

# Modelling multilevel nonlinear treatment-by-covariate interactions in cluster randomized controlled trials using a generalized additive mixed model

Sun-Joo Cho , Kristopher J. Preacher, Haley E. Yaremych, Matthew Naveiras, Douglas Fuchs and Lynn S. Fuchs

Vanderbilt University, Nashville, Tennessee, USA

A cluster randomized controlled trial (C-RCT) is common in educational intervention studies. Multilevel modelling (MLM) is a dominant analytic method to evaluate treatment effects in a C-RCT. In most MLM applications intended to detect an interaction effect, a single interaction effect (called a *conflated* effect) is considered instead of level-specific interaction effects in a multilevel design (called *unconflated multilevel interaction* effects), and the linear interaction effect is modelled. In this paper we present a generalized additive mixed model (GAMM) that allows an unconflated multilevel interaction to be estimated without assuming a prespecified form of the interaction. R code is provided to estimate the model parameters using maximum likelihood estimation and to visualize the nonlinear treatment-by-covariate interaction. The usefulness of the model is illustrated using instructional intervention data from a C-RCT. Results of simulation studies showed that the GAMM outperformed an alternative approach to recover an unconflated logistic multilevel interaction. In addition, the parameter recovery of the GAMM was relatively satisfactory in multilevel designs found in educational intervention studies, except when the number of clusters, cluster sizes, and intraclass correlations were small. When modelling a linear multilevel treatment-by-covariate interaction in the presence of a nonlinear effect, biased estimates (such as overestimated standard errors and overestimated random effect variances) and incorrect predictions of the unconflated multilevel interaction were found.

## 1. Introduction

### 1.1. Study motivation

Intervention studies in education (e.g., concerning curriculum, policy, or instructional programmes) have become more common over the past two decades. The evaluation setting is often a randomized controlled trial (RCT). One popular RCT design in education research is the cluster RCT (C-RCT) with control and treatment groups. In the C-RCT design, clusters (e.g., schools) are randomized to the control or treatment group, and the inferential goal is to test hypotheses related to treatment effects.

\*Correspondence should be addressed to Sun-Joo Cho, Peabody Hobbs 213a, 230 Appleton Place, Nashville, TN 37203, USA (email: sj.cho@vanderbilt.edu).

The effectiveness of the intervention measures whether a programme, policy, or approach improves outcomes. In education, outcomes are commonly continuous variables (e.g., continuous scale scores of students' achievement). Multilevel modelling (MLM) is a dominant analytic method because it allows researchers to model complex patterns of variability within and across different levels of analysis (e.g., students in schools) to evaluate whether outcomes are improved due to the effects of covariates. Education researchers often report an average treatment effect (ATE), sometimes conditional on a pre-intervention covariate as a moderator (MOD; e.g., pretest scores, demographic information) (hereafter denoted  $TRT \times MOD$ , indicating that treatment effects (TRT) differ based on the values of MOD). MODs have an important role in understanding variability in treatment effects, and they are commonly used as covariates in MLM to explain variability in treatment effects. In multilevel settings commonly found in education, MODs can be assessed either at the individual level or at the cluster level. In addition, MODs can be categorical (e.g., grade levels, school types) or continuous (e.g., student pretest scores, teaching experience in years). In this study, we focus on continuous MODs assessed at the individual level and a categorical TRT at the cluster level in a C-RCT.

## 1.2. Current issues

In multilevel designs, the  $TRT \times MOD$  effect represents *multilevel interaction* or *moderation* (e.g., Preacher, Curran, & Bauer, 2006; Raudenbush & Bryk, 2002). In testing hypotheses of multilevel interaction, several conceptual and statistical problems have been discussed (Preacher, Zhang, & Zyphur, 2016). Preacher et al. (2016) noted that problems occur because most applications testing multilevel interaction do not separate level-1 and level-2 effects into their orthogonal components (in a two-level design as an example) and instead combine them into a single coefficient (called *conflation*). As a conceptual problem, conflation results in insensitivity to the theoretically meaningful ways in which multilevel interaction can occur. As a statistical problem, conflation leads to estimates of a weighted average of within- and between-cluster effects in the presence of the level-specific interaction effects. As a solution, Preacher and Sterba (2019) suggested modelling level-specific interaction ( $TRT \times MOD$ ) for fixed effects by centring a covariate (e.g., pretest) at its cluster mean (called the *unconflated* solution).

When a continuous MOD (e.g., pretest scores) and a treatment variable (e.g., control versus treatment groups) are considered in detecting  $TRT \times MOD$  interaction, the linear effect of the  $TRT \times MOD$  interaction is often modelled in MLM to detect interaction effects (Preacher & Sterba, 2019). When the treatment variable is dummy-coded (control group = 0, treatment group = 1), the MOD effect is the linear effect of MOD in the control group and the  $MOD + TRT \times MOD$  effect is the linear effect of MOD in the treatment group. As such, the expected difference between the control and the treatment condition is the marginal TRT effect plus the  $TRT \times MOD$  effect for each value of MOD. If the relationship between an outcome and MOD is incorrectly assumed to be linear, estimates of treatment effects are expected to be inaccurate (Harrell, 2015). To model nonlinear  $TRT \times MOD$  interaction, Preacher and Sterba (2019) presented a logistic function in MLM. The parametric logistic function may work well when there are floor and ceiling effects. As a more flexible approach, *smooth functions* can be used for MODs that are known to predict an outcome nonlinearly. To the best of our knowledge, smooth functions have not been applied to MLM in the context of detecting unconflated  $TRT \times MOD$  interaction effects.

### 1.3. Study purpose

The purpose of this study is to illustrate modelling of nonlinear multilevel TRT  $\times$  MOD interaction with unconfounded effects in intervention studies from C-RCT designs. With multilevel TRT  $\times$  MOD interaction with unconfounded effects, the TRT effect is estimated as a function of MOD at each level. A modelling framework to detect the nonlinear multilevel interaction effect is the generalized additive mixed model (GMM; Lin & Zhang, 1999; Wood, 2017) with an identity link function. A GMM can be considered to be a multilevel mixed model (having fixed and random effects) with an identity link in which the linear predictor partly depends on some unknown smooth functions. The *nonlinear* multilevel TRT  $\times$  MOD interaction using GMM is illustrated using an instructional intervention data set in a C-RCT and compared with the *linear* multilevel TRT  $\times$  MOD interaction from MLM. For parameter estimation, we utilize the `gmm` function in the *mgcv* package (Wood, 2019) in R (R Core Team, 2020) for maximum likelihood estimation. In the `gmm` function, smooth functions in GMM are reformulated as random effects, and parameters of GMM are estimated as parameters of generalized linear mixed-effects models (GLMMs) (Wood, 2019). In this study, the key derivations on the reformulation of smooth functions as random effects by Wood (2004, 2006, 2017, p. 239) in the statistics literature are illustrated for researchers in the social and behavioral sciences. Furthermore, the R code is provided to visualize nonlinear multilevel TRT  $\times$  MOD based on results from the `gmm` function. In addition, the accuracy and precision of parameter estimates is evaluated and the consequences of modelling *linear* multilevel TRT  $\times$  MOD are presented in the presence of nonlinear multilevel TRT  $\times$  MOD via simulation.

The remainder of this paper is organized as follows. In Section 2 we present the GMM specification, provide the estimation method using R, and describe model checking and testing. In Section 3 we illustrate the model using an empirical data set. In Section 4 we present the design of simulation studies and their results. In Section 5 we conclude with a summary and a discussion.

## 2. Methods

In this section, the GMM is specified with a comparison to MLM, and its parameter estimation method in the `gmm` function is described. In addition, testing for nonlinear TRT  $\times$  MOD interaction is explained.

### 2.1. The generalized additive mixed model

A GMM with an identity link and univariate smooth functions is written as

$$\mu = \mathbf{X}\boldsymbol{\gamma} + \mathbf{Z}\mathbf{u} + \sum_b^H f_b(x_b), \mathbf{u} \sim MN(\mathbf{0}, \boldsymbol{\Sigma}), y \sim N(\mu, \sigma^2), \quad (1)$$

where  $b$  is an index for the smooth function ( $b = 1, \dots, H$ );  $y$  is an outcome variable;  $\mathbf{X}$  is a design matrix for fixed effects;  $\mathbf{Z}$  is a design matrix for random effects;  $\boldsymbol{\gamma}$  is the vector of fixed parameters;  $\mathbf{u}$  is the vector of random parameters;  $f_b$  is the univariate smooth function for covariate  $x_b$ ;  $\boldsymbol{\Sigma}$  is a covariance matrix of the random parameters in a multivariate normal (*MN*) distribution;  $\sigma^2$  is an error variance; and  $N(\mu, \sigma^2)$  denotes a normal distribution with mean  $\mu$  and an error variance  $\sigma^2$ . Here, one can see that the

GAMM is a GLMM in which the linear predictor partly depends linearly on some unknown smooth functions ( $f_b$ ).

A general form of GAMM with an identity link (equation (1)) can be presented using an MLM specification with a smooth function for a nonlinear interaction of categorical TRT  $\times$  continuous MOD. To illustrate MLM specifications as GAMM in testing a nonlinear multilevel interaction, a level-1 MOD ( $x_{ij}$ ) and a level-2 TRT ( $z_j$ ; focal covariate) in a C-RCT are considered for a level-1 outcome ( $y_{ij}$ ) in a two-level nested design in which an individual  $i$  is nested within a cluster  $j$ . For an unconfated solution for a nonlinear multilevel interaction of categorical TRT  $\times$  continuous MOD, a level-1 MOD  $x_{ij}$  can be decomposed into uncorrelated level-1 and level-2 components by subtracting the cluster average  $x_j$  from  $x_{ij}$  (i.e.,  $x_{ij} - x_j$ ) and using  $x_j$  as a level-2 MOD:  $x_{ij} = (x_{ij} - x_j) + x_j$ . In this design, researchers can test whether the level-1 part of the MOD ( $x_{ij} - x_j$ ) or level-2 part of the MOD ( $x_j$ ) moderates a level-2 TRT ( $z_j$ ) effect on an outcome variable ( $y_{ij}$ ). Below, with the unconfated level-1 MOD ( $x_{ij} - x_j$ ) and the level-2 MOD ( $x_j$ ), we first present a multilevel model specification for a *linear* multilevel interaction of categorical TRT  $\times$  continuous MOD for comparison purposes and then present a multilevel model specification for a *nonlinear* multilevel interaction with smooth functions as a special case of the GAMM.

Following the multilevel model specification, notation, and symbols in Raudenbush and Bryk (2002), the multilevel model with a random intercept  $\beta_{0j}$ , a random slope  $\beta_{1j}$ , and a *linear* interaction of dummy-coded TRT  $\times$  continuous MOD is written as follows for a two-level nested design. A level-1 model is given by

$$y_{ij} = \beta_{0j} + \beta_{1j}(x_{ij} - x_j) + r_{ij},$$

a level-2 model is presented as

$$\beta_{0j} = \gamma_{00} + \gamma_{01}x_j + \gamma_{02}z_j + \gamma_{03}x_jz_j + u_{0j},$$

$$\beta_{1j} = \gamma_{10} + \gamma_{11}z_j + u_{1j},$$

and the reduced form is expressed as

$$y_{ij} = \gamma_{00} + \gamma_{10}(x_{ij} - x_j) + \gamma_{01}x_j + \gamma_{02}z_j + \gamma_{03}x_jz_j + \gamma_{11}(x_{ij} - x_j)z_j + u_{0j} + (x_{ij} - x_j)u_{1j} + r_{ij}, \tag{2}$$

where  $\gamma_{00}$  is a fixed intercept;  $\gamma_{10}$  is a fixed effect of a level-1 component ( $x_{ij} - x_j$ ) of MOD  $x_{ij}$  where  $z_j = 0$ ;  $\gamma_{01}$  is a fixed effect of a level-2 component  $x_j$  of MOD  $x_{ij}$  where  $z_j = 0$ ;  $\gamma_{02}$  is a fixed effect of a dummy coded level-2 TRT  $z_j$  (a conditional treatment effect where  $x_j = 0$ );  $\gamma_{03}$  is a fixed linear interaction of a level-2 component of MOD  $x_{ij}$  ( $x_j$ ) and a dummy coded level-2 TRT  $z_j$ ;  $\gamma_{11}$  is a fixed linear interaction of a level-1 component of MOD  $x_{ij}$  ( $x_{ij} - x_j$ ) and a dummy coded level-2 TRT  $z_j$ ;  $u_{0j}$  is a random intercept;  $u_{1j}$  is a random slope of a level-1 component ( $x_{ij} - x_j$ ) of MOD  $x_{ij}$ ; and  $r_{ij}$  is random error. The random effects,  $[u_{0j}, u_{1j}]'$ , are assumed to follow a multivariate normal distribution,  $[u_{0j}, u_{1j}]' \sim MN(0, \Sigma)$ , with a random intercept variance  $\tau_{00}$ , a random slope variance  $\tau_{11}$ , and a covariance  $\tau_{01}$  in  $\Sigma$ . The random error,  $r_{ij}$ , is assumed to follow a normal distribution,  $r_{ij} \sim N(0, \sigma^2)$ .

