}^{(t)} | \boldsymbol{\gamma}, \boldsymbol{\eta}^{(c)}, \boldsymbol{\eta}^{(t)}) := \pi(\boldsymbol{\beta} | \boldsymbol{\gamma}) \cdot \pi(\mathbf{p}^{(c)} | \boldsymbol{\eta}^{(c)}) \cdot \pi(\mathbf{p}^{(t)} | \boldsymbol{\eta}^{(t)}),$$

where  $\boldsymbol{\gamma} = \{\mu_k, \Sigma_k\}_k$  denotes the collection of group-specific parameters determining the effects of class size on student test scores.

## 4 Estimation

Let  $\boldsymbol{\alpha} := (\boldsymbol{\beta}, \mathbf{p}^{(c)}, \mathbf{p}^{(t)})$ ,  $\boldsymbol{\beta}_s := \{\beta_{isk} : s(i) = s\}$ ,  $\boldsymbol{\alpha}_s := (\boldsymbol{\beta}_s, p_s^{(c)}, p_s^{(t)})$ ,  $\boldsymbol{\rho} := (\boldsymbol{\gamma}, \boldsymbol{\eta}^{(c)}, \boldsymbol{\eta}^{(t)})$  and  $\boldsymbol{\rho}_k := (\boldsymbol{\gamma}_k, \boldsymbol{\eta}_k^{(c)}, \boldsymbol{\eta}_k^{(t)})$ . Also, let  $(\mathbf{y}_s, \mathbf{n}_s^{(c)}, \mathbf{n}_s^{(t)})$  denote the set of test scores and class sizes pertaining to school  $s$ . The GRE problem is to maximize the likelihood of the data jointly over both the group assignments and the common and group-specific parameters:

$$\begin{aligned} & \max_{\{w_s(k)\}, \theta, \{\boldsymbol{\rho}_k\}} \ln \int p(\mathbf{y}, \mathbf{n}^{(c)}, \mathbf{n}^{(t)} | \boldsymbol{\alpha}, \mathbf{x}, \boldsymbol{\sigma}_\epsilon^2, \theta) \pi(\boldsymbol{\alpha} | \boldsymbol{\rho}) d\boldsymbol{\alpha} \\ & = \max_{\{w_s(k)\}, \theta, \{\boldsymbol{\rho}_k\}} \sum_{s=1}^S \sum_{k=1}^K w_s(k) \ln \int p(\mathbf{y}_s, \mathbf{n}_s^{(c)}, \mathbf{n}_s^{(t)} | \boldsymbol{\alpha}_s, \mathbf{x}, \boldsymbol{\sigma}_{\epsilon,k}^2, \theta) \pi(\boldsymbol{\alpha}_s | \boldsymbol{\rho}_k) d\boldsymbol{\alpha}_s. \end{aligned} \quad (10)$$

Following [Adusumilli \(2020\)](#), we solve the above maximization problem using the Expectation, Assignment, Maximization, and Propagation (EAMP) algorithm. To use the algorithm, we first use the Donsker-Varadhan variational formula to rewrite the likelihood in (10) as follows:

$$\max_{\{w_s(k)\}, \theta, \{\boldsymbol{\rho}_k\}} \sum_{s=1}^S \sum_{k=1}^K w_s(k) \ln \int p(\mathbf{y}_s, \mathbf{n}_s^{(c)}, \mathbf{n}_s^{(t)} | \boldsymbol{\alpha}_s, \mathbf{x}, \boldsymbol{\sigma}_{\epsilon,k}^2, \theta) \pi(\boldsymbol{\alpha}_s | \boldsymbol{\rho}_k) d\boldsymbol{\alpha}_s$$

$$= \max_{\{q_{sk}(\cdot)\}, \{w_s(k)\}, \theta, \{\rho_k\}} \sum_{s=1}^S \sum_{k=1}^K w_s(k) \left\{ E_{q_{sk}(\cdot)} \left[ \ln p(\mathbf{y}_s, \mathbf{n}_s^{(c)}, \mathbf{n}_s^{(t)} | \boldsymbol{\alpha}_s, \mathbf{x}, \boldsymbol{\sigma}_{\epsilon, k}^2, \theta) \right] - \text{KL}(q_{sk}(\boldsymbol{\alpha}_s) || \pi(\boldsymbol{\alpha}_s | \rho_k)) \right\}. \quad (11)$$

In the above equation,  $q_{sk}(\cdot)$  denotes a group-specific distribution over  $\boldsymbol{\alpha}_s$  and the maximization is carried out over the space of all possible distributions  $q_{sk}(\cdot)$ . The EAMP Algorithm proceeds by repeatedly maximizing over each of  $\{q_{sk}(\cdot)\}$ ,  $\{w_s(k)\}$ ,  $\theta$ , and  $\{\rho_k\}$  holding other quantities fixed. This results in a sequence of four steps – Expectation, Assignment, Maximization, and Propagation – that are repeated in an iterative process until the algorithm converges. Note that we first demean all variables by their school-specific means to eliminate the school fixed effects before running the estimation algorithm. A detailed description of the algorithm’s steps is provided in Appendix A.

## 5 Empirical Findings

In this section we present the estimates of the model with controls for child gender, race (an indicator for being neither white nor Asian), and free lunch status. The EAMP algorithm takes the number of groups as given. In choosing the group size, we aim to balance the competing attractions of parsimony and allowing for richer patterns of heterogeneity. The EAMP framework allows for the number of groups to be as many as the number of schools. However, allowing for too many groups may cause the model to capture more noise than signal and render the estimates difficult to interpret. In practice, we use the Bayesian Information Criterion (BIC) to select the number of groups. BIC selects three groups as optimal.<sup>10</sup>

Altogether, our analysis consists of kindergarten students in the seventy-nine participating schools in the Tennessee STAR experiment who were randomly assigned to either a small or a regular class, who do not lack information on race, gender, or eligibility for free/reduced price lunch, and who have at least one exam score at the end of kindergarten. These restrictions result in an estimation sample of 3813 students.

---

<sup>10</sup>Appendix D presents estimation results with four groups. These estimates are similar to the estimates from the three-group model.

**Table 2: Class Size Marginal Effects by Group**

Group	$\mu_k$	$\Sigma_k$	$\sigma_{\epsilon,k}^2$	Schools	Students
1	-0.068 (0.061)	0.042 (0.065)	4.933 (1.55)	31	1425
2	-0.339 (0.062)	0.235 (0.15)	6.442 (1.86)	23	1137
3	0.106 (0.059)	0.001 (0.07)	11.852 (2.07)	25	1251
Avg. Effect	-0.092 (0.033)	0.118 (0.048)	7.653 (0.68)		
N	3813				

The table shows the estimated model parameters governing the effect of class size on test scores. Regressions include controls for gender, race, and free lunch status. Controls are constrained to be equal across groups. The class size effects represent the effect of a one-unit increase in class size on a simple mean of three test scores. Bootstrapped standard errors from 95 bootstrapped data sets are in parentheses. Bootstrapping involves first sampling schools from the set of 79 schools that participated in the experiment and then sampling individuals within each school.

Summary statistics for the estimation sample are presented in Appendix Table B-1. The exams, which were originally measured on a scale of 0 to 1000, have been converted to a 0-100 scale. The table shows that the mean exam score is 45.30 points, and the standard deviation is 3.59 points. The minimum and maximum test scores in the sample are 28.80 and 61.53, respectively. The average student sits in a classroom with nearly 19 students including herself. 49% of students are female, 47% of students are eligible for free lunch, and 32% are neither white nor Asian, meaning that they are either Black, Hispanic, or Native American. Slightly fewer than half of students, 46.5%, are in small (treated) classrooms.

Table 2 presents the estimated parameters of the model for test scores with three

groups. It presents the mean ( $\mu_k$ ) and variance ( $\Sigma_k$ ) of the marginal effects of class size on test scores for each group, alongside the variance of the unobserved shocks in test scores by group,  $\sigma_{\epsilon,k}^2$ . Our results reveal striking differences in the marginal effects of class size on test scores across groups. Group two, comprising twenty-three schools, exhibits substantial gains from reducing class size, evidenced by a prominently negative marginal effect of class size on test scores. In contrast, group three, encompassing twenty-five schools, demonstrates an insignificant, but positive effect, implying that, on average, reducing class size is detrimental to students in these schools. Class size has a moderate negative effect at schools in group one.

The test score performance of students in group two responds strongly to changes in class size. The mean test score increase for students in group two schools caused by reducing class size by one unit is 0.339 points, representing an improvement of 0.094 standard deviations. When averaged across all groups and students, the marginal effect of a one-student reduction in class size translates to an improvement of 0.092 points, equivalent to 0.026 standard deviations.

These heterogeneous effects of class-size reductions are accompanied by heterogeneous compliance across schools when determining class sizes. Table 3 displays the estimated Dirichlet prior means for the class size model by group. The Dirichlet prior mean can be interpreted in our context as a prior on the fraction of classes with each size. A higher value, all else equal, indicates that schools within that group are more likely to create a classroom of that size. The table shows that schools in Group 1 choose relatively large control class sizes, while Group 2 has relatively small control classes, and Group 3 has smaller treated class sizes. Variation in targeted treatment and control class sizes means that compliance with the STAR protocol differed across schools, with schools in Group 1 creating the largest expected difference between treated and control classrooms.

