# The Impact of the 1999 Education Reform in Poland

*Maciej Jakubowski*
*Harry Anthony Patrinos*
*Emilio Ernesto Porta*
*Jerzy Wiśniewski*

## Abstract

Increasing the share of vocational secondary schooling has been a mainstay of development policy for decades, perhaps nowhere more so than in formerly socialist countries. The transition, however, led to significant restructuring of school systems, including a declining share of vocational students. Exposing more students to a general curriculum could improve academic abilities. This paper analyzes Poland's significant improvement in international achievement tests and the restructuring of the education system that expanded general schooling to test the hypothesis that delayed vocational streaming improves outcomes. Using propensity score matching and differences-in-differences estimates, the authors show that delayed vocationalization had a positive and significant impact on student performance on the order of one standard deviation.

# The Impact of the 1999 Education Reform in Poland[1]

Maciej Jakubowski
Organisation for Economic Co-operation and Development (OECD)

Harry Anthony Patrinos
World Bank

Emilio Ernesto Porta
World Bank

Jerzy Wiśniewski
Center for Social and Economic Research (CASE), Poland

---

# 1.      Introduction

The vocationalization of the education curriculum has been a major feature of education plans since the post war period.  It is often argued that vocational skills are needed for job creation, employment and productivity.  The common sense view is that vocational education is necessary for a country to modernize and acquire the technical skills needed for economic development.  A number of reasons have been put forward to argue for increasing the proportion of students in vocational education programs, neatly summarized by Psacharopoulos (1997):

1.  Youth unemployment.  The argument is that with one action, policymakers can take the youth off the streets, and at the same time equip them with skills that could be used later in the labor market.

2.  Instilling technological knowledge.  Starting from the Industrial Revolution it is a common belief that economic progress heavily depends on technological knowhow. The next logical step in this reasoning is that vocational education must expand.

3.  Academically less able students.  There has always been concern with a great number of students who 'are not able' to advance through the school system, especially the academic curriculum of secondary education.  Hence, the provision of vocational education to them would allegedly equip them with something useful to do later in life.

4.  Lack of middle level technicians.  There is no country in the world where a number of specialties are not in 'scarce' supply (for example, plumbers and nurses).  So, it would appear logical to reason that such skills in short supply should be created by the country's vocational schools and training institutions.

5.  Poverty among urban dwellers.  A more modern variant of the above themes is that, given the increased poverty of urban dwellers, the provision of vocational education would give useful skills to the unemployed and make them find productive employment and thus raise their income.

6.  Economic globalization.  Given the fact that frontiers have expanded by means of freer trade and multinationals, this development has to have implications on the nature of vocational education received by the labor force.

Since the post war period many countries have developed vocational education systems, at times diversified (for example, Colombia and Tanzania).  Socialist countries integrated vocational schooling into the overall economic planning system, assigning them to different ministries.  The upshot was that employment was guaranteed in those models.  After the transition began, however, the link between vocational education and employment was broken, leaving vocational students lacking in job opportunities and skills demanded by the workforce.

Indeed, vocationalization has been under attack for many decades.  Psacharopoulos (1987) argues that the social costs of vocational education may not match the social benefits associated with it.  The argument that vocational education would bring industrialization and jobs was challenged early on by Foster (1965), naming it the "vocational school fallacy."  More importantly, the "vocational skills" – what is needed in the world of work, what students must learn to compete – of today are not the old traditional skills linked to specific jobs; rather, they are the critical thinking and "learning to learn" skills (see Murnane et al. 1995), exemplified by success in math, reading and science, for example.

Despite its prominent place in school policy, there has been little rigorous evaluation of vocational school efforts. Much more work has been undertaken on financing, arguing that general skills are a public priority while specific vocational skills should be privately financed or financed by employers (Becker 1964). Even more work has gone into estimating wage effects or returns to schooling for vocational tracks and comparing them to general or academic tracks. Overall, cost-benefit studies show that returns are lower and costs higher (Psacharopoulos and Patrinos 2004).

A small empirical literature does suggest that that there are advantages of targeted vocational *training* – that is not school based – programs (Karlan and Valdivia 2006). Evaluations from the randomized training programs in the United States show modest effects, at best (see, for example, Heckman, Lalonde and Smith 1999). Evidence on the effectiveness of training in developing countries is more limited. Betcherman, Olivas and Dar (2004) review 69 impact evaluations of unemployed and youth training programs, only 19 of which are in developing countries. They find that the impacts in developing countries are more positive than the impacts of programs in the United States and Europe. Most of those programs, however, are not experimental. Card et al. (2007) report on the first randomized evaluation of a job training program in Latin America. The subsidized program in the Dominican Republic showed no impact on employment, a marginally significant impact on hourly wages and on the probability of health insurance coverage, conditional on employment. Attanasio, Kugler and Meghir (2009) evaluate the impact of a randomized training program for disadvantaged youth in Colombia on employment and earnings outcomes. They find that the program raises earnings and employment for both men and women, with larger effects on women. Cost-benefit analysis of these results suggests that the program generates a large net gain, especially for women.

Fewer evaluations – randomized or otherwise – have been undertaken on the impacts of vocational education. Earlier assessments of vocational education programs in a number of countries (for example, Colombia and Tanzania) have shown that most graduates of such schools go to university rather than entering manual occupations (Psacharopoulos and Loxley 1985). In 1991, the upper secondary school two-year vocational programs were transformed into three-year programs as a pilot before the reform was implemented all over the country in 1995. This "natural experiment" was evaluated in terms of years of upper secondary education, university enrolment, and the rate of inactivity. Results suggest positive effects on years of upper secondary education for those who lived in a pilot municipality in 1990. One of the important changes was that the third year in upper secondary vocational education gave individuals general eligibility to continue to higher education. However, the third year did not have a statistically significant effect on the probability to continue to higher education, at least not within six years after completing upper secondary education (Ekström 2002). No rigorous study, to our knowledge, has been undertaken on the learning outcomes associated with vocational secondary schooling.

