



## LJMU Research Online

Ashworth, E, Humphrey, N and Hennessey, A

**Game Over? No Main or Subgroup Effects of the Good Behavior Game in a Randomized Trial in English Primary Schools**

<http://researchonline.ljmu.ac.uk/id/eprint/11997/>

### Article

**Citation** (please note it is advisable to refer to the publisher's version if you intend to cite from this work)

**Ashworth, E, Humphrey, N and Hennessey, A (2020) Game Over? No Main or Subgroup Effects of the Good Behavior Game in a Randomized Trial in English Primary Schools. Journal of Research on Educational Effectiveness. ISSN 1934-5739**

LJMU has developed **LJMU Research Online** for users to access the research output of the University more effectively. Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. Users may download and/or print one copy of any article(s) in LJMU Research Online to facilitate their private study or for non-commercial research. You may not engage in further distribution of the material or use it for any profit-making activities or any commercial gain.

The version presented here may differ from the published version or from the version of the record. Please see the repository URL above for details on accessing the published version and note that access may require a subscription.

For more information please contact [researchonline@ljmu.ac.uk](mailto:researchonline@ljmu.ac.uk)

<http://researchonline.ljmu.ac.uk/>

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

2

## Abstract

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

1 The current study aimed to examine the impact of a universal, school-based intervention, the  
2 Good Behavior Game (GBG), on children's behavior, and to explore any subgroup moderator  
3 effects among children at varying levels of cumulative risk (CR) exposure. A two-year  
4 cluster-randomized controlled trial was conducted comprising 77 primary schools in England.  
5 Teachers in intervention schools delivered the GBG, while their counterparts in control  
6 schools continued their usual provision. Behavior (specifically disruptive behavior,  
7 concentration problems, and pro-social behavior) was assessed via the checklist version of  
8 the Teacher Observation of Classroom Adaptation (TOCA-C). A CR index was calculated by  
9 summing the number of risk factors to which each child was exposed. Multi-level models  
10 indicated that no main or subgroup effects were evident. These findings were largely  
11 insensitive to the modeling of CR, although a small intervention effect on disruptive behavior  
12 was found when the curvilinear trend was used. Further sensitivity analyses revealed no  
13 apparent influence of the level of program differentiation. In sum, our findings indicate that  
14 the GBG does not improve behavior when implemented in this sample of English schools.

**Keywords**

17 *Good Behavior Game; externalizing behavior; differential effects; universal intervention;*  
18 *cumulative risk*

37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 19 Game over? No Main or Subgroup Effects of the Good Behavior Game in a Randomized  
4  
5 20  
6 Trial in English Primary Schools

7  
8 21 Persistent low-level disruptive behavior (e.g. chatting, calling-out without permission)  
9  
10 22 was highlighted as an area of concern by the Chief Inspector for schools in England (Office  
11  
12 23 for Standards in Education [Ofsted], 2013). Such behavior is proposed to have a detrimental  
13  
14 24 impact on life chances of students, and reduce retention among teachers (Ofsted, 2014).  
15  
16 25 Teachers also share these concerns, with 69% of the 1048 surveyed by Ofsted identifying  
17  
18 26 talking and chatting as a key problem in every lesson. Additional concerns in most lessons  
19  
20 27 included calling-out without permission, failing to follow instructions, and fidgeting. Over 25%  
21  
22 28 reported that the impact of this on learning was high (Ofsted, 2014). While these behaviors  
23  
24 29 are not as severe as the aggressive and anti-social behaviors that characterize conduct  
25  
26 30 disorders (National Institute for Health and Care Excellence, 2013), all behavior problems are  
27  
28 31 likely to impact on the learning, participation and achievement of students, and it is estimated  
29  
30 32 that up to an hour of learning is lost each day as a direct consequence of low-level disruption  
31  
32 33 in classrooms (Ofsted, 2014).

33  
34  
35 34 Korpershoek and colleagues' recent meta-analysis (2016) produced a useful  
36  
37 35 taxonomy of different approaches to classroom management, namely *teachers' behavior-*  
38  
39 36 *focused, teacher-student relationship-focused, students' behavior-focused, and students'*  
40  
41 37 *social-emotional development-focused* interventions. One such *students' behavior-focused*  
42  
43 38 intervention is the Good Behavior Game (GBG; see Intervention section in Method), an  
44  
45 39 interdependent group-contingency behavior management strategy (Lastrapes, 2013) that aims  
46  
47 40 to target low-level disruptive behaviors that interfere with learning, in order to allow more  
48  
49 41 time to teach (Chan, Foxcroft, Smurthwaite, Coomes, & Allen, 2012). It was originally  
50  
51 42 developed by Barrish, Saunders, and Wolf (1969) in the United States of America (USA) and  
52  
53 43 is designed to be used by teachers alongside the curriculum in elementary schools.  
54  
55  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

4

1  
2  
3 44 Evidence spanning several decades across many countries worldwide (e.g. Dion et al.,  
4  
5 45 2011; Leflot, Van Lier, Onghena, & Colpin, 2010; Ruiz-Olivares, Pino, & Herruzo, 2010;  
6  
7 46 van Lier, Muthén, van der Sar, & Crijnen, 2004) attests to the impact of the GBG on  
8  
9 47 behavioral outcomes (Flower, Mckenna, Bunuan, Muething, & Vega, 2014). This research  
10  
11 48 can be broadly categorized into three domains: low-level disruptive behaviors, conduct  
12  
13 49 problems, and aggressive behaviors, with disruptive behavior the outcome most commonly  
14  
15 50 examined (Flower et al., 2014). However, while other positive effects are outlined in the  
16  
17 51 program logic model (Chan et al., 2012), less evidence exists regarding the effectiveness of  
18  
19 52 the GBG in improving these outcomes. For instance, while increased on-task behavior is  
20  
21 53 theorized as an immediate outcome of GBG implementation, few studies have examined this  
22  
23 54 explicitly. One exception found that an adapted version of the GBG in Canada improved  
24  
25 55 students' attentional focus (Dion et al., 2011). The program logic model also outlines  
26  
27 56 increased pro-social behaviors and social awareness as immediate impacts of the GBG, in  
28  
29 57 addition to increased positive peer interactions and reduced anti-social behavior as short- and  
30  
31 58 medium- term impacts (Chan et al., 2012). Indeed, a pilot study of the GBG in Oxfordshire,  
32  
33 59 England found that pro-social behavior improved, along with social competence and  
34  
35 60 decreased social isolation. Qualitative data also supported this, with teachers reporting more  
36  
37 61 effective interpersonal communication among students, and increased sociability (Coombes,  
38  
39 62 Chan, Allen, & Foxcroft, 2016).

40  
41  
42  
43  
44  
45  
46  
47 63 However, a recent meta-analysis (Flower et al., 2014) highlighted some  
48  
49 64 inconsistencies in findings regarding the impact of the GBG on children's behavior. For  
50  
51 65 example, while several studies have reported a positive impact of the GBG on disruptive  
52  
53 66 behaviors (e.g., Barrish et al., 1969; Kleinman & Saigh, 2011; Saigh & Umar, 1983), Leflot  
54  
55 67 and colleagues' (2010) study in Belgium found no significant impact of the program on out-  
56  
57 68 of-seat behaviors. These null findings may be the result of adaptations that were made to the  
58  
59  
60

69 GBG in the Dutch<sup>1</sup> version of the game. Indeed, Coombes and colleagues (2016) emphasize  
70 the importance of adhering to the manualized procedures specified by the program developers.

### 71 **The GBG in England: Issues of Cultural Transferability and Program Differentiation**

72 The inconsistent findings noted above may also be due to cultural incompatibility of  
73 the GBG, as the vast majority of the research has been conducted in the USA (Humphrey et  
74 al., 2016; Lendrum & Humphrey, 2012). Indeed, a recent meta-analysis of school-based  
75 interventions found that larger effect sizes for some outcomes are evidenced when a program  
76 is implemented in its country of origin (Wigelsworth et al., 2016). It is thought that both local  
77 needs and fit with the new cultural context are factors that can influence the success of these  
78 interventions (Castro, Barrera, & Martinez, 2004), and so some aspects of the English school  
79 system may impact on the delivery of the GBG. For example, the National Curriculum and  
80 priorities outlined by Ofsted affect the time and resources that teachers have to successfully  
81 implement optional programs. Furthermore, some teachers have previously noted the  
82 prohibition of teacher-student interaction during GBG gameplay sessions as problematic  
83 (Ashworth, Demkowicz, Lendrum, & Frearson, 2018; Chan et al., 2012). This may reduce the  
84 social validity of the intervention (e.g., its acceptability, feasibility and utility), which in turn  
85 is likely to influence the extent to which teachers in English schools adhere to the guidelines  
86 provided by the intervention developers, and thus the likelihood that implementation will be  
87 high and will be sustained (Wehby, Maggin, Partin, & Robertson, 2011).

88 Conversely, it is possible that the GBG may not be sufficiently distinct from teachers'  
89 usual practice to evidence significant gains in students' outcomes. In the decades since the  
90 intervention was first established, many of the procedures embodied within it have become  
91 standard behavior management practices (e.g., provision of rewards, classroom rules, group-  
92 based contingencies, monitoring behavior). This speaks to the concept of program

---

<sup>1</sup> The GBG was implemented in a Flemish-speaking (Dutch) area of Belgium

1  
2  
3 93 differentiation, defined as, “the extent to which a program’s theory and practices can be  
4  
5 94 distinguished from other programs” (Durlak & DuPre, 2008, p. 329). This aspect of  
6  
7  
8 95 implementation has been sorely neglected in research (e.g., 0/59 studies documented in  
9  
10 96 Durlak and DuPre’s (2008) seminal review), despite its potential importance as a moderator  
11  
12 97 of outcomes. Low levels of program differentiation may be advantageous because a given  
13  
14 98 intervention will feel more familiar to staff and can be assimilated more easily into existing  
15  
16 99 processes and practices; conversely, high levels may be desirable in that the intervention will  
17  
18 100 be seen as more distinctive, adding value to what is already in place (Humphrey, 2013).  
19  
20  
21 101 Either way, it is important to take into account the “uniqueness” of an intervention when  
22  
23 102 evaluating its impact, in order to understand what led to any change in students’ outcomes (or  
24  
25 103 lack thereof; Humphrey et al., 2016).

#### 28 104 **Differential Gains**

30 105 The intention-to-treat (ITT) principle, whereby analyses include every subject who is  
31  
32 106 randomized, regardless of anything that happens after randomization (e.g., noncompliance,  
33  
34 107 protocol deviations, withdrawal), dominates analyses of randomized trials as it provides an  
35  
36 108 unbiased effect estimate (Gupta, 2011). However, it is well established that students do not  
37  
38 109 respond in a uniform manner to universal interventions; natural heterogeneity exists within  
39  
40 110 student populations, and these interventions can differentially affect various strata of the  
41  
42 111 population (Greenberg & Abenavoli, 2017). Indeed, students deemed to be “at-risk” typically  
43  
44 112 evidence greater benefits (Farrell, Henry, & Bettencourt, 2013). For example, an evaluation  
45  
46 113 of Second Step, a universal preventive intervention, found evidence of differential gains  
47  
48 114 among students from socio-economically disadvantaged backgrounds in terms of social  
49  
50 115 competence, school performance and life satisfaction (Holsen, Iversen, & Smith, 2009).