The EAMP algorithm thus uncovers distinct patterns of learning gains and class size reductions in the data. Figure 2 plots the treatment-control difference in test scores against the treatment-control difference in class sizes across schools. The areas inhabited by Groups 2 and 3 scarcely overlap. Group 2 consists of schools where students attending small classes saw large benefits, while students

**Table 3:** Dirichlet Class Size Parameters by Group

		Treatment Class Size					
Group		12	13	14	15	16	17
1		0.06	0.15	0.09	0.22	0.25	0.22
2		0.03	0.14	0.29	0.13	0.21	0.20
3		0.08	0.17	0.16	0.17	0.26	0.15

		Control Class Sizes											
Group		16	17	18	19	20	21	22	23	24	25	26	27
1		0	0.03	0.03	0.11	0	0	0.14	0.17	0.32	0.11	0.09	0
2		0	0	0	0.03	0.12	0.33	0.32	0.07	0.1	0.03	0	0
3		0.03	0.09	0	0.06	0.06	0.12	0.15	0.25	0.15	0.03	0.03	0.03

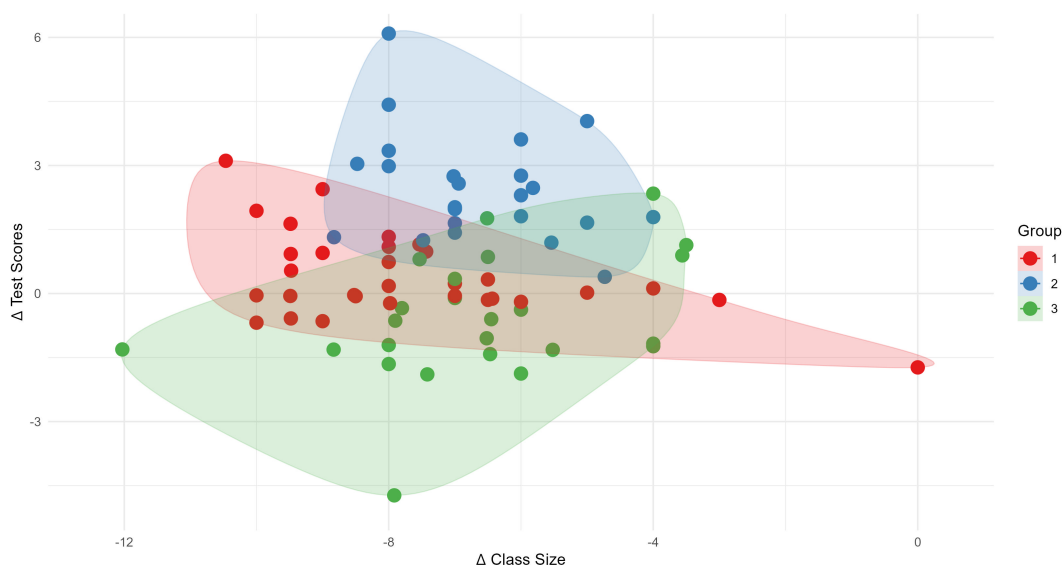
The top panel of the table displays the Dirichlet prior means for each treatment class size by group. The bottom panel reports the Dirichlet prior means for each control class size by group. Each reported number can be interpreted as a prior on the fraction of observations in each cell. Zeros indicate that a particular group does not generate classes of that size.

in Group 3 schools benefited little from smaller classes. Group 1, which has a moderate relationship between class size and test scores, is relevant for policies that reduce class size marginally from the status quo, because it is the only group with any support on small reductions in class size.

In Appendix Table B-2, we show how students in different groups differ by observable characteristics. We repeat the same analysis for teacher-specific variables in Table B-3 and for school-specific variables in Table B-4. Observable differences across groups are small, although a few are of particular interest. Group 2, which has the greatest marginal effect of class size, has the highest fraction of nonwhite students. Students in Group 2 also attend school for between four and five fewer days of the school year relative to students in the other groups, a result due to both absences and the length of the school year. Group 3, which is the only group with a positive relationship between class size and student performance, has by contrast the highest fraction of white and Asian students. Table B-3 shows that schools in Groups 2 and 3, which have opposite responses to class size reductions, display similar populations of teachers. These teachers are less



Figure 2: Test Scores and Class Size Reductions



The figure plots the average difference in test scores between treated and control classrooms along the Y-axis against the average difference in class size between treated and control classrooms along the X-axis. Each dot represents a school. The groups are determined by the EAMP algorithm described in section 3.

likely to have a masters degree but have more years of experience than Group 1 teachers.

Table B-4 shows little evidence of geographic differences in the distribution of schools across the three groups. Group 2 schools tend to have the largest cohorts, while schools in Group 1, which complied the most with the policy by creating large class size dose reductions, had small cohorts, the smallest average cohort size, and the least race-segregated schools.<sup>11</sup> These findings suggest that the schools in Tennessee that deliver the greatest marginal returns to class size have more nonwhite students and larger cohorts. The results also suggest that compliance was lowest at the most segregated schools.

<sup>11</sup>A school is deemed to be race-segregated if over 80% of its student body belongs to a single race. According to this metric, 88% of schools in the STAR experiment were segregated.

## 5.1 Individual Heterogeneity in the Impact of Class Size

These findings lend support to the subgroup analysis in [Krueger \(1999\)](#) that found that Black students and students qualifying for a free or reduced price lunch benefited more from smaller class sizes than did their white and wealthier peers. In [Table 5](#), we show that differential effects by race and income are driven both by the type of school attended and by individual heterogeneity across race. Columns (1) and (2) of [Table 5](#) display separate 2SLS regressions of test scores on class size for nonwhite and white and Asian students. Columns (3) and (4) instead run OLS regressions after interacting class size with group membership as recovered by the EAMP algorithm.<sup>12</sup> The comparison reveals that differences in the returns to class size by race identified by the 2SLS regressions are entirely driven by differences in group one. White and Asian students in Group 1 schools do not benefit from reduced class sizes, while nonwhite students experience significant gains. There is no racial difference in the return to class size at schools in Groups 2 and 3, suggesting that individual heterogeneity in the response to class size interventions is driven at least in part by the type of school attended. This is relevant for policymakers, because if a goal of a particular policy is to reduce the Black-white test score gap, then the policy must include a sizeable number of Group 1 schools.

Columns (5)-(8) of [Table 5](#) repeat the same exercise for students qualifying for free lunch and those who do not. While columns (5) and (6) seem to suggest that the relationship between test score and class size does not differ by free lunch status, columns (7) and (8) reveal that this is actually due to offsetting effects in different groups. Students in Group 1 schools who qualify for free lunch benefit significantly from a reduction in class size, while students from richer households do not benefit. But, in the remaining schools, richer students react more to class size reductions than poorer students.

Table 5: Heterogeneity in Marginal Effects of Class Size on Test Scores by Race and Income

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	By Race				By Income			
	Nonwhite	White & Asian	Nonwhite	White & Asian	Free lunch	No free lunch	Free lunch	No free lunch
Class size	-0.134 (0.025)	-0.085 (0.018)			-0.110 (0.020)	-0.093 (0.021)		
Class size × Group 1			-0.169 (0.035)	-0.023 (0.026)			-0.116 (0.028)	-0.032 (0.029)
Class size × Group 2			-0.349 (0.050)	-0.346 (0.032)			-0.311 (0.038)	-0.357 (0.036)
Class size × Group 3			0.114 (0.038)	0.115 (0.033)			0.098 (0.031)	0.123 (0.038)
Constant	-0.282 (0.086)	0.132 (0.063)	-0.271 (0.142)	0.116 (0.099)	-0.705 (0.069)	0.632 (0.072)	-0.564 (0.112)	0.477 (0.113)
Observations	1,216	2,597	1,216	2,597	1,803	2,010	1,803	2,010
R <sup>2</sup>	0.014	0.008	0.063	0.048	0.011	0.009	0.051	0.053

Columns (1)-(2) and (5)-(6) show 2SLS regression results for four subsamples of interest: nonwhite, white & Asian, students who do not qualify for free or reduced price lunch, and those who do qualify. The excluded instrument is a binary indicator for attending a small classroom. Columns (3)-(4) and (7)-(8) show the results from OLS regressions of test scores on class size interacted with group membership, as uncovered by the EAMP algorithm, for the same subsamples. The dependent variable is a simple mean of scores on three cognitive tests. All models include school fixed effects.

