

# One Curriculum to Rule Them All

Experimental Evidence from Utah's Jordan School District

Sam Lee

Gavin Hatch

Brigg Tredler\*

December 2025

## Abstract

Elementary curriculum choices in early grades can shape later academic outcomes. However, in our evaluation, causal evidence on alternative instructional materials remains understudied. We design a randomized crossover experiment in two second-grade classrooms in the Jordan School District in northern Utah. 33 students received either a district-mandated reading curriculum or a teacher-curated curriculum in spelling and phonics, with pre- and post-tests in each of two weeks. We estimate a bivariate fixed effects model that uses within-student changes to identify the causal effect of curated instruction on both outcomes. We then study the finite sample behavior of Wald and  $F$ -tests for the joint hypothesis of no effect using a simulation design that matches the structure of the experiment. To allow for student-specific heterogeneity in baseline achievement, we also fit a bivariate Bayesian random effects model. Frequentist joint tests do not reject the null at conventional levels in this small sample. However, posterior draws from the random effects model indicate a substantial positive treatment effect. The posterior probability that curated instruction improves at least one of the two scores is estimated at 0.98. Because test scores are censored at 100%, these estimates likely understate the magnitude of the true gains. Notwithstanding small sample limitations, we suggest that a simple curated curriculum can generate meaningful short-run improvements in reading skills for students in this setting.

## 1 Introduction

Elementary curriculum is an important driver of early academic success. When instructional materials fail to engage students, they may lose interest in school at a young age and have difficulty recovering. In practice, districts often adopt a single mandated curriculum, while individual teachers develop their own materials in an effort to keep students motivated. Despite the prevalence of these informal alternatives, there is little direct causal evidence that compares district-mandated curricula with teacher-curated instruction in real classrooms.

We address this gap with a randomized crossover experiment in two second-grade classrooms in the Jordan School District in northern Utah. A total of 33 students followed either a district-mandated reading curriculum or a teacher-curated curriculum in spelling and phonics over two weeks, with pre- and post-tests in each week. The experiment is designed with crossover assignment at the classroom week level and with repeated measures for each student. This yields within-student variation that supports a clear causal interpretation for short-run gains in test scores. Under standard assumptions, the design identifies the effect of curated instruction on learning outcomes for this class and time period.

The experiment is multivariate by construction, since each observation consists of a pair of outcomes for spelling and phonics scores. We use this bivariate structure to obtain tractable results in a multivariate setting. Working with two related outcomes allows us to study joint treatment effects and multivariate test statistics without the complexity of higher-dimensional multivariate methods. The bivariate model provides a natural setting to compare Wald and  $F$ -tests for a joint hypothesis and to examine how these tests behave in small samples that reflect our realistic classroom experiment.

The remainder of the paper is organized as follows. Section 2 describes the experimental design, the timing of the crossover assignment, and the resulting dataset. Section 3 introduces our methods. We first develop a bivariate fixed effects model that uses within-student changes to identify the causal effect of curated instruction, then study the finite sample behavior of Wald and  $F$ -tests through a simulation calibrated to our

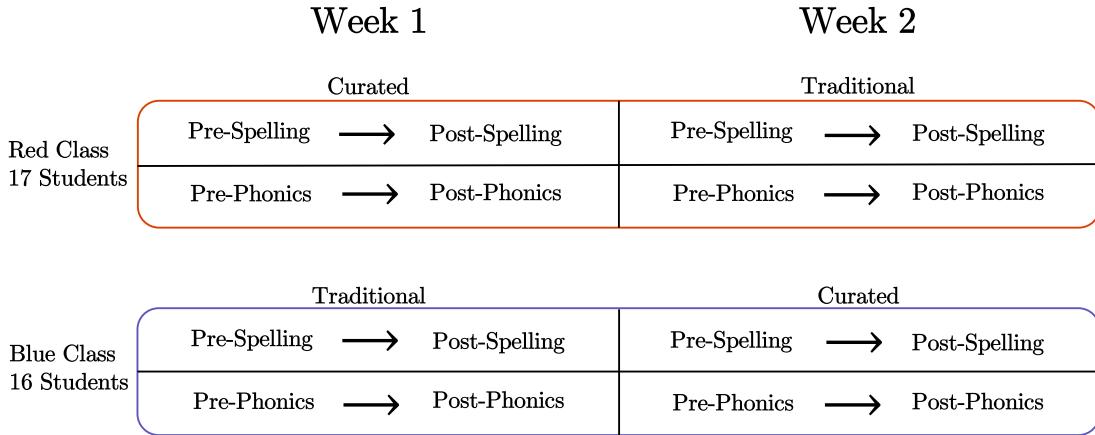
---

\*Department of Statistics, Brigham Young University. We acknowledge the generous support that made this research possible. We thank Ella Hatch, who secured permission for data collection. We thank Michael Christensen for educating us in multivariate methods.

design, and finally extend the framework to a Bayesian random effects model that allows for heterogeneity between students. Section 4 reports the empirical results from both the fixed effects and random effects analyses, and interprets the joint posterior distribution of the treatment effect vector. Section 4.1 discusses how the bounded nature of test scores affects the interpretation of estimated gains. Section 5 concludes, outlines limitations of the study, and describes threats to identification and directions for future work.

## 2 Data and Experimental Design

The data were collected on 33 students across two second-grade classes (Blue and Red) taught by Mrs. Ella Hatch over the course of two weeks at Fox Hollow Elementary (See Figures 7 and 8 in the appendix for raw grades). She used two different curricula: traditional and teacher-curated. The traditional curriculum was written by the private company, Really Great Reading, and chosen by the district. The curated curriculum covered the same content but was taught by Mrs. Hatch in a way that she thought would be more engaging for the students. The response variable is the recorded test scores for two subjects: spelling and phonics. Each class was randomly assigned an instruction type in the first week and assigned the opposite in the following week. Each week, a pre-test and post-test for each subject were administered on Monday and Friday, respectively. Figure 1 shows this experiment design. This design results in a total of eight scores per student (two tests for two subjects over two weeks). In addition to the designed experiment, we also include a covariate for student age. Since we implement our experiment in two classes of second-graders, each student will be either seven or eight years old. While less than one year seems like an insignificant difference, numerous studies have shown that small differences in age in early childhood development can have a significant effect on learning ability [Sharp et al., 2009, Bedard and Dhuey, 2006, Navarro et al., 2016].



*Notes:* This figure summarizes the randomized crossover design for the two second-grade classrooms. Each row represents a classroom (Blue or Red), and each column represents a week of instruction. Boxes indicate the assigned curriculum in that classroom week (traditional or teacher-curated). Within each cell, the pre-test is given on Monday, and the post-test is given on Friday for both spelling and phonics, so each student contributes four scores per week and eight scores in total across the two-week study.

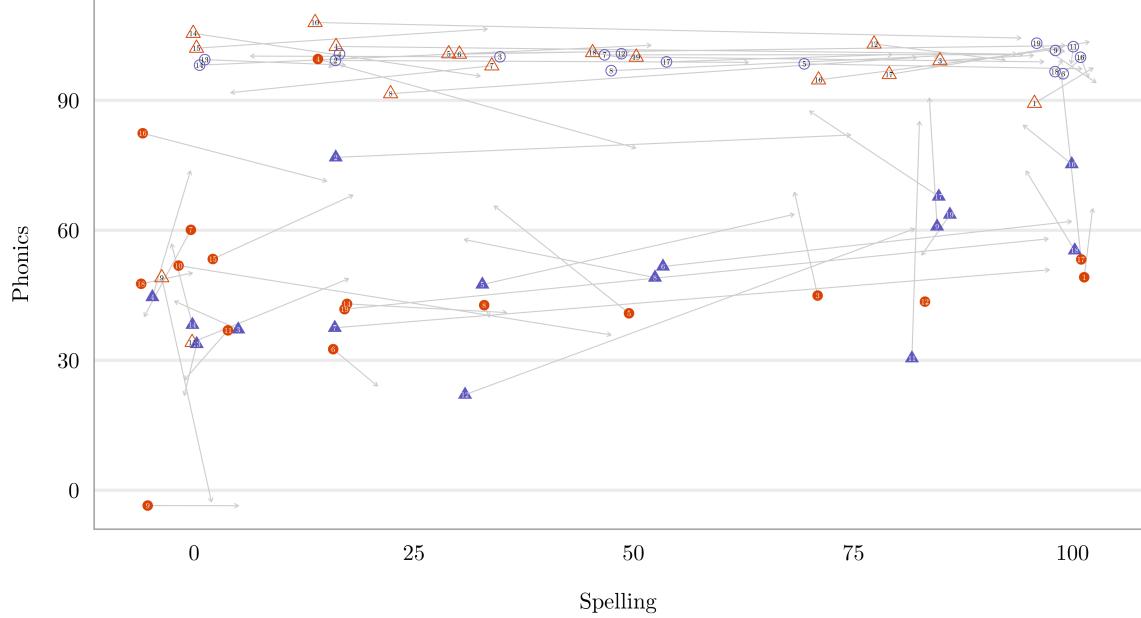
Figure 1: Experimental design

Exploratory analysis in Figure 2 shows each student’s progression in score from pre-tests to post-tests broken up by class, week, and instruction type. Most arrows show improvement by pointing to the right and upwards, as we expect by virtue of any amount of instruction. In the subsequent section, we estimate which curriculum has a greater influence on score improvement. We also aggregate<sup>1</sup> these trends in Figure 6 (see appendix) for a more simplified view of the data.

Some students were absent on test days and do not have a recorded score. There were a total of 5 missing values. We reached out to the teacher, Mrs. Hatch, and asked if she had any thoughts on how these students

<sup>1</sup>Note that in our analysis, the aggregate trends may ignore within-student improvement. Hence, we account for this heterogeneity in Section 3.

would have scored had they been present. She was confident in her assessment of the missing scores due to her prior experience with those particular students, so we decided to impute<sup>2</sup> the grades based on “teacher’s intuition.”



*Notes:* This figure displays the observed bivariate outcomes from the experiment. Each point represents a student’s spelling-phonics pre-test score, with arrows depicting within-student changes from pre-test to post-test in each week. For the purpose of presentation, we jitter each pre-test score according to a Gaussian distribution with variance 9. Triangles denote curated-instruction weeks and circles denote traditional-instruction weeks; blue markers correspond to the blue classroom and red markers to the red classroom.

Figure 2: Observed Spelling and Phonics Scores

### 3 Methods

Our empirical strategy proceeds in three steps. We first develop a bivariate fixed-effects model that leverages the crossover structure of the experiment to identify the short-run causal effect of curated instruction on spelling and phonics performance. We then assess the small-sample behavior of the resulting multivariate hypothesis tests through a simulation study calibrated to the structure of our data. Finally, to allow for student-level heterogeneity in baseline proficiency, we extend the framework to a bivariate random-effects specification and estimate it using a Bayesian multilevel model.

