Education
Endowment
Foundation

| INTERVENTION | Best Practice in Setting |
|---|---|
| DEVELOPER | UCL Institute of Education |
| EVALUATOR | National Foundation for Educational Research (NFER) |
| TRIAL REGISTRATION NUMBER | ISRCTN17963123 |
| TRIAL STATISTICIAN | Palak Roy |
| TRIAL CHIEF INVESTIGATOR | Ben Styles |
| SAP AUTHOR | Palak Roy & Ben Styles |
| SAP VERSION | 1 |
| SAP VERSION DATE | 15/10/2017 |
| EEF DATE OF APPROVAL | 18/10/2017 |
| DEVELOPER DATE OF APPROVAL | 18/10/2017 |

## Protocol and SAP changes

No changes since updated protocol was published.

# Table of contents

## Introduction

The Education Endowment Foundation (EEF) has commissioned UCL Institute of Education to investigate best practice in grouping students by attainment. The project is led by Professor Becky Francis and consists of two randomised controlled trials. The first trial tests an intervention which trains schools in a best practice approach to setting (BPS). The second trial is a feasibility study and pilot RCT exploring the use of mixed attainment teaching in secondary schools (BPMA).

This SAP refers to the first trial (BPS); the SAP for BPMA can also be found on the EEF website. The intervention will help schools address poor practices, which include misallocation, low expectations, less demanding curricula and fixed positioning in low groups. The trial will focus on teaching within English and mathematics in years 7 and 8.

## Study design

The evaluation is as a cluster-randomised controlled trial (RCT), that started in September 2015 and follows children through years 7 and 8. The trial has two arms: intervention and control. The planned sample size for the evaluation was 120 secondary schools randomised to receive either the intervention or to be part of a control group.

### Description of population including eligibility criteria

The population for this trial is all state-funded English secondary schools. At recruitment, it was assumed that most schools were employing setting for both key stage 3 English and mathematics. However, to support recruitment any school using ability grouping in year 7 was eligible to take part regardless of their prior grouping arrangements, including schools that were undertaking streaming. Schools that stream (i.e. allocate their pupils to fixed ability groups across subjects) were eligible if they were prepared to amend their grouping arrangements to setting. The following eligibility rules were used to recruit schools i.e. the trial population cannot be considered to include schools that employ mixed ability grouping for the subject in question.

**Table 1: Eligibility criteria to recruit Secondary schools for BPS trial**

| Year 7 | | Year 8 | | Eligible? |
|---|---|---|---|---|
| English | Mathematics | English | Mathematics | |
| Setting | Setting | Setting | Setting | Yes |
| Mixed | Mixed | Setting | Setting | No |
| Mixed | Setting | Mixed | Setting | Yes (for mathematics outcome only) |
| Mixed | Mixed | Mixed | Setting | No |
| Mixed | Mixed | Mixed | Mixed | No |
| Streaming | Streaming | Streaming | Streaming | Yes, if prepared to set |

As seen in above table, during recruitment it was assumed that schools always prefer setting in mathematics to setting in English. Therefore, not all schools would be willing or able to participate in both English and mathematics. In addition to this, the sample included schools that took part only in English. i.e. only applying the intervention in English (if randomised to intervention group) and therefore eligible for only the English trial. Therefore, separate sample sizes and power calculations were considered for both the subjects and outcome measures.

Intervention group: Heads of Maths and English departments1 along with teachers delivering Maths and English to the Year 7 cohort in intervention schools to attend four workshops provided across the two-year intervention. These were focused on the KS3 leadership developing a department-wide approach to addressing the factors identified: increasing fluidity of movement between sets, raising teacher expectations and pedagogy in lower sets, improving access to the whole curriculum and higher status qualifications, improving pupil engagement and attitudes and tackling the 'self-fulfilling' prophecy of ending up in a lower set.

Control group: Schools from this group continued with their grouping practices as usual. Once they have completed the year 8 tests at the end of the trial, they will receive £1,000.

### Number and timing of measurement points

The trial will span over two academic years (2015/16 to 2016/17) where the outcome measurement will take place once at the end of academic year 2016/17.

