

Review of Educational Research

The Effect of School Tracking on Student Achievement and Inequality: A Meta-Analysis

Journal:	<i>RER</i>
Manuscript ID	RER-21-Feb-MS-103.R2
Manuscript Type:	Manuscript

SCHOLARONE™
Manuscripts

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

The Effect of School Tracking on Student Achievement and Inequality: A Meta-Analysis

For Peer Review

Abstract

This meta-analysis examines the effects of sorting secondary students into different tracks (“between-school” tracking) or classrooms (“within-school” tracking) on the efficiency and inequality levels of an educational system. Efficiency is related to the overall learning achievement of students, while inequality can refer to “inequality of achievement” (i.e., the dispersion of outcomes) or “inequality of opportunity” (i.e., the strength of the influence of family background on student achievement). The selected publications are 53 analyses performed in the period from 2000 to 2021, yielding 213 estimates on efficiency and 230 estimates on inequality. The results show that the mean effect size (Hedge’s G) of tracking on efficiency is not statistically significant ($G = -.063$), whereas it is significantly positive ($G = .117$) on equality. We further set out to explain variation in effect sizes by (a) policy characteristics, (b) the operationalization of main variables, (c) the research design, (d) the set of control variables included in the statistical analyses, and (e) the quality of the study, year of publication, and publication status (peer-reviewed or not peer-reviewed).

Keywords: secondary school tracking, achievement, efficiency, inequality, meta-analysis

The Effect of School Tracking on Student Achievement and Inequality: A Meta-Analysis

Almost all industrialized societies provide unified initial schooling in which all students receive the same general education and are exposed to a relatively homogeneous school environment. Nevertheless, in many countries, students are allocated to different types of education at some point in their educational career (Blossfeld et al., 2016). Indeed, most school systems practice some form of tracking, through which they provide a differentiated learning environment to students (Dupriez et al., 2008).

In a broad sense, tracking (also known as streaming, sorting, or ability grouping) refers to the assignment of students to different types of education—kind of school, curricula, subjects, classes—according to their ability, interests, or attitudes (Betts, 2001; Brunello & Checchi, 2007; Woessmann, 2009). School tracking is considered by social scientists to be one of the most important features of school systems because it has important consequences for students' school performance, educational pathways, entrance into higher education, and subsequent labor market outcomes (Brunello & Checchi, 2007; van de Werfhorst & Mijs, 2010).

National educational systems differ in the way they organize secondary education. Some offer different pathways to students or group them based on their ability and interests relatively early in their educational career, while others opt for later student sorting or a more comprehensive model of education (Dupriez et al., 2008; Blossfeld et al., 2016). It has long been argued that the shape of secondary education has consequences for the efficiency of an educational system, that is, the overall achievement of students (Betts, 2001; Clifford & Heath, 1984; Gamoran & Mare, 1984).

INTRODUCTION

4

1
2
3 The underlying argument supporting student tracking is that more homogeneous groups
4 of students allow teachers to tailor their pedagogical strategies to students' abilities and interests,
5 making the learning process more efficient. This is called the "specialization effect" and,
6 following this argument, tracking can allow students to learn more. Nonetheless, the possibility
7 that school tracking may affect educational inequality must also be considered (van de Werfhorst
8 & Mijs, 2010). Some authors argue that tracking increases inequality among students (Gamoran
9 & Mare, 1989; Lucas, 1999/2001; Oakes, 1985) not only in terms of dispersion of achievement
10 (Galindo-Rueda & Vignoles, 2005; Hanushek & Woessmann, 2006; Manning & Pischke, 2006)
11 but also in terms of inequality of opportunities, namely, the strength of the influence of family
12 background on student outcomes. These studies usually quantify the extent to which economic,
13 cultural, and social resources in the home environment affect student achievement across
14 educational systems characterized by different institutional arrangements (Bol & van de
15 Werfhorst, 2013; Brunello & Checchi, 2007; Jackson, 2013). From this perspective, gains in
16 achievement in a tracked system concern students at more advanced levels or with higher
17 socioeconomic status, whereas students from socio-economically disadvantaged families face
18 deterioration in their learning progress (Hallinan, 1994).

19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41 Studies on the effects of tracking are important not only for the academic debate but also
42 because of their informative potential in shaping educational policies. Three scenarios arise
43 regarding the effect of tracking on both dimensions (efficiency and inequality). First, if tracking
44 increases efficiency as it exacerbates inequality, policy-makers will be dealing with an
45 "efficiency-equity trade-off" (Skopek et al., 2019, p. 224): is it more important to maximize
46 student performance at all costs, or is it preferable to distribute education more equally within the
47 system? Second, if the effects of tracking are negative or null in terms of efficiency, would
48
49
50
51
52
53
54
55
56
57
58
59
60

INTRODUCTION

5

1
2
3 increasing inequality only for the benefit of more privileged students be legitimate? Third, if the
4 effects on both efficiency and inequality are null, would tracking still be a valid policy choice for
5 providing students with a more specialized education?
6
7
8
9

10 While a fully comprehensive assessment of the effects of tracking would also include its
11 longer-term consequences for employability and occupational outcomes, in this analysis, we will
12 focus on the consequences of tracking for short-term outcomes only, that is, educational
13 achievement. This homogenizes tracking outcomes both in conceptual terms and in their
14 measurements, thus posing fewer problems in the design of the meta-analytical review.
15
16 Furthermore, thanks to the existence of several international pre-harmonized surveys using
17 student test scores, a far greater amount of consolidated empirical literature is available regarding
18 the effects of tracking on educational achievement than about its longer-term consequences for
19 occupational prospects.
20
21
22
23
24
25
26
27
28
29
30

31
32 As suggested by Woessmann (2009), since there are plausible contrasting arguments at a
33 theoretical level about the role of tracking in contemporary educational systems, it is of
34 paramount importance to address this issue at the empirical level. However, there is no
35 consensus in the empirical literature. Indeed, some studies report either positive (Ariga &
36 Brunello, 2007; Galindo-Rueda & Vignoles, 2005; Horn, 2013; Lassibile & Gomes, 2010) or
37 null effects of tracking on the overall level of student performance, and others report negative
38 (Ammermüller, 2005; Bol & van de Werfhorst, 2013; Brunello & Checchi, 2007; Dunne, 2010;
39 Hanushek & Woessmann, 2006; Piopiunik, 2014) or null effects (Jakubowski, 2010; Ruhose &
40 Shwerdt, 2016; Waldinger, 2007) on social inequalities. Differences in the qualitative direction
41 of the results and the magnitude of the effects might stem from the heterogeneity of samples, the
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

INTRODUCTION

6

1
2
3 definition of tracking, the type of research design, and several other characteristics (e.g.,
4
5 Ammermüller, 2005; Jakubowski, 2010; Waldinger, 2007).
6
7

8
9 The last extensive meta-analysis of secondary tracking was performed in the 1980s by
10
11 Kulik and Kulik (1982). It focused exclusively on within-school tracking and the problem of
12
13 efficiency, with only side concerns for inequality. Their results pointed to small yet significant
14
15 effects of ability grouping on overall achievement. Given that the empirical literature has since
16
17 flourished, especially in the 2000s, due to growing attention to the topic and the increased
18
19 availability of suitable data and analytical techniques, it is important to understand the findings
20
21 of recent research. The few existing reviews have summarized the key elements of the recent
22
23 literature. In their assessment, van de Werfhorst & Mijs (2010) suggest that tracking increases
24
25 inequality and lowers efficiency. Meanwhile, Woessmann (2016) shows that tracking is a crucial
26
27 institutional feature for explaining differential performances across countries in international
28
29 standardized tests such as the Program for International Student Assessment (PISA) and that
30
31 evidence points to a deterioration of the performance of students placed in low and middle
32
33 tracks. Skopek et al. (2019) review the plurality of approaches employed to analyze the effects of
34
35 tracking, suggesting that the heterogeneous results of empirical studies might depend partly on
36
37 the diversity of analytical strategies. However, these reviews are non-systematic and cannot
38
39 provide an exhaustive description of the phenomenon or explain the most important drivers of
40
41 the variation in findings across empirical studies with proper statistical analysis. Older reviews
42
43 include the best-evidence synthesis performed by Slavin (1990) and the summary of survey and
44
45 ethnographic research compiled by Gamoran and Berends (1987), both published by this journal.
46
47
48
49
50
51
52 The former found no effect of ability grouping on overall achievement in studies that treated
53
54 tracking as a systemic variable and in studies that evaluated the effect of placement in different
55
56
57
58
59
60

INTRODUCTION

7

1
2
3 ability classes. The latter stressed the difficulties of the literature in disentangling the influence
4
5 of tracks from the effect of students' social-class, calling for more robust quantitative
6
7 longitudinal research, which only flourished recently and will be reviewed in this meta-analysis.
8
9

10
11 Therefore, we believe that the field would benefit from a systematization of these
12
13 research findings capable of answering two key questions in the current debate. First, to what
14
15 extent does tracking increase or decrease efficiency and inequality in the educational system?
16
17 Second, to what extent can the variation in results across studies be explained by the latter's
18
19 characteristics?
20
21

22
23 We propose a meta-analysis of studies that operationalized tracking as a
24
25 systemic/institutional macro-level variable (i.e., the sorting of students within the educational
26
27 system) and measured its impact on efficiency and/or inequality for the educational system as a
28
29 whole. The meta-analytic method offers some advantages. Through its concept of "mean effect
30
31 size," it allows us to take the estimates reported in different statistical forms and produced based
32
33 on different research designs and make them comparable through a standardization process that
34
35 mainly takes into account sample size and variance (Borenstein et al., 2009; Lipsey & Wilson,
36
37 2001). Second, via the systematic coding of studies, it is possible to explain variation in effect
38
39 sizes by study and policy features. We (meta-)regress effect sizes on five blocks of explanatory
40
41 variables: (a) policy characteristics, (b) the operationalization of the tracking and outcome
42
43 variables, (c) the research design, (d) the set of control variables included in the statistical
44
45 analyses, and (e) the quality of the study, year of publication, and publication status (peer-
46
47 reviewed or not peer-reviewed). All of these elements received diverse treatment in the studies
48
49 selected for analysis, each possibly contributing to variation in effect sizes. We run separate
50
51 analyses for efficiency and inequality.
52
53
54
55
56
57
58
59
60

Theoretical framework

In the following section, we discuss the main mechanisms highlighted in the literature that could explain the effects of tracking on efficiency and inequality in educational systems. They include (a) institutional characteristics (e.g., curricula), (b) characteristics of the learning environment (e.g., classroom composition and peer effects), and (c) the unequal distribution of educational resources (e.g., teachers' characteristics, class size). We will review the main theoretical arguments in favor of the tracking policy before considering the arguments that highlight its potentially adverse effects.

Potential benefits of tracking

The proponents of tracking argue that early tracking can have a positive effect on the overall level of learning and school performance because, in a tracked system, students attend a school environment that is tailored to their needs and receive instruction adequate for their skill level, which in turn maximizes each student's potential. The most important pedagogic rationale for tracking is that students differ in their academic potential and the environments in which they learn best.

A tracked school system is also helpful for teachers, who can develop their teaching approach based on the ability level of the classroom, which is on average more homogeneous than in a comprehensive school system because of the prior selection and streaming (Lassibile & Gomez, 2010; Manning & Pischke, 2006). From a systemic point of view, if the allocation is efficient—that is, students are appropriately assigned to the track that best fits their ability level—a tracked system should produce an overall higher level of student performance than a

THEORETICAL FRAMEWORK

9

1
2
3 non-tracked system. In other terms, by raising average outcomes, efficient tracking should
4
5 enhance educational productivity (Betts, 2001).
6
7

8 According to Gamoran and Mare (1989), this argument can be traced back to old studies
9
10 that analyzed the effects of classroom homogeneity on educational achievement (Cook, 1924;
11
12 Keliher, 1931; Whipple, 1936) and concluded that a more homogeneous environment improves
13
14 the efficiency of the learning process. An underlying assumption of this view is that the process
15
16 of sorting students should be based solely, or at least largely, on students' academic abilities. In
17
18 this perspective, an appropriately designed tracking system could not only positively affect the
19
20 efficiency of the system but also compensate for inequalities in achievement between low- and
21
22 high-performing students. If, after track allocation, the achievement of low-track students
23
24 increases more than that of high-track students, inequality will decline (Gamoran & Mare, 1989).
25
26
27
28
29

30 More recent works suggest that, in many educational systems, the process of sorting
31
32 students occurs not only on the basis of their abilities but also by taking into account their
33
34 attitudes and aspirations, leaving space for families to choose their child's academic destination
35
36 independently from their demonstrated academic competence (Jackson, 2013). Some students
37
38 aim to enter the labor market early and develop occupation-specific and applied skills. On the
39
40 contrary, others are oriented to continue to study at university either because they like studying
41
42 and do not have well-defined occupational expectations or because their target job requires a
43
44 university degree (e.g., doctor, lawyer). Many European educational systems differentiate
45
46 between schools offering academic-oriented curricula and others providing vocationally oriented
47
48 instruction and training (Shavit & Müller, 1998). A school with an efficient tracking system
49
50 should provide study courses appropriate for different types of students to better address their
51
52 needs. As suggested by Gamoran and Mare (1989), "ideally, a system of academic tracking
53
54
55
56
57
58
59
60

THEORETICAL FRAMEWORK

10

1
2
3 matches students' aptitudes with the objectives and learning environments to which they are best
4 suited" (p. 1148). If student sorting can effectively assign students to the most appropriate track,
5
6
7 it could lead to higher levels of learning and academic performance, greater student satisfaction,
8
9
10 and lower dropout rates.