Poland is a good case for such an evaluation. In 1999, Poland reformed its basic education system in order to raise the level of education in society, increase educational opportunities, and improve the quality of education. The new government at that time restructured basic education by converting the old 8-year primary school that was followed by early vocational tracking, into a 6-year primary education followed by three years of lower

2

general secondary education. Only after 9 years of schooling would a decision about what type of upper secondary education – academic or vocational – be undertaken. In other words, the new system postponed the choice of type of curriculum at the secondary level (general or vocational) for one year. This structural change was accompanied by curricular reform. A concept of core curricula was developed which aimed to provide schools with extensive scope of autonomy and responsibility. A system of examination and tests at the end of primary and lower secondary were introduced.

We use the variation created by the policy change in 1999 to test the impact on test scores over time. Specifically, we estimate a difference in difference model that compares the change in test scores of the likely vocational school students that were able to study in the general, academic track because of the change in school policy. The purpose of our paper is to explain the significant improvement in international achievement tests by Poland in recent years. We find that, on average, that the reform was associated with significant improvements.

In math, Poland improved its score by 0.25 of a standard deviation. In reading the increase is 0.28 of a standard deviation. In science, the scores increased by 0.16 of a standard deviation. We confirm these results using our evaluation model – propensity score matching and differences in differences to create counterfactual scores for the group of likely vocational students in subsequent years – and the OECD's Program for International Student Assessment (PISA), an internationally comparable standardized student test conducted every three years to test reading, mathematics and science achievement of 15-year-olds, data for 2000, 2003 and 2006, using 2000 as the baseline since most of the existing students were continuing their lower secondary schooling under the old system. We explore threats to identification using among other things decomposition analysis We conclude that the reform is associated with an improvement in likely vocational students' scores of about 100 points, or a whole standard deviation. We explore the implications using as well a 2006 special application of PISA in Poland to 16 and 17 year-olds, and warn of the dangers of early vocationalization.

The organization of the paper is as follows. Section 2 describes the policy change in Poland. Section 3 describes the increase in test scores over time. Our hypotheses are presented in Section 4. Section 5 describes our empirical methods and data. Section 6 presents the average impact results. Additional analyses are presented in section 7. Finally, we summarize our conclusions and discuss the implications in section 8.

## 2. Reform of 1998-1999

In 1998, the Minister of Education presented the outline of the reform. The following goals were set (Ministry of National Education 1998):

1. Raising the level of education in society by increasing the number of people with secondary and higher education qualifications;
2. Ensuring equal educational opportunities; and
3. Supporting improvements in the quality of education.

The reform was envisaged to cover the following areas:

- the structure of the education system ranging from nursery school to doctoral studies, including particularly the introduction of a new structure of the school system;

- changes in the methods of administration and supervision;

- a curricular reform comprising the introduction of core curricula and changes in the organization and methods of teaching;

- the establishment of an assessment and examination system independent of the school;

- school finance; and

- the identification of qualification requirements for teachers, which would also be linked with their promotion paths and, the system of remuneration at an adequately high level.

The structural reform led to the introduction of the lower secondary school "gymnasium" as a new type of school. It was most visible to society, thus becoming a symbol of the reform. It was decided that the previous structure of education, comprising the eight-year primary school followed by the four-year secondary school or the three-year vocational school as an option to be chosen, would be replaced with a system described in brief as 6+3+3 (Figure 1). This meant that the duration of education in the primary school would be reduced to 6 years. Following the educational cycle in the primary school, the pupil would continue his/her education in a three-year gymnasium and only upon completion of education in the gymnasium would he/she move on to a three-year secondary school (specialized lyceum) or a two-year vocational school. The structural reform postponed the choice of the direction of education at the secondary level (general or vocational curriculum) for one year. With the clear division of the education system into stages, pupil achievements could be assessed reliably through tests and examinations.

**Figure 1: Structure of the Polish Education System**

| Before the reform of 1999 | After the reform of 1999 |
|---|---|

**Before the reform of 1999**

| age | | grade |
|---|---|---|
| 6 | Zero class (primary schools or kindergartens) | 0 |
| 7 | Comprehensive primary schools | I |
| 8 | | II |
| 9 | | III |
| 10 | | IV |
| 11 | | V |
| 12 | | VI |
| 13 | | VII |
| 14 | | VIII |

**Entrance exam**

| age | | | | grade |
|---|---|---|---|---|
| 15 | General secondary schools (*liceum*) | Secondary vocational schools (*technikum*) | Basic vocational schools | I |
| 16 | | | | II |
| 17 | | | | III |
| 18 | | | | IV |
| 19 | **Matura** | | | V |

**Matura**

**After the reform of 1999**

| age | | grade |
|---|---|---|
| 6 | Zero class (primary schools or kindergartens) | 0 |
| 7 | Comprehensive primary schools | I |
| 8 | | II |
| 9 | | III |
| 10 | | IV |
| 11 | | V |
| 12 | | VI |

**Final test**

| age | | grade |
|---|---|---|
| 13 | Comprehensive lower secondary schools (*gimnazjum*) ISCED 2A | I |
| 14 | | II |
| 15 | | III |

**Final exam**

| age | | | | | grade |
|---|---|---|---|---|---|
| 16 | General secondary schools ISCED 3A | Profiled general secondary schools ISCED 3B | Secondary vocational schools ISCED 3B | Basic vocational schools ISCED 3C | I |
| 17 | | | | | II |
| 18 | | | | | III |
| 19 | **Matura** | **Matura** | | | IV |

**Matura**

The reformers assumed that the gymnasia would allow Poland to raise the level of education particularly in rural areas where the schools were small. The new lower secondary schools would be larger (at least 150 pupils), well equipped, and employ teachers with adequate qualifications. Since the number of pupils in the school varies with the school-catchment area, the establishment of gymnasia involved, unavoidably, the reorganization of the school network. The structural reform did not cover nursery schools and did not result in the lowering of the compulsory school age (7 years).