# Protocol amendments

As planned in the initial protocol, it was not possible to collect year 7 pupil data such as pupil names, DOBs and UPNs before randomisation as schools were required to administer parental consent to opt out from the data collection. Therefore, this commenced after the randomisation in academic term autumn 2015. However, the baseline measure for the primary outcome is the pupil attainment at Key Stage 2 which took place prior to randomisation.

# Randomisation

The intervention is a whole-school approach, i.e., allocation of teachers across sets. Therefore, school-level randomisation was conducted. Due to the difficulty in recruiting enough schools that wished to set in English, three different types of schools were identified- those taking part in both subjects, those taking part in Maths only and those taking part in English only. Therefore, it was necessary to stratify the randomisation by setting practice (English/maths/both) to allow a lower powered analysis to be run on English outcomes. Randomisation was carried out by a statistician at NFER using a full syntax audit trail. This was done in five blocks due to staggered school recruitment and intervention workshops running concurrently. Randomisation was conducted in June-July 2015. Details on the blocked randomisation is included in Table 2.

There were 129 Secondary schools randomised for this trial. One school withdrew participation without knowledge of group allocation and another school was randomised due to an administrative error. Of the 127 schools, 75 schools took part in both maths and English, 46 schools took part in maths only and 6 schools took part in English only. This meant there were 121 schools taking part in maths that were randomly assigned to one of the groups (61 intervention and 60 control). Of the 79 schools taking part in English, 43 were assigned to intervention and 36 to control. Overall, there is an imbalance in the group allocation for schools taking part in the English trial. This occured as a result of not correcting the group imbalance that arose at each block. We deliberately adopted this

---

[1] If a school is taking part only in one subject, they will have a representative only from the one department which attends the workshops.

approach to ensure that it was not possible to predict the sequence of group allocation. Table 2 describes the unplanned blocks that were required due to slow recruitment to the trial.

**Table 2: Number and proportion of schools randomised**

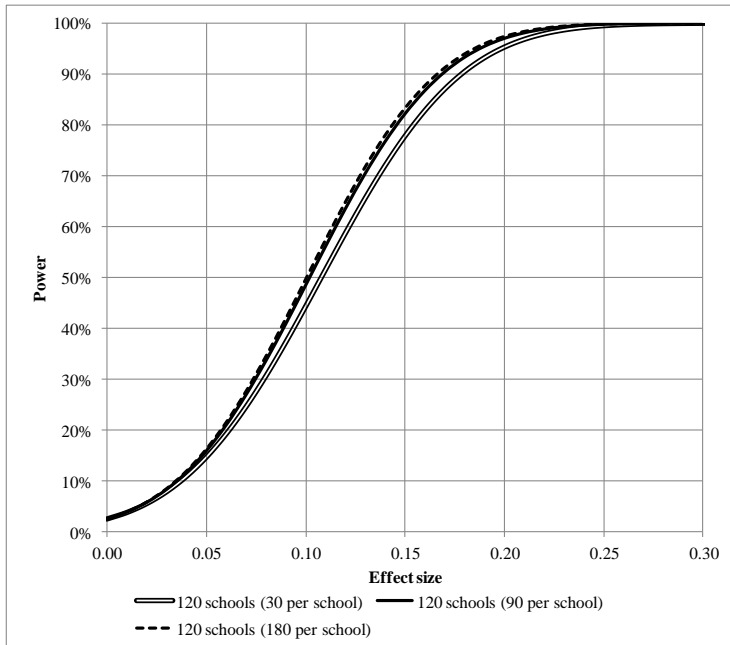|  |  | Maths n (%) | English n (%) |
|---|---|---|---|
| Block 1 | Intervention | 47 (49%) | 30 (52%) |
|  | Control | 49 (51%) | 28 (48%) |
| Block 2 | Intervention | 11 (50%) | 10 (56%) |
|  | Control | 11 (50%) | 8 (45%) |
| Block 3 | Intervention | 1 (100%) | 1 (100%) |
|  | Control | 0 (0%) | 0 (0%) |
| Block 4 | Intervention | 1 (100%) | 1 (100%) |
|  | Control | 0 (0%) | 0 (0%) |
| Block 5 | Intervention | 1 (100%) | 1 (100%) |
|  | Control | 0 (0%) | 0 (0%) |
| Total | Intervention | 61 (50%) | 43 (54%) |
|  | Control | 60 (50%) | 36 (46%) |
|  | **Total** | **121** | **79** |