#### 3.1 Fixed Effects Model

We begin by outlining a simple bivariate fixed-effects regression model to evaluate the treatment effect on a student’s spelling and phonics scores. We model the student’s outcome vector jointly as,

$$\mathbf{y}_{iwt} = \begin{pmatrix} \text{spelling}_{iwt} \\ \text{phonics}_{iwt} \end{pmatrix} \in \mathbb{R}^2, \quad (1)$$

for each student  $i = 1, \dots, N$ , week  $w \in \{1, 2\}$ , and time  $t \in \{\text{pre}, \text{post}\}$  indicating the time periods before and after (respectively) a student takes the quiz during a given week. Recall in Section 2 that treatment

---

<sup>2</sup>We are confident that the imputation of the five missing values over 264 observations, or 132 bivariate observations, has little effect on our resulting analysis.

occurs at the classroom level. Hence, we first propose modeling the outcome vector through the following specification.

$$\begin{aligned} \mathbf{y}_{iwt} &= \boldsymbol{\alpha}_i + \boldsymbol{\lambda}_w + \boldsymbol{\kappa}_c + \mathbf{x}'_{iwt} \boldsymbol{\gamma} \\ &+ \mathbb{1}\{t = \text{post}\} \boldsymbol{\beta}_1 + \mathbb{1}\{\text{instruction}_{cw} = \text{curated}\} \boldsymbol{\beta}_2 \\ &+ \mathbb{1}\{t = \text{post}\} \mathbb{1}\{\text{instruction}_{cw} = \text{curated}\} \boldsymbol{\beta}_3 \\ &+ \mathbf{u}_{iwt}, \end{aligned} \quad (2)$$

$$\text{where } \mathbf{u}_{iwt} \sim \mathcal{N}_2 \left( \mathbf{0}, \boldsymbol{\Sigma}_u = \begin{bmatrix} \sigma_s^2 & \rho \\ \rho & \sigma_p^2 \end{bmatrix} \right) \quad (3)$$

where we include bivariate fixed effects<sup>3</sup> for each student ( $\boldsymbol{\alpha}_i$ ), week ( $\boldsymbol{\lambda}_w$ ), and each class (either blue or red) ( $\boldsymbol{\kappa}_c$ ). The vector  $\mathbf{x}_{iwt}$  represents additional controls such as the age of the student. We are particularly interested in  $\boldsymbol{\beta}_3$ , which represents how the student performs at the end of the week on their quiz after receiving the curriculum treatment. Stacking score vectors into the corresponding  $n \times 2$  matrix  $\mathbf{Y}$ , using  $\mathbf{X}$  to represent  $n \times k$  the matrix of fixed effect, time, and treatment indicators, and using the  $k \times 2$  matrix  $\mathbf{B}$  to represent the matrix of regression coefficients to be estimated ( $\boldsymbol{\alpha}, \boldsymbol{\lambda}, \boldsymbol{\kappa}, \boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \boldsymbol{\beta}_3$ ), the standard maximum likelihood estimators hold. Namely<sup>4</sup>,

$$\hat{\mathbf{B}} := (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}, \quad (4)$$

$$\hat{\boldsymbol{\Sigma}}_u := \frac{1}{n} (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}})' (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}}). \quad (5)$$

Thus, it follows<sup>5</sup> that,

$$\text{Var}(\text{vec}(\hat{\mathbf{B}})) := \boldsymbol{\Sigma}_u \otimes (\mathbf{X}' \mathbf{X})^{-1}. \quad (6)$$

We are particularly interested in whether at least one score (either that of spelling or phonics) is improved using the curated curriculum. Formally, we state this hypothesis as,

$$\begin{aligned} \mathcal{H}_0 : \mathbf{B}' \mathbf{c} &= 0, & \mathbf{c} &= (0 \ 0 \ \cdots \ 1)' \\ \mathcal{H}_1 : \mathbf{B}' \mathbf{c} &\neq 0, \end{aligned} \quad (7)$$

where the last entry in  $\mathbf{c}$  indicates the row position of  $\boldsymbol{\beta}_3$  in  $\mathbf{B}$ . Using a consistent estimator<sup>6</sup> for Eq. (6), we compute the corresponding Wald statistic as,

$$W := (\hat{\mathbf{B}}' \mathbf{c})' [(\mathbf{X}' \mathbf{X})_{kk}^{-1} \hat{\boldsymbol{\Sigma}}_u]^{-1} (\hat{\mathbf{B}}' \mathbf{c}) \xrightarrow{d} \chi_2^2. \quad (8)$$

It is not difficult to see that because  $\mathbf{B}' \mathbf{c} \sim \mathcal{N}_2(\mathbf{0}, \mathbf{V})$  by the Normality assumption in Eq. (3), where  $\mathbf{V} = (\mathbf{X}' \mathbf{X})_{kk}^{-1} \hat{\boldsymbol{\Sigma}}_u$ ,  $W$  can be decomposed into two statistics that converge to bivariate standard Normal distributions by constructing  $\mathbf{Z}$  such that  $\mathbf{Z} = \mathbf{V}^{-1/2} (\mathbf{B}' \mathbf{c})$ . Thus,  $E[\mathbf{Z}] = \mathbf{V}^{-1/2} E[\mathbf{B}' \mathbf{c}] = \mathbf{0}$ , and  $\text{Var}(\mathbf{Z}) = \mathbf{V}^{-1/2} \mathbf{V} \mathbf{V}^{-1/2} = \mathbf{I}_2$ . Hence,

$$W = (\hat{\mathbf{B}}' \mathbf{c})' \hat{\mathbf{V}}^{-1/2} \hat{\mathbf{V}}^{-1/2} (\hat{\mathbf{B}}' \mathbf{c}) = (\hat{\mathbf{V}}^{-1/2} \hat{\mathbf{B}}' \mathbf{c})' (\hat{\mathbf{V}}^{-1/2} \hat{\mathbf{B}}' \mathbf{c}) = \hat{\mathbf{Z}}' \hat{\mathbf{Z}}.$$

We show this, frankly, obvious result to motivate our use of Hotelling's correction, which we outline as follows. While  $\hat{\mathbf{Z}} \xrightarrow{d} \mathbf{Z} \sim \mathcal{N}_2(\mathbf{0}, \mathbf{I}_2)$ , it may be the case that  $\hat{\mathbf{V}}$  is not the best estimator for  $\mathbf{V}$  to the extent that  $\hat{\mathbf{Z}}$  differs significantly from the Normal distribution in finite samples. To this end, we remark that  $\hat{\mathbf{V}}$

<sup>3</sup>In this case, since we only have two classrooms in our experiment, the fixed effects given by  $\boldsymbol{\lambda}_w$  render  $\mathbf{X}' \mathbf{X}$  rank deficient and are thus excluded in empirical estimation. In other words, we can always perfectly identify all the parameters in the model if we are given a student and his or her corresponding classroom. We include these fixed effects in Eq. (3) simply for tractability.

<sup>4</sup>In practice, we use the unbiased estimator for  $\boldsymbol{\Sigma}_u$ , that is,  $\frac{1}{n-k} (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}})' (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}})$ .

<sup>5</sup>This is a standard result. However, for the interested reader, we reproduce it in Section B.

<sup>6</sup>Specifically, we use  $\hat{V}(\text{vec}(\hat{\mathbf{B}})) = \hat{\boldsymbol{\Sigma}}_u \otimes (\mathbf{X}' \mathbf{X})^{-1}$ , where  $\hat{\boldsymbol{\Sigma}}_u$  is given in Eq. (5).

follows a Wishart distribution [Wishart, 1928]. Therefore, the statistic  $W = \hat{\mathbf{Z}}'\hat{\mathbf{Z}}$  follows a Hotelling  $T^2$  distribution. Using the result from Hotelling [1931], it follows that

$$F := \frac{n-p-1}{2(n-p)} W \sim F_{2,n-p-1}. \quad (9)$$

We evaluate the performance of both the Wald and  $F$ -statistics for testing Eq. (7) under our experimental design in the subsequent section. We do not consider eigenvalue-based MANOVA tests, since in our setting they collapse to the same underlying quadratic form; see Section A.1.

The selling feature of this specification is its causal implications within the experimental design. The inclusion of student and week fixed effects absorbs all time-invariant heterogeneity across students and all common shocks within a given week of instruction and thus enables us to recover an unbiased treatment effect. The bivariate student effects control for persistent differences in ability, motivation, or home environment that jointly influence spelling and phonics scores; likewise, the week effects  $\lambda_w$  remove any systematic differences in testing conditions or instructional content that are shared by all students in a given week. Conditional on these fixed effects, the remaining within-student variation across the pre- and post-tests isolates the short-run gain (linearly) attributable to instruction during that week<sup>7</sup>. Under this design, given usual model assumptions,  $\beta_3$  is identified from a clean difference-in-differences comparison of within-student gains. The estimator of  $\beta_3$  unbiasedly recovers the causal effect of the curated curriculum on student performance.

### 3.1.1 Simulation Study

The joint hypothesis in Eq. (7) concerns a bivariate treatment effect. The curated curriculum may increase spelling performance, phonics performance, or both. Because our experiment is small by design, large-sample approximations supporting multivariate Wald tests may perform poorly in finite samples. Moreover, the power of multivariate hypothesis tests in short panels can depend sensitively on the covariance structure of the data-generating process. To justify our inferential strategy and to evaluate the operating characteristics of our test procedure in a setting commensurate with our empirical application, we conduct a simulation study calibrated to our fixed-effects model described in the previous section.

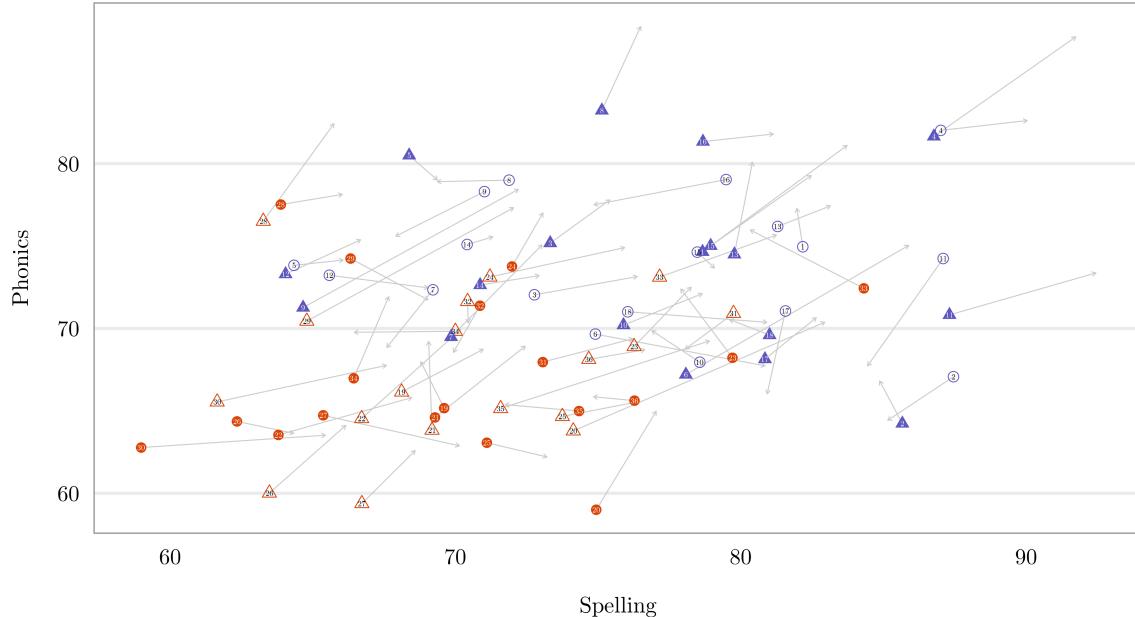
For each simulated dataset, we preserve the exact crossover structure of our experiment: two classrooms, two weeks, pre- and post-measurements in each week, and a single bivariate observation for each  $(i, w, t)$ . We generate synthetic data under the multivariate fixed-effects model in Eq. (3). Specifically, each student receives independent bivariate intercepts, and scores evolve within each week according to idiosyncratic bivariate error terms with residual covariance  $\rho$ . We vary the true treatment effect vector over a  $11 \times 11$  grid of values spanning both positive and negative effect sizes, yielding 121 distinct simulation configurations. We assess both the size of the test under the null hypothesis and the power using this design. For each configuration of  $\beta_3$ , we generate 2,000 Monte Carlo replications for each configuration. Figure 3 displays one such realization.