51  
52 116 There has been a strong focus in GBG research on the benefits for certain “at-risk”  
53  
54 117 groups, namely boys and those from low socio-economic backgrounds. Results have  
55  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

7

1  
2  
3 118 consistently shown that these groups of students, and in particular boys who were found to be  
4  
5 119 highly aggressive at baseline, evidenced greater gains from the intervention (Dolan et al.,  
6  
7 120 1993; Kellam et al., 2011; Kellam, Rebok, Ialongo, & Mayer, 1994; Kellam, Ling, Merisca,  
8  
9 121 Brown, & Ialongo, 1998). Studies have also examined the effects of the intervention on  
10  
11 122 individuals displaying “early risk behaviors” for substance abuse, depression and antisocial  
12  
13 123 behavior, reporting positive results (e.g. Ialongo et al., 1999). Indeed, the emphasis on  
14  
15 124 differential gains of the GBG is so strong, some studies *only* report subgroup effects (e.g.,  
16  
17 125 effects reported by gender; Dolan et al., 1993); although this brings into question the overall  
18  
19 126 effectiveness of the intervention at the ITT level.  
20  
21  
22

23  
24 127         However, extant GBG research that utilizes the term “at-risk” generally uses this as a  
25  
26 128 proxy for “highly aggressive” (Dolan et al., 1993). This implies that risk is binary, whereas  
27  
28 129 evidence suggests that risk factors cluster together and are not independent of each other, as  
29  
30 130 proposed in cumulative risk (CR) theory (Rutter, 1979). Studies that have adopted a CR  
31  
32 131 model of risk when examining outcomes for children have typically found that the *number* of  
33  
34 132 risk factors present is a superior predictor of negative outcomes than the *nature* of the  
35  
36 133 individual risks (Evans, Li, & Whipple, 2013). It is thought that measuring the effects of risk  
37  
38 134 factors in isolation – as has been the case in existing studies of the GBG – can over-estimate  
39  
40 135 the importance of a given factor by not accounting for the complex relationships between  
41  
42 136 them (Gerard & Buehler, 1999; Sameroff, Gutman, & Peck, 2003). It is the confluence of  
43  
44 137 various risk factors, rather than any particular factor, that leads to dysfunction (Flouri &  
45  
46 138 Kallis, 2007).  
47  
48  
49

50  
51 139         Thus, research which adopts a CR perspective represents a potentially important step  
52  
53 140 forward in examining potential subgroup effects of preventive interventions, as it more  
54  
55 141 accurately represents individual differences and the multitude of factors influencing a child’s  
56  
57 142 development. However, to the authors’ knowledge, only one study to date has utilized a CR  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

8

1  
2  
3 143 approach when determining the effectiveness of a school-based intervention. An evaluation  
4  
5 144 of the “Guiding Responsibility and Expectations in Adolescents Today and Tomorrow”  
6  
7  
8 145 (GREAT) student curriculum (The Multisite Violence Prevention Project, 2008) found that  
9  
10 146 short- and long-term effects on social-cognitive factors varied as a function of students’ pre-  
11  
12 147 intervention level of risk; while high-risk students evidenced gains in self-efficacy and  
13  
14 148 attitudes towards aggression and non-violent behavior, effects for low-risks students were in  
15  
16 149 the opposite direction. The authors argued that this differential pattern of intervention effects  
17  
18 150 may explain why main effects are not typically found in evaluations of universal  
19  
20 151 interventions in middle schools, thus highlighting the importance of looking beyond the ITT  
21  
22 152 approach.

23  
24  
25  
26 153 It is theorized that the variation in outcomes of universal interventions for different  
27  
28 154 subgroups is due to the extent to which the individuals within them display deficiencies in the  
29  
30 155 skills targeted by the intervention (Farrell, Henry, & Bettencourt, 2013; Greenberg &  
31  
32 156 Abenavoli, 2017). It therefore follows that those children at the higher levels of risk for the  
33  
34 157 intended outcomes of the intervention, and are thus the most in need of it, will evidence the  
35  
36 158 greatest gains. To date, however, “no research has yet examined the effectiveness of the GBG  
37  
38 159 in relation to baseline risk profiles reflecting different constellations of risks across  
39  
40 160 developmental domains” (2013, p. 480), although there is conjecture. For example, Muthén  
41  
42 161 and colleagues (2002) hypothesized: “GBG may have its largest effect for those who are in  
43  
44 162 the middle trajectory class, showing milder forms of problems, while not being strong enough  
45  
46 163 to affect the most seriously aggressive children and not needed for members of the stable  
47  
48 164 non-aggressive group” (p.461).

### 165 **The Current Study**

166 The aforementioned pilot study of the GBG conducted in six English primary schools  
167 for one year found positive effects on various aspects of students’ behavior (Chan et al., 2012;  
168



## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

9

1  
2  
3 168 Coombes et al., 2016). However, this study did not utilize a control group, thus limiting the  
4  
5  
6 169 extent to which these improvements could be securely attributed to the intervention.  
7  
8 170 Therefore, the current study, a cluster-randomized controlled trial (RCT) of the GBG in  
9  
10 171 England, intended to extend this previous research by conducting a more rigorous evaluation  
11  
12 172 utilizing a larger, more representative sample. In addition to examining the impact of the  
13  
14 173 GBG on student behavior (namely disruptive behavior, concentration problems, and pro-  
15  
16 174 social behavior) in English schools, a sensitivity analysis was also conducted to establish  
17  
18 175 whether students' behavioral outcomes varied as a function of levels of program  
19  
20 176 differentiation. Furthermore, pre-specified subgroup analyses were conducted to explore any  
21  
22 177 potential differential effects of the intervention on children at low, medium and high levels of  
23  
24 178 CR exposure. To determine the sensitivity of our findings to changes in the modeling of  
25  
26 179 children's risk status (Evans et al., 2013), additional analyses were conducted, whereby CR  
27  
28 180 exposure (both the linear and quadratic terms, the latter of which represents the curvilinear  
29  
30 181 trend of the continuous CR measure) was treated as a continuous variable.  
31  
32  
33  
34

35 182 Thus, our research questions were as follows:

- 36  
37 183 1. Do children in primary schools implementing the GBG over a two-year period  
38  
39 184 demonstrate significant improvements in behavior (specifically, a) disruptive behavior, b)  
40  
41 185 concentration problems and c) pro-social behavior) compared to those children attending  
42  
43 186 usual practice schools?  
44  
45  
46 187 2. Are any findings in relation to RQ1 sensitive to varying levels of program differentiation?  
47  
48 188 3. Are there differential intervention gains in behavioral outcomes among children at  
49  
50 189 different levels of cumulative risk exposure?  
51  
52 190 4. Are any findings in relation to RQ3 sensitive to changes in the way risk exposure is  
53  
54 191 modeled?  
55  
56  
57

## 192 Method

### 193 Design

1  
2  
3 194 The current study utilizes data from an efficacy trial of the GBG in England and so  
4  
5 195 methods have been published previously. In brief, a two-year cluster-randomized design was  
6  
7 196 utilized, with participating schools as the unit of randomization. A local trials unit randomly  
8  
9 197 allocated schools to one of two trial arms: (1) GBG (intervention arm); or (2) usual provision  
10  
11 198 (UP arm). To ensure balance across the arms of the trial, a minimization algorithm (adaptive  
12  
13 199 stratification) was applied regarding the proportion of children eligible for free school meals  
14  
15 200 (FSM) and school size. Teachers in schools allocated to the intervention arm were trained and  
16  
17 201 supported to implement the GBG during the two-year trial period (2015/16 and 2016/17). The  
18  
19 202 trial protocol is available here [masked for peer review].  
20  
21  
22  
23

24 203 Schools were recruited between March and July 2015. Eligible schools were  
25  
26 204 mainstream, state-maintained primary schools (serving children aged four-11 years).  
27  
28 205 Participation required informed consent from the schools' Head Teachers. Child assent and  
29  
30 206 parental opt-out consent were also sought. In total, 68 parents (2.2%) exercised their right to  
31  
32 207 opt their children out of the trial, and no children declined assent or exercised their right to  
33  
34 208 withdraw from the study. The study received approval from the ethics committee of the  
35  
36 209 authors' host institution.  
37  
38  
39

#### 40 210 **Participants**

41  
42 211 The trial sample were N=3084 children aged six-seven in 77 schools in three regions  
43  
44 212 across England (see supplemental files for CONSORT diagram). The composition of  
45  
46 213 participating schools mirrored that of primary schools in England regarding size and the  
47  
48 214 proportion of students speaking English as an Additional Language (EAL). However, trial  
49  
50 215 schools typically had significantly larger proportions of children with special educational  
51  
52 216 needs and disabilities (SEND) and those eligible for FSM, in addition to lower rates of  
53  
54 217 absence and attainment. At the student level, the trial sample were also generally above the  
55  
56 218 national average regarding the proportion who were identified as having an SEND, eligible  
57  
58  
59  
60

219 for FSM, and speaking EAL; while they were typically below average with regards to  
220 attainment (DfE, 2015; Table 1). Differences between schools and students across the trial  
221 arms were negligible, indicating good balance and successful randomization.

### 222 **Sample Size**

223 Sample size calculations were carried out using Optimal Design Software at the point  
224 of randomization. With an intra-cluster correlation co-efficient (ICC) of 0.08 for the primary  
225 outcome measure (disruptive behavior) at baseline, a pre-post correlation of 0.63, an average  
226 cluster size of 40, and standard Power and Alpha thresholds of 0.80 and 0.05 respectively, the  
227 minimum detectable effect size (MDES) for an ITT analysis was determined to be 0.16.

228 Given this, the trial was considered to be well powered.

### 229 **Intervention**

230 The GBG is an “interdependent group-oriented contingency management procedure”  
231 (Tingstrom, Sterling-Turner, & Wilczynski, 2006, p. 225) designed to be integrated into the  
232 existing curriculum without taking up any additional teaching time. It is underpinned by  
233 behaviorism (i.e., reinforcement of desired behaviors; Skinner, 1948), social learning theory  
234 (i.e., vicarious learning through the modeling of appropriate behaviors; Bandura, 1977), and  
235 life course/social field theory (i.e., an individuals’ ability to meet the social demands of a  
236 particular environment; Kellam, Branch, Agrawal, & Ensminger, 1975). Core components  
237 are (1) *classroom rules*, (2) *team membership*, (3) *monitoring behavior*, and (4) *positive*  
238 *reinforcement*. It is suggested that the GBG should initially be implemented three times a  
239 week, for ten minutes each time, with this increasing to everyday for up to 30 minutes over  
240 the course of the year. It should also be played at varying points in the day, during an  
241 assortment of lessons and activities. The logic model for the program is available in Chan et  
242 al. (2012).

1  
2  
3 243           When playing the game, students are divided into teams of up to seven that are  
4  
5 244           gender-balanced and heterogeneous in behavior and academic ability. Teams are expected to  
6  
7 245           follow four rules during the game: (1) *we will work quietly*<sup>2</sup>, (2) *we will be polite to others*, (3)  
8  
9 246           *we will get out of our seats with permission*, and (4) *we will follow directions* (Kellam et al.,  
10  
11 247           2011); teachers monitor behavior and record any infractions that occur as a result of a team  
12  
13 248           member failing to follow these. Other than when recording an infraction, teachers should not  
14  
15 249           interact with students during the game. In order to win the GBG, and thus access agreed  
16  
17 250           rewards or privileges, teams need to have four or fewer infractions at the end of the game. At  
18  
19 251           the beginning of the year, it is recommended that these rewards are tangible (e.g., stickers)  
20  
21 252           and given immediately after the game ends; as the school year progresses, the rewards should  
22  
23 253           become intangible (e.g., free time) and their receipt should be delayed (e.g., end of day).

24  
25  
26  
27  
28 254           As the trial took place over a two-year period, different teachers delivered the GBG in  
29  
30 255           the first and second years. All teachers attended two days of training prior to implementation,  
31  
32 256           in the September or October of their delivery year, with a further day of top-up training a few  
33  
34 257           months later. Training focused on the theoretical underpinnings of the GBG and procedures  
35  
36 258           for implementation. Trained GBG coaches visited schools approximately once per month  
37  
38 259           throughout the trial to support teachers' implementation (e.g. observation and feedback,  
39  
40 260           modeling delivery; Ashworth et al., 2018). Coaches had all previously worked as teachers or  
41  
42 261           education professionals and were trained by the program developers.

## 43 262   **Implementation**

44  
45  
46  
47 263           Implementation fidelity/quality, participant responsiveness and reach were assessed  
48  
49 264           via annual structured observations in the second term of each year of the trial (January-April).  
50  
51 265           These were developed and piloted using video footage of GBG delivery in English schools  
52  
53 266           recorded in the aforementioned UK pilot (Chan et al., 2012). Inter-rater reliability was found

---

54  
55  
56  
57  
58  
59  
60  
<sup>2</sup> Adherence to “quietly” is defined as working at a “voice level” set by the teacher that is deemed to be appropriate for a particular activity.

267 to be very high. Dosage was measured via an online scoreboard developed for teachers to  
268 record infractions during the game; this also included details of length and frequency of play.