## Calculation of sample size

A number of within school sample sizes were considered for analysis. In order to reduce testing burden per school without sizeable impact on the power, it was decided that NFER will randomly select 60 pupils from the year 8 school roll from each of the recruited schools. Half of the pupils will sit the mathematics test and half will sit the English test. For schools that were randomised for only one subject, 30 pupils will sit the test in that subject[2]. (Some of these options are illustrated in power curves in figure 1).

The power curves in figure 1 use the following assumptions: intra-cluster correlation of 0.15 (lowered from 0.2 through the use of key stage 2 as a covariate) and correlation between key stage 2 and year 8 test of 0.7. With a statistical power of 80%, the MDES for the maths outcome will be 0.16 for 30 pupils sitting the test in maths. However, the sample of schools for the English outcome measure is slightly smaller and with the same assumptions, the MDES for the English outcome will be 0.19 for 30 pupils sitting the test in English.
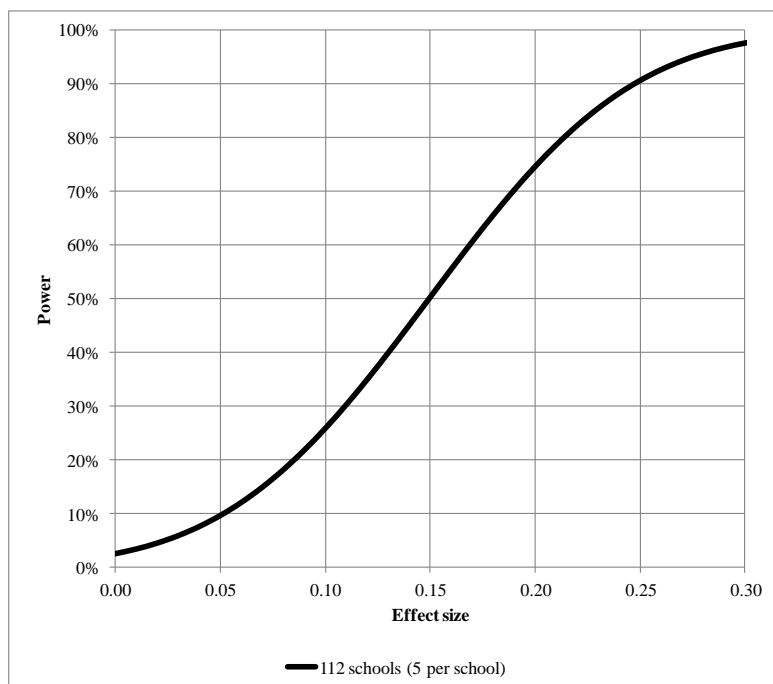
---

[2] Pupils being tested in English will be different from those being tested in maths. Sampled pupils will not be replaced in any case even if a pupil was no longer available for testing.

**Figure 1: Power curves for 60 versus 60 schools (recruitment target as per the protocol)**



These power curves clearly illustrate how testing burden per school can be greatly reduced through within-school sampling with minimal impact on power. For such strategies to work, the within school sampling has to be random to ensure unbiased cluster mean estimates. Within-school sampling has an impact on the power of sub-group analysis. As FSM-eligible pupils represent a particularly important subgroup the power of a separate FSM analysis was considered at this stage. In order to achieve sufficient number of pupils with FSM eligibility, this within-school sampling will be stratified by pupil FSM status. Based on this calculation, we can expect an average of 5 FSM pupils to be sampled in each school cohort and at least one FSM pupil to be sampled in each school in 93% of recruited schools. As we are just estimating regression coefficients some small cluster sizes will not compromise the multi-level models (Snijders et al., 2005). The minimum detectable effect size (MDES) was then calculated for FSM only analysis at 0.22  for 112 schools (93% of 120) and with the same assumptions as above. This is illustrated in the following figure.