Why tracking can be harmful

11
12
13
14
15
16 At the theoretical level, tracking can also have negative consequences for both efficiency
17
18 and equity. First, it has been highlighted that the expected benefits outlined in the previous
19
20 section depend on the extent to which institutions and families are able to place the students in
21
22 the type of school or study program that is most appropriate for their skill and aptitude levels. In
23
24 this respect, it is often claimed that the earlier the separation of students occurs, the greater the
25
26 likelihood of making an inappropriate assignment (Betts, 2001; Brunello & Checchi, 2007;
27
28 Dustmann, 2004). This is because the pace of cognitive development varies across students, as
29
30 well as the degree of maturation and development of non-cognitive skills that are conducive to
31
32 higher achievement (Cohn, 1991).
33
34
35

36
37
38 Several authors have pointed out possible negative outcomes of tracking in terms not
39
40 only of dispersion of achievement levels but also inequality in educational opportunity (Triventi
41
42 et al., 2019; van de Werfhorst & Mijs, 2010). First, it is likely that the stratification introduced by
43
44 the practice of tracking particularly penalizes students from disadvantaged socio-economic
45
46 backgrounds when the placement in lower-level study programs is affected by pupils' social
47
48 origin, beyond school achievement. If enrolment in different types of schools is free and depends
49
50 on families' decisions (as in Italy), the process of segregation in different types of education is
51
52 mainly based on students' self-selection. In those school systems where allocation to different
53
54
55
56
57
58
59
60

THEORETICAL FRAMEWORK

11

1
2
3 tracks is based on teacher recommendations (as in some German *Länder*), the resulting social
4
5 stratification is mainly due to institutional practices (Contini & Scagni, 2011).
6
7

8 In a free-choice system, for culturally embedded reasons (Bourdieu & Passeron, 1970) or
9
10 because of risk aversion (Breen & Goldthorpe, 1997), lower-status parents may decide to enroll
11
12 their children at a vocational school even if they have a relatively good school record. On the
13
14 contrary, children from higher socio-economic status (SES) families with a similar level of
15
16 achievement are more likely to attend the academic track to increase their chances of admission
17
18 to university in the future (Gambetta, 1987).
19
20
21
22

23 In a system taking into account teacher recommendations, students of lower social origin
24
25 are more likely to be directed toward vocational schools or study programs, for two reasons: (a)
26
27 lower-class pupils may have lower school achievement in primary school because of differences
28
29 between social classes in the development of cognitive abilities in early age (Brunello &
30
31 Checchi, 2007); (b) teacher recommendations can be biased in favor of upper-class children due
32
33 to teachers' beliefs about students' "natural" attitudes as well as pressures from upper-class
34
35 parents to place their children in the higher track (Barg, 2013; Brantlinger, 2013; Laureau, 1987).
36
37
38
39

40 Tracking can exacerbate social inequality in school achievement and performance
41
42 through several mechanisms. The first one is peer-group effects: if highly motivated and high-
43
44 achieving pupils are grouped, the low-achieving students are segregated and cannot benefit from
45
46 proximity to such peers (Betts & Shkolnik, 2000; Deci et al., 2001; Esser, 2016). Additionally,
47
48 attending school with more motivated peers could create a favorable learning environment,
49
50 allowing teachers to devote more time to effective teaching and less time to managing the
51
52 classroom or dealing with students' misconduct. The second is teacher sorting (Bonesronning et
53
54 al., 2005): it is possible that the ablest and most motivated teachers prefer to teach the brightest
55
56
57
58
59
60

THEORETICAL FRAMEWORK

12

1
2
3 students in the academic track (Brunello & Checchi, 2007, p. 795). If this is the case, teachers
4
5 engage in positive self-selection into the higher-level tracks, resulting in an increased learning
6
7 gap between students placed in high and low tracks. For instance, until a few years ago, teachers'
8
9 official preparation for the academic track in Germany lasted longer and was more complete than
10
11 for other tracks (Roloff et al., 2020, p. 3). The third mechanism refers to differences in the
12
13 quality of curricula and teachers' expectations. In higher-level tracks, educational standards are
14
15 stricter and teachers often have higher expectations regarding their students' academic potential.
16
17 This could contribute to improving student performance in this type of school because positive
18
19 teacher expectations (Rosenthal & Jacobson, 1968) and rigorous educational standards seem to
20
21 foster school achievement (Betts & Grogger, 2003). The fourth potential mechanism concerns
22
23 the educational resources invested in different tracks, such as average expenditure per student or
24
25 student-teacher ratio (Brunello & Checchi, 2007). One example of potentially heterogeneous
26
27 resources across tracks is class size. A student in a class with a small number of students is more
28
29 likely to receive personal teacher support during lessons than students in more crowded
30
31 classrooms. Some international evidence points to lower student-teacher ratios in general tracks
32
33 than in vocational tracks (Brunello & Checchi, 2007, p. 795). However, the extent to which this
34
35 can exacerbate inequalities or compensate for them seems to vary across countries (Betts, 2001).
36
37
38
39
40
41
42

43 In summary, since tracking policy can intensify inequalities among students in terms of
44
45 the quality of their regular educational environment (educational resources, peers), we expect it
46
47 to lead to heterogeneous outcomes across tracks or ability groups, canceling out any possible
48
49 improvement in efficiency produced by specialization. Therefore, we hypothesize that *tracking*
50
51 *does not increase the overall efficiency (achievement) of an educational system (H1a)* and that
52
53
54
55
56
57
58
59
60

1
2
3 *tracking increases the level of inequality of achievement/opportunity within an educational*
4
5 *system (H1b).*
6
7

8 If tracking increases efficiency as it increases inequality, support is provided for the
9
10 “equality-efficiency trade-off” (Skopek et al., 2019, p. 224). If its effect is null on efficiency but
11
12 positive on inequality, support is provided for studies that are more critical of tracking. Finally, if
13
14 its effects are null on both dimensions, support is provided for the organizational and transitional
15
16 arguments, that is, the idea that tracking is still valid for attending to students’ specific interests
17
18 and facilitating their transition to the labor market.
19
20
21
22

23 **Explaining variations across studies: hypotheses on the role of meta-regressors**

24
25

26 As anticipated, research on the effect of tracking as an organizational feature of the
27
28 educational system used a variety of tracking measures, research designs, datasets, and country
29
30 or regional samples. We set out to explain variation in effect sizes based on five blocks of
31
32 explanatory variables: (a) policy characteristics, (b) the operationalization of tracking and
33
34 outcome variables, (c) the research design, (d) the set of control variables included in the
35
36 statistical analyses, and (e) the quality of the study, year of publication, and publication status.
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

HYPOTHESES

14

Policy characteristics

Type of tracking. A basic distinction in tracking policy is that between “within-school” and “between-school” tracking. While the first is more common in anglophone countries, the latter is usually observed in Europe (Betts, 2000; Blossfeld et al., 2016). Within-school tracking, frequently referred to in the literature as “ability grouping,” consists in assigning students to different groups or classrooms (e.g., basic, intermediate, or advanced mathematics classes) within the same school. Between-school tracking assigns students to different school types (e.g., academic, technical, or vocational tracks). This distinction is important because the type of tracking adopted by an educational system will affect the change in school/class composition and the distribution of educational resources differently.

For instance, within-school tracking affects student group composition only at the classroom level, where the learning process mainly takes place, whereas between-school tracking alters it more drastically, re-shaping peer composition at both the classroom and school levels. Consequently, we believe that the effects of tracking should be different for each of these modalities. Students that are sorted into different schools will learn in a fully specialized environment, while students that are sorted only into different classes or courses will have a part-time specialization, still sharing the school environment and potentially other courses with less advanced peers. Thus, we expect that *between-school tracking can increase efficiency more than within-school tracking* (H2a) and that *between-school tracking can increase inequality more than within-school tracking* (H2b).

HYPOTHESES

Outcome variables

Subject domain. Several studies operationalize the dependent variable by pooling scores from different domains in one single overall score (Ayalon, 2006; Kerr et al., 2013; Koerselman, 2013; Piopiunik, 2014). Nevertheless, numerous others maintain the separation between subject domains (Hanushek & Woessmann, 2006; Waldinger, 2007), assuming that tracking may have differentiated effects across school subjects (Dammrich & Triventi, 2018).

The bulk of reading and mathematical competencies is developed at school. However, social background seems to play a more important role in the development of reading skills than in that of mathematical competencies (Bol et al., 2014; Cooper et al., 1996). Verbal and reading skills are also acquired via interactions with parents and peers and reading at home, whereas mathematical competencies are mostly developed at school (Cooper et al., 1996). Following this reasoning, on theoretical grounds, we can expect student achievement in mathematics to be more sensitive to the characteristics of the educational system than student achievement in reading. Thus, *the effect of tracking should affect mathematics scores more markedly than reading scores* (H3a) and *create more inequality in mathematics than in reading* (H3b).

Type of outcome. As anticipated, we divide studies into two main groups. The first analyzes the effect of tracking on the average level of student achievement, which is thought to capture the consequences of school tracking for efficiency. The outcome is usually measured based on test scores constructed by applying item-response theory or related techniques to a set of student answers to a standardized test. The second group focuses on the effect of tracking on achievement inequality. Given that achievement inequality has been measured in two very distinct ways in tracking studies, a specific hypothesis for the operationalization of the dependent variable is needed for research focusing on this outcome.

HYPOTHESES

16

1
2
3 The first subgroup of studies measures achievement inequality in terms of dispersion of
4 student performance, that is, the standard deviations of scores, interquartile range, score gap
5 between top and low performers, and ratio of low and top performers, inter alia (Bol & van de
6 Werfhorst, 2013; Galindo-Rueda & Vignoles, 2005; Hanushek & Woessmann, 2006;
7
8 Jakubowski, 2010; Vandenberghe, 2006). The second subgroup of studies, in their basic setting,
9
10 measures inequality of opportunities by interacting an indicator of family background (e.g.,
11 number of books at home, parental education, parental occupation or income) with a macro-
12 variable capturing the tracking regime (e.g., Ammermüller, 2005; Le Donné, 2014). A positive
13 interaction indicates that the effect of social background on student achievement (i.e., inequality
14 of opportunity) is larger in tracked than in untracked systems. While achievement dispersion is
15 an inter-individual measure of inequality, the second is a group-based measure of inequality in
16 which categories of students are identified on the basis of an ascriptive trait (Breen & Jonsson,
17 2005) that can lead to scholastic advantages or disadvantages for which students do not bear any
18 responsibility (Roemer, 1998). Many authors agree that both measures of achievement inequality
19 are needed to develop a more comprehensive view of the social consequences of tracking
20 (Schütz et al., 2008; Waldinger, 2007; Pfeffer, 2008; Schlicht et al., 2010; Burger, 2016; Ruhose
21 & Schwerdt, 2016).
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

43 It is not easy to formulate a hypothesis on whether we should expect a higher level of
44 inequality using one measure or the other because empirical research has not engaged in this
45 comparison explicitly. To this end, we should consider that since inequality of opportunity in
46 these studies is measured in a more indirect way (in a moderation analysis) than achievement
47 dispersion, we can tentatively expect that *measurements of inequality of opportunity tend to*
48 *report smaller estimates than measurements of inequality in achievement (H3b.2).*
49
50
51
52
53
54
55
56
57
58
59
60

HYPOTHESES

17

Tracking indicator

The independent variable (tracking) is operationalized diversely across the studies. A first basic distinction concerns whether tracking is treated as a dummy variable (“tracked” vs. “comprehensive” or “highly tracked” vs. “moderately tracked” systems) or as a continuous variable (e.g., age at first tracking, number of years exposed to tracking, or number of tracks offered). If tracking is operationalized as a dummy variable, the study reflects the full potential effect of a tracking policy. Conversely, if it is operationalized as a continuous variable, the study reports the effect of one additional unit of a certain tracking component.

Another important distinction in this sense is when tracking is operationalized based on reform dummies. Systems that have implemented educational reforms—such as Finland, the United Kingdom, Sweden, Poland, and Bavaria in Germany—offer interesting cases for measuring the impact of tracking because of its longitudinal nature. Yet, one important drawback of studies that analyze tracking reforms is that if the reform includes any institutional or organizational change other than tracking (e.g., interventions on time of instruction, curricula modifications), its effects are very difficult to isolate (Koerselman, 2013). Additionally, predicting the direction of the potential bias is tricky since it depends on the nature of concomitant interventions. In any case, we believe that by focusing on longitudinal variations within countries, the studies analyzing reforms are not affected by time-constant unobserved heterogeneity at the macro-level, which is more likely to plague estimates in international studies, especially those based on cross-sectional designs. Following this consideration, we hypothesize that *studies that analyze national educational reforms tend to report smaller effects of tracking both for efficiency (H4a) and inequality (H4b)*.

HYPOTHESES

18

Research design

Choice of counterfactual. One important feature of tracking studies is the type of reference category or counterfactual to which tracking is compared. Many studies compare tracking not to comprehensive systems but to tracking regimes with less strict rules. The most obvious example is studies comparing regimes that track at a very early age to regimes that track later. Bauer and Riphan (2013, p. 112), for instance, take advantage of the variation in tracking rules across Swiss cantons: the first tracking in Switzerland occurs mostly between grade 5 (approximately 10–11 years of age) and grade 7 (12–13 years). Since their outcome variable comes from a dataset that measured the ability of students at the age of 17, all students had already been tracked by this time. Thus, in this study, the estimated impact was the consequences of tracking one or two years earlier, not the impact of attending a tracked rather than a comprehensive system. Other examples of this type of comparison can be found in Dronkers et al. (2012) and Dunne (2010), who contrasted “between-school” tracking with the “within-school” modality.

We set out to test whether studies that compare tracking with comprehensive systems (Felouzis & Charmillot, 2008; Hanushek & Woessmann, 2006; Hoffer, 1992; Horn, 2013; Jakubowski, 2010; Pfeffer, 2008) report larger estimates for tracking than those that use other types of reference categories. Since the contrast is starker when tracking regimes are compared to non-tracked ones, we hypothesize that *studies that use comprehensive education as their counterfactual tend to report larger effects of tracking on achievement (H5a) and inequality (H5b) than those whose counterfactuals represent less-tracked systems.*

HYPOTHESES

19

1
2
3 *Analytical design.* When looking at comparative studies investigating the role of school tracking
4 systems on inequalities in educational outcomes, two main types of research designs can be
5 identified (Skopek et al., 2019): cross-sectional multilevel-analysis (e.g., Ayalon & Gamoran,
6 2000; Burger, 2016; Schütz et al., 2008) and difference-in-differences designs (e.g., Hoffer,
7 1992; Zimmer, 2003).

