

Ground your practice in evidence, because thinking matters: A national evaluation of the impact of the Thinking Schools approach on the achievement of primary and secondary age pupils in England (2016, 2017).

Dr Dave Walters – Honorary Research Fellow, Exeter University Graduate School of Education.

December 2017.

Abstract

The main purpose of this paper is to provide a lens through which the impact of taking a whole school approach to the development and embedding of cognitive education can be evaluated. Specifically operationalized, this takes the form of examining the progress outcomes for 2016 and 2017 of primary and secondary age pupils in England from schools who have (i) adopted a whole school approach to the teaching of thinking and have subsequently been accredited as a Thinking School, and (ii) schools who have either been accredited as a Thinking School *or* who have registered and started the Thinking Schools journey. Using the widely acknowledged impact measure of *effect size*, these progress outcomes are set against key benchmarks so that a relative measure of impact can be presented. The process follows the author's preferred style of action research based on a form of illuminative evaluation and his adaptation / application of the SPARE 'Wheel' model initially developed by Burden (1998). 'Very High' impact (equivalent to over a whole grade extra growth) is shown for both 2016 and 2017 in the secondary progress measure of Progress 8 (P8). 'Moderate' to 'High' impact (equivalent to 0.5 – 0.6 of a grade extra growth at GCSE) is shown for both 2016 and 2017 in the primary progress measures relating to reading, writing and mathematics. In addition, 2017 shows marked impact improvements for secondary P8 and primary reading progress. Further, marked improvement is evident in all progress measures, including the primary progress measure relating to writing, as schools move to successful accreditation. In conclusion, a shift towards evidence based practice, for developing a 'systems' approach to organisational growth using the Thinking Schools accreditation criteria as a blueprint, is proposed as a vehicle for a self-improving system to more fully develop the learning potential of all children. This is recommended in addition to the more general support provided by Thinking Matters and 'hub' Thinking Schools.

Keywords: cognitive education; evaluative research; effect size; Thinking Matters; Grounded Practice; self-improving system; Thinking Schools

Introduction

As the number of schools undertaking the Thinking Schools journey increases, opportunities to research and evaluate are also developing. Through the Thinking Schools accreditation already undertaken, there is much anecdotal evidence demonstrating the benefits of taking a whole school approach to the teaching of thinking. These benefits are detailed in the evaluation reports of accredited Thinking Schools across the globe and also in Ofsted / Estyn inspection reports in the UK. The growth of the Thinking Matters training support across the globe, for schools aspiring to become accredited as Thinking Schools by Exeter University's Cognitive Education Development Unit, makes this school development journey a truly international venture. Schools and school systems in Ethiopia, Lithuania, Malaysia, Norway, South Africa, New Zealand, the USA, Ireland, Egypt, Nigeria, Dubai, India, Thailand, Malaysia, Australia, and the UK have all taken up the challenge of developing a cognitive approach to education.

In order to add to the evidence base of impact, this report attempts to focus more robustly on the impact of whole school approaches to the teaching of thinking on the academic achievement of pupils in UK state primary and secondary schools. Specifically, this evaluation focuses on the progress outcome measures of reading, writing and mathematics at the end of the primary phase, and the 'new' secondary progress outcome of progress 8 (P8) at the end of the secondary phase (a progress measure encompassing 8 subjects with English and Mathematics being 'weighted' more strongly in the overall collective measure). A more detailed rationale for this can be found in the Setting and Plan sections of this report. Further, it is the authors view that attainment is a measure with no origin or starting point from which true orientation of destination can be judged and so measures of progress add this missing dimension. However, although what may be described as 'narrow' measures of achievement have been used (P8, Progress in Reading, Writing and Mathematics), it would be unwise to fight these accountability measures in the short-term due to the high status they have in the UK. Therefore, whilst being appreciative of other indicators of achievement (some attitudinal and dispositional), and that 'excellence' may be open to multiple interpretations that go beyond the academic achievement measures presented here, this evaluation plays the high profile accountability 'ball' straight back and uses the key measures that UK schools are held to account over.

The 'dreaded' questions for any innovation to field are, as Ellis (2010) coins, the 'So what?' and 'So, has it had an impact?' questions. Randomised controlled trials (RCTs) have often been used to provide definitive scientifically acceptable results in order to answer these 'dreaded' questions. However, this evaluation takes the stance RCTs, however robust their application, only provide a limited insight to those interested in 'real-world' applications. This evaluation takes an alternative approach based on the notion that innovative projects in 'real-world' contexts require a multi-dimensional approach more attuned to a socio-cultural perspective where social and historical influences are acknowledged. Specifically, Burden's (1998) SPARE

'wheel' model of 'Illuminative Evaluation' is used as a working application of evaluative research that places measurement in the context of the social, historical and cultural influences of the time. Thus, the SPARE 'wheel' model takes the Setting, the Plans, the Actions and Reactions of participants into account, as well as any 'robust' quantitative measure, as part of the Evaluation process. Although presented in a rather linear format, it has to be stressed that the process follows a *cycle* of enquiry. What follows may be better described as being a holistic narrative of how the author addressed the dreaded 'So what?' and 'So, has it had an impact?' questions.