The arguments given to justify the need for change were twofold. First, the new division of the school career stages would allow for better adjustment of teaching methods and curricula to the specific needs of pupils of various ages. Second, a structural reform would need to be linked with a curricular reform. Otherwise, conservative teachers resisting the reform may continue to teach their pupils – as they had for many years – following exactly the same patterns. However, a structural reform could hardly pass unnoticed. It would be difficult to teach the old
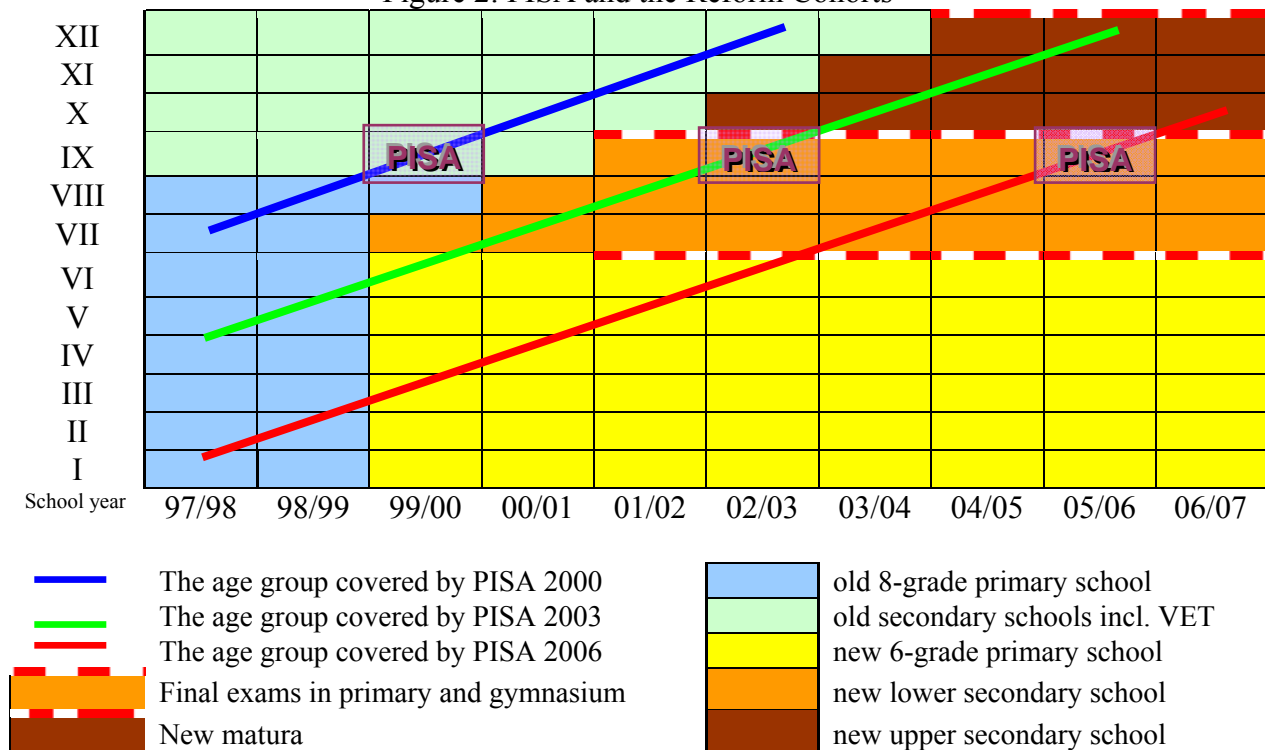
way in a new school. Thus, the idea was to provide an impulse to deep reflection in the teacher community and to bring about actual changes in teaching contents and style.

After years of complaints of an overloaded curricula and disputes about possible ways forward, the decision was finally made to implement a concept of core curricula. The concept aimed to provide schools with an extensive scope of autonomy and responsibility; schools were to build their own curriculum within a pre-determined general framework while balancing three dimensions of education: acquiring knowledge, developing skills, and shaping attitudes. The curricular reform was designed not only to bring about change in the contents of school education and to encourage the introduction of innovative teaching methods, but above all to change the teaching philosophy and culture of schools. Instead of following instructions of the educational authorities passively, teachers were expected to come up with their own teaching styles that would be best suited to the needs of their pupils. Teachers were thus faced with entirely new tasks.

The introduction of a curricular reform based on a deep decentralization required that a system for the collection of information and the monitoring of the education system be implemented simultaneously. It was, therefore, decided that common compulsory tests assessing pupil achievements should be organized at the end of the primary cycle and at the end of the lower secondary cycle (both administered for the first time in 2002). Schooling would culminate with the *matura* examination taken upon completion of the upper secondary education. All these examinations were to be organized, set and corrected by the central examination board and regional examination boards, new institutions to be set up as part of the reform. The *matura* was administered for the first time in 2005. The results of the primary school test do not affect the students' schooling career, as the completion of the cycle does not depend on the results. The score earned on the gymnasium final exam is taken into account together with the final marks in the selection process for upper secondary schools.

It is important to note that the age cohorts covered by PISA in 2000, 2003 and 2006 have been affected by the introduction of the reform in different ways (Figure 2). The first group (2000) was not affected by the reform. The group that was 15 years of age in 2003 and was covered by the second cycle of PISA started their education in primary school in the former system but attended the gymnasium, which was part of the new structure (the flagship of the reform). They did not take the final test in the sixth grade of primary school. The test was for the first time administered in 2002, when they were already gymnasium students. Finally, the group covered by PISA 2006 has attended reformed schools. They took the primary school final test in 2003 and were prepared for the final gymnasium exams a few weeks after PISA was administered in 2006.

6

Figure 2: PISA and the Reform Cohorts

Grade (left axis): XII, XI, X, IX, VIII, VII, VI, V, IV, III, II, I

PISA

School year: 97/98  98/99  99/00  00/01  01/02  02/03  03/04  04/05  05/06  06/07

Legend:
- The age group covered by PISA 2000
- The age group covered by PISA 2003
- The age group covered by PISA 2006
- Final exams in primary and gymnasium
- New matura
- old 8-grade primary school
- old secondary schools incl. VET
- new 6-grade primary school
- new lower secondary school
- new upper secondary school

The group covered by PISA 2000 consisted of the first grade students of the pre-reform secondary schools: general lyceum (one had to pass an entrance exam to enter), secondary vocational school and basic vocational school (with very low prestige). The results of PISA 2000 in Poland showed huge variations among schools, which was not surprising at all, as the pre-reform system was based on a strong selection of primary school leavers. However, the group covered by PISA 2003 (and PISA 2006) consisted of students of the last (third) grade of compulsory gymnasium. Not surprisingly, the results showed smaller variations among schools and larger ones among students within schools.