269 Average scores for fidelity/quality (2015/16: 69.79%; 2016/17: 70.11%) were high,  
270 suggesting that teachers followed most of the prescribed steps in the manual and did so in an  
271 enthusiastic and engaging manner. Participant responsiveness was also high (2015/16:  
272 74.51%; 2016/17: 69.07%), as was participant reach (2015/16: 95.26%; 2016/17: 95.98%),  
273 indicating that students responded favorably to the GBG and were largely present for delivery.  
274 The intervention was implemented on average twice per week in the first year of the trial  
275 (2015/16: 1.93 games per week), although this reduced slightly in the second year (2016/17:  
276 1.55 games per week); the average game lasted approximately 15 minutes in both years. Thus,  
277 while average *duration* was well within the range of the only other GBG trials that have  
278 reported dosage data, the *frequency* of game play was somewhat lower (Dimitrovich et al.,  
279 2015; Hagermoser Sanetti & Fallon, 2011; Kellam et al., 1998; Pas et al., 2015).

## 280 **Measures**

281 **Behavior.** Behavior was assessed using the checklist version of the Teacher  
282 Observation of Classroom Adaptation (TOCA-C; Koth, Bradshaw, & Leaf, 2009). This 21-  
283 item scale assesses students' concentration problems (inattentive and off-task behavior; seven  
284 items), disruptive behavior (disobedient, disruptive and aggressive behaviors; nine items) and  
285 pro-social behavior (positive social interactions; five items). Statements are provided about  
286 the child (e.g. gets angry when provoked by other children), which teachers read and endorse  
287 on a six-point scale (Never/Rarely/Sometimes/Often/Very Often/Almost Always). Children's  
288 scores are then summed and averaged (one-six), with higher scores indicating more  
289 maladaptive behaviors for concentration and disruptive behavior, and lower scores indicative  
290 of poorer pro-social behavior (Kourkounasiou & Skordilis, 2014).

291 The TOCA has been frequently used in previous research on the GBG (e.g. Bradshaw

et al., 2015; Chan et al., 2012; Kellam et al., 1994). The checklist version has good psychometric properties, including high internal consistency (all subscales  $\alpha > 0.86$ ), and a factor structure that is invariant across gender, race and age (Koth et al., 2009; Bradshaw, Waasdorpe, & Leaf, 2015). Internal consistency of the TOCA-C subscales in the trial was excellent (all  $\alpha > 0.87$  at baseline).

**Program differentiation.** As a means through which to determine the level of program differentiation of the GBG, usual practice surveys (based on existing measures of classroom management strategies; Reupert & Woodcock, 2010) were administered to all teachers at baseline. Twelve items deemed to reflect the presence or absence of key GBG procedures and practices (e.g. establishing and maintaining a set of classroom rules; observing and monitoring students' behavior in the classroom; use of prizes as rewards for good behavior; use of group rewards; use of a warning/strike system) were extracted to create a program differentiation index (PDI), with scores ranging from 0-12; higher scores were indicative of lower levels of program differentiation. PDI scores were transformed to percentages for ease of interpretation and aggregated to the school-level. Finally, the school-level PDI score was converted to a binary variable, with schools categorized into either low or moderate<sup>3</sup> program differentiation (utilizing the 50<sup>th</sup> percentile as a cut-point). 38 schools (17 GBG; 21 UP) were designated as moderate PDI and 38 (21 GBG; 17 UP) as low PDI. The remaining school's PDI score could not be calculated as they failed to complete usual practice surveys.

**Cumulative risk.** Previous research by the authors utilizing the same dataset identified the student- and school-level risk factors that were significant predictors of baseline disruptive behavior scores. These analyses were subsequently extended for the present study to incorporate concentration problems and pro-social behavior (see Table 2). As CR theory

---

<sup>3</sup> Given the distribution of scores, which ranged from 56-92%, 'moderate' and 'low PDI was deemed to be more accurate than 'high' and 'low'.

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

15

1  
2  
3 316 states that the number of risk factors is more important than their nature, these risk factors  
4  
5 317 were then dichotomized (coded as either '0' for absent or '1' for present) and summed for  
6  
7 318 each of the behavioral outcomes, creating three CR scores for each child that represented the  
8  
9 319 number of risk factors to which they were exposed. This is consistent with previous work in  
10  
11 320 the field (e.g., Ashworth & Humphrey, 2018; Gerard & Buehler, 2004; Hebron, Oldfield, &  
12  
13 321 Humphrey, 2016; Oldfield, Humphrey, & Hebron, 2015). Calculations of effect size were  
14  
15 322 conducted to determine where differences in mean behavior scores between risk levels lay;  
16  
17 323 line graphs were also plotted to provide a visual representation of the risk-outcome  
18  
19 324 relationship (see supplemental files). Risk groups were then created in accordance with the  
20  
21 325 elbow points present on these graphs, and where differences between risk levels were notable.  
22  
23 326 Children were therefore categorized into one of three groups for each of the three measures of  
24  
25 327 behavior: *low-risk*, *medium-risk*, or *high-risk* (see Table 2).

**328 Analysis**

329 ITT analyses (controlling for school-level FSM and school size, and child-level  
330 gender, FSM and SEND status) were conducted and subsequently extended to incorporate an  
331 analysis of subgroup moderator effects. Multi-level modeling (in MLwiN 2.36) was used to  
332 account for the clustered and hierarchical nature of the data (students nested within schools;  
333 Twisk, 2006). A fixed effects random intercepts model was utilized, which assumes that  
334 baseline scores will have different explanatory power for different schools, and that the  
335 relationship between baseline and follow-up will not vary by school. Prior to analysis,  
336 behavior scores were standardized by converting them to z scores, meaning that the  
337 coefficients reported can be interpreted as effect sizes akin to *Cohen's d*, thus facilitating  
338 interpretation across models (Bierman et al., 2014).

339 First, two-level models were fitted (one for each measure of behavior) with  
340 intervention group allocation at the school-level (with *UP allocation* utilized as the reference

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

16

1  
2  
3 341 category) and the relevant baseline behavior scores and student-level covariates (gender,  
4  
5 342 SEND, FSM) at the student-level as explanatory variables, to establish any main effects of  
6  
7 343 the GBG on the outcome variables (post-test behavior scores). Second, PDI group was added  
8  
9 344 at the school-level as an explanatory variable, and interaction terms between trial group and  
10  
11 345 PDI (with *low PDI* utilized as the reference category) were specified using dummy coding  
12  
13 346 (e.g. 0 = low, 1 = moderate) to establish whether our substantive findings were sensitive to  
14  
15 347 levels of program differentiation.  
16  
17  
18

19 348 For research question 2, risk group categorization was added to the original model at  
20  
21 349 the student-level. Cross-level interaction terms between the intervention group and the three  
22  
23 350 risk groups (with the *no risk* group utilized as the reference category) were specified using  
24  
25 351 dummy coding, to establish any subgroup moderator effects. Finally, to determine the  
26  
27 352 sensitivity of our substantive findings to changes to the modeling of children's risk status,  
28  
29 353 two additional multi-level models were fitted. The original form of the continuous CR score  
30  
31 354 was added to the first, along with an interaction term between CR score and intervention  
32  
33 355 group. For the second, consistent with previous literature (Ashworth & Humphrey, 2018;  
34  
35 356 Oldfield et al., 2015), the continuous CR score was mean-centered and squared, to generate a  
36  
37 357 quadratic CR score, and was then added to the model, along with an interaction term between  
38  
39 358 quadratic CR score and intervention group.  
40  
41  
42  
43

44 359 18.5% of participants in the sample had incomplete data, in cases where they had left  
45  
46 360 the school (12.6%) or teachers had failed to provide post-test behavior data (5.9%). Missing  
47  
48 361 value analysis was conducted through binary logistic regression to identify the variables that  
49  
50 362 predicted partially observed data. Missingness was predicted by school size, school-level  
51  
52 363 absence, school-level behaviour, student SEND status, and student looked-after status. Thus,  
53  
54 364 data were likely to be missing at random. Therefore, in order to maintain the sample size,  
55  
56 365 multiple imputation (MI) procedures were implemented. This reduces the bias associated  
57  
58  
59  
60



with attrition and allows for the use of statistical techniques designed for complete datasets (Pampaka, Hutcheson, & Williams, 2016). MI was conducted in REALCOM-Impute with demographic variables and trial group allocation added as auxiliary (where data were fully observed) and response variables. REALCOM-Impute default settings of 10 datasets, 1000 iterations, a burn-in of 100, and a refresh of 10 were utilized, in accordance with guidance produced by Carpenter and colleagues (2011) for multi-level imputation with mixed response types.

## Results

Descriptive statistics pertaining to both the main and subgroup analyses are presented in Table 3.

### Research Question 1

ITT analyses (Table 4) indicated that the GBG had no overall effect on children's a) disruptive behavior ( $d = 0.056$ ,  $p = .235$ ), b) concentration problems ( $d = 0.022$ ,  $p = .400$ ) or c) pro-social behavior ( $d = -0.108$ ,  $p = .160$ ).

### Research Question 2

There was no significant interaction found between moderate levels of program differentiation and GBG trial group allocation for a) disruptive behavior ( $d = 0.019$ ,  $p = .451$ ), b) concentration problems ( $d = 0.053$ ,  $p = .378$ ), or c) pro-social behavior ( $d = -0.282$ ,  $p = .097$ ) (see Table 5).

### Research Question 3

There were no statistically significant subgroup effects of the GBG on a) disruptive behavior (medium-risk:  $d = 0.055$ ,  $p = .217$ ; high-risk:  $d = -0.234$ ,  $p = .100$ ), b) concentration problems (medium-risk:  $d = -0.086$ ,  $p = .137$ ; high-risk:  $d = -0.098$ ,  $p = .179$ ), or c) pro-social behavior (medium-risk:  $d = -0.027$ ,  $p = .362$ ; high-risk:  $d = 0.191$ ,  $p = .165$ ) for students in any of the risk groups, relative to the low-risk group (Table 6).

**391 Research Question 4**

392 The sensitivity analysis (Table 7), whereby CR was modeled as a linear continuous  
393 variable, confirmed previous findings, with no significant interaction found between trial  
394 group and risk level in predicting post-test a) disruptive behavior ( $d = 0.004$ ,  $p = .453$ ), b)  
395 concentration problems ( $d = -0.048$ ,  $p = .057$ ), or c) pro-social behavior ( $d = 0.034$ ,  $p = .187$ ).  
396 With regards to the quadratic term (Table 8), a small but significant interaction was found  
397 between trial group and risk level in predicting post-test disruptive behavior only ( $d = -0.092$ ,  
398  $p = <.001$ ). No such effects were found for concentration problems ( $d = 0.002$ ,  $p = .462$ ) or  
399 pro-social behavior ( $d = 0.020$ ,  $p = .213$ ).

**400 Discussion**

401 The results of this RCT demonstrate that the GBG had no main effect on students'  
402 disruptive behavior, concentration problems, or pro-social behavior. Allocation to the  
403 intervention group was also not a statistically significant predictor of outcomes for students  
404 exposed to varying levels of CR. In other words, exposure to the GBG for two years did not  
405 result in significant improvements in students' behavior, irrespective of risk status. These  
406 findings were insensitive to levels of program differentiation, and were largely unaffected by  
407 the way that risk exposure was modeled, although a significant effect was found for  
408 disruptive behavior when CR exposure was modeled using the quadratic term. However,  
409 given the large number of comparisons conducted, the latter finding may well be due to  
410 familywise error as opposed to a genuine effect. Furthermore, the size of the subgroup  
411 moderator effect was very small, and thus this finding was not considered to be practically  
412 meaningful.