Setting

The setting is the growing consensus amongst educators that there is a great need to place cognitive education at the heart of pedagogical innovation and reform (Burden and Nichols, 2000). The need for schools to take a whole school approach to the teaching of thinking is no longer seriously in dispute. The extensive work of Hattie (2009) and Higgins *et al.* (2013) both draw together data through meta-analyses that clearly place meta-cognitive strategies (cognitive education) and feedback (assessment) at the top of 'high-impact', 'low-cost', 'based on extensive evidence' factors that underpin high achievement for children in schools. Whilst Assessment for Learning (AfL) has long been acknowledged as a highly effective approach for teachers (Stobart, 2006; Clarke, 2005), the development of the Thinking Schools educational approach, devised by Professor Bob Burden and Thinking Matters (TM), has also gained wide acceptance and acknowledgment as impacting positively on the achievement and cognitive growth of children. The approach allows schools the opportunity to benefit from TM's training in the use of cognitive tools, together with the Cognitive Education Development Unit evaluation/accreditation as a 'Thinking School'. As the number of schools undertaking the Thinking Schools journey has increased, much *anecdotal* evidence demonstrating the benefits have been detailed in the evaluation reports of accredited Thinking Schools and also in Ofsted / Estyn inspection reports.

Set against this backdrop of evidence supporting the impact of cognitive education, are three influential contextual factors relating to an evaluation of this kind:

- The increasing use of effect sizes to quantify the effectiveness of a particular innovation or intervention.
- The move in England from using measures of attainment as key student outcome data for accountability purposes to measures of progress / Value-Added (VA).
- The crackdown on schools in England who are perceived to be 'gaming the system' to gain an advantage in league table accountability measures, in response to the second bullet point above.

Each of these factors merits consideration here, to set the scene for how this evaluation evolved.

The increased use of effect sizes

The use of effect sizes has moved steadily from the routine use for meta-analysis (combining and comparing estimates from different studies) to gradually becoming a simple way of quantifying the difference between two groups (Coe, 2002; Higgins *et al.* 2013). It is a way of measuring the *extent* of the difference between two groups and allows an evaluation of impact to move beyond the simplistic ‘Does it work?’ to the more valuable insight of ‘How *well* does it work across a *range* of contexts?’. As Higgins *et al.* point out: ‘For these reasons, effect size is the most important tool in reporting and interpreting effectiveness, particularly when drawing comparisons about relative effectiveness of different approaches’ (2013, 6). The technical aspects relating to the use of effect sizes in this evaluation appear in the subsequent sections relating to ‘Plan’ and ‘Action’.

The changes in school accountability measures in England

Recent changes in school accountability measures first came to prominence in the DfE performance tables for 2016. Rather than a focus on attainment, the key pupil outcome measure has become progress / Value-Added. Specifically in primary schools this takes the form of:

- Average progress in mathematics
- Average progress in reading; and
- Average progress in writing.

(see DfE 2016)

Progress scores for primary schools are centred around 0, with most schools in the range of -5 to +5. A score of 0 means pupils in this school on average do about as well at KS2 (the end of the primary phase – age 11) as those with similar prior attainment nationally. A positive score means pupils in this school do better at KS2 as those with similar prior attainment nationally. A negative score means pupils in this school on average do worse at KS2 as those with similar prior attainment nationally.

That said , when we look at the more complicated progress measure used for secondary schools, ‘Progress 8’ (P8), it is the primary school progress measures of mathematics and

reading that come to dominate as these measures form the baseline from which progress in secondary schools is judged (see DfE 2017). Progress 8 aims to capture the progress of a pupil from the end of primary school to the end of secondary school. It is a type of value-added measure, where pupils' results are compared to the actual achievements of pupils with the same prior attainment (average of **mathematics** and **reading** scores at the end of KS2). P8 is based on a calculation of pupils' performance across 8 qualifications. These qualifications are:

- A double weighted **mathematics** element that will contain the point score of the pupil's English Baccalaureate (EBacc) mathematics qualification.
- An English element based on the highest point score in a pupil's EBacc **English language** or **English literature** qualification. This will be double weighted provided a pupil has taken both qualifications.
- An element which includes the three highest point scores from any of the **EBacc** qualifications in science subjects, computer science, history, geography, and languages. The qualifications can count in any combination and there is no requirement to take qualifications in each of the 'pillars' of the EBacc.
- The remaining element contains the three highest point scores in any three **other** subjects, including English language or literature (if not counted in the English slot), further GCSE qualifications (including EBacc subjects) or any other technical awards from the DfE approved list.

If a pupil has not taken the maximum number of qualifications that count in each group then they will receive a point score of zero where a slot is empty. No legacy GCSEs (A*-G), International GCSEs or level 1/level 2 certificates in these subjects will count in performance tables once new GCSEs (9-1) in that subject are introduced. A score of zero means pupils in this school on average do about as well at key stage 4 (end of the secondary education phase – age 16) as other pupils across England who got similar results at the end of key stage 2. A score above zero means pupils made more progress, on average, than pupils across England who got similar results at the end of key stage 2. A score below zero means pupils made less progress, on average, than pupils across England who got similar results at the end of key stage 2. **Given the scope and comprehensive nature of P8, it is clearly a more encompassing measure than the legacy measure of attainment, 5A*-C including English and mathematics.**

'Gaming the system'

In a letter to inspectors, Harford (2017), Ofsted's national director of education (UK), draws the attention of inspectors to the importance of following up unusual examination entry patterns, what is frequently termed as 'gaming the system'. Unfortunately, when an accountability measure is introduced (P8 for example) there is a risk that 'game theory' emerges (see Waters, 2013). Harford (2017) raises particular concerns in relation to what he views as increasingly common practices, namely:

- Schools which enter large numbers of pupils for qualifications that are not core subjects or do not reflect a school’s specialisms – often subjects of a technical or vocational nature not suited to the majority of pupils.
- Double entry in qualifications that overlap in content. For example, statistics and free-standing mathematics qualifications; GCSE English and IGCSE English as a second language qualification for pupils who have English as a first language.
- Schools which enter pupils for GCSEs in English literature, without teaching the latter properly. Pupils sit the exam purely to ensure that the language result is counted doubly towards P8.
- Moving underperforming pupils into alternative provision so that they would not bring down results – a practice known as ‘off-rolling’.