Among the PISA 2000 participants, only students of lyceums and some secondary vocational schools had a previous experience of taking a written entrance exam. The others had no experience at all. Actually, the lyceum entrance exam was not a test. It consisted of two parts: Polish language (written essay) and mathematics (solving five slightly complicated but standard problems). The first national final tests after primary school and gymnasium were carried out in 2002. At that time, the group of PISA 2003 were in the second grade of gymnasium, so they did not take the final primary school test; however, the PISA 2006 group were then still in primary school (the fifth grade) so they took the full set of the new external exams. PISA 2000 was, for most Polish students covered by the survey, the first experience in writing a test-item exam. PISA 2003 participants, although not having written an actual test-item exam, had had some previous experience due to the preparation (mock exams introduced by their teachers) for their upcoming final gymnasium exams. PISA 2006 participants were very well acquainted with doing tests. They took the final primary school test and had three years of preparation for the gymnasium exam. Konarzewski (2004) shows that teachers considered seriously the results of the 2002 final exams (the first ones ever). In fact, one-third of a

representative sample of teachers declared to have introduced changes in their teaching to acquaint students with test requirements. This aspect was also taken into account while choosing textbooks and other supporting teaching materials. To find out if these changes could have influenced the PISA results one should compare the scope (content) of both PISA and gymnasium exams. It should be noted that 26 percent of the teachers expressed the opinion that the unsatisfactory results of their students were not caused by their poor knowledge or low skills, but by their lack of experience in taking such tests. Therefore, teachers concluded that it was important to practice taking tests. Konarzewski (2008) shows that in literally all gymnasia, a substantial amount of time is devoted to solving test type tasks and doing mock exams. Approximately 5 percent of the respondents have changed their assessment schemes making them more test-like. In his conclusion, Konarzewski (2008) writes: "The test exam, being so predictable as ours, each year less and less measures the competences of gymnasium leavers but more and more the effort and time spent by schools on training students to do the exams."

## 3. Relative Increase in Scores

Improvements in student outcomes in Poland measured by PISA have been impressive. In math, Poland improved its score from 470 points in 2000, to 490 in 2003, and to 495 in 2006 (see Table 1). Reading has seen a steady improvement over time, from 479, to 490, to 508 in the latest round. In fact, in the first assessment Poland ranked below the OECD country average in reading. In 2003, Poland reached the OECD average. Moreover, by 2006, the latest assessment, Poland scores above average, ranking 9th among all countries in the world. In science, the scores are 483, 498 and 498.