As a special case of the GAMM, the multilevel model with a random intercept, a random slope of a level-1 covariate, and a smooth function for a *nonlinear* multilevel interaction of dummy-coded TRT  $\times$  continuous MOD is written as follows:

$$y_{ij} = \gamma_{00} + \gamma_{02}z_j + f_1(x_j)(z_j = 0) + f_1(x_j)(z_j = 1) + f_2(x_{ij} - x_j)(z_j = 0) + f_2(x_{ij} - x_j)(z_j = 1) + u_{0j} + (x_{ij} - x_j)u_{1j} + r_{ij}, \quad (3)$$

where  $\gamma_{02}$  is a mean of all smooth functions when  $z_j$  ( $z_j = 0$  for a control group;  $z_j = 1$  for a treatment group) is specified as a factor in R;  $f_1(x_j)(z_j = 0)$  is a smooth function of a level-2 component  $x_j$  of MOD  $x_{ij}$  where  $z_j = 0$  (i.e., nonlinear level-2 interaction for a control group);  $f_1(x_j)(z_j = 1)$  is a smooth function of a level-2 component  $x_j$  of MOD  $x_{ij}$  where  $z_j = 1$  (i.e., nonlinear level-2 interaction for a treatment group);  $f_2(x_{ij} - x_j)(z_j = 0)$  is a smooth function of a level-1 component  $x_{ij} - x_j$  of MOD  $x_{ij}$  where  $z_j = 0$  (i.e., nonlinear level-1 interaction for a control group); and  $f_2(x_{ij} - x_j)(z_j = 1)$  is a smooth function of a level-1 component  $x_{ij} - x_j$  of MOD  $x_{ij}$  where  $z_j = 1$  (i.e., nonlinear level-1 interaction for a treatment group).

In MLM (equation (2)) and GAMM (equation (3)), the same fixed intercept terms for TRT ( $z_j$ ),  $\gamma_{00} + \gamma_{02}z_j$ , are specified. However, different slope terms of MOD ( $x_{ij}$ ) for TRT ( $z_j$ ) are specified in MLM and GAMM. In MLM,  $\gamma_{03}x_jz_j$  and  $\gamma_{11}(x_{ij} - x_j)z_j$  terms are for modelling *linear* multilevel interaction: the  $\gamma_{03}x_jz_j$  is for the linear TRT  $\times$  level-2 MOD interaction and  $\gamma_{11}(x_{ij} - x_j)z_j$  is for the linear TRT  $\times$  level-1 MOD interaction. In GAMM, the  $f_1(x_j)(z_j = 0)$ ,  $f_1(x_j)(z_j = 1)$ ,  $f_2(x_{ij} - x_j)(z_j = 0)$ , and  $f_2(x_{ij} - x_j)(z_j = 1)$  are for modelling *nonlinear* multilevel interactions: the  $f_1(x_j)(z_j = 0)$  and  $f_1(x_j)(z_j = 1)$  are for the nonlinear TRT  $\times$  level-2 MOD interaction, and  $f_2(x_{ij} - x_j)(z_j = 0)$  and  $f_2(x_{ij} - x_j)(z_j = 1)$  are for the nonlinear TRT  $\times$  level-1 MOD interaction.

### 2.1.1. Smooth functions for categorical TRT $\times$ Continuous MOD

The univariate smooth function  $f_b(x_b)$  of a covariate  $x_b$  is specified as a weighted sum of a set of basis functions over the covariate  $x_b$ :

$$f_b(x_b) = \sum_{k=1}^K \delta_{bk} b_{bk}(x_b), \quad (4)$$

where  $k$  is an index for a basis function ( $k = 1, \dots, K$ ),  $x_b$  is a covariate for a smooth function  $b$ ,  $\delta_{bk}$  is a basis coefficient, and  $b_{bk}(x_b)$  is the  $k$ th basis function for smooth function  $b$ . The basis functions ( $\mathbf{b}_b = [b_{b1}, \dots, b_{bK}]'$ ) are a set of known curves to represent  $f_b(x_b)$  and they are functioned as covariates to estimate basis coefficients ( $\delta_b = [\delta_{b1}, \dots, \delta_{bK}]'$ ). In the *mgcv* package, a smooth function is estimated with an identification constraint such that  $f_b$  sums to 0 over the observed covariate values (i.e.,  $\sum f_b(x_{bv}) = 0$  for each  $b$ , where  $v$  is an index for observations); otherwise,  $f_b(x)$  can be confounded with the intercept. When the TRT is specified as a factor, the *mgcv* package automatically computes a separate smooth function for the MOD effect, for every level in TRT (Wieling, 2018).

In GAMM applications using the *mgcv* package, a cubic regression spline (CRS; Wood, 2017, Section 5.3.1) and a thin plate regression spline (TPRS; Wood, 2017, Section 5.5.1) are commonly used splines for the univariate smooth function ( $f_b(x_b)$ ). The CRS is a smooth curve made up of sections of cubic polynomials. The sections are joined together

at locations referred to as *knots*. At each knot, the joined sections of the cubic polynomials have equivalent values, first and second derivatives (Wood, 2017, Section 5.3.1). In the *mgcv* package, the default is for the knots to be equally spaced over the entire range of the observed covariate, and the number of knots is the same as the number of basis functions ( $K$ ). The CRS and the TRPS yield comparable results for the univariate smooth function (e.g., Finch & Finch, 2018), although the CRS yields better computational efficiency; therefore, we use the CRS in the current study.

As shown in equation (3), a nonlinear categorical TRT  $\times$  continuous MOD interaction is specified by including different smooth functions of a continuous MOD multiplied by a dummy-coded TRT,  $f_1(x_j)z_j$  and  $f_2(x_{ij} - x_j)z_j$ . In Appendix S1, the CRS for a smooth function of  $z_j = 0$  or  $z_j = 1$  is illustrated.

For the selected basis functions for smooth functions, the number of basis functions ( $K$ ) should be selected to obtain a good fit. The dimensionality of the basis expansion is determined by  $K$ . When  $K$  is too small, oversmoothing will occur, and when  $K$  is too large, computation time is increased.  $K = 10$  is the default in *mgcv*, and is often sufficient in generalized additive modelling (e.g., Bringmann et al., 2017). Thus, we set  $K = 10$  in the current study. To determine whether a selected  $K$  is large enough, the value of the  $k$ -index can be assessed. The  $k$ -index is a measure of the remaining pattern in the residuals. Let  $\tilde{\mathbf{r}}$  denote the vector of residuals  $r_{ij}$ , ordered according to the value of covariate  $x_b$  and define *differencing residuals* that are near neighbours according to the covariate of the smooth as  $\Delta_{ij} = \tilde{r}_{(ij)+1} - \tilde{r}_{ij}$ . The  $k$ -index is calculated as the ratio of (a) an estimate of the mean of the squared differencing residuals ( $\sigma_\Delta^2 = E[\Delta_{ij}^2]$ ) to (b) an estimate of residual variance from a model fit ( $\sigma^2$ ) (Wood, 2017, pp. 243, 330). A  $k$ -index below 1 indicates that there is a missed pattern left in the residuals with a specified  $K$ , and a larger  $K$  should be considered in this case. The  $k$ -index can be obtained through the *gam.check* function in *mgcv*.

In addition to the  $k$ -index, the corrected Akaike information criterion (correctedAIC; Wood, Pya, & Säfken, 2016) is considered to select a model with an adequate amount of smoothing from the data among candidate models differing in  $K$ , as a commonly used model selection criterion in generalized additive modelling (Ruppert, Wand, & Carroll, 2003, p. 120). The correctedAIC for GAMM uses the effective degrees of freedom (*edf*) as the number of parameters needed to represent smooth functions in the penalty term of the Akaike information criterion (AIC; Akaike, 1974). The correctedAIC is specified as follows:

$$\text{CorrectedAIC} = -2ll + (2 \times edf), \tag{5}$$

where  $ll$  is the log-likelihood. The  $ll$  and *edf* in the correctedAIC for GAMM can be extracted using the function *logLik.gam* for a fitted model in the *mgcv* package.

**2.2. Parameter Estimation**

The *glm* function in the *mgcv* package was used for maximum likelihood estimation. Below, we describe the details of the implementation of the *glm* function for the specified GAMM (equation (3)).

The ‘wiggleness’ of the smooth function  $f_b(x_b)$  is controlled less by  $K$  (the number of basis functions) than by a quadratic smoothing penalty (e.g., Wood, 2017). The quadratic smoothing penalty for the model can be written as

$$\lambda_b \delta_b^T \mathbf{S}_b \delta_b, \quad (6)$$

where  $\lambda_b$  is a smoothing parameter,  $\delta_b$  is a vector of basis coefficients, and  $\mathbf{S}_b$  is a penalty matrix embedded as a diagonal block in a matrix. For smooth functions, the elements of  $\mathbf{S}_b$  are known and are determined by the chosen basis functions. The parameter  $\lambda_b$  controls the trade-off between goodness of fit and model smoothness.

For the identity link, the glmm function uses a GLMM formulation to fit a GAMM. Wood (2004, 2006, 2017, p. 239) showed how a smooth function in a GAMM can be reformulated into fixed and random effects in a GLMM. Key derivations in Wood (2004, 2006, 2017, p. 239) are explained and illustrated in Appendix S2.

### 2.3. Testing for nonlinear TRT×MOD interactions

To determine whether or not the smooth function  $f_b(x_b)$  is distinguishable from zero, the null hypothesis  $H_0 : f_b(x_b) = 0$  for all  $x_b$  in the range of interest can be tested. A test statistic for  $f_b(x_b)$  is as follows:

$$T_r = \hat{\mathbf{f}}_b^T \mathbf{V}_{f_b}^- \hat{\mathbf{f}}_b, \quad (7)$$

where  $r$  is the rounded effective degrees of freedom (*edf*) for  $f_b(x_b)$  (integer; e.g.,  $r = 1$  in the case of  $edf = 1.45$ ),  $\hat{\mathbf{f}}_b$  is the vector of  $f_b(x_b)$  evaluated at the observed predictor values, and  $\mathbf{V}_{f_b}^-$  is a pseudo-inverse of  $\mathbf{V}_{f_b}$  of rank  $r$  ( $\mathbf{V}_{f_b} = \mathbf{X}\mathbf{V}_\delta\mathbf{X}^T$ , where  $\mathbf{X}$  are basis functions and  $\mathbf{V}_\delta$  is the covariance matrix of basis coefficient estimates) (Wood, 2017, pp. 305-306). Under  $H_0$ , the test statistic  $T_r$  follows a chi-square distribution ( $T_r \sim \chi_r^2$ ) (Wood, 2013).

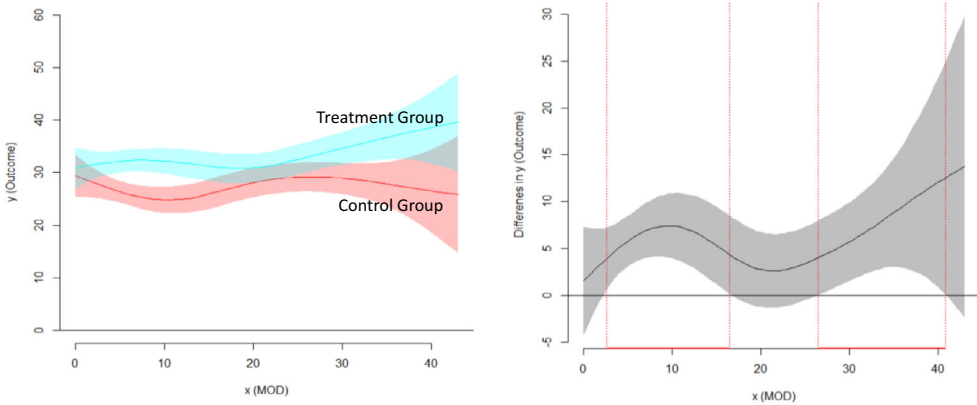
Smooth functions have confidence intervals around them, which are obtained by taking the quantiles from the posterior distribution of the  $f_b(x_b)$  (Marra & Wood, 2012). To calculate the distribution of the  $f_b(x)$ , a large number (e.g., 1,000) of basis coefficient parameters ( $\delta_b$ ) are simulated from the posterior distributions of basis coefficients using a multivariate normal distribution:

$$\delta_b \sim MN(\hat{\delta}_b, \mathbf{V}_\delta), \quad (8)$$

where  $\hat{\delta}_b$  contains basis coefficient estimates. Then a large number of the  $f_b(x_b)$  can be calculated using sampled basis coefficients and basis functions using equation (4). The .025 and .975 quantiles of the posterior distribution can be used for the lower and upper bounds of a 95% confidence interval of the smooth functions.

