

# Using Complier Average Causal Effect Estimation to Determine the Impacts of the Good Behavior Game Preventive Intervention on Teacher Implementers

Juliette K. Berg<sup>1,4</sup> · Catherine P. Bradshaw<sup>1</sup> · Booil Jo<sup>2</sup> · Nicholas S. Ialongo<sup>3</sup>

Published online: 20 May 2016  
© Springer Science+Business Media New York 2016

**Abstract** Complier average causal effect (CACE) analysis is a causal inference approach that accounts for levels of teacher implementation compliance. In the current study, CACE was used to examine one-year impacts of PAX good behavior game (PAX GBG) and promoting alternative thinking strategies (PATHS) on teacher efficacy and burnout. Teachers in 27 elementary schools were randomized to PAX GBG, an integration of PAX GBG and PATHS, or a control condition. There were positive overall effects on teachers' efficacy beliefs, but high implementing teachers also reported increases in burnout across the school year. The CACE approach may offer new information not captured using a traditional intent-to-treat approach.

**Keywords** Whole school intervention · Implementation · Causal inference · Elementary school · Teacher effects

## Introduction

There is increased focus on the use of universal school-based interventions to promote a range of academic and behavioral outcomes for students and possibly staff;

however, it is common for there to be variation in the extent to which teachers fully comply with the intended implementation model (Domitrovich et al. 2009). Most intervention studies examining the effects of school-based programs have taken an intent-to-treat (ITT) approach, whereby the researchers estimate the effect of being assigned to the treatment condition (Schochet et al. 2014). However, there are typically small or null effects for those participants who do not fully implement, and thus the ITT estimates may understate the effect of intervention (Stuart et al. 2008).

An alternative to traditional ITT analysis is the complier-average causal effect (CACE) analysis approach, which has been successfully used to estimate treatment effects accounting for compliance (see Angrist et al. 1996; Little and Yau 1998; Jo 2002; Stuart et al. 2008). Both ITT and CACE approaches are useful for gaining a more complete understanding of intervention effects (Jo 2002; Stuart et al. 2008). The CACE method has been applied to the estimation of treatment effects with noncompliance within several randomized interventions serving families and youth (e.g., Barnard et al. 2003; Connell et al. 2007; Stanger et al. 2011); however, there has been less focus on classroom-based preventive interventions implemented by teachers. In the current study, we used the CACE method to estimate the impacts of a commonly used classroom-based preventive intervention called the good behavior game (GBG) (Bradshaw et al. 2009b; Ialongo et al. 1999, 2001; Kellam et al. 2008) on teachers' self-efficacy and burnout over the course of a school year. This universal, behavioral management model was combined with a social-emotional learning intervention and implemented in elementary schools. The overall goal of the current study was to provide an example of CACE analysis as applied to a classroom-based intervention. This study represents a novel extension and application of

✉ Juliette K. Berg  
berg.juliette@gmail.com

<sup>1</sup> Curry School of Education, University of Virginia, Charlottesville, VA 22904, USA

<sup>2</sup> Department of Psychiatry and Behavioral Sciences, School of Medicine, Stanford University, Stanford, CA 94305, USA

<sup>3</sup> Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD 21218, USA

<sup>4</sup> American Institutes for Research, 1000 Thomas Jefferson Street, NW, Washington, DC 20007, USA

the CACE analytic approach to better understand the impact of the intervention, which had varying levels of implementation across trained teachers.

### Teachers' Compliance with School-Based Program Implementation

Teachers are often the primary implementers of classroom-based preventive interventions, yet the degree to which they opt to implement the various components of the intervention often varies (Domitrovich et al. 2009). In fact, many of the efficacy and effectiveness studies of school-based prevention programs have noted considerable variation in implementation quality, which in turn attenuates program impacts on student and staff outcomes (Durlak and DuPre 2008; Durlak et al. 2011; Ringeisen et al. 2003). There is emerging evidence that certain characteristics of implementers, including teacher characteristics of attitudes and beliefs about themselves or the school environment, are associated with variation in implementation (Domitrovich et al. 2015; Payne and Eckert 2010). Much of the exploration into the effects associated with variation in implementation has been descriptive and post hoc, with limited use of causal inference approaches and a lack of clarity about the direction of effects. Nevertheless, this line of research suggests that poor implementation and characteristics systematically associated with variation in implementation are typically unmeasured or not accounted for in traditional randomized trials employing an ITT approach. This source of bias may in turn result in an under-estimation of intervention effects. Additional research is needed to demonstrate a causal impact of programs, while taking into consideration program implementation. In fact, the assumption has been that increased compliance with intervention implementation translates into better outcomes (Botvin et al. 1995; Derzon et al. 2005; Durlak and DuPre 2008; Rohrbach et al. 1993). The current proposal focused on implementation, or compliance, which is defined as “the discrepancy between what is planned and what is actually delivered when an intervention is conducted” (Domitrovich et al. 2008, also see Chen 1998; Hulleman and Cordray 2009; O'Donnell 2008). A common indicator of compliance of school-based interventions is program dosage, which includes the frequency with which the program is implemented, with the expectation that a higher dosage is associated with better outcomes.

### Background on the Interventions

The current study tested two evidence-based elementary school prevention programs: the PAX version of the good behavior game (PAX GBG; Embry et al. 2003) and

Promoting Alternative Thinking Strategies (PATHS; Greenberg, Kusché and Conduct Problems Prevention Research Group 2011; Kusché, Greenberg, and Conduct Problems Prevention Research Group 2011). Specifically, a randomized controlled trial (RCT) design was used to compare these two intervention models against a control group. The first intervention model was the PAX GBG alone. The second model was the integration of the PAX GBG and the PATHS program (Domitrovich et al. 2010). PAX GBG focuses on providing teachers with an efficient way of reinforcing the inhibition of aggressive/disruptive and off-task behavior in a “game” like context (Embry et al. 2003). The PATHS curriculum trains teachers to provide explicit instruction to students to promote the development of emotional awareness and communication, self-regulation, social problem solving, and relationship management skills (e.g., interpersonal skills, conflict management) through didactic lessons that take place weekly across the school year (Greenberg and Kusché 2006). Several large RCTs of GBG have demonstrated positive effects on student peer relations, aggressive/off-task behavior, substance use, and academic outcomes (e.g., Bradshaw et al. 2009b; Jalongo et al. 1999, 2001; Kellam et al. 2008). Similarly, prior RCTs of PATHS have yielded positive effects on student social-emotional skills, peer relations, prosocial cognitive functioning, socially-competent behaviors, and behavioral adjustment (e.g., Conduct Problems Prevention Research Group 1999; Greenberg and Kusché 2006; Greenberg et al. 1995).

Although most of the focus has been on student impacts of PAX GBG and PATHS, as well as other whole-school social-emotional programs, teachers implementing these programs may experience benefits such as increased efficacy in managing their classrooms and reduced emotional exhaustion and other forms of burnout (Bradshaw et al. 2008; Bradshaw et al. 2009a; Han and Weiss 2005). On the other hand, the additional burden placed on teachers to implement the program may unintentionally cause some teachers to experience increased burnout and stress. Impacts on teachers, whether positive or negative, may be secondary effects of the program's impacts on students, may stem from teachers' involvement in the training component of the intervention, or may be a function of the supports accompanying the intervention (see Domitrovich et al. 2015 for a more extensive discussion). Further, impacts on students may be a function of how engaged teachers are in implementing the components of the intervention. Greater involvement could result in positive effects as teachers learn to better manage their classrooms, or in negative effects as teachers' burden increases. The effects of such school-based interventions, either positive or negative, on the teachers who implement them likely have important implications for how effective the

interventions are at producing positive effects on students. Yet few studies have specifically tested the impacts of student-focused classroom-based interventions on teacher outcomes (Domitrovich et al. 2016).