**413 No Main Effects of the GBG on Behavior**

414 As noted previously, studies of the GBG typically find positive effects for students'  
415 behavior; therefore the findings from the present study are incongruent with much of the

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

19

1  
2  
3 416 existing literature in the field. However, of the three measures of behavior utilized in the  
4  
5 417 present study, only disruptive behavior has previously been examined in-depth. Although the  
6  
7 418 program logic model also outlines improvements in on-task behavior and pro-social behavior  
8  
9 419 as immediate impacts of the GBG (Chan et al., 2012), Flower and colleagues' meta-analysis  
10  
11 420 (2014) found that only two studies had tested the effects of the intervention on social  
12  
13 421 outcomes (antisocial behaviors and social interactions), and none had examined concentration  
14  
15 422 explicitly (only attention). Thus, while these outcomes are hypothesized to occur as a result  
16  
17 423 of the GBG by the program developers, they have not been extensively tested. Furthermore,  
18  
19 424 studies of the GBG often find greater subgroup than main effects on behavior, namely for  
20  
21 425 boys and students from low socio-economic backgrounds, and some studies only report  
22  
23 426 subgroup effects (e.g. Bradshaw, Zmuda, Kellam, & Ialongo, 2009; Dolan et al., 1993). Thus,  
24  
25 427 when looking specifically at ITT analyses, and more infrequently tested outcomes, it was  
26  
27 428 unclear if any main effects would be found.  
28  
29  
30  
31  
32

33 429 Alternatively, the null results in the present study may be explained by a lack of  
34  
35 430 cultural transferability of the GBG. Although previous research has found the intervention to  
36  
37 431 be effective outside of its country of origin, it was adapted prior to implementation to suit the  
38  
39 432 culture of the countries in which the studies were located (e.g. France, Spain, the Netherlands;  
40  
41 433 Dion et al., 2011; Ruiz-Olivares et al., 2010; van Lier et al., 2004). However, in the present  
42  
43 434 study, the GBG was implemented in its original format. Thematic analyses of the qualitative  
44  
45 435 data from both the Oxfordshire pilot (Chan et al., 2012) and the IPE associated with the  
46  
47 436 current study indicated that teachers had concerns regarding several aspects of the  
48  
49 437 intervention, namely the lack of teacher-student interaction permitted, and the inflexibility of  
50  
51 438 the program not allowing teachers to adapt it to suit their classes' needs. Thus, it is possible  
52  
53 439 that the intervention was not compatible with the school culture in England, and hence its  
54  
55 440 effects were diluted (Wigelsworth et al., 2016).  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

20

1  
2  
3 441 It is therefore also possible that the null results were due to implementation failure,  
4  
5 442 whereby teachers failed to adhere to the steps outlined in the manual due to perceived  
6  
7 443 incompatibility of the intervention. Indeed, almost one-quarter of schools ceased  
8  
9 444 implementation over the course of the trial (although they still complied with data collection  
10  
11 445 protocols, and there were no significant differences between teachers who ceased and  
12  
13 446 sustained implementation). Furthermore, dosage (frequency and duration of game play) was  
14  
15 447 also lower than is recommended (Ford, Keegan, Poduska, Kellam, & Littman, 2014),  
16  
17 448 meaning that the game was not implemented as often or for as long as is specified by the  
18  
19 449 program developers. Thus, it may be that school timetables and other demands placed on  
20  
21 450 teachers (Education Committee, 2017) mean that the GBG does not fit well in English  
22  
23 451 schools. However, the other studies of the GBG in the USA that have measured dosage have  
24  
25 452 reported similar issues (e.g. Domitrovich et al., 2015; Hagermoser Sanetti & Fallon, 2011),  
26  
27 453 meaning that the time to implement the GBG is unlikely to be an issue specific to the English  
28  
29 454 schooling system. Furthermore, while a certain level of dosage is specified in the  
30  
31 455 intervention's manual, these dosage benchmarks have not been empirically validated, and so  
32  
33 456 the levels of dosage necessary to bring about student gains are unknown (Becker, Darney, &  
34  
35 457 Domitrovich, 2013).

36  
37 458 As the present study was an efficacy trial, in which significant resources were made  
38  
39 459 available to optimize implementation (e.g. developer support, subsidized costs), the levels of  
40  
41 460 implementation reported here are likely the 'best case' scenario (that is, we could reasonably  
42  
43 461 expect further dilution of implementation under 'real world' conditions). Furthermore,  
44  
45 462 procedural fidelity, quality, participant responsiveness and reach recorded in our structured  
46  
47 463 observations were high, indicating that with the exception of suboptimal dosage,  
48  
49 464 implementation failure was likely not a key issue. Nevertheless, it is important that future  
50  
51 465 studies examine whether null results are still evident once implementation variability is  
52  
53  
54  
55  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

21

1  
2  
3 466 accounted for using appropriately robust methods (e.g. complier average causal effect  
4  
5 467 estimation; CACE). Indeed, a CACE analysis would not only help to determine whether the  
6  
7 468 null results found in the present study were a result of suboptimal dosage, but would also  
8  
9 469 remove the potential bias introduced in ITT (wherein full compliance with the intervention is  
10  
11 470 assumed) analyses (Peugh, Strotman, McGrady, Rausch, & Kashikar-Zuck, 2017). This will  
12  
13 471 be formally examined in a future paper.  
14  
15

16  
17 472         Alternatively, it is possible that the GBG was not required in English schools. First,  
18  
19 473 behavior management strategies akin to key aspects of the GBG were commonplace in  
20  
21 474 classrooms at baseline, and so it may not have been distinct enough to bring about additional  
22  
23 475 gains in students' behavior. Thus, it may not be that the GBG was ineffective, but that it was  
24  
25 476 simply no *more* effective than the existing behavior management strategies already in use  
26  
27 477 both at baseline and in the usual provision group during the trial. Our sensitivity analysis  
28  
29 478 demonstrated that the interaction between trial group allocation and program differentiation  
30  
31 479 level did not predict post-test outcomes, indicating that the intervention was equally  
32  
33 480 unsuccessful in improving behavior in contexts in which it was more distinct from existing  
34  
35 481 practice. However, it should be borne in mind that, such was the proliferation of GBG-like  
36  
37 482 strategies in use among all participating teachers, it was not possible to model outcomes in  
38  
39 483 the context of high program differentiation. Furthermore, as the collection of usual practice  
40  
41 484 data is not standard procedure when evaluating the effectiveness of an intervention, this  
42  
43 485 finding cannot be compared to previous studies of the GBG that have identified a positive  
44  
45 486 effect. This therefore highlights the importance of collecting this type of data when  
46  
47 487 conducting an RCT.  
48  
49  
50  
51  
52

53 488         Second, contrary to widely publicized concerns noted at the outset of the current study,  
54  
55 489 other data challenges the deficit view of student behavior in English schools. For instance,  
56  
57 490 Ofsted's most recent report identified that 92.3% of all schools in England were judged Good  
58  
59  
60

1  
2  
3 491 or Outstanding for behavior standards (DfE, 2012). Indeed, students' average behavior scores  
4  
5 492 on the TOCA-C were low at baseline across the trial arms (e.g., disruptive behavior = 1.71/6),  
6  
7  
8 493 leaving little room for improvement.

9  
10 494 **No Subgroup Effects Among Students Exposed to Varying Levels of Cumulative Risk**

11  
12 495 The current study was among the first to apply CR theory to subgroup moderator  
13  
14 496 analyses of behavior in a preventive intervention trial. Although no subgroup effects were  
15  
16  
17 497 found in the present study for students at any risk level, the analyses in the wider trial using  
18  
19 498 'traditional' subgroups (boys at-risk of conduct problems and children eligible for FSM) also  
20  
21 499 found no evidence of statistically significant differential gains. This contradicts the majority  
22  
23  
24 500 of GBG studies, which typically find such effects for students using a single risk factor  
25  
26 501 marker (e.g. Dolan et al., 1993; Kellam et al., 2011; Kellam, Ling, Merisca, Brown, &  
27  
28 502 Ialongo, 1998; Kellam, Rebok, Ialongo, & Mayer, 1994). Similarly to the above discussion, it  
29  
30  
31 503 is possible that cultural incompatibility or implementation failure provide an explanation for  
32  
33 504 these incongruous results. Thus, the null findings were likely not a reflection of the choice to  
34  
35 505 utilize CR theory to measure subgroup effects, but instead suggest that the GBG simply did  
36  
37  
38 506 not have any impact on at-risk children, regardless of the approach taken to the modeling of  
39  
40 507 risk.

41  
42 508 This, therefore, does not rule out the utility of CR indices when examining subgroup  
43  
44 509 effects of preventive interventions. Indeed, there is evidence from other studies to suggest its  
45  
46  
47 510 utility (The Multisite Violence Prevention Project, 2008). In line with ecological systems  
48  
49 511 theory (Bronfenbrenner, 1986), it is thought that CR theory provides a more accurate  
50  
51 512 representation of a child's experiences regarding risk exposure, accounting for the clustering  
52  
53  
54 513 of risk factors and interactions between them that are likely to occur (Flouri & Kallis, 2007).  
55  
56 514 However, as the subgroup moderator analysis was exploratory in nature, the small sample  
57  
58 515 sizes of the different risk groups mean that the study may have been under-powered. Thus,  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

23

1  
2  
3 516 larger-scale explorations of subgroup effects utilizing CR indices warrant further attention in  
4  
5 517 future research.

7  
8 518 **Implications**

9  
10 519 The results from this study highlight the importance of reporting null results of trials  
11  
12 520 (Fiennes, 2018), helping to reduce the disconnect between scientific worth and culture  
13  
14 521 (Matosin, Frank, Engel, Lum, & Newell, 2014). In addition to tackling the widely publicized  
15  
16 522 publication bias (which favors statistically significant, ‘positive’ results; Fanelli, 2010), these  
17  
18 523 null results also advance knowledge in the field and are important for several stakeholders.  
19  
20 524 For instance, this finding is beneficial to both funding bodies and schools that need to be  
21  
22 525 aware of the utility of an intervention in order to make an informed decision to implement it  
23  
24 526 with students. This is particularly important when a wealth of evidence already exists  
25  
26 527 regarding the intervention’s efficacy in its country of origin, but it has not previously been  
27  
28 528 rigorously tested within the local culture, as decreases in effectiveness are often identified  
29  
30 529 once an intervention is exported (Wigelsworth et al., 2016). Furthermore, it is vital that  
31  
32 530 intervention developers know if their program is ineffective outside of its country of origin;  
33  
34 531 this information allows them to make decisions regarding necessary adaptations to ensure the  
35  
36 532 intervention’s viability in different countries. Indeed, the finding that the GBG is not  
37  
38 533 effective in achieving its primary intended outcome is a key issue that needs to be addressed.  
39  
40  
41  
42  
43

44 534 The findings from the present study also highlight the broader need to evaluate the  
45  
46 535 cultural transferability of previously successful interventions before they are exported and  
47  
48 536 implemented on a large scale. Significant adaptations may be necessary if the expected  
49  
50 537 outcomes are to be achieved, as they can enhance ownership and commitment, support  
51  
52 538 ‘goodness-of-fit’ to the local culture, and improve sustainability (Lendrum & Humphrey,  
53  
54 539 2012). Indeed, previous research suggests that a major factor in the successful transferability  
55  
56 540 of interventions is their adaptability (Castro et al., 2004). Thus, future research should seek to  
57  
58  
59  
60

1  
2  
3 541 identify the cultural adaptations of both the GBG and other school-based interventions that  
4  
5 542 ensure that they are suitable to the English context, while also ensuring the program's critical  
6  
7 543 components are still in place (Sharples, Albers, & Fraser, 2018).  
8  
9

#### 10 544 **Strengths and Limitations**

11  
12 545 There are several factors that increase the security of the findings noted above. The  
13  
14 546 use of a cluster-randomized design with well-balanced trial arms at the school- and student-  
15  
16 547 levels means that the likelihood of diffusion or contamination effects was minimized  
17  
18 548 (Campbell, Mollison, & Grimshaw, 2001). In addition, the trial was well powered to detect  
19  
20 549 effects (MDES of 0.16), and although attrition was at 18.5% at the student-level, this was  
21  
22 550 within acceptable limits (Dumville, Torgerson, & Hewitt, 2006) and was appropriately  
23  
24 551 addressed using MI. The use of sensitivity analyses to establish that our substantive findings  
25  
26 552 did not vary by program differentiation levels or the way in which CR status was modeled  
27  
28 553 further increases the security of these results.  
29  
30  
31

32  
33 554 Nevertheless, it is possible that null results were the result of research design  
34  
35 555 limitations. For instance, while all schools in the desired regions were invited to participate in  
36  
37 556 the trial, the schools in the sample were typically larger than average, with higher rates of  
38  
39 557 students with an SEND, eligible for FSM, and speaking EAL (DfE, 2015). Indeed, the  
40  
41 558 majority of schools were situated in one densely populated, ethnically diverse region with  
42  
43 559 high levels of socio-economic deprivation. In addition, schools that chose to participate in the  
44  
45 560 trial were likely those that had a greater perceived need for a behavior management  
46  
47 561 intervention. As such, the schools participating in the trial may not have been fully  
48  
49 562 representative of those in England overall. Furthermore, as class composition was not known  
50  
51 563 for schools in the UP arm of the trial, it was not possible to include teachers as a level of  
52  
53 564 analysis in the model. Thus, we therefore missed the opportunity to model teacher  
54  
55 565 characteristics as a potentially strong source of variance.  
56  
57  
58  
59  
60



1  
2  
3 566 While our primary outcome measure, the TOCA-C, has previously been validated,  
4  
5  
6 567 psychometric validation studies are limited, and have typically been conducted by developers  
7  
8 568 of either the TOCA or GBG (e.g. Bradshaw et al., 2015; Kellam et al., 1994). As the TOCA  
9  
10 569 was designed by a developer of the GBG, and has primarily been used in GBG studies, there  
11  
12 570 may be concerns that the measure is ‘inherent to treatment’ (Slavin & Madden, 2011).  
13  
14 571 However, as this issue would have *avored* positive outcomes in the intervention arm of the  
15  
16 572 trial, and no effects of the GBG were identified, it is unlikely that this was an issue. In  
17  
18 573 addition, the TOCA-C is a teacher informant-report measure. Its developers advise that a  
19  
20 574 variety of factors (e.g. the demographic characteristics of the child) can influence teachers’  
21  
22 575 reports of behavior problems. It is also thought that the timing of administration can influence  
23  
24 576 reports, with a notable difference in scores between the beginning and end of the school year  
25  
26 577 (Dolan et al., 1993; Koth et al., 2009). However, such issues apply equally to both trial arms  
27  
28 578 and are therefore unlikely to have biased our findings. Finally, it is acknowledged that no  
29  
30 579 single informant can provide a comprehensive picture of a student’s behavior. Inter-rater  
31  
32 580 correlations between teacher- and self-report behavior measures have previously been found  
33  
34 581 to be weak (Goodman, Meltzer, & Bailey, 1998) and so collecting similar additional data  
35  
36 582 from other informants may have provided a more comprehensive and valid assessment (De  
37  
38 583 Los Reyes et al., 2015).