In addition to significance testing for the smooth function, we can visualize over what ranges of MOD the smooth functions differ significantly (called *the region of significance*). As an example, varying treatment effects depending on the levels of  $x$  can be estimated for each level of a TRT variable (one for a control group and another for a treatment group) using the GAMM specification (equation (3)), as shown in Figure 1 (left). Based on the result of the smooth functions, *differences* in the smooth functions in Figure 1 (left) (a smooth function for a treatment group minus a smooth function for a

<sup>1</sup> In the `gamm` output, the degrees of freedom used in computing test statistics and  $p$ -values are presented as `Ref.df`.



**Figure 1.** Varying effects of MOD  $x$  on outcome  $y$  by the TRT  $z$  modelled with smooth functions with confidence bands (left) and differences in outcome  $y$  between the two smooth functions (the smooth function for the treatment group minus the smooth function for the control group) with confidence bands (right). Vertical lines in Figure 1 (right) indicate windows of significant differences.

control group) can be presented as in Figure 1 (right). Windows of significant differences are found in ranges between 2.606 and 16.505, and between 26.495 and 40.828 (noted with vertical bars in Figure 1 (right)) in  $x$ .

### 3. Empirical study

In this section a GAMM specification is illustrated using an empirical data set to detect a nonlinear multilevel interaction of categorical level-2 TRT and continuous level-1 MOD in a C-RCT design. The data set comes from Fuchs et al. (2021). The purpose of the study was to evaluate the efficacy of two revised versions of first-grade Peer-Assisted Learning Strategies (PALS) called the PALS-Only and PALS+Fluency programmes. Using a subset of the data from Fuchs et al. (2021), an analysis goal in the present study is to test whether the contrast of both PALS conditions against control moderates the effect of students' pre-treatment scores of phonological awareness (prePA) on students' post-treatment scores of phonological awareness (postPA).

#### 3.1. Data description

##### 3.1.1. Participants

Teachers from 33 first-grade classrooms in eight elementary schools and their 491 students participated. In the C-RCT, the 33 teachers were assigned randomly within schools to a control group (11 teachers and their 171 students) and two treatment groups – the PALS-Only (11 teachers and their 168 students) and PALS+Fluency groups (11 teachers and their 152 students). Cluster size (the number of students per classroom) ranged from 1 to 18 (median = 15, semi-interquartile range = 1). One classroom had one student in a PALS+Fluency group. Fuchs et al. (2021) reported that there were no statistically significant differences among the three study groups on student demographics, teacher demographics, or pre-treatment reading performance.



### 3.1.2. Measures

Students were tested before and immediately following the 22-week treatment period. The same measures, used by Yopp (1988) for segmenting sounds in words and by Fuchs, Fuchs, Hosp, and Jenkins (2001) for blending sounds in words, were used to assess phonological awareness at two time points. In this study, the same continuous scores of phonological awareness were used in GAMM and MLM analyses as were used in Fuchs et al. (2021). There are no missing data in prePA or postPA. As descriptive information, the normality assumption of the postPA scores is tested. A Shapiro–Wilk test indicated that the postPA scores were significantly non-normal ( $W = 0.950$ ,  $p < .0001$ ). However, the deviance from normality is not large in the quantile–quantile plot and it is localized in the tails of the distribution. As employed in Fuchs et al. (2021), zero-centring of the prePA and postPA scores was used. The mean and standard deviation of prePA scores are 0 and 0.911, respectively, and the mean and standard deviation of postPA scores are 0 and 0.842, respectively.

## 3.2. Analyses and results

Below, we describe how to test whether the contrast of both PALS conditions against the control group moderates the effect of prePA on students' postPA. The R code for the empirical data analyses is presented in Appendix S3.

### 3.2.1. Step 1: Fitting the unconditional GAMM

The postPA results are from a nested data structure: 491 students (level 1) nested within 33 classes (level 2), nested within eight schools (level 3). Dependencies in the postPA scores due to clusters (classes and schools) can be accounted for in the three-level model. However, there was a convergence problem with estimating the variance of the random intercept for schools in the unconditional three-level random intercept model.<sup>1</sup> This problem may be due to having only eight schools, which is fewer upper-level units than recommended in MLM (Snijders & Bosker, 2012). In this circumstance, it is recommended to replace a random intercept with  $L-1$  dummy codes (where  $L$  is the number of clusters) for cluster membership (e.g., McNeish & Stapleton, 2016). The unconditional GAMM with the  $L-1$  dummy codes for eight schools is specified as follows (the model can be called an MLM because smooth functions have not yet been introduced):

$$y_{ij} = \gamma_{00} + \sum_{l=2}^L \alpha_l D_l + u_{0j} + r_{ij}, \quad (9)$$

where  $l$  is an index for a school ( $l = 2, \dots, L$ ;  $L = 8$  in this example),  $D_l$  is a dummy code for school membership with the first school as the reference school, and  $\alpha_l$  is the fixed effect of  $D_l$ . The intraclass correlation coefficient (ICC), calculated based on the results of equation (9), was .107 ( $= 0.066/[0.066 + 0.552]$  where 0.066 is  $\hat{\tau}_{00}$  and 0.552 is  $\hat{\sigma}^2$ ), which suggests that there is non-ignorable dependency in postPA scores due to class clustering.

<sup>1</sup> The warning message from the *lmer* function in the *lme4* package is 'Model failed to converge with max|grad| = 0.00206106 (tol = 0.002, component 1)'.

### 3.2.2. Step 2: Adding covariates (TRT and MOD) and comparing and checking models

In Step 2, a dummy-coded level-2 TRT ( $z_j = 0$  for a control group;  $z_j = 1$  for PALS-Only and PALS+Fluency groups) and a cluster (class)-centred level-1 MOD (prePA  $x_{ij}$ ;  $x_{ij} - x_j$  and  $x_j$ ) were considered as covariates. Prior to modelling, the relationship between postPA scores ( $y_{ij}$ ) and prePA scores ( $x_{ij} - x_j$  or  $x_j$ ) was explored by TRT groups ( $z_j$ ) using scatter plots. Figure 2 (top) presents a scatter plot of  $y_{ij}$  against  $x_{ij} - x_j$  by  $z_j$ , and Figure 2 (bottom) presents a scatter plot of  $y_{ij}$  against  $x_j$  by  $z_j$ . This figure shows nonlinear relationships presented with smooth lines deviant from the linear dotted lines at each level. In addition, the figure shows that the differences in postPA between the two groups (presented with 95% confidence bands) differ depending on the levels of prePA. Given the patterns identified in Figure 2, smooth functions of  $f_1(x_j)(z_j = 0)$ ,  $f_1(x_j)(z_j = 1)$ ,  $f_2(x_{ij} - x_j)(z_j = 0)$ , and  $f_2(x_{ij} - x_j)(z_j = 1)$  (along with  $z_j$ ) were added to equation (9).

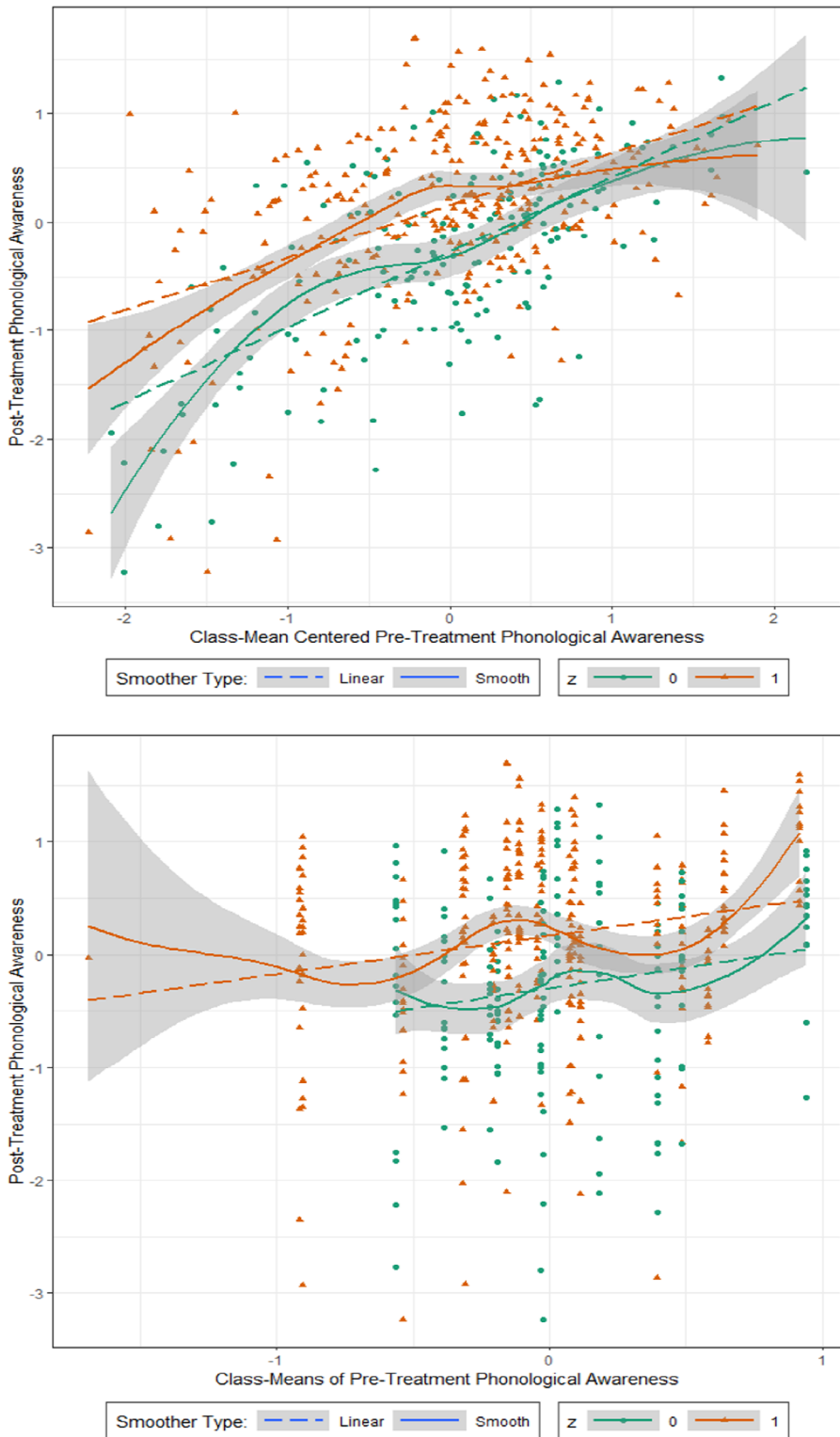
The random-intercept-and-slope GMM and the random-intercept GMM were fitted with 10 basis functions ( $K = 10$  as the *mgcv* default). The  $k$ -index was close to 1 and the correctedAIC differs in the first decimal place for the models with  $K = 10$  (correctedAIC = 897.706), 12 (correctedAIC = 897.683), 14 (correctedAIC = 897.606), and 16 (correctedAIC = 897.453). These results indicate that  $K = 10$  is adequate to obtain a good fit in both models.

Results of the random-intercept-and-slope GMM were compared with those of the random-intercept GMM. In the random-intercept-and-slope GMM,  $\hat{\tau}_{11}$  was  $3.699829 \times 10^{-9}$ . In addition, estimates and standard errors of fixed effects differed in the second or third decimal places and patterns in the smooth functions were similar between the two models. Furthermore, the AIC and the Bayesian information criterion (BIC; Schwarz, 1978) suggested that the random-intercept GMM fits better than the random-intercept-and-slope GMM (see the AIC and BIC values in Table 1 (top)). Thus, the random-intercept GMM was chosen for result interpretations. Residual analysis of the model indicates that there is evidence of good model-data fit (see Appendix S4).

### 3.2.3. Step 3: Interpreting results

Table 1 (bottom) presents results of the random-intercept GMM, compared with the MLM results we will discuss in the following subsection. The GMM results are interpreted below.