## 5.2 Assessing Linearity

An important assumption of our model is that class size has a linear effect on test scores in each group. We now examine this assumption more closely.<sup>13</sup> To do this, we estimate the following semiparametric regressions of test scores on class size separately by each group  $k$ :

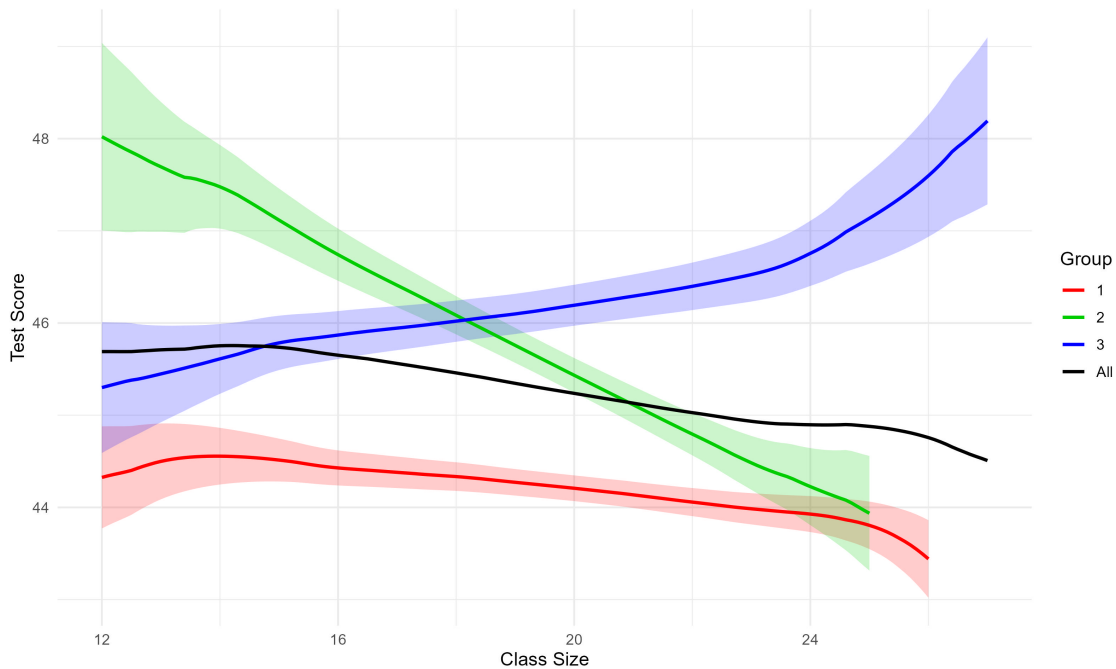
$$y_{isk} = \phi_k(n_{isk}) + x'_{isk}\theta_k + \varepsilon_{isk}. \quad (12)$$

We estimate the model in equation (12) according to the [Robinson \(1988\)](#) semi-parametric estimator, which first entails partialing out the group-specific means from both test scores and class sizes for the six cells created by  $Female \times Nonwhite \times Freelunch$  before nonparametrically regressing demeaned test scores on demeaned class size using local linear regression. We use an Epanechnikov kernel and a

<sup>12</sup>Under the assumption of no endogenous compliance within groups, this OLS specification delivers unbiased estimates of group-specific returns to class size.

<sup>13</sup>[Bandiera, Larcinese, and Rasul \(2010\)](#) find evidence of nonlinear effects of class size on academic performance among university students.

Figure 3: Nonparametric Regressions of Test Score on Class Size



The figure plots the semiparametric relationship between class size and test scores by group, as uncovered by the EAMP algorithm. All regressions use an Epanechnikov kernel and a bandwidth of 7.5. The regressions control for gender, race, and free lunch status. The black curve labeled *All* is a kernel-weighted average of the group-specific curves.

bandwidth of 7.5 students, which lets in about 50% of the data.

Figure 3 displays the results of the nonparametric regressions of test scores on class size over the support of class sizes within each group. The shaded regions represent asymptotic 95% confidence intervals. The sharp downward relationship between class size and test scores in Group 2 is noteworthy. The difference in size (of 12 students) between the largest and smallest classes in Group 2 produces a difference of approximately 1.25 sd on the exam. This relationship dwarfs the effects seen in other groups. Although the slopes vary noticeably across groups, there is limited indication of nonlinearities within any of the groups. We therefore use estimates from the linear model in Table 2 to conduct subsequent analysis.

## 6 Lessons for Scaling Class-Size Reduction Policies

We now discuss the extent to which it is possible to learn about the scalability of class-size reduction initiatives from the STAR experimental data. Prior literature has emphasized how supply-side factors, particularly a limited pool of high quality educators, pose a major challenge in scaling the STAR experiment. For example, a state-wide reduction in class sizes in California in 1996 caused schools to hire less-qualified teachers, which contributed to the disappointing “voltage effects,” whereby the policy’s effects were much smaller compared to those observed in the experimental setting (List 2022).

Our paper provides a new framework for understanding the link between experimental interventions and widespread policy adoption. We argue that the first step in scaling an intervention is aligning the econometric model with the experimental design. Even after this is done, additional challenges remain, but we show how the GRE framework can provide tools to assess the likelihood that the intervention will scale in new contexts. We conclude by discussing when targeting resources may emerge as a useful scaling alternative to universal policies.

### 6.1 Assessing the Internal Validity of the STAR Evaluation

Following the discussion in section 2.2, we argue that the 2SLS estimator is not well-aligned with the STAR experimental design.<sup>14</sup> Table B-6 compares the results from our estimated model with 2SLS that uses a treatment indicator as an instrument for class size in an IV model. The coefficient obtained by 2SLS (−0.101) does not differ much from the average marginal effect uncovered by the EAMP algorithm (−0.092), but this is more by chance than by design. Recall that in section 2.2 we showed how the 2SLS estimand is given by

$$\beta_1^{2sls} = \sum_{s \in \mathcal{S}} \beta_s \underbrace{\frac{\phi_s \bar{z}_s (1 - \bar{z}_s) (\mathbb{E}[\Delta n_{is} | s, Z_{i,s} = 1] - \mathbb{E}[\Delta n_{is} | s, Z_{i,s} = 0])}{\sum_{s \in \mathcal{S}} \phi_s \bar{z}_s (1 - \bar{z}_s) (\mathbb{E}[\Delta n_{is} | s, Z_{i,s} = 1] - \mathbb{E}[\Delta n_{is} | s, Z_{i,s} = 0])}}_{weight_s},$$

<sup>14</sup>The same is true for the reduced-form estimator, as dose heterogeneity renders the estimates difficult to interpret.

Table 6: IV Weights

Group	Marginal Effect	$\phi_g$	$\bar{Z}_g(1 - \bar{Z}_g)$	Dose Reduction	2SLS Weight
1	-0.068	0.374	0.249	-7.404	0.406
2	-0.339	0.298	0.247	-6.528	0.283
3	0.106	0.328	0.249	-6.465	0.311

The table displays instrumental variables weights computed according to equation (3). The three rightmost columns – corresponding to the fraction of individuals in each group who are treated, the fraction of the sample in each group, and the differential between treated and control class sizes in each group – are all individual components of the 2SLS weight formula.

where the school-specific weights on  $\beta_s$  depend on each schools’ endogenous dose reduction. Table 6 aggregates these weights by group and breaks them down further into their individual components.<sup>15</sup> Schools in Group 1 have the highest weight, because they created the largest dose reductions and had the highest fraction of treated observations. Lower dose reductions and treated proportions cause schools in Groups 2 and 3 to be underweighted relative to their size. 2SLS does not differ significantly from the EAMP estimate, because endogenous compliance in STAR took a particular form that caused 2SLS to overweight the schools with the most typical impact effects (Group 1) and underweight the schools in the tails of the distribution (Groups 2 and 3).

Krueger (1999) documents how dissatisfaction among parents whose children had been randomly assigned to larger classrooms led to attrition from control classes and re-randomization in later years of the experiment. An attempt to forestall parental concerns like these may have led some principals to target similar sizes for treatment and control classrooms at the program’s inception, potentially leading schools with higher impact effects to implement smaller dose reductions.<sup>16</sup> An important takeaway from this episode is that interventions are often more likely to scale successfully when they receive approval and support

<sup>15</sup>Appendix Figure B-1 plots the distribution of weights across schools within each group.

<sup>16</sup>School infrastructure may also have affected class size targets. The state of Tennessee covered the costs of hiring teachers and aides, but schools were required to provide additional classroom space if needed (Word et al. 1990).

from the local community (Agostinelli, Avitabile, and Bobba 2023). Such approval tends to be easier to obtain when the distribution of resources is perceived as equitable (Rawls 1971).

This may explain why all three groups of schools undercomplied with the STAR experimental protocols. To show this, we simulate the Tennessee STAR experiment as it was originally intended with small classes ranging from 13 to 17 students and regular classes between 22 and 25 students. We impose equal compliance by setting the Dirichlet distributions that govern class size to place equal weight on all sizes, so that  $(\eta_k^{(t)}, \eta_k^{(c)}) = (1, 1)$  for all groups  $k = 1, \dots, K$ . We then simulate class sizes and test scores and compute moments. These moments, in the third column of Table 7, show that the intended Tennessee STAR experiment would have generated a larger dose reduction, of over one student per classroom, between treated and control classrooms than was actually implemented. Endogenous compliance in STAR thus manifested itself in two ways: under-compliance overall and particularly large under-compliance at schools in Groups 2 and 3 that are associated with larger impact effects. If STAR had been implemented as intended, the reduced form effect would have been 18% higher (0.762 instead of 0.645).

Alternative forms of noncompliance could have generated still larger ATEs. We simulate an implementation of STAR that maximizes the difference between treatment and control class sizes in Groups 1 and 2 and minimizes this difference for schools in Group 3. This version of STAR would have generated an average difference in test scores between treated and control classrooms of over ten students, and a huge reduced form effect of 1.971 points (or 0.55 s.d.). The 2SLS estimate obtained from this experiment would be  $-0.199$ , more than twice the average marginal effect of  $-0.092$  estimated by our model. This occurs because 2SLS places negative weights on schools in Group 3, but these are schools which have a positive effect of class size on test scores. Researchers relying on 2SLS would therefore have erroneously concluded that reducing class size was more than twice as effective as it actually is.