Figure 4 reports the empirical rejection probabilities of both the Wald test and the  $F$ -test across the  $11 \times 11$  grid of treatment effect values. We also report the Monte Carlo errors in Table 2 in the appendix. We discuss several important implications for the purposes of our experimental study. First, along the null line  $\beta_3 = (0, 0)'$ , rejection rates for the Wald test concentrate around 0.06. Thus, while the Wald test may slightly over-reject in a small experimental setting such as ours, we remark that it is generally appropriately sized. The  $F$ -test, in contrast, improves on this performance and is precisely the nominal 5 percent level under the null (subject to Monte Carlo error).

Second, the power of the procedure increases smoothly as both spelling and phonics effects move in a single direction. We suggest that this is a desirable property for our experiment. Since we hypothesize that the curated curriculum will move both spelling and phonics scores above the counterfactual expectation, if

---

<sup>7</sup>We give a brief outline of this intuition here. Because the curriculum is assigned at the classroom-week level through a deterministic crossover design, the indicator  $\mathbb{1}(\text{instruction}_{cw} = \text{curated})$  varies for each student across weeks but is orthogonal to their latent ability and all other time-invariant characteristics by virtue of the inclusion of fixed effects. This follows from the experimental structure: each classroom receives both instructional regimes, but in different weeks, thereby generating quasi-experimental within-student treatment variation that is unconfounded by individual heterogeneity. Thus, the interaction term  $\mathbb{1}(t = \text{post})\mathbb{1}(\text{instruction}_{cw} = \text{curated})$  compares post-test improvements that come by way of the curated curriculum *relative* to the improvement that would've otherwise come by the traditional instruction, differencing out all baseline differences captured by  $\alpha_i$  and all week-specific shocks captured by  $\lambda_w$ .



*Notes:* This figure displays a single simulated realization of spelling and phonics scores. Each point represents an  $(i, w, t)$  observation for a student, where student IDs are denoted by the numbers inscribed in each point. Arrows from each node point to the end-of-week score vector. Treatment status is indicated by the marker shape: circles denote traditional instruction, while triangles denote the curated curriculum treatment. Red shapes indicate observations from the red class, and blue shapes indicate observations from the blue class. Open shapes indicate quiz scores in week 1, and closed shapes indicate quiz scores in week 2. We expect that some instruction, regardless of type, will improve scores at the end of the week. However, on average, we hypothesize that the idiosyncratic slopes (the measure of improvement) for each student will be larger in magnitude than the non-treated groups. Compare Figure 2.

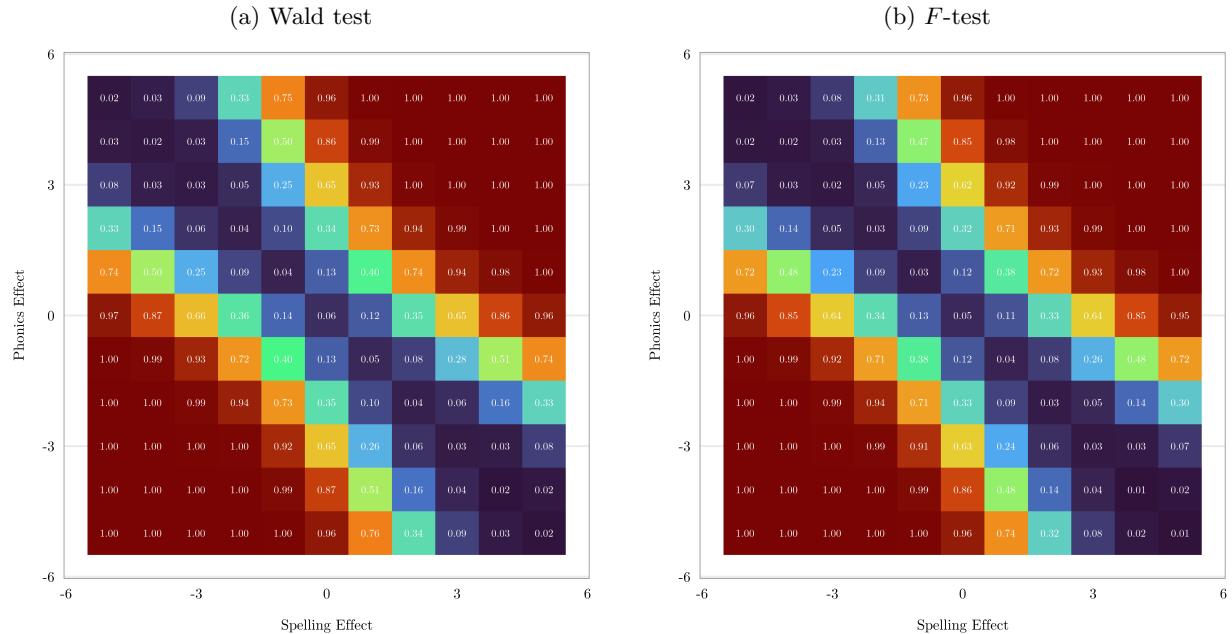
Figure 3: Simulated Bivariate Outcomes

we detect a positive treatment effect vector, we maintain that the effect we see is recovered with high power. The Wald and  $F$ -tests perform similarly in this region.

Third, when the two components of  $\beta_3$  have opposite signs, rejection probabilities fall dramatically, even when the magnitudes of the effects are large. We suspect this occurs because both the Wald statistic and the  $F$  statistic depend on the same quadratic form  $\hat{\beta}_3' \hat{V}^{-1} \hat{\beta}_3$ , which is sensitive to the direction of the true effect vector. Alternatives that move away from the origin but point into different quadrants relative to the null do not produce large test statistics. In this sense, opposite-signed effects “cancel” each other out in both procedures. We remark that this behavior is consistent with the substantive hypothesis of interest. The tests are designed to detect whether curated instruction produces simultaneous improvements in both spelling and phonics, so improvements in only one outcome, or improvements in one paired with deterioration in the other, do not provide evidence against the joint null.

### 3.2 Random Effects Model

A natural alternative to the fixed-effects specification in Eq. (3) is to treat student-specific heterogeneity as arising from a random sampling process. In multilevel settings, a random-effects formulation offers two advantages. First, by shrinking individual intercepts toward a common population distribution, the random-effects estimator can be substantially more efficient than the fixed-effects estimator when the student effect is uncorrelated with the regressors. We defend this assumption in our setting because the curriculum assignment varies only at the classroom–week level and thus is orthogonal to latent student ability by design. Second, modeling the dependence structure through a random-effects model allows us to quantify heterogeneity in baseline spelling and phonics proficiency across the population. To incorporate the random-effects design, we make a slight adjustment to Eq. (3) by replacing the student fixed effects  $\alpha_i$  with a bivariate random



*Notes:* Each panel reports Monte Carlo rejection probabilities for the joint null  $\mathcal{H}_0 : \beta_3 = (0, 0)'$  in our bivariate fixed-effects model. The horizontal and vertical axes index the true spelling and phonics components of  $\beta_3$  on an  $11 \times 11$  grid ranging from  $-5$  to  $5$ . Cell colors and overlaid numbers give the proportion of 2,000 simulated samples in which the null is rejected at the nominal  $\alpha = 5\%$  level. Panel (a) shows the performance of the  $\chi^2$  Wald test; panel (b) displays the corresponding  $F$ -test. Selected results are also presented in Table 2.

Figure 4: Empirical rejection rates for joint tests over a grid of treatment effects

intercept:

$$\begin{aligned} \mathbf{y}_{iwt} &= \boldsymbol{\mu}_0 + \boldsymbol{\lambda}_w + \boldsymbol{\kappa}_c + \mathbf{x}'_{iwt} \boldsymbol{\gamma} \\ &\quad + \mathbf{z}'_{iwt} \mathbf{b}_i \\ &\quad + \mathbb{1}\{t = \text{post}\} \boldsymbol{\beta}_1 + \mathbb{1}\{\text{instruction}_{cw} = \text{curated}\} \boldsymbol{\beta}_2 \\ &\quad + \mathbb{1}\{t = \text{post}\} \mathbb{1}\{\text{instruction}_{cw} = \text{curated}\} \boldsymbol{\beta}_3 \\ &\quad + \mathbf{u}_{iwt}, \end{aligned} \tag{10}$$

$$\text{where } \mathbf{u}_{iwt} \sim \mathcal{N}_2 \left( \mathbf{0}, \boldsymbol{\Sigma}_u = \begin{bmatrix} \sigma_s^2 & \rho \\ \rho & \sigma_p^2 \end{bmatrix} \right), \tag{11}$$

and we include an overall intercept  $\boldsymbol{\mu}_0$  so that the student random effects  $\mathbf{b}_i$  may be parameterized with mean zero. We set

$$\begin{aligned} \mathbf{z}_{iwt} &= \mathbf{I}_2, \\ \mathbf{b}_i &\sim \mathcal{N}_2 \left( \mathbf{0}, \boldsymbol{\Omega} = \begin{bmatrix} \eta_s^2 & \eta_{sp} \\ \eta_{sp} & \eta_p^2 \end{bmatrix} \right), \end{aligned}$$

so that each student has a pair of latent baseline abilities (one for spelling and one for phonics). A priori, these random intercepts may be correlated through the off-diagonal parameter  $\eta_{sp}$ , but in estimation we regularize this correlation toward zero via a shrinkage prior. To that end, we opt to switch to a Bayesian paradigm for ease of estimation. While it is possible to estimate our bivariate mixed model with a somewhat complex Newton-Raphson optimization routine, we note that computational and inferential limitations under such an approach motivate a Bayesian treatment of the random-effects specification. We include a discussion of our decision to switch to a Bayesian framework in the appendix.

We give a descriptive outline of our Bayesian estimation strategy as follows. Conditional on the random

effects, the likelihood factorizes from Eq. (11) as,

$$\mathbf{y}_{iwt} \mid \mathbf{b}_i, \mathbf{B}, \boldsymbol{\Sigma}_u \sim \mathcal{N}_2(\boldsymbol{\mu}_0 + \boldsymbol{\lambda}_w + \boldsymbol{\kappa}_c + \mathbf{x}'_{iwt}\boldsymbol{\gamma} + \mathbf{z}'_{iwt}\mathbf{b}_i, \boldsymbol{\Sigma}_u),$$

with  $\mathbf{b}_i \sim \mathcal{N}_2(\mathbf{0}, \boldsymbol{\Omega})$  as above. For the fixed effects we place weakly informative Gaussian priors of the form

$$\beta_{jk} \sim \mathcal{N}(0, s_{jk}^2),$$

where the prior scales  $s_{jk}$  are chosen automatically<sup>8</sup> by our statistical software [Goodrich et al., 2025].

For the covariance matrices  $\boldsymbol{\Omega}$  and  $\boldsymbol{\Sigma}_u$  we adopt the decomposed covariance prior. Each covariance matrix can be written as  $\boldsymbol{\Sigma} = \mathbf{D}\mathbf{R}\mathbf{D}$ , where  $\mathbf{D}$  is diagonal with nonnegative standard deviations on the diagonal and  $\mathbf{R}$  is a  $2 \times 2$  correlation matrix. Independent weakly informative priors are placed on the standard deviations.<sup>9</sup> Posterior inference is based on draws from the joint posterior distribution

$$\pi(\mathbf{B}, \boldsymbol{\Omega}, \boldsymbol{\Sigma}_u, \{\mathbf{b}_i\}_{i=1}^N \mid \{\mathbf{y}_{iwt}\}),$$

obtained via Hamiltonian Monte Carlo as implemented in Stan. Again, our primary interest lies in the bivariate treatment-effect vector  $\boldsymbol{\beta}_3$ . Rather than relying on an asymptotic Wald test as we did for our fixed effect model, we summarize the joint posterior distribution of  $\boldsymbol{\beta}_3$  directly. The posterior probability that the curated curriculum improves performance in both subjects is then estimated by

$$\widehat{\Pr}(\beta_{3,\text{spelling}} > 0, \beta_{3,\text{phonics}} > 0 \mid \text{data}) = \frac{1}{M} \sum_{m=1}^M \mathbb{1}\left\{\beta_{3,\text{spelling}}^{(m)} > 0, \beta_{3,\text{phonics}}^{(m)} > 0\right\}, \quad (12)$$

for each posterior draw,  $m$ . We are also interested in the following posterior probability.

$$1 - \widehat{\Pr}(\beta_{3,\text{spelling}} < 0, \beta_{3,\text{phonics}} < 0 \mid \text{data}), \quad (13)$$

which is the posterior probability that it is *not* the case that curated instruction harms both spelling and phonics simultaneously. Or in other words, this represents the posterior probability that at least one test score is improved relative to the traditional curriculum, even at the expense of the improvement of the other score metric. These posterior probabilities provide a Bayesian equivalent of a joint test of  $H_0 : \boldsymbol{\beta}_3 = \mathbf{0}$ , with two clear advantages. First, it does not rely on large-sample  $\chi^2$  approximations and is straightforward to interpret as the posterior mass assigned to the region of  $(\beta_{3,\text{spelling}}, \beta_{3,\text{phonics}})$  space corresponding to positive treatment effects in both outcomes.