8
9
10
11
12
13
14
15 In the first stream of research, the increased availability of international large-scale
16 student surveys – e.g., PISA, Progress in International Reading Literacy Study (PIRLS), Trends
17 in International Mathematics and Science Study (TIMSS) – has promoted the development of
18 comparative quantitative research on educational inequalities. Multilevel cross-sectional designs
19 rely on the variation in educational policies across countries, taking into account the hierarchical
20 nature of educational data and correcting for its dependencies (e.g., students nested in schools,
21 nested in countries). Since individual achievement is explained by variation at the macro-level, it
22 is rather unlikely that confounders at the micro-level are correlated with these macro-features. A
23 problem with this approach, however, is that it does not rule out confounders at the macro-level
24 (Hanushek & Woessmann, 2006). Variation in efficiency and inequality across national
25 educational systems could be attributed not only to tracking but also to differences between
26 countries regarding other features, such as GDP, educational investment, or even national
27 culture. While some of these macro-level confounders can be taken into consideration, the
28 limited number of countries available makes it impossible to account simultaneously for all
29 possible confounders, leaving open the possibility that the effect of tracking estimated by these
30 studies is biased by unobserved factors.

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52 The second strategy, which was introduced in this field by the pioneering work of
53 Hanushek and Woessmann (2006), intends to overcome some of the aforementioned limitations.

HYPOTHESES

20

1
2
3 Applied to international data on student test scores, the difference-in-differences strategy aims to
4 enhance the internal validity of the estimated effect of tracking by relying on (at least) two points
5 in time. The key idea is to compare student outcomes (e.g., PISA scores measured around the age
6 of 15) in countries that have already tracked students and countries that have not yet tracked
7 them, comparing this difference with a previous point in time at which tracking had not yet
8 occurred in all countries (e.g., PIRLS and TIMSS scores measured around the age of 10). By
9 relating changes in outcomes to the changing policy (from untracked to tracked in a subsample
10 of countries), this design can remove time-constant heterogeneity across countries. While this is
11 a valuable improvement on the previous strategy, the causal interpretation of the tracking effect
12 will also be undermined in this design if there are unaccounted changes in other institutional
13 features between the two data points, which correlate with changes in tracking practice and
14 influence student performance (e.g., expenditure on education, quality of teachers).

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31 In any case, in line with the literature (Skopek et al., 2019; van de Werfhorst, 2019), we
32 believe that longitudinal designs are better at granting internal validity and ruling out possible
33 alternative explanations. Thus, we hypothesize that, *compared to the cross-sectional multilevel*
34 *design, the difference-in-differences design tends to report smaller effects of tracking on overall*
35 *student achievement (H6a) and inequality (H6b).*

36
37
38
39
40
41
42
43 ***Inclusion of educational mechanisms/resources.*** In discussing our theoretical framework, we
44 have highlighted that tracking could affect student outcomes by altering the school/class
45 composition and educational resources to which children are exposed in secondary education.
46 Relying on the literature on counterfactual causal inference, we distinguish between studies that
47 estimate the total effect of tracking and those that estimate an over-controlled model by including
48 potential mediators as well, that is, variables that intervene in the relationship between tracking
49
50
51
52
53
54
55
56
57
58
59
60

HYPOTHESES

21

1
2
3 and outcomes (Angrist & Pischke, 2008; Pearl, 2009). Regarding the set of controls, we classify
4
5 studies into four groups: (a) studies that control for educational resources (e.g., student-teacher
6
7 ratio, time of instruction, teacher qualifications), (b) studies that control for school/class
8
9 composition (e.g., school average SES or proportion of minority status among students), (c)
10
11 studies that account for both sets of controls (i.e., composition and resources), and (d) studies
12
13 that do not control for either of these mechanisms. Consequently, we hypothesize that *when a*
14
15 *study controls for educational resources and/or school/class composition indicators, the effect*
16
17 *sizes of tracking tend to be smaller for efficiency (H7a) and inequality (H7b).*
18
19
20
21
22

Quality, year, and publication status

23
24
25 **Quality score.** In addition to the choice of an adequate research design, the quality and rigor with
26
27 which they are applied can explain variation in effect sizes. We propose an evaluation matrix that
28
29 takes into account general and more specific criteria of evaluation (adapted from Conn, 2017).
30
31 The first two criteria are general and apply to all research designs. The first criterion (a) takes
32
33 into account the quality of the language and whether basic quantitative information is missing
34
35 (e.g., descriptive statistics, sample size, standard errors). The second criterion (b) considers how
36
37 well presented the research design is and whether the authors show awareness of its limitations.
38
39 The third and fourth criteria are design-specific and focus on the application of best practices to
40
41 reduce threats to internal validity. The third criterion (c) only concerns studies that use a
42
43 difference-in-differences design and reflects whether they included at least one previous point in
44
45 time to check for pre-existing trends and whether they balance treatment and control groups on
46
47 important observable variables (i.e., for which variables they control). The fourth criterion (d)
48
49 only applies to studies that employ a cross-sectional design and verifies whether a study respects
50
51
52
53
54
55
56
57
58
59
60

HYPOTHESES

22

1
2
3 the hierarchical structure of the educational data and controls for confounders at the macro-level.

4
5 Table 1 presents the structure of the evaluation matrix.

6
7
8 Scores can range from 0 to a maximum of 6 per study. Here, we follow the general
9
10 pattern found in other meta-analyses and hypothesize that *the higher the quality of the study, the*
11
12 *smaller the tracking effects will be on efficiency (H8a) and inequality (H8b)*. To test these
13
14 hypotheses, we use the overall quality score but also design-specific items that focus on internal
15
16 validity. We believe that if the first two criteria contain any sort of subjectivity, using only
17
18 criteria (c) and (d) will largely eliminate it, offering a more objective test.
19
20
21
22

23 ***Year of publication.*** We aim to understand whether the size of the estimates reported by the
24
25 studies under analysis follows any systematic trend over time, from 2000 to 2021.
26
27

28 ***Peer-review.*** Another potentially relevant feature is the publication status of the studies. In
29
30 particular, we distinguish between peer-reviewed articles and all other types of products, such as
31
32 policy reports and working papers. Although it can be assumed that peer-reviewed studies
33
34 underwent more rigorous quality control, the literature highlights many issues in the current
35
36 publication system in academic journals. Among the most serious concerns is the “file-drawer
37
38 problem” (e.g., Rosenthal, 1979), namely, the fact that the system discourages the publication of
39
40 studies that do not report statistically significant estimates. This occurs directly in the peer-
41
42 review system, in which reviewers are more likely to approve papers with statistically significant
43
44 findings, and indirectly via authors’ reticence to submit to academic journals papers whose
45
46 results provide no or weak support for their hypotheses. Thus, we anticipate that *peer-reviewed*
47
48 *studies tend to report larger tracking estimates than policy and discussion reports both for*
49
50 *efficiency (H9a) and inequality (H9b)*. If supported, this will provide evidence of a publication
51
52 bias in the current scientific publication system.
53
54
55
56
57
58
59
60

Method

Meta-analysis

Meta-analysis is characterized by clear-cut procedures and transparency. It offers an objective and mostly unbiased way to review a large body of literature on a certain topic and, especially in the social sciences, is a useful tool for standardizing results presented in the most diverse statistical forms to make them comparable. More specifically, we rely on random-effect meta-analysis, which is capable of systematizing study estimates by calculating effect sizes and can explain their variation based on study and policy features. This technique is particularly well suited when studies analyze heterogeneous populations and use different research designs, as is the case for the empirical literature on school tracking. Furthermore, estimates are weighted by the inverse of their squared variance, giving more weight to studies characterized by greater precision and statistical power (Lipson & Wilson, 2001; Siddaway et al., 2018). For our calculations, we used the “metafor” package (Viechtbauer, 2010), which is designed specifically to perform meta-analyses with the R software.

Search procedure

A systematic search in the Education Resources Information Center (ERIC) and Web of Science (WoS) databases was performed in September 2021 using three blocks of terms combined by the operator “AND”: (1) educational-field terms (“educ*” OR “student*”); (2) tracking-related terms (“track*” OR “ability group*” OR “stream*” OR “sort*”); and (3) terms associated with identification strategies (“experiment” OR “quasi-experiment” OR “difference-in-difference*” OR “multilevel” OR “regression*” OR “propensity score” OR “instrumental variable”). This procedure yielded 3,655 results, which were filtered in a first-level screening by

METHOD

24

1
2
3 scanning the publications' title and abstract to identify whether they concerned educational
4 tracking and employed any sort of quantitative analysis. A total of 222 reports were assessed for
5 eligibility, and 24 were retained. Figure 1 shows a PRISMA diagram describing the screening
6 process and reasons for exclusion. To complement this initial sample, we snowballed Skopek et
7 al. (2019), which was found to be the most recent non-systematic qualitative literature review on
8 the topic. This procedure produced 168 additional reports for screening, of which 29 remained
9 after the assessment based on the eligibility criteria. Therefore, the final sample is composed of
10 53 publications.
11
12
13
14
15
16
17
18
19
20
21
22

Selection criteria

23
24
25 All the selection criteria can be found in Appendix A. This set of criteria first aimed to
26 detect a basic distinguishing feature: the central theme of the study should be school tracking in
27 secondary education. As the independent variable, we considered "tracking" measured variously
28 by, inter alia, the presence or absence of tracking (dummy, educational reforms), the age at first
29 tracking, the number of tracks, and the length of tracking. One important reminder is that only
30 studies that treated tracking as a systemic organizational variable were coded for the analysis.
31 This excludes numerous studies on "within-school" tracking performed mainly in the United
32 States, in which the individual-level effect of attending specific tracks is estimated, taking into
33 account selection into different tracks or ability groups (Betts, 2001; Burks, 1994; Gamoran &
34 Mare, 1989). These studies measure the individual impact of assignment to a certain track
35 instead of the impact of tracking as a system-level policy, which is a linked but analytically
36 distinct research question (Triventi et al., 2021).
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52

53 As dependent variables, for educational efficiency, we selected a range of possible
54 indicators of student performance: standardized proficiency tests (PISA, TIMSS, PIRLS), the
55
56
57
58
59
60

METHOD

25

1
2
3 number of correct answers in a test, and non-standardized tests.¹ For achievement inequality, we
4
5 considered measurements of dispersion (e.g., standard deviation, interquartile range, the share of
6
7 top/bottom performers, etc.). For inequality of opportunity, we essentially considered the
8
9 interaction between tracking and the family-background effect.
10
11

12
13 We reviewed only studies written in English and available online. The publications were
14
15 mainly of two types: scientific journal articles and policy/discussion reports. At first, we chose to
16
17 include all studies published since 1982, the year of publication of the last extensive meta-
18
19 analysis on the topic, performed by Kulik and Kulik. Nonetheless, we eventually settled on a
20
21 final sample of only one study from 1992; the rest was published from 2000 onward. We believe
22
23 that this is due to the use of research designs such as multilevel analysis and difference-in-
24
25 differences, which have gained popularity in the social sciences in the last two decades.
26
27 Therefore, we adjusted the time span to include only studies from 2000 onward, excluding
28
29 Hoffer (1992). Other examples of exclusion concern studies that do not cover the population of
30
31 interest (secondary school students) but focus on primary school students (e.g., Duflo et al.,
32
33 2011; Sterbinsky et al., 2006), studies that measure the impact of attending a specific track (e.g.,
34
35 Livingston, 2010; Pop-Eleches & Urquiola, 2011; Van Houtte et al., 2013) instead of the overall
36
37 effect of tracking on the educational system, studies whose dependent variable was not cognitive
38
39 achievement or student performance but, for instance, non-cognitive skills or educational
40
41 transitions (e.g., Elfers, 2011; Guyon et al., 2012; Sampermans et al., 2021), and studies that
42
43
44
45
46
47
48
49
50
51

52
53 ¹ IQ was reported as a dependent variable in only one study: Koerselman (2013). This represents only two statistically
54
55 non-significant estimates for efficiency. Although IQ can express the ability of a student to learn, justifying its use
56
57 only as an independent variable, Kenneth et al. (2015) contend that it can vary over time depending on the amount of
58
59 education received and the social environment. We considered this debate to be unresolved and respected the author's
60
choice to treat IQ as a dependent variable.

METHOD

26

1
2
3 address different research questions and/or employ a different analytical design (Banerjee, 2017;
4
5 Beattie, 2017; Mickelson & Everett, 2008).
6
7

Mean effect-size calculations

8
9
10
11 The studies that analyzed the effects of tracking, in general, applied various forms of
12 regression models and reported β coefficients and standard errors for these estimates. However,
13
14 the statistical theory for meta-analysis was mainly developed for experimental designs that
15
16 presuppose independent treatment and control groups and the reporting of means and standard
17
18 deviations for each. Even though Lipsey and Wilson (2001) consider that, when a non-
19
20 experimental design uses a dummy variable as an independent variable, this could be interpreted
21
22 as the mean difference between two groups (i.e., *tracked system = 1*), this might not be true
23
24 when continuous variables such as the age at first tracking or the number of tracks are employed.
25
26 Thus, the calculation of effect sizes here considers the effect of tracking for one unit, be it a
27
28 dummy or a continuous variable. Consequently, to interpret the results, we must bear in mind
29
30 that they measure the difference between tracked and comprehensive systems but also between
31
32 tracked and less-tracked systems. All multivariate models presented in the results section are
33
34 adjusted for whether tracking is operationalized as a dummy or a continuous variable.
35
36
37
38
39
40
41

42 The formulas used to calculate effect sizes are displayed in Appendix B and were taken
43
44 mainly from Borenstein et al. (2009), who propose specific calculations for non-experimental
45
46 designs. Two blocks of formulas for effect-size calculations are shown in two sections of
47
48 Appendix B, (a) for studies that reported already standardized regression coefficients and (b) for
49
50 studies that reported unstandardized regression coefficients.
51
52
53
54
55
56
57
58
59
60

1
2
3 A potential limitation of retrieving β estimates from different multivariate-regression
4 models for meta-analysis is comparability. Borenstein et al. (2009) argue that because the set of
5 controls might vary from one model to the next, the magnitude of the effect would depend on
6 what is being controlled for. We decided to always take the fullest model specification presented
7 by each study, that is, the model with the most controls. Furthermore, we also theorized the
8 mechanisms through which tracking works to account for this problem. We classified studies
9 according to whether they control for crucial intervening variables in the relationship between
10 tracking and student outcomes, such as school/class composition and educational resources.
11
12
13
14
15
16
17
18
19
20
21

22 We elected Hedge's G as our measure to quantify the mean effect size, which is a
23 corrected version of Cohen's D for samples that are too small ($n < 20$). Since some studies use
24 small samples of countries (sometimes as little as eight countries), we correct for this upward
25 Cohen's D bias (Borenstein et al., 2009). The signs of the retrieved β estimates were consistently
26 coded for a positive sign to indicate that tracking improves student achievement and increases
27 inequality. Hedge's G unit offers an intuitive interpretation as it measures effect sizes in the form
28 of standardized mean differences between groups, which can be read as standard deviations.
29
30
31
32
33
34
35
36
37
38

39 **Coding of studies**

40
41
42 It is important to stress that in a meta-analysis, one publication does not always equal one
43 study. For instance, a publication may run the same model for different outcome variables, such
44 as "efficiency" and "inequality." In this case, we ran separate analyses for each of these
45 dimensions. Furthermore, a publication can be split into two or more studies when the sample
46 changes. Here, we follow the exact definition proposed by Lipsey and Wilson (2001, p. 76),
47 which states that a study is "a set of data analyzed under a single research plan from a designated
48 [specific] sample of respondents." This means that one study does not correspond to one
49
50
51
52
53
54
55
56
57
58
59
60

METHOD

28

1
2
3 publication but to a unique combination of author and sample. Since meta-analysis also requires
4
5 binding different populations together for generalizability, by doing this, we do not treat different
6
7 populations as equal simply because they are included in the same publication. Furthermore, we
8
9 gain in sample size at the macro-level (number of studies) and, consequently, in degrees of
10
11 freedom, not restricting our main model excessively. Thus, our sample of 53 publications
12
13 provides 213 estimates for efficiency that can be grouped into 55 studies, as well as 230
14
15 estimates for inequality that can be grouped into 46 studies. Table 2 provides the descriptive
16
17 statistics.
18
19
20
21

22 Both authors scanned all the results retrieved during the search procedure, screened all
23
24 studies identified based on the eligibility criteria, and reviewed all 53 articles to build the dataset.
25
26 We reached 85% agreement on the coding process, favoring the main author's coding in cases of
27
28 disagreement. Table S.1 in the online Supplemental Material contains the references of all 53
29
30 coded publications.
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

Results

Publication bias

A common concern in meta-analyses is publication bias (Lipson & Wilson, 2001; Borenstein et al., 2009). This results from the well-known trend whereby studies that report more statistically significant estimates have a higher probability of being published.