It appears that the motive for some schools is to boost their league table position rather than act in the best interests of the children. Clearly, there is a need to preserve **authentic education** for all our children. This authenticity lies at the heart of the values of ‘Thinking Schools’.

The ‘setting’ or context, outlined in this section, provided the springboard for the initial ‘plan’. It is this ‘plan’ that features next.

Plan

Given the prominent use of effect sizes to judge the impact of educational innovations, the author decided to use this metric so that a degree of consistency was retained such that the results of the evaluation could be compared to research outcomes already available, particularly the findings of Hattie (2009) and Higgins (2013). Thus, the original calculation of effect size followed that of Hattie (2009) and is illustrated below:

$$\text{Effect size} = \frac{[\text{Mean of experimental group}] - [\text{Mean of control group}]}{\text{Average spread (standard deviation, or sd)}}$$

(The experimental group being accredited state ‘Thinking Schools’ in England and the control group being other schools in England)

English state schools were chosen, as opposed to non-state schools (or independent schools), to retain consistency in terms of statutory curriculum expectations. Due to the changing nature of accountability measures in England, it was decided that this metric be applied to progress scores in mathematics, reading and writing for primary schools, and progress 8 (P8) scores for secondary schools. Primary reading scores were chosen alongside mathematics progress scores a main focus due to these outcomes being used as the baseline from which P8 is calculated – again the idea was to retain consistency and also the ability to analyse how the primary measures might influence P8 going forward. On that basis, the original plan was to conduct three separate effect size calculations to provide four measures of impact as follows:

- Impact on mathematics progress in primary schools (key focus)
- Impact on reading progress in primary schools (key focus)
- Impact on writing progress in primary schools (to provide a supplementary measure for comparative purposes)
- Impact on progress 8 (P8) scores in secondary schools (key focus)

In addition to using the common metric of effect size for this summative evaluation, the plan was to provide, through this evaluation, guidance to schools as to how they might use effect sizes formatively in an on-going way in order judge their own impact. The idea here was to provide schools with a robust mechanism by which they could advance pupils' learning and in doing so undertake authentic school improvement without having to resort to 'gaming the system'.

Action

Although the essence of the original plan was actually applied, three refinements were made. The first two refinements related to the actual effect size metric. Although using the pooled standard deviation (SD) to calculate the effect size generally provides a better estimate than the control group SD, Coe (2002) points to the issue of bias in that it generally gives a value slightly larger than the true population value. Hedges and Olkin (1985) give a formula which provides an approximate correction to this bias which was subsequently applied to the effect size formula used in the evaluation. This correction had the added benefit of allowing the calculated effect sizes to be compared to The Sutton Trust research (Higgins *et al.*, 2013) on an 'even', 'like for like' basis, as this research also corrected for this bias. In addition, due to the increased rigour allowed by this correction, high impact results could be viewed with

increased status as results not corrected for this bias would tend to yield higher scores. The second refinement, also relating to the effect size metric, came in response to Coe's (2002) recommendation that effect sizes should be calculated and reported with **confidence intervals**. The confidence interval for effect size is a measure of the significance of the effect size taking into account the spread of the data and also the number of observations. A confidence interval that doesn't include zero indicates that there is a significant difference between the two groups. Another way of looking at this issue is that statistical significance does *not* tell you the key feature: the size of the effect. To overcome this effect sizes are reported together with an estimate of its likely 'margin for error' or 'confidence interval'. In keeping with common research protocols (Higgins *et al.*, 2013), a 95% confidence interval was calculated and reported alongside the effect size. A 95% confidence interval also features in England's Department for Education (DfE) performance tables for both primary and secondary school progress data. Confidence intervals are presented as two numbers – the lower and upper limits within which we are 95% confident the true score may lie. This is a loose interpretation, but is useful as a rough guide. The strictly-correct interpretation of a confidence interval is based on the hypothetical notion of considering the results that would be generated if the study were repeated many times. If a study were repeated infinitely often, and each time a 95% confidence interval calculated, then 95% of these intervals would contain the true effect.

The third refinement to the original plan related to increasing the original 'pool' of schools to include schools which had 'registered' with TM as embarking on the Thinking Schools journey. Thus two separate illustrations of impact were selected namely:

- Accredited Thinking Schools
- All registered Thinking Schools (including accredited Thinking Schools)

The idea behind this extra dimension was two-fold. Firstly, it allowed the impact of taking a whole school approach to cognitive education to be evaluated more widely by including schools which were well on the journey (accredited), together with schools who had committed to developing their practice in this way but who had yet to fully embed the innovation. Secondly, by adding this dimension the 'progress' and 'gain' of moving to full accreditation could also be judged by comparing all registered Thinking Schools with those which had successfully become accredited. In short, it would provide a 'proxy' measure of the impact of pursuing accreditation.