## 4 Results

In this section, we present the empirical findings from our experiment. We begin by examining the estimated effects of our main parameter vector of interest, that is, the effect of the curated instruction relative to the effects of the traditional curriculum. The estimated coefficient vector ( $\boldsymbol{\beta}_3$ ) and other parameters of interest are reported in Table 1. The fixed effects estimates indicate that the point estimates of the curated curriculum interaction are positive for both spelling and phonics. Under our frequentist estimation, neither component is individually significant, and we fail to reject the joint null  $H_0 : \boldsymbol{\beta}_3 = \mathbf{0}$  and  $\alpha = 0.05$ . We suggest that this is consistent with the limited size of the experiment and the resulting variability in the estimated effects. Beyond the treatment effects, Table 1 also reveals interesting patterns in the covariates. Age is positively associated with both spelling and phonics scores in all specifications. This supports the evidence given in Section 2 that small age differences yield the strongest marginal effects in scholastic performance for younger students. Students in the red classroom score lower on average in both subjects, which indicates baseline differences across classes that the fixed-effects specification is intended to absorb. As

<sup>8</sup>Specifically,  $s_{jk}$  is proportional to the ratio of the marginal standard deviation of outcome  $j$  to that of predictor  $k$ , multiplied by a fixed constant. This yields a weakly informative prior that is invariant to linear rescaling of the predictors and outcomes.

<sup>9</sup>These are exponential-type priors that downweight unrealistically large values. Additionally, an LKJ(2) prior with concentration parameter is placed on the correlation matrix. The LKJ prior centers  $\mathbf{R}$  at the identity and mildly shrinks correlations toward zero while still allowing substantial posterior correlation when supported by the data. We specify the same prior family for both  $\boldsymbol{\Omega}$  and  $\boldsymbol{\Sigma}_u$ , with a regularization parameter tuned to avoid overfitting in this small sample.

a result, it may be reasonable to hypothesize that potential peer spillover effects may differ by classroom, since students surrounded by higher-performing peers may experience stronger motivational or strategic incentives to improve. We view this as a natural direction for future work for researchers who seek to replicate or extend a similar experimental design. We also point out that the main effect of curated instruction ( $\beta_2$ ) is near zero in both models. This suggests that differences across instruction types arise primarily through within-week improvements.

Table 1: Frequentist and Bayesian Estimates for Curriculum Effect

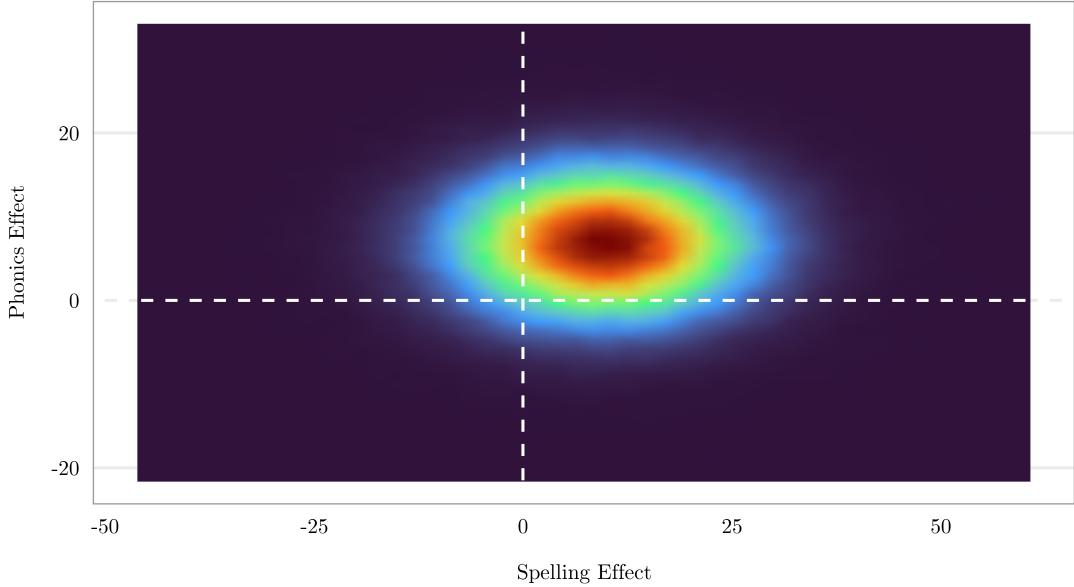
Panel A. Coefficient Summary					
	Frequentist (FE model)		Bayesian (RE model)		
	Estimate	SE	Mean	Median	95% HPD
<i>Spelling</i>					
Age	22.92	9.19	17.37	17.38	[2.34, 32.40]
Class Red	-19.03	6.11	-17.40	-17.42	[-28.79, -5.90]
$\beta_1$	13.09	7.83	13.12	13.14	[-1.81, 28.09]
$\beta_2$	-1.59	7.84	-1.56	-1.56	[-16.71, 13.22]
$\beta_3$	9.64	11.08	9.60	9.63	[-11.36, 30.83]
<i>Phonics</i>					
Age	16.54	4.72	13.69	13.69	[4.37, 22.65]
Class Red	-7.74	5.40	-8.00	-7.99	[-14.12, -1.88]
$\beta_1$	1.00	4.02	1.00	0.98	[-6.91, 8.74]
$\beta_2$	-3.41	4.02	-3.40	-3.40	[-11.32, 4.26]
$\beta_3$	6.82	5.69	6.81	6.82	[-4.10, 17.92]

Panel B. Inference Summary			
	Test	Statistic	p-value
<i>Frequentist</i>	Wald	$W = 1.44$	0.48
	$F$	$F = 0.72$	0.49
<i>Bayesian</i>	$\Pr(\beta_{3,\text{spell}} > 0, \beta_{3,\text{phon}} > 0   \text{data})$	0.72	
	$1 - \Pr(\beta_{3,\text{spell}} < 0, \beta_{3,\text{phon}} < 0   \text{data})$	0.98	

*Notes:* Parameters estimated with  $n = 132$  observations with 33 total students from two classes. Frequentist estimates come from the bivariate fixed-effects model estimation in Section 3.1. The Wald statistic in Panel B is computed using the quantiles of a  $\chi^2_2$  distribution, and the  $F$ -stat is computed using the quantiles of the  $F_{2,108}$  distribution. Bayesian estimates come from the random-effects model in Section 3.2 and reflect estimates on the posterior draws. HPD represents the 95% highest posterior density intervals. Joint inference evaluates the null hypothesis  $H_0 : \beta_3 = \mathbf{0}$  for the spelling and phonics treatment effects as given in Eq. (7). Posterior probabilities summarize the mass assigned to the positive-positive and non-negative quadrants of the joint posterior distribution.

However, we contend that our Bayesian random effects analysis provides a more informative description of the joint parameter space. Figure 5 displays the posterior density of  $(\beta_{3,\text{spelling}}, \beta_{3,\text{phonics}})$ , which places most of its mass in the region where both effects are positive. Posterior means, medians, and highest posterior density intervals in Table 1 reflect the same pattern. The two joint posterior probabilities computed from Eq. (13) and (12) summarize this evidence. The probability that both effects are positive is computed as 0.72. Hence, we suggest that the effect we estimate is practically significant even if it does not show statistical significance from a frequentist standpoint. The posterior probability that the curated curriculum instruction improves *at least one* score relative to the traditional curriculum is computed as 0.98 and grants our most substantive result of the analysis. This shows that the data place almost all of their weight on outcomes in which students benefit.



*Notes:* This figure displays the estimated posterior joint density of the bivariate treatment-effect vector from the Bayesian random-effects model. Warmer colors indicate regions of higher posterior probability. The dashed lines mark the axes corresponding to zero effects on spelling and phonics. The posterior mass is concentrated in the positive-positive quadrant, indicating strong posterior support that curated instruction improves both skills simultaneously.

Figure 5: Posterior Joint Density of the Bivariate Treatment Effect

#### 4.1 Attenuation Bias

The astute reader will notice that we have ignored the censored nature of our data. So far, we have not considered that by construction, quiz scores are bounded between 0 and 100; we ignore this constraint in both our modeling and simulation studies. Several students achieve perfect scores before and after instruction, as can be seen in Figure 2. Since post-treatment progress cannot be expressed beyond the score cap, measured gains are compressed relative to true gains. Hence, this kind of censoring induces attenuation bias in the estimated effects and shifts mass in the posterior toward more conservative values. Under this structure, the estimates reported here should be interpreted as lower bounds on the magnitude of the true treatment effect.

### 5 Conclusion

We study a randomized crossover experiment in two second-grade classrooms in the Jordan School District and evaluate the effect of a teacher-curated reading curriculum on spelling and phonics scores. A bivariate fixed effects model uses within-student changes to identify the short-run causal effect of curated instruction. In this framework, the multivariate Wald and  $F$ -tests are approximately well-sized in finite samples and exhibit strong power when effects move in the same direction. The quadratic form of these statistics limits hypothesis testing when effects diverge in opposite directions. A Bayesian random effects model complements this analysis by providing a richer summary of the joint treatment effect. Posterior draws place most of their mass on regions where curated instruction improves at least one outcome relative to the effects of the traditional curriculum. The posterior probability that at least one of the two scores improves relative to the traditional curriculum is estimated at 0.98, which we interpret as strong evidence that the curated curriculum benefits students in this setting.

## 5.1 Limitations

The empirical design and data impose several constraints on external validity. The sample is small, with 33 students drawn from a single school and a single instructor. The experiment also covers only two weeks, so we do not study the longer-run persistence of gains or the possibility of dynamic treatment effects. As discussed in earlier sections, these features limit the generalizability of our conclusions and suggest that our estimates should be viewed as local to this population, teacher, and time horizon. Moreover, we conclude that any effects that are present represent a lower bound on the true treatment effect due to the attenuation bias that occurs by the nature of the data.

Despite these limitations, we conclude that our study provides a concrete example of how elementary schools can use simple experimental designs and multivariate methods to evaluate instructional choices. The evidence indicates that a curated curriculum can generate meaningful improvements in reading skills over a short horizon and that Bayesian multilevel models can extract more information from small educational experiments than standard large sample tests. Future work could extend this design to more classrooms, grade levels, and schools with different demographic and socioeconomic profiles. Larger studies that randomize curricula across multiple teachers in a double-blind fashion, track students over longer periods, and measure peer interactions would help separate teacher effort from curricula effects that researchers intend to assess and clarify how alternative curricula shape the joint development of reading skills.

### 5.1.1 Threats to Identification

A central threat to identification arises from potential experimenter bias. Because the instructor in our experiment expressed a strong preference against the traditional curriculum, differential effort itself may account for at least part of the observed treatment effect. In this case, our estimates would capture a compound effect of instructional content and teacher motivation, and hence, we cannot disentangle these channels within the confines of this study. A more credible design would randomize curricula at the classroom level without disclosing assignments to teachers, thereby holding teacher beliefs and effort constant while isolating curricular effects. Nevertheless, our findings provide motivation for such a design. Even if the estimated gains reflect effort, the result is substantively meaningful. We conclude that variation in teacher engagement can generate measurable differences in short-run learning outcomes for their students.