## Study Design

Data for this study came from 27 elementary schools in a large urban, east coast public school district. Schools were recruited and principals agreed to participate in a randomized controlled trial of two intervention models and to potentially receive one year of training and coaching. Schools were then randomized (i.e., cluster randomized trial) to one of three conditions: the PAX GBG only (nine schools), the integration of PAX GBG and PATHS (referred to as PATHS to PAX [P2P]) (nine schools), and a control condition (nine schools) where teachers conducted their usual practice. The study took place over the course of one school year. A novel aspect of the design of the current study was the plan to contrast the PAX GBG classroom management model when implemented alone with an integrated training, combining PAX GBG with the PATHS social-emotional learning program. In the current study, we were particularly interested in impacts on teacher outcomes, rather than the traditional impacts solely on students. In fact, as noted above, both PAX GBG and PATHS have the potential to positively impact teacher outcomes of burnout and efficacy, as a function of their positive impact on classroom management and student behavior; however, no studies have taken into consideration compliance when examining these effects. Specifically, our prior analysis of data from this trial using an ITT approach suggested that teachers in the integrated condition reported feeling more efficacious and feeling more personal accomplishment relative to control teachers; however, they did not report reduced levels of emotional exhaustion or depersonalization (Domitrovich et al. 2015).

In the current study, we operationalized implementation compliance as the teachers' use of the PAX GBG "games" in the classroom using records of how many games they played throughout the school year and for how long they played each game. More specifically, compliance was defined as being above a cut point on both the number of games played and the total number of minutes of games played. Based on our prior ITT findings (Domitrovich et al. 2016), we expected to find stronger effects on teacher efficacy and personal accomplishment among intervention teachers who sufficiently complied with the program components because these were the teachers who stood to gain the most from the intervention. Furthermore, we anticipated that these effects would be most pronounced in the integrated PATHS to PAX condition. We also expected

to find intervention effects on emotional exhaustion and depersonalization, although the direction of effects was less clear to us. On the one hand, teachers who were provided the tools to handle behavior management challenges and to improve children's social skills and who felt more efficacious in doing so could in turn experience reductions in emotional exhaustion and depersonalization. However, the potential burden of implementing a new program could put additional strains on teachers and increase burnout (Domitrovich et al. 2010; Han and Weiss 2005), particularly among teachers who spent the most time integrating program components into their daily practice. Thus it is especially important to examine the program impacts on teachers using a CACE analysis, in light of the potential added burden of implementing the multicomponent PATHS to PAX program (Domitrovich et al. 2010).

## Overview of the CACE Approach

The overall goal of the current study was to estimate the effects of the interventions on teachers while accounting for compliance with assigned intervention. In order to do this, we needed to compare outcomes for teachers in the treatment group who implemented the intervention to the outcomes for teachers in the control group who would potentially do the same if assigned to the intervention group. Angrist et al. (1996) provided a framework for this approach, which outlined a process for a two-arm trial with binary compliance in the potential outcomes framework (also see Frangakis and Rubin 2002; Holland 1986). They defined four compliance types on the basis of individuals' treatment assignment status (1 = treatment, 0 = control) and potential treatment receipt status (1 = received/participated, 0 = not received/not participated). These groups are important because we assume that the treatment and control groups are likely to have the same proportion of each compliance type due to the fact that the groups were randomly assigned. Therefore, the difference between the treatment and control condition within each compliance type can be interpreted as a causal effect (Frangakis and Rubin 2002). A nuance of the application of CACE to the current study is that the participants here are the teachers who are in effect *delivering* the intervention, rather than *receiving* it from another source. Specifically, in the current application of the CACE framework, compliers are those who participate in the treatment (i.e., implement) when assigned to the treatment group and do not participate/ implement when assigned to the control group. In an ITT analysis, the effect of treatment assignment is the same as the effect of full participation for the compliers. Always-takers are those who will always implement the treatment, no matter what group they are assigned to. Never-takers are

those who will never implement, regardless of the treatment assignment. In contrast, defiers are those who will not implement if assigned to the treatment group and will implement if assigned to the control group. Similar to Angrist et al. (1996), our primary interest in this study was the causal treatment effect for compliers (i.e., CACE). This focus on implementation, coupled with the use of a continuous compliance indicator for which we set a threshold of high and low compliance (as compared to a traditional categorical approach to compliance) make this application of CACE particularly unique.

### Assumptions to Identify CACE

Given that implementation behavior of each participating teacher can be observed only under the condition the teacher is assigned to, CACE cannot be calculated directly comparing the same teacher's outcomes under the treatment and under the control conditions. Holland (1986) called this the fundamental problem of causal inference. However, under certain conditions, we are able to identify the causal effect at the average level. In particular, the set of conditions (assumptions) used in Angrist et al. (1996) have been widely used to identify CACE. One core assumption is *ignorable treatment assignment*, which provides the basis for causal inference as it guarantees the comparability between treatment arms. In our case, this assumption is automatically satisfied as schools are randomized to intervention and control conditions. A second assumption is the *stable unit treatment value*, which means that the potential outcome of each individual is not affected by the treatment assignment status of other individuals. This is a questionable assumption in school settings because teachers in the same school are highly likely to interact with one another. To minimize this interaction across different treatment arms, we employed cluster randomization, where the unit of randomization is school. Previous studies suggested that by employing cluster randomized trials, interaction or contamination among individuals becomes a more manageable problem (Jo et al. 2008; Sobel 2006). That is, now we only need to worry about interactions among teachers within schools, which can be handled using statistical techniques such as multilevel analysis or generalized estimating equations. We cannot prevent interactions among teachers across different intervention arms, but the likelihood of these interactions will remain about the same as that observed when no nesting exists. A third assumption is *monotonicity*, which assumes that there are no defiers. This is a reasonable assumption in our case because teachers in the control group did not have access to the intervention. It is also assumed that there are at least some compliers, meaning that the offer of the intervention induces at least some teachers to participate. This is a reasonable assumption in our study.

Finally, the *exclusion restriction* assumes that always-takers and never-takers in the control group will not benefit from the program and therefore the distribution of outcomes is the same in the treatment and control groups for these two types. In our context, this means that there is no effect of assignment for never-takers. Since teachers in the control group did not have access to the intervention, the stratum of always-takers does not apply to our study. Given our simplified setting with only compliers and never-takers, we will use non-compliers to refer to never-takers. The assumption of the exclusion restriction may need to be relaxed in school-based interventions where it is quite possible that teachers are affected by the intervention assignment even if they do not participate. For example, teachers may be affected by the training at the beginning of the year even if they do not end up implementing the program in their classrooms according to our definition of implementation. Since compliance in our case is not a dichotomous variable (i.e., teachers can vary in their frequency and quality of program implementation) for whether a teacher participates or not, the cut-off for determining participation will affect the likelihood that the exclusion restriction is met. Given the possible deviation from the exclusion restriction, we additionally conducted CACE estimation assuming an alternative assumption that the intervention effects are additive (Jo et al. 2008), meaning that the intervention effects do not change depending on the values of covariates. In addition, we conducted CACE estimation using two different cut points of our original continuous implementation indicator of program dosage. Comparing the CACE estimates with the exclusion restriction assumption compared to with the additivity assumption, and with two different cut points, served as sensitivity analyses because we cannot ensure that we met the exclusion restriction.