#### 39 584 **Conclusions**

40 585 The present study is, to the authors’ knowledge, the largest RCT of the GBG  
41  
42 586 worldwide to date, and is the first in an English setting, thus providing a significant  
43  
44 587 contribution regarding the efficacy of the intervention and its cultural transferability when  
45  
46 588 delivered in its original format. The current study also contributes to the evidence base  
47  
48 589 regarding the effects of the GBG on concentration problems and pro-social behavior,  
49  
50 590 examining the previously infrequently tested claims outlined in the program logic model.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

26

1  
2  
3 591 Furthermore, the study has advanced the examination of differential gains by testing  
4  
5 592 subgroup effects for students at varying levels of CR exposure, an area highlighted as a  
6  
7 593 priority for future research in Durlak and colleagues' (2011) meta-analysis, and the influence  
8  
9 594 of program differentiation, an aspect of implementation sorely neglected in prior literature  
10  
11 595 (Durlak & DuPre, 2008).  
12  
13

**Acknowledgements**

14 596  
15  
16  
17 597 This work was supported by funding from [masked for peer review].  
18

**Declaration of Interest Statement**

19 598  
20  
21 599 The authors declare they have no conflicts of interest.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 600 References

- 601 Ashworth, E. & Humphrey, N. (2018). More than the sum of its parts: Cumulative risk effects  
602 on school functioning in middle childhood. *British Journal of Educational Psychology*.  
603 <https://doi.org/10.1111/bjep.12260>
- 604 Ashworth, E., Demkowicz, O., Lendrum, A., & Frearson, K. (2018). Coaching Models of  
605 School-Based Prevention and Promotion Programmes: A Qualitative Exploration of UK  
606 Teachers' Perceptions. *School Mental Health, 10*. [https://doi.org/10.1007/s12310-018-](https://doi.org/10.1007/s12310-018-9282-3)  
607 [9282-3](https://doi.org/10.1007/s12310-018-9282-3)
- 608 Bandura, A. (1977). *Social learning theory*. Englewood Cliffs, NJ: Prentice Hall.
- 609 Barrish, H., Saunders, M., & Wolf, M. (1969). Good Behavior Game: Effects of individual  
610 contingencies for group consequences on disruptive behavior in a classroom. *Journal of*  
611 *Applied Behavior Analysis, 2*, 119–124. <https://doi.org/10.1901/jaba.1969.2-119>
- 612 Becker, K., Darney, D., & Domitrovich, C. (2013). Supporting universal prevention programs:  
613 A two-phased coaching model. *Clinical Child and Family Psychology Review, 16*, 213–  
614 228. <https://doi.org/10.1007/s10567-013-0134-2>. Supporting
- 615 Bierman, K., Nix, R., Heinrichs, B., Domitrovich, C., Gest, S., Welsh, J., & Gill, S. (2014).  
616 Effects of Head Start REDI on children's outcomes 1 year later in different kindergarten  
617 contexts. *Child Development, 85*, 140–159. <https://doi.org/10.1111/cdev.12117>
- 618 Bradshaw, C., Waasdorp, T., & Leaf, P. J. (2015). Examining variation in the impact of  
619 school-wide positive behavioral interventions and supports: Findings from a randomized  
620 controlled effectiveness trial. *Journal of Educational Psychology, 107*, 546–557.  
621 <https://doi.org/10.1037/a0037630>
- 622 Bradshaw, C., Zmuda, J., Kellam, S., & Ialongo, N. (2009). Longitudinal impact of two  
623 universal preventive interventions in first grade on educational outcomes in high school.  
624 *Journal of Educational Psychology, 101*, 926–937. <https://doi.org/10.1037/a0016586>

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

28

- 1  
2  
3 625 Bronfenbrenner, U. (1986). Ecology of the family as a context for human development:  
4  
5 626 Research perspectives. *Developmental Psychology*, 22, 723–742.  
6  
7  
8 627 <https://doi.org/10.1037/0012-1649.22.6.723>  
9  
10 628 Campbell, M., Mollison, J., & Grimshaw, J. (2001). Cluster trials in implementation research:  
11  
12 629 Estimation of intraclass correlation coefficients and sample size. *Statistics in Medicine*,  
13  
14 630 20, 391–399.  
15  
16 631 Carpenter, J., Goldstein, H., & Kenward, M. (2011). REALCOM-IMPUTE software for  
17  
18 632 multilevel multiple imputation with mixed response types. *Journal of Statistical*  
19  
20 633 *Software*, 45, 1–14. <https://doi.org/http://dx.doi.org/10.18637/jss.v045.i05>  
21  
22  
23 634 Castro, F., Barrera, M., & Martinez, C. (2004). The cultural adaptation of prevention  
24  
25 635 interventions: Resolving tensions between fit and fidelity. *Prevention Science*, 5, 41–45.  
26  
27  
28 636 Chan, G., Foxcroft, D., Smurthwaite, B., Coomes, L., & Allen, D. (2012). *Improving child*  
29  
30 637 *behaviour Management: An evaluation of the Good Behaviour Game in UK primary*  
31  
32 638 *schools*.  
33  
34  
35 639 Coombes, L., Chan, G., Allen, D., & Foxcroft, D. (2016). Mixed-methods evaluation of the  
36  
37 640 Good Behaviour Game in English primary schools. *Journal of Community & Applied*  
38  
39 641 *Social Psychology*, 26, 369–387. <https://doi.org/10.1002/casp>  
40  
41  
42 642 De Los Reyes, A., Augenstein, T., Wang, M., Thomas, S., Drabkci, D., Burgers, D., &  
43  
44 643 Rabinowitz, J. (2015). The validity of the multi-informant approach to assessing child  
45  
46 644 and adolescent mental health. *Psychological Bulletin*, 14, 858–900.  
47  
48 645 <https://doi.org/10.1158/1940-6207.CAPR-14-0359.Nrf2-dependent>  
49  
50  
51 646 DfE. (2012). *Pupil behaviour in schools in England*.  
52  
53 647 DfE. (2015). *Schools, pupils, and their characteristics: January 2015*. London, UK: DfE.  
54  
55  
56 648 Dion, E., Roux, C., Landry, D., Fuchs, D., Wehby, J., & Dupéré, V. (2011). Improving  
57  
58 649 attention and preventing reading difficulties among low-income first-graders: A  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

29

- 1  
2  
3 650 randomized study. *Prevention Science*, 12, 70–79. <https://doi.org/10.1007/s11121-010->  
4  
5 651 0182-5  
6  
7  
8 652 Dolan, L., Kellam, S., Hendricks Brown, C., Werthamer-Larsson, L., Rebok, G., Mayer,  
9  
10 653 L., ... Wheeler, L. (1993). The short-term impact of two classroom-based preventive  
11  
12 654 interventions of aggressive and shy behaviors and poor achievement. *Journal of Applied*  
13  
14 655 *Developmental Psychology*, 14, 317–345.  
15  
16  
17 656 Domitrovich, C., Pas, E., Bradshaw, C., Becker, K., Keperling, J., Embry, D., & Ialongo, N.  
18  
19 657 (2015). Individual and School Organizational Factors that Influence Implementation of  
20  
21 658 the PAX Good Behavior Game Intervention. *Prevention Science*, 16, 1064–1074.  
22  
23 659 <https://doi.org/10.1007/s11121-015-0557-8>  
24  
25  
26 660 Dumville, J., Torgerson, D., & Hewitt, C. (2006). Reporting attrition in randomised  
27  
28 661 controlled trials. *BMJ*, 332, 969–971. <https://doi.org/10.1136/bmj.332.7547.969>  
29  
30  
31 662 Durlak, J., & DuPre, E. (2008). Implementation matters: A review of research on the  
32  
33 663 influence of implementation on program outcomes and the factors affecting  
34  
35 664 implementation. *American Journal of Community Psychology*, 41, 327–350.  
36  
37 665 <https://doi.org/10.1007/s10464-008-9165-0>  
38  
39  
40 666 Durlak, J., Weissberg, R., Dymnicki, A., Taylor, R., & Schellinger, K. (2011). The impact of  
41  
42 667 enhancing students' social and emotional learning: A meta-analysis of school-based  
43  
44 668 universal interventions. *Child Development*, 82, 405–432.  
45  
46 669 <https://doi.org/10.1111/j.1467-8624.2010.01564.x>  
47  
48  
49 670 Education Committee. (2017). *Recruitment and retention of teachers*. London, UK.  
50  
51 671 Evans, G., Li, D., & Whipple, S. (2013). Cumulative risk and child development.  
52  
53 672 *Psychological Bulletin*, 139, 1342–1396. <https://doi.org/10.1037/a0031808>  
54  
55  
56 673 Fanelli, D. (2010). Do pressures to publish increase scientists' bias? An empirical support  
57  
58 674 from US states data. *PLoS ONE*, 5. <https://doi.org/10.1371/journal.pone.0010271>  
59  
60

- 1  
2  
3 675 Farrell, A., Henry, D., & Bettencourt, A. (2013). Methodological challenges examining  
4  
5 676 subgroup differences: Examples from universal school-based youth violence prevention  
6  
7 677 trials. *Prevention Science, 14*, 121–133. <https://doi.org/http://doi.org/10.1007/s11121->  
8  
9 678 011-0200-2
- 10  
11  
12 679 Fiennes, C. (2018). Charity begins with admitting we got it wrong. Retrieved from  
13  
14 680 <https://www.ft.com/content/49e715b6-8458-11e8-9199-c2a4754b5a0e>
- 15  
16  
17 681 Flouri, E., & Kallis, C. (2007). Adverse life events and psychopathology and prosocial  
18  
19 682 behavior in late adolescence: Testing the timing, specificity, accumulation, gradient, and  
20  
21 683 moderation of contextual risk. *Journal of the American Academy of Child and*  
22  
23 684 *Adolescent Psychiatry, 46*, 1651–1659. <https://doi.org/10.1097/chi.0b013e318156a81a>
- 24  
25  
26 685 Flower, A., Mckenna, J., Bunuan, R., Muething, C., & Vega, R. (2014). Effects of the Good  
27  
28 686 Behavior Game on challenging behaviors in school settings. *Review of Educational*  
29  
30 687 *Research, 84*, 546–571. <https://doi.org/10.3102/0034654314536781>
- 31  
32  
33 688 Ford, C., Keegan, N., Poduska, J., Kellam, S., & Littman, J. (2014). *Implementation manual*.  
34  
35 689 Washington, DC: American Institutes for Research.
- 36  
37  
38 690 Gerard, J., & Buehler, C. (1999). Multiple risk factors in the family environment and youth  
39  
40 691 problem behaviors. *Journal of Marriage and the Family, 61*, 343–361.
- 41  
42  
43 692 Gerard, J., & Buehler, C. (2004). Cumulative environmental risk and youth maladjustment:  
44  
45 693 The role of youth attributes. *Child Development, 75*, 1832–1849.  
46  
47 694 <https://doi.org/10.1111/j.1467-8624.2004.00820.x>
- 48  
49  
50 695 Goodman, R., Meltzer, H., & Bailey, V. (1998). The strengths and difficulties questionnaire:  
51  
52 696 A pilot study on the validity of the self-report version. *European Child & Adolescent*  
53  
54 697 *Psychiatry, 7*, 125–130. <https://doi.org/10.1080/0954026021000046137>
- 55  
56  
57 698 Greenberg, M., & Abenavoli, R. (2017). Universal Interventions: Fully Exploring Their  
58  
59 699 Impacts and Potential to Produce Population-Level Impacts. *Journal of Research on*  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