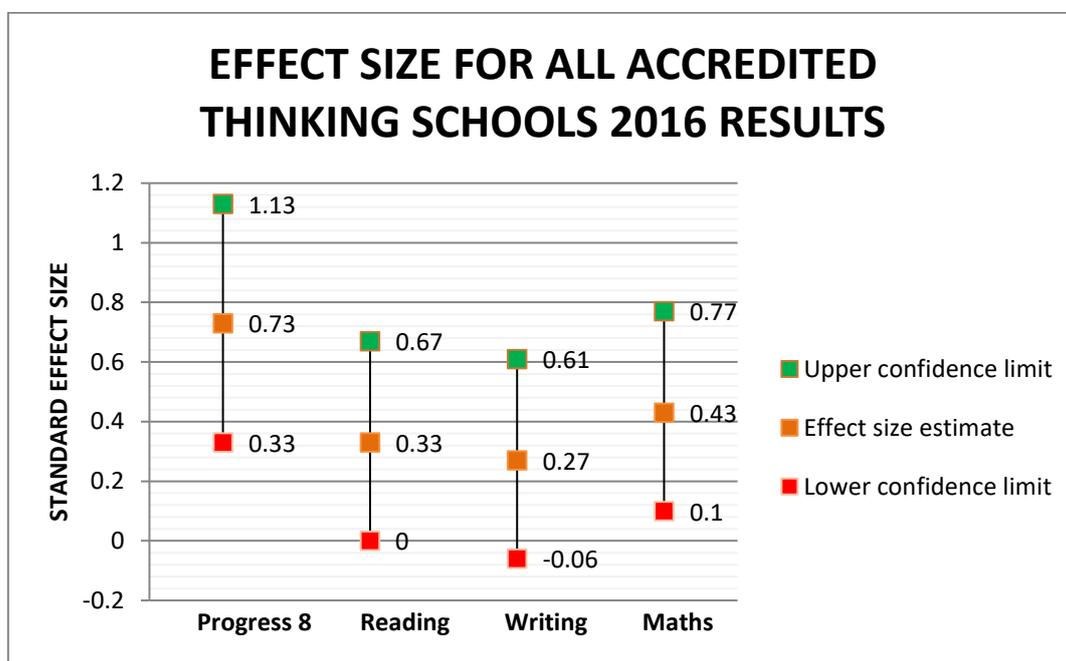
What follows next are the results (or reaction) together with guidance in the form of a lens through which to view the outcomes.

Results / Reaction

Table 1 presents the effect sizes for accredited primary and secondary Thinking Schools (2016) and figure 1 represents the effect sizes graphically.

Table 1: Effect sizes for all accredited Thinking Schools (2016) - (n = 34 primary, n = 24 secondary)

| STANDARDISED EFFECT SIZE | Progress 8 | Reading progress | Writing progress | Maths progress |
|-----------------------------|-------------|------------------|------------------|----------------|
| Upper confidence limit | 1.13 | 0.67 | 0.61 | 0.77 |
| Effect size estimate | 0.73 | 0.33 | 0.27 | 0.43 |
| Lower confidence limit | 0.33 | 0 | -0.06 | 0.1 |



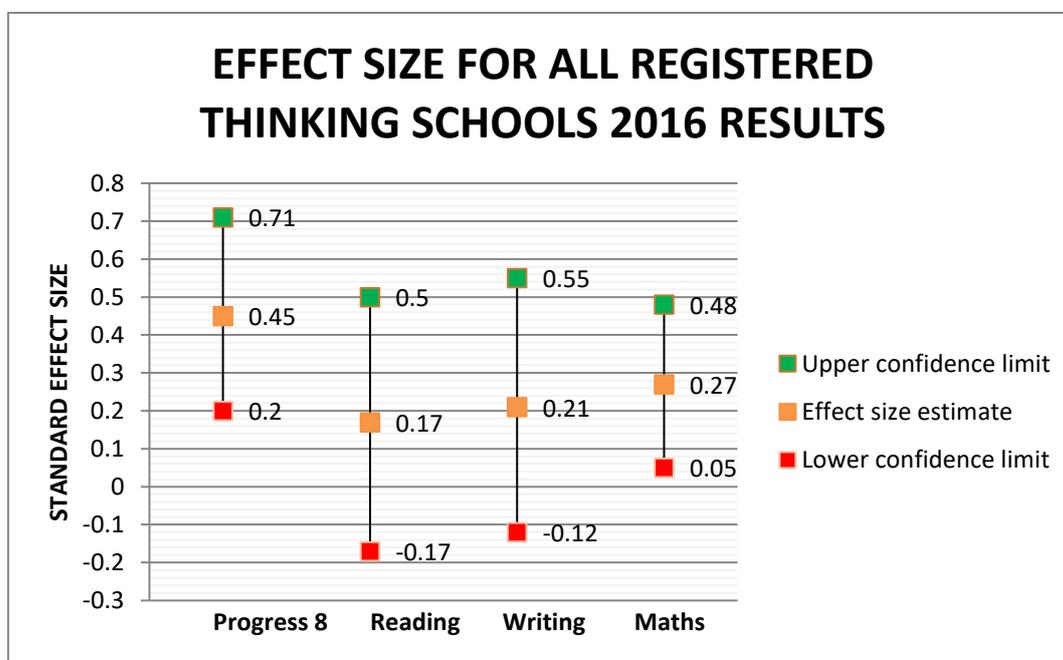
(Average effect size across all measures = 0.44)

Figure 1. Effect sizes for all accredited Thinking Schools (2016).