## References

- K. Bedard and E. Dhuey. The persistence of early childhood maturity: International evidence of long-run age effects. *Quarterly Journal of Economics*, 121(4):1437–1472, 2006. doi: 10.1093/qje/121.4.1437. URL <https://academic.oup.com/qje/article-abstract/121/4/1437/1855234>.
- Ben Goodrich, Jonah Gabry, Imad Ali, and Sam Brilleman. *rstanarm: Bayesian Applied Regression Modeling via Stan*, 2025. URL <https://mc-stan.org/rstanarm/>. R package version 2.32.2.
- Harold Hotelling. The generalization of student’s ratio. *The Annals of Mathematical Statistics*, 2(3):360–378, 1931.
- R. M. Navarro, M. C. Aguilar, and D. Alonso. The relative age effect and its influence on academic performance in primary school. *PLoS One*, 10(10):e0141895, 2016. doi: 10.1371/journal.pone.0141895. URL [journals.plos.org](https://journals.plos.org).
- Caroline Sharp, Nalia George, Claire Sargent, Sharon O’Donnell, and Maureen Heron. The influence of relative age on learner attainment and development: A review of international evidence. Technical report, Department for Education, February 2009. URL [https://assets.publishing.service.gov.uk/media/5a82a07de5274a2e87dc2281/0209\CarolineSharp\\_et\\_al\\_RelativeAgeReviewRevised.pdf](https://assets.publishing.service.gov.uk/media/5a82a07de5274a2e87dc2281/0209\CarolineSharp_et_al_RelativeAgeReviewRevised.pdf). A UK government review confirming effects in primary school.
- John Wishart. The generalised product moment distribution in samples from a normal multivariate population. *Biometrika*, 20A(1/2):32–52, 1928.

## A Discussion

### A.1 Eigenvalue Tests

In principle, one could replace the Wald statistic with classical MANOVA tests that are functions of the eigenvalues of the matrix  $\mathbf{E}^{-1}\mathbf{H}$ , such as Wilks's  $\Lambda$ , Pillai's trace, the Hotelling–Lawley trace, or Roy's largest root. However, in our setting, the linear hypothesis  $H_0 : \boldsymbol{\beta}_3 = \mathbf{0}$  has rank one, so the associated hypothesis sum of squares and products matrix  $\mathbf{H}$  has rank at most one and  $\mathbf{E}^{-1}\mathbf{H}$  has at most a single nonzero eigenvalue, say  $\lambda_1$ . All four MANOVA test statistics become monotone functions of this same scalar quantity. In particular, the Hotelling–Lawley trace reduces to  $T = \lambda_1$ , and the Wald statistic for  $H_0 : \boldsymbol{\beta}_3 = \mathbf{0}$  is proportional to  $T$  and hence to the usual Hotelling  $T^2$  statistic. Eigenvalue-based tests, therefore, compress the sample information into the same quadratic form that underlies our Wald and  $F$  tests, so they cannot remedy the low power against mixed-sign alternatives documented in our simulations.

### A.2 A Case for a Bayesian Approach

To illustrate why estimating the multivariate random effects model in Eq. (11) may be more tractable in a Bayesian paradigm, we first consider a motivating result from statistical theory. Observe that, given our design, by conditioning on the week and classroom effects and the observed covariates, the marginal variance of an individual observation can be decomposed as,

$$\begin{aligned}\text{Var}(\mathbf{y}_{iwt} | \boldsymbol{\lambda}_w, \boldsymbol{\kappa}_c, \mathbf{x}_{iwt}, \mathbf{z}_{iwt}) &= \text{Var}(\mathbf{z}'_{iwt} \mathbf{b}_i) + \text{Var}(\mathbf{u}_{iwt}) \\ &= \boldsymbol{\Omega} + \boldsymbol{\Sigma}_u.\end{aligned}\tag{14}$$

Stacking the four  $(w, t)$  observations for student  $i$  into the vector  $\mathbf{y}_i \in \mathbb{R}^8$ , and letting  $\mathbf{J}$  and  $\mathbf{I}$  denote  $4 \times 4$  matrices of ones and identity matrices respectively, the implied covariance structure for the multivariate repeated-measures vector is

$$\text{Var}(\mathbf{y}_i) = (\mathbf{J} \otimes \boldsymbol{\Omega}) + (\mathbf{I} \otimes \boldsymbol{\Sigma}_u),\tag{15}$$

where  $\otimes$  denotes the Kronecker product. We continue to assume that students are independent across  $i$ , so the full covariance matrix of all stacked outcomes is block diagonal with blocks given by  $\text{Var}(\mathbf{y}_i)$ . In principle, under the Gaussian specification above, the generalized least squares estimator of the coefficient matrix  $\mathbf{B}$  that collects all fixed effects across both outcomes is

$$\hat{\mathbf{B}}(\boldsymbol{\theta}) = (\mathbf{X}' \mathbf{V}(\boldsymbol{\theta})^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}(\boldsymbol{\theta})^{-1} \mathbf{Y},$$

where  $\boldsymbol{\theta}$  stacks the variance parameters in  $\boldsymbol{\Sigma}_u$  and  $\boldsymbol{\Omega}$ . The corresponding asymptotic covariance of  $\text{vec}(\hat{\mathbf{B}})$  has the form

$$\text{Var}(\text{vec}(\hat{\mathbf{B}}(\boldsymbol{\theta}))) = \boldsymbol{\Sigma}_u \otimes (\mathbf{X}' \mathbf{V}(\boldsymbol{\theta})^{-1} \mathbf{X})^{-1},$$

so in theory one could obtain a Wald-type test for the contrast  $\mathbf{B}'\mathbf{c}$  by plugging in a REML estimate  $\hat{\boldsymbol{\theta}}$  of the variance parameters and a consistent estimate of  $\boldsymbol{\Sigma}_u$ . This follows as,

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) \equiv \arg \max_{\boldsymbol{\theta}} \left[ -\frac{1}{2} \log |\mathbf{V}(\boldsymbol{\theta})| - \frac{1}{2} \log |\mathbf{X}' \mathbf{V}(\boldsymbol{\theta})^{-1} \mathbf{X}| - \frac{1}{2} (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}}(\boldsymbol{\theta}))' \mathbf{V}(\boldsymbol{\theta})^{-1} (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}}(\boldsymbol{\theta})) \right].$$

In practice, however, this approach is cumbersome for two reasons. First, standard frequentist mixed-effects software typically fits separate univariate mixed models for each outcome and does not expose the joint multivariate covariance structure needed to evaluate the expressions above. Implementing a full multivariate REML estimator for  $\boldsymbol{\theta}$  and the corresponding Wald covariance would require jointly optimizing over a high-dimensional, non-convex parameter space. Moreover, the REML estimates of the fixed effects and the variance components are mutually dependent, so obtaining the Wald covariance entails solving a system with no closed-form expression. Second, for our empirical purposes, the available sample size is small (33 students with four repeated measurements per student), and the large-sample  $\chi^2_2$  approximation underlying a multivariate Wald test is unlikely to be reliable, as demonstrated in Section 3.1.1.

## B Additional Derivations

To obtain the sampling variance of  $\text{vec}(\hat{\mathbf{B}})$  as given in Eq. (6), we write the multivariate regression model as  $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{U}$  with  $\text{vec}(\mathbf{U}) \sim \mathcal{N}(\mathbf{0}, \Sigma_u \otimes \mathbf{I}_n)$ . The maximum likelihood estimator satisfies

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \mathbf{B} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U},$$

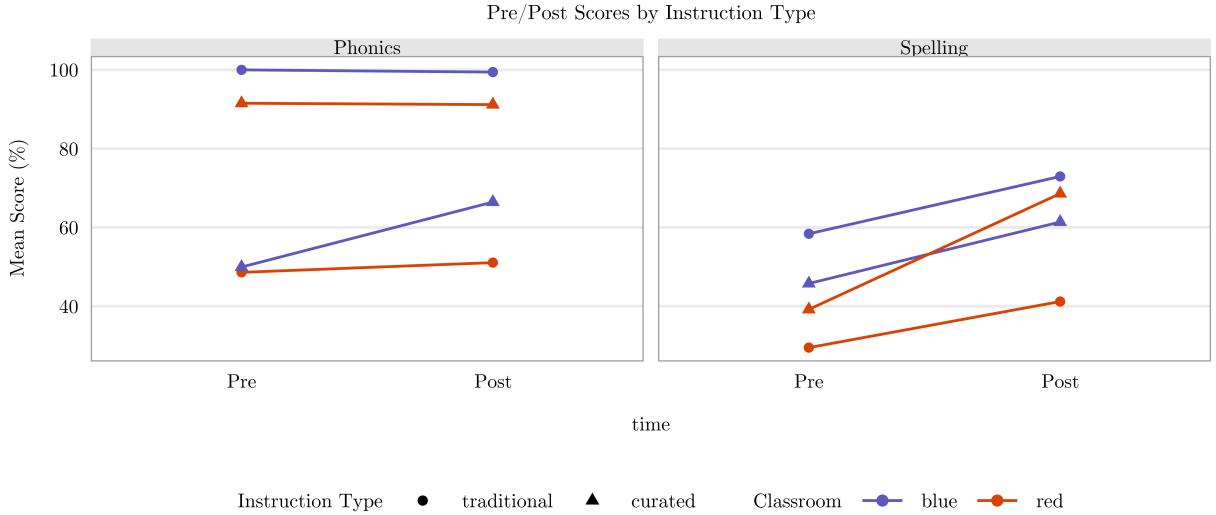
so vectorization yields

$$\text{vec}(\hat{\mathbf{B}}) = \text{vec}(\mathbf{B}) + (\mathbf{I}_2 \otimes (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \text{vec}(\mathbf{U}).$$

Hence, it follows directly,

$$\begin{aligned} \text{Var}(\text{vec}(\hat{\mathbf{B}})) &= (\mathbf{I}_2 \otimes (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \text{Var}(\text{vec}(\mathbf{U})) (\mathbf{I}_2 \otimes \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) \\ &= (\mathbf{I}_2 \otimes (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') (\Sigma_u \otimes \mathbf{I}_n) (\mathbf{I}_2 \otimes \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) \\ &\quad (\text{using } (\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})) \\ &= (\mathbf{I}_2\Sigma_u\mathbf{I}_2) \otimes ((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{I}_n\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) \\ &= \Sigma_u \otimes ((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) \\ &= \Sigma_u \otimes (\mathbf{X}'\mathbf{X})^{-1}. \end{aligned}$$

## C Additional Tables & Figures



*Notes:* This figure aggregates the student-level data from Figure 2 to show mean pre- and post-test scores for spelling and phonics by instruction type (traditional vs. curated) and classroom (blue vs. red). Circles represent traditional instruction and triangles represent the curated curriculum. Lines connect pre/post means within each classroom-instruction cell. The figure summarizes the same underlying observations as Figure 2 but collapses them to group means for a comparison across instructional regimes.