31

- 1  
2  
3 700 *Educational Effectiveness*, 10, 40–67. <https://doi.org/10.1080/19345747.2016.1246632>
- 4  
5 701 Gupta, S. (2011). Intention-to-treat concept: A Review. *Perspectives in Clinical Research*, 2,  
6  
7 702 109–112.
- 8  
9 703 Hagermoser Sanetti, L., & Fallon, L. (2011). Treatment integrity assessment: How estimates  
10  
11 704 of adherence, quality, and exposure influence interpretation of implementation. *Journal*  
12  
13 705 *of Educational and Psychological Consultation*, 21, 209–232.  
14  
15 706 <https://doi.org/10.1080/10474412.2011.595163>
- 16  
17 707 Hebron, J., Oldfield, J., & Humphrey, N. (2016). Cumulative risk effects in the bullying of  
18  
19 708 children and young people with autism spectrum conditions. *Autism*, 21, 291–300.  
20  
21 709 <https://doi.org/10.1177/1362361316636761>
- 22  
23 710 Holsen, I., Iversen, A. C., & Smith, B. H. (2009). Universal social competence promotion  
24  
25 711 programme in school: Does it work for children with low socio-economic background?  
26  
27 712 *Advances in School Mental Health Promotion*, 2, 51–60.  
28  
29 713 <https://doi.org/10.1080/1754730X.2009.9715704>
- 30  
31 714 Humphrey, N. (2013). *Social and emotional learning: A critical appraisal*. London, UK:  
32  
33 715 Sage.
- 34  
35 716 Humphrey, N., Lendrum, A., Ashworth, E., Frearson, K., Buck, R., & Kerr, K. (2016).  
36  
37 717 *Implementation and process evaluation (IPE) for interventions in educational settings:*  
38  
39 718 *A synthesis of the literature*. EEF. London, UK.  
40  
41 719 <https://doi.org/10.1017/CBO9781107415324.004>
- 42  
43 720 Ialongo, N., Werthamer, L., Kellam, S., Brown, C., Wang, S., & Lin, Y. (1999). Proximal  
44  
45 721 impact of two first-grade preventive interventions on the early risk behaviors for later  
46  
47 722 substance abuse, depression, and antisocial behavior. *American Journal of Community*  
48  
49 723 *Psychology*, 27, 599–641. <https://doi.org/10.1023/A:1022137920532>
- 50  
51 724 Kellam, S., Branch, J., Agrawal, K., & Ensminger, M. (1975). *Mental health and going to*  
52  
53  
54  
55  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

32

- 1  
2  
3 725 *school: The woodlawn program of assessment, early intervention, and evaluation.*  
4  
5 726 Chicago: University of Chicago Press.  
6  
7 727 Kellam, S., Ling, X., Merisca, R., Brown, C. H., & Ialongo, N. (1998). The effect of the level  
8  
9  
10 728 of aggression in the first grade classroom on the course and malleability of aggressive  
11  
12 729 behavior into middle school. *Development and Psychopathology, 10*, 169–170.  
13  
14 730 <https://doi.org/10.1017/S0954579498001564>  
15  
16 731 Kellam, S., Mackenzie, A., Brown, C., Poduska, J., Wang, W., Petras, H., & Wilcox, H.  
17  
18 732 (2011). The Good Behavior Game and the future of prevention and treatment. *Addiction*  
19  
20 733 *Science & Clinical Practice, 6*, 73–84.  
21  
22 734 Kellam, S., Mayer, L., Rebok, G., & Hawkins, W. (1998). Effects of improving achievement  
23  
24 735 on aggressive behavior and of improving aggressive behavior on achievement through  
25  
26 736 two preventive interventions: An investigation of causal paths. In B. P. Dohrenwend  
27  
28 737 (Ed.), *Adversity, stress, and psychopathology* (pp. 591–600). New York, NY: Oxford  
29  
30 738 University Press. <https://doi.org/10.1038/nrm3860.RNAi>  
31  
32 739 Kellam, S., Rebok, G., Ialongo, N., & Mayer, L. (1994). The course and malleability of  
33  
34 740 aggressive behavior from early first grade into middle school: Results of a  
35  
36 741 developmental epidemiologically-based preventive trial. *Journal of Child Psychology*  
37  
38 742 *and Psychiatry, 35*, 259–281.  
39  
40 743 Kleinman, K., & Saigh, P. (2011). The effects of the Good Behavior Game on the conduct of  
41  
42 744 regular education New York city high school students. *Behavior Modification, 35*, 95–  
43  
44 745 105. <https://doi.org/10.1177/0145445510392213>  
45  
46 746 Korpershoek, H., Harms, T., de Boer, H., van Kuijk, M., & Doolaard, S. (2016). A meta-  
47  
48 747 analysis of the effects of classroom management strategies and classroom management  
49  
50 748 programs on students' academic, behavioral, emotional, and motivational outcomes.  
51  
52 749 *Review of Educational Research, 86*, 643–680.  
53  
54  
55  
56  
57  
58  
59  
60



## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

33

- 1  
2  
3 750 <https://doi.org/10.3102/0034654315626799>  
4  
5 751 Koth, C., Bradshaw, C., & Leaf, P. (2009). Teacher observation of classroom adaptation-  
6  
7 752 checklist: Development and factor structure. *Measurement and Evaluation in*  
8  
9 753 *Counseling and Development, 42*, 15–30. <https://doi.org/10.1177/0748175609333560>  
10  
11 754 Kourkounasiou, M., & Skordilis, E. (2014). Validity and reliability evidence of the TOCA-C  
12  
13 755 in a sample of Greek students. *Psychological Reports, 115*, 766–783.  
14  
15 756 Lastrapes, R. E. (2013). Using the Good Behavior Game in an inclusive classroom.  
16  
17 757 *Intervention in School and Clinic, 49*, 225–229.  
18  
19 758 <https://doi.org/10.1177/1053451213509491>  
20  
21 759 Leflot, G., Van Lier, P., Onghena, P., & Colpin, H. (2010). The role of teacher behavior  
22  
23 760 management in the development of disruptive behaviors: An intervention study with the  
24  
25 761 Good Behavior Game. *Journal of Abnormal Child Psychology, 38*, 869–882.  
26  
27 762 <https://doi.org/10.1007/s10802-010-9411-4>  
28  
29 763 Lendrum, A., & Humphrey, N. (2012). The importance of studying the implementation of  
30  
31 764 interventions in school settings. *Oxford Review of Education, 38*, 635–652.  
32  
33 765 <https://doi.org/10.1080/03054985.2012.734800>  
34  
35 766 Matosin, N., Frank, E., Engel, M., Lum, J., & Newell, K. (2014). Negativity towards negative  
36  
37 767 results: A discussion of the disconnect between scientific worth and scientific culture.  
38  
39 768 *Disease Models & Mechanisms, 7*, 171–173. <https://doi.org/10.1242/dmm.015123>  
40  
41 769 Muthén, B., Brown, C. H., Booil, K., Khoo, S., Yang, C., Wang, C., ... Liao, J. (2002).  
42  
43 770 General growth mixture modeling for randomized preventive interventions. *Biostatistics,*  
44  
45 771 *3*, 459–475.  
46  
47 772 National Institute for Health and Care Excellence. (2013). *Antisocial behaviour and conduct*  
48  
49 773 *disorders in children and young people: recognition, intervention and management.*  
50  
51 774 London: NICE.  
52  
53  
54  
55  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

34

- 1  
2  
3 775 Office for Standards in Education. (2013). *The report of Her Majesty's Chief Inspector of*  
4  
5 776 *Education, Children's Services and Skills*. London: Ofsted.  
6  
7  
8 777 Office for Standards in Education. (2014). *Below the radar: Low-level disruption in the*  
9  
10 778 *country's classrooms*. London: Ofsted.  
11  
12 779 Oldfield, J., Humphrey, N., & Hebron, J. (2015). Cumulative risk effects for the development  
13  
14 780 of behaviour difficulties in children and adolescents with special educational needs and  
15  
16 781 disabilities. *Research in Developmental Disabilities, 41*, 66–75.  
17  
18 782 <https://doi.org/10.1016/j.ridd.2015.05.010>  
19  
20  
21 783 Pampaka, M., Hutcheson, G., & Williams, J. (2016). Handling missing data: Analysis of a  
22  
23 784 challenging data set using multiple imputation. *International Journal of Research &*  
24  
25 785 *Method in Education, 39*, 19–37. <https://doi.org/10.1080/1743727X.2014.979146>  
26  
27  
28 786 Pas, E., Bradshaw, C., Becker, K., Domitrovich, C., Berg, J., Musci, R., & Ialongo, N. (2015).  
29  
30 787 Identifying patterns of coaching to support the implementation of the good behavior  
31  
32 788 game: The role of teacher characteristics. *School Mental Health, 7*, 61–73.  
33  
34 789 <https://doi.org/10.1007/s12310-015-9145-0>  
35  
36  
37 790 Peugh, J., Strotman, D., McGrady, M., Rausch, J., & Kashikar-Zuck, S. (2017). Beyond  
38  
39 791 intent to treat (ITT): A complier average causal effect (CACE) estimation primer.  
40  
41 792 *Journal of School Psychology, 60*, 7–24. <https://doi.org/10.1016/j.jsp.2015.12.006>  
42  
43  
44 793 Reupert, A., & Woodcock, S. (2010). Success and near misses: Pre-service teachers' use,  
45  
46 794 confidence and success in various classroom management strategies. *Teaching and*  
47  
48 795 *Teacher Education, 26*, 1261–1268. <https://doi.org/10.1016/j.tate.2010.03.003>  
49  
50  
51 796 Ruiz-Olivares, R., Pino, M., & Herruzo, J. (2010). Reduction of disruptive behaviors using an  
52  
53 797 intervention based on the Good Behavior Game and the Say-Do-Report correspondence.  
54  
55 798 *Psychology in the Schools, 47*, 1046–1058. <https://doi.org/10.1002/pits.20523>  
56  
57  
58 799 Rutter, M. (1979). Maternal deprivation, 1972-1978: New findings, new concepts, new  
59  
60

- 1  
2  
3 800 approaches. *Child Development*, 50, 283–305.  
4  
5 801 Saigh, P., & Umar, A. (1983). The effects of a Good Behavior Game on the disruptive  
6  
7 802 behavior of Sudanese elementary school students. *Journal of Applied Behavior Analysis*,  
8  
9 803 16, 339–344.  
10  
11 804 Sameroff, A., Gutman, L., & Peck, S. (2003). Adaptation among youth facing multiple risks:  
12  
13 805 Protective research findings. In S. Luthar (Ed). *Resilience and Vulnerability* (pp. 364–  
14  
15 806 391). Cambridge, UK: Cambridge University Press.  
16  
17 807 <https://doi.org/10.1176/appi.ajp.162.8.1553-a>  
18  
19 808 Sharples, J., Albers, B., & Fraser, S. (2018). *Putting evidence to work: A school's guide to*  
20  
21 809 *implementation*. London, UK: EEF.  
22  
23 810 Skinner, B. F. (1948). "Superstition" in the pigeon. *Journal of Experimental Psychology*, 38,  
24  
25 811 168–172.  
26  
27 812 Slavin, R., & Madden, N. (2011). Measures inherent to treatments in program effectiveness  
28  
29 813 reviews. *Journal of Research on Educational Effectiveness*, 4, 370–380.  
30  
31 814 <https://doi.org/10.1080/19345747.2011.558986>  
32  
33 815 Spilt, J., Koot, J., & van Lier, P. (2013). For whom does It work? Subgroup differences in the  
34  
35 816 effects of a school-based universal prevention program. *Prevention Science*, 14, 479–  
36  
37 817 488. <https://doi.org/10.1007/s11121-012-0329-7>  
38  
39 818 The Multisite Violence Prevention Project. (2008). The multisite violence prevention project:  
40  
41 819 Impact of a universal school-based violence prevention program on social-cognitive  
42  
43 820 outcomes. *Prevention Science*, 9, 231–244. <https://doi.org/10.1007/s11121-008-0101-1>  
44  
45 821 Tingstrom, D., Sterling-Turner, H., & Wilczynski, S. (2006). The Good Behavior Game:  
46  
47 822 1969-2002. *Behavior Modification*, 30, 225–253.  
48  
49 823 <https://doi.org/10.1177/0145445503261165>  
50  
51 824 Twisk, J. (2006). *Applied multilevel analysis*. Cambridge, UK: Cambridge University Press.  
52  
53  
54  
55  
56  
57  
58  
59  
60