Table 2 presents the effect sizes for all *registered* primary and secondary Thinking Schools (2016) and figure 2 represents this graphically.

Table 2: Effect sizes for all registered primary and secondary Thinking Schools (2016) - (n = 86 primary, n = 59 secondary)

| STANDARDISED EFFECT SIZE | Progress 8 | Reading progress | Writing progress | Maths progress |
|-----------------------------|-------------|------------------|------------------|----------------|
| Upper confidence limit | 0.71 | 0.5 | 0.55 | 0.48 |
| Effect size estimate | 0.45 | 0.17 | 0.21 | 0.27 |
| Lower confidence limit | 0.2 | -0.17 | -0.12 | 0.05 |



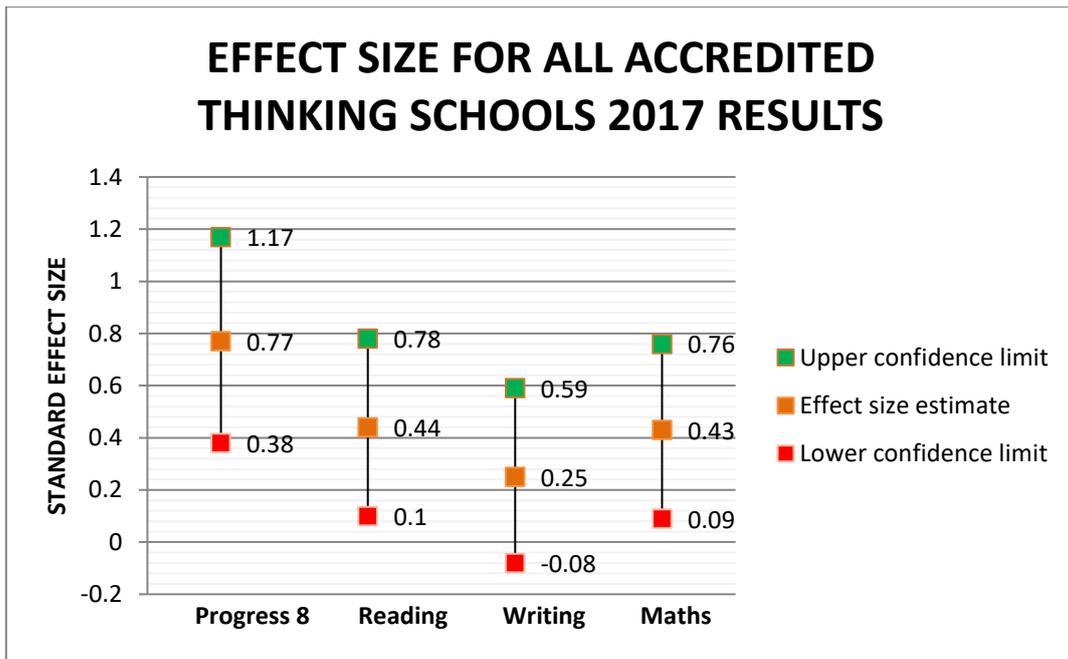
(Average effect size across all measures = 0.28)

Figure 2. Effect sizes for all registered Thinking Schools (2016).

Table 3 presents the effect sizes for accredited primary and secondary Thinking Schools (2017) and figure 3 represents the effect sizes graphically.

Table 3: Effect sizes for all accredited Thinking Schools (2017) - (n = 34 primary, n = 25 secondary)

| STANDARDISED EFFECT SIZE | Progress 8 | Reading progress | Writing progress | Maths progress |
|-----------------------------|-------------|------------------|------------------|----------------|
| Upper confidence limit | 1.17 | 0.78 | 0.59 | 0.76 |
| Effect size estimate | 0.77 | 0.44 | 0.25 | 0.43 |
| Lower confidence limit | 0.38 | 0.1 | -0.08 | 0.09 |



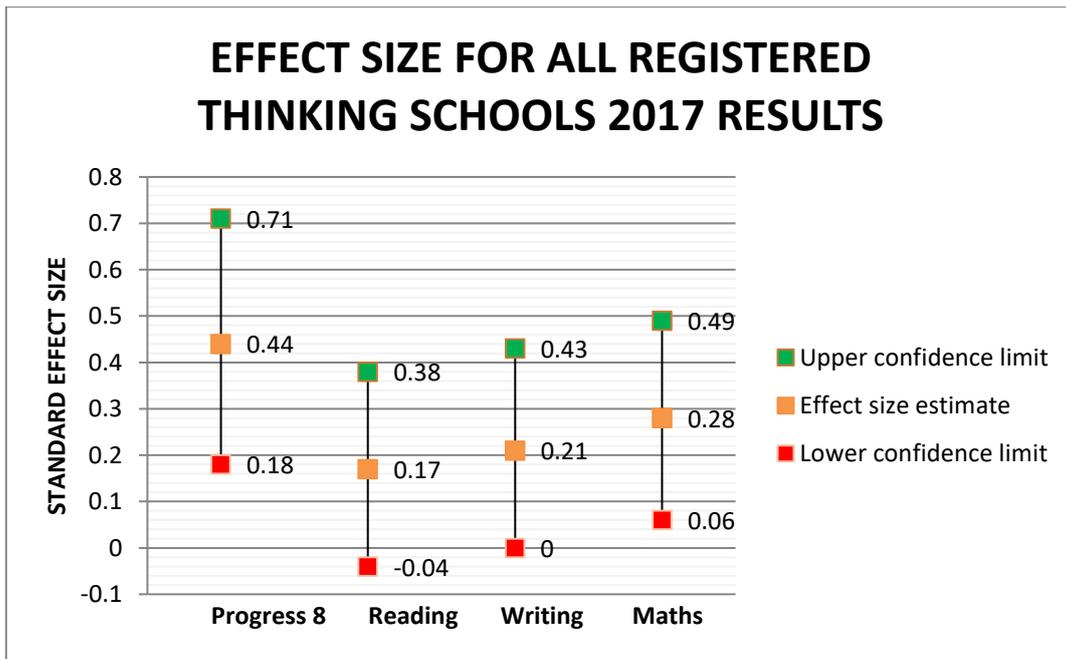
(Average effect size across all measures = 0.47)