These simulations show how the twin features of compliance and impact heterogeneity mean that different implementations of the STAR experiment on the same

**Table 7:** Model Fit and Counterfactual Simulations

	Data	Model	Model's Counterfactual Simulation	
			Intended Experiment	Extreme Noncompliance
Reduced Form	0.727	0.645	0.762	1.971
2SLS	-0.101	-0.089	-0.090	-0.199
Avg. class size	18.991	19.02	19.764	20.026
S.D. class size	4.067	4.081	4.371	6.420
Dose Reduction	-7.228	-7.282	-8.417	-10.384

The table presents several moments from the actual data and simulations. The first row is the reduced form effect of randomization into a small class on test scores, estimated using linear regression model with school fixed effects. The 2SLS model employs randomization into a small class as an instrument for class size in a two-stage least squares regression, incorporating school fixed effects. All regressions control for race, free lunch status, and gender. Simulated and counterfactual estimates are averages across 100 simulations. Details regarding the simulations are provided in section 6.

underlying population will generate different 2SLS estimates. While all estimates in Table 7 indicate a negative relationship between class size and test scores, the marginal effects estimated by 2SLS can differ considerably. These variable estimates can generate different rationales for action by policymakers, particularly when policies are justified on the basis of the marginal value of public funds (Mayshar 1990; Kline and Walters 2016). The 2SLS estimator is therefore poorly-aligned with an experimental design that allows for dose heterogeneity and endogenous dose responses. Furthermore, 2SLS tells researchers nothing about the heterogeneity in impact effects and compliance behavior, both of which affect scalability of the experimental intervention.

## 6.2 Assessing the Policy-Relevance of the STAR Evaluation for Scaling

GRE, by contrast, is robust to different implementations of the experiment and informs us about the distributions of compliance and impact heterogeneity. In the GRE framework, the policy-relevant treatment effect (PRTE) for moving from



policy  $a$  to policy  $a'$  can be decomposed as follows:

$$PRTE_{a',a} = \sum_{k=1}^K \rho_k \mu_k \underbrace{(\mathbb{E}[n_{isk}|k, Policy = a'] - \mathbb{E}[n_{isk}|k, Policy = a])}_{Compliance_k}, \quad (13)$$

Equation (13) shows that the expected effect of a policy intended to reduce class size hinges on schools' compliance with the policy, the structural impact effects  $\mu_k$ , and the proportion of schools belonging to each group in the population, given by  $\rho_k$ . [Al-Ubaydli et al. \(2021\)](#) identify nonrepresentative populations and nonrepresentative behavioral responses as two main threats to the scalability of experimental interventions. In our framework,  $\rho_k$  corresponds to representativeness of the population and  $Compliance_k$  corresponds to representativeness of the behavioral response. When scaling an intervention, it is important to assess whether  $\rho_k$  and  $Compliance_k$  will differ in the new context.

The history of the Tennessee STAR experiment sheds light on the potential obstacles to implementing a statewide reduction in class size. [Rockoff \(2009\)](#) explains that "only about one in five eligible schools volunteered to participate."<sup>17</sup> This reluctance to engage in the experiment may have stemmed from various factors, including disinterest or practical limitations in reducing class size. These schools, which opted out of the experiment, will likely be involved in any statewide expansion of the class-size reduction policy, which may diminish the PRTE.

When extrapolating from the STAR experiment to new contexts, it is important to assess the representativeness of the experimental sample for the new population. This corresponds to determining whether the group proportions,  $\rho_k$ , are the same in the experimental and scaled environments. The GRE model groups schools based on their compliance and impact effect heterogeneity, and we suggest that researchers use this discrete grouping to estimate multinomial models of group membership as functions of site-specific covariates. Researchers can then use such a model to predict the proportion of groups in new environments.

The estimated coefficients for two multinomial logit models of group membership are presented in Table 8. None of the coefficients are statistically significant,

---

<sup>17</sup>See also [Boyd-Zaharias et al. 2007](#) for details on the STAR experimental design.

partly due to the sample size of 79 schools. This may also reflect the fact that treatment effect heterogeneity in STAR stems less from observable school-specific factors than from unobservables. Observable characteristics may, however, be more predictive of impact effect heterogeneity in other multi-site experiments. In these settings, researchers can use GRE to classify sites based on impact and compliance heterogeneity, estimate group membership as a function of site-specific covariates, and use these multinomial models to predict group membership in new sites where they hope to scale the intervention. When researchers can accurately predict impact heterogeneity based on site-specific covariates, resource targeting becomes a promising alternative to universal policy adoption.

### 6.3 Targeting Small-Class Policies

Targeted policies are less costly than universal ones and can be better monitored, which makes them less susceptible to a “voltage drop” in implementation. If patterns in treatment effect heterogeneity are known, a targeted policy may be particularly efficient. To demonstrate this, we simulate a universal class-size reduction of five students across all schools and examine its effect on test scores and on the test-score gaps between Black and white students and between students who qualify for free lunch and those who do not. We then ask: What fraction of these effects could be achieved by a targeted policy that reduces class size by five students in the single group of schools with the largest marginal returns to class size?

Table 9 presents estimates of these counterfactual treatment effects. The first row shows that reducing class size by five for all students in the experiment would have raised test scores by 0.458 points overall, but the same intervention applied to only the 29% of schools in Group 2 would have generated gains that are 10.3% larger. The second row shows that the universal reduction of five students would have lowered the Black-white test score gap by 0.137 points, but an intervention targeted only to Group 2 schools would have reduced the gap by 0.105 points, or 76.7% of the original effect. The effects of both the universal and targeted interventions on the gap between students growing up in richer and poorer households (as measured by free lunch status) are small. These simulations suggest

Table 8: Predictive Models

	Model 1		Model 2	
	Group 2	Group 3	Group 2	Group 3
Suburban	0.262 (0.714)	-0.143 (0.691)	-0.701 (1.034)	-0.470 (0.972)
Urban	-0.431 (0.963)	-1.529 (1.180)	-0.738 (1.007)	-1.788 (1.216)
Inner City	-0.208 (0.708)	-1.124 (0.776)	-2.456 (1.833)	-2.999 (1.873)
Cohort Size			0.157 (1.109)	0.257 (1.149)
Teacher Experience			0.177 (0.119)	0.127 (0.112)
Teacher: Master's Degree			-1.427 (1.240)	-1.399 (1.151)
Days in School			-2.181 (2.933)	-0.051 (3.063)
% Black			2.819 (2.300)	1.677 (2.284)
% Free lunch			-1.038 (2.313)	0.653 (2.112)
AIC	184.170		199.428	
Observations	79		79	

The table shows the estimated coefficients from two multinomial logit models predicting group membership as a function of school-specific covariates. Group one is the reference group. School-specific averages of individual- and teacher-specific variables are used. Standard errors are in parentheses.

Table 9: Comparing Universal and Targeted Interventions

	Universal Intervention			Targeted Intervention	
	Data	Effect on Scores	Percent Reduction	Targeted Effect on Scores	Fraction of Universal Effect
Overall		0.458		0.506	1.103
Black-white gap	-1.395	0.137	9.8 %	0.105	0.767
Free lunch gap	-1.982	0.013	0.7 %	0.023	1.748

The table shows the effect of reducing class size by 5 students for everyone (universal intervention) with the effect of reducing class size by 5 students for schools in group two only (targeted intervention) on three separate outcomes. The first row simulates the effect on average test scores, while the second and third rows simulate the effects on the gaps in test scores between Black and white students and between students who qualify for free lunch and those who do not. Estimates are averages across 100 simulations. Details regarding the implementation of these counterfactual interventions are presented in section 6.

that there may be large cost efficiencies in policy implementation if resources can be targeted where they are likely to have the largest effects.

## 7 Conclusion

This paper develops and implements a novel empirical framework that can be applied to large families of randomized controlled trials with potentially endogenous compliance and heterogeneous treatment effects. We make use of the Grouped Random Effects approach of [Adusumilli \(2020\)](#), a generalization of the Grouped Fixed Effects framework of [Bonhomme and Manresa \(2015\)](#), to identify patterns in treatment effect heterogeneity and compliance behavior across multi-site experiments. We then show how this framework can help researchers use the experimental data to learn about the scalability of the policy.

We apply this framework to data from the Tennessee STAR experiment to evaluate the causal effects of class size on student learning. The grouping algorithm uncovers marked differences in the effectiveness of class size reductions across schools. We find that 29% of schools in the experiment generated an increase in test scores of 0.09 sd on average for each one-student reduction in class size. However, the effect of class size on test scores was small in the remaining schools.

Our findings contribute to the debate on the value of universal versus targeted interventions. We show how a well-targeted investment to reduce class size can generate similar overall effects on test scores as a universal program. The same program would also reduce the black-white test score gap by 76.7% as much as

a universal program. Targeted interventions may also have other benefits not studied in this paper, the most important being that they can be scaled more easily while retaining quality.

By presenting evidence of greater heterogeneity in both impact effects and compliance behavior than previously acknowledged, we posit that the STAR experiment can reconcile contrasting conclusions from previous studies of class size effects on children’s learning. We show how different implementations of STAR could have generated a wide range of estimates regarding class size effectiveness when using common regression estimators like two-stage least squares (2SLS). Our findings underscore the significance of designing pilot interventions capable of revealing the extent of treatment effect heterogeneity before advocating for the scaling of interventions across different contexts, situations, and populations.