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

36

- 1  
2  
3 825 van Lier, P., Muthén, B., van der Sar, R., & Crijnen, A. (2004). Preventing disruptive  
4  
5 826 behavior in elementary school children: Impact of a universal classroom-based  
6  
7 827 intervention. *Journal of Consulting and Clinical Psychology*, *72*, 467–478.  
8  
9  
10 828 <https://doi.org/10.1037/0022-006X.72.3.467>  
11  
12 829 Wehby, J., Maggin, D., Partin, T., & Robertson, R. (2011). The impact of working alliance,  
13  
14 830 social validity, and teacher burnout on implementation fidelity of the Good Behavior  
15  
16 831 Game. *School Mental Health*, *4*, 22–33. <https://doi.org/10.1007/s12310-011-9067-4>  
17  
18  
19 832 Wigelsworth, M., Lendrum, A., Oldfield, J., Scott, A., ten Bokkel, I., Tate, K., & Emery, C.  
20  
21 833 (2016). The impact of trial stage, developer involvement and international transferability  
22  
23 834 on universal social and emotional learning programme outcomes: A meta-analysis.  
24  
25 835 *Cambridge Journal of Education*, *46*, 347–376.  
26  
27  
28 836 <https://doi.org/10.1080/0305764X.2016.1195791>  
29  
30  
31 837  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 1***Mean baseline characteristics at individual and school levels*

Demographic	School			Student		
	Overall	GBG	UP	Overall	GBG	UP
Size – number of pupils on roll	306.9	298.2	315.4	-	-	-
Sex – proportion of male students	-	-	-	52.6	50.4	54.9
FSM – proportion of pupils eligible for FSM	26.0	27.6	24.5	24.8	27.4	22.8
EAL – proportion of pupils speaking EAL	22.6	22	23.2	27.3	26.2	31
Ethnic Minority – proportion of ethnic minority pupils	32.9	32.4	33.3	33.5	32.8	34.2
SEND – proportion of pupils with SEND	19.5	20.9	18.2	20.3	23.1	17.6

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

**Table 2**

*Cumulative risk index development*

Significant risk factors for behavioral outcomes				Number of students by risk group and trial allocation					
Risk factor (RF)	Disruptive	Concentration	Pro-social	Risk Group	N RF	N students	% sample	N GBG	N UP
School-level				Disruptive behavior					
High urbanicity		✓		Low-risk	0+1	1638	56	762	876
% EAL students	✓		✓	Medium-risk	2+3	1181	40	644	537
% students	✓	✓	✓	High-risk	4+	119	4	82	37
conduct problems				Concentration problems					
				Low-risk	0+1	568	19.9	285	283
Student-level				Medium-risk	2+3	1780	62.2	837	943
Male gender	✓	✓	✓	High-risk	4+5	513	17.9	303	210
Summer-born				Pro-social behavior					
FSM eligible	✓	✓	✓	Low-risk	0+1	1680	55.3	782	898
SEND	✓	✓	✓	Medium-risk	2+3	1228	40.4	674	554
Looked-after child	✓	✓	✓	High-risk	4+	129	4.2	88	41

**Table 3***Mean (standard error) behavior scores by trial group allocation*

	GBG arm		UP arm	
	Pre-test score	Post-test score	Pre-test score	Post-test score
Disruptive behavior	1.66 (.022)	1.74 (.025)	1.59 (.022)	1.64 (.023)
Low-risk	1.46 (0.21)	1.59 (0.03)	1.40 (0.20)	1.46 (0.02)
Medium-risk	1.91 (0.37)	1.89 (0.04)	1.91 (0.42)	1.90 (0.05)
High-risk	2.37 (0.10)	2.01 (0.11)	2.32 (0.19)	2.46 (0.20)
Concentration problems	2.57 (.033)	2.54 (.033)	2.53 (.032)	2.49 (.031)
Low-risk	1.94 (.050)	2.06 (.061)	1.91 (.052)	1.88 (.055)
Medium-risk	2.58 (.037)	2.47 (.043)	2.55 (.037)	2.53 (.039)
High-risk	3.33 (.064)	3.08 (.193)	3.40 (.078)	3.26 (.080)
Pro-social behavior	4.94 (.025)	4.81 (.027)	4.97 (.025)	4.93 (.026)
Low-risk	5.16 (.027)	4.96 (.035)	5.21 (.027)	5.15 (.030)
Medium-risk	4.66 (.036)	4.67 (.042)	4.63 (.044)	4.63 (.046)
High-risk	4.25 (.093)	4.53 (.111)	3.93 (.136)	4.13 (.157)

**Table 4***Main effects of the GBG on behavior*

	Coefficient $\beta$	SE	p value
Disruptive behavior			
$\beta_{0ij} = -1.491(0.120); df = 6$			
School-level	0.093	0.018	<.001**
Trial group (if GBG)	0.056	0.077	.235
Proportion FSM	0.001	0.003	.370
Size	0.000	0.000	.500
Student-level	0.522	0.015	<.001**
FSM	0.098	0.036	.004*
SEN	0.132	0.040	<.001**
Gender (male)	0.150	0.031	<.001**
Baseline disruptive behavior scores	0.786	0.021	<.001**
Concentration problems			
$\beta_{0ij} = -1.726(0.130); df = 6$			
School-level	0.115	0.021	<.001**
Trial group (if GBG)	0.022	0.084	.400
Proportion FSM	0.152	0.035	<.001**
Size	0.000	0.000	.500
Student-level	0.498	0.014	<.001**
FSM	0.004	0.003	.093
SEN	0.256	0.041	<.001**
Gender (male)	0.232	0.030	<.001**
Baseline concentration	0.506	0.015	<.001**



## problems scores

## Pro-social behavior

$$\beta_{0ij} = -1.963(0.193); df = 6$$

School-level	0.195	0.035	<.001**
Trial group (if GBG)	-0.108	0.108	.160
Proportion FSM	-0.003	0.004	.160
Size	-0.000	0.000	.500
Student-level	0.604	0.017	<.001**
FSM	-0.136	0.039	<.001**
SEN	-0.278	0.044	<.001**
Gender (male)	-0.171	0.032	<.001**
Baseline pro-social	0.476	0.021	<.001**

## behavior scores

\* p<.05; \*\* p<.01

FSM = free school meals; SEN = special educational needs

**Table 5***Interaction between trial group allocation and program differentiation*

	Coefficient $\beta$	SE	p value
Disruptive behavior			
$\beta_{0ij} = -1.369(0.087); df = 4$			
School-level	0.092	0.018	<.001**
Trial group (if GBG)	0.042	0.109	.351
GBG* high differentiation	0.028	0.154	.428
High differentiation	-0.022	0.107	.419
Student-level	0.534	0.016	<.001**
Baseline disruptive behavior scores	0.832	0.020	<.001**
Concentration problems			
$\beta_{0ij} = -1.455(0.099); df = 4$			
School-level	0.127	0.024	<.001**
Trial group (if GBG)	-0.007	0.125	
High differentiation	-0.078	0.123	.478
GBG* high differentiation	0.065	0.177	.357
Student-level	0.524	0.015	<.001**
Baseline concentration problems scores	0.577	0.014	<.001**
Prosocial skills			
$\beta_{0ij} = -2.690(0.149); df = 4$			
School-level	0.194	0.035	<.001**
Trial group (if GBG)	0.012	0.153	.469
High differentiation	0.115	0.152	.226
GBG* high differentiation	-0.269	0.217	.109
Student-level	0.626	0.018	<.001**
Baseline prosocial skills scores	0.545	0.020	<.001**

\*  $p < .05$ ; \*\*  $p < .01$ 

FSM = free school meals; SEN = special educational needs

**Table 6***Subgroup effects of the GBG on behavior for students at-risk*

	Coefficient $\beta$	SE	p value
Disruptive behavior			
$\beta_{0ij} = -1.487(0.119); df = 11$			
School-level	0.090	0.018	<.001**
Trial group (if GBG)	0.049	0.082	.476
Proportion FSM	0.001	0.003	.370
Size	0.000	0.000	.500
Student-level	0.523	0.015	<.001**
FSM	0.118	0.044	.005*
SEN	0.159	0.049	<.001**
Gender (male)	0.170	0.039	<.001**
Baseline disruptive behavior scores	0.785	0.021	<.001**
Risk group:			
Low-risk $\diamond$	$\diamond$	$\diamond$	$\diamond$
Medium-risk	-0.062	0.066	.175
High-risk	-0.005	0.167	.488
GBG*medium-risk	0.055	0.070	.217
GBG*high-risk	-0.234	0.181	.100
Concentration problems			
$\beta_{0ij} = -1.766(0.135); df = 11$			
School-level	0.115	0.022	<.001**
Trial group (if GBG)	0.090	0.104	.195

1				
2				
3	Proportion FSM	0.004	0.003	.093
4				
5	Size	0.000	0.000	.500
6				
7	Student-level	0.497	0.014	<.001**
8				
9				
10	FSM	0.137	0.041	<.001**
11				
12	SEN	0.237	0.046	<.001**
13				
14	Gender (male)	0.215	0.037	<.001**
15				
16	Baseline concentration			
17				
18	problems scores	0.507	0.015	<.001**
19				
20				
21	Risk group:			
22				
23				
24	Low-risk	◇	◇	◇
25				
26	Medium-risk	0.068	0.060	.130
27				
28	High-risk	0.108	0.100	.142
29				
30	GBG*medium-risk	-0.086	0.078	.137
31				
32	GBG*high-risk	-0.098	0.106	.179
33				
34				

---

Prosocial skills

$\beta_{0ij} = -1.962(0.194); df = 11$

---

35				
36				
37				
38				
39				
40	School-level	0.193	0.035	<.001**
41				
42	Trial group (if GBG)	-0.108	0.112	.169
43				
44	Proportion FSM	-0.003	0.004	.228
45				
46	Size	-0.000	0.000	.500
47				
48	Student-level	0.603	0.017	<.001**
49				
50				
51	FSM	-0.155	0.048	<.001**
52				
53	SEN	-0.301	0.054	<.001**
54				
55	Gender (male)	-0.185	0.042	<.001**
56				
57	Baseline prosocial	0.476	0.021	<.001**
58				
59				
60				

1  
2  
3 skills scores  
4

5 Risk group:  
6

7	Low-risk	◇	◇	◇
8				
9	Medium-risk	0.034	0.073	.321
10				
11	High-risk	0.019	0.185	.459
12				
13	GBG*medium-risk	-0.027	0.076	.362
14				
15	GBG*high-risk	0.191	0.195	.165
16				

17  
18  
19 \* p<.05; \*\* p<.01

20  
21 ◇ reference category

22  
23 FSM = free school meals; SEN = special educational needs  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 7**

*Sensitivity analysis for effects of the GBG on behavior for students at-risk (continuous CR score)*

	Coefficient $\beta$	SE	p value
Disruptive behavior			
$\beta_{0ij} = -1.457(0.120); df = 9$			
School-level	0.087	0.017	<.001**
Trial group (if GBG)	0.070	0.090	.230
Proportion FSM	0.001	0.003	.370
Size	0.000	0.000	.500
Student-level	0.522	0.015	<.001**
FSM	0.221	0.066	<.001**
SEN	0.254	0.059	<.001**
Gender (male)	0.277	0.064	<.001**
Baseline disruptive behavior scores	0.789	0.022	<.001**
CR score	-0.129	0.060	.017*
GBG*CR score	0.004	0.034	.453
Concentration problems			
$\beta_{0ij} = -1.757(0.140); df = 9$			
School-level	0.111	0.021	<.001**
Trial group (if GBG)	0.142	0.112	.104
Proportion FSM	0.004	0.003	.093
Size	0.000	0.000	.500
Student-level	0.497	0.014	<.001**
FSM	0.172	0.046	<.001**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