A significant fixed conditional TRT ( $z_j$ ) effect was found ( $\hat{\gamma}_{02} = 0.439$ , SE = 0.071). This result means that students in PALS-Only and PALS+Fluency programs together outperformed control students. The corresponding effect size (Hedges's  $g$ ) is 0.537, following the guideline suggested by What Works Clearinghouse (2017). When TRT is a focal covariate at the class level (level 2) and MOD is a moderator, the level-2 TRT effect ( $z_j$ ) on class means ( $y_j$ ) at any chosen value of the level-2 part of MOD ( $x_j$ ) is of interest to interpret. Figure 3(a) presents the effect of  $z_j$  on  $y_j$  by quantiles (.1, .25, .5, .75, .9) of  $x_j$ , and Figure 3(b) shows the level-2 TRT effects ( $\hat{\gamma}_{02} + \{f_1(x_j)(z_j = 1) - f_1(x_j)(z_j = 0)\}$ ) against  $x_j$  from the GMM. The region of significance of class-level prePA ( $x_j$ ) was  $[-0.494, 0.568]$ , presented in the vertical lines of Figure 3(b). PostPA scores were higher for the treatment groups than for the control group in the range of  $[-0.494, 0.568]$ . However, the difference at the extremes was not significant (see Figure 3(b)).



**Figure 2.** Empirical study: Scatter plots of  $y_{ij}$  against  $x_{ij} - x_j$  by  $z_j$  (top) and  $y_{ij}$  against  $x_j$  by  $z_j$  (bottom) (raw data).

**Table 1.** Empirical study: Results of model selection (top) and results (bottom) of GMM and MLM

Model	GMM		MLM	
	AIC	BIC	AIC	BIC
Random-intercept	895.59	975.32	911.98	974.93
Random-intercept-and-slope	899.59	987.71	*	*

Fixed effects	GMM		MLM	
	EST	SE	EST	SE
Intercept[ $\gamma_{00}$ ]	<b>0.384</b>	0.133	<b>0.318</b>	0.149
$x_{ij} - x_j[\gamma_{10}]$	-		<b>0.693</b>	0.056
$x_j[\gamma_{01}]$	-		<b>0.537</b>	0.189
$z_j[\gamma_{02}]$	<b>0.439</b>	0.071	<b>0.463</b>	0.081
$x_j z_j[\gamma_{03}]$	-		-0.183	0.209
$(x_{ij} - x_j)z_j[\gamma_{11}]$	-		<b>-0.212</b>	0.070
$D_2[\alpha_2]$	-0.216	0.173	-0.304	0.206
$D_3[\alpha_3]$	<b>-1.123</b>	0.150	<b>-1.056</b>	0.180
$D_4[\alpha_4]$	<b>-0.353</b>	0.155	<b>-0.381</b>	0.183
$D_5[\alpha_5]$	<b>-0.928</b>	0.139	<b>-0.877</b>	0.161
$D_6[\alpha_6]$	<b>-0.680</b>	0.152	<b>-0.689</b>	0.181
$D_7[\alpha_7]$	<b>-0.713</b>	0.220	<b>-0.506</b>	0.232
$D_8[\alpha_8]$	<b>-0.679</b>	0.146	<b>-0.594</b>	0.168

Random effects	EST	EST
$\tau_{00}$	0.010	0.023
$\sigma^2$	0.316	0.337

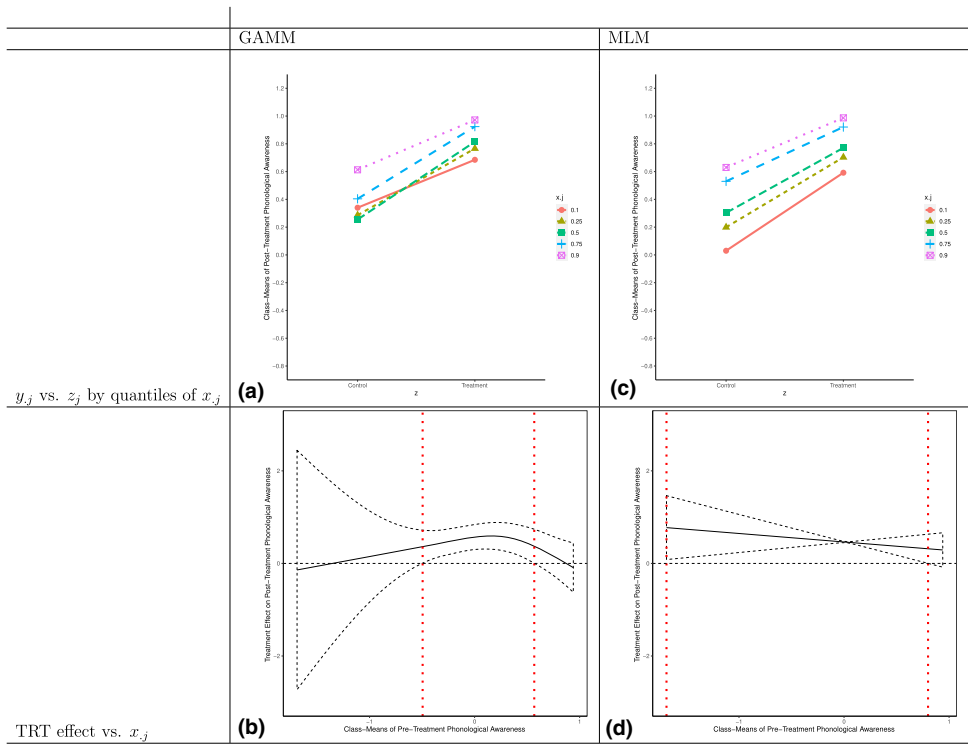
Smooth functions	Ref. <i>edf</i>	$T_r$ ( <i>p</i> -value)	
$f_2(x_{ij} - x_j)(z_j = 0)$	3.630	47.747 ( $< 2 \times 10^{-16}$ )	--
$f_2(x_{ij} - x_j)(z_j = 1)$	3.153	51.731 ( $< 2 \times 10^{-16}$ )	--
$f_1(x_j)(z_j = 0)$	2.270	7.728 (.000165)	--
$f_1(x_j)(z_j = 1)$	1.000	3.981 (.046566)	--

Note. - indicates a parameter which was not considered under GMM; -- indicates a smooth function which was not considered under MLM; Significance for fixed effects in bold based on *t*-test at  $\alpha = .05$ ; \* indicates that AIC and BIC were not reported because there was a convergence problem with estimating the random-intercept-and-slope model.

### 3.3. Comparisons between GMM and MLM

Because MLM is a dominant analytic method to evaluate TRT effects in a C-RCT, the results of GMM and MLM are compared for instructive purposes in Table 1 (bottom). Similarly to the GMM, there was a convergence problem with estimating a random intercept variance for classes in the unconditional three-level MLM.  $L-1$  dummy codes for school memberships (where  $L$  is the number of schools) were considered to account for dependency due to school clustering. In addition, there was a convergence problem with estimating  $\tau_{11}$  in the random-intercept-and-slope model (equation (2)).<sup>2</sup> Thus, the

<sup>2</sup>The error message from the lme function was 'nlminb problem, convergence error code = 1 message = iteration limit reached without convergence (10)'.



**Figure 3.** Empirical study: Probing an interaction between a level-2 TRT (focal covariate) and a level-2 part of MOD for GAMM and MLM. Vertical lines in (b) and (d) indicate windows of significant differences.

random-intercept MLM with dummy codes for school memberships was considered for comparison with the random-intercept GAMM.

Regarding model selection between the random-intercept GAMM and the random-intercept MLM as shown in Table 1 (top), the random-intercept GAMM fits better than the random-intercept MLM based on the AIC, whereas the random-intercept MLM fits better than the random-intercept GAMM based on the BIC. However, the differences in the BIC between the two models were small (975.32 for GAMM and 974.93 for MLM, yielding a difference of 0.39), indicating that there is no strong evidence that the random-intercept MLM fits better than the random-intercept GAMM based on the BIC.

There were similar patterns in the effects of school memberships and variances of random effects between GAMM and MLM, as shown in Table 1 (bottom). The standard errors of fixed effects were larger in MLM than in GAMM. For the comparison with an interaction between a level-2 TRT ( $z_j$ ) and the level-2 part of MOD ( $x_j$ ) from GAMM, Figure 3(c) presents the effect of  $z_j$  on  $y_j$  by quantiles (.1,.25,.5,.75,.9) of  $x_j$ , and Figure 3(d) shows the level-2 TRT effects ( $\hat{\gamma}_{02} + \hat{\gamma}_{03}x_j$ ) against  $x_j$  from MLM. In Figure 3(d), 95% confidence bands were calculated as  $(\hat{\gamma}_{02} + \hat{\gamma}_{03}x_j) \pm z_{crit} \sqrt{\text{Var}(\hat{\gamma}_{02}) + \text{Var}(\hat{\gamma}_{03})x_j^2 + 2x_j\text{Cov}(\hat{\gamma}_{02}, \hat{\gamma}_{03})}$ . Unlike GAMM (Figure 3(a)), postPA scores were higher for the treatment groups than for the control group

over all quantiles of  $x_j$ , as shown in Figure 3(c). The region of significance of class-level prePA ( $x_j$ ) was  $[-1.690, 0.800]$ , presented in the vertical lines of Figure 3(d).

### 4. Simulation study

A simulation study was designed to investigate the relative performance of GAMM in detecting a nonlinear multilevel TRT  $\times$  MOD interaction compared with an alternative approach (MLM with a logistic function [MLM-Logistic]) in the presence of the level-specific logistic (parametric) form of the interaction; and the accuracy of GAMM (equation (3)) parameter estimates and their precision (standard errors). For both foci, the results of modelling a *nonlinear* TRT  $\times$  MOD in the GAMM are compared with those of modelling a *linear* TRT  $\times$  MOD in the MLM.

#### 4.1. Simulation design

A two-level nested design (e.g., students nested within schools) was chosen for the simulation study. As is common in education intervention studies, a balanced design for control and treatment groups was used. The following simulation conditions were varied because they are expected to affect parameter recovery and precision in multilevel designs (Lüdtke, Marsh, Robitzsch, & Trautwein, 2011; Preacher, Zhang, & Zyphur, 2011): the number of clusters, cluster sizes (i.e., the number of individuals within a cluster), and ICC for outcomes. The levels of these three simulation conditions were chosen based on literature reviews of study designs for educational intervention studies in 30 papers published in the *Journal of Educational Psychology*, *American Educational Research Journal*, and *Exceptional Children*, and in other related MLM work. The number of clusters in 30 papers we reviewed ranged from 28 to 225 (median = 70, semi-interquartile range = 25). To mimic these numbers of clusters, the number of clusters was selected as 30 (small), 70 (medium), and 200 (large). Cluster sizes in 30 papers ranged from 13 to 35 (median = 16, semi-interquartile range = 7). Balanced cluster sizes were selected as 15 and 30 in the simulation study. Eight of 30 papers reported ICCs ranging from .08 to .25. The ICC of outcomes was selected to be at .05, .10, or .30. ICC values are rarely greater than .30 in educational and organizational studies (e.g., Fox, 2010).

For the first focus (the performance of GAMM in the presence of the logistic form of the TRT  $\times$  MOD), the data-generating model is MLM-Logistic for the level-specific TRT  $\times$  MOD. The MLM with a logistic function as a data-generating model is written as.

$$y_{ij} = \frac{1}{1 + \exp[-\mu]} + u_{0j} + (x_{ij} - x_j)u_{1j} + r_{ij}, \tag{10}$$

where  $\mu = \gamma_{00} + \gamma_{10}(x_{ij} - x_j) + \gamma_{01}x_j + \gamma_{02}z_j + \gamma_{03}x_jz_j + \gamma_{11}(x_{ij} - x_j)z_j$ , and the error term  $r_{ij}$  is assumed to be distributed independently as  $N(0, \sigma^2)$ . Fixed parameters were selected to generate no differences in logistic functions between control and treatment groups at level 1 and to have large differences in logistic functions between control and treatment groups at level 2 as expected in a C-RCT:  $\gamma_{00} = 0$ ,  $\gamma_{10} = 1.5$ ,  $\gamma_{01} = 7.2$ ,  $\gamma_{02} = 1.7$ ,  $\gamma_{03} = 5.1$ , and  $\gamma_{11} = 0.1$ . The three levels of ICC were manipulated using the ‘true’ variances of random intercept and random errors. That is, given the error variance  $\sigma^2 = 0.6$ , the three levels of  $\tau_{00}$  were calculated as 0.032, 0.067, and 0.257 for ICC = .05, .10, and .30, respectively. The slope variance  $\tau_{11}$  was set as 0.1 and the covariance  $\tau_{01}$  was

set to 0. In generating the logistic functions for control and treatment groups at level 2, it was assumed that students with low and high values of MOD would not benefit from TRT and students in the middle values of MOD might gain from TRT, but to some extent conditional on TRT (Preacher & Sterba, 2019). The generated logistic function in MLM with  $K = 10$  is shown for one condition (number of clusters = 200, cluster size = 30, ICC = .30) in Appendix S5 as an example.