## Method

### Sample

The current study sample included 350 K-5 teachers across 27 schools. Schools, and therefore teachers, were enrolled in three cohorts (i.e., for 1 year each, in three consecutive years) and provided consent for their voluntary participation. The sample was generally evenly split across the three cohorts (31 % cohort 1, 34 % cohort 2, and 35 % cohort 3) and across the three conditions (25 % PAX GBG, 29 % PATHS to PAX, 37 % control). The majority of students in the schools was African American (88 % on average) and received free and reduced meals (i.e., FARMS; 85 %). The vast majority of the teacher sample was female (i.e., 88 %). Less than half was 30 or younger (41.4 %), and taught

**Table 1** Descriptive information on teacher participants and schools

	Total	PATHS to PAX	PAX GBG	Control
Teacher characteristics (%)				
Female	88.7	90.2	88.0	88.3
Taught grades 3–5	41.5	45.1	38.0	42.2
Age $\leq$ 30 years	38.7	39.0	45.4	32.8
Has graduate degree	52.8	52.4	56.5	50.0
Cohort				
Cohort 1	32.1	32.9	37.0	27.3
Cohort 2	34.3	41.5	31.5	32.0
Cohort 3	33.6	25.6	31.5	40.6
Teacher self-report				
	Mean (SD)			
Openness to innovation Time 1	3.75 (0.86)	3.79 (0.86)	3.79 (0.84)	3.70 (0.89)
Mindfulness Time 1	3.98 (0.43)	3.97 (0.37)	3.99 (0.46)	3.99 (0.44)
Behavioral management efficacy				
Time 1	3.84 (0.63)	3.78 (0.58)	3.93 (0.61)	3.82 (0.68)
Time 4	4.03 (0.60)	4.16 (0.55)	4.00 (0.59)	3.96 (0.63)
Social-emotional learning efficacy				
Time 1	3.60 (0.67)	3.51 (0.60)	3.72 (0.64)	3.56 (0.73)
Time 4	3.76 (0.68)	3.92 (0.65)	3.78 (0.63)	3.63 (0.73)
Emotional exhaustion				
Time 1	3.39 (1.39)	3.38 (1.34)	3.39 (1.44)	3.40 (1.40)
Time 4	3.17 (1.48)	3.07 (1.45)	3.14 (1.45)	3.27 (1.54)
Personal accomplishment				
Time 1	5.90 (0.87)	5.77 (0.91)	5.97 (0.82)	5.92 (0.88)
Time 4	5.96 (0.83)	6.09 (0.76)	5.91 (0.79)	5.91 (0.91)
Depersonalization				
Time 1	2.21 (1.32)	2.08 (1.13)	2.35 (1.35)	2.20 (1.42)
Time 4	2.39 (1.38)	2.36 (1.30)	2.39 (1.32)	2.42 (1.49)
School-level variable				
Mobility rate	35.38 (8.20)	37.70 (9.30)	34.04 (7.53)	34.38 (8.14)

students in grades three through five (44.1 %). Just over half of the teachers had a graduate degree (56.4 %). See Table 1 for additional details on the sample as well as average scores on the key measures administered in this study.

## Measures

All outcome measures were assessed using a teacher self-report measure administered four times (i.e., fall baseline and three follow-ups) over the course of the school year.

### Teacher Burnout

Teachers were asked to report on their level of *emotional exhaustion* (nine items, e.g., I feel used up at the end of the workday,  $\alpha = 0.92$ ), *personal accomplishment* (eight

items, e.g., I deal very effectively with the problems of my students,  $\alpha = 0.85$ ), and *depersonalization* (3 items, e.g., I've become more callous towards people since I took this job,  $\alpha = 0.64$ ) from the *Maslach Burnout Inventory* (Maslach et al. 1997). Responses were rated on a 7-point Likert scale from *never* to *every day*, with higher scores indicating greater burnout.

### Teacher Efficacy

Teachers reported on a 5-point scale their self-efficacy in two domains. The *Behavior Management Self-Efficacy Scale* (Main and Hammond 2008) assessed teachers' self-efficacy in promoting classroom behavior management (14 items; e.g., I am able to use a variety of behavior management techniques;  $\alpha = 0.94$ ). The *Social-Emotional*



*Learning Efficacy Scale* (Domitrovich and Poduska 2008) assessed teachers' self-efficacy in promoting social-emotional skills in students (eight items; e.g., I am able to teach children to show empathy and compassion for each other;  $\alpha = 0.93$ ).

### Compliance

Teachers in the intervention groups completed weekly logs in which they recorded the number of games they played and the number of minutes they spent on each game. These were the indicators of compliance. The number of games and the number of minutes played were each summed, for a total score for each measure across the school year. The compliance cut point will affect the exclusion restriction (Stuart et al. 2008); therefore there is a trade-off when deciding where to set the cut point. For example, if the cut point is five games, the assumption is that teachers who led less than five games would not be affected by the intervention. Setting the cut point too low may lead to great variation in the degree to which compliers implemented the program. However, with a higher cut point CACE is assumed to be larger, but the sample size among compliers becomes small thereby reducing the quality of CACE estimates. Therefore, participation, or compliance, was defined in two ways: a *medium compliance* cut point for teachers who fell above the 50th percentile on both the number of games played and the minutes played ( $n = 81$  total treatment teachers), and a *high compliance* cut point was defined as teachers who fell above the 75th percentile on both number of games and minutes played ( $n = 29$  total treatment teachers).

### Covariates

Several baseline variables that were correlated with whether or not treatment teachers were classified as compliers were included as covariates in all models (Domitrovich et al. 2015). A teacher information form was completed at baseline to collect information on teacher demographics (e.g., gender, age, education, degree attained), professional development experiences, and information regarding other social-emotional and classroom management interventions being used by the teacher. Teacher gender, age, graduate degree attainment, the grade level taught, cohort, and school mobility were included in the current study. In addition, several other baseline scales were included as covariates. A total mindfulness score was computed as the mean of 20 items (e.g., When I am in the classroom I have difficulty staying focused on what is happening in the present;  $\alpha = 0.84$ ) from the *Mindfulness in Teaching Scale* (Frank et al. 2016). The *Openness to Innovation* subscale from the *Trust in Schools* measure (Bryk and Schnieder

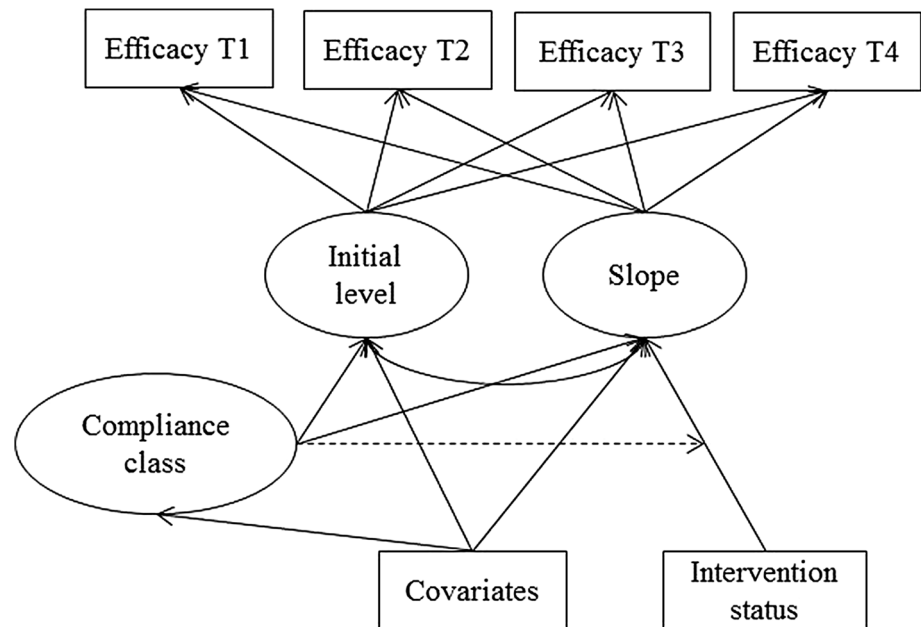
2002) was computed as the mean score of three items (e.g., take responsibility for improving the school;  $\alpha = 0.84$ ). Baseline depersonalization and emotional exhaustion were also included as covariates in all models where they were not the outcome.

### Estimation of CACE

CACE models were estimated separately for each of the treatment conditions relative to control (i.e., PATHS to PAX integrated vs. control and PAX GBG vs. control). Linear growth curve models with intercept and slope parameters were used to estimate the initial level and change of each outcome over the school year (see Fig. 1). In this longitudinal framework, CACE was defined as the effect of intervention assignment for compliers on the change (slope) in each outcome (Jo and Muthén 2003; Jo et al. 2009). As described above, compliance was defined in the current study using two different cut points (50th and 75th percentile), as a sensitivity analysis for both the cut point and the potential deviation from the exclusion restriction. Specifically, the CACE models were identified in two different ways. First, we assumed the exclusion restriction. In this model, the slope was regressed on treatment assignment in the complier class but not in the non-complier class. In this case, we are assuming that non-compliers are not affected by treatment assignment. However, as discussed earlier, this assumption might have been violated in our trial. Given these possibilities of deviation from the exclusion restriction, we additionally conducted CACE estimation assuming that the intervention effects are additive (Jo et al. 2008). That is, we assumed that the intervention effects do not change depending on the covariates. In the model with the additive treatment effect assumption (instead of the exclusion restriction), the slope was regressed on treatment in both the complier and non-complier classes. In both models, the intercept and slope were regressed on the pre-treatment covariates.