Figure 6: Mean pre/post scores by instruction type and classroom

	Spelling Unit 11 Pre-test	Phonics Unit 11 Pre-test	Spelling Unit 11	Unit 11 Phonics Pre-test	Unit 12 Spelling Pre-test	Unit 12 Phonics Pre-test	Unit 12 Spelling Test	Unit 12 Phonics Test	Unit 12 Phonics Test
1	100%.	91%.	100%.	100%.	100%.	90%.	100%.	67%.	67%.
2	100%.	100%.	100%.	100%.	100%.	100%.	100%.	100%.	100%.
3	83%.	100%.	100%.	100%.	67%.	42%.	67%.	67%.	67%.
4	17%.	100%.	100%.	100%.	17%.	100%.	50%.	83%.	83%.
5	33%.	100%.	100%.	100%.	50%.	42%.	33%.	67%.	67%.
6	33%.	100%.	50%.	100%.	17%.	33%.	17%.	25%.	25%.
7	33%.	91%.	0%.	100%.	0%.	67%.	0%.	42%.	42%.
8	17%.	91%.	83%.	100%.	33%.	42%.	33%.	33%.	33%.
9	0%.	50%.	0%.	0%.	0%.	0%.	0%.	0%.	0%.
10	17%.	100%.	100%.	100%.	0%.	50%.	50%.	33%.	33%.
11	0%.	33%.	17%.	50%.	0%.	33%.	0%.	25%.	25%.
12	83%.	100%.	100%.	100%.	83%.	42%.	83%.	42%.	42%.
13	5%	41%	7%	11%	1%	23%	1%	17%	17%
14	0%.	100%.	33%.	100%.	17%.	42%.	33%.	42%.	42%.
15	0%.	100%.	33%.	100%.	0%.	58%.	17%.	67%.	67%.
16	67%.	100%.	100%.	100%.	0%.	75%.	17%.	67%.	67%.
17	83%.	100%.	100%.	100%.	100%.	58%.			
18	50%.	100%.	50%.	100%.	0%.	50%.	0%.	50%.	50%.
19	50%.	100%.	100%.	100%.	17%.	42%.	100%.	58%.	58%.

Red

Figure 7: Gradebook for Red Class

	Unit 11 Spelling Pre test	11 phonics Pre test	Unit 11 Spelling test	Unit 11 phonics	Unit 12 Spelling Pre	Unit 12 phonics Pre	Unit 12 Spelling test	Unit 12 phonics	10 - 83 11 - 91 12 - 100	a c b a b c
1	87%	97%	97%	97%	97%	97%	97%	97%		
2	17%   100%   67%   100%   17%	88%   83%   83%								
3	33%   100%   0%   100%   0%   33%   0%   47%									
4	17%   100%   17%   100%   0%   50%   0%   75%									
5	67%   100%   100%   100%   33%   42%   67%   67%									
6	100%   100%	50%   50%   100%   50%								
7	50%   100%   100%   100%   17%   42%   100%   50%									
8	50%   100%   83%   100%   50%   50%   33%   50%									
9	100%   100%   100%   100%   83%   58%   83%   91%									
10	77%   17%   97%   100%   77%   88%   100%									
11	100%   100%   100%   100%   83%   33%   83%									
12	50%   100%   50%   100%   33%   25%   83%   58%									
13	0%   100%   17%   100%   0%   33%   0%   25%									
14	0%   100%   50%   100%   0%   50%   0%   50%									
15	18%   17%   87%   84%   85%   77%   100%									
16	100%   100%   100%   100%   100%   75%   100%   83%									
17	50%   100%   83%   100%   83%   67%   67%   91%									
18	100%   100%   100%   100%   100%   58%   100%   83%									
19	100%   100%   100%   91%   83%   58%   83%   58%									

Blue

Figure 8: Gradebook for Blue Class

Table 2: Selected empirical rejection rates ( $\alpha = 0.05$ ) for joint tests of  $\mathcal{H}_0 : \beta_3 = (0, 0)'$

True effects		Wald test		F test	
Spelling	Phonics	Rej. prob.	MCSE	Rej. prob.	MCSE
<b>Panel A: Positive-positive effects</b>					
0	0	0.057	0.005	0.050	0.005
1	1	0.396	0.011	0.375	0.011
2	2	0.936	0.005	0.930	0.006
3	3	0.999	0.001	0.999	0.001
<b>Panel B: Opposite-sign effects</b>					
4	-4	0.017	0.003	0.015	0.003
2	-2	0.035	0.004	0.029	0.004
-2	2	0.038	0.004	0.034	0.004
-4	4	0.024	0.003	0.022	0.003

*Notes:* Each entry reports the empirical rejection frequency across 2,000 Monte Carlo replications for tests of  $\mathcal{H}_0 : \beta_3 = (0, 0)'$  in the bivariate fixed-effects model of Section 3.1. The data-generating process uses two classrooms (blue and red) with  $n = 18$  students per class; baseline ability is generated with  $\alpha_i \sim \mathcal{N}(70, 5^2)$ , and subject-specific test shifts of 2 points in spelling and -2 points in phonics. Week effects are  $\lambda_1 = \lambda_2 = \mathbf{0}$ , class effects are  $\kappa_{\text{blue}} = (5, 5)'$  and  $\kappa_{\text{red}} = \mathbf{0}$ , and main effects satisfy  $\beta_1 = \beta_2 = \mathbf{0}$ . Idiosyncratic errors are bivariate normal with standard deviations  $\sigma_s = \sigma_p = 2$  and correlation  $\rho = 0.9$ . The true treatment-effect vector in each row is  $\beta_3 = (\beta_{3,\text{spelling}}, \beta_{3,\text{phonics}})'$ . The bold entry marks the rejection probability closest to the nominal size under the null.

## D R Code

All code is also documented and publicly available here: <https://github.com/SamLeeBYU/south-jordan-curriculum-experiment>.

### D.1 Data Cleaning

```
1 ## data_cleaning.R
2
3 library(tidyverse)
4
5 red <- read.csv("data/raw_data/red_scores.csv")
6 blue <- read.csv("data/raw_data/blue_scores.csv")
7
8 # Red class curated first
9 # Blue class curated second
10 # Red student 17 -> 100 both post assessments
11 # Blue
12 #   6 -> 100 for both assessments
13 #   10, 11 -> 83 for both
14
15 ### 1. Impute data based on Ella's feedback
16 red$Unit.12.Spelling.Post[17] = "100.00%"
17 red$Unit.12.Phonics.Post[17] = "100.00%"
18
19 blue$Unit.11.Spelling.Post[6] = 100
20 blue$Unit.11.Phonics.Post[6] = 100
21
22 blue$Unit.12.Phonics.Post[10] = 83
23 blue$Unit.12.Phonics.Post[11] = 83
24
25 ### 2. Drop rows where Student.Number is NA ----
26 red <- red %>% filter(!is.na(Student.Number))
27 blue <- blue %>% filter(!is.na(Student.Number))
28
29 ### 3. Add classroom labels ----
30 red <- red %>% mutate(classroom = "red")
31 blue <- blue %>% mutate(classroom = "blue")
32
33 ### 4. Make score columns numeric ----
34 pct_to_num <- function(x) {
35   x %>% str_replace("%", "") %>% as.numeric()
36 }
37
38 red <- red %>%
39   mutate(across(starts_with("Unit"), pct_to_num))
40
41 blue <- blue %>%
42   mutate(across(starts_with("Unit"), as.numeric()))
43
44
45 ### 5. Combine into one dataset ----
46 all <- bind_rows(red, blue)
47
48 ### 6. Pivot to long format ----
49 long <- all %>%
50   pivot_longer(
51     cols = matches("^Unit\\\\.\\d+"),
```

```

52     names_to = c("unit", "test", "time"),
53     names_pattern = "^(Unit\\\\.\\d+)\\.(Spelling|Phonics)\\.(Pre|Post)$",
54     values_to = "score"
55   )
56
57 ### 7. Recode variables ----
58 long <- long %>%
59   mutate(
60     id = Student.Number,
61     week = ifelse(unit == "Unit.11", 1L, 2L),
62     time = factor(time, levels = c("Pre", "Post")),
63     test = tolower(test), # "Spelling" -> "spelling", "Phonics" -> "phonics"
64     classroom = factor(classroom)
65   )
66
67 ### 8. Assign instruction type based on classroom times week ----
68 instruction_fun <- function(classroom, week) {
69   if (classroom == "blue" && week == 1) {
70     return("traditional")
71   }
72   if (classroom == "blue" && week == 2) {
73     return("curated")
74   }
75   if (classroom == "red" && week == 1) {
76     return("curated")
77   }
78   if (classroom == "red" && week == 2) return("traditional")
79 }
80
81 long <- long %>%
82   rowwise() %>%
83   mutate(instruction = instruction_fun(classroom, week)) %>%
84   ungroup() %>%
85   mutate(
86     instruction = factor(instruction, levels = c("traditional", "curated"))
87   )
88
89 ### 9. Final cleanup ----
90 final_data <- long %>%
91   dplyr::select(id, Age, classroom, week, time, instruction, test, score) %>%
92   arrange(id, week, test, time)
93
94 glimpse(final_data)
95 saveRDS(final_data, "data/final_project_scores.rds")
96
97 ## data.R
98 library(dplyr)
99
100 ellas.class <- readRDS("data/final_project_scores.rds") %>%
101   pivot_wider(names_from = "test", values_from = "score")

```

## D.2 DGP

```

1 ##dgp.R
2
3 library(dplyr)
4 library(tidyr)

```

```

5 #This DGP was originally written by Gavin Hatch,
6 # recoded by Sam Lee to match regression model in the paper
7
8 library(dplyr)
9 library(MASS)
10
11 generate_classroom_data <- function(
12   n_per_class = 18,
13   classes = c("blue", "red"),
14
15   # student ability
16   ability_mean = 70,
17   ability_sd = 6,
18
19   # test-specific shift
20   test_effect = c(spelling = 2, phonics = -2),
21
22   # week and class fixed effects
23   lambda_week = c('1' = 0, '2' = 1),
24   kappa_class = c(blue = 0, red = 0),
25
26   # regression parameters
27   beta1 = c(spelling = 4, phonics = 4), # effect of post indicator
28   beta2 = c(spelling = 0, phonics = 0), # main effect of curated
29   beta3 = c(spelling = 3, phonics = 3), # post x curated gain
30
31   # bivariate error covariance
32   sigma_s = 4, # spelling SD
33   sigma_p = 4, # phonics SD
34   rho = 0, # covariance term
35   clip_lo = 0,
36   clip_hi = 100,
37   seed = NULL
38 ) {
39   if (!is.null(seed)) {
40     set.seed(seed)
41   }
42
43   # error covariance matrix
44   Sigma_u <- matrix(
45     c(sigma_s^2, rho, rho, sigma_p^2),
46     nrow = 2,
47     byrow = TRUE
48   )
49   L.u <- chol(Sigma_u)
50
51   # cross-over assignment
52   instruction_fun <- function(classroom, week) {
53     if (week == 1 && classroom == "blue") {
54       return("traditional")
55     }
56     if (week == 1 && classroom == "red") {
57       return("curated")
58     }
59     if (week == 2 && classroom == "blue") {
56       return("curated")
57     }
58     if (week == 2 && classroom == "red") {
59

```

```

64     return("traditional")
65   }
66 }
67
68 # students and design
69 students <- data.frame(
70   id = 1:(length(classes) * n_per_class),
71   classroom = rep(classes, each = n_per_class)
72 )
73
74 design <- expand.grid(
75   id = students$id,
76   week = 1:2,
77   time = c("pre", "post"),
78   KEEP.OUT.ATTRS = FALSE,
79   stringsAsFactors = FALSE
80 ) %>%
81   left_join(students, by = "id") %>%
82   mutate(
83     instruction = mapply(instruction_fun, classroom, week),
84     week = factor(week),
85     time = factor(time, levels = c("pre", "post")),
86     instruction = factor(instruction, levels = c("traditional", "curated")),
87     classroom = factor(classroom)
88   )
89
90 # student abilities
91 ability_student_s <- rnorm(
92   nrow(students),
93   mean = ability_mean + test_effect["spelling"],
94   sd = ability_sd
95 )
96 ability_student_p <- rnorm(
97   nrow(students),
98   mean = ability_mean + test_effect["phonics"],
99   sd = ability_sd
100 )
101 names(ability_student_s) <- students$id
102 names(ability_student_p) <- students$id
103
104 clip01 <- function(x, lo = clip_lo, hi = clip_hi) pmin(pmax(x, lo), hi)
105
106 alpha.mat <- matrix(
107   c(ability_student_s[design$id], ability_student_p[design$id]),
108   ncol = 2
109 )
110 lambda <- matrix(
111   c(lambda_week[design$week], lambda_week[design$week]),
112   ncol = 2
113 )
114 kappa <- matrix(
115   c(kappa_class[design$classroom], kappa_class[design$classroom]),
116   ncol = 2
117 )
118 post = 1 * (design$time == "post")
119 cur = 1 * (design$instruction == "curated")
120
121 mu = alpha.mat +
122   lambda +