**Figure 2: Power curves for FSM analysis (as per the protocol)**



## Follow-up

As stated earlier, pupil data could not be collected prior to randomisation. Pupil data collection started in autumn 2015 and is ongoing at present. Out of 127 schools, 94 BPS schools have provided pupil data. There are 29 schools that did not provide pupil data and have refused to participate in the trial. The attrition rate is different across the randomisation groups and is presented in the following table.

**Table 3: Number and proportion of schools withdrawn from the primary outcome**

| Group allocation | Randomised (n schools) | Withdrawn[3] (n schools) | Withdrawn (%) |
|---|---|---|---|
| Intervention | 65 | 22 | 34 |
| Control | 62 | 7 | 11 |
| Total | 127 | 29 | 23 |

At present, NFER is arranging the end-point test administration. The final sample with follow-up test data will be comprised of all schools that take part in this test administration[4].

## Outcome measures

### Primary outcome

As mentioned previously, some schools took part in the trial for both the subjects and other schools took part only for one subject. This means there are different schools taking part in maths than those taking part in English. Therefore, it is necessary to undertake two separate

---

[3] We couldn't locate a physical copy of the signed MoU (memorandum of understanding) for one school that withdrew participation from the trial. However, the school withdrew from the trial after the knowledge of group allocation, therefore this school will be retained in the missing data analysis.
[4] Note that there are four BPS schools that have not provided pupil data yet. If this is delayed any further, it will not be possible to match FSM status from NPD in time for test administration. Pupils from these schools will be randomly selected without the FSM stratifier.

analyses, one for each subject. It is further decided (jointly by EEF and NFER) that the pupils taking English and maths tests are different and therefore adjustment for multiple testing is not required in this case.

As the best practice intervention is aimed at pupils in years 7 and 8, testing is necessary as there is no statutory assessment in these years. Testing will take place in June-July 2017 at the end of year 8. GL Assessment's Progress in English (PTE13)[5] and Progress in Mathematics (PTM13) tests will be used.

The primary outcome measures of attainment will answer the following research questions:

1. What is the impact of best practice in setting on pupils' attainment in mathematics?
2. What is the impact of best practice in setting on pupils' attainment in English?

As these tests have a broad coverage of the curriculum, we will use the raw total score for each subject that covers all curriculum content. Maths total score (maximum possible score 70) will consist of fluency in facts and procedures, fluency in conceptual understanding, mathematical reasoning and problem solving. English total score (maximum possible score 66) will consist of spelling, grammar and punctuation, reading comprehension: narrative and non-narrative.

NFER will take responsibility for collecting and delivering PTE13 and PTM13 in paper form using its test administrators while the tests will be marked by GL assessment blind to treatment allocation.

## Secondary outcomes- pupil attitudes

As outlined in the protocol, the secondary research questions are:

1. What is the impact of best practice in setting on pupils' self-confidence in mathematics?
2. What is the impact of best practice in setting on pupils' self-confidence in English?

These are being measured by administering a pupil survey at the start of year 7 in September 2015 (baseline survey administered post randomisation) and at the end of year 8 in summer 2017 (follow-up survey). UCL Institute of Education are responsible for administration of these surveys. The surveys will be administered with an entire cohort from participating schools. However, for trial purposes, pupil survey data will be analysed based on the original randomisation group and subject participation. E.g. if a school is taking part in the mathematics only trial, pupils' self-confdience in English from this school will not be considered in the English analysis.

In partnership with Queen's University Belfast, UCL Institute of Education have developed pupil self-confidence measures in maths and English.

Self-confidence measures were developed using factor analysis on selected items from the baseline pupil survey data (combined dataset for both the trials, BPS and BPMA). These items were drawn from several instruments previously used (SDQII from Marsh, 1990; TIMSS questions from IEA, 2011 and PISA questions from OECD, 2012). Please see below table for the list of items included in the principal axis factor analysis.

---

[5] At the time of the protocol, these tests were being developed by GL assessment. They were being called New Progress in English (NPiE) and New Progressm in Mathematics (NPiM). After the development, these tests are called Progress in English and Progressm in Mathematics.