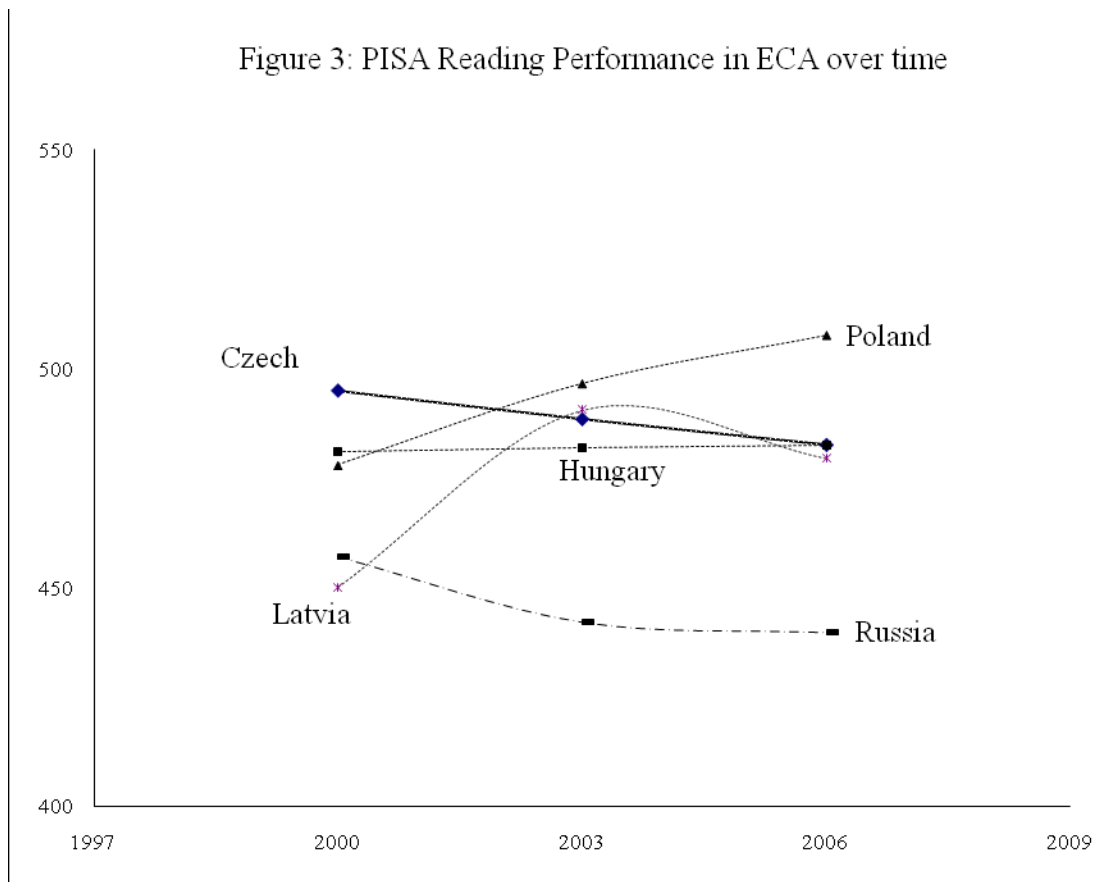Table 1: Top 10 Reading over Time, PISA

|   | 2000 | | 2003 | | 2006 | |
|---|------|---|------|---|------|---|
| 1 | Finland | 549 | Finland | 543 | Korea | 556 |
| 2 | Netherlands | 537 | Korea | 534 | Finland | 547 |
| 3 | Canada | 535 | Canada | 528 | Hong Kong | 536 |
| 4 | Hong Kong | 532 | Australia | 525 | Canada | 527 |
| 5 | Australia | 528 | Liechtenstein | 525 | New Zealand | 521 |
| 6 | Ireland | 528 | New Zealand | 522 | Ireland | 517 |
| 7 | New Zealand | 526 | Ireland | 515 | Australia | 513 |
| 8 | Japan | 525 | Sweden | 514 | Liechtenstein | 510 |
| 9 | United Kingdom | 524 | Netherlands | 513 | Poland | 508 |
| 10 | Korea | 522 | Hong Kong | 510 | Sweden | 507 |

## 4. Hypotheses for Explaining Change over Time

Several factors could explain these changes; however, it is hard to say whether there any causal relationships. In other words, in most cases validity of given explanations is doubtful. The claim of this paper is that many statistics produced from PISA are not fully comparable across years if specific empirical questions are to be answered. If one is interested in international comparisons to assess the effectiveness of countries' educational policies then only samples equivalent in a distribution of important student and family characteristics could be

compared. To give an example, if countries differ in a distribution of parental education which strongly affects student outcomes, then it is not valid to compare mean performance in these two countries to conclude whether one has more effective education policy than the other. Most likely the difference in mean performance depends more heavily on the difference in parental education than on the policy itself. Thus, any comparison on unadjusted samples could be policy irrelevant or unhelpful. Similarly, if one wants to compare achievement levels in a particular country in different years, then one needs to adjust samples to make them fully comparable. While PISA organizers try to maintain sampling schemes that are the same in all countries and years, it is difficult to preserve similar samples across time, especially when the school system changes.

*Not all transition countries improved over time*. Among participating countries in Eastern Europe, Poland is one of the best performers, with a solid improvement over time. Figure 3 shows the performance of the five Eastern European countries that participated in all three rounds of PISA. Poland is the only country with consistent improvement over time. In fact, among the five that participated in all three rounds of PISA, only Latvia and Poland improved over time. Latvia started at a lower point than did Poland and its performance over time is impressive. However, unlike Poland, in reading Latvia improved between 2000 and 2003, but slightly declined between 2003 and 2006.



Figure 3: PISA Reading Performance in ECA over time

*Reform led to improvement – through delay of vocational, more education-relevant inputs (hours, motivation, better teachers, etc.)*.  This paper concentrates on Poland.  It tries to recognize differences between the country's samples collected in 2000, 2003 and 2006.  Unadjusted score distribution is compared across years.  After this purely descriptive analysis some adjustments to make samples more comparable are proposed and semi-parametric methods are employed to produce equivalent score distributions.  Finally, the differences-in-differences method is applied to test whether extension of general comprehensive education was the main reason for score improvement over time.

*Students are more accustomed to taking tests and teachers are preparing students for tests*.  Rigorous academic testing was not the norm prior to the 1999 reforms.  Soon after the reforms tests became more important and regular.  This exposure to assessments may have prepared students, thus making them better "test takers."

## 5.      Empirical Methods and Data

We test whether the reform – more specifically the change in the structure of the school system – led to the improvement in test scores, through the delay of vocational education.  Our main approach is based on propensity score matching and reweighting.  Assume that one wants to compare survey results that are directly non-comparable because of differences in the distribution of observable characteristics.  Then one can calculate conditional expectations based on these characteristics and use them to calculate the difference of interest.  However, when the number of distinct values of important covariates is high or some of them are continuous, then any comparisons of this kind become problematic.  This is known as the "curse of dimensionality."  To resolve this problem propensity score matching methods were proposed by Rosenbaum and Rubin (1983).  They showed that instead of using multiple covariates one can use the *propensity score* that reflects the probability of being sampled to one of the groups conditional on covariates.  Originally, propensity score matching methods were applied to solve selection problems but in recent applications they were also used to adjust statistics across datasets (see Tarozzi 2007).  Similar methods were also applied earlier to compare whole outcome distributions before and after reweighting based on observable individual characteristics (DiNardo, Fort and Lemieux 1996).  In this paper, when comparing whole distributions of student achievement, we use simple propensity score weight adjustment.  The counterfactual outcome distribution is obtained using kernel density estimators with weights given by:

$$w = \frac{1 - \Pr(\text{Depvar} = 1)}{\Pr(\text{Depvar} = 1)}$$

Tarozzi (2007) argues that such reweighting produces comparable outcome distributions.  *Depvar*=1 is defined as being in a sample of interest ("target" sample), which in our case means the sample of PISA students in 2000.  *Depvar* equals 0 for students sampled in 2003 or 2006, depending on a comparison made.  Conditional probabilities are estimated using logit regression with a set of student and family characteristics defined in the same way in all waves of the PISA survey (recoded to have similar categories).  Additionally, we considered sample weights that are of importance when one wants to make inferences about population effects. PISA survey design was accounted for by multiplying propensity score weights and survey weights.

As covariates we used gender, age, mothers and father's education, the highest value of the International Socio-Economic Index among parents, number of books at home, and grade. Usually, researchers also control for immigrant status; however, there is a negligible number of migrants in the Polish sample. Missing data were imputed using the multiple imputation approach (Royston 2004). Results without any imputation were qualitatively similar, while less precise because of smaller sample sizes.

Estimates of score change for students in different tracks

Reweighting produces factual and counterfactual distributions, which are balanced in observable characteristics and can be compared across survey cycles. However, it is clear that the performance of Polish students could change for other reasons besides the introduction of comprehensive schooling. The education reform of 1999/2000 modified not only school structure but also curriculum, teacher compensation and many other things. Thus, the test score change cannot be solely attributed to the replacement of old-type tracks in secondary schools by lower secondary schools for 15-year-olds.