To check for publication bias, we plotted the effect sizes on a funnel plot where the x-axis represents Hedges' G effect size and the y-axis represents the standard error for each effect (see Figure 2 for efficiency and Figure 3 for inequality). Both outcomes revealed asymmetry. In general, this could indicate that the selection of studies for the meta-analysis concentrated on published articles and that cases from unpublished papers were missing (Borenstein et al., 2009, p. 277).

In the present meta-analysis, we rely on a more balanced sample because 100 of the 213 estimates for efficiency (46.9%) and 66 of the 230 estimates for inequality (28.7%) are taken from working and discussion papers, policy reports, book chapters, or doctoral dissertations, that is, works that are not peer-reviewed. Furthermore, we can address this concern later in our meta-regression models using the variable "peer-review" to gather additional insight into the strength of this bias.

RESULTS

30

Mean effect sizes

Educational efficiency: level of student achievement. Table 3 reports the mean effect size of tracking on the average student's achievement. The effect size is not statistically significant at the 95% confidence level, providing support for the hypothesis that tracking does not increase the overall efficiency of an educational system (H1a).² However, this result must be assessed together with the results on inequality, since a 0 value could be due to the fact that students in the higher tracks improve their results, but this improvement is canceled out by the deterioration of performance in the lower tracks. Another explanation for the null effect of tracking might be related to the fact that many studies use international standardized tests such as PISA, which restricts the analysis to student outcomes in grade 9 (15 years old), when tracking starts in some countries. In this case, for many country units in the sample, the effects of tracking have not had time to fully develop yet.

Heterogeneity tests (Cochran's Q) returned statistically significant results, indicating that the studies included in the analysis contain too much variation and, thus, cannot be treated as belonging to the same population (Hak et al., 2016; Borenstein et al., 2009). This test provides the first evidence that, for the studies included in this meta-analysis, a random-effect analysis would be more appropriate than a fixed-effect one.

The tau² parameter provides the magnitude of this variation, such that if we were to draw infinite samples of studies, their variance would be .116 on average. Tests for the detection of

² Although some of the hypotheses outlined in the theoretical section are directional, we opted to rely on nondirectional statistical tests of the null hypothesis in the empirical analysis, for several reasons: 1) for some issues, there are plausible theoretical predictions in both directions (Woessman, 2009); 2) for other aspects, there is little previous theoretical elaboration, and our hypotheses are more explorative; 3) given that the critical rejection values are higher in a nondirectional test than in a directional test, this is a more conservative approach and the most commonly used in practice (Pillemer, 1991; Salkind, 2010); 4) this choice allows us to maintain homogeneity in the type of statistical test used throughout the paper.

RESULTS

31

1
2
3 outliers revealed four out-of-bound estimates for student achievement (see Supplemental Table
4 S.2, available in the online version of the journal, for a graphic representation of the outlier tests
5 performed). We run a model without outliers to determine how much the magnitude of the effect
6 size would change; it gets closer to zero and remains statistically non-significant. I^2 is the
7 proportion of between-study variance that reflects real differences among effect sizes, namely,
8 the amount of variation that cannot be explained by chance only. Hence, a high value of I^2
9 ($>25\%$) means that there is enough non-random variation to be explored in a multivariate meta-
10 regression at a second stage (Hak et al., 2016). Model 2 reports a high value of I^2 (99.6%),
11 justifying the multivariate analysis in the results section.
12
13
14
15
16
17
18
19
20
21
22
23

24 Model 3 applies multilevel techniques to account for the hierarchical structure of the data
25 and the non-independence of estimates taken from the same studies. The 213 estimates of the
26 effect of tracking on educational efficiency (level 1) were nested within 55 studies (level 2). The
27 estimate of the mean effect size remains statistically non-significant. Thus, all in all, the results
28 of three different models indicate that the effect of tracking on educational efficiency (i.e.,
29 whether overall student achievement is greater in a tracked system), oscillates around 0. This
30 does not mean that the variation in results cannot be explained by moderators but merely that
31 tracking studies report neutral, positive, and negative estimates in a rather balanced way.
32
33
34
35
36
37
38
39
40
41
42

43 ***Inequality in educational achievement.*** The same strategy was used to assess the mean effect
44 size of tracking on student-achievement inequality (Table 4). This estimate is positive in sign and
45 statistically significant at the 99.9% confidence level (in all specifications), providing support for
46 the hypothesis that tracking increases the level of inequality in the educational system (H1b).
47
48
49
50
51

52 Model 1 reports a statistically significant estimate of .097 standard deviation (SD)
53 increase in educational inequality. After excluding three outliers (see Supplemental Table S.3
54
55
56
57
58
59
60

RESULTS

32

1
2
3 available online), model 2 reveals a smaller effect of .076 SD. When grouping the 230 estimates
4
5 into 46 studies in model 3, the mean effect size of tracking on inequality grows to .117 SD
6
7 (99.9% CI: .008; .225). The mean effect size increases after grouping estimates into studies
8
9 because all studies now operate with equal, balanced weight in the calculations, which are, thus,
10
11 not biased by the number of estimates that each study reports. Even though the effect size is
12
13 statistically significant with a 99.9% CI, it can be considered small on a meta-analytic scale
14
15 (Cohen, 1988).³
16
17
18
19

20 Still, it is worth considering whether its magnitude has any practical social significance
21
22 (Bernardi et al., 2017). Woessmann (2016) reported that in the PISA scale, for instance, students
23
24 are capable of progressing, on average, from 0.25 to 0.33 SD in a school year. If all studies in our
25
26 sample were to use PISA scores (and the majority does), this would mean that tracked systems
27
28 exacerbate existing inequalities in student achievement from a third to almost half of a school
29
30 year. Whether this increase is remarkable or not will depend on the already existing levels of
31
32 inequality in elementary and primary schools. The top four countries in inequality growth
33
34 between primary (PIRLS 2001) and secondary education (PISA 2003) reported by Hanushek and
35
36 Woessmann (2006, p. 69) are Germany (0.71 SD), Greece (0.3 SD), the Czech Republic (0.25
37
38 SD), and Italy (0.22 SD), all considered early trackers by the authors. Thus, an increase of .117
39
40 SD for the whole variety of tracking systems reported in the literature (not only early trackers)
41
42 seems to be relevant and capable of explaining much of this increase. This magnitude can
43
44 explain, for example, the whole increase in achievement inequality in the Netherlands after
45
46 tracking occurs.
47
48
49
50
51
52
53
54

55 ³ Cohen (1988) has proposed that effect sizes should be considered small between .0 and .29, medium between .3 and
56 .49, and large above .5.
57
58
59
60

Multivariate meta-regressions

The following multivariate meta-regressions set out to explain the variation in effect sizes observed in the population of studies by a set of explanatory variables related to the tracking policy and the features of the studies.

We specified four statistical models with stepwise inclusion of independent variables in addition to the “intercept-only” models provided in the previous section. The first specification (model 1) in Tables 5 (efficiency) and 6 (inequality) include the independent variables related to the policy characteristics and aspects of substantive interest (i.e., type of tracking, subject of student competencies, and reform analysis). Model 2 adds the variables related to the methodological aspects of the research designs (i.e., counterfactual reference category, identification strategy, and whether the study controls for potential mediators). Model 3 adds the overall score for publication quality, and model 4 tests for the year of publication, peer-review, and quality of the research design.

Educational efficiency. Model 1 in Table 5 reveals a statistically significant positive effect of between-school tracking compared to within-school tracking. This indicates that studies focusing exclusively on measuring the impact of between-school tracking tend to report estimates .023 SD larger than studies of within-school tracking (see Table S.4 in the online Supplemental Material for an additional analysis considering only between-school tracking estimates). This confirms our expectation that between-school tracking can increase overall student achievement more than within-school tracking (H2a).

Model 1 also checks whether achievement in one specific cognitive domain (reading, mathematics, or science) is more sensitive to tracking. The conditional differences across subject

RESULTS

34

1
2
3 domains are not statistically significant and are small in size. This does not support the
4
5 hypothesis that mathematics performance is more sensitive to the effects of tracking (H3a).
6
7 Even though not statistically significant, studies that use pooled scores across domains tend to
8
9 report effects much smaller in magnitude (-.047) than studies focusing exclusively on
10
11 mathematics. The estimates related to reading, on the contrary, are slightly larger (.016) than
12
13 those related to mathematics.
14
15

16
17
18 In Model 1, we also investigate whether studies concentrating on the impact of national
19
20 reforms tend to report smaller estimates of the effects of tracking on student achievement. The
21
22 effect size is not statistically significant, which does not confirm the hypothesis that studies of
23
24 national educational reforms tend to report smaller effects of tracking (H4a).
25
26

27
28 Model 2 includes methodological and research-design characteristics. The choice of the
29
30 counterfactual situation (reference category) in a study appears to influence the reported effect of
31
32 tracking but in the opposite direction than expected. Studies that use “comprehensive” systems as
33
34 a reference (instead of “less-tracked” systems) tend to report estimates .071 SD smaller,
35
36 contradicting hypothesis H5a. Looking at the analytical strategy, when a study applies a
37
38 difference-in-differences design, it tends to report smaller estimates (-.041 SD) than those using
39
40 cross-sectional multilevel strategies, thus providing support for hypothesis H6a. This result could
41
42 suggest that this analytical design is more effective in ruling out confounders that may inflate the
43
44 size of tracking estimates. This estimate remains significant until model 4, increasing its
45
46 magnitude to -.154.
47
48
49
50

51
52 Model 2 further indicates that studies adopting an extensive model specification that
53
54 controls for educational resources report on average -.022 SD smaller effects on the level of
55
56 student achievement. Studies that control both for educational resources and school/class
57
58
59
60

RESULTS

35

1
2
3 composition also produce smaller estimates than studies that control for neither of these
4
5 mechanisms. However, studies that control only for school/class composition do not report
6
7 significantly smaller effect sizes, suggesting that a major mechanism of the tracking policy
8
9 works through the unequal distribution of educational resources across different tracks and/or
10
11 classrooms. These estimates provide support for hypothesis H7a.
12
13
14

15 Model 3 tests whether the reported effects of tracking on efficiency depend on the overall
16
17 quality of the studies (i.e., the scores given to each study according to the evaluation matrix and
18
19 presented in Table 1). Model 4 instead tests whether the effects of tracking are associated with
20
21 the year of publication (from 2000 to 2021), the status of the study (peer-reviewed or not), and
22
23 whether tracking effects depend on the quality of the research design (a sub-component of the
24
25 overall quality score), which we believe to be the most objective item in the evaluation matrix
26
27 and could isolate possible subjectivity contained in the other items. The results indicate that more
28
29 recent studies report smaller estimates for the gain in achievement associated with tracking. This
30
31 is in line with the literature that questions the gain in achievement due to specialization of
32
33 tracked educational systems. Thus, the hypothesis that the higher the quality of the study, the
34
35 smaller the effects will be (H8a) is contradicted as higher-quality studies report larger estimates.
36
37 Furthermore, the hypothesis that peer-reviewed studies tend to report larger estimates (H9a) is
38
39 not supported.
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

RESULTS

36

1
2
3 ***Educational inequality.*** The same models were applied almost identically to explain the
4
5 variation in the effect sizes of tracking on student achievement inequality (Table 6). The only
6
7 difference is the inclusion of a dummy variable that accounts for the fact that studies measured
8
9 inequality in two ways, through the dispersion of student performance (“achievement
10
11 inequality,” used as reference category) or through the effect of family background on
12
13 achievement (“inequality of opportunity”).
14
15

16
17 Model 1 reports a significant negative effect of -.40 SD of between-school tracking on
18
19 inequality. Thus, studies on the consequences of between-school tracking for inequality tend to
20
21 report smaller estimates than studies focusing on within-school tracking, contradicting
22
23 hypothesis H2b (see Table S.5 in the online Supplemental Material for an additional analysis
24
25 considering only between-school tracking estimates). Nonetheless, the importance of informal
26
27 and within-school forms of tracking for inequality of opportunities in education has been
28
29 highlighted in some recent cross-national studies (Blossfeld et al., 2016; Triventi et al., 2020). A
30
31 possible explanation for this result can be found in the stigmatization and self-fulfilling prophecy
32
33 theory (Merton, 1968). Students that are placed into less-advanced classes still have to share the
34
35 school environment with students in the more advanced classes, possibly creating a feeling of
36
37 inferiority. Furthermore, teachers within the school deal with students who are labeled as less
38
39 advanced beforehand, sometimes via decisions taken by the teachers themselves, possibly
40
41 lowering teaching standards. Additionally, between-school tracking, in principle, is mostly
42
43 related to the optimal allocation of students according to vocational, technical, and academic
44
45 profiles to facilitate the school-to-work transition. Instead, ability grouping is mostly related to
46
47 dividing students based on achievement, concentrating high-performers in more advanced
48
49 courses.
50
51
52
53
54
55
56
57
58
59
60