For the second focus (parameter recovery of GAMM), the data-generating model is a random-intercept-and-slope GAMM (equation (3)). Fixed parameters,  $\gamma_{00} = 0.4$  and  $\gamma_{02} = 0.4$ , were selected to mimic the results of the empirical study. The same  $\sigma^2$ ,  $\tau_{00}$ ,  $\tau_{11}$ , and  $\tau_{01}$  used for the first focus were also used in the second focus. For 'true' smooth functions, increasing nonlinear functions were generated using equation (4) with  $K = 10$ . According to literature reviews of educational intervention studies in the 30 papers mentioned earlier, educational interventions are implemented mainly to improve learning for students with low achievement levels. Thus, larger nonlinear treatment effects were generated at the lower end and smaller differences were generated in the other ranges of a covariate, assuming that the intervention is the most effective for them. Parameters of basis coefficients ( $\delta_b$ ) and generated smooth functions are shown for one condition (number of clusters = 200, cluster size = 30, ICC = .30) in Appendix S5 for illustrative purposes.

In the two data-generating models, the  $x_{ij}$  was generated from a standard normal distribution and then  $x_j$  and  $x_{ij} - x_j$  were calculated. For each simulation condition, the same MOD ( $x_{ij} - x_j$  and  $x_j$ ) and generated functions for TRT  $\times$  MOD were used across replications, and random effects were generated at each replication.

The simulation conditions regarding multilevel designs were fully crossed, yielding 18 (= 3 numbers of clusters  $\times$  2 cluster sizes  $\times$  3 ICCs) conditions. Five hundred replications were simulated for each condition. For MLM-Logistic and GAMM as data-generating models, GAMM (equation (3)) was fitted to the generated data sets. In addition, for MLM-Logistic and GAMM as data-generating models, MLM (equation (2)) was fitted to the same generated data sets to demonstrate the consequences of modelling *linear* categorical TRT  $\times$  continuous MOD interactions in the presence of nonlinear categorical TRT  $\times$  continuous MOD interactions. And for MLM-Logistic as a data-generating model, MLM-Logistic (see Appendix S5 for estimation in R) was fitted to the same generated data sets to compare its results with GAMM's results. In addition, five candidate models for different values of  $K$  ( $K = 6, 8, 10, 12, 14$ ) were fitted to the generated data sets for each replication in a condition to check whether the  $K$  used in generated smooth functions was adequate based on the correctedAIC. Thus, the total number of fitted models is 225,000 (18 multilevel designs  $\times$  500 replications  $\times$  3 models [MLM-Logistic, GAMM, and MLM]  $\times$  5 models for different values of  $K$  for MLM-Logistic as a data-generating model = 135,000; 18 multilevel designs  $\times$  500 replications  $\times$  2 models [GAMM and MLM]  $\times$  5 models for different values of  $K$  for GAMM as a data-generating model = 90,000).

#### 4.2. Analysis

For the first focus, the 'true' level-specific TRT  $\times$  MOD generated using MLM-Logistic is compared with the predicted level-specific TRT  $\times$  MOD generated using MLM-Logistic, MLM, and GAMM, respectively. As an evaluation measure for the TRT  $\times$  MOD, the root mean squared difference (RMD) between predicted values (calculated based on estimates of fixed effects for TRT  $\times$  MOD) and true values (calculated based on parameters of fixed effects for TRT  $\times$  MOD) was obtained. The RMD is interpreted as the standard deviation of

the differences between predicted and true values. Equations to calculate the level-specific TRT  $\times$  MOD in the MLM-Logistic, GAMM, and MLM are presented in Table 2, and equations to calculate the RMD are presented in Table 3 (top). As a summary of the RMD, the mean of RMDs over 500 replications was reported. For variance and covariance estimates of random intercept and slope ( $\hat{\tau}_{00}, \hat{\tau}_{11}, \hat{\tau}_{01}$ ), and the variance estimate of random residuals ( $\hat{\sigma}^2$ ) in the MLM-Logistic, GAMM, and MLM, the bias (calculated  $\sum_{rep=1}^{500} (\hat{\tau}_{00,rep} - \tau_{00})/500$  where rep denotes a replication number as an example) and the

root mean square error (RMSE; calculated as  $\sqrt{\sum_{rep=1}^{500} (\hat{\tau}_{00,rep} - \tau_{00})^2/500}$  as an example)

was used to evaluate overall accuracy (i.e., bias and variability) patterns with respect to levels of simulation conditions.

To evaluate the accuracy of estimates in the parametric part of the GAMM ( $\hat{\gamma}$ ,  $\hat{\Sigma}$ , and  $\hat{\sigma}^2$ ) in the second focus, the bias<sup>3</sup> and RMSE were calculated and compared across simulation conditions. To evaluate the accuracy of standard errors, the mean standard error of the estimates (M(SE)) across 500 replications was compared with the standard deviation of the estimates (SD) across 500 replications. A ratio of M(SE) to SD close to 1 suggests that the estimated standard errors are approximately correct. Because smooth functions are generated based on basis functions and basis coefficients, the accuracy of basis coefficient estimates ( $\hat{\delta}$ , calculated based on  $\hat{\gamma}$  and  $\hat{\lambda}$ ) was evaluated using the bias and RMSE for the smooth functions. As shown in equation (6), the precision of the smooth functions depends on the standard errors of basis coefficient estimates. Thus, for the smooth functions, the standard errors of the basis coefficient estimates were evaluated using the ratio of M(SE) to SD. In addition, the RMD between predicted values of the level-specific TRT  $\times$  MOD under GAMM or MLM and ‘true’ smooth functions was obtained (see the RMD calculations in Table 3 (bottom)). To summarize the results of the RMD, its mean over 500 replications was obtained. For GAMM, the  $k$ -index for  $K = 10$  used in generating smooth functions was close to 1 for all smooth functions in all conditions and a model with  $K = 10$  was selected among candidate models with different values of  $K$  ( $K = 6, 8, 10, 12, 14$ ) based on the correctedAIC. These results indicate that  $K = 10$  is adequate. In addition, to show the effects of modelling *linear* TRT  $\times$  MOD interactions on MLM parameter estimates in the presence of nonlinear TRT  $\times$  MOD interactions, the bias, RMSE, and the ratio of M(SE) to SD were calculated for MLM estimates which are not part of the linear interaction effects and are comparable with GAMM estimates: the intercept ( $\hat{\gamma}_{00}$ ), the effect of TRT ( $\hat{\gamma}_{02}$ ), the covariance matrix of random effects ( $\hat{\Sigma} = [\hat{\tau}_{00}, \hat{\tau}_{01}, \hat{\tau}_{10}, \hat{\tau}_{11}]'$ ), and the residual variance ( $\hat{\sigma}^2$ ) in equation (2).

### 4.3. Results

The results below are summarized by data-generating models. No convergence problems occurred in any simulation condition for MLM-Logistic, GAMM, or MLM.

<sup>3</sup> Relative percentage bias was not considered because it leads to scaling problems in the case of parameters close to 0 as in our simulation study.



**Table 2.** Simulation study: Comparisons of level-specific TRT  $\times$  MOD among a 'true' model (MLM-Logistic), GAMM, and MLM

	Function form	Level 1	Level 2
MLM-Logistic	Logistic	$\frac{1}{1 + \exp[-\{(\hat{\gamma}_{10} + \hat{\gamma}_{11}z_j)(x_{ij} - x_j)\}]}$	$\frac{1}{1 + \exp[-\{(\hat{\gamma}_{00} + \hat{\gamma}_{02}z_j) + (\hat{\gamma}_{01} + \hat{\gamma}_{03}z_j)x_j\}]}$
GAMM	Nonlinear	$\begin{aligned} & \tilde{f}_2(x_{ij} - x_j)(z_j = 0) + \tilde{f}_2(x_{ij} - x_j)(z_j = 1) \\ &= \sum_{k=1}^9 \hat{\delta}_{2k(z_j=0)} b_{2k}(x_{ij} - x_j) + \sum_{k=1}^9 \hat{\delta}_{2k(z_j=1)} b_{2k}(x_{ij} - x_j) \end{aligned}$	$\begin{aligned} & (\hat{\gamma}_{00} + \hat{\gamma}_{02}z_j) + \tilde{f}_1(x_j)(z_j = 0) + \tilde{f}_1(x_j)(z_j = 1) \\ &= (\hat{\gamma}_{00} + \hat{\gamma}_{02}z_j) + \sum_{k=1}^9 \hat{\delta}_{1k(z_j=0)} b_{1k}(x_j) + \sum_{k=1}^9 \hat{\delta}_{1k(z_j=1)} b_{1k}(x_j) \end{aligned}$
MLM	Linear	$(\hat{\gamma}_{10} + \hat{\gamma}_{11}z_j)(x_{ij} - x_j)$	$(\hat{\gamma}_{00} + \hat{\gamma}_{02}z_j) + (\hat{\gamma}_{01} + \hat{\gamma}_{03}z_j)x_j$

Note.  $\hat{\delta}_{1k(z_j=0)}$  is an estimate of a basis coefficient for  $\tilde{f}_1(x_j)(z_j = 0)$ ;  $\hat{\delta}_{1k(z_j=1)}$  is an estimate of a basis coefficient for  $\tilde{f}_1(x_j)(z_j = 1)$ ;  $\hat{\delta}_{2k(z_j=0)}$  is an estimate of a basis coefficient for  $\tilde{f}_2(x_{ij} - x_j)(z_j = 0)$ ;  $\hat{\delta}_{2k(z_j=1)}$  is an estimate of a basis coefficient for  $\tilde{f}_2(x_{ij} - x_j)(z_j = 1)$ .

**Table 3.** Simulation study: Root mean squared differences (RMDs) between predicted values and true values for MLM-Logistic (top) and GMM (bottom) as data-generating models

Fitting model	Level	RMD
MLM-Logistic	Level 1	$\sqrt{\sum_{j=1}^J \sum_{i=1}^{n_j} d_{ij,ML1}^2 / J n_j}$ where $d_{ij,ML1} = \frac{1}{1 + \exp[-\{(\hat{\gamma}_{10} + \hat{\gamma}_{11}z_j)(x_{ij} - x_j)\}]} - \frac{1}{1 + \exp[-\{(\gamma_{10} + \gamma_{11}z_j)(x_{ij} - x_j)\}]}$
	Level 2	$\sqrt{\sum_{j=1}^J d_{j,ML2}^2 / J}$ where $d_{j,ML2} = \frac{1}{1 + \exp[-\{(\hat{\gamma}_{00} + \hat{\gamma}_{02}z_j) + (\hat{\gamma}_{01} + \hat{\gamma}_{03}z_j)x_j\}]} - \frac{1}{1 + \exp[-\{(\gamma_{00} + \gamma_{02}z_j) + (\gamma_{01} + \gamma_{03}z_j)x_j\}]}$
	Level 1	$\sqrt{\sum_{j=1}^J \sum_{i=1}^{n_j} d_{ij,G1}^2 / J n_j}$ where $d_{ij,G1} = \{f_2(x_{ij} - x_j)(z_j = 0) + \tilde{f}_2(x_{ij} - x_j)(z_j = 1)\} - \frac{1}{1 + \exp[-\{(\gamma_{10} + \gamma_{11}z_j)(x_{ij} - x_j)\}]}$
	Level 2	$\sqrt{\sum_{j=1}^J d_{j,G2}^2 / J}$ where $d_{j,G2} = (\hat{\gamma}_{00} + \hat{\gamma}_{02}z_j) + \{f_1(x_j)(z_j = 0) + \tilde{f}_1(x_j)(z_j = 1)\} - \frac{1}{1 + \exp[-\{(\gamma_{00} + \gamma_{02}z_j) + (\gamma_{01} + \gamma_{03}z_j)x_j\}]}$
MLM	Level 1	$\sqrt{\sum_{j=1}^J \sum_{i=1}^{n_j} d_{ij,M1}^2 / J n_j}$ where $d_{ij,M1} = (\hat{\gamma}_{10} + \hat{\gamma}_{11}z_j)(x_{ij} - x_j) - \frac{1}{1 + \exp[-\{(\gamma_{10} + \gamma_{11}z_j)(x_{ij} - x_j)\}]}$
	Level 2	$\sqrt{\sum_{j=1}^J d_{j,M2}^2 / J}$ where $d_{j,M2} = (\hat{\gamma}_{00} + \hat{\gamma}_{02}z_j) + (\hat{\gamma}_{01} + \hat{\gamma}_{03}z_j)x_j - \frac{1}{1 + \exp[-\{(\gamma_{00} + \gamma_{02}z_j) + (\gamma_{01} + \gamma_{03}z_j)x_j\}]}$
Fitting model	Level	RMD
GMM	Level 1	$\sqrt{\sum_{j=1}^J \sum_{i=1}^{n_j} d_{ij,G10}^2 / J n_j}$ where $d_{ij,G10} = f_2(x_{ij} - x_j)(z_j = 0) - f_2(x_{ij} - x_j)(z_j = 0)$
		$z_j = 1$
	Level 2	$\sqrt{\sum_{j=1}^J \sum_{i=1}^{n_j} d_{ij,G11}^2 / J n_j}$ where $d_{ij,G11} = f_2(x_{ij} - x_j)(z_j = 1) - f_2(x_{ij} - x_j)(z_j = 1)$
		$z_j = 0$
	Level 2	$\sqrt{\sum_{j=1}^J d_{j,G20}^2 / J}$ where $d_{j,G20} = \tilde{f}_1(x_j)(z_j = 0) - f_1(x_j)(z_j = 0)$
		$z_j = 1$
		$\sqrt{\sum_{j=1}^J d_{j,G21}^2 / J}$ where $d_{j,G21} = \tilde{f}_1(x_j)(z_j = 1) - f_1(x_j)(z_j = 1)$