Missing data on the compliance measure caused 32 cases (19 in the integrated condition and 13 cases in the PAX GBG condition) to drop out (through listwise deletion), resulting in a total sample size of 318 teachers. An additional set of teachers was excluded in the analysis phase due to missing data on the covariates. Specifically, in the models comparing the PATHS to PAX condition to the control condition, 25 teachers were dropped, resulting in a sample size of 185 teachers. In the models comparing the PAX GBG condition to the control condition, 34 teachers were dropped, resulting in a total sample size of 202. In principle, we could incorporate all cases including the ones with incomplete information. However, we employed listwise deletion, given that there is little research which provides guidelines for handling simultaneous

**Fig. 1** Compliance average causal effect model with covariate



complications of noncompliance, clustering, and missing data. In this study, we focused on handling of noncompliance and clustering. We ignored biases introduced by dropping teachers with missing data, which is a limitation of the study. We used maximum likelihood (ML) estimation with the expectation maximization (EM) algorithm (Little and Rubin 2002) for CACE estimation, conveniently implemented using the mixture modeling feature in Mplus (Muthén & Muthén, 1998–2015). In this framework, compliance status was defined by a categorical latent variable, with one class referring to the compliers and the other class referring to the non-compliers. Given our simplified setting where there are only compliers and never-takers, the compliance class membership was completely observed in the treatment group whereas it was completely unobserved in the control group. The unknown compliance type of individuals in the control condition was handled as missing data via the EM algorithm. Characteristics of teachers were used as predictors of the latent complier class membership (Domitrovich et al. 2015).

In principle, between school and within school level parameters can be formally modeled taking into account compliance in the context of cluster randomized trials (i.e., multilevel modeling). However, in practice, the number of clusters is often small (nine schools per condition in our study). Fairly large numbers of clusters (preferably 50 or more) are necessary to yield accurate CACE estimates when taking a formal multilevel approach (Jo et al. 2008). Instead, we used the sandwich estimator in conjunction with the ML-EM mixture (type = mixture complex) to adjust the

standard errors for the clustering of teachers within schools. In order to interpret the magnitude of effects across the different models, effect sizes (ES) were calculated by dividing the outcome difference across the two conditions by the square root of the total variance obtained from a fully unconditional model. We also applied a Bonferroni correction to the  $p$  value for multiple tests, which set the  $p$ -value for significance at 0.01. Because of the conservative nature of Bonferroni correction, effects between alpha levels of 0.05 and 0.10 were also noted and referred to as trend-level effects.

## Results

Descriptive statistics on the study variables for the two sample conditions are reported in Table 1, whereas Table 2 indicates baseline differences on the covariates between compliers (i.e., high implementer) and non-compliers (i.e., low implementer). Complier teachers in the integrated PATHS to PAX condition were less burnt out at baseline compared to non-compliers, using either the medium or high fidelity/compliance cut point. Complier teachers in the integrated PATHS to PAX condition also had higher mindfulness scores at baseline (high compliance cut point only). Complier teachers in the PAX GBG condition were less likely to have a graduate degree compared to non-compliers. Complier teachers in both conditions were in schools with less mobility regardless of the compliance cut point.

**Table 2** Baseline differences between compliers and non-compliers randomized to the treatment condition

Covariate	With medium compliance cut-off			With high compliance cut-off		
	Non-complier (%)	Complier (%)	Omnibus test ( $\chi^2$ or t-test)	Non-complier (%)	Complier (%)	Omnibus test ( $\chi^2$ or t-test)
<b>PATHS to PAX vs. control</b>						
Female	91.7	88.2	0.27	92.8	76.9	3.11
Late elementary school	48.9	46.9	0.03	46.9	53.8	0.21
New teacher	44.4	37.5	0.37	42.2	38.5	0.06
Graduate degree	52.3	62.5	0.79	58.7	46.2	0.69
Cohort			6.17*			0.70
Cohort 1	43.8	17.6		34.8	23.1	
Cohort 2	35.4	50.0		40.6	46.2	
Cohort 3	20.8	32.4		24.6	30.8	
Mobility (school)	39.1	34.8	2.51*	37.90	34.50	1.46
Mindfulness	3.91 (0.40)	4.05 (0.32)	−1.69	3.92 (0.37)	4.21 (0.28)	−2.63*
Depersonalization	2.34 (1.16)	1.71 (0.99)	2.51*	2.20 (1.18)	1.46 (0.52)	3.59**
Openness to innovation	3.71 (0.85)	3.90 (0.88)	−0.92	3.74 (0.85)	4.05 (0.88)	−1.21
Emotional exhaustion	3.69 (1.33)	2.94 (1.25)	2.46*	3.56 (1.36)	2.52 (0.86)	2.64*
<b>Pax vs. control</b>						
Female	86.9	89.4	0.15	88.0	87.5	0.00
Late E.S.	42.4	34.0	0.77	37.8	43.8	0.20
New teacher	39.7	56.5	2.93	45.5	56.3	0.63
Graduate degree	69.0	45.7	5.75*	59.1	56.3	0.05
Cohort			5.26			7.75*
Cohort 1	45.9	25.5		42.4	6.3	
Cohort 2	29.5	34.0		28.3	50.0	
Cohort 3	24.6	40.4		29.3	43.8	
Mobility (school)	36.2	32.3	3.09**	35.20	30.60	2.57*
Mindfulness	3.98 (0.45)	4.00 (0.447)	−0.18	3.98 (0.46)	4.03 (0.44)	−0.37
Depersonalization	2.43 (1.35)	2.25 (1.37)	0.65	2.38 (1.44)	2.19 (0.78)	0.76
Openness to innovation	3.68 (0.80)	3.91 (0.88)	−1.37	3.79 (0.80)	3.81 (1.04)	−0.12
Emotional exhaustion	3.51 (1.54)	3.26 (1.31)	0.83	3.36 (1.49)	3.55 (1.16)	−0.47

\*  $p < 0.05$ ; \*\*  $p < 0.01$ 

### ITT Estimates

Teachers in the integrated PATHS to PAX condition reported increases in SEL efficacy ( $ES = 0.15$ ,  $p < 0.001$ ) and BM efficacy ( $ES = 0.11$ ,  $p < 0.001$ ) relative to teachers in the control condition. In addition, teachers in the integrated PATHS to PAX condition reported increases in personal accomplishment, one dimension of burnout, relative to teachers in the control condition ( $ES = 0.09$ ,  $p < 0.01$ ). There were no significant impacts on change in depersonalization ( $ES = 0.01$ ,  $p > 0.01$ ) or emotional exhaustion ( $ES = 0.02$ ,  $p > 0.01$ ). The PAX GBG condition did not impact teachers' SEL efficacy ( $ES = 0.04$ ,  $p > .01$ ), BM efficacy ( $ES = 0.02$ ,  $p > 0.01$ ), personal accomplishment ( $ES = 0.03$ ,  $p > 0.01$ ), depersonalization ( $ES = 0.01$ ,  $p > 0.01$ ), or emotional exhaustion

( $ES = 0.01$ ,  $p > 0.01$ ). Results are presented in the right-hand column of Table 3.