More research is still needed to understand why educational interventions cause such heterogeneous effects across schools. This paper provides a first step at documenting the wide variability in the production function relating schooling investments to knowledge. Our hope is that future researchers use these tools to generate further insights on who responds endogenously to investments, who benefits from them, and why.

## References

- Achilles, C. M., J. B. Zaharias, B. A. Nye, and D. Fulton. 1995. “Analysis of Policy Application of Experimental Results: Project Challenge.” Nashville, TN: Tennessee State University, Center of Excellence for Research in Basic Skills. Retrieved from <http://eric.ed.gov/?id=ED393151>.
- Adusumilli, Karun. 2020. “Unobserved Heterogeneity, Grouped Random Effects and the EAMP Algorithm.” *Unpublished Manuscript*.
- Agostinelli, Francesco, Ciro Avitabile, and Matteo Bobba. 2023. “Enhancing Human Capital in Children: A Case Study on Scaling.” Technical Report, National Bureau of Economic Research.
- Al-Ubaydli, Omar, Min Sok Lee, John A List, Claire L Mackevicius, and Dana Suskind. 2021. “How Can Experiments Play a Greater Role in Public Policy?”

- Twelve Proposals from an Economic Model of Scaling." *Behavioural Public Policy* 5 (1): 2–49.
- Al-Ubaydli, Omar, John A. List, and Dana Suskind. 2020. "2017 Klein Lecture: The Science of Using Science: Toward an Understanding of the Threats to Scalability." *International Economic Review* 61 (4): 1387–1409.
- Angrist, Joshua D., and Guido W. Imbens. 1995. "Two-Stage Least Squares Estimation of Average Causal Effects in Models with Variable Treatment Intensity." *Journal of the American Statistical Association* 90 (430): 431–442.
- Angrist, Joshua D., and Victor Lavy. 1999. "Using Maimonides' Rule to Estimate the Effect of Class Size on Scholastic Achievement." *The Quarterly Journal of Economics* 114 (2): 533–575.
- Angrist, Joshua D, Victor Lavy, Jetson Leder-Luis, and Adi Shany. 2019. "Maimonides rule redux." *American Economic Review: Insights* 1 (3): 309–324.
- Bandiera, Oriana, Valentino Larcinese, and Imran Rasul. 2010. "Heterogeneous Class Size Effects: New Evidence from a Panel of University Students." *The Economic Journal* 120 (549): 1365–1398.
- Bergman, Peter, Raj Chetty, Stefanie DeLuca, Nathaniel Hendren, Lawrence F Katz, and Christopher Palmer. 2019. "Creating Moves to Opportunity: Experimental Evidence on Barriers to Neighborhood Choice." Technical Report, National Bureau of Economic Research.
- Blei, David M, Alp Kucukelbir, and Jon D McAuliffe. 2017. "Variational inference: A review for statisticians." *Journal of the American statistical Association* 112 (518): 859–877.
- Bohrnstedt, George W, and Brian M Stecher. 1999. "Class Size Reduction in California 1996-1998: Early Findings Signal Promise and Concerns. Report Summary."
- . 2002. "What We Have Learned about Class Size Reduction in California. Capstone Report."
- Bonhomme, Stéphane, Thibaut Lamadon, and Elena Manresa. 2022. "Discretizing unobserved heterogeneity." *Econometrica* 90 (2): 625–643.

- Bonhomme, Stéphane, and Elena Manresa. 2015. "Grouped patterns of heterogeneity in panel data." *Econometrica* 83 (3): 1147–1184.
- Borghesan, Emilio, and Gabrielle Vasey. 2024. "The Marginal Returns to Distance Education: Evidence from Mexico's Telesecundarias." *American Economic Journal: Applied Economics* 16 (1): 253–285.
- Boyd-Zaharias, J, J Finn, R Fish, and S Gerber. 2007. "Project STAR and Beyond: Database User's Guide." *HEROS Inc. and University of New York at Buffalo technical report*.
- Card, David, and Alan B. Krueger. 1992a. "Does School Quality Matter? Returns to Education and the Characteristics of Public Schools in the United States." *Journal of Political Economy* 100 (1): 1–40.
- . 1992b. "School Quality and Black-White Relative Earnings: A Direct Assessment." *The Quarterly Journal of Economics* 107 (1): 151–200.
- Carneiro, Pedro, James J Heckman, and Edward J Vytlačil. 2011. "Estimating Marginal Returns to Education." *American Economic Review* 101 (6): 2754–2781.
- Chetty, Raj, John N. Friedman, Nathaniel Hilger, Emmanuel Saez, Diane Whitmore Schanzenbach, and Danny Yagan. 2011. "How Does Your Kindergarten Classroom Affect Your Earnings? Evidence from Project STAR." *The Quarterly Journal of Economics* 126 (4): 1593–1660.
- Chetty, Raj, John N Friedman, and Jonah E Rockoff. 2014. "Measuring the Impacts of Teachers II: Teacher value-Added and Student Outcomes in Adulthood." *American Economic Review* 104 (9): 2633–2679.
- DellaVigna, Stefano, Elizabeth Linos, and Woojin Kim. 2023. "Bottlenecks for evidence adoption." Working Paper.
- Finn, Jeremy D, and Charles M Achilles. 1990. "Answers and Questions about Class Size: A Statewide Experiment." *American Educational Research Journal* 27 (3): 557–577.
- Folger, John, and Carolyn Breda. 1989. "Evidence from Project STAR about Class Size and Student Achievement." *Peabody Journal of Education* 67 (1): 17–33.

- Gechter, Michael, Keisuke Hirano, Jean Lee, Mahreen Mahmud, Orville Mondal, Jonathan Morduch, Saravana Ravindran, and Abu S Shonchoy. 2023. "Site Selection for External Validity: Theory and an Application to Mobile Money in South Asia."
- Hanushek, Eric A. 1986. "The Economics of Schooling: Production and Efficiency in Public Schools." *Journal of Economic Literature* 24 (3): 1141–1177.
- . 1997. "Assessing the Effects of School Resources on Student Performance: An Update." *Educational Evaluation and Policy Analysis* 19 (2): 141–164.
- Heckman, James. 1992. "Randomization and Social Policy Evaluation." In *Evaluating Welfare and Training Programs*. Edited by C. F. Manski and I. Garfinkel. Harvard University Press.
- Heckman, James, Anne Layne-Farrar, and Petra Todd. 1995, September. "Does Measured School Quality Really Matter? An Examination of the Earnings-Quality Relationship." Working paper 5274, National Bureau of Economic Research.
- Heckman, James J., Sergio Urzua, and Edward Vytlacil. 2006. "Understanding Instrumental Variables in Models with Essential Heterogeneity." *The Review of Economics and Statistics* 88 (3): 389–432.
- Heckman, James J, and Edward J Vytlacil. 2001. "Policy-Relevant Treatment Effects." *American Economic Review* 91 (2): 107–111.
- Hippel, Paul von, and Chandi Wagner. 2018. "Does a Successful Randomized Experiment Lead to Successful Policy? Project Challenge and What Happened in Tennessee After Project STAR (March 31, 2018)." Available at SSRN: <https://ssrn.com/abstract=3153503>.
- Hoxby, Caroline M. 2000. "The Effects of Class Size on Student Achievement: New Evidence from Population Variation\*." *The Quarterly Journal of Economics* 115 (4): 1239–1285 (11).
- Jepsen, Christopher, and Steven Rivkin. 2009. "Class Size Reduction and Student Achievement: The Potential Tradeoff between Teacher Quality and Class Size." *The Journal of Human Resources* 44 (1): 223–250.



- Kline, Patrick, Evan K Rose, and Christopher R Walters. 2022. "Systemic Discrimination Among Large US Employers." *The Quarterly Journal of Economics* 137 (4): 1963–2036.
- Kline, Patrick, and Christopher Walters. 2021. "Reasonable Doubt: Experimental Detection of Job-Level Employment Discrimination." *Econometrica* 89 (2): 765–792.
- Kline, Patrick, and Christopher R Walters. 2016. "Evaluating Public Programs with Close Substitutes: The Case of Head Start." *The Quarterly Journal of Economics* 131 (4): 1795–1848.
- Krueger, Alan B. 1999. "Experimental Estimates of Education Production Functions." *The Quarterly Journal of Economics* 114 (2): 497–532.
- Krueger, Alan B, and Diane M Whitmore. 2001. "The Effect of Attending a Small Class in the Early Grades on College-test Taking and Middle School Test Results: Evidence from Project STAR." *The Economic Journal* 111 (468): 1–28.
- List, John A. 2022. *The Voltage Effect: How to Make Good Ideas Great and Great Ideas Scale*. Currency.
- Mayshar, Joram. 1990. "On Measures of Excess Burden and Their Application." *Journal of Public Economics* 43 (3): 263–289.
- Mishel, Lawrence, and Richard Rothstein. 2002. In *The class size debate*, 1–102. Economic Policy Institute.
- Neal, Radford M, and Geoffrey E Hinton. 1998. "A View of the EM Algorithm that Justifies Incremental, Sparse, and other Variants." In *Learning in Graphical Models*, 355–368. Springer.
- Rawls, John. 1971. *A Theory of Justice: Original Edition*. Harvard University Press.
- Rivkin, Steven G., Eric A. Hanushek, and John F. Kain. 2005. "Teachers, Schools, and Academic Achievement." *Econometrica* 73 (2): 417–458.
- Robinson, Peter M. 1988. "Root-N-Consistent Semiparametric Regression." *Econometrica: Journal of the Econometric Society*, pp. 931–954.