To deal with that, our strategy is to assess how the extension of obligatory comprehensive education by one year affected the performance of students in different tracks. More specifically, we are interested in whether students who were in old-type vocational schools in 2000 would have similar scores in 2003 or 2006 in newly established lower secondary comprehensive schools. That could be investigated by matching vocational school students from 2000 with their counterparts in 2003 and 2006. This way an estimate of performance change for students sharing characteristics common in each track can be obtained. Having that, we look at the differential impact of the reform for students who were in different tracks in 2000. The change for vocational school students minus the change for general (or mixed vocational-general) school students could be attributed mainly to the introduction of lower secondary schools. The point is that without the reform 15-year-old students in vocational schools would not have the opportunity to study in general programs; however, students in other tracks had this opportunity despite the reform. Students from general tracks can serve as a control group and the difference in a simulated score change for them and for the former vocational schools students could be attributed to the postponing of vocational education by one year.

Our approach to estimate the differential score change is similar to the differences-in-differences (DD) method. This method compares outcome change in the group of interest (treatment group) with similar change in the control group. DD estimates of treatment effect take into account trends in the whole population that equally affect both groups. In our case, we calculate the difference between the achievement of students in vocational schools in 2000 and similar students in 2003 or 2006, and we subtract it from the difference between scores of secondary general track students in 2000 and their counterparts in 2003 or 2006. Assuming that we are able to match similar students across waves of the PISA study, we can estimate how the reform affected students who without the reform would still be in vocational schools.

To define it formally, we use treatment evaluation nomenclature (see Lee 2005). The treatment is defined as being a 15-year-old student in vocational secondary school (*szkoła zawodowa*) in 2000. The control group is defined as 15-year-olds in general (*liceum ogólnokształcące*) or mixed general-vocational (*technikum*) secondary schools. We have to construct counterfactual groups of students from 2003 or 2006 samples based on their observable characteristics. A crucial assumption is that these observable characteristics constitute the main

11

factors explaining differences in student achievement across treatment groups. This assumption is called "selection on observables" in the econometric evaluation literature. Having in mind that PISA collects a rich set of background characteristics, which are strong predictors of student performance, we believe that this assumption is fulfilled and our approach is valid.

Let $Y_{it}$ be an outcome of an $i$-th individual in time $t$=0,1. We assume that some individuals were exposed to the treatment between $t$=0 and $t$=1, and write $D_{it}$=1 if an $i$-th individual was exposed to the treatment. In the rest of the paper we drop individual argument $i$ for simplicity. The differences-in-differences model is given by:

$$\alpha = \left\{ E(Y_1 \mid D_1 = 1) - E(Y_0 \mid D_1 = 1) \right\} - \left\{ E(Y_1 \mid D_1 = 0) - E(Y_0 \mid D_1 = 0) \right\}$$

A crucial assumption in this model is that a difference between transitory shocks in time $t$=0 and $t$=1 is mean independent of the treatment (see Abadie 2005; Heckman, Ichimura and Todd 1998). It means that without the treatment the average outcome for the treated would experience the same change as the average outcome for the controls (or untreated). That assumption could be demanding if groups differ in important characteristics. Thus, usually conditional differences-in-differences estimator is employed which controls for the set of covariates:

$$\alpha_X = \left\{ E(Y_1 \mid X, D_1 = 1) - E(Y_0 \mid X, D_1 = 1) \right\} - \left\{ E(Y_1 \mid X, D_1 = 0) - E(Y_0 \mid X, D_1 = 0) \right\}$$

The crucial assumption here is that quasi-experimental groups differ only by observable covariates and this condition eliminates any bias caused by this. Typically, the differences-in-differences model is estimated using simple regression analysis when any characteristic one wants to control for could be entered into the equation and interacted with time and treatment (Meyer 1995; Gruber 1994). Another approach is to balance covariates across groups to make them more comparable, which can be achieved through matching methods (Rosenbaum and Rubin 1983; Heckman, Ichimura and Todd 1998).

In our case, we need to find counterparts for the treatment and control groups in 2000 among students in lower secondary schools in 2003 or 2006. This can be achieved with matching methods where counterfactual $t$=1 scores are constructed using scores of students with similar characteristics to those observed in $t$=0. Usually, matching methods are used to make control and treatment groups more comparable assuming that we have the same observations in each group in $t$=0 and $t$=1. In our case, we do not want to adjust for dissimilarities among treatment and control groups. We know that students who were in vocational schools differed from those in general schools but we are interested whether moving students from different tracks, who differ by assumption, into the one-type comprehensive lower secondary schools affected them similarly. Matching in our case is used to adjust in time by drawing comparable groups from 2003 or 2006 samples, not for adjustments across quasi-experimental groups.

As already mentioned, when dimension of $X$ is high, then exact matching on covariates is not possible (the "curse of dimensionality"). In this case, individuals can be matched on one-dimensional propensity score $P = P(D=1|X)$, where $D$ again indicates treatment and $P$ reflects the conditional probability of being treated (see Rosenbaum and Rubin 1983). However, as we noted above we have to balance covariates not between treatment and control groups, which differ by assumption, but between waves of the survey. Only in 2000 were students *treated*, which means they were separated into different types of secondary schools. After the reform, in PISA 2003 and PISA 2006, all students were in lower secondary comprehensive schools. Nevertheless, one can draw from 2003 and 2006 samples to find good matches and construct

reference groups for students tested in 2000. We match using propensity score $P^{2000} = P(T=2000|X)$, reflecting the propensity to be in the PISA 2000 sample. Two propensity scores need to be estimated. One measuring a propensity of being in a vocational school in 2000 for students tested in 2003 or 2006, and the other for being in a general (or mixed vocational-general) school in 2000 for students tested in 2003 or 2006. Thus, we have the propensity score for treated units (vocational school students) $P_T^{2000}$ and the propensity score for controls $P_C^{2000}$ (students in other tracks), both reflecting the propensity of being sampled in 2000 for students sampled in 2003 or 2006.