### CACE Estimates

Effects of each intervention condition relative to the control condition with compliance are reported in the left-hand columns of Table 3, with the left-most columns reporting the medium compliance cut point with and without the exclusion restriction. Using the medium compliance cut point, complier teachers in the integrated PATHS to PAX condition showed statistically significant increases in SEL efficacy ( $ES = 0.13$ ,  $p < 0.01$ ) and trend-level increases in depersonalization ( $ES = 0.11$ – $0.13$ ,  $p < 0.05$ ) across the school year (with and without the exclusion restriction). Teachers in the PAX GBG condition also showed trend-



**Table 3** CACE effects with covariates

	CACE									ITT		
	With medium compliance cut point											
	With ER			Without ER								
	Compliers			Compliers			Non-compliers					
	Slope	SE	ES	Slope	SE	ES	Slope	SE	ES			
SEL efficacy												
P2P vs. control	0.09	0.04 <sup>†</sup>	0.13	0.09	0.04 <sup>†</sup>	0.13	0.10	0.04 <sup>†</sup>	0.14			
PAX vs. control	0.00	0.04	0.00	-0.02	0.06	0.02	0.07	0.07	0.10			
BM efficacy												
P2P vs. control	0.04	0.04	0.06	0.04	0.04	0.07	0.04	0.06	0.06			
PAX vs. control	0.05	0.75	0.07	0.04	0.03	0.06	-0.04	0.03	0.06			
Depersonalization												
P2P vs. control	0.15	0.07 <sup>†</sup>	0.11	0.16	0.07 <sup>†</sup>	0.12	-0.23	0.12 <sup>†</sup>	0.17			
PAX vs. control	0.16	0.10	0.11	0.19	0.08 <sup>†</sup>	0.13	-0.18	0.12	0.13			
Personal accomplishment												
P2P vs. control	-0.06	0.08	0.06	-0.07	0.06	0.07	0.23	0.06**	0.25			
PAX vs. control	-0.07	0.04	0.07	-0.07	0.04	0.07	0.00	0.06	0.00			
Emotional exhaustion												
P2P vs. control	0.08	0.10	0.06	0.18	0.08 <sup>†</sup>	0.13	-0.44	0.12**	0.32			
PAX vs. control	0.11	0.07	0.07	0.15	0.07 <sup>†</sup>	0.10	-0.32	0.14 <sup>†</sup>	0.22			

	CACE									ITT		
	With high compliance cut point											
	With ER			Without ER								
	Compliers			Compliers			Non-compliers					
	Slope	SE	ES	Slope	SE	ES	Slope	SE	ES			
SEL efficacy												
P2P vs. control	-0.32	0.26	0.48	-0.36	0.28	0.54	0.07	0.02*	0.10	0.10	0.02**	0.15
PAX vs. control	0.12	0.08	0.11	0.02	0.04	0.02	0.01	0.03	0.02	0.03	0.02	0.04
BM efficacy												
P2P vs. control	0.25	0.07**	0.39	0.24	0.07**	0.39	0.02	0.02	0.04	0.07	0.02*	0.11
PAX vs. control	0.21	0.04**	0.32	0.21	0.05**	.32	-0.03	0.02	0.03	0.01	0.02	0.02
Depersonalization												
P2P vs. control	0.16	0.12	0.12	0.16	0.13	0.12	-0.13	0.10	0.09	0.02	0.05	0.01
PAX vs. control	0.08	0.17	0.06	0.12	0.08	0.08	-0.11	0.08	0.07	-0.01	0.05	0.01
Personal accomplishment												
P2P vs. control	-0.11	0.07	0.12	-0.18	0.04**	0.20	0.12	0.04*	0.14	0.08	0.03 <sup>†</sup>	0.09
PAX vs. control	0.17	0.08 <sup>†</sup>	0.19	0.61	0.12**	0.69	-0.07	0.03 <sup>†</sup>	0.08	-0.02	0.03	0.03
Emotional exhaustion												
P2P vs. control	0.36	0.16 <sup>†</sup>	0.26	0.32	0.15 <sup>†</sup>	0.24	-0.25	0.11 <sup>†</sup>	0.18	-0.02	0.05	0.02
PAX vs. control	-0.42	0.44	0.29	-0.42	0.54	0.29	0.01	0.06	0.01	-0.02	0.04	0.01

ES effect size, P2P integrated PATHS to PAX, ER exclusion restriction

\*  $p < 0.01$ ; \*\*  $p < 0.001$ ; <sup>†</sup>  $p < 0.05$

level increases in depersonalization, but only without the exclusion restriction (ES = 0.13,  $p < 0.05$ ). Complier teachers showed trend-level increases in emotional

exhaustion in both conditions, but only without the exclusion restriction (PAX GBG condition ES = 0.10,  $p < 0.05$ ; integrated PATHS to PAX condition ES = 0.13,

$p < 0.05$ ). Without the exclusion restriction, the effects of treatment assignment on the slopes of the outcomes were also estimated for the non-compliers. Non-complier teachers in the integrated PATHS to PAX condition reported decreases in emotional exhaustion ( $ES = 0.32$ ,  $p < 0.001$ ). Non-compliers in the PAX GBG condition reported trend-level decreases in emotional exhaustion (PAX GBG  $ES = 0.22$ ,  $p < 0.05$ ). Non-compliers in the integrated condition reported increases in personal accomplishment ( $ES = 0.25$ ,  $p < 0.001$ ). The effect of being assigned to treatment among non-compliers on these outcomes was stronger than the effect of treatment among compliers. In addition, non-compliers in the integrated condition reported trend-level decreases in depersonalization ( $ES = 0.17$ ,  $p < 0.05$ ), whereas the effect among compliers was positive, indicating a trend towards an increase in depersonalization.

The right-hand columns show results from the models using the high compliance cut point with and without the exclusion restriction. In most instances the results were in the same direction, with the exception of personal accomplishment when comparing PAX GBG to the control condition. Using the high compliance cut point, complier teachers in both conditions showed increases in BM efficacy (PAX GBG  $ES = 0.32$ ,  $p < 0.01$ ; integrated  $ES = 0.39$ ,  $p < 0.001$ ) across the school year with and without the exclusion restriction. Those in the integrated PATHS to PAX condition showed trend-level increases in emotional exhaustion with and without the exclusion restriction ( $ES = 0.24$ – $0.26$ ,  $ps < 0.05$ ). Personal accomplishment increased among high complier teachers in the PAX GBG condition ( $ES = 0.69$ ,  $p < 0.001$ ), and decreased among those in the integrated condition ( $ES = 0.20$ ,  $p < 0.001$ ) without the exclusion restriction. In addition, in contrast to compliers, non-compliers in the integrated condition reported increases in SEL efficacy ( $ES = 0.10$ ,  $p < 0.01$ ) and personal accomplishment ( $ES = 0.14$ ,  $p < 0.01$ ). Finally, results were similar when models were estimated using a sandwich estimator, and were therefore not sensitive to adjusting the standard errors to account for the clustering of teachers in schools.

### Covariate Associations with Compliance

Table 4 shows results from the logistic regression models predicting high compliance from the models without the exclusion restriction. When comparing the integrated condition to the control condition, gender (odds ratio [OR] = 5.38), cohort (OR = 3.94), mobility (OR = 1.15), and mindfulness (OR = 0.03) significantly predicted compliance when personal accomplishment was the outcome (all  $ps < 0.05$ ). Gender (OR = 21.19), cohort (OR = 13.10), mobility (OR = 1.30), and depersonalization (OR = 2.59) predicted compliance when emotional exhaustion was the

outcome (all  $ps < 0.05$ ). Mindfulness predicted compliance with regard to BM efficacy (OR = 0.24), such that compliers were more likely to be higher on mindfulness at baseline. The covariates did not significantly predict SEL efficacy or depersonalization. When comparing the PAX GBG condition to the control condition, mobility predicted compliance for the SEL efficacy outcome (OR = 1.14). Cohort predicted compliance for the BM efficacy (OR = 0.55) and depersonalization (OR = 0.38) outcomes (all  $ps < 0.05$ ).