```

```

123  kappa +
124  post * beta1 +
125  cur * beta2 +
126  post * cur * beta3
127
128  Z = matrix(rnorm(nrow(design) * 2), ncol = 2)
129
130  Y <- Z %*% L.u + mu
131
132  long_data <- design %>%
133    mutate(
134      spelling = clip01(Y[, 1]),
135      phonics = clip01(Y[, 2])
136    )
137
138  long_data
139 }
```

### D.3 Model Code

```

1 #Fixed effects model
2 fit.lm <- function(dat) {
3   X <- model.matrix(
4     ~ 0 +
5       as.factor(id) +
6       as.factor(week) +
7       as.factor(time) * as.factor(instruction),
8     data = dat
9   )
10  m <- lm(
11    cbind(spelling, phonics) ~ 0 + X,
12    data = dat
13  )
14  Uhat <- residuals(m)
15  Sigma.u.hat <- crossprod(Uhat) / df.residual(m)
16  Var.B <- Sigma.u.hat * solve(crossprod(X)) [ncol(X), ncol(X)]
17  list(
18    est = coef(m) [ncol(X), ],
19    var = Var.B,
20    n = nrow(X),
21    p = ncol(X)
22  )
23}
24
25 wald.test <- function(estimates, alpha = 0.05) {
26   v.i <- solve(estimates$var)
27   w <- t(estimates$est) %*% v.i %*% estimates$est
28   (1 - pchisq(w, df = 2)) < alpha
29 }
30
31 f.test <- function(estimates, alpha = 0.05) {
32   v.i <- solve(estimates$var)
33   w <- t(estimates$est) %*% v.i %*% estimates$est
34
35   n = estimates$n
36   p = estimates$p
37   nu <- (n - p - 1)
```

```

38     f <- nu / (2 * (n - p)) * w
39
40     (1 - pf(f, df1 = 2, df2 = nu)) < alpha
41 }

```

## D.4 Bayesian Model Code

```

1 library(rstanarm)
2
3 source("dgp.R")
4
5 ellas.class <- generate_classroom_data(
6   rho = 0.9,
7   sigma_s = 2,
8   sigma_p = 2,
9   ability_sd = 5,
10
11 #This is effectively deviations in mu_0
12 test_effect = c(spelling = 2, phonics = -2),
13
14 lambda_week = c('1' = 0, '2' = 0),
15 kappa_class = c(blue = 5, red = 0),
16
17 beta1 = c(spelling = 0, phonics = 0), # effect of post indicator
18 beta2 = c(spelling = 0, phonics = 0), # main effect of curated
19 beta3 = c(spelling = 3, phonics = 3), # post x curated gain
20
21 seed = 666
22 )
23
24 source("data.R")
25
26 fit_mv <- stan_mvmer(
27   formula = list(
28     spelling ~ 0 +
29       Age +
30       factor(week) +
31       factor(classroom) +
32       factor(time) * factor(instruction) +
33       (1 | id),
34
35     phonics ~ 0 +
36       Age +
37       factor(week) +
38       factor(classroom) +
39       factor(time) * factor(instruction) +
40       (1 | id)
41   ),
42   data = ellas.class,
43   family = gaussian(),
44   prior = normal(0, 5, autoscale = TRUE),
45   prior_intercept = normal(0, 5, autoscale = TRUE),
46   prior_covariance = decov(regularization = 2),
47   chains = 4,
48   iter = 100000,
49   cores = 4,
50   refresh = T

```

```

51 }
52
53 allsamps <- as.matrix(fit_mv)
54 saveRDS(allsamps, "mcmc-samples.rds")

```

## D.5 Simulation Code

```

1 source("dgp.R")
2 source("models.R")
3
4 #Simulation size and power
5 plan <- expand.grid(
6   #different beta 3
7   spelling = -5:5,
8   phonics = -5:5
9 )
10 simulate <- function(plan, B = 1000, seed = NULL) {
11   if (!is.null(seed)) {
12     set.seed(seed)
13   }
14   plan$phat.w <- plan$phat.f <- 0
15   for (i in 1:nrow(plan)) {
16     rejections.w <- matrix(0, nrow = B, ncol = 1)
17     rejections.f <- matrix(0, nrow = B, ncol = 1)
18     plan.i <- plan[i, ]
19     for (b in 1:B) {
20       ellas.class <- generate_classroom_data(
21         rho = 0.9,
22         sigma_s = 2,
23         sigma_p = 2,
24         ability_sd = 5,
25
26         #This is effectively deviations in mu_0
27         test_effect = c(spelling = 2, phonics = -2),
28
29         lambda_week = c('1' = 0, '2' = 0),
30         kappa_class = c(blue = 5, red = 0),
31
32         beta1 = c(spelling = 0, phonics = 0), # effect of post indicator
33         beta2 = c(spelling = 0, phonics = 0), # main effect of curated
34         beta3 = c(spelling = plan.i$spelling, phonics = plan.i$phonics), # post x
35         curated gain
36
37         seed = NULL
38       )
39
40       fe <- fit.lm(ellas.class)
41       rejections.w[b, 1] <- wald.test(fe)
42       rejections.f[b, 1] <- f.test(fe)
43     }
44     plan$phat.w[i] = mean(rejections.w[, 1])
45     plan$phat.f[i] = mean(rejections.f[, 1])
46     print(round(100 * i / nrow(plan)) / 100)
47   }
48   plan$MCSE.w <- sqrt(plan$phat.w * (1 - plan$phat.w) / B)
49   plan$MCSE.f <- sqrt(plan$phat.f * (1 - plan$phat.f) / B)
50   plan

```

```

50 }
51
52 rejection.rates <- simulate(
53   plan,
54   B = 2000
55 )
56 saveRDS(rejection.rates, "rejection-rates-sim.RDS")

```

## D.6 Code for Analysis

```

1 #Frequentist analysis
2 source("data.R")
3 source("models.R")
4
5 X <- model.matrix(
6   ~ 0 +
7     as.factor(id) +
8     as.factor(classroom) +
9     Age +
10    as.factor(time) * as.factor(instruction),
11   data = ellas.class
12 )
13 m <- lm(
14   cbind(spelling, phonics) ~ 0 + X,
15   data = ellas.class
16 )
17 summary(m)
18
19 #Spelling | Coef | SE
20 #Age | 22.917 | 9.186
21 #Class Red | -19.025 | 6.108
22 #\beta_1 | 13.091 | 7.834
23 #\beta_2 | -1.586 | 7.835
24 #\beta_3 | 9.636 | 11.078
25
26 #Phonics | Coef | SE
27 #Age | 16.542 | 4.715
28 #Class Red | -7.738 | 5.403
29 #\beta_1 | 1.000 | 4.021
30 #\beta_2 | -3.410 | 4.022
31 #\beta_3 | 6.818 | 5.687
32
33 estimates <- fit.lm(ellas.class)
34 #F-test
35 v.i <- solve(estimates$var)
36 w <- t(estimates$est) %*% v.i %*% estimates$est
37
38 #Chisq test
39 w
40 # 1.44
41 1 - pchisq(w, 2)
42 #0.48
43
44 n = estimates$n
45 p = estimates$p
46 nu <- (n - p - 1)
47 f <- nu / (2 * (n - p)) * w

```

```

48 | f
49 | #0.7169982
50 | (1 - pf(f, df1 = 2, df2 = nu))
51 | #0.4905246
52 |
53 | #Bayesian analysis
54 | allsamps <- readRDS("mcmc-samples.rds")
55 | relevant.coefs <- allsamps[, c(1:14)]
56 |
57 | #Posterior Estimates
58 | cbind(
59 |   colMeans(relevant.coefs),
60 |   apply(relevant.coefs, 2, median),
61 |   apply(relevant.coefs, 2, \((x) HDInterval::hdi(x)) |> t()
62 | )
63 |
64 | #Spelling | Mean | Median | 95% HPD
65 | #Age | 17.37 | 17.38 | 2.34, 32.398
66 | #Class Red | -17.40 | -17.42 | -28.79, -5.898
67 | #\beta_1 | 13.124 | 13.139 | -1.8146, 28.0868
68 | #\beta_2 | -1.5608 | -1.5621 | -16.706, 13.2158
69 | #\beta_3 | 9.599 | 9.625 | -11.36, 30.83
70 |
71 | #Phonics | Mean | Median | 95% HPD
72 | #Age | 13.687 | 13.69 | 4.373, 22.65
73 | #Class Red | -8.00 | -7.99 | -14.11697, -1.88
74 | #\beta_1 | 0.996 | 0.982 | -6.911, 8.735574
75 | #\beta_2 | -3.39778 | -3.3988075 | -11.3178, 4.25631
76 | #\beta_3 | 6.810 | 6.8184 | -4.1028, 17.918
77 |
78 | #Bayesian hypothesis testing
79 | beta3 <- relevant.coefs[, c(7, 14)]
80 | #Pr(beta_3, spelling > 0, beta_3, phonics > 0)
81 | mean(beta3[, 1] > 0 & beta3[, 2] > 0)
82 | # = 0.72305
83 | #1-Pr(beta_3, spelling < 0, beta_3, phonics < 0)
84 | 1 - mean(beta3[, 1] < 0 & beta3[, 2] < 0)
85 | # = 0.9785

```

## D.7 Reproduce Figures

```

1 ##themes.R
2 library(ggplot2)
3 library(patchwork)
4 library(sysfonts)
5
6 font_add("cm", regular = "fonts/cmunrm.ttf")
7 showtext::showtext_auto()
8
9 theme_paper <- function(base_size = 14, text_size = 14, base_family = "cm") {
10   theme_minimal(base_size = base_size, base_family = base_family) %+replace%
11   theme(
12     # Text
13     plot.title = element_text(
14       face = "bold",
15       size = text_size + 2,
16       hjust = 0.5,

```

```

17     margin = margin(b = 6)
18   ),
19   plot.subtitle = element_text(
20     size = text_size,
21     hjust = 0.5,
22     margin = margin(b = 8)
23   ),
24   plot.caption = element_text(
25     size = text_size - 2,
26     hjust = 1,
27     margin = margin(t = 6)
28   ),
29   axis.title.x = element_text(size = text_size, margin = margin(t = 8)),
30   axis.title.y = element_text(
31     size = text_size,
32     angle = 90,
33     margin = margin(r = 8)
34   ),
35   axis.text = element_text(size = text_size - 1),
36
37   # Panel & grid
38   panel.grid.major.x = element_blank(),
39   panel.grid.minor.x = element_blank(),
40   panel.grid.minor.y = element_blank(),
41   #panel.grid.major.y = element_line(linewidth = 0.3, color = "grey85"),
42   panel.border = element_rect(fill = NA, color = "grey60", linewidth = 0.4),
43
44   # Legend
45   legend.position = "bottom",
46   legend.title = element_text(size = base_size - 1, face = "bold"),
47   legend.text = element_text(size = base_size - 2),
48   legend.key.width = unit(1.2, "lines"),
49
50   # Strips (for facets)
51   strip.background = element_rect(fill = "grey90", color = NA),
52   strip.text = element_text(face = "bold", size = base_size - 1),
53
54   # Margins
55   plot.margin = margin(t = 8, r = 8, b = 8, l = 8)
56 )
57 }
58
59 paper_palette <- c(
60   "#5e57be",
61   "#d94202",
62   "#1b9e77",
63   "#e7298a",
64   "#66a61e",
65   "#e6ab02"
66 )
67
68 scale_color_paper <- function(...) {
69   scale_color_manual(values = paper_palette, ...)
70 }
71
72 scale_fill_paper <- function(...) {
73   scale_fill_manual(values = paper_palette, ...)
74 }
75