Continued

**Table 3.** (Continued)

Fitting model	Level	Group	RMD
MLM	Level 1	$z_j = 0$	$\sqrt{\sum_{j=1}^J \sum_{i=1}^{n_j} d_{ij,M10}^2 / J n_j}$ where $d_{ij,M10} = \hat{\gamma}_{10}(x_{ij} - x_j) - f_2(x_{ij} - x_j) (z_j = 0)$
		$z_j = 1$	$\sqrt{\sum_{j=1}^J \sum_{i=1}^{n_j} d_{ij,M11}^2 / J n_j}$ where $d_{ij,M11} = \{ \hat{\gamma}_{10}(x_{ij} - x_j) + \hat{\gamma}_{11}(x_{ij} - x_j) \} - f_2(x_{ij} - x_j) (z_j = 1)$
	Level 2	$z_j = 0$	$\sqrt{\sum_{j=1}^J d_{j,M20} / J}$ where $d_{j,M20} = \hat{\gamma}_{01} x_j - f_1(x_j) (z_j = 0)$
		$z_j = 1$	$\sqrt{\sum_{j=1}^J d_{j,M21} / J}$ where $d_{j,M21} = (\hat{\gamma}_{01} x_j + \hat{\gamma}_{03} x_j z_j) - f_1(x_j) (z_j = 1)$

#### 4.3.1. Results for MLM with a logistic function as a data-generating model

Table 4 shows the average RMD across 500 replications for level-specific TRT  $\times$  MOD interactions, and the bias and RMSE of  $\hat{\tau}_{00}$ ,  $\hat{\tau}_{11}$ ,  $\hat{\tau}_{01}$ , and  $\hat{\sigma}^2$  in MLM-Logistic, GMM, and MLM. In Table 4, the averaged results are reported by the levels of simulation conditions in order to understand the main effects of each condition. The results of all 18 simulation conditions are presented in the figures of Appendix S6.

*Performance of GMM compared with MLM-Logistic for Predictions of Level-Specific TRT  $\times$  MOD Interactions.* When the data-generating model is MLM-Logistic, the RMDs in GMM are smaller than the RMDs in MLM-Logistic as an alternative approach to GMM in all simulation conditions (the RMD ranged from 0.042 to 0.126 at level 1 and ranged from 0.106 and 0.281 at level 2 in GMM; it from 0.082 to 0.301 at level 1 and ranged from 0.250 and 0.410 at level 2 in MLM-Logistic). This result indicates that a logistic form of the level-specific TRT  $\times$  MOD interactions can be recovered better using smooth functions in GMM than using nonlinear fitting in MLM. Regarding patterns in the RMD by the levels of simulation conditions, the following is observed. First, in MLM-Logistic, the RMD decreased with increasing number of clusters ( $J$ ) and decreasing cluster size ( $n_j$ ) at level 1 RMD, whereas it decreased with decreasing number of clusters ( $J$ ) and increasing cluster size ( $n_j$ ) at level 2 RMD. In the GMM, the RMD decreased with increasing number of clusters ( $J$ ) and cluster size ( $n_j$ ) at both levels. Second, the RMD decreased with increasing ICCs at level 1 and at level 2 (from ICC = .1 to ICC = .3) in MLM-Logistic, and it decreased with decreasing ICCs at level 2 in GMM; in GMM, ICC had no effect on RMD at level 1.

*Effects of modelling linear TRT  $\times$  MOD on predictions for level-specific TRT  $\times$  MOD interactions.* Under MLM-Logistic as the data-generating model, the RMDs in MLM with linear TRT  $\times$  MOD are larger than the RMDs in MLM-Logistic or in GMM in all simulation conditions (the RMD ranged from 0.499 to 0.597 at level 1 and from 0.512 and 0.663 at level 2 in MLM). In the MLM, there were small differences (to two decimal places) in RMDs across simulation conditions. These results suggest that misspecifying the functional forms for TRT  $\times$  MOD interactions in MLM leads to biased predictions of the interactions in all multilevel designs we considered.

*Random effects comparisons.* Overall, the bias and RMSE of  $[\hat{\tau}_{00}, \hat{\tau}_{11}, \hat{\sigma}^2]'$  were smaller in GMM than in MLM-Logistic. For  $\hat{\tau}_{01}$ , smaller bias was observed in GMM than in MLM-Logistic, while larger RMSE was found in GMM than MLM-Logistic. In addition, the bias and RMSE of  $[\hat{\tau}_{00}, \hat{\tau}_{11}, \hat{\sigma}^2]'$  in MLM were similar to those in GMM because misspecification in MLM is not for these variances but for the level-specific TRT  $\times$  MOD interactions. With respect to simulation conditions, the bias and RMSE of  $[\hat{\tau}_{00}, \hat{\tau}_{11}, \hat{\sigma}^2]'$  tended to decrease as the number of clusters ( $J$ ) and cluster size ( $n_j$ ) increased in all three models, except for a few cases: the bias and RMSE of  $\hat{\tau}_{00}$  in MLM-Logistic regarding  $n_j$ ; the bias of  $\hat{\tau}_{11}$  and  $\hat{\sigma}^2$  regarding  $n_j$  in MLM-Logistic; and the bias of  $\hat{\sigma}^2$  with respect to  $J$  and  $n_j$  in MLM. As noticeable patterns regarding ICCs, the bias and RMSE of  $\hat{\tau}_{00}$  tended to decrease with decreasing ICCs in MLM-Logistic and GMM. And the bias and RMSE of  $\hat{\tau}_{11}$  decreased with increasing ICCs in MLM-Logistic and those of  $\hat{\tau}_{01}$  decreased with increasing ICCs mainly in GMM.

#### 4.3.2. Results for GMM as a data-generating model

The averaged bias, RMSE, and the ratio of M(SE) to SD by levels of simulation conditions are reported in Table 5 for fixed and random effects of GMM and MLM, and in Table 6 for averaged basis coefficients across nine basis coefficients of each smooth function in

**Table 4.** Simulation study: Results for predictions of level-specific TRT  $\times$  MOD interactions (top) and for random effects (bottom) of MLM-Logistic, GAMM, and MLM under MLM-Logistic as a data-generating model

		MLM-Logistic		GAMM	MLM		
Conditions		RMD		RMD	RMD		
Prediction of interactions							
Level 1	$J = 30$	0.205	0.107	0.519			
	$J = 70$	0.170	0.073	0.526			
	$J = 200$	0.162	0.050	0.528			
	$n_j = 15$	0.173	0.089	0.522			
	$n_j = 30$	0.187	0.070	0.527			
	ICC = .05	0.190	0.080	0.531			
	ICC = .1	0.189	0.079	0.526			
	ICC = .3	0.159	0.081	0.517			
	Level 2	$J = 30$	0.312	0.208	0.554		
$J = 70$		0.318	0.161	0.546			
$J = 200$		0.323	0.127	0.554			
$n_j = 15$		0.347	0.180	0.554			
$n_j = 30$		0.286	0.156	0.548			
ICC = .05		0.329	0.150	0.547			
ICC = .1		0.333	0.156	0.574			
ICC = .3		0.290	0.196	0.532			
Random effects							
Random effects	Conditions	MLM-Logistic		GAMM		MLM	
		Bias	RMSE	Bias	RMSE	Bias	RMSE
$\tau_{00}$	$J = 30$	0.131	0.155	-0.018	0.041	-0.007	0.038
	$J = 70$	0.116	0.126	-0.006	0.026	0.002	0.026
	$J = 200$	0.101	0.105	-0.001	0.016	0.006	0.018
	$n_j = 15$	0.109	0.122	-0.011	0.031	0.001	0.029
	$n_j = 30$	0.123	0.134	-0.006	0.025	0.000	0.025
	ICC = .05	0.099	0.104	-0.002	0.011	0.005	0.013
	ICC = .1	0.105	0.112	-0.005	0.018	0.002	0.018
	ICC = .3	0.143	0.169	-0.017	0.054	-0.006	0.051
	$\tau_{11}$	$J = 30$	0.044	0.068	-0.007	0.034	-0.007
$J = 70$		0.045	0.056	-0.001	0.023	-0.001	0.023
$J = 200$		0.042	0.046	0.000	0.014	0.000	0.014
$n_j = 15$		0.043	0.058	-0.004	0.026	-0.003	0.026
$n_j = 30$		0.045	0.056	-0.002	0.021	-0.002	0.021
ICC = .05		0.052	0.063	-0.003	0.024	-0.002	0.024
ICC = .1		0.047	0.059	-0.002	0.024	-0.002	0.024
ICC = .3		0.032	0.048	-0.004	0.024	-0.004	0.024
$\tau_{01}$		$J = 30$	0.074	0.086	0.019	0.348	0.022
	$J = 70$	0.070	0.074	0.004	0.189	0.004	0.174
	$J = 200$	0.063	0.064	0.003	0.105	0.001	0.100
	$n_j = 15$	0.068	0.075	0.009	0.256	0.014	0.216
	$n_j = 30$	0.070	0.075	0.008	0.173	0.004	0.163
	ICC = .05	0.071	0.074	0.023	0.260	0.019	0.214
	ICC = .1	0.070	0.074	-0.003	0.208	0.001	0.188

*Continued*

**Table 4.** (Continued)

Random effects	Conditions	MLM-Logistic		GAMM		MLM	
		Bias	RMSE	Bias	RMSE	Bias	RMSE
$\sigma^2$	ICC = .3	0.067	0.077	0.006	0.174	0.006	0.167
	$J = 30$	0.003	0.036	-0.004	0.036	0.001	0.036
	$J = 70$	0.003	0.023	-0.002	0.023	0.001	0.023
	$J = 200$	0.004	0.015	0.000	0.014	0.003	0.014
	$n_j = 15$	0.003	0.029	-0.003	0.029	0.000	0.029
	$n_j = 30$	0.004	0.020	-0.002	0.020	0.003	0.020
	ICC = .05	0.004	0.025	-0.002	0.024	0.002	0.024
	ICC = .1	0.004	0.024	-0.002	0.023	0.002	0.024
	ICC = .3	0.002	0.025	-0.003	0.025	0.001	0.025

Note. RMD is the mean the root mean squared difference between predicted values and true values across 500 replications.

GAMM. The results of all 18 simulation conditions are presented in the figures of Appendix S6.