Figure 3. Effect sizes for all accredited Thinking Schools (2017).

Table 4 presents the effect sizes for all *registered* primary and secondary Thinking Schools (2017) and figure 4 represents this graphically.

Table 4: Effect sizes for all registered primary and secondary Thinking Schools (2017) – (n = 86 primary, n = 55 secondary)

| STANDARDISED EFFECT SIZE | Progress 8 | Reading progress | Writing progress | Maths progress |
|-----------------------------|-------------|------------------|------------------|----------------|
| Upper confidence limit | 0.71 | 0.38 | 0.43 | 0.49 |
| Effect size estimate | 0.44 | 0.17 | 0.21 | 0.28 |
| Lower confidence limit | 0.18 | -0.04 | 0 | 0.06 |



(Average effect size across all measures = 0.28)

Figure 4. Effect sizes for all registered Thinking Schools (2017).

Before analysing and interpreting these data, it is important to understand some key benchmarks for effect sizes. What follows, is a comprehensive lens through which to view the effect sizes in this evaluation based on these established benchmarks.

Effect Size Comparison / Interpretation Data

So, we have a calculated effect size. How should we interpret this? To come up with an 'objective' and widely acknowledged benchmark we need to use three main considerations:

- When we look at many major longitudinal databases – the Progress in International Literacy Study (PIRLS); the Program for International Student Assessment (PISA); the Trends in International Mathematics and Science Study (TIMSS); the National Assessment of Educational Progress (NAEP); the National Assessment Programme – Literacy and Numeracy (NAPLAN) – they all lead to a similar estimate of an effect size of **0.40** for a year's input of schooling. For example, using NAPLAN (Australia's national assessments) reading, writing, and maths data for students moving from one year to the next, the average effect size across all students is **0.40**.

- The average of 900+ meta-analyses carried out by Professor John Hattie (2009), based on over 250 million students show an average intervention effect size of **0.40**.
- The Sutton Trust Education Endowment Foundation also indicates the overall distribution of effects found in education research with an average around **0.40** (see Higgins *et al.*, 2013).

Therefore, an effect greater than **0.40** is seen as above the norm, and indicates a more than anticipated impact. In other words, the innovation is working better than expected.

Further, Cohen (1988, p. 25) describes an effect size of **0.29** as not being perceptible to the naked eye and equal to the difference between the height of a 5' 11" and a 6' 0" person. In addition, Cohen (1988, p. 26) describes an effect size of 0.50 as being perceptible to the naked eye and therefore 'medium'. Cohen (1988, p. 27) goes on to describe an effect size of **0.8** as 'grossly perceptible and therefore large', equating it to the difference between the heights of 13 year old and 18 year old girls.

Converting Effect Size to Grade Growth

If we then consider how effect sizes relate to actual improvements in examination grades, the following interpretation may prove useful by way of an illustration. If we take Coe's (2002) analysis that the distribution of GCSE grades in compulsory subjects (ie. English and Mathematics) have standard deviations of between 1.5 – 1.8 grades, so an improvement of one GCSE grade represents an effect size of 0.5 – 0.7. In a secondary school therefore, introducing an innovation whose effect size was known to be 0.6 would be likely to yield an improvement of about a grade for each pupil in each subject. For a school in which 50% of pupils were previously gaining a grade 5 or more in English and Mathematics (or indeed overall), this percentage (other factors being equal, and assuming that the effect applied equally across the range of subjects offered) would rise to 73%! For a school with a prior progress 8 score of zero (the national average – pupils grades were in line with national grades for pupils with the same starting points), progress 8 would move up to 1. This would place the school in the 'Outstanding' inspection category in England for this significant measure. For a school with a prior progress 8 score of -0.5 (pupils grades being half a grade lower than national grades for pupils with the same starting points) and in an 'Inadequate' inspection category in England, progress 8 would move to 0.5 (pupils grades being half a grade higher than national grades for pupils with the same starting points) and place the school in the 'Good' inspection category in England. Further, for a school with a prior progress 8 score of 0.5 (pupils grades being half a grade higher than national grades for pupils with the same starting points) and in a 'Good' inspection category in England, progress 8 would move to 1.5 (pupils grades being a grade and a half higher than national grades for pupils with the same

starting points) and place the school in a notionally 'Exceptional' inspection category in England. This pattern of increases in attainment and progress would apply for primary assessment measures also as well as international student outcome measures in compulsory subjects or subjects with large percentage entries.