Let us now define $Y^1$ as the score of students separated into tracks in secondary schools in 2000 and $Y^0$ as the score for students tested in 2003 or 2006. Now, the DD estimator could be defined by:

$$\alpha_{DD} = \left\{ E(Y^1 \mid D=1) - E(Y^0 \mid P_T^{2000}, D=1) \right\} - \left\{ E(Y^1 \mid D=0) - E(Y^0 \mid P_C^{2000}, D=0) \right\}$$

In this equation $E(Y^1 \mid D=1)$ and $E(Y^1 \mid D=0)$ are directly observed in the data, but $E(Y^0 \mid P_T^{2000}, D=1)$ and $E(Y^0 \mid P_C^{2000}, D=0)$ have to be constructed from 2003 or 2006 PISA samples using propensity scores. We first estimate the performance change for students in each type of secondary school in 2000 and their matched counterparts in 2003 or 2006. Then we compare these performance changes across students from different tracks. The difference between performance gains of students in the former vocational track and for students in other tracks is the differences-in-differences estimator of the impact of abolishing the vocational curriculum for 15-year-olds. This estimator reflects the causal impact of the reform under the crucial assumption that the score change for students in the general track would be the same without the reform. This assumption is not directly testable, however. For general track students the curriculum did not change in a fundamental way, while other changes affected them as much as they did other students.

Propensity scores were estimated using logit regressions. Two kinds of propensity score matching were then employed: nearest neighbor matching 1-to-1 matching and kernel matching. The first method matches to each treated observation one control observation with the closest value of the propensity score. The kernel method constructs values for matched counterparts by weighting control observations by their proximity in the propensity score to the treated observation, using a kernel function (we used Epanechnikov kernel with bandwidth 0.6; see Becker and Ichino 2002 for details of the Stata procedure used). In both methods a common support restriction was imposed, which means that if propensity score distribution does not overlap at the bottom or top of the distribution, then observations with extreme propensity score values will not be considered. This restriction rarely affects the results in our case, but guarantees that proper matches were drawn from the 2003 and 2006 samples.

Finally, we need to decide which covariates to balance across surveys or use to draw counterparts of 2000 students in different tracks from 2003 and 2006 data. An obvious limitation here is availability of control variables which are identically defined across waves of PISA. Fortunately, PISA collects crucial variables reflecting student socio-economic background, including the HISEI index (highest of mother or father international socio-economic index), mother and father ISCED education level, and number of books at home; in addition, student gender, age, grade they attended at the time of PISA survey, and family structure, were also used as covariates. Some of these indicators, mainly HISEI index, parental education levels, and

13

family structure, have a small number of missing observations. To keep the sample size and performance distribution untouched by the matching exercise, missing values for matching covariates were imputed through multiple imputation models (Royston 2004).

We use the Poland's data from the OECD's Program for International Student Assessment (PISA), an internationally comparable standardized student test conducted every three years to test reading, mathematics and science achievement of 15-year-olds. PISA assesses and compares student achievement based on a standardized and highly reliable framework. It assesses the knowledge and skills essential for full participation in society and the world of work. It is conducted near the end of compulsory education. In all cycles, the domains of reading, mathematical and scientific literacy are covered not merely in terms of mastery of the school curriculum, but in terms of important knowledge and skills needed in adult life. This makes the PISA framework attractive as a basis for international comparisons.

The detailed description of the assessment framework, cautious procedures of translation and supervision of country specific implementation, and careful calibration of student scores based on response item theory and collected background variables all support the view of PISA as a valid framework for international comparisons of student achievement. Depending on population and country involvement there are approximately 2,500 to nearly 30,000 students tested in each participating country. The survey is realized as a representative to the population of interest with a two-stage stratified sample with random sampling of schools and within each school. Survey weights are provided by the organizers that reflect the different probabilities of schools and students to be sampled (OECD 2001).

The PISA survey has a complex structure, similar to methods commonly used in other educational surveys such as the International Association for the Evaluation of Educational Achievement's (IEA) Trends in International Mathematics and Science Study (TIMSS) and Progress in International Reading Literacy Study (PIRLS), or the United States' National Assessment of Educational Progress (NAEP), with sampling conducted with different probabilities in two stages within separate strata. This complexity should be taken into account by using probability weights when calculating point estimates and by adjusting for clustering and strata design when estimating standard errors. However, there is little advice in the literature on how to account for survey design in matching methods (see Zanutto 2006, for example, of analysis with survey weights and stratification matching). In our case, we used survey weights when calculating average outcomes for the treated students in PISA 2000. This way the results are representative for the population of 15-year-olds in 2000. Also, students are answering a randomly assigned groups of test items (so-called booklets), but responses are put into one common scale using psychometric models. The performance of each student is reflected by five plausible values which give equally probable performance scores for individuals. Plausible values should not be used to judge individual performance but provide unbiased estimates of achievement for whole populations of interest. Each analysis should be repeated five times with each plausible value used once to take into consideration measurement error in student performance. We follow this strategy in this paper by repeating each analysis five times, once with each plausible value of reading performance. When using the multiple imputation method, we imputed missing values once for each plausible value and then repeat any estimation five times, once with each dataset containing one plausible value and imputations obtained with this

plausible value. That should guarantee that all imputation errors, one in plausible values and the others in imputed covariates, will be taken into account (see OECD 2002, 2005). The final set of variables from the PISA dataset used in this analysis are re-sampling replicate weights used in the calculation of standard errors. Intra-cluster correlation violates an assumption needed for the unbiasedness of the analytical method of calculating standard errors based on the variation of the sample. Re-sampling methods such as bootstrapping, Jackknifed Repeated Replication and Balanced Repeated Replication serve as alternative means of calculating standard errors. These methods calculate sampling variance by re-sampling the same sample to mimic re-sampling of the original population. Replicate weights are alternative sample weights which represent a sub-sample based on the original sampling design. PISA provides replicate weights compatible with Fay's adjusted Balanced Repeated Replication. These weights were constructed to reflect the sampling design including any country specific modifications, as well as non-response by students or schools (OECD 2002: 89-98). Standard errors were obtained by the BRR method. In our case, the additional benefit of using BRR weights is that these were produced by survey organizers who used confidential information not available for external users.