```

```

76 scale_fill_gradient2_paper <- function(
77   low = "#7570b3",
78   mid = "white",
79   high = "#d95f02",
80   midpoint = 0,
81   ...
82 ) {
83   scale_fill_gradient2(
84     low = low,
85     mid = mid,
86     high = high,
87     midpoint = midpoint,
88     ...
89   )
90 }
91
92 ##figures.R
93 source("themes.R")
94 source("dgp.R")
95
96 ellas.class <- generate_classroom_data(
97   rho = 0.9,
98   sigma_s = 2,
99   sigma_p = 2,
100  ability_sd = 5,
101
102  #This is effectively deviations in mu_0
103  test_effect = c(spelling = 2, phonics = -2),
104
105  lambda_week = c('1' = 0, '2' = 0),
106  kappa_class = c(blue = 5, red = 0),
107
108  beta1 = c(spelling = 0, phonics = 0), # effect of post indicator
109  beta2 = c(spelling = 0, phonics = 0), # main effect of curated
110  beta3 = c(spelling = 3, phonics = 3), # post x curated gain
111
112  seed = 666
113 )
114 delta <- 0.04
115
116 arrows.df <- ellas.class %>%
117   mutate(time = tolower(time)) %>%
118   dplyr::select(id, week, instruction, time, spelling, phonics) %>%
119   tidyr::pivot_wider(
120     names_from = time,
121     values_from = c(spelling, phonics)
122   ) %>%
123   mutate(
124     x_end = spelling_pre + (1 - delta) * (spelling_post - spelling_pre),
125     y_end = phonics_pre + (1 - delta) * (phonics_post - phonics_pre)
126   )
127
128 base.dat <- ellas.class %>%
129   filter(time == "pre") %>%
130   mutate(
131     shape_class = ifelse(
132       instruction == "traditional",
133       ifelse(week == 1, "open_circle", "closed_circle"),
134       ifelse(week == 1, "open_triangle", "closed_triangle"))

```

```

135      )
136  )
137
138 ellas.class.dgp <- ggplot(
139   base.dat,
140   aes(
141     x = spelling,
142     y = phonics
143   )
144 ) +
145   # arrows
146   geom_segment(
147     data = arrows.df,
148     aes(
149       x = spelling_pre,
150       y = phonics_pre,
151       xend = x_end,
152       yend = y_end
153     ),
154     arrow = arrow(length = unit(0.05, "cm")),
155     color = "#CCC",
156     linewidth = 0.15,
157     inherit.aes = FALSE
158   ) +
159   geom_point(
160     aes(
161       color = classroom,
162       shape = shape_class
163     ),
164     stroke = 0.2,
165     show.legend = FALSE,
166   ) +
167   # labels: white for week 2
168   geom_text(
169     data = dplyr::filter(base.dat, week == 2),
170     aes(label = id),
171     color = "white",
172     size = 8,
173     vjust = 0.5,
174     hjust = 0.5,
175     show.legend = FALSE
176   ) +
177   # labels: black for week 1
178   geom_text(
179     data = dplyr::filter(base.dat, week == 1),
180     aes(label = id),
181     color = "black",
182     size = 8,
183     vjust = 0.5,
184     hjust = 0.5,
185     show.legend = FALSE
186   ) +
187   scale_color_paper() +
188   scale_shape_manual(
189     values = c(
190       open_circle = 1, # open circle
191       open_triangle = 2, # open triangle
192       closed_circle = 16, # filled circle
193       closed_triangle = 17 # filled triangle

```

```

194      )
195  ) +
196  labs(
197    x = "Spelling",
198    y = "Phonics",
199    shape = "Instruction",
200    color = "Class"
201  ) +
202  theme_paper(base_size = 10, text_size = 64) +
203  theme(
204    legend.title = element_text(size = 64),
205    legend.text = element_text(size = 64)
206  )
207
208 ggsave(
209   "Figures/dgp.png",
210   ellas.class.dgp,
211   width = 5,
212   height = 2.8,
213   dpi = 900
214 )
215
216 ##### Real Data #####
217
218 source("data.R")
219
220 jitter.sd <- 3
221 ellas.class <- ellas.class %>%
222   group_by(time == "Pre") %>%
223   mutate(
224     spelling = spelling + rnorm(n(), sd = jitter.sd),
225     phonics = phonics + rnorm(n(), sd = jitter.sd)
226   ) %>%
227   ungroup()
228
229 delta <- 0.04
230 arrows.df <- ellas.class %>%
231   mutate(
232     time = tolower(time)
233   ) %>%
234   dplyr::select(id, week, instruction, time, spelling, phonics) %>%
235   tidyr::pivot_wider(
236     names_from = time,
237     values_from = c(spelling, phonics)
238   ) %>%
239   mutate(
240     x_end = spelling_pre + (1 - delta) * (spelling_post - spelling_pre),
241     y_end = phonics_pre + (1 - delta) * (phonics_post - phonics_pre)
242   )
243 base.dat <- ellas.class %>%
244   filter(time == "Pre") %>%
245   mutate(
246     shape_class = ifelse(
247       instruction == "traditional",
248       ifelse(week == 1, "open_circle", "closed_circle"),
249       ifelse(week == 1, "open_triangle", "closed_triangle")
250     )
251   )
252

```

```

253 ellas.class.plt <- ggplot(
254   base.dat,
255   aes(
256     x = spelling,
257     y = phonics
258   )
259 ) +
260   # arrows
261   geom_segment(
262     data = arrows.df,
263     aes(
264       x = spelling_pre,
265       y = phonics_pre,
266       xend = x_end,
267       yend = y_end
268     ),
269     arrow = arrow(length = unit(0.05, "cm")),
270     color = "#CCC",
271     linewidth = 0.15,
272     inherit.aes = FALSE
273   ) +
274   geom_point(
275     aes(
276       color = classroom,
277       shape = shape_class
278     ),
279     stroke = 0.2,
280     show.legend = FALSE,
281   ) +
282   # labels: white for week 2
283   geom_text(
284     data = dplyr::filter(base.dat, week == 2),
285     aes(label = id),
286     color = "white",
287     size = 8,
288     vjust = 0.5,
289     hjust = 0.5,
290     show.legend = FALSE
291   ) +
292   # labels: black for week 1
293   geom_text(
294     data = dplyr::filter(base.dat, week == 1),
295     aes(label = id),
296     color = "black",
297     size = 8,
298     vjust = 0.5,
299     hjust = 0.5,
300     show.legend = FALSE
301   ) +
302   scale_color_paper() +
303   scale_shape_manual(
304     values = c(
305       open_circle = 1, # open circle
306       open_triangle = 2, # open triangle
307       closed_circle = 16, # filled circle
308       closed_triangle = 17 # filled triangle
309     )
310   ) +
311   labs(

```

```

312     x = "Spelling",
313     y = "Phonics",
314     shape = "Instruction",
315     color = "Class"
316   ) +
317   theme_paper(base_size = 10, text_size = 64) +
318   theme(
319     legend.title = element_text(size = 64),
320     legend.text = element_text(size = 64)
321   )
322
323 ggsave(
324   "Figures/observed.png",
325   ellas.class.plt,
326   width = 5,
327   height = 2.8,
328   dpi = 900
329 )
330
331 ##### Simulation Study Results #####
332
333 rr <- readRDS("rejection-rates-sim.RDS")
334 lims <- range(c(rr$phat.w, rr$phat.f), na.rm = TRUE)
335
336 rr.plt.W <- ggplot(rr, aes(spelling, phonics, fill = phat.w)) +
337   geom_raster(interpolate = FALSE) +
338   geom_text(aes(label = sprintf("%.2f", phat.w)), color = "white", size = 16) +
339   scale_fill_viridis_c(option = "turbo", limits = lims) +
340   coord_equal() +
341   labs(fill = "Rejection Rates", x = "Spelling Effect", y = "Phonics Effect") +
342   theme_paper(base_size = 10, text_size = 64) +
343   theme(legend.position = "none")
344
345 rr.plt.F <- ggplot(rr, aes(spelling, phonics, fill = phat.f)) +
346   geom_raster(interpolate = FALSE) +
347   geom_text(aes(label = sprintf("%.2f", phat.f)), color = "white", size = 16) +
348   scale_fill_viridis_c(option = "turbo", limits = lims) +
349   coord_equal() +
350   labs(fill = "Rejection Rates", x = "Spelling Effect", y = "Phonics Effect") +
351   theme_paper(base_size = 10, text_size = 64) +
352   theme(legend.position = "none")
353
354 ggsave(
355   "Figures/rrw.png",
356   rr.plt.W,
357   width = 4,
358   height = 4,
359   dpi = 900
360 )
361
362 ggsave(
363   "Figures/rrf.png",
364   rr.plt.F,
365   width = 4,
366   height = 4,
367   dpi = 900
368 )
369
370 ##### BAYESIAN ESTIMATION #####

```

```

371 allsamps <- readRDS("mcmc-samples.rds")
372
373 curriculum.effects <- allsamps[, c(7, 14)]
374 x <- allsamps[, 7] #Spelling
375 y <- allsamps[, 14] #Phonics
376 colMeans(curriculum.effects)
377 # mean(curriculum.effects[,1] > 0 | curriculum.effects[,2] > 0)
378
379 dens <- kde2d(x, y, n = 50)
380 df <- with(
381   dens,
382   expand.grid(x = x, y = y) |>
383     mutate(z = as.vector(z))
384 )
385
386
387 bayes.plt <- ggplot(df, aes(x, y, fill = z)) +
388   geom_raster(interpolate = T) +
389   scale_fill_viridis_c(option = "turbo") +
390   coord_equal() +
391   geom_vline(xintercept = 0, color = "white", linetype = 2) +
392   geom_hline(yintercept = 0, color = "white", linetype = 2) +
393   labs(fill = "Density", x = "Spelling Effect", y = "Phonics Effect") +
394   theme_paper(base_size = 10, text_size = 64) +
395   theme(
396     legend.position = "none"
397   )
398
399 ggsave(
400   "Figures/bayes.png",
401   bayes.plt,
402   width = 5,
403   height = 2.8,
404   dpi = 900
405 )
406
407 final_data <- readRDS("data/final_project_scores.rds")
408 post.scores.by.instruction <-
409   final_data %>%
410   group_by(instruction, test, time, classroom) %>%
411   summarize(mean_score = mean(score), .groups = "drop") %>%
412   ggplot(aes(
413     time,
414     mean_score,
415     group = interaction(instruction, classroom),
416     color = classroom
417   )) +
418   geom_line() +
419   geom_point(aes(shape = instruction)) +
420   scale_color_paper() +
421   facet_wrap(
422     ~test,
423     labeller = labeller(
424       test = c(
425         "phonics" = "Phonics",
426         "spelling" = "Spelling"
427       )
428     )
429   ) +

```

```
430 labs(
431   title = "Pre/Post Scores by Instruction Type",
432   y = "Mean Score (%)",
433   color = "Classroom",
434   shape = "Instruction Type"
435 ) +
436 theme_paper(base_size = 10, text_size = 64) +
437 theme(
438   legend.title = element_text(size = 64),
439   legend.text = element_text(size = 64),
440   strip.text.x = element_text(size = 64),
441   strip.text.y = element_text(size = 64)
442 )
443 ggplot(
444   post.scores.by.instruction,
445   width = 6,
446   height = 2.8,
447   dpi = 900
448 )
449 }
```