### Discussion

Many district, school, structural, training, and teacher factors can facilitate or impede the implementation of school-based prevention programs, particularly those that are dependent on teachers' use in the classroom (Domitrovich et al. 2009; Han and Weiss 2005). Comparing average outcomes of schools randomized to an intervention group and a control group should produce unbiased estimates of a program's impacts, assuming that randomization was successful in creating equivalent groups, but the contrasts between the conditions are diminished as a result of variation in treatment received (Weiss et al. 2013). In the current study, average outcomes across treatment conditions should have produced unbiased estimates of the impacts on teacher efficacy and burnout of two evidence-based prevention programs—PAX GBG and PATHS—intended to build students' social-emotional skills, reduce aggressive behaviors, and help teachers manage their classrooms. But teachers varied in the degree to which they implemented the programs, which is a common occurrence in school-based programming (Domitrovich et al. 2008; Fixsen et al. 2005). Ignoring this lower implementation (i.e., compliance) may result in decreased power to detect average effects (Jo et al. 2008). This source of variation likely diminishes the contrasts between treatment conditions and, in turn, attenuates program impacts. CACE estimation is one approach to accounting for this source of variation. In this study, we applied the CACE framework to account for teachers' implementation dosage of the PAX GBG program. The interventions tested within this study are in fact similar in many ways to other classroom-based prevention programs that largely rely on teachers for implementation (e.g., 4R's, Second Step).

Overall, the CACE estimation approach was helpful in understanding treatment-control differences when accounting for variation in treatment conditions due to teacher implementation. This approach revealed impacts on teachers that were different than those produced using an ITT approach. First, some intervention effects on teacher efficacy and burnout were stronger among higher implementing teachers compared to teachers overall. Specifically, we found positive effects on social-emotional and behavioral management

**Table 4** Logistic regression of high compliance on baseline covariates (compliers vs. never-takers)

Covariate	SEL Efficacy		Behavior management efficacy		Depersonalization		Personal accomplishment		Emotional exhaustion	
	Estimate	SE	Estimate	SE	Estimate	SE	Estimate	SE	Estimate	SE
PATHS to PAX vs. Control										
Intercept	4.99	4.802	-10.555	4.396	-4.095	4.27	-4.11	3.988	6.892	5.545
Female	-0.91	0.81	-1.294	0.78	-1.734	1.304	<b>-1.683</b>	<b>0.766</b>	<b>-3.054</b>	<b>1.512</b>
Late elementary school	-0.26	0.609	0.274	0.58	-0.326	0.563	-0.658	0.662	-0.029	0.683
New teacher	0.23	0.676	-0.045	0.534	0.542	0.918	0.32	0.847	0.197	0.676
Graduate degree	-0.39	0.631	-0.049	0.497	-0.711	0.749	-0.023	0.689	-0.901	0.825
Cohort	0.34	0.519	0.676	0.419	-0.543	0.949	<b>-1.371</b>	<b>0.512</b>	<b>-2.569</b>	<b>0.874</b>
Mobility (school)	-0.09	0.052	0.077	0.045	-0.117	0.109	<b>-0.141</b>	<b>0.05</b>	<b>-0.264</b>	<b>0.074</b>
Mindfulness	-0.11	1.04	<b>1.445</b>	<b>0.632</b>	2.561	1.606	<b>3.391</b>	<b>1.004</b>	2.578	1.768
Depersonalization	-0.33	0.426	0.39	0.304	N/A	N/A	-0.397	0.337	<b>-0.95</b>	<b>0.398</b>
Openness to innovation	0.31	0.312	0.001	0.323	0.403	0.442	0.041	0.368	0.213	0.428
Emotional exhaustion	-0.45	0.286	-0.329	0.224	-0.449	0.329	-0.481	0.255	N/A	N/A
PAX vs. control										
Intercept	3.40	3.736	-6.736	3.643	-6.667	3.731	0.509	4.468	-2.754	3.691
Female	-1.21	0.666	0.505	0.86	0.802	0.751	0.251	0.786	-0.075	0.786
Late elementary school	0.60	0.615	-0.086	0.475	0.097	0.468	-0.341	0.524	-0.329	0.651
New teacher	-0.05	0.572	-0.285	0.558	-0.447	0.682	-0.209	0.679	-0.509	0.716
Graduate degree	-0.73	0.575	-0.645	0.614	-1.323	0.722	-0.655	0.635	-0.758	0.86
Cohort	0.04	0.478	<b>0.599</b>	<b>0.290</b>	<b>0.971</b>	<b>0.445</b>	-0.025	0.358	1.009	0.652
Mobility (school)	<b>-0.13</b>	<b>0.044</b>	0.033	0.033	0.015	0.034	-0.045	0.046	0.012	0.033
Mindfulness	-0.39	0.871	0.624	0.628	0.909	0.65	-0.096	0.734	0.031	1.033
Depersonalization	-0.04	0.249	-0.056	0.209	N/A	N/A	-0.254	0.236	-0.271	0.258
Openness to innovation	0.63	0.511	-0.144	0.311	0.182	0.272	-0.097	0.315	-0.092	0.679
Emotional exhaustion	-0.05	0.235	0.264	0.177	-0.291	0.189	0.261	0.232	N/A	N/A

The logistic regression represents the prediction of the covariates on the compliance class. Bolded estimates were statistically significant at  $p < 0.05$ . N/A = covariate was left out in the case where it was the outcome

efficacy among higher implementing teachers in both of the intervention conditions. Teachers on average had greater increases in efficacy in the integrated condition than in the control condition. In the case of behavior management efficacy, these effects seemed to be concentrated among higher implementing teachers. In the case of social-emotional efficacy, the effects among higher implementers and among the teacher sample overall were similar.

The estimation of program impacts on burnout (i.e., personal accomplishment, depersonalization, and emotional exhaustion) while accounting for implementation revealed a different story than the estimation of average impacts across all teachers. Specifically, program effects on personal accomplishment were stronger among teachers most likely to comply in both intervention conditions. On average, there were no significant differences between treatment conditions and the control condition in growth or change of emotional exhaustion or depersonalization across the school year. However, accounting for implementation showed some

evidence for opposing findings among higher and lower implementers and for some increases in burnout among higher implementers. Specifically, being in the higher implementing group in the integrated condition led to greater reports of emotional exhaustion, whereas being in the lower implementing group within an intervention school was associated with somewhat reduced emotional exhaustion. Furthermore, being a higher implementer led to slightly greater reports of depersonalization in both conditions, while being a lower implementer in an integrated school was associated with slightly reduced depersonalization. Overall program impacts on emotional exhaustion and depersonalization may have been masked by these opposing findings. In addition, the trends using the medium and high implementation cut points were somewhat similar, but there were several differences. The effects using the high implementation cut point were notably stronger for all the outcomes except depersonalization. In addition, there were a few cases where the direction of the effect was different (e.g., SEL

efficacy in the integrated condition and emotional exhaustion in the PAX GBG condition).

### Limitations

The initial sample size was small, and the sample of teachers became smaller when split into implementation groups. Estimation of program impacts using the high implementation cut point yielded an especially small sample size and likely rendered the estimates unstable. In addition, despite the fact that schools and not teachers were randomized, the small sample of schools in each condition prevented us from employing a multilevel modeling approach. Multilevel mixture modeling using the EM estimation approach is computationally demanding and treatment effects accounting for compliance are poorly estimated when the number of clusters are small (Jo et al. 2008). Therefore, we were not able to accommodate both implementation and the clustering of teachers in schools.