RESULTS

37

1
2
3 Still, in model 1, similarly to the efficiency results, the effects of tracking on inequality
4 appear to be similar across subject domains. This provides no support for the hypothesis that
5 inequality in mathematics scores is more sensitive to the effects of tracking (H3b.1). Model 1 still
6 reports a non-statistically significant estimate for studies that operationalized their measurement
7 of inequality through the influence of family background on achievement, compared to measures
8 of achievement dispersion. This does not confirm the hypothesis that the effects of tracking tend
9 to be weaker on inequality of opportunity than on measurements of achievement dispersion
10 (H3b.2).
11
12
13
14
15
16
17
18
19
20
21

22 Reform studies tend to report more conservative estimates for inequality in the magnitude
23 of .05 SD, keeping its magnitude and significance until model 4. Therefore, the hypothesis that
24 studies of national educational reforms tend to report smaller effects of tracking on inequality
25 (H4b) is supported.
26
27
28
29
30
31

32 Model 2 contradicts the expectation that the effects of tracking tend to be larger when
33 studies use a comprehensive educational system as their counterfactual (H5b). Moreover, this
34 model does not corroborate the hypothesis that difference-in-differences designs tend to produce
35 smaller estimates than cross-sectional multilevel designs (H6b) because the estimates are not
36 statistically significant, even though their signs point in the hypothesized direction. Nevertheless,
37 it provides evidence that the few studies that employed very simple identification strategies (i.e.,
38 one-level regression models) tend to report larger estimates for tracking effects, ranging from
39 .242 (model 2) to .305 SD (model 4). Also, no support is provided for the hypothesis that when a
40 study controls for educational resources and/or school/class composition indicators, the effect
41 sizes of tracking tend to be smaller for inequality (H7b).
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

RESULTS

38

1
2
3 Model 3 presents a significant positive estimate for the overall quality of the study.
4
5 Better-evaluated studies tend to report estimates .084 SD larger, contradicting our expectation
6
7 that the higher the quality of the study, the lower the estimates. In model 4, when we test this
8
9 hypothesis, substituting the overall evaluation score for its sub-component focused on the quality
10
11 of the research design, the estimate increases to .094 and is still significant at the 95% CI. An
12
13 explanation is that the quality of the study and the peer-review process could be intertwined.
14
15 Articles that are published tend to be evaluated as higher quality. This would once again lend
16
17 credence to the file-drawer problem, contradicting hypothesis H8b (i.e., the higher the quality of
18
19 the study, the smaller the effects) and not supporting hypothesis H9b (i.e., peer-reviewed studies
20
21 tend to report larger estimates). Table S.5 in the online Supplemental Material gives an overview
22
23 of the testing of all the hypotheses.
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

Discussion

We analyzed the effect of tracking on two important outcomes of contemporary educational systems: efficiency and inequality. By “efficiency,” we mean the overall (average) level of student achievement, whereas “inequality” refers to achievement dispersion (an indicator of interpersonal inequality) and inequality of opportunity (measuring the strength of the relationship between social background and achievement). Our main finding, which we obtained by applying meta-analytical techniques to studies from the last two decades, is that the mean effect size of tracking on efficiency is null, whereas it is positive for inequality. This evidence provides no support for the existence of an “equality-efficiency trade-off” (Skopek et al., 2019, p. 224), that is, the need to sacrifice equality to improve the overall performance of the educational system. Instead, this finding suggests that the stream of literature that emphasizes the role of tracking in enhancing both student achievement dispersion and inequality of opportunity relies on more solid empirical evidence than the theoretical arguments suggesting that tracking increases efficiency. Thus, our results indicate that de-tracking reforms, which postpone tracking, reduce the number of tracks, or smooth out the distinctions across tracks, have the potential to reduce inequality in educational opportunities based on social background without harming overall student achievement. Of course, our finding only considers cognitive-related achievement, and it should be complemented with considerations about the labor-market effects of school tracking and the institutional specificities of each country (Shavit & Muller, 1998). Moreover, future studies should start considering the possible effects of tracking on students’ socioemotional development, a clear gap in the tracking literature.

As regards the choice of tracking policy, the between-school modality seems to be more effective in terms of specialization while producing less inequality than the within-school type,

DISCUSSION

40

1
2
3 but a caveat is warranted. Van de Werfhorst (2019) analyzed educational reforms in 21 European
4 countries and found that de-tracking educational reforms reduce inequality of opportunity at the
5 expense of lowering the achievement of more advantaged students. This means that there could
6 be a risk that children from higher socioeconomic backgrounds will face deterioration in their
7 learning progress if a de-tracking reform is implemented. Alternative policies should be put in
8 place to avoid this drawback, such as stimulating students to take on a more proactive role, for
9 example, by making them tutor low-performing students so that they may not be affected
10 negatively (Hattie, 2009).
11
12
13
14
15
16
17
18
19
20
21

22 Important methodological remarks concern the fact that longitudinal designs, in general,
23 report more conservative estimates than cross-sectional multilevel designs. This was confirmed
24 for studies that analyzed the effect of educational reforms on inequality and studies that used a
25 difference-in-differences design to identify the effects of tracking on efficiency. Multilevel
26 techniques, for their part, are capable of correcting for part of the bias when compared to
27 standard one-level multivariate regressions. Additionally, we provided evidence from two
28 different tests that publication bias is present in the tracking literature and needs to be addressed
29 by the scientific publication system to minimize the file-drawer problem.
30
31
32
33
34
35
36
37
38
39
40

41 Another contribution of this meta-analysis is theoretical. The debate in the field of
42 sociology of education has moved from an input-output black-box type of research to one that
43 investigates the mechanisms through which macro-systemic variables work. We found evidence
44 that the presence of a specialization effect of tracking on efficiency depends on the type of
45 controls included in the analysis (i.e., educational resources, school/class composition, or both).
46 This reveals a need for more exchanges between the literature on tracking and the literature
47 discussing the distribution of educational resources between schools as well as the effects of
48
49
50
51
52
53
54
55
56
57
58
59
60

DISCUSSION

1
2
3 school/class composition and peers on student achievement. Tracking likely works partly
4
5 through these factors, and its effects should not be discussed without taking into account these
6
7 mechanisms.
8
9

10
11 Lastly, our analysis showed no evidence of a differential effect of tracking across
12
13 subjects. This probably would require more statistical power. Another limitation of this meta-
14
15 analysis is that it focuses on the short-term effects of tracking on cognitive outcomes only, that
16
17 is, how tracking affects educational achievement during the school phase, ignoring the transition
18
19 to higher education and the job market. Future research should also concentrate on systematizing
20
21 the effects of tracking on longer-term outcomes when more studies are available. This would
22
23 expand the debate to other relevant questions, such as whether it is worth tolerating some level of
24
25 short-term inequality if the professional placement of students is assured later, and whether the
26
27 inequalities reproduced and intensified by tracked educational systems remain in access to higher
28
29 education and higher-status positions in the labor market. Moreover, future investigations should
30
31 look at the effects of various forms of tracking on students' socio-emotional competencies and
32
33 civic engagement, which have been covered by a limited number of studies (e.g., Witschge &
34
35 van de Werfhorst, 2020; Korthals et al., 2021; Österman, 2021) even though they are
36
37 increasingly considered important for individuals' success in life (OECD 2015). Lastly, the
38
39 literature seems to have somehow overlooked the economic aspects related to tracking policies.
40
41 The financial cost of maintaining a tracked system instead of a comprehensive one is rarely
42
43 discussed. Taking into consideration all these aspects is crucial to improving our overall
44
45 understanding of how educational arrangements work and helping policymakers make informed
46
47 decisions based on solid empirical evidence.
48
49
50
51
52
53
54
55
56
57
58
59
60

REFERENCES

References

- 1
2
3
4
5
6 Ammermüller, A. (2005). *Educational opportunities and the role of institutions* (ZEW Discussion Paper
7
8 No. 05-44). Retrieved from the Econstor website:
9
10 <https://www.econstor.eu/bitstream/10419/24135/1/dp0544.pdf>
11
12
13
14 Angrist, J. D., & Pischke, J. S. (2008). *Mostly harmless econometrics: An empiricist's companion*.
15
16 Princeton, NJ: Princeton University Press. doi:10.2307/j.ctvc4j72
17
18
19 Ariga, K., & Brunello, G. (2007). *Does Secondary School Tracking Affect Performance? Evidence from*
20
21 *IALS* (Discussion paper n. 630). Retrieved from the Research Gate network:
22
23 [https://www.researchgate.net/publication/5136980_Does_Secondary_School_Tracking_Affect_P](https://www.researchgate.net/publication/5136980_Does_Secondary_School_Tracking_Affect_Performance_Evidence_from_IALS)
24
25 [erformance_Evidence_from_IALS](https://www.researchgate.net/publication/5136980_Does_Secondary_School_Tracking_Affect_Performance_Evidence_from_IALS)
26
27
28
29 Ayalon, H. (2006). Nonhierarchical curriculum differentiation and inequality in achievement: A
30
31 different story or more of the same? *Teachers College Record*, 108(6), 1186-1213.
32
33 doi:10.1111/j.1467-9620.2006.00690.x
34
35
36
37 Ayalon, H., & Gamoran, A. (2000). Stratification in academic secondary programs and educational
38
39 inequality in Israel and the United States. *Comparative Education Review* 44(1), 54-80.
40
41 doi:10.1086/447591
42
43
44 Banerjee, N. (2017). Student–Teacher Ethno-Racial Matching and Reading Ability Group Placement in
45
46 Early Grades. *Education and Urban Society*. V. 51 (3), pp.: 395-422.
47
48 doi:10.1177/0013124517721948
49
50
51
52
53
54
55
56
57
58
59
60

REFERENCES

- 1
2
3 Barg, K. (2013). The influence of students' social background and parental involvement on teachers'
4 school track choices: Reasons and consequences. *European Sociological Review*, 29(3), 565–
5 579. doi:[10.1093/esr/jcr104](https://doi.org/10.1093/esr/jcr104)
6
7
8
9
10
11 Bauer, P., & Riphahn, R. T. (2006). Timing of school tracking as a determinant of intergenerational
12 transmission of education. *Economics Letters*, 91(1), 90–97. doi:[10.1016/j.econlet.2005.11.003](https://doi.org/10.1016/j.econlet.2005.11.003)
13
14
15
16 Beattie, I. R. (2014). Tracking Women's Transitions to Adulthood: Race, Curricular Tracking, and
17 Young Adult Outcomes. *Youth and Society*. First Published March 20, 2014.
18
19
20
21
22
23
24 Bernardi, F., Chakhaia, L., & Leopold, L. (2017). 'Sing me a song with social significance': The
25 (mis)use of statistical significance testing in European sociological research. *European*
26
27
28
29
30
31
32 Betts, J. R. & Grogger, J. (2003). The impact of grading standards on student achievement, educational
33 attainment, and entry-level earnings. *Economics of Education Review*, 22(4), 343-352.
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
- Betts, J. R. (2001). The economics of tracking in education. In Hanushek, Eric A., Machin, S. & Woessmann, L. (Eds.), *Handbook of the economics of education*, Vol. 3 (pp. 341-381), Amsterdam, The Netherlands: North-Holland.
- Betts, J. R., & Shkolnik, J. L. (2000). The effects of ability grouping on student achievement and resource allocation in secondary schools. *Econ. Educ. Rev.* 19(1), 1-15. doi:[10.1016/S0272-7757\(98\)00044-2](https://doi.org/10.1016/S0272-7757(98)00044-2)

REFERENCES

- 1
2
3 Blossfeld, H. P., Bucholz, S., Skopek, J., & Triventi, M. (2016). *Models of secondary education and*
4 *social inequality: An international comparison*. Cheltenham, UK and Northampton, MA, USA:
5 Edward Elgar Publishing. doi:[10.4337/9781785367267](https://doi.org/10.4337/9781785367267)
6
7
8
9
10
11 Bol, T., & Van de Werfhorst, H. (2013). The measurement of tracking, vocational orientation, and
12 standardization of educational systems. *Comparative Education Review*, 57(2), 285-308.
13
14
15
16 Bol, T., Witschge, J., Van de Werfhorst, H. G., & Dronkers, J. (2014). Curricular tracking and central
17 examinations: Counterbalancing the impact of social background on student achievement in 36
18 countries. *Social Forces* 92(4), 1545-1573. doi:[10.1093/sf/sou003](https://doi.org/10.1093/sf/sou003)
19
20
21
22
23
24 Bonesronning, Hans & Falch, Torberg & Strom, Bjarne, 2005. Teacher sorting, teacher quality, and
25 student composition. *European Economic Review*, 49(2), 457-483. doi:[10.1016/S0014-](https://doi.org/10.1016/S0014-2921(03)00052-7)
26 [2921\(03\)00052-7](https://doi.org/10.1016/S0014-2921(03)00052-7)
27
28
29
30
31 Borenstein, M., Hedges, L. V., Higgins, J. P. T., & Rothstein H. R. (2009). *Introduction to Meta-*
32 *Analysis*. West Sussex, UK: John Wiley & Sons Ltd. doi:[10.1002/9780470743386](https://doi.org/10.1002/9780470743386)
33
34
35
36 Bourdieu, P., & Passeron, J. C. (1970). *La Reproduction*. Paris : Éditions de Minuit.
37
38
39
40 Brantlinger, A. (2013). Between politics and equations: Teaching critical mathematics in a remedial
41 secondary classroom. *American Educational Research Journal*, 50(5), 1050-1080.
42
43
44
45
46
47
48 Breen, R., & Goldthorpe, J. H. (1997). Explaining educational differentials: Towards a formal rational
49 action theory. *Rationality and Society*, 9(3), 275-305. doi:[10.1177/104346397009003002](https://doi.org/10.1177/104346397009003002)
50
51
52
53
54
55
56
57
58
59
60

REFERENCES

- 1
2
3 Brunello, G., & Checchi, D. (2007). Does school tracking affect equality of opportunity? New
4 international evidence. *Economic Policy*, 22(52), 781–861.
5
6 doi:[10.1111/j.14680327.2007.00189.x](https://doi.org/10.1111/j.14680327.2007.00189.x)
7
8
9
10
11 Burger, K. (2016). Intergenerational transmission of education in Europe: Do more comprehensive
12 education systems reduce social gradients in student achievement? *Research in Social*
13 *Stratification and Mobility*, 44, 54–67. doi:[10.1016/j.rssm.2016.02.002](https://doi.org/10.1016/j.rssm.2016.02.002)
14
15
16
17
18
19 Burks, L. C. (1994). Ability group level and achievement. *School Community Journal*, 4(1), 11-24.
20
21
22 Clifford, P., & Heath, A. (1984). Selection does make a difference. *Oxford Review of Education*, 10(1),
23 85-97. doi:[10.1080/0305498840100108](https://doi.org/10.1080/0305498840100108)
24
25
26
27 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ:
28 Lawrence Earlbaum Associates.
29
30
31
32
33 Cohn, L. D. (1991). Sex differences in the course of personality development: A meta-analysis.
34 *Psychological bulletin*, 109(2), 252-266. doi:[10.1037/0033-2909.109.2.252](https://doi.org/10.1037/0033-2909.109.2.252)
35
36
37
38 Conn, K. M. (2017). Identifying effective education interventions in Sub-Saharan Africa: A meta-
39 analysis of impact evaluations. *Review of Educational Research*, 87(5), 863-898.
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- Cook, R. R. (1924). A study of the results of homogeneous grouping of abilities in high-school classes.
In Whipple, G. M. (Ed.) *The education of gifted children* (pp. 302-12). Bloomington, Ill.