- Rockoff, Jonah. 2009. "Field Experiments in Class Size from the Early Twentieth Century." *Journal of Economic Perspectives* 23 (4): 211–30 (December).
- Rose, Evan K., and Yotam Shem-Tov. 2021. "How Does Incarceration Affect Reoffending? Estimating the Dose-Response Function." *Journal of Political Economy* 129 (12): 3302–3356.
- Schanzenbach, Diane Whitmore. 2006. "What Have Researchers Learned from Project STAR?" *Brookings Papers on Education Policy*, no. 9:205–228.
- . 2016. *Long-term Impacts of Class Size Reduction*. In Peter Blatchford, et al., eds., *International Perspectives on Class Size*. Routledge.
- Stecher, Brian M, and George W Bohrnstedt. 2000. "Class Size Reduction in California: Summary of the 1998-99 Evaluation Findings."
- Walters, Christopher R. 2018. "The Demand for Effective Charter Schools." *Journal of Political Economy* 126 (6): 2179–2223.
- Word, Elizabeth, John Johnston, Helen P Bain, B DeWayne Fulton, Jayne B Zaharias, Charles M Achilles, Martha N Lintz, John Folger, and Carolyn Breda. 1990. "The State of Tennessee's Student/Teacher Achievement Ratio (STAR) Project." *Tennessee Board of Education*.

# Appendices

## A Details of the EAMP algorithm

This appendix describes the steps of the EAMP algorithm used to estimate the GRE model.

### Step E: Expectation

By the Donsker-Varadhan variational formula, the optimal value of  $q_{sk}(\boldsymbol{\alpha}_s)$  is just the posterior distribution of  $\boldsymbol{\alpha}_s$ , as implied by the likelihood  $p(\mathbf{y}_s, \mathbf{n}_s^{(c)}, \mathbf{n}_s^{(t)} | \boldsymbol{\alpha}_s, \mathbf{x}, \boldsymbol{\sigma}_\epsilon^2, \theta)$  and the prior  $\pi(\boldsymbol{\alpha}_s | \boldsymbol{\rho}_k)$ . Since the prior is conjugate to the likelihood, the posterior can be computed very quickly. To characterize the posterior, we first note that due to the structure of the model, the posterior is separable:

$$q_{sk}(\boldsymbol{\alpha}_s) = q_{sk}(\boldsymbol{\beta}_s) \cdot q_{sk}(\mathbf{p}_s^{(c)}) \cdot q_{sk}(\mathbf{p}_s^{(t)}).$$

We can then update each of these quantities separately as follows. The update to the posterior distribution of  $\boldsymbol{\beta}_s$  is given by

$$q_{sk}(\boldsymbol{\beta}_s) \equiv \prod_{i:s(i)=s} q_{sk}(\beta_{isk}) \leftarrow \prod_{i:s(i)=s} N(\beta_{isk} | \mu_{isk}, \Sigma_{isk})$$

where, for each  $i$  such that  $s(i) = s$ , we update

$$\begin{aligned} \Sigma_{isk} &\leftarrow \left( \Sigma_k^{-1} + \frac{n_{isk} \mathbf{n}_{isk}^\top}{\sigma_{\epsilon,k}^2} \right)^{-1}, \\ \mu_{isk} &\leftarrow \Sigma_{isk} \left( \Sigma_k^{-1} \mu_k + \frac{n_{isk} (y_{isk} - x'_{isk} \theta)}{\sigma_{\epsilon,k}^2} \right) \end{aligned} \tag{A-1}$$

The update to  $q_{sk}(p_s^{(c)})$  is given by

$$q_{sk}(p_{sk}^{(c)}) \leftarrow \text{Dirichlet}(p_{sk}^{(c)} | \eta_k^{(c)}),$$

where  $\eta_k^{(c)} \equiv \{\eta_{k1}^{(c)}, \dots, \eta_{kL}^{(c)}\}$  denotes the posterior values of  $\eta_k^{(c)}$ , and is given by

$$\eta_{kl}^{(c)} = \eta_{kl} + \sum_{g=1}^{N^{(c)}(s)} \mathbb{I}\{n_{skg}^{(c)} \equiv 16 + m\}, \text{ for each } m = 0, \dots, M.$$

The update to  $q_{sk}(p_{sk}^{(t)})$  is analogous.

### Step A: Assignment

We assign each observation to one of the  $K$  groups by maximizing (11) with respect to  $w_s(k)$ . As in Adusumilli (2020), group assignments are obtained as the solution to the following problem:

$$k(s) \leftarrow \arg \max_k I_{sk}; \quad I_{sk} := \ln \frac{p(\mathbf{y}_s, \mathbf{n}_s^{(c)}, \mathbf{n}_s^{(t)} | \boldsymbol{\alpha}_s, \mathbf{x}, \sigma_{\epsilon}^2, \theta) \pi(\boldsymbol{\alpha}_s | \boldsymbol{\rho}_k)}{q_{sk}(\boldsymbol{\alpha}_s)}. \quad (\text{A-2})$$

Since the posterior distribution,  $q_{sk}(\boldsymbol{\alpha}_s)$ , is known from Step E, we can obtain an analytical expression for  $I_{sk}$ . To obtain this expression, note that

$$\ln I_{sk} = \frac{p(\mathbf{y}_s | \boldsymbol{\beta}_s, \mathbf{x}, \sigma_{\epsilon, k(s)}^2, \theta) \cdot p(\mathbf{n}_s^{(c)} | \mathbf{p}_s^{(c)}) \cdot p(\mathbf{n}_s^{(t)} | \mathbf{p}_s^{(t)}) \pi(\boldsymbol{\beta}_s | \boldsymbol{\rho}_k) \cdot \pi(\mathbf{p}_s^{(c)} | \boldsymbol{\eta}_k^{(c)}) \cdot \pi(\mathbf{p}_s^{(t)} | \boldsymbol{\eta}_k^{(t)})}{q_{sk}(\boldsymbol{\beta}_s) \cdot q_{sk}(\mathbf{p}_s^{(c)}) \cdot q_{sk}(\mathbf{p}_s^{(t)})}, \quad (\text{A-3})$$

which can be grouped into three separate terms:

$$\begin{aligned} \ln \frac{p(\mathbf{y}_s | \boldsymbol{\beta}_s, \mathbf{x}, \sigma_{\epsilon, k(s)}^2, \theta) \pi(\boldsymbol{\beta}_s | \boldsymbol{\rho}_k)}{q_{sk}(\boldsymbol{\beta}_s)} &= -\frac{1}{2} \sum_{i:s(i)=s} \left\{ \frac{(y_{isk} - x'_{isk} \theta)^2}{\sigma_{\epsilon, k}^2} + \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k - \boldsymbol{\mu}_{k_{si}}^\top \boldsymbol{\Sigma}_{k_{si}}^{-1} \boldsymbol{\mu}_{k_{si}} \right\} \\ &+ \frac{1}{2} \sum_{i:s(i)=s} \{ \ln |\boldsymbol{\Sigma}_{k_{si}}| - \ln |\boldsymbol{\Sigma}_k| \} - \frac{1}{2} \ln \sigma_{\epsilon, k}^2 + \text{const} \end{aligned} \quad (\text{A-4})$$

$$\begin{aligned} \ln \frac{p(\mathbf{n}_s^{(c)} | \mathbf{p}_s^{(c)}) \pi(\mathbf{p}_s^{(c)} | \boldsymbol{\eta}_k^{(c)})}{q_{sk}(\mathbf{p}_s^{(c)})} &= \ln \frac{N_s^{(c)}!}{\prod_{j=1}^J n_{skj}^{(c)}!} \frac{B(\eta_k^{(c)} + n_{sk}^{(c)})}{B(\eta_k^{(c)})} \\ &= \ln B(\eta_k^{(c)} + n_{sk}^{(c)}) - \ln B(\eta_k^{(c)}) + \text{const} \end{aligned} \quad (\text{A-5})$$

$$\ln \frac{p(\mathbf{n}_s^{(t)} | \mathbf{p}_s^{(t)}) \pi(\mathbf{p}_s^{(t)} | \boldsymbol{\eta}_k^{(t)})}{q_{sk}(\mathbf{p}_s^{(t)})} = \ln \frac{N_s^{(t)}!}{\prod_{j=1}^J n_{skj}^{(t)}!} \frac{B(\eta_k^{(t)} + n_{sk}^{(t)})}{B(\eta_k^{(t)})}$$

$$= \ln B(\eta_k^{(t)} + n_{sk}^{(t)}) - \ln B(\eta_k^{(t)}) + \text{const} \quad (\text{A-6})$$

where  $n_{sk}^{(c)}$  and  $n_{sk}^{(t)}$  are the count vectors containing the number of classes of each type in school  $s$ ,  $N_s^{(c)}$  and  $N_s^{(t)}$  are the number of control and treatment classrooms in school  $s$ ,  $B(\eta) = \frac{\prod_{j=1}^J \Gamma(\eta_j)}{\Gamma(\sum_{j=1}^J \eta_j)}$ , and  $\Gamma(n) = (n-1)!$ . We compute  $I_{sk}$  as the sum of (A-4) (A-5) and (A-6) separately for each school and group, and then we assign each school to the group with the greatest value of  $I_{sk}$ .