The interpretation of our findings is limited by our measure of implementation. In most prior applications of the CACE approach, full implementation was known, either based on a theoretical idea of what level of implementation is needed to benefit from the program (e.g., Stuart et al. 2008) or because implementation was defined as electing to receive the intervention or not (e.g., Connell et al. 2007; Cowen 2008). In our case, we did not have a target level of dosage to measure perfect implementation and implementation was measured on a continuum from low to high. As evidenced by the findings, the cut point we used for implementation made a difference in terms of the pattern of findings. Thus, when using the CACE approach to account for service delivery by teachers rather than program uptake of participants, as the approach has most often been used, it would be useful to have a more precise threshold for full implementation. In addition, CACE estimation relies on a set of assumptions, some of which are difficult to meet when applied to school-randomized trials. As discussed earlier, a violation of the exclusion restriction is imaginable in the current study because teachers are likely affected by the intervention even if they did not participate and because our implementation measure is a continuous variable from which we established artificial cut points. The models assuming additivity that relax the exclusion restriction are more likely to suffer from a violation of normality, however. Given these trade-offs, we conducted the CACE estimation with two different cut points and with and without the assumption of the exclusion restriction. We gained more confidence in our findings because our results generally held through sensitivity testing in which we relaxed the exclusion restriction (Jo 2002). However, the results were somewhat sensitive to the cut point used for implementation. This is not

surprising given that implementation was higher and more concentrated using the high implementation cut point rather than the medium implementation cut point. On the other hand, the sample size was significantly reduced using the high implementation cut point. Further tests of violations of the identifying assumptions are not possible. It is important to keep in mind that bias from violation of these assumptions will be problematic in any application of CACE estimation (Jo 2002). Although a strength of this study was the relatively large number of outcomes we examined, it did result in a rather large number of tests conducted. We applied a Bonferroni correction to adjust for the multiple tests; interestingly, many of the findings which dropped to trend-level significance after applying this correction tended to be those which were suggestive of an iatrogenic effect of high implementation. As a result, we are cautious in our interpretation of this finding.

### Conclusion and Implications

Most universal school-based interventions are tested using an ITT approach. Average treatment effects are useful for understanding whether school-based interventions can work under real world conditions. Estimating variation in implementation and impacts can help unpack under what conditions and for whom programs are effective. This can be helpful in targeting interventions and informing design and implementation of evidence-based interventions (Schochet et al. 2014). CACE estimation is one approach for taking into account implementation or compliance when estimating causally estimated treatment impacts. Applied to a case study of two evidence-based interventions implemented in the classroom, we believe that this approach was helpful in uncovering effects that differed from the traditional ITT approach. The findings suggest the possibility that the highest implementing teachers benefited from the interventions in that they felt more efficacious in their instruction across the school year than teachers in the control group. On the other hand, the results raise the possibility that the increased demand put on teachers in the intervention schools may have increased burnout for some teachers over the year. The current findings suggest the possibility that the implementation of an intervention may increase stress and burnout for certain teachers, even as the design intends for the intervention to be integrated into the regular curriculum and seeks to minimize the amount of additional burden placed on the teacher. Another possibility is that certain teachers who are engaging most in the intervention are becoming more aware and learning to recognize their own emotional responses. As a result of this increased emotional awareness, they may be perceiving and reporting greater feelings of burnout. However, we do not know the extent to which the level of burnout reported

may translate into significant or clinical impairment. Regardless, the extent to which these emotional symptoms may be affecting students still needs to be addressed. Given previous research on the negative associations between teacher burnout and student outcomes (Maslach et al. 1997; Pas and Bradshaw 2014), it is possible that any increased burden placed on teachers could attenuate the effects of the interventions on children (for further discussion of the effects of teacher burnout on students, see Abenavoli et al. 2013). Additional research is needed to better understand what potential elements of programs teachers find stressful to implement. With widespread concern for the effective dissemination of evidence-based programs, these data provide some insight into why some teachers may experience resistance to their implementation.

The question also becomes how we can provide the necessary supports for teachers to implement classroom-based interventions without the generation or perception of increased stress. It is possible that some teachers are better equipped to implement these types of programs with high compliance and experience better outcomes than others. Some teachers respond to and benefit from programs without additional supports; other teachers may benefit from implementation supports such as coaching. In this case, program implementation should be somewhat tailored to the implementing teacher, including adapting a program to decrease implementing burden when burden is a concern. This line of research also provides some insight regarding why we might see variation in program impacts. Variation in impacts might be somewhat attributable to variation in teacher response.

Finally, teachers' implementation of social-emotional learning programs involves the development of a similar set of skills among adults. One possibility is that prevention programs that foster social-emotional skill-building in children could also provide the supports and skill-building for teachers to manage and develop their own emotional responses. This is particularly important as greater demands are placed on teachers to incorporate these lessons into their everyday classroom routines. Additional research is needed to further explore the predictive validity of teacher compliance as it relates to student outcomes achieved, over and above the teacher impacts examined in the current study.

**Acknowledgments** This research was supported in part by grants from the Institute of Education [R305A080326; R305A130060] and the National Institute of Mental Health [P30 MH08643]. We also thank Celene Domitrovich for her contribution to the project.

## References

- Abenavoli, R. M., Jennings, P. A., Greenberg, M. T., Harris, A. R., & Katz, D. A. (2013). The protective effects of mindfulness against burnout among educators. *The Psychology of Education Review*, 37, 57–69.
- Angrist, J. D., Imbens, G. W., & Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association*, 91, 444–455.
- Barnard, J., Frangakis, C. E., Hill, J. L., & Rubin, D. (2003). Principal stratification approach to broken randomized experiments. *Journal of American Statistical Association*, 98, 299–323. doi:10.1198/016214503000071.
- Botvin, G. J., Baker, E., Dusenbury, L., Botvin, E. M., & Diaz, T. (1995). Long-term follow-up results of a randomized drug abuse prevention trial in a white middle-class population. *JAMA*, 273, 1106–1112.
- Bradshaw, C. P., Koth, C. W., Bevans, K. B., Ialongo, N., & Leaf, P. J. (2008). The impact of school-wide Positive Behavioral Interventions and Supports (PBIS) on the organizational health of elementary schools. *School Psychology Quarterly*, 23, 462–473. doi:10.1037/a0012883.
- Bradshaw, C. P., Koth, C. W., Thornton, L. A., & Leaf, P. J. (2009a). Altering school climate through school-wide Positive Behavioral Interventions and Supports: Findings from a group-randomized effectiveness trial. *Prevention Science*, 10(2), 100–115. doi:10.1007/s11121-008-0114-9.
- Bradshaw, C. P., Zmuda, J. H., Kellam, S. G., & Ialongo, N. S. (2009b). Longitudinal impact of two universal preventive interventions in first grade on educational outcomes in high school. *Journal of Educational Psychology*, 101(4), 926–937. doi:10.1037/a0016586.
- Bryk, A. S., & Schneider, B. (2002). *Trust in schools: a core resource for improvement*. New York: Russell Sage.
- Chen, H. T. (1998). Theory-driven evaluations. *Advances in Educational Productivity*, 7, 15–34.
- Conduct Problems Prevention Research Group. (1999). Initial impact of the fast track prevention trial for conduct problems: II. Classroom effects. *Journal of Consulting and Clinical Psychology*, 67, 648–657.
- Connell, A. M., Dishion, M. Y., & Kavanagh, K. (2007). An adaptive approach to family intervention: Linking engagement in family-centered intervention to reductions in adolescent behavior. *Journal of Consulting and Clinical Psychology*, 75, 568–579. doi:10.1037/0022-006X.75.4.568.
- Cowen, J. M. (2008). School choice as a latent variable: Estimating the “complier average causal effect” of vouchers in Charlotte. *The Policy Studies Journal*, 36, 301–315.
- Derzon, J. H., Sale, E., Springer, J. F., & Brounstein, P. (2005). Estimating intervention effectiveness: Synthetic projection of field evaluation results. *Journal of Primary Prevention*, 26, 321–343.
- Domitrovich, C. E., Bradshaw, C. P., Berg, J., Pas, E. T., Becker, K., Musci, R., Ialongo, N. (2016). How do school-based prevention programs impact teachers? Findings from a randomized trial of combined classroom management and social-emotional programs. *Journal of Prevention Science*, 1–13.
- Domitrovich, C., Bradshaw, C., Greenberg, M., Embry, D., Poduska, J., & Ialongo, N. (2010). Integrated models of school-based prevention: Logic and theory. *Psychology in the Schools*, 47, 71–88.
- Domitrovich, C. E., Bradshaw, C. P., Poduska, J. M., Hoagwood, K. E., Buckley, J. A., Olin, S., & Ialongo, N. S. (2008). Maximizing the implementation quality of evidence-based preventive interventions in schools: A conceptual framework. *Advances in School Mental Health Promotion*, 1, 6–28.
- Domitrovich, C. E., Gest, S. D., Gill, S., Jones, D. J., & DeRouise, R. S. (2009). Teacher factors related to the professional development process of the Head Start REDI intervention. *Early Education and Development*, 20, 402–430.