Effect Size of Metacognitive Strategies

- Hattie (2009) – Effect size = **0.69**
- Higgins *et al.* (2013) - Sutton Trust EEF Teacher Toolkit – Effect size = **0.62 to 0.69**

Converting Effect Size to Months' Development

| Months' Development | Effect Size from... | ...to | Description |
|---------------------|---------------------|-------|-----------------------|
| 0 | -0.01 | 0.01 | Very low or no effect |
| 1 | 0.02 | 0.09 | Low |
| 2 | 0.10 | 0.18 | Low |
| 3 | 0.19 | 0.26 | Moderate |
| 4 | 0.27 | 0.35 | Moderate |
| 5 | 0.36 | 0.44 | Moderate |
| 6 | 0.45 | 0.52 | High |
| 7 | 0.53 | 0.61 | High |
| 8 | 0.62 | 0.69 | High |
| 9 | 0.70 | 0.78 | Very high |
| 10 | 0.79 | 0.87 | Very high |
| 11 | 0.88 | 0.95 | Very high |
| 12 | 0.96 | >1.0 | Very high |

(From Higgins *et al.*, 2013, p. 8, reproduced with kind permission)

And so how might we interpret the results of this study given the aforementioned lenses of interpretation?

So, to interpret the calculated effect sizes for this evaluation the following key features are offered. Tables 1 and 3, together with Figures 1 and 3, clearly indicate what would be termed 'Very High' impact (over a whole grade extra growth) relating to P8 in both 2016 and 2017 for accredited secondary Thinking Schools (0.73 and 0.77 respectively). Further, the primary impact measures relating to reading, writing and mathematics would be viewed as 'Moderate' to 'High' (equivalent to between 0.5 – 0.6 extra grade growth) for accredited primary Thinking Schools (0.33/0.44, 0.27/0.25, 0.43/0.43 respectively). In the case of secondary P8, primary reading progress and primary mathematics progress, the confidence intervals for accredited Thinking Schools *do not* include zero and so we can reasonably assume that these effects are also *significant*.

For all registered Thinking Schools (inclusive of accredited and those pursuing accreditation), Tables 2 and 4, together with Figures 2 and 4, clearly indicate what would be termed 'Moderate' to 'High' impact for P8 in both 2016 and 2017 (0.45 and 0.44 respectively). Further, the primary impact measures relating to reading, writing and mathematics would be viewed as 'Moderate' on the whole (around 0.3 extra grade growth) for registered primary Thinking Schools (0.17/0.17, 0.21/0.21, 0.27/0.28 respectively). Once again, in the case of secondary P8 and primary mathematics progress, the confidence intervals *do not* include zero and so we can reasonably assume that these effects are *significant*. For primary reading progress and primary writing progress in 2016 and 2017, the confidence intervals indicate that these effects are very close to being significant.

If we compare the data for accredited Thinking Schools with that of all registered Thinking Schools for both 2016 and 2017, we can see that the picture is one of clear and significant growth in terms of impact (all effects show a marked increase from registered to accredited status). In addition, primary reading and mathematics progress impact measures show growth from registered to accredited status. In addition, secondary P8 impact and primary reading impact shows clear growth from 2016 to 2017.

The following section adds more of a story to these initial interpretations and offers some overall evaluative conclusions.

Evaluation

By referring back to each of the previous sections of the SPARE wheel model, it is possible to draw some tentative conclusions and recommendations for further consideration.

Without attempting to draw *absolute* cause-effect claims, the evidence suggests that taking the Thinking Schools approach impacts greatly on the progress of pupils in both primary and secondary schools in England as measured by P8, reading progress, writing progress and mathematics progress. Indeed, the effect sizes are consistent with the high impact on pupils' achievement of meta-cognitive strategies illustrated by the research of Hattie (2009) and Higgins *et al.* (2013). The very high impact on P8 is of particular note as this measure spans a *wide* range of subject disciplines, not just the traditional canon of achievement represented by mathematics and English.

Comparing the data for accredited Thinking Schools with registered Thinking Schools illustrates clear impact growth in all progress areas examined once the criteria for accreditation has been met. This would tend to reflect the importance of fully embedding cognitive education, a process that would normally take at least 3 years given the profile of the schools in this study. The very high impact on P8 would tend to suggest that further research into possible 'latent' achievement development as pupils move through primary and secondary phases may further add to the body of knowledge in the areas covered by this study.

Given the importance of the context for specific schools and that schools naturally apply the Thinking Schools approaches in different ways, there is a need to not only share 'what works' *generally* across all schools but also enable schools to develop 'what works' for their own *particular* context. For example, under the original model of school training in the use of cognitive tools provided by TM, schools were provided with three, interrelated pathways to pursue:

- Visual Tools for Thinking – tools that explicitly support thinking processes
- Dispositions for Mindfulness – intelligent learning behaviours
- Questioning for Enquiry – skills for effective questioning and enquiry

The selection of pathways, including the order of progression, was left for the schools to decide. Some would start with Questioning for Enquiry and then move onto Visual Tools for Thinking. Others would start with Visual Tools for Thinking and then move onto Dispositions for Mindfulness. The corollary of the pathway model has been that schools develop in a rather dis-jointed manner where the cognitive tools are seen as discrete components. Further, schools began to describe themselves as a Visual Tools school or Questions for Enquiry Schools. In short, schools were showing good impact as illustrated by this evaluation, but were adopting a narrow approach to the teaching of thinking. TM, having recognised this and in their pursuit of building further on the already impressive impact of their approach, have

created a more integrated model of teaching thinking that is based on the dynamic and interrelated nature of the thinking process and the use of Grounded Practice principles (see Walters, 2014).