**Table 4: List of items included in the secondary outcome measures of self-confidence**

| Composite measure | Constituent items | Source |
|---|---|---|
| Self-confidence in mathematics | Work in Maths lessons is easy for me | Adapted from Marsh (1990) verbal [sic] self-concept |
| | I am not very good at Maths | Adapted from Marsh (1990) verbal [sic] self-concept |
| | Maths is one of my best subjects | Adapted from Marsh (1990) |
| | I hate maths | Adapted from Marsh (1990) |
| | I do well at maths | Adapted from Marsh (1990) school [sic] self-concept |
| | I get good marks in maths | Adapted from PISA self-concept in mathematics and Marsh (1990) verbal [sic] self-concept |
| | I learn things quickly in maths lessons | Adapted from TIMSS self-confidence in learning mathematics and Marsh (1990) verbal [sic] self-concept |
| Self-confidence in English | Work in English lessons is easy for me | Adapted from Marsh (1990) verbal self-concept |
| | I am not very good at English | Adapted from Marsh (1990) verbal self-concept |
| | English is one of my best subjects | Marsh (1990) verbal self-concept |
| | I hate English | Marsh (1990) verbal self-concept |
| | I do well at English | Adapted from Marsh (1990) school [sic] self-concept |
| | I get good marks in English | Adapted from Marsh (1990) verbal self-concept |
| | I learn things quickly in English lessons | Adapted from Marsh (1990) verbal self-concept |

Subsequently, all the items were retained and the composite measures were created as an average of all constituent items. Self-confidence in maths had an internal reliability (Cronbach's α) of 0.88 and self-confidence in English had an internal reliability (Cronbach's α) of 0.86. Scores for these composite measures will range from 1 to 5 with higher scores reflecting higher self-confidence in the given subject.

## Analysis

The trial analysis will follow EEF Analysis Policy[6].

### Primary intention-to-treat (ITT) analysis

The primary outcome analysis will be 'intention-to-treat'.

There will be two separate analyses- one for each subject. Pupils from schools that took part in a given subject trial (maths trial and English trial) will be included in each analysis irrespective of whether or not the schools implemented the intervention. If a school took part in both the subjects, their pupil data will be included in both the analyses- maths and English.

---

[6]https://educationendowmentfoundation.org.uk/public/files/Evaluation/Writing_a_Research_Report/2015_Analysis_for_EEF_evaluations.pdf

The analyses will determine whether the Best Practice in Setting initiative had an overall effect on Year 8 pupils' mathematics and English attainment. This will be determined by fitting two separate models, one for each subject. Multilevel models with two levels (school and pupil) will be used for the analysis to account for the cluster randomisation.

The mathematics model will include data for all schools that took part in the mathematics trial and the dependent variable for this model will be the raw total score in mathematics for PTM13 with the following covariates:

- an indicator of whether the pupil is in an intervention school
- pupil prior attainment as measured by KS2 Maths point score (KS2_KS2MATPS variable on NPD)
- an indicator of whether the school took part in one subject or both (representing the stratification variable used at randomisation)

The English model will include data for all schools that took part in the English trial and the dependent variable for this will be the raw total score in English for PTE13 with the following covariates:

- an indicator of whether the pupil is in an intervention school
- pupil prior attainment as measured by KS2 English point score (KS2_KS2READPS variable)
- an indicator of whether the school took part in one subject or both (representing the stratification variable used at randomisation)

In addition to the above models, we will also report a point estimate (without a confidence interval) from similar models which don't include the stratification variable. This will be reported for the purposes of cross-study comparisons.

### Imbalance at baseline for analysed groups

Although we expect no systematic bias to have arisen from randomisation, it will be important to examine bias due to high attrition. In the absence of named pupil data from schools that withdrew participation from the trial, we will use de-identified NPD for these schools to examine imbalance in the samples using background characteristics such as pupil FSM status and prior attainment at key stage 2. We will use multilevel modelling to examine imbalance for prior attainment.

### Missing data

We will run two multilevel logistic models (one for each subject) with two levels (school and pupil) on whether or not a pupil is missing at follow-up, regressed on the covariates of the main model.