### Step M: Maximization

The maximization step updates the estimates of the nonrandom parameters,  $\theta$  and  $\sigma_\epsilon^2$ . We compute  $\theta$  by solving

$$\begin{aligned} & \max_{\theta} \sum_{s=1}^S \sum_{k=1}^K w_s(k) E_{q_{sk}(\cdot)} [\ln p(\mathbf{y}_s, \mathbf{n}_s^{(c)}, \mathbf{n}_s^{(t)} | \boldsymbol{\alpha}_s, \mathbf{x}, \sigma_{\epsilon,k}^2, \theta)] = \\ & \max_{\theta} \sum_{s=1}^S \sum_{k=1}^K w_s(k) E_{q_{sk}(\cdot)} [\ln p(\mathbf{y}_s | \boldsymbol{\alpha}_s, \mathbf{x}, \sigma_{\epsilon,k}^2, \theta) + \ln p(\mathbf{n}_s^{(t)} | \boldsymbol{\alpha}_s) + \ln p(\mathbf{n}_s^{(c)} | \boldsymbol{\alpha}_s)] = \\ & \max_{\theta} \sum_{s=1}^S \sum_{k=1}^K w_s(k) E_{q_{sk}(\cdot)} [\ln p(\mathbf{y}_s | \boldsymbol{\alpha}_s, \mathbf{x}, \sigma_{\epsilon,k}^2, \theta)] \end{aligned} \quad (\text{A-7})$$

Since  $y_{isk}$  is normally distributed conditional on covariates, class size, and the random coefficients, the solution to (A-7) is a linear projection:

$$\theta = \left( \sum_{i=1}^N x_{isk} x'_{isk} \right)^{-1} \left( \sum_{i=1}^N (y_{isk} - \sum_{k=1}^K w_{s(i)}(k) n_{isk} E_{q_{sk}(\cdot)}[\beta_{isk}]) \right), \quad (\text{A-8})$$

where  $w_{s(i)}(k)$  is an indicator for whether student  $i$ 's school ( $s$ ) belongs to group  $k$ , and  $E_{q_{sk}(\cdot)}[\beta_{isk}]$  is the posterior mean of  $\beta_{isk}$ , specifically  $\mu_{isk}$  from equation (A-1). Note that because we are taking expectation with respect to the posterior distribution conditional on observing the data (including class size), class size is in the conditioning set so that

$$E_{q_{sk}(\cdot)}[w_{s(i)}(k)n_{isk}\beta_{isk}] = w_{s(i)}(k)n_{isk}E_{q_{sk}(\cdot)}[\beta_{isk}] ,$$

which delivers the formula in (A-8).

We then compute  $\sigma_{\epsilon,k}^2$  for  $k = 1, \dots, K$  by solving

$$\begin{aligned} \sigma_{\epsilon,k}^2 &= \frac{1}{n_k} \sum_{i=1}^N w_{s(i)}(k) E_{q_{isk}(\cdot)} [(y_{isk} - \beta_{isk}n_{isk} - x'_{isk}\theta)^2] , \\ &= \frac{1}{n_k} \sum_{i=1}^N w_{s(i)}(k) [(y_{isk} - \mu_{isk}n_{isk} - x'_{isk}\theta)^2 + \Sigma_{isk}n_{isk}^2] , \end{aligned}$$

where  $n_k$  is the number of students in group  $k$ , and  $E_{q_{isk}(\cdot)}(\beta_{isk})$  and  $Var(\beta_{isk})$  are posterior means and variances computed in (A-1).

### Step P: Propagation

The prior is from the exponential family. Hence, as in Adusumilli (2020), updating the prior parameters involves matching the sufficient statistics of the exponential family between the prior and average posterior. Due to separability of both the prior and posterior, we can separately update the prior parameters  $\gamma_k$ ,  $\eta_k^{(c)}$ , and  $\eta_k^{(t)}$ .

We update the mean and variance for each group as follows:

$$\mu_k \leftarrow \frac{1}{n_k} \sum_s w_s(k) \sum_{i:s(i)=s} \mu_{isk}, \quad (\text{A-9})$$

$$\Sigma_k \leftarrow \frac{1}{n_k} \sum_s w_s(k) \sum_{i:s(i)=s} \{\Sigma_{isk} + \mu_{isk}\mu_{isk}^\top\} - \mu_k\mu_k^\top, \quad (\text{A-10})$$

where  $n_k$  is the number of observations (students) in group  $k$ . If, in the process of optimization, a group turns out to be empty, we do not update the posterior for that group.

To update  $\eta_k^{(c)}$ , we match the posterior average and prior moments of  $\ln p_{skl}^{(c)}$  for each  $l$  as these are the sufficient statistics of the Dirichlet family. This implies that

$\eta_k^{(c)} \equiv (\eta_{k1}^{(c)}, \dots, \eta_{kL}^{(c)})$  can be obtained as the solution to the system of  $L$  equations. Denote by  $\tilde{\eta}_{kj}^{(c)}$  the updated parameter for control class type  $j$  in group  $k$ . Then the system of  $M$  equations in  $M$  unknowns for control group  $k$  is given by

$$\psi\left(\tilde{\eta}_{km}^{(c)}\right) - \psi\left(\sum_{m=1}^M \tilde{\eta}_{km}^{(c)}\right) = \frac{1}{N_{sk}} \sum_s w_s(k) \left\{ \psi\left(\eta_{skm}^{(c)}\right) - \psi\left(\sum_{m=1}^M \eta_{skm}^{(c)}\right) \right\}, \text{ for each } m = 1, \dots, M \quad (\text{A-11})$$

where  $\psi(\cdot)$  denotes the Digamma function,  $N_{sk}$  is the number of schools in group  $k$ , and the right hand side variables,  $\eta_{sk}^{(c)}$ , were obtained in the E-step as the sum of the Dirichlet prior and the vector of class size counts for school  $s$ :  $\eta_{sk}^{(c)} = \eta_k^{(c)} + \mathbf{n}_{sk}$ .

The system of equations for treatment group  $k$  is analogous:

$$\psi\left(\tilde{\eta}_{kl}^{(t)}\right) - \psi\left(\sum_{l=1}^L \tilde{\eta}_{kl}^{(t)}\right) = \frac{1}{N_{sk}} \sum_s w_s(k) \left\{ \psi\left(\eta_{skl}^{(t)}\right) - \psi\left(\sum_{l=1}^L \eta_{skl}^{(t)}\right) \right\}, \text{ for each } l = 1, \dots, L \quad (\text{A-12})$$

Note that, because the test score model generates scores for students, but the class size model generates counts of classes,  $N_{sk}$  in equations (A-11) and (A-12) refer to the number of schools.

These systems of equations are solved for each group and treatment/control status to obtain  $2XK$  posterior parameter vectors,  $\tilde{\eta}_k^{(c)}$  and  $\tilde{\eta}_k^{(t)}$  for each group,  $k = 1, \dots, K$ .

## B Additional Tables and Figures - Main Sample

Table B-1: Summary Statistics

	Mean	Standard Deviation
Test Score	45.30	3.586
Class Size	18.991	4.067
Female	0.489	0.5
Not White/Asian	0.319	0.466
Eligible for Free Lunch	0.473	0.499
Treated	0.465	0.499

The table presents descriptive statistics for the sample of 3813 school children used in estimation. Test Score is the average of student scores on the math, reading, and word skills SAT exams administered at the end of kindergarten. Female, Nonwhite, and Eligible for Free Lunch are all binary variables.

Table B-2: Student Characteristics by Group

Group	Female	Nonwhite	Free lunch	Avg. Days Present	S.D. Days Present
1	0.48	0.32	0.48	158.0	25.7
2	0.49	0.37	0.48	153.4	27.6
3	0.50	0.30	0.48	157.3	24.6

The table shows average student characteristics for students attending schools in each group. Avg. Days Present is determined by both the length of the school year and student absences.



**Table B-3: Teacher Characteristics by Group**

Group	Masters Degree	Experience	Nonwhite
1	0.42	8.80	0.17
2	0.31	9.52	0.17
3	0.31	9.49	0.14

The table shows average teacher characteristics for teachers at schools in each group. Experience is measured in years.

**Table B-4: School Characteristics by Group**

Group	Mean Class Size	S.D. Class Size	Cohort Size	S.D. Cohort Size	% Race-Segregated	% SES-Segregated
1	20.92	4.23	78.5	19.3	0.81	0.29
2	20.04	3.52	82.3	29.1	0.96	0.39
3	19.92	4.00	79.9	31.3	0.92	0.36
Group	Inner City	Rural	Suburban	Urban		
1	8	13	6	4		
2	5	10	6	2		
3	3	15	6	1		

The table shows school characteristics by group membership. A school is classified as race-segregated if over 80% of the student body belongs to a single racial or ethnic group. A school is classified as ses-segregated if either over 80% or under 20% of the student body qualifies for free or reduced price lunch.