- Domitrovich, C., Pas, E., Bradshaw, C. P., Becker, K., Keperling, J. P., Embry, D., & Ialongo, N. (2015). Individual and school organizational factors that influence implementation of the Pax Good Behavior Game intervention. *Prevention Science*, 6(8), 1064–1074.
- Domitrovich, C. & Poduska, J. (2008). The social-emotional learning efficacy scale. *Unpublished technical report*.
- Durlak, J. A., & DuPre, E. P. (2008). Implementation matters: A review of research on the influence of implementation on program outcomes and the factors affecting implementation. *American Journal of Community Psychology*, 41, 327–350.
- Durlak, J. A., Weissberg, R. P., Dymnicki, A. B., Taylor, R. D., & Schellinger, K. B. (2011). The impact of enhancing students' social and emotional learning: A meta-analysis of school-based universal interventions. *Child Development*, 82, 405–432.
- Embry, D., Staatemeier, G., Richardson, C., Lauger, K., & Mitich, J. (2003). *The PAX Good Behavior Game* (1st ed.). Center City : Hazelden.
- Fixsen, D., Naoom, S., Blasé, K., Friedman, R., & Wallace, F. (2005). *Implementation research: A synthesis of the literature*. Tampa: University of South Florida, Louis de la Parte Florida Mental Health Institute. The National Implementation Research Network (FMHI Publication #231).
- Frangakis, C. E., & Rubin, D. B. (2002). Principal stratification in causal inference. *Biometrics*, 58, 21–29.
- Frank, J. L., Jennings, P. A., & Greenberg, M. T. (2016). Validation of the mindfulness in teaching scale. *Mindfulness*, 1, 155–163.
- Greenberg, M., & Kusché, C. (2006). Building social and emotional competence: The PATHS Curriculum. In S. R. Jimerson & M. J. Furlong (Eds.), *Handbook of school violence and school safety: From research to practice* (pp. 395–412). Mahwah: Erlbaum.
- Greenberg, M. T., Kusché, C. A., & Conduct Problems Prevention Research Group. (2011). *Grade level PATHS (Grades 3-5)*. South Deerfield: Channing-Bete Co.
- Greenberg, M., Kusché, C., Cook, E., & Quamma, J. (1995). Promoting emotional competence in school-aged children: The effects of the PATHS curriculum. *Development and Psychopathology*, 7, 117–136.
- Han, S., & Weiss, B. (2005). Sustainability of teacher implementation of school-based mental health programs. *Journal of Abnormal Child Psychology*, 33, 665–679.
- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American Statistical Association*, 81, 945–970.
- Hulleman, C. S., & Cordray, D. S. (2009). Moving from the lab to the field: The role of fidelity and achieved relative intervention strength. *Journal of Research on Educational Effectiveness*, 2, 88–110.
- Ialongo, N., Poduska, J., Werthamer, L., & Kellam, S. (2001). The distal impact of two first-grade preventive interventions on conduct problems and disorder in early adolescence. *Journal of Emotional and Behavioral Disorders*, 9, 146–160.
- Ialongo, N., Werthamer, L., Kellam, S., Brown, C., Wang, S., & Lin, Y. (1999). Proximal impact of two first-grade preventive interventions on the early risk behaviors for later substance abuse, depression, and anti-social behavior. *American Journal of Community Psychology*, 27, 599–641.
- Jo, B. (2002). Estimation of intervention effects with noncompliance: Alternative model specifications. *Journal of Educational and Behavioral Statistics*, 27, 385–409.
- Jo, B., Asparouhov, T., Muthén, B. O., Ialongo, N. S., & Brown, C. H. (2008). Cluster randomized trials with treatment noncompliance. *Psychological Methods*, 13, 1–18.
- Jo, B., & Muthén, B. (2003). Longitudinal studies with intervention and noncompliance: Estimation of causal effects in growth mixture modeling. In N. Duan & S. Reise (Eds.), *Multilevel modeling: Methodological advances, issues, and applications* (pp. 112–139). Mahwah: Lawrence Erlbaum Associates.
- Jo, B., Wang, C. P., & Ialongo, N. S. (2009). Using latent outcome trajectory classes in causal inference. *Statistics Interface*, 2, 403–412.
- Kellam, S., Brown, C. H., Poduska, J., Ialongo, N., Wang, W., Toyinbo, P., & Wilcox, H. (2008). Effects of a universal classroom behavior management program in first and second grades on young adult behavioral, psychiatric, and social outcomes. *Drug and Alcohol Dependence*, 95S, S5–S28.
- Kusché, C. A., Greenberg, M. T., & Conduct Problems Prevention Research Group. (2011). *Grade level PATHS (Grades 1-2)*. South Deerfield: Channing-Bete Co.
- Little, R. J. A., & Rubin, D. B. (2002). *Statistical analysis with missing data*. Hoboken: Wiley.
- Little, R. J. A., & Yau, L. H. Y. (1998). Statistical techniques for analyzing data from prevention trials: Treatment of no-shows using Rubin's causal model. *Psychological Methods*, 3, 147–159.
- Main, S., & Hammond, L. (2008). Best practice or most practiced? Pre-service teachers' beliefs about effective behavior management strategies and reported self-efficacy. *Australian Journal of Teacher Education*, 33, 28–39.
- Maslach, C., Jackson, S. E., & Leiter, M. P. (1997). Maslach burnout inventory. In R. J. Wood (Ed.), *Evaluating stress: A book of resources* (3rd ed., pp. 191–218). Lanham: Scarecrow Education.
- O'Donnell, C. L. (2008). Defining, conceptualizing, and measuring fidelity of implementation and its relationship to outcomes in K-12 curriculum intervention research. *Review of Educational Research*, 78, 33–84.
- Pas, E., & Bradshaw, C. P. (2014). What affects teacher ratings of student behaviors? The potential influence of teachers' perceptions of the school environment and experiences. *Prevention Science*, 15, 940–950. doi:10.1007/s11121-013-0432-4.
- Payne, A. A., & Eckert, R. (2010). The relative importance of provider, program, school, and community predictors of the implementation quality of school-based prevention programs. *Prevention Science*, 11(2), 126–141.
- Ringeisen, H., Henderson, K., & Hoagwood, K. (2003). Context matters: Schools and the “research to practice gap” in children's mental health. *School Psychology Review*, 32(2), 153–169.
- Rohrbach, L. A., Graham, J. W., & Hansen, W. B. (1993). Diffusion of a school-based substance abuse prevention program: Predictors of program implementation. *Preventive Medicine*, 22, 237–260.
- Schochet, P. Z., Puma, M., & Deke, J. (2014). *Understanding variation in treatment effects in education impact evaluations: An overview of quantitative methods* (NCEE 2014-4017). Washington, DC: U.S. Department of Education, Institute of Education Sciences, national Center for Education Evaluation and Regional Assistance, Analytic Technical Assistance and Development. Retrieved from <http://ies.ed.gov/ncee/edlabs>.
- Sobel, M. E. (2006). What do randomized studies of housing mobility demonstrate: Causal inference in the face of interference. *Journal of the American Statistical Association*, 101, 1398–1407.
- Stanger, C., Ryan, S. R., Fu, H., & Budney, A. J. (2011). Parent training plus contingency management for substance abuse families: A Complier Average Causal Effects (CACE) analysis. *Drug and Alcohol Dependence*, 118, 119–126. doi:10.1016/j.drugalcdep.2011.03.007.
- Stuart, E. A., Perry, D. F., Le, H., & Ialongo, N. S. (2008). Estimating intervention effects of prevention programs: Accounting for noncompliance. *Prevention Science*, 9, 288–298. doi:10.1007/s11121-008-0104-y.
- Weiss, M. J., Bloom, H. S., & Brock, T. (2013). *A conceptual framework for studying the sources of variation in program effects*. MDRC. New York: MDRC.