SEN	0.279	0.051	<.001**
Gender (male)	0.253	0.042	<.001**
Baseline concentration	0.507	0.015	<.001**
problems scores			
CR score	0.003	0.034	.460
GBG*CR score	-0.048	0.030	.057

---

 Prosocial skills

 $\beta_{0ij} = -1.951(0.196); df = 9$ 


---

School-level	0.189	0.034	<.001**
Trial group (if GBG)	-0.165	0.120	.087
Proportion FSM	-0.003	0.004	.228
Size	-0.000	0.000	.500
Student-level	0.604	0.017	<.001**
FSM	-0.184	0.087	.019*
SEN	-0.329	0.089	<.001**
Gender (male)	-0.220	0.085	<.001**
Baseline prosocial	0.576	0.022	<.001**
skills scores			
CR score	0.033	0.081	.343
GBG*CR score	0.034	0.038	.187

\* p&lt;.05; \*\* p&lt;.01

FSM = free school meals; SEN = special educational needs

**Table 8***Sensitivity analysis for effects of the GBG on behavior for students at-risk (quadratic)*

	Coefficient $\beta$	SE	p value
<b>Disruptive behavior</b>			
$\beta_{0ij} = -1.541(0.119); df = 9$			
School-level	0.089	0.017	<.001**
Trial group (if GBG)	0.162	0.080	.023*
Proportion FSM	0.001	0.003	.370
Size	0.000	0.000	.500
Student-level	0.519	0.015	<.001**
FSM	0.104	0.036	2.889
SEN	0.141	0.042	<.001**
Gender (male)	0.162	0.031	<.001**
Baseline disruptive behavior scores	0.786	0.021	<.001**
Quadratic CR	0.045	0.017	.005*
GBG*Quadratic CR	-0.092	0.023	<.001**
<b>Concentration problems</b>			
$\beta_{0ij} = -1.712(0.132); df = 9$			
School-level	0.115	0.022	<.001**
Trial group (if GBG)	0.022	0.088	.402
Proportion FSM	0.004	0.003	.093
Size	0.000	0.000	.500
Student-level	0.500	0.015	<.001**
FSM	0.151	0.036	<.001**
SEN	0.257	0.044	<.001**



Gender (male)	0.232	0.030	<.001**
Baseline concentration	0.504	0.016	<.001**
problems scores			
Quadratic CR	-0.005	0.017	.385
GBG*Quadratic CR	0.002	0.021	.462

---

 Prosocial skills

 $\beta_{0ij} = -1.963(0.194); df = 9$ 


---

School-level	0.194	0.035	<.001**
Trial group (if GBG)	-0.132	0.111	.119
Proportion FSM	-0.003	0.004	.228
Size	-0.000	0.000	.500
Student-level	0.603	0.017	<.001**
FSM	-0.145	0.039	<.001**
SEN	-0.293	0.046	<.001**
Gender (male)	-0.171	0.033	<.001**
Baseline prosocial	0.476	0.021	<.001**
skills scores			
Quadratic CR	0.004	0.018	.413
GBG*Quadratic CR	0.020	0.025	.213

\* p<.05; \*\* p<.01

FSM = free school meals; SEN = special educational needs

### Supplemental Files

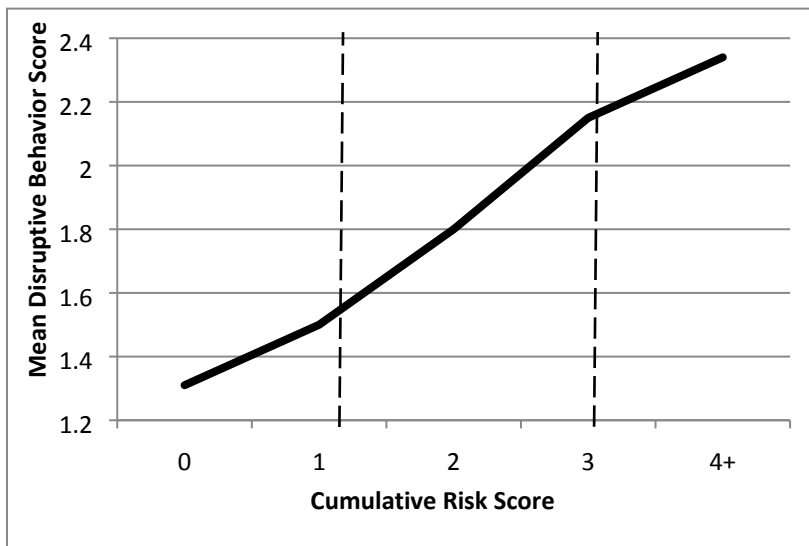
**Table A.1**

*Change in behavior scores between risk levels*

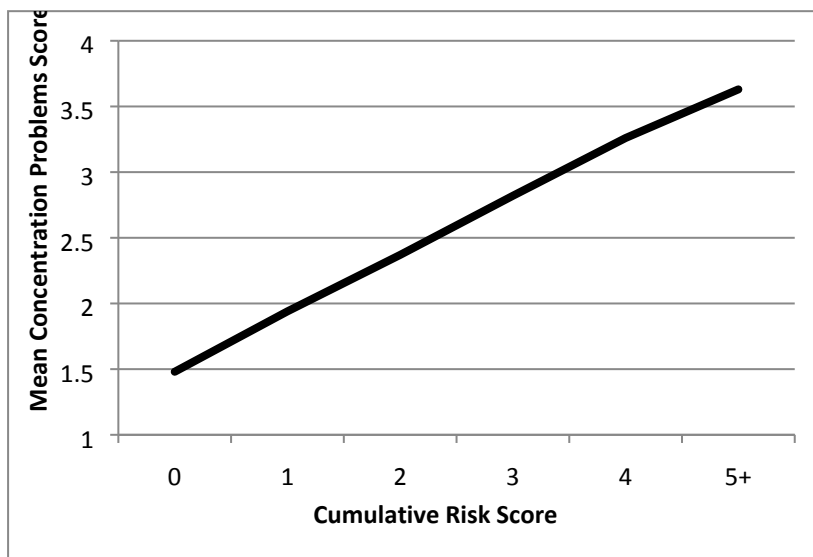
Risk group	Effect size (Cohen's d)
<b>Disruptive behavior</b>	
0-1	0.35
1-2	0.39
2-3	0.36
3-4+	0.19
<b>Concentration problems</b>	
0-1	0.61
1-2	0.45
2-3	0.42
3-4	0.39
4-5+	0.35
<b>Pro-social behavior</b>	
0-1	0.28
1-2	0.45
2-3	0.35
3-4+	0.25

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

3



**Fig A.1.** Line graph demonstrating cumulative risk effect – disruptive behavior



**Fig A.2.** Line graph demonstrating cumulative risk effect - concentration problems

RANDOMIZED TRIAL OF THE GBG IN ENGLAND

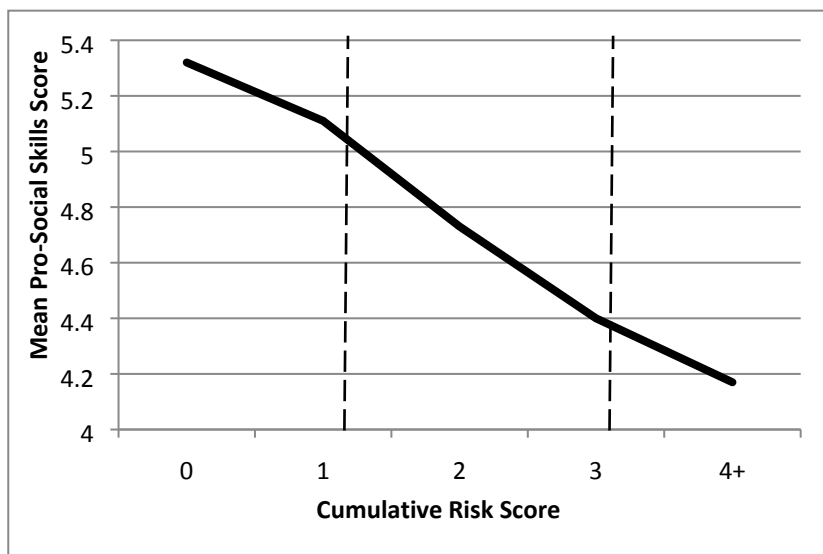


Fig A.3. Line graph demonstrating cumulative risk effect - pro-social behavior

## RANDOMIZED TRIAL OF THE GBG IN ENGLAND

5

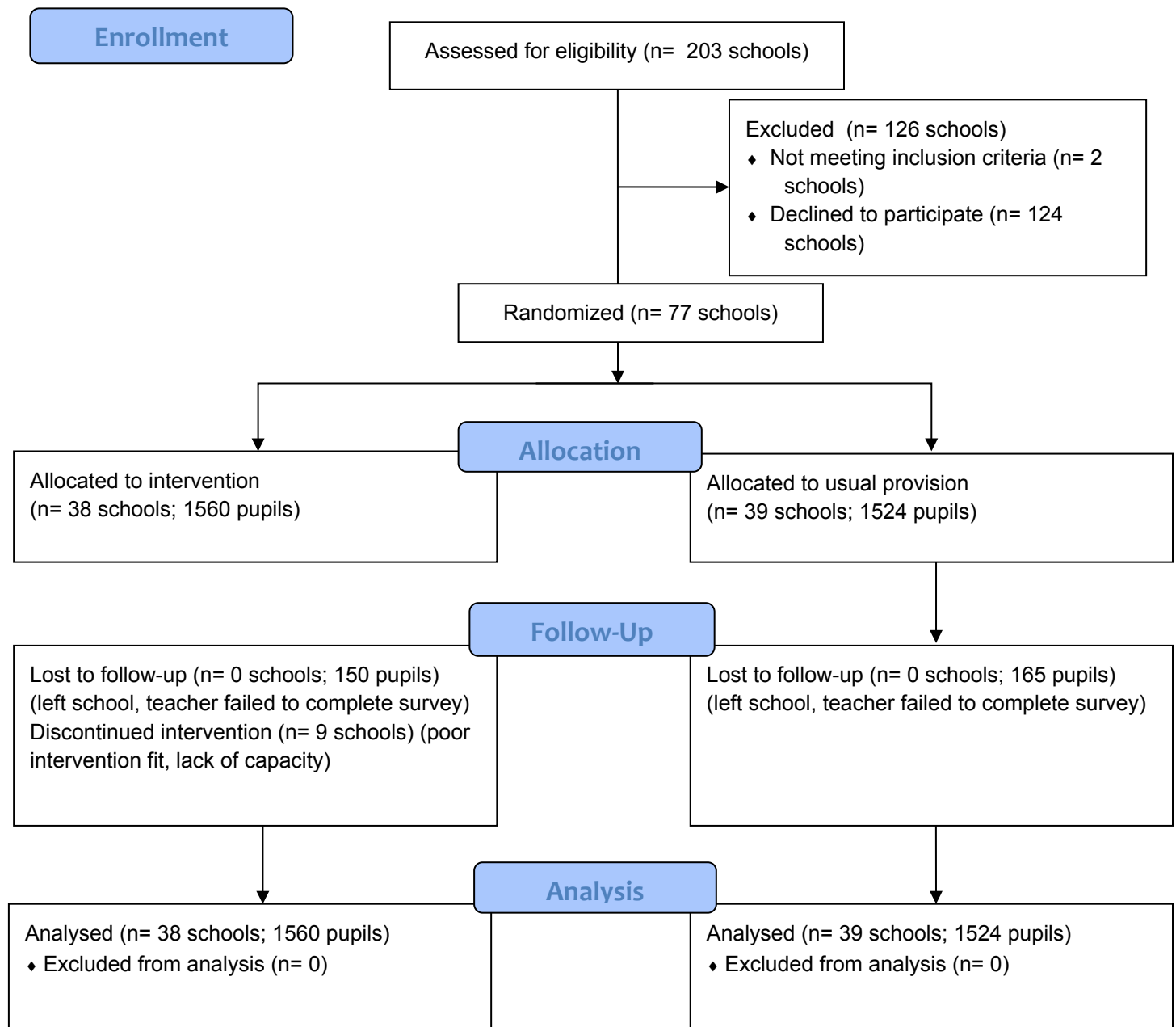


Fig A.4 CONSORT diagram