Since there are many schools that withdrew participation from the trial and the primary outcome measurement, it is important to explore the level of missing data and the extent of bias. As we are unlikely to be in a situation where a school has complete follow-up data and missing baseline, multiple imputation may not be useful. Instead, under the 'missing at random' given baseline assumption, we would expect a completers analysis to be unbiased. However, since we already know that the extent of school dropout was unequal between randomised groups, the 'missing not at random' assumption is likely to hold and we will need to conduct sensitivity analyses. This will be done by initially running multilevel multiple imputation and then extending this model using a weighting approach according to

Carpenter et al. (2007). This approach works by replacing a simple average by a weighted average where estimates from the imputations that are more likely under 'missing not at random' are upweighted relative to the others.

After adjusting for the observed variables, the chance of observing the outcome measure per unit change in that measure has log-odds ratio of δ. If data are 'missing at random', δ will be zero. If δ is positive, the chance of observing the outcome measure is higher for higher values of the outcome measure. Thus after imputation under the 'missing at random' assumption, there is under representation of imputation with small values of the outcome measure. The weights correct this by upweighting the estimates from the imputed data. We will use values of δ between -0.5 and 0.5 as indicated by the literature as being suitable.

## Secondary outcome analyses- self-confidence models

Models for the secondary outcomes of self-confidence will be run similar to the primary outcomes of attainment. There will be two dependant variables in two separate multilevel models- each including pupils from schools that took part in the relevant subject trial . The covariates for these models will be similar to the secondary attainment models wherein pupil self-confidence measures in given subject at baseline will be one of the covariates instead of prior attainment measures.

It is anticipated that the models of pupil self-confidence in mathematics and English will encounter some attrition. Therefore, it will be important to determine the extent of bias. If there is more than 5% data missing at baseline where we have outcome measures at follow-up, multilevel multiple imputation will be used to impute the missing values at baseline. As a check for missing at random assumption, an imputation model will be run for both baseline and follow-up. This model will be compared with the completers model.

Data manipulation will be carried out in SPSS while the multilevel models will be run in R package nlme and imputation macros available from missingdata.org.uk.

## Non-compliance with intervention

Fidelity analysis will be carried out on the primary outcome measure only. The developer collected data on the level of school engagement throughout the two-year delivery period using a number of pre-defined variables as described in below table 4 . They sent us data on these individual variables and we will summarise them according to pre-agreed categorisation. This categorisation will yield three measures that are listed in the following table. Measures for English are provided below as an example. Similar information was collected for mathematics.

**Table 5: Dosage variables for primary and secondary outcomes on English:**

| Combined measure | Variable | Level of measurement |
|---|---|---|
| 1. Effectiveness of training practices | 1. English department represented at each training session | Binary. Did the expected number and type of staff attend each session? |
| | | 0 = No |
| | | 1= Yes |
| | | Binary. Has some form of cascading/internal training taken place? |

| | | |
|---|---|---|
| | 2. Training is cascaded to members of the English department | 0= No |
| | | 1= Yes (if one or more departmental members concur) |
| | 3. Setting arrangements follow BP principles (only 3-4 sets) | Binary. |
| | | 0 = 5 or more sets |
| | | 1= 3/4 sets or fewer |
| 2. Effectiveness of setting/allocation practices | 4. Teachers are randomly allocated to classes | Were BP principles followed? |
| | | 0 = No, teachers not randomised or allocated with reference to BP principles |
| | | 1= Partial, teachers allocated with reference to BP principles |
| | | 2= Yes, teachers randomised to classes |
| | 5. Students are allocated to classes according to KS2 results | Binary – 95% or more of students are allocated on the basis of KS2 results |
| | | 0= No |
| | | 1= Yes |
| | 6. Students are re-set no more than three times in two years. | Binary – 95% or more of students are re-set no more than three times in two years (i.e. 95% or more of students are re-set 'at the most' three times in two years. Schools that re-set students more than three times in two years will get a value of 0) |
| | | 0= No |
| | | 1= Yes |
| 3. High expectations | 7. Teachers have high expectations for all students | Binary |
| | | 0= No |
| | | 1= Yes |

In order to obtain a more accurate measure of the 'pure' dosage effect of the intervention on pupil attainment and self-confidence, the Complier Average Causal Effect (CACE) impact estimate will be calculated. Measure 2 from the above table (effectiveness of setting/allocation practices) will be used for this purpose. Because schools may potentially have unobserved characteristics that have an influence on both compliance with the trial and academic attainment a two stage least squares model will be used to calculate the CACE estimate (Angrist and Imbens, 1995).