REFERENCES

- 1
2
3 Cooper, H., Nye, B., Charlton, K., & Lindsay, J. (1996). The effects of summer vacation on achievement
4 test scores: A narrative and meta-analytic review. *Review of Educational Research*, 66(3), 227–
5 268. doi:[10.2307/1170523](https://doi.org/10.2307/1170523)
6
7
8
9
10
11 Deci, E. L, Koestner, R., & Ryan, R. M. (2001). Extrinsic rewards and intrinsic motivation in education:
12 Reconsidered once again. *Review of Educational Research*, 71(1), 1-27.
13
14
15 doi:[10.3102/00346543071001001](https://doi.org/10.3102/00346543071001001)
16
17
18 Dronkers, J., Van der Velden, R. & Dunne, A. (2012). Why are migrant students better off in certain
19 types of educational systems or schools than in others? *European Educational Research Journal*.
20
21
22
23
24
25
26
27 Duflo, E., Dupas, P., & Kremer, M. (2011). Peer effects, teacher incentives, and the impact of tracking:
28 Evidence from a randomized evaluation in Kenya. *American Economic Review*, 101(5), 1739–
29 1774. Doi: [0.1257/aer.101.5.1739](https://doi.org/0.1257/aer.101.5.1739)
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
- Dunne, A. (2010). *Dividing lines: Examining the relative importance of between- and within-school differentiation during lower secondary education* (Doctoral dissertation). Retrieved from the European University Institute website: <https://cadmus.eui.eu/handle/1814/14497>.
- Dupriez, V., Dumay, X., & Vause, A. How do school systems manage pupils heterogeneity? *Comparative Education Review* 52(2), 245-273. doi:[10.1086/528764](https://doi.org/10.1086/528764)
- Dustmann, C. (2004). Parental background, secondary school track choice, and wages. *Oxford Economic Papers*, 56(2), 209-230. doi:[10.1093/oep/gpf048](https://doi.org/10.1093/oep/gpf048)
- Esser, H. (2016). The model of ability tracking—Theoretical expectations and empirical findings on how educational systems impact on educational success and inequality. In *Models of secondary education and social inequality*. Edward Elgar Publishing.

REFERENCES

- 1
2
3 Felouzis, G., & Charmillot, S. (2013). School tracking and educational inequality: A comparison of 12
4 education systems in Switzerland. *Comparative Education*, 49(2), 181-205.
5
6 doi:[10.1080/03050068.2012.706032](https://doi.org/10.1080/03050068.2012.706032)
7
8
9
10 Galindo-Rueda, F., & Vignoles, A. F. (2005). The heterogeneous effect of selection in secondary
11 schools: Understanding the changing role of ability (IZA Discussion Paper No. 1245). Retrieved
12 from the Ideas website: <https://ideas.repec.org/p/cep/ceedps/0052.html>
13
14
15
16
17
18 Gambetta, D. (1987). *Studies in rationality and social change. Were they pushed or did they jump?*
19
20 *Individual decision mechanisms in education*. Cambridge University Press.
21
22
23
24 Gamoran, A., & Mare, R. D. (1989). Secondary school tracking and educational inequality:
25
26 Compensation, reinforcement, or neutrality? *American Journal of Sociology* 94(5), 1146-1183.
27
28 doi:[10.1086/229114](https://doi.org/10.1086/229114)
29
30
31 Gamoran, A., Berends, M. The Effects of Stratification in Secondary Schools: Synthesis of Survey and
32
33 Ethnographic Research *Review of Educational Research*, vol. 57, 4: pp. 415-435. , First
34
35 Published Dec 1, 1987. doi:[10.3102/00346543057004415](https://doi.org/10.3102/00346543057004415)
36
37
38
39 Guyon, N., Maurin, E., & McNally, S. (2012). *The Effect of Tracking Students by Ability into Different*
40
41 *Schools A Natural Experiment*. doi:<http://statline.cbs.nl/Statweb/publication>
42
43
44 Hak, T., Van Rhee, H. J., & Suurmond, R. (2016). *How to interpret results of meta-analysis* (Version
45
46 1.0) Rotterdam, The Netherlands: Erasmus Rotterdam Institute of Management.
47
48 doi:[10.2139/ssrn.3241367](https://doi.org/10.2139/ssrn.3241367)
49
50
51 Hallinan, M. T. (1994). School differences in tracking effects on achievement. *Social Forces*, 72(3),
52
53 799–820. doi:[10.1093/sf/72.3.799](https://doi.org/10.1093/sf/72.3.799)
54
55
56
57
58
59
60

REFERENCES

- 1
2
3 Hanushek, E. A., & Woessmann, L. (2006). Does educational tracking affect performance and
4 inequality? Differences-in-differences evidence across countries. *Economic Journal*, 116(510),
5 C63-C76. doi:[10.1111/j.1468-0297.2006.01076.x](https://doi.org/10.1111/j.1468-0297.2006.01076.x)
6
7
8
9
10
11 Hoffer, T. B. (1992). Middle school ability grouping and student achievement in science and
12 mathematics. *Educational Evaluation Policy Analysis*, 14(3), 205-227.
13
14
15 doi:[10.3102/01623737014003205](https://doi.org/10.3102/01623737014003205)
16
17
18
19 Horn, D. (2013). Diverging performances: the detrimental effects of early educational selection on
20 equality of opportunity in Hungary. *Research in Social Stratification and Mobility*, 32, 25-43.
21
22
23 doi:[10.1016/j.rssm.2013.01.002](https://doi.org/10.1016/j.rssm.2013.01.002)
24
25
26 Horn, D., Balázsi, I., Takács, S., & Zhang, Y. (2007). Tracking and inequality of learning outcomes in
27 Hungarian secondary schools. *PROSPECTS*, 36, 433-446. doi:[10.1007/s11125-006-9003-9](https://doi.org/10.1007/s11125-006-9003-9)
28
29
30
31 Huang, M-H. (2009). Classroom homogeneity and the distribution of student math performance: A
32 country-level fixed-effects analysis. *Social Science Research*, 38,81–791.
33
34
35
36 doi:[10.1016/j.ssresearch.2009.05.001](https://doi.org/10.1016/j.ssresearch.2009.05.001)
37
38
39 Jackson, M. (2013). *Determined to succeed? Performance versus choice in educational attainment*.
40
41
42
43 Stanford: Stanford University Press. doi:[10.1515/9780804784481](https://doi.org/10.1515/9780804784481)
44
45
46
47 Jakubowski, M. (2010). Institutional tracking and achievement growth: Exploring difference-in-
48 differences approach to PIRLS, TIMSS, and PISA data. In J. Dronkers (Ed.), *Quality and*
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

REFERENCES

- 1
2
3 Kendler K. S., Turkheimer, E., Ohlsson, H., Sundquist, J., & Sundquist, K. (2015). Family environment
4 and the malleability of cognitive ability: A Swedish national home-reared and adopted-away co-
5 sibling control study. *PNAS* 112 (15) pp. 4612-4. doi:[10.1073/pnas.1417106112](https://doi.org/10.1073/pnas.1417106112)
6
7
8
9
10 Kerr, S. P., Pekkarinen, T., & Uusitalo, R. (2013). School tracking and development of cognitive skills.
11
12 *Journal of Labour Economics*, 31(3), 577–602. doi:[10.1086/669493](https://doi.org/10.1086/669493)
13
14
15
16 Koerselman, K. (2013). Incentives from curriculum tracking. *Economics of Education Review* 32, 140–
17
18 150. doi:[10.1016/j.econedurev.2012.08.003](https://doi.org/10.1016/j.econedurev.2012.08.003)
19
20
21 Korthals, R., Schils, T., & Borghans, L. (2021). Track placement and the development of cognitive and
22
23 non-cognitive skills. *Education Economics*, 1-20.
24
25
26
27 Kulik, C.-L. C., & Kulik, J. A. (1982). Effects of ability grouping on secondary school students: A meta-
28
29 analysis of evaluation findings. *American Educational Research Journal*, 19(3), 415–428.
30
31 doi:[10.3102/00028312019003415](https://doi.org/10.3102/00028312019003415)
32
33
34
35 Lareau, A. (1987). Social class differences in family-school relationships: The importance of cultural
36
37 capital. *Sociology of Education*, 60(2), pp. 73-85. doi:[10.2307/2112583](https://doi.org/10.2307/2112583)
38
39
40 Lassibile, G., & Gómez, L. N. (2000). Organization and efficiency of education systems: Some
41
42 empirical findings. *Comparative Education*, 36(1), 7-19. doi:[10.1080/03050060027737](https://doi.org/10.1080/03050060027737)
43
44
45
46 Le Donné, N. (2014). European variations in socioeconomic inequalities in students' cognitive
47
48 achievement: The role of educational policies. *European Sociological Review*, 30(3), 1–15.
49
50 doi:[10.1093/esr/jcu040](https://doi.org/10.1093/esr/jcu040)
51
52
53
54 Lipsey, M. W., & Wilson, D. B. (2001). *Practical Meta-Analysis*. Thousand Oaks, CA: Sage
55
56 Publications, Inc.
57
58
59
60

REFERENCES

- 1
2
3 Lucas, S. R. (1999). *Tracking inequality: stratification and mobility in American high schools*. New
4
5 York: Teachers College Press.
6
7
8
9 Lucas, S. R. (2001) Effectively maintained inequality: Education transitions, track mobility, and social
10
11 background effects. *American Journal of Sociology* 106(6):1642-90. doi:[10.1086/321300](https://doi.org/10.1086/321300)
12
13
14 Manning, A., & Pischke, J.S. (2006). *Comprehensive versus selective schooling in England in Wales:
15
16 What do we know?'* (IZA Discussion Paper No. 2072). doi:[10.3386/w12176](https://doi.org/10.3386/w12176)
17
18
19 Marín-Martínez, F., & Sánchez-Meca, J. (2010). Weighting by inverse variance or by sample size in
20
21 random-effects meta-analysis. *Educational and Psychological Measurement*, 70(1), 56–73.
22
23
24 doi:[10.1177/0013164409344534](https://doi.org/10.1177/0013164409344534)
25
26
27 Merton, R. (1968). The self-fulfilling prophecy. In *Social Theory and Social Structure*. New York: The
28
29 Free Press.
30
31
32 Mickelson, R. A., & Everett, B. J. (2008). Neotracking in North Carolina: How high school courses of
33
34 study reproduce race and class-based stratification. *Teachers College Record*, 110(3), 535–570.
35
36
37 Doi: <https://psycnet.apa.org/record/2009-00714-001>
38
39
40 Oakes, J. (1985). *Keeping track: How schools structure inequality*. New Haven: Yale University Press.
41
42
43 OECD (2015), *Skills for Social Progress: The Power of Social and Emotional Skills*, Paris: OECD
44
45 Publishing.
46
47
48 Österman, M. (2021). Can We Trust Education for Fostering Trust? Quasi-experimental Evidence on the
49
50 Effect of Education and Tracking on Social Trust. *Social Indicators Research*, 154(1), 211-233.
51
52
53 Pearl, J. (2009). *Causality*. Cambridge University Press. doi:[10.1017/CBO9780511803161](https://doi.org/10.1017/CBO9780511803161)
54
55
56
57
58
59
60

REFERENCES

- 1
2
3 Pfeffer, F. T. (2008). Persistent inequality in educational attainment and its institutional context.
4
5 *European Sociological Review*, 24(5): 543-565. doi:[10.1093/esr/jcn026](https://doi.org/10.1093/esr/jcn026)
6
7
8 Pillemer, D. B. (1991). One-versus two-tailed hypothesis tests in contemporary educational research.
9
10 *Educational Researcher*, 20(9), 13-17.
11
12
13 Piopiunik, M. (2014). The effects of early tracking on student performance: Evidence from a school
14
15 reform in Bavaria. *Economics of Education Review*, 42, 12-33.
16
17 doi:[10.1016/j.econedurev.2014.06.002](https://doi.org/10.1016/j.econedurev.2014.06.002)
18
19
20 Roloff, J., Klusmann, U., Lüdtke, O. & Trautwein U. (2020) The Predictive Validity of Teachers'
21
22 Personality, Cognitive and Academic Abilities at the End of High School on Instructional
23
24 Quality in Germany: A Longitudinal Study. *AERA Open*, vol. 6, 1, First Published January 8.
25
26 doi: [10.1177/2332858419897884](https://doi.org/10.1177/2332858419897884)
27
28
29
30 Rosenthal, R. (1979). The file drawer problem and tolerance for null results. *Psychological bulletin*,
31
32 86(3), 638-641. doi:[10.1037/0033-2909.86.3.638](https://doi.org/10.1037/0033-2909.86.3.638)
33
34
35
36 Rosenthal, R., & Jacobson, L. (1968). Pygmalion in the classroom. *The urban review*, 3(1), 16-20.
37
38 doi:[10.1007/BF02322211](https://doi.org/10.1007/BF02322211)
39
40
41 Ruhose, J., & Schwerdt, G. (2016). Does early educational tracking increase migrant-native achievement
42
43 gaps? Differences-in-differences evidence across countries. *Economics of Education Review*, 52,
44
45 134-154. doi:[10.1016/j.econedurev.2016.02.004](https://doi.org/10.1016/j.econedurev.2016.02.004)
46
47
48
49 Salkind, N. J. (2010). *Encyclopedia of research design* (Vols. 1-0). Thousand Oaks, CA: SAGE
50
51 Publications, Inc. doi: [10.4135/9781412961288](https://doi.org/10.4135/9781412961288)
52
53
54 Sampermans D., Claes, E. & Janmaat J. G. (2021). Back on track? How civic learning opportunities
55
56 widen the political knowledge gap in a tracked education system. *School Effectiveness and*
57
58
59
60