*Accuracy of parameter estimates and precision of GAMM.* As shown in the GAMM columns in Table 5 and the figures of Appendix S6, the bias of the intercept estimate ( $\hat{\gamma}_{00}$ ), the TRT estimate ( $\hat{\gamma}_{02}$ ), the variance and covariance estimates of random intercept and slope ( $\hat{\tau}_{00}, \hat{\tau}_{11}, \hat{\tau}_{01}$ ), and the variance estimate of random residuals ( $\hat{\sigma}^2$ ) was close to 0 (ranging from  $-0.086$  to  $0.109$  for all fixed and random estimates in all 18 conditions). Overall, the bias and RMSE of these estimates decreased with increasing number of clusters ( $J$ ) and cluster size ( $n_j$ ). For  $\hat{\gamma}_{00}, \hat{\gamma}_{02}$ , and  $\hat{\tau}_{00}$ , the bias and RMSE of the estimates decreased with smaller ICCs. However, this pattern was not observed for the other parameter estimates. Except for two of the conditions with the smallest number of clusters and smallest cluster size ( $J = 30, n_j = 15, ICC = .05$ ; and  $J = 30, n_j = 15, ICC = .1$ ), the ratios of M(SE) to SD for both  $\hat{\gamma}_{00}$  and  $\hat{\gamma}_{02}$  were close to 1 (ranging from 0.950 to 1.067 and from 0.964 to 1.024 across 16 conditions, respectively). The ratio approached 1 as the number of clusters and cluster size increased and as ICC decreased.

As presented in Table 6, the bias of the average basis coefficient estimate (across nine basis coefficient estimates) for each smooth function was relatively small (ranging from  $-0.012$  to  $0.278$  across 18 conditions and four smooth functions). For all level-1 ( $f_2(x_{ij} - x_j)(z_j = 0); f_2(x_{ij} - x_j)(z_j = 1)$ ) and level-2 ( $f_1(x_j)(z_j = 0); f_1(x_j)(z_j = 1)$ ) smooth functions, the bias and RMSE decreased with increasing the number of clusters and cluster size. However, different patterns were found regarding ICCs at level 1 and level 2. The bias and RMSE decreased with increasing ICCs for the smooth functions at level 1, whereas they decreased with decreasing ICCs for the smooth functions at level 2. Because a larger ICC corresponds with greater between-cluster variability, the pattern at level 2 indicates that the accuracy of basis coefficients for level-2 smooth functions can decrease when there is greater between-cluster variability. The ratio of M(SE) to SD for level-1 smooth functions ranged from 0.900 to 1.003 and the ratio of M(SE) to SD for level-2 smooth functions ranged from 0.920 to 1.033. In addition, as shown in Table 6, the mean RMDs across 500 replications are close to 0 (ranged from 0.037 to 0.326) across all simulation conditions, indicating that the predicted smooth functions are close to the true smooth functions. The RMDs decreased with increasing number of clusters ( $J$ ) and cluster

**Table 5.** Simulation study: Results for fixed and random effects of GAMM ('true' model) and MLM (misspecified model) under GAMM as a data-generating model

Parameters	Conditions	GAMM			MLM		
		Bias	RMSE	Ratio	Bias	RMSE	Ratio
<b>Fixed effects</b>							
$\gamma_{00}$	$J = 30$	-0.004	0.094	1.091	-0.028	0.137	1.721
	$J = 70$	-0.005	0.070	0.980	0.054	0.112	1.633
	$J = 200$	-0.002	0.033	1.009	0.075	0.105	1.612
	$n_j = 15$	-0.005	0.079	1.045	0.022	0.144	1.546
	$n_j = 30$	-0.002	0.053	1.008	0.045	0.092	1.765
	ICC = .05	-0.002	0.048	1.016	0.011	0.114	1.951
	ICC = .1	0.004	0.058	1.056	0.026	0.105	1.687
	ICC = .3	-0.013	0.093	1.009	0.064	0.136	1.328
	$\gamma_{02}$	$J = 30$	0.031	0.113	1.109	-0.011	0.145
$J = 70$		0.014	0.077	0.981	0.024	0.066	1.151
$J = 200$		0.006	0.043	0.979	0.030	0.058	1.250
$n_j = 15$		0.031	0.107	1.054	0.012	0.125	1.381
$n_j = 30$		0.003	0.049	0.992	0.017	0.055	1.173
ICC = .05		0.014	0.063	1.026	0.056	0.091	1.397
ICC = .1		0.016	0.076	1.051	-0.015	0.085	1.272
ICC = .3		0.020	0.095	0.992	0.001	0.092	1.161
<b>Random effects</b>							
$\tau_{00}$	$J = 30$	0.031	0.074		0.218	0.233	
	$J = 70$	0.005	0.034		0.253	0.259	
	$J = 200$	-0.002	0.016		0.257	0.259	
	$n_j = 15$	0.022	0.052		0.209	0.218	
	$n_j = 30$	0.000	0.030		0.276	0.283	
	ICC = .05	0.004	0.018		0.260	0.262	
	ICC = .1	0.018	0.038		0.239	0.243	
	ICC = .3	0.011	0.067		0.229	0.246	
	$\tau_{11}$	$J = 30$	-0.003	0.035		0.001	0.038
$J = 70$		-0.002	0.023		0.005	0.026	
$J = 200$		0.000	0.014		0.006	0.016	
$n_j = 15$		-0.002	0.026		0.003	0.028	
$n_j = 30$		-0.001	0.022		0.005	0.025	
ICC = .05		-0.001	0.023		0.005	0.027	
ICC = .1		-0.002	0.024		0.004	0.026	
ICC = .3		-0.001	0.025		0.004	0.027	
$\tau_{01}$		$J = 30$	-0.011	0.261		-0.094	0.278
	$J = 70$	0.009	0.115		-0.035	0.155	
	$J = 200$	0.002	0.097		-0.058	0.106	
	$n_j = 15$	-0.001	0.164		-0.040	0.185	
	$n_j = 30$	0.000	0.151		-0.084	0.174	
	ICC = .05	0.003	0.167		-0.075	0.184	
	ICC = .1	-0.008	0.159		-0.065	0.179	
	ICC = .3	0.004	0.147		-0.047	0.175	
	$\sigma^2$	$J = 30$	0.008	0.039		0.191	0.198
$J = 70$		0.004	0.025		0.132	0.136	
$J = 200$		0.000	0.014		0.133	0.134	

*Continued*

**Table 5.** (Continued)

Parameters	Conditions	GAMM			MLM		
		Bias	RMSE	Ratio	Bias	RMSE	Ratio
	$n_j = 15$	0.005	0.031		0.111	0.117	
	$n_j = 30$	0.002	0.021		0.194	0.196	
	ICC = .05	0.004	0.026		0.175	0.178	
	ICC = .1	0.004	0.026		0.146	0.151	
	ICC = .3	0.004	0.026		0.135	0.140	

Note. Ratio of M(SE) to SD was considered for fixed effects;  $J$  is the number of clusters;  $n_j$  is the cluster size; ICC is the intraclass correlation coefficient.

size ( $n_j$ ) for all smooth functions. They decreased with decreasing ICC for the level-2 smooth functions, whereas they were not affected by ICC for the level-1 smooth functions.

*Effects of modelling linear TRT × MOD on estimates of MLM and on prediction for level-specific TRT × MOD interactions.* As presented in Table 5, a larger bias, RMSE, and ratio of M(SE) to SD were observed in MLM estimates than in GAMM estimates. First, for all parameter estimates reported in Table 5, the bias decreased with decreasing number of clusters and cluster sizes, except for  $\hat{\tau}_{01}$  and  $\hat{\sigma}^2$  regarding the number of clusters. Second, for all parameter estimates, the RMSE decreased with increasing number of clusters and cluster sizes, with a few exceptions for  $\hat{\tau}_{00}$  with respect to the number of clusters and cluster sizes, and for  $\hat{\sigma}^2$  with respect to cluster sizes. Third, bias tended to be larger with decreasing levels of ICCs for all parameter estimates except  $\hat{\gamma}_{00}$ . Fourth,  $\hat{\tau}_{00}$ ,  $\hat{\tau}_{11}$ , and  $\hat{\tau}_{01}$  from MLM were overestimated. Fifth, the ratio of M(SE) to SD in MLM ranged from 1.210 to 2.241 for  $\hat{\gamma}_{00}$  and from 1.110 to 2.233 for  $\hat{\gamma}_{02}$  across 18 simulation conditions, indicating that standard errors of  $\hat{\gamma}_{00}$  and  $\hat{\gamma}_{02}$  were overestimated. The degree of overestimation of standard errors increased mainly with decreasing number of clusters and ICCs. To conclude, these results suggest that misspecifying the functional forms for TRT × MOD interactions leads to biased estimates of MLM parameters. In addition, as shown in Table 6, larger RMDs were found in MLM than in GAMM across all simulation conditions, indicating that modelling linear TRT × MOD interactions leads to biased predictions of the interactions in the presence of nonlinear interactions.

### 5. Summary and discussion

In this paper we presented a GAMM specification to model a nonlinear multilevel TRT × MOD interaction with unconfounded effects in intervention studies from C-RCT designs. The nonlinear multilevel TRT × MOD interaction was modelled using smooth functions in GAMM. Maximum likelihood estimation was implemented using the `gamm` function in the `mgcv` package. Because the smooth functions are reformulated as random effects in the `gamm` function, it may be challenging for researchers in the social and behavioural sciences to interpret results from the software output. Thus, core derivations from the statistical literature were explained.

The GAMM specification and its estimation were illustrated using instructional intervention data from a C-RCT. We provided the R code to visualize the nonlinear multilevel TRT × MOD interaction with unconfounded effects. Because MLM is a dominant



**Table 6.** Simulation study: Results (bias, RMSE, and ratio) of basis coefficient estimates in GMM, and RMD in GMM and MLM under GMM as a data-generating model

Smooth functions	Conditions	GMM				MLM
		Bias	RMSE	Ratio	RMD	RMD
$f_2(x_{ij} - x_j)(z_j = 0)$	$J = 30$	0.155	0.511	0.938	0.126	0.398
	$J = 70$	0.147	0.397	0.945	0.084	0.361
	$J = 200$	0.032	0.192	0.980	0.045	0.439
	$n_j = 15$	0.154	0.405	0.940	0.110	0.365
	$n_j = 30$	0.069	0.329	0.968	0.071	0.433
	ICC = .05	0.116	0.381	0.964	0.089	0.380
	ICC = .1	0.110	0.369	0.952	0.095	0.410
	ICC = .3	0.108	0.350	0.945	0.095	0.408
	$f_2(x_{ij} - x_j)(z_j = 1)$	$J = 30$	0.177	0.444	0.936	0.130
$J = 70$		0.133	0.348	0.991	0.084	0.297
$J = 200$		0.037	0.207	0.985	0.045	0.359
$n_j = 15$		0.157	0.396	0.966	0.110	0.305
$n_j = 30$		0.074	0.270	0.976	0.077	0.350
ICC = .05		0.123	0.344	0.980	0.095	0.308
ICC = .1		0.115	0.334	0.967	0.095	0.334
ICC = .3		0.109	0.320	0.965	0.095	0.340
$f_1(x_j)(z_j = 0)$		$J = 30$	0.105	0.355	0.953	0.217
	$J = 70$	0.101	0.278	0.967	0.126	0.555
	$J = 200$	0.075	0.236	0.983	0.071	0.633
	$n_j = 15$	0.142	0.389	0.962	0.179	0.589
	$n_j = 30$	0.045	0.191	0.974	0.118	0.654
	ICC = .05	0.067	0.223	0.996	0.095	0.689
	ICC = .1	0.081	0.293	0.960	0.158	0.578
	ICC = .3	0.133	0.353	0.947	0.184	0.597
	$f_1(x_j)(z_j = 1)$	$J = 30$	0.083	0.387	0.977	0.217
$J = 70$		0.051	0.322	0.992	0.138	0.416
$J = 200$		0.003	0.169	0.996	0.077	0.430
$n_j = 15$		0.084	0.380	0.985	0.173	0.464
$n_j = 30$		0.007	0.205	0.992	0.134	0.457
ICC = .05		0.036	0.231	0.998	0.118	0.465
ICC = .1		0.043	0.272	0.990	0.138	0.487
ICC = .3		0.058	0.375	0.978	0.195	0.430

Note. For each smooth function, the bias, RMSE, and ratio reported are averaged across nine basis coefficient estimates; RMD is the mean RMD across 500 replications;  $J$  is the number of clusters;  $n_j$  is the cluster size; ICC is the intraclass correlation coefficient.

analytic method to detect a multilevel TRT  $\times$  MOD interaction in education, the GMM results were contrasted with MLM results. In GMM, TRT effects were different depending on the values of MOD (pretest scores). However, this pattern was obscured when the linear multilevel TRT  $\times$  MOD interaction was modelled in MLM. In addition, simulation studies were implemented to evaluate the performance of GMM in recovering the level-specific logistic (parametric) form of the TRT  $\times$  MOD interaction, compared with MLM-Logistic as an alternative approach and MLM as a misspecification approach. We found that GMM outperformed MLM-Logistic to recover the level-specific logistic form of TRT  $\times$  MOD interaction and MLM led to incorrect prediction of the

interaction. Simulation studies were also conducted to evaluate parameter recovery in GMM and to show consequences of modelling a nonlinear multilevel TRT  $\times$  MOD as a linear multilevel TRT  $\times$  MOD. The parameter recovery in GMM was relatively satisfactory in most multilevel designs typical of educational intervention studies except designs with small number of clusters, small cluster size, and small ICC. When ignoring the nonlinear multilevel TRT  $\times$  MOD interaction, biased estimates such as overestimated standard errors and overestimated variance estimates of random effects were found. These bias patterns were also observed in the empirical study.