In addition, through the sharing of best practice offered by ‘hub’ Thinking Schools, and application of formative evaluative processes to compliment the more summative evaluation offered by accreditation (although it is advisable to use this more ‘summative’ evaluation in a formative style on a three year cycle), schools can monitor what works best in their *own* specific context. This new wave of development brings cognitive education and assessment together in order to secure self-improving organisations with a common conception of impact, independent but embracing of, any accountability measure. Given the higher impact on achievement for accredited Thinking Schools, where formal training in cognitive education has been undertaken, this also needs consideration if schools are to fully realise the potential of taking a whole school approach to the teaching of thinking.

Finally, the ‘snap-shot’ picture of impact provided by this evaluation of 2016 and 2017 national student outcomes in England shows a relatively stable positive picture for accredited and registered Thinking Schools. This would tend to suggest that this positive impact is not cohort dependent as the data is representative of different cohorts of pupils. Also, given the increased rigour and demand of the new 9 – 1 GCSE grading system in secondary schools in England, and that most secondary schools saw their results drop in 2017 (Treadway, 2017), the stability apparent in the Thinking Schools would tend to indicate that they were not impacted negatively by this change of the examination system. Together, these factors point to positive resilience to large-scale changes in the educational landscape.

Future avenues of enquiry

Although this study presents a snap-shot of 2016 and 2017 outcomes in England and does not track individual school development as they move through the Thinking Schools process, continuing this evaluative approach in an ongoing way would extend our understanding of how Thinking Schools develop in a changing national and international educational landscape. Further, it would also allow the new TM model of teaching thinking to be monitored and evaluated as it gains ground.

In addition to continuing with an ongoing snapshot evaluative approach, a future research focus in England would seem to lie in a longitudinal study of how pupils and schools develop over time. Evidence of long lasting ‘far transfer’, in terms of accelerated pupil growth, as indicated by multiple indicators (not *just* academic achievement), would be of particular interest here.

On a more international level, the challenge would be to replicate the evaluation offered by this study more globally using the key outcome measures specific to particular countries. Again, once a snap-shot of impact has been illustrated, countries across the globe may wish

to fund more longitudinal studies in order to grow their Thinking Schools further in their own particular contexts.

Finally, if a vision for the teaching profession is truly one of developing researching professionals on order to create a self-improving system, then teachers need support and training in order to research their work themselves. This would extend the general 'what works best' at the macro level to more fully realise the potential of taking a whole school approach to the teaching of thinking in specific school contexts at the micro level.

References

- Burden, R.L. 1998. Illuminative evaluation. *Educational and Child Psychology* 15, no.3: 15-23.
- Burden, R.L. and Nichols, L. 2000. Evaluating the process of introducing a thinking skills programme into the secondary school curriculum. *Research Papers in Education* 15, no.3: 293-306
- Clarke, S. 2005. *Formative Assessment in the Classroom* London: Hodder Murray
- Coe, R. 2002. It's the effect size, stupid: What effect size is and why it is important. *Paper presented at the Annual Conference of the British Educational Research Association, University of Exeter, England, 12-14 September 2002.*
- Cohen, J. 1988. *Statistical Power Analysis for the Behavioural Sciences* (2nd ed.). Hillsdale, New Jersey: Lawrence Erlbaum Associates.
- DfE, 2016. *Primary Progress Measures: How the Primary Progress Measures are Calculated.* London: DfE.
- DfE, 2017. *Progress 8 and Attainment 8: Guide for Maintained Secondary Schools, Academies and Free Schools.* London: DfE.
- Ellis, P.D. 2010. *The Essential Guide to Effect Sizes: Statistical Power, Meta-Analysis and the Interpretation of Research Results.* London: Cambridge.
- Harford, S. 2017. School inspection update. March 2017, Issue 9. London: Ofsted.
- Hattie, J. A. 2009. *Visible Learning – A Synthesis of over 800 Meta-Analyses Relating to Achievement.* London: Routledge.
- Hedges, L. and Olkin, I. 1985. *Statistical Methods for Meta-Analysis.* New York: Academic Press.
- Higgins, S. et al. 2013. *The Sutton Trust – Education Endowment Foundation Teaching and Learning Toolkit: Technical Appendices.* London: The Sutton Trust.
- Stobart, G. 2006. 'The validity of formative assessment' in Gardner, J. (ed) 2006. *Assessment and Learning* London: Sage pp. 133-146.

Treadway, M. 2017. *Key stage 4 performance tables: Closing the gap just got harder*. (<https://educationdatalab.org.uk/2018/01/key-stage-4-performance-tables-2017-closing-the-gap-just-got-harder/>). London: FFT Datalab.

Walters, D. 2014. Grounded practice: putting the 'self' back into self-evaluation. *Educational Action Research* Vol. 22 No. 1

Waters, M. 2013. *Thinking Allowed on Schooling*. London: Crown House Publishing.