Table B-5: Estimates of the Model for Test Scores

Group	$\mu_k$	$\Sigma_k$	$\sigma_{\epsilon,k}^2$	Schools	Students
1	-0.068 (0.061)	0.042 (0.065)	4.933 (1.55)	31	1425
2	-0.339 (0.062)	0.235 (0.15)	6.442 (1.86)	23	1137
3	0.106 (0.059)	0.001 (0.07)	11.852 (2.07)	25	1251
Covariates					
Female	0.666 (0.14)				
Nonwhite	-1.109 (0.29)				
Free Lunch	-1.683 (0.20)				
N	3,813				

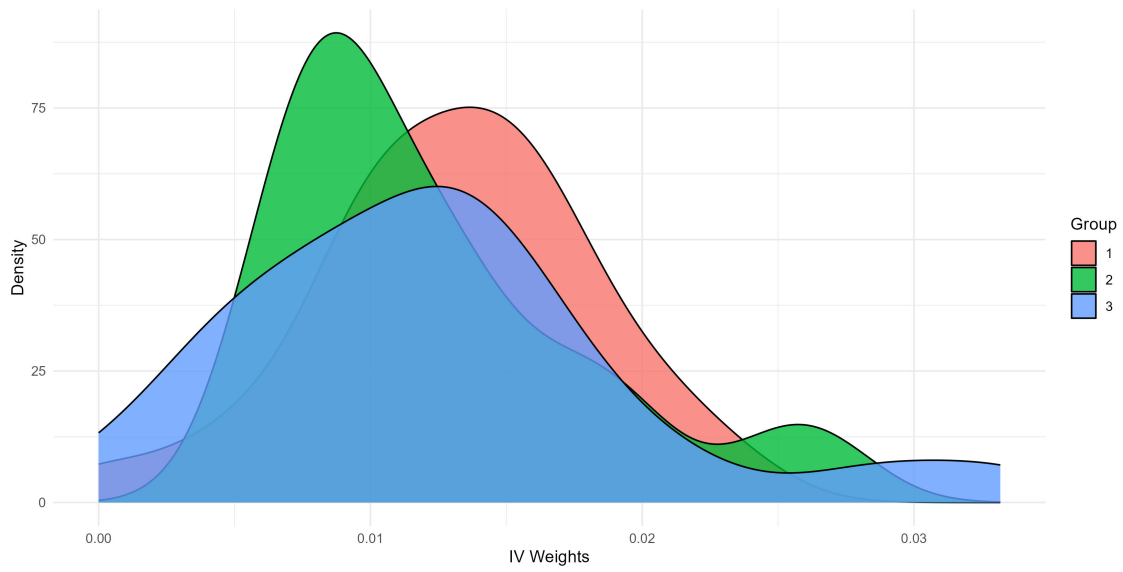
The table shows the estimated model parameters governing the effect of class size on test scores. Regressions include controls for gender, race, and free lunch status. Controls are constrained to be equal across groups. The class size effects represent the effect of a one-unit increase in class size on test score performance. Bootstrapped standard errors from 95 bootstrapped data sets are in parentheses. Bootstrapping involves first sampling schools from the set of 79 schools that participated in the experiment and then sampling individuals within each school.

Table B-6: 2SLS Model

	2SLS
Class Size	-0.101 (0.014)
Female	0.650 (0.099)
Nonwhite	-1.057 (0.209)
Free Lunch	-1.776 (0.119)
First Stage F-Statistic	42237.32
Observations	3823

The table presents estimates from a two-stage least squares regression that uses the binary randomization into a small classroom as the excluded instrument. Heteroskedasticity-robust standard errors are in parentheses.

Figure B-1: Density of School-Specific 2SLS Weights



The figure plots the density of weights across schools implied by the 2SLS estimand in equation (3). Separate densities are estimated for each group as uncovered the EAMP algorithm.

## C Results From Sample With Regular + Aide Classrooms

Table C-1: Class Size Marginal Effects by Group

Group	$\mu_k$	$\Sigma_k$	$\sigma_{\epsilon,k}^2$	Schools	Students
1	-0.094	0.019	5.32	28	2046
2	-0.313	0.192	8.364	21	1609
3	0.056	0.001	10.652	30	2247
Avg. Effect	-0.097	0.081	8.180		
N	5902				

The table shows the estimated model parameters governing the effect of class size on test scores. Regressions include controls for gender, race, and free lunch status. Controls are constrained to be equal across groups. The class size effects represent the effect of a one-unit increase in class size on test score performance.

Table C-2: Treatment Class Size Support by Group

Group	Class Size (Number of Students)					
	12	13	14	15	16	17
1	0.05	0.15	0.08	0.20	0.26	0.26
2	0.03	0.24	0.17	0.15	0.27	0.14
3	0.09	0.11	0.24	0.17	0.22	0.16

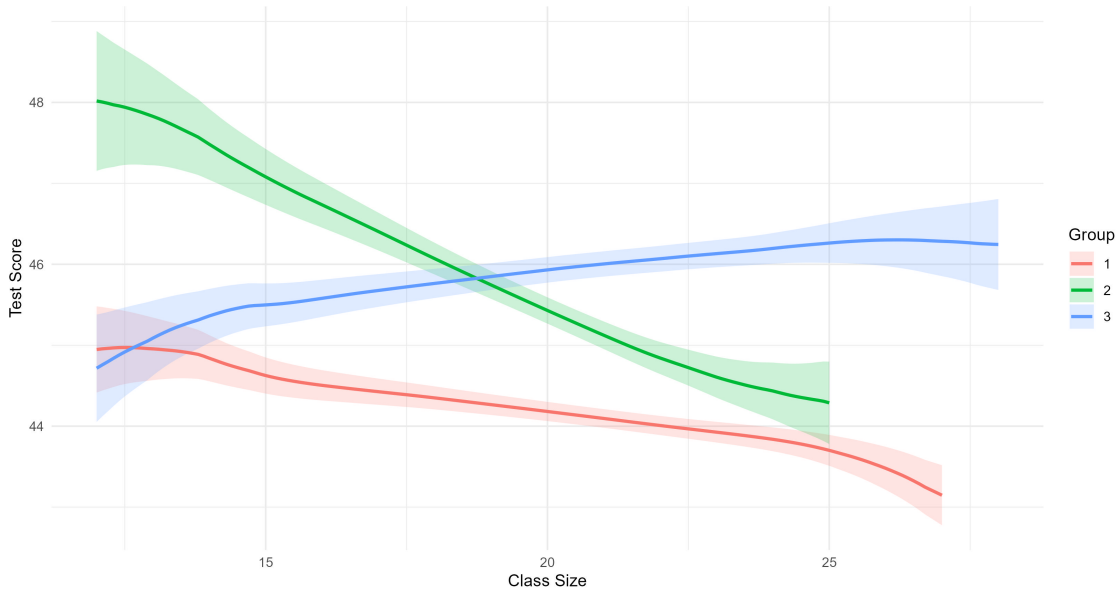
The table shows the Dirichlet prior means for each treatment class size by group. Each reported number can be interpreted as the fraction of observations in each cell. Zeros indicate that a particular group does not generate classes of that size.

Table C-3: Control Class Size Support by Group

Group	Class Size (Number of Students)													
	15	16	17	18	19	20	21	22	23	24	25	26	27	28
1	0.02	0	0.02	0.03	0.05	0.05	0.09	0.17	0.19	0.23	0.08	0.06	0.02	0
2	0	0	0.07	0.02	0.12	0.05	0.31	0.18	0.13	0.07	0.05	0	0	0
3	0	0.01	0	0	0.06	0.06	0.04	0.26	0.23	0.20	0.04	0.03	0.07	0.01

The table shows the Dirichlet prior means for each control class size by group. Each reported number can be interpreted as the fraction of observations in each cell. Zeros indicate that a particular group does not generate classes of that size.

Figure C-1: Nonparametric Regressions of Test Score on Class Size



The figure plots the semiparametric relationship between class size and test scores by group, as uncovered by the EAMP algorithm estimated on a sample that combines regular and regular + aide classrooms in the control group. All regressions use an Epanechnikov kernel and a bandwidth of 7.5. The regressions control for gender, race, and free lunch status.

## D Estimation Results with $K = 4$ Groups

Table D-1: Class Size Marginal Effects by Group

Group	$\mu_k$	$\Sigma_k$	$\sigma_{\epsilon,k}^2$	Schools	Students
1	-0.09 (0.073)	0.028 (0.162)	4.642 (3.116)	27	1179
2	0.132 (0.085)	0.001 (0.096)	11.249 (2.585)	14	1132
3	-0.349 (0.07)	0.248 (0.121)	6.425 (1.844)	20	693
4	0.032 (0.063)	0.001 (0.178)	10.996 (3.45)	18	809
Avg. Effect	-0.094 (0.032)	0.105 (0.051)	8.276 (0.676)		
	3813				
N	3813				

The table shows the estimated parameters governing the effect of class size on test scores for the model with  $K = 4$  groups. Regressions include controls for gender, race, and free lunch status. Controls are constrained to be equal across groups. The class size effects represent the effect of a one-unit increase in class size on test score performance. Bootstrapped standard errors from 95 bootstrapped data sets are in parentheses. Bootstrapping involves first sampling schools from the set of 79 schools that participated in the experiment and then sampling individuals within each school.