REFERENCES

- 1
2
3 *School Improvement An International Journal of Research, Policy and Practice* V. 32 (2).
4
5 doi:[10.1080/09243453.2020.1830125](https://doi.org/10.1080/09243453.2020.1830125)
6
7
8 Schlicht, R., Stadelmann-Steffen, I., & Freitag, M. (2010). Educational inequality in the EU: The
9
10 effectiveness of the national education policy. *European Union Politics*, 11(1), 29–59.
11
12 doi:[10.1177/1465116509346387](https://doi.org/10.1177/1465116509346387)
13
14
15
16 Schütz, G., Ursprung, H., & Woessmann, L. (2008). Education policy and equality of opportunity.
17
18 *Kyklos*, 61(2), 279–308. doi:[10.1111/j.1467-6435.2008.00402.x](https://doi.org/10.1111/j.1467-6435.2008.00402.x)
19
20
21 Shavit, Y., & Müller, W. (1998). *From school to work: A comparative study of educational*
22
23 *qualifications and occupational destinations*. Oxford, UK: Oxford: University Press.
24
25
26
27 Siddaway, A., Wood, A., & Hedges, L. (2019). How to do a systematic review: A best practice guide for
28
29 conducting and reporting narrative reviews, meta-analyses, and meta-syntheses. *Annual Review*
30
31 *of Psychology*. 70, 747-770. doi:[10.1146/annurev-psych-010418-102803](https://doi.org/10.1146/annurev-psych-010418-102803)
32
33
34 Skopek, J., Triventi, M., & Buchholz, S. (2019). How do educational systems affect social inequality of
35
36 educational opportunities? The role of tracking in comparative perspective. In Becker, R. (Ed.),
37
38 *Research Handbook on the Sociology of Education* (pp. 214-232). Cheltenham, UK and
39
40 Northampton, MA, USA: Edward Elgar Publishing. doi:[10.4337/9781788110426.00022](https://doi.org/10.4337/9781788110426.00022)
41
42
43
44 Slavin, R. E. (1990) Achievement Effects of Ability Grouping in Secondary Schools: A Best-Evidence
45
46 Synthesis. *Review of Educational Research*, vol. 60, 3: pp. 471-499. , First Published Sep 1,
47
48 1990. doi:[10.3102/00346543060003471](https://doi.org/10.3102/00346543060003471)
49
50
51
52 Snijder, T., & Bosker, R. (1999). *Multilevel analysis: An introduction to basic and advanced multilevel*
53
54 *Modeling*. SAGE Publications London Thousand Oaks New Delhi.
55
56
57
58
59
60

REFERENCES

- 1
2
3 Triventi, M., Skopek, J., Kulic, N., Buchholz, S., & Blossfeld, H. P. (2020). Advantage 'finds its way':
4
5 How privileged families exploit opportunities in different systems of secondary education.
6
7 *Sociology*, 54(2), 237–257. doi:[10.1177/0038038519874984](https://doi.org/10.1177/0038038519874984)
8
9
10
11 Van de Werfhorst, H. G. (2019). Early tracking and social inequality in educational attainment:
12
13 Educational reforms in 21 European countries. *American Journal of Education* 126(1), 65-99.
14
15 doi:[10.1086/705500](https://doi.org/10.1086/705500)
16
17
18 Van de Werfhorst, H. G., & Mijs, J. J. B. (2010). Achievement inequality and the institutional structure
19
20 of educational systems: A comparative perspective. *Annual Review of Sociology*, 36(1), 407-428.
21
22 doi:[10.1146/annurev.soc.012809.102538](https://doi.org/10.1146/annurev.soc.012809.102538)
23
24
25
26 Van Houtte, M., Demanet, J. & Stevens, P.A.J. (2013). Curriculum tracking and teacher evaluations of
27
28 individual students: selection, adjustment or labeling?. *Soc Psychol Educ* 16, 329–352 (2013).
29
30 doi:[10.1007/s11218-013-9216-8](https://doi.org/10.1007/s11218-013-9216-8)
31
32
33
34 Vandenberghe, V. (2006). Achievement effectiveness and equity: the role of tracking, grade repetition
35
36 and inter-school segregation. *Applied Economics Letters*, 13(11), 685-693.
37
38 doi:[10.1080/13504850500404944](https://doi.org/10.1080/13504850500404944)
39
40
41
42 Viechtbauer, W. (2010). Conducting meta-analyses in R with the metafor package. *Journal of Statistical*
43
44 *Software*, 36(3), 1-48. doi:[10.18637/jss.v036.i03](https://doi.org/10.18637/jss.v036.i03)
45
46
47 Waldinger, F. (2007). Does ability tracking exacerbate the role of family background for students' test
48
49 scores? (mimeo).
50
51
52 Whipple, G. M. (1936). *The Grouping of Pupils* (35, pt. 1). Bloomington, Ill.: Public School Publishing,
53
54 National Society for the Study of Education Yearbook,
55
56
57
58
59
60

REFERENCES

- 1
2
3 Witschge, J., & van de Werfhorst, H. G. (2020). Curricular tracking and civic and political engagement:
4
5 Comparing adolescents and young adults across education systems. *Acta Sociologica*, 63(3),
6
7 284-302.
8
9
10
11 Woessmann, L. (2009). *International evidence on school tracking: A review*. CESifo DICE Research
12
13 Report 1. Retrieved from <https://www.ifo.de/DocDL/dicereport109-rr1.pdf>
14
15
16 Woessmann, L. (2009). International evidence on school tracking: A review. CESifo DICE Report, 7(1),
17
18 26-34.
19
20
21 Woessmann, L. (2016). The importance of school systems: Evidence from international differences in
22
23 student achievement. *Journal of Economic Perspectives*. 30(3), 3-32. doi:10.1257/jep.30.3.3
24
25
26 Zimmer, R. (2003). A new twist in the educational tracking debate. *Economics of Education Review*.
27
28 22(3), 307-315. doi:10.1016/S0272-7757(02)00055-9
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

TABLES

Table 1

Evaluation matrix^a

	Criteria	Research design	0 point	1 point	2 points
1	Presentation of results	All	Poor: unclear language, incomplete tables, crucial information missing (sample size, SE or SD, etc.).	Adequate: although not very well textually presented, results are clear in the tables, or vice-versa.	Professional: results presented in clear language with complete and self-explaining tables, graphs, figures, etc.
2	Presentation of research design	All	Vague explanation of procedures and limitations not discussed.	Fair explanation of procedures and/or some limitations discussed.	Adequate explanation of procedures and limitations discussed.
3	Trend/balances	Difference-in-differences	Does not check for previous trend and does not control for balancing variables.	Does not check for previous trend but controls for balancing variables.	Checks for previous trend and controls for balancing variables.
4	Macro-level heterogeneity	Cross-sectional models	Does not account for macro-level unobserved heterogeneity or applies one-level regression.	Applies multilevel regression and controls for at least one macro-level confounder (e.g., GDP).	Applies multilevel and controls for more than one macro-level confounder (e.g., GDP and investment per student).

^a Adapted from Conn (2017).

TABLES

56

Table 2 *Descriptive Statistics*

Variables	Efficiency					Inequality				
	N.	%	Mean	Min.	Max.	N.	%	Mean	Min.	Max.
<i>Policy characteristics</i>										
Within school (ref.)	35	16.4%	0.16	0	1	25	10.9%	0.11	0	1
Between school	178	83.6%	0.84	0	1	205	89.1%	0.89	0	1
Achiev. Dispersion (ref.)	-	-	-	-	-	86	37.4%	0.37	0	1
Opportunity	-	-	-	-	-	144	62.6%	0.63	0	1
<i>Dep. variable oper.</i>										
Mathematics (ref.)	94	44.1%	0.44	0	1	75	32.6%	0.33	0	1
Reading	61	28.6%	0.29	0	1	62	27.0%	0.27	0	1
Science	25	11.7%	0.12	0	1	22	9.6%	0.10	0	1
Pooled	33	15.5%	0.15	0	1	71	30.9%	0.31	0	1
<i>Indep. variable oper.</i>										
Non-reform (ref.)	174	81.7%	0.82	0	1	211	91.7%	0.92	0	1
Reform	39	18.3%	0.18	0	1	19	8.3%	0.08	0	1
Dummy	171	80.30%	0.8	0	1	97	45.50%	0.46	0	1
Continuous	42	19.70%	0.2	0	1	133	62.40%	0.62	0	1
<i>Study characteristics</i>										
Less tracked (ref.)	103	48.4%	0.48	0	1	127	55.2%	0.55	0	1
Comprehensive	110	51.6%	0.52	0	1	103	44.8%	0.45	0	1
Multilevel (ref.)	88	41.3%	0.41	0	1	102	44.3%	0.44	0	1
Diff-in-diff	111	52.1%	0.52	0	1	110	47.8%	0.48	0	1
One level regression	14	6.6%	0.07	0	1	18	7.8%	0.08	0	1
Controls: neither (ref.)	24	11.3%	0.11	0	1	8	3.5%	0.03	0	1
Resources	100	46.9%	0.47	0	1	124	53.9%	0.54	0	1
School/class composition	69	32.4%	0.32	0	1	30	13.0%	0.13	0	1
Both	20	9.4%	0.09	0	1	68	29.6%	0.30	0	1
<i>Publication quality</i>										
Publication year	213	100%	12.07	0	21	230	100%	10.45	0	21
Study overall quality	213	100%	4.67	1	6	192	100%	4.74	1	6
Not peer-reviewed (ref.)	100	46.9%	0.46	0	1	66	28.7%	0.29	0	1
Peer-reviewed	113	53.1%	0.39	0	1	164	71.3%	0.71	0	1
Quality of research design	213	100%	0.93	0	2	230	100%	1.17	0	2

TABLES

57

Table 3 Tracking mean effect size on student achievement level (educational efficiency)

Dep. Variable: Hedge's G	Random-effect models								
	1. Full sample			2. No outliers			3. Nested		
	β	SE	p-val	β	SE	p-val	β	SE	p-val
Mean effect-size	-.024	.025	.320	-.002	.019	.903	-.063	.042	.138
Heterogeneity test	1,492.86	***		1,412.00	***		1,492.18	***	
Tau ²	.116	.014		.065	.009		.087	.024	
I ²	99.76%			99.59%			99.69%		
N. of estimates (studies)	213			209			213 (55)		
Degrees of freedom (df)	212			208			212		

† $p < 0.10$, * $p < .05$, ** $p < .01$, *** $p < 0.001$

Note. Model 1: estimation by "Hedges and Vevea (HE)" method according to recommendations from Marín-Martínez and Sánchez-Meca (2010) for standardized mean difference effect sizes. Model 2: estimation by "Restricted Maximum Likelihood (REML)". Estimates are nested within studies.

TABLES

Table 4*Tracking mean effect sizes on achievement inequality*

Dep. Variable: Hedge's G	Random-effect models											
	1. Full sample			2. No outliers			3. No outliers / Nested					
	β	SE	p-val	β	SE	p-val	β	SE	p-val			
Mean effect-size	.097	***	.022	<.0001	.076	***	.018	<.0001	.117	***	.033	<.001
Heterogeneity test	1782.12	***			1647.51	***			1782.12	***		
Tau ²	.100		.011		.063		.008		.045			
I ²	99.95%				99.93%				99.89%			
N. of estimates (studies)	230				227				230 (46)			
Degrees of freedom (df)	229				226				229			

† $p < 0.10$, * $p < .05$, ** $p < .01$, *** $p < 0.001$

Note. Models 1 and 2: estimation by "Hedges and Vevea (HE)" method according to recommendations from Marín-Martínez and Sánchez-Meca (2010) for standardized mean difference effect sizes. Model 3: estimation by "Restricted Maximum Likelihood (REML)". Estimates are nested within studies

TABLES

Table 5*Meta-regressions on the effects of tracking on educational efficiency with robust standard errors (nested within studies)*

<i>Dep. Variable: Hedge's G</i>	<i>Model 1</i>			<i>Model 2</i>			<i>Model 3</i>			<i>Model 4</i>		
	β	<i>SE</i>	<i>p-val</i>	β	<i>SE</i>	<i>p-val</i>	β	<i>SE</i>	<i>p-val</i>	β	<i>SE</i>	<i>p-val</i>
Intercept	-.086 †	.048	.076	-.011	.065	.867	-.196	.124	.115	.214 †	.125	.087
<i>Policy characteristics</i>												
Between VS within school	.023 ***	.003	<.0001	.022 ***	.003	<.0001	.022 ***	.003	<.0001	.022 ***	.003	<.0001
<i>Outcome subject</i>												
Reading VS mathematics	.016	.014	.254	.015	.014	.270	.014	.014	.301	.014	.014	.319
Science VS mathematics	.002	.003	.384	.002	.003	.388	.002	.003	.384	.002	.003	.383
Pooled VS mathematics	-.047	.079	.553	-.046	.079	.565	-.046	.079	.562	-.047	.079	.556
<i>Institutional variation</i>												
Reform VS other	.097	.102	.341	.126	.104	.227	.104	.105	.322	.125	.101	.215
<i>Research design</i>												
Comprehensive VS less tracked				-.071 ***	.010	<.0001	-.071 ***	.010	<.0001	-.072 ***	.010	<.0001
Method: Diff-in-diff VS multilevel				-.041 ***	.011	.000	-.154 *	.066	.019	-.154 **	.052	.003
One level regression VS multilevel				-.082	.167	.626	-.054	.168	.748	-.137	.163	.402
Controls: Resources VS neither				-.222 *	.104	.032	-.247 *	.105	.018	-.314 **	.105	.003
Composition VS neither				-.063	.119	.597	-.107	.122	.382	-.055	.127	.666
Both VS neither				-.199 †	.103	.054	-.225 *	.104	.031	-.293 **	.105	.005
<i>Publication quality</i>												
Study overall quality							.057 †	.033	.081			
Publication year										-.020 *	.009	.021
Peer-reviewed										-.040	.097	.678
Quality of research design										.115 *	.052	.026
N. estimates (n. studies)	213 (55)			213 (55)			213 (55)			213 (55)		

† $p < 0.10$, * $p < .05$, ** $p < .01$, *** $p < 0.001$.