The following methodological limitations remain because this paper is the first attempt to apply the GMM to model a nonlinear multilevel TRT  $\times$  MOD interaction with unconfounded effects for educational intervention studies. First, we presented the GMM specification for two groups (control and treatment groups) in a two-level nested design. As in Fuchs et al. (2021), there can be three groups (control, treatment 1, and treatment 2 groups). In this example, the two contrasts (e.g., between control and treatment 1 + treatment 2; and between treatment 1 and treatment 2) can be created for the nonlinear multilevel TRT  $\times$  MOD interaction with unconfounded effects. In addition, there are more complex multilevel designs than the two-level nested design, such as three levels with cross-classified units, for example, student (level 1) nested in a cross-classification of rater and classroom (both level 2) nested in school (level 3). Further studies are needed to apply the GMM to detect a nonlinear multilevel TRT  $\times$  MOD interaction with unconfounded effects for more than two groups and in more complex multilevel designs. Second, the simulation study results are limited to the selected simulation conditions and the selected parameters and nonlinear functions in this study. More extensive simulations that vary these limited conditions should be conducted to make solid generalizations. Third, when newly specified GMMs are presented to researchers in substantive areas, it is important to plan sample sizes to ensure high power for detecting hypothesized magnitudes of ATEs and variability in treatment effects. In a C-RCT design, it is important to have a large number of clusters for inferences about the ATE and to have a large number of clusters and large cluster size for inferences about TRT  $\times$  MOD (Raudenbush & Liu, 2000). Equations for power calculation have been provided for TRT  $\times$  MOD. For example, Raudenbush and Liu (2000) derived a non-central *F*-statistic for the *confounded* fixed effects and variances of random effects of site-level TRT  $\times$  MOD in a multisite randomized trial (MRT) in which individuals are randomly assigned within sites. Dong, Kelcey, and Spybrook (2020) provided power calculation formulas for level-1 TRT  $\times$  binary and continuous MOD in MRTs. Bloom (2005) presented power calculation formulas for TRT  $\times$  binary level-1 or level-2 MOD in two-level C-RCTs. Spybrook, Kelcey, and Dong (2016) provided power calculation formulas for level-2 TRT  $\times$  Level-1 binary MOD and in C-RCTs. Dong, Kelcey, and Spybrook (2018) presented power calculation formulas for level-3 TRT  $\times$  level-1 binary and continuous MOD in C-RCTs. However, existing formulas for power calculation have not been designed for detecting *unconfounded* fixed and variances of random effects for TRT  $\times$  MOD in the C-RCT design. Further studies are needed to provide equations of power calculations to detect such effects.

MLM is frequently used to detect a linear multilevel TRT  $\times$  MOD interaction with confounded effects in educational intervention studies. However, conflation results in insensitivity to theoretically meaningful interactions, and estimates a weighted average of within- and between-cluster effects in the presence of level-specific interaction effects. In addition, a linear interaction is a misspecification in the presence of a more complex nonlinear interaction. The main goal of this study is to illustrate the applicability of the GMM to detect a nonlinear multilevel TRT  $\times$  MOD interaction with unconfounded effects.

We hope that this paper can serve as an example of modelling nonlinear effects using smooth functions in the GAMM for educational intervention research.

### Conflicts of interest

All authors declare no conflict of interest.

### Author contribution

**Sun-Joo Cho:** Conceptualization (equal); Formal analysis (equal); Investigation (equal); Methodology (equal); Project administration (equal); Software (equal); Supervision (equal); Validation (equal); Visualization (equal); Writing – original draft (equal). **Kristopher J. Preacher:** Conceptualization (equal); Investigation (equal); Methodology (equal); Validation (equal); Visualization (equal); Writing – review & editing (equal). **Haley E. Yaremych:** Validation (equal); Writing – review & editing (equal). **Matthew Naveiras:** Validation (equal); Writing – review & editing (equal). **Douglas Fuchs:** Funding acquisition (equal); Resources (equal). **Lynn S. Fuchs:** Funding acquisition (equal); Resources (equal).

### Data availability statement

Data sharing is not available. However, all code is available as supplementary material.

### References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, *19*, 716–723. <https://doi.org/10.1109/TAC.1974.1100705>
- Bloom, H. S. (2005). Randomizing groups to evaluate place-based programs. In H. S. Bloom (Ed.), *Learning more from social experiments: Evolving analytic approaches* (pp. 115–172). New York: Russell Sage Foundation.
- Bringmann, L. F., Hamaker, E. L., Vigo, D. E., Aubert, A., Borsboom, D., & Tuerlinckx, F. (2017). Changing dynamics: Time-varying autoregressive models using generalized additive modeling. *Psychological Methods*, *22*, 409–425. <https://doi.org/10.1037/met0000085>
- Dong, N., Kelcey, B., & Spybrook, J. (2018). Power analyses of moderator effects in three-level cluster randomized trials. *Journal of Experimental Education*, *86*, 489–514. <https://doi.org/10.1080/00220973.2017.1315714>
- Dong, N., Kelcey, B., & Spybrook, J. (2020). Design considerations in multisite randomized trials probing moderated treatment effects. *Journal of Educational and Behavioral Statistics*, *46*, 527–559. <https://doi.org/10.3102/1076998620961492>
- Finch, W. H., & Finch, M. H. (2018). A simulation study evaluating the generalized additive model for assessing intervention effects with small samples. *Journal of Experimental Education*, *86*, 652–670. <https://doi.org/10.1080/00220973.2017.1339010>
- Fox, J.-P. (2010). *Bayesian item response modeling: Theory and applications*. New York: Springer. <https://doi.org/10.1007/978-1-4419-0742-4>
- Fuchs, D., Cho, E., Toste, J. R., Fuchs, L. S., Gilbert, J. K., McMaster, K. L., . . . Thompson, A. (2021). A quasiexperimental evaluation of two versions of first-grade PALS: One with and one without repeated reading. *Exceptional Children*, *87*, 141–162. <https://doi.org/10.1177/0014402920921828>

- Fuchs, L. S., Fuchs, D., Hosp, M. K., & Jenkins, J. R. (2001). Oral reading fluency as an indicator of reading competence: A theoretical, empirical, and historical analysis. *Scientific Studies of Reading*, 5, 239–256. <https://doi.org/10.1207/S1532799XSSR05033>
- Harrell, Jr, F. E. (2015). *Regression modeling strategies: With applications to linear models, logistic and ordinal regression, and survival analysis* (2nd ed.). New York: Springer.
- Lin, X., & Zhang, D. (1999). Inference in generalized additive mixed models by using smoothing splines. *Journal of the Royal Statistical Society, Series B*, 61, 381–400. <https://doi.org/10.1111/1467-9868.00183>
- Lüdtke, O., Marsh, H. W., Robitzsch, A., & Trautwein, U. (2011). A  $2 \times 2$  taxonomy of multilevel latent contextual models: Accuracy-bias trade-offs in full and partial error correction models. *Psychological Methods*, 16, 444–467. <https://doi.org/10.1037/a0024376>
- Marra, G., & Wood, S. N. (2012). Coverage properties of confidence intervals for generalized additive model components. *Scandinavian Journal of Statistics*, 39(1), 53–74. <https://doi.org/10.1111/j.1467-9469.2011.00760.x>
- McNeish, D. M., & Stapleton, L. M. (2016). The effect of small sample size on two-level model estimates: A review and illustration. *Educational Psychology Review*, 28, 295–314. <https://doi.org/10.1007/s10648-014-9287-x>
- Preacher, K. J., Curran, P. J., & Bauer, D. J. (2006). Computational tools for probing interactions in multiple linear regression, multilevel modeling, and latent curve analysis. *Journal of Educational and Behavioral Statistics*, 31, 437–448. <https://doi.org/10.3102/10769986031004437>
- Preacher, K. J., & Sterba, S. K. (2019). Aptitude-by-treatment interactions in research on educational interventions. *Exceptional Children*, 85, 248–264. <https://doi.org/10.1177/0014402918802803>
- Preacher, K. J., Zhang, Z., & Zyphur, M. J. (2011). Alternative methods for assessing mediation in multilevel data: The advantages of multilevel SEM. *Structural Equation Modeling*, 18, 161–182. <https://doi.org/10.1080/10705511.2011.557329>
- Preacher, K. J., Zhang, Z., & Zyphur, M. J. (2016). Multilevel structural equation models for assessing moderation within and across levels of analysis. *Psychological Methods*, 21(2), 189–205. <https://doi.org/10.1037/met0000052>
- R Core Team (2020). *R: A language and environment for statistical computing*. Vienna: R Foundation for Statistical Computing. <https://www.R-project.org/>
- Raudenbush, S. W., & Bryk, A. S. (2002). *Hierarchical linear models: Applications and data analysis methods*. Thousand Oaks, CA: Sage.
- Raudenbush, S. W., & Liu, X. (2000). Statistical power and optimal design for multisite randomized trials. *Psychological Methods*, 5, 199–213. <https://doi.org/10.1037/1082-989x.5.2.199>
- Ruppert, D., Wand, M. P., & Carroll, R. J. (2003). *Semiparametric regression*. New York: Cambridge University Press. <https://doi.org/10.1017/CBO9780511755453>
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6, 461–464.
- Snijders, T. A., & Bosker, R. J. (2012). *Multilevel analysis: An introduction to basic and advanced multilevel modeling*. Thousand Oaks, CA: Sage.
- Spybrook, J., Kelcey, B., & Dong, N. (2016). Power for detecting treatment by moderator effects in two and three-level cluster randomized trials. *Journal of Educational and Behavioral Statistics*, 41, 605–627. <https://doi.org/10.3102/1076998616655442>
- What Works Clearinghouse (2017). *What Works Clearinghouse: Procedures handbook (Version 4.0)*. US Department of Education, Institute of Education Sciences. Retrieved from [https://ies.ed.gov/ncee/wwc/docs/referenceresources/wwc\\_procedures\\_handbook\\_v4.pdf](https://ies.ed.gov/ncee/wwc/docs/referenceresources/wwc_procedures_handbook_v4.pdf)
- Wieling, M. (2018). Analyzing dynamic phonetic data using generalized additive mixed modeling: A tutorial focusing on articulatory differences between L1 and L2 speakers of English. *Journal of Phonetics*, 70, 86–116. <https://doi.org/10.1016/j.wocn.2018.03.002>
- Wood, S. N. (2004). Stable and efficient multiple smoothing parameter estimation for generalized additive models. *Journal of the American Statistical Association*, 99, 673–686. <https://doi.org/10.1198/016214504000000980>

- Wood, S. N. (2006). Low-rank scale-invariant tensor product smooths for generalized additive mixed models. *Biometrics*, 62, 1025–1036. <https://doi.org/10.1111/j.1541-0420.2006.00574.x>
- Wood, S. N. (2013). On  $p$ -values for smooth components of an extended generalized additive model. *Biometrika*, 100(1), 221–229. <https://doi.org/10.1093/biomet/ass048>
- Wood, S. N. (2017). *Generalized additive models: An introduction with R* (2nd. ed.). Boca Raton: Chapman & Hall/CRC.
- Wood, S. N. (2019). Package ‘mgcv’: Mixed GAM computation vehicle with automatic smoothness estimation. Retrieved from <https://cran.r-project.org/web/packages/mgcv/mgcv.pdf>
- Wood, S. N., Pya, N., & Säfken, B. (2016). Smoothing parameter and model selection for general smooth models. *Journal of the American Statistical Association*, 111, 1548–1563. <https://doi.org/10.1080/01621459.2016.1180986>
- Yopp, H. K. (1988). The validity and reliability of phonemic awareness tests. *Reading Research Quarterly*, 23, 159–177. <https://doi.org/10.2307/747800>

Received 26 February 2021

### Supporting Information

The following supporting information may be found in the online edition of the article:

**Appendix S1.** Illustration of the CRS for a smooth function.

**Appendix S2.** Illustration for the reformulation of a smooth function as a random effect.

**Appendix S3.** R code used for the empirical study.

**Appendix S4.** Diagnostic plots of the selected random-intercept GAMM.

**Appendix S5.** Examples of generated functions and estimation code of the simulation study.

**Appendix S6.** Results of the simulation study.