The first stage of the model will be engagement level regressed on all covariates that are used in the main primary outcome model and in addition will include, as an instrumental variable, a binary variable that indicates a pupil's pre-intervention treatment allocation. The second stage of the model will regress the primary and the secondary outcomes on the covariates used in the main model and will also include a covariate representing the pupil's estimated dosage level from the first stage of the model and an interaction term between the estimated dosage and the pupil's pre-intervention treatment allocation. The coefficient of the

interaction term is the CACE estimate of the dosage effect. In the event that there are no confounding factors affecting compliance and attainment the CACE estimate will be equal to the intention-to-treat estimate.

A further factor that must be taken into account is the hierarchical nature of the data. To ensure that this factor of the data is accounted for correctly the R package ivpack, which has the functionality to correctly handle hierarchical data when using instrumental variables, will be used to perform the CACE analysis.

### Subgroup analyses

Sub-group analyses on the primary outcomes will be carried out as per the protocol and the most recent EEF analysis guidelines. As per the protocol, we will explore the differential effect for different pupil ability levels. An interaction term will be added to the main models. The intervention indicator will be interacted with pupil ability as measured by above mentioned prior attainment measures at KS2.

As per the EEF guidance, there will also be another interaction model of whether a pupil has ever received free school meals (as measured by EVERFSM_6 variable from the Autumn School census 2015/16). This will be done using a model identical to the primary outcome model but including EVERFSM_6 and EVERFSM_6 interacted with the intervention indicator as covariates. Analysis shall proceed as per the original primary outcome modelling i.e. the first model shall be identical to the primary outcome model but with EVERFSM_6 as a covariate.

A separate analysis of FSM only pupils will also be carried out as per the EEF analysis guidance. These models will be similar to the main models of overall effect but will only include pupils who were eligible for FSM as measured by EVERFSM_6 variable.

### Effect size calculation

The numerator for the effect size calculation will be the coefficient of the intervention group from the multilevel model. All effect sizes will be calculated using total variance from a multilevel model, without covariates, as the denominator i.e. equivalent to Hedges' g. Confidence intervals for each effect size will be derived by multiplying the standard error of the intervention group model coefficient by 1.96. These will be converted to effect size confidence intervals using the same formula as the effect size itself. Parallel to this, we will also use the R package for analysing education trials (eefanalytics) which employs a slight correction to the calculation of effect size as referenced in the analysis guidance.

## Report tables

All the tables will be structured according to the EEF trial report template[7].

## References

Angrist, J.D. and Imbens, G.W. (1995). 'Two-stage least squares estimation of average causal effects in models with variable treatment intensity'. *Journal of the American Statistical Association.* 90**,** 430, 431-442. [online]. Available: http://amstat.tandfonline.com/doi/abs/10.1080/01621459.1995.10476535 [19 May, 2017].

---

[7] https://educationendowmentfoundation.org.uk/evaluation/resources-centre/writing-a-research-report/

Borenstein M, Hedges LV, Higgins JPT and Rothstein HR. (2009) Introduction to Meta-Analysis. Wiley.

Carpenter, J.R., Kenward, M.G. and White, I.R. (2007) Sensitivity analysis after multiple imputation under missing at random – a weighting approach. Statistical Methods in Medical Research 16:259-275.

IEA (2011). *TIMSS 2011 Student Questionnaire*. IEA: Boston.

Marsh, H. W. (1990). The Self-Description Questionnaire II *Manual*. Australia: University of Western Sydney.

OECD (2012). *PISA 2012 Student Questionnaire*. OECD Publishing, Paris.

Snijders, T. A. B. (2005). 'Power and Sample Size in Multilevel Linear Models'. In: Everitt, B. S. and Howell, D. C. (Eds.)  Encyclopedia of Statistics in Behavioral Science. 3, 1570–1573. Chicester (etc.): Wiley, 2005 [online]. Available: https://pdfs.semanticscholar.org/a769/1eb67c5806e154da58a74f7b1a1bc9ccb58a.pdf [19 May, 2017].