Note: all models control for whether tracking was measured as a dummy or a continuous variable.

TABLES

60

Table 6*Meta-regressions on the effects of tracking on educational inequality with robust standard errors (nested within studies)*

<i>Dep. Variable: Hedge's G</i>	<i>Model 1</i>			<i>Model 2</i>			<i>Model 3</i>			<i>Model 4</i>		
	β	<i>SE</i>	<i>p-val</i>	β	<i>SE</i>	<i>p-val</i>	β	<i>SE</i>	<i>p-val</i>	β	<i>SE</i>	<i>p-val</i>
Intercept	.138 **	.045	.002	.148 *	.060	.014	-.219	.142	.123	-.023	.102	.822
<i>Policy characteristics</i>												
Between VS within school	-.040 ***	.007	<.0001	-.038 ***	.007	<.0001	-.038 ***	.007	<.0001	-.038 ***	.007	<.0001
<i>Outcome</i>												
Opportunity VS achv	.020	.047	.664	.046	.050	.365	.026	.050	.609	.007	.052	.891
Reading VS mathematics	-.001	.002	.683	-.001	.002	.681	-.001	.002	.691	-.001	.002	.686
Science VS mathematics	-.001	.003	.678	-.001	.003	.674	-.001	.003	.679	-.001	.003	.680
Pooled VS mathematics	.026	.019	.157	.026	.019	.159	.021	.019	.261	.022	.019	.250
<i>Institutional variation</i>												
Reform VS other	-.050 ***	.013	<.0001	-.044 ***	.013	.001	-.045 ***	.013	.000	-.045 ***	.013	.001
<i>Research design</i>												
Comprehensive VS less tracked				-.004 *	.002	.045	-.004 *	.002	.047	-.004 *	.002	.048
Method: Diff-in-diff VS multilevel				-.016	.034	.646	-.025	.034	.461	-.014	.034	.688
One level regression VS multilevel				.244 *	.107	.023	.382 **	.116	.001	.309 **	.111	.005
Controls: Resources VS neither				-.055	.092	.550	-.120	.091	.184	-.081	.102	.425
Composition VS neither				-.098	.131	.454	-.046	.126	.713	-.165	.129	.202
Both VS neither				-.070	.106	.506	-.047	.101	.642	-.069	.105	.513
<i>Publication quality</i>												
Study overall quality							.084 **	.030	.005			
Publication year										.008	.008	.306
Peer-reviewed										.030	.071	.673
Quality of research design										.094 †	.053	.078
N. estimates (n. studies)	230 (46)			230 (46)			230 (46)			230 (46)		

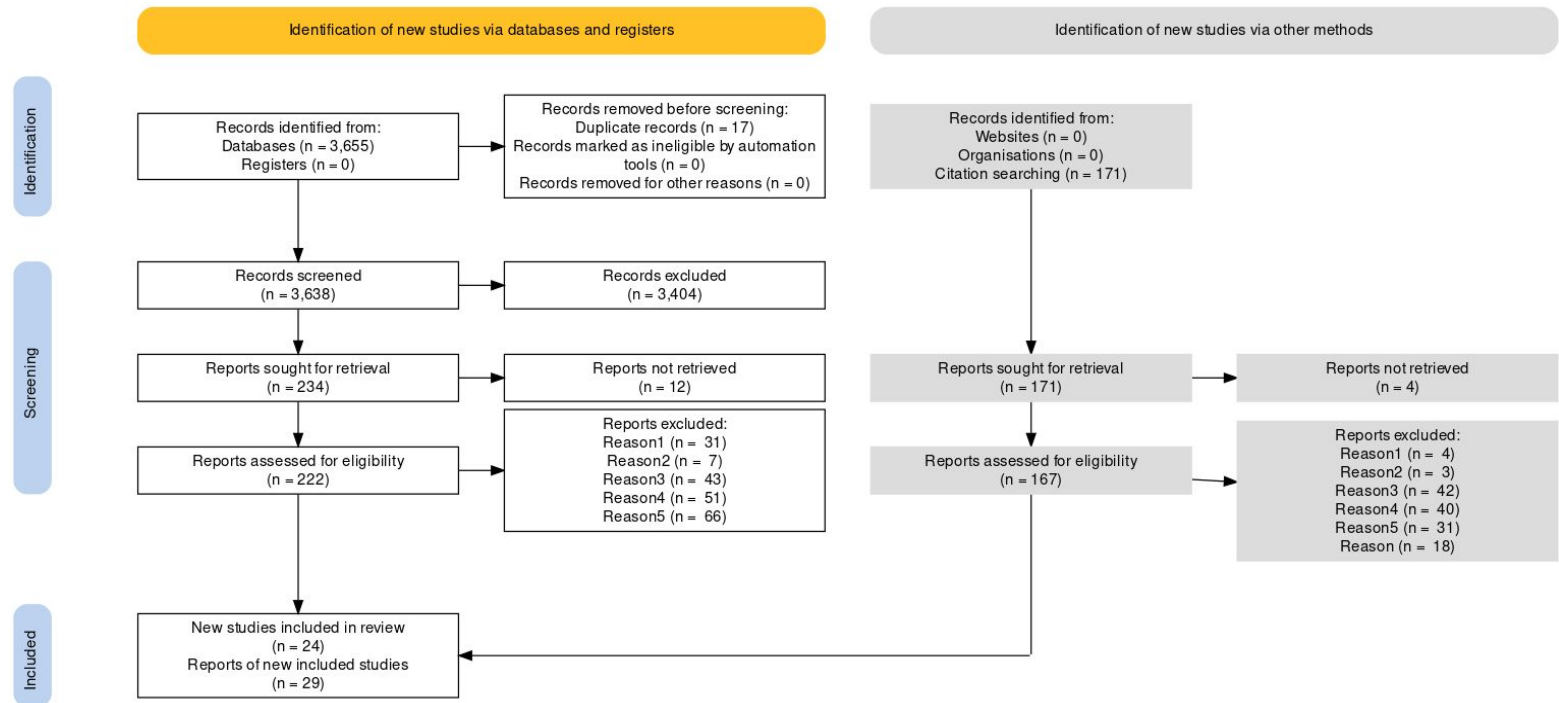
† $p < 0.10$, * $p < .05$, ** $p < .01$, *** $p < 0.001$

Note: all models control for whether tracking was measured as a dummy or a continuous variable.

FIGURES

Figure 1

PRISMA diagram



Reason 1: educational level considered is not secondary education (e.g. primary education or higher education).

Reason 2: out of publication year range (2000-2021).

Reason 3: treatment is not tracking as a macro/institutional characteristic (in most cases it refers to studies focused on comparing students' outcomes across tracks).

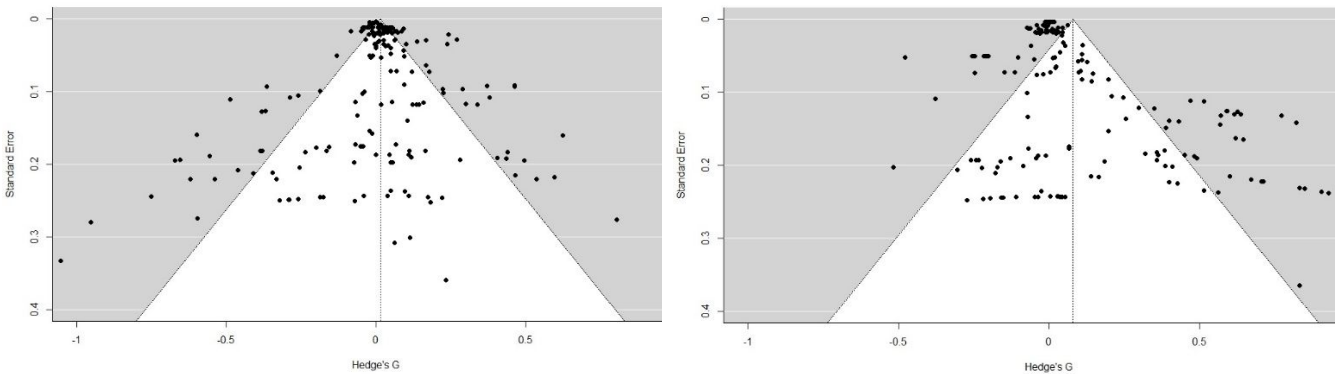
Reason 4: design (it addresses a different research question and employs a different analytical design).

Reason 5: outcome is not achievement/student performance (it might be non-cognitive skills or educational transitions).

Reason 6: duplicated.

FIGURES

62

Figure 2*Funnel plot for educational efficiency (left) and inequality (right)*

*Note. Regression Test for Funnel Plot
Asymmetry: $z = -4.9638$, $p < .0001$
Limit Estimate (as $sei \rightarrow 0$): $b = 0.0854$ (CI: $0.0254, 0.1454$)*

*Test for Funnel Plot Asymmetry: $z = 6.9189$, $p < .0001$
Limit Estimate (as $sei \rightarrow 0$): $b = -0.0330$ (CI: $-0.0852, 0.0192$)*

APPENDICES

Appendix A

Eligibility criteria

Independent variable: tracking (dummy, educational reforms, age of tracking, n. of tracks, length of tracking, share of vocational track, interaction tracking x family background, etc.)

Dependent variable (efficiency): achievement: standardized students proficiency (PISA, TIMSS, PIRLS), IQ tests, n. of correct questions, non-standardized tests, etc.

Dependent variable (inequality of achievement): inequality in achievement: gap between top/bottom performers, standard deviation, interquartile range, the share of top/bottom performers.

Dependent variable (inequality of opportunity): effect of family background on achievement.

Measurement: β coefficients.

Research respondents (population): Secondary students

Research methods: Experimental and quasi-experimental designs.

Language: English

Time frame: 1980 – 2020 (obs.: a meta-analysis on the topic was published in 1982).

Publication type: published journals, books, dissertations, technical reports, conference presentations, policy reports, etc.

Availability: online.

Search terminology

Terms for educational field:

“educ*” OR “student*”

AND

Terms for tracking:

track OR “ability group” OR stream OR sort

AND

Methodological terms:

experiment OR quasi-experiment OR difference-in-difference OR multilevel OR regression OR “propensity score” OR “instrumental variable”

Appendix B

Formulae: From Borenstein et al. (2009), chapter 4 (p. 21).

For regression coefficients reported in a standardized form:

Effect size (Cohen's D):

$$D = \bar{X}_1 - \bar{X}_2 \quad \text{Eq. (B.1)}$$

Where D is the mean difference between two independent groups, \bar{X}_1 is the standardized mean of group 1 and \bar{X}_2 is the standardized mean of group 2, i.e., $\bar{X}_1 - \bar{X}_2$ is equal to the standardized regression coefficient.

Standard deviation of the difference:

$$S_{diff} = SE \times \sqrt{n-1} \quad \text{Eq. (B.2)}$$

Where SE is the standard error of the regression coefficient and n is the sample size at the considered level of the analysis, i.e., n of macro-units for macro-level variables and n of micro-level units for interactions between micro and macro-level variables.

Variance of effect size:

$$V_D = \frac{S_{diff}^2}{n} \quad \text{Eq. (B.3)}$$

Where S_{diff}^2 is the standard deviation of the difference squared and n is the sample size.

Standard error of the effect size:

$$SE_D = \sqrt{V_D} \quad \text{Eq. (B.4)}$$

Where $\sqrt{V_D}$ is the squared root of the variance of the effect size.

For regression coefficients reported in an unstandardized form:

Effect size (Cohen's D):

$$D = \frac{\bar{Y}_{diff}}{S_{within}} = \frac{\bar{Y}_1 - \bar{Y}_2}{S_{within}} \quad \text{Eq. (B.5)}$$

Where $\bar{Y}_1 - \bar{Y}_2$ is the unstandardized regression coefficient and S_{within} is the pooled standard deviation.

Standard deviation of the difference:

$$S_{diff} = SE \times \sqrt{n-1} \quad \text{Eq. (B.6)}$$

Where SE is the standard error of the regression coefficient and n is the sample size for the respective level of analysis.

Pooled standard deviation:

$$S_{within} = \frac{S_{diff}}{\sqrt{2(1-r)}} \quad \text{Eq. (B.7)}$$

Where S_{diff} is the standard deviation of the difference and r is the correlation of the dependent variable in times 1 and 2. No study reported r , so it was set to .5 in order to make the divisor neutral = 1.

Variance of effect size:

$$V_D = \left(\frac{1}{n} + \frac{D^2}{2n} \right) \times 2(1-r) \quad \text{Eq. (B.8)}$$

Where D is the regression coefficient, n is the sample size for the respective level of analysis and r is the correlation of the dependent variable between times 1 and 2 and was set to .5.

Standard error of the effect size:

$$SE_D = \sqrt{V_D} \quad \text{Eq. (B.9)}$$

Bias correction (transformations of Cohen's D to Hedge's G):

Hedge's bias correction factor:

$$J = 1 - \frac{3}{4df - 1} \quad \text{Eq. (B.10)}$$

Where df is equal to $n-2$.

Effect size (Hedge's G):

$$G = J \times D \quad \text{Eq. (B.11)}$$

Where J is Hedge's correction factor and D is Cohen's D .

Variance of effect size (Hedge's G):

$$V_G = J^2 \times V_D \quad \text{Eq. (B.12)}$$

Where J is Hedge's correction factor and V_D is the variance of Cohen's D .

Standard error of effect size (Hedge's G):

$$SE_G = \sqrt{V_G} \quad \text{Eq. (B.13)}$$

Where V_G is the variance of Hedge's G