Decompose Change over Time

A simple decomposition analysis is undertaken in order to attempt to explain one of the possible pathways be which the reform may have led to improved student achievement. We decompose reading scores over time PISA 2000 and 2006 to explain to what extent the increase in scores is due to changes in characteristics and what proportion is due to changes in returns to characteristics. A simple education production function is estimated (Hanushek 1986, 2002; Todd and Wolpin 2003; Glewwe 2002), which relates various inputs affecting student learning to measured outputs, the PISA standardized reading test score. Past empirical research does not always agree on which school and family inputs improve children's achievement. Examples are the disagreements found on the role of schooling inputs such as class-size, teacher experience, teacher education and mother's employment. Nevertheless, although a child's achievement is inherently individual in nature, a large body of evidence points to the existence of persistence effects in educational achievement across generations (Fertig 2003; Fertig and Schmidt 2002; Currie and Thomas 1999). Consequently, one must control for individual pupil characteristics as well as family background. Finally, one needs to control for characteristics on school environment as well as institutional arrangements. Evidence also suggests that socioeconomic and family background variables, such as parent's education and the number of books a child has, are very important determinants of test scores at early ages (Fryer and Levitt 2002). Therefore, we specify and estimate education production functions that relate students' achievement to individual, family and school inputs. We then proceed to decompose the over time test score gap into an explained component (accounting for student, family, and school characteristics) and an "unexplained" – or returns, or the efficiency by which the country is able to convert characteristics into student learning outcomes as measured by test scores – component, using the traditional Oaxaca (1973)-Blinder (1973) decomposition method.

The model specification for the estimation of the production function for cognitive achievement is as follows:

$$T_{ija} = T_a(A_{ija}, F_{ija}, S_{ija}) + \epsilon_{ija}$$

where $T_{ij}$ is the observed test score (from PISA reading) of student $i$ in household $j$ at time $a$ (time of the test), $A_{ija}$ is a vector of individual, student, characteristics, $F_{ija}$ is a vector of parent inputs, $S_{ija}$ is a vector of school-related inputs, and $\epsilon_{ija}$ is an additive error, which includes all the omitted variables including those which relate to the history of past inputs, endowed mental capacity and measurement error. The linear specification (after dropping subscript $a$) of the production function is given by:

$$T_{ij} = \beta_0 + \beta_1 A_{ij} + \beta_2 F_{ij} + \beta_3 S_{ij} + \epsilon_{ij}$$

where $\beta_0$ to $\beta_4$ are coefficients to be estimated. The standard procedure for analyzing the determinants of the test score differences over time is to fit equations between test scores and observed characteristics. The observed test score differential can be decomposed as:

$$T_{2006} - T_{2000} = (X_{2006} - X_{2000})\beta_{2006} + X_{2006}(\beta_{2006} - \beta_{2000})$$

where $T$ is the standardized test score, $X_i$ is a vector of student, family and school characteristics for the $i$th individual, $\beta$ is a vector of coefficients and 2006, 2000 subscripts are identifiers of the PISA test score in reading in years 2000 and 2006; evaluated at 2006 prices.

The overall test-score increase can, therefore, be decomposed into two components: one is the portion attributed to differences in characteristics ($X_{2006} - X_{2000}$) evaluated with the 2006 prices, or 2006 group performance ($\beta_{2006}$); the other portion is attributable to differences in effects on performance ($\beta_{2006} - \beta_{2000}$) of 2000 and 2006 students derived from the same characteristics. This second (unexplained) component, while more difficult to interpret in the present context compared to an earnings gap decomposition framework, can be assigned more than one interpretation. An obvious one is that the unexplained portion of the test score increase may reflect certain unobserved family characteristics that are correlated with achievement over time, possibly relating to household wealth. In addition, it may be that the different cohorts of students do not reap the same benefits from equivalent school and classroom resources. Finally, the differences in the returns may reflect the impact of changes over time based on past reforms that both increased school enrollments in Poland and contributed towards improving the quality of school inputs. Certain of the above coefficient estimates may be subject to biases. For example, if a school characteristic is correlated with unobserved family characteristics that influence achievement (such as family wealth and parents' motivation), the effect of attending a school with such characteristics may be biased. While test scores and individual and family information are at the individual level, school resources and other school-related inputs are at the school level. In choosing the estimation method, we recognize that observed test scores are expected to be correlated at the school level due to clustering effects. Therefore, the assumption that disturbances are independently and identically distributed with fixed conditional variance does not hold. The estimation method of OLS by cluster at the school level is used.

## 6.    Results

In what follows, we focus on reading literacy as only performance in this domain is fully comparable across PISA cycles. Performance in mathematics can be compared across 2003 and 2006 only because the 2000 assessment framework was later modified. Science performance in 2006 cannot be related to previous cycles as the framework was completely changed in 2006. The results are presented for the whole sample and for the modal grade only, which is the 9[th] grade in Poland. In PISA 2000, only the 9[th] grade was sampled while in PISA 2003 and 2006 students from 7[th], 8[th] and 10[th] grade were also sampled. The results suggest that students in non-

modal grades slightly affect the estimates and should be taken into account.  In the regression and matching analysis we simply adjust for student grade to take into account these differences.

Reweighting clearly lowers the mean scores of students in 2003 and 2006 (Table 2).  On the other hand, scores for students in the modal grade are slightly higher.  These effects, which influence results in opposite ways, when combined are positive, suggesting that overall student performance increased between 2000 and 2003 or 2006.  For example, the change in factual scores (weighted only with survey weights) from 2000 to 2003 is 17.5, and from 2000 to 2006 it is 28.5; but it diminishes after reweighting to 6.1 and 23.7, respectively.  However, after reweighting and taking students from the modal grade only, the gains are equal to 13.5 and 30.6, respectively.  Thus, there is no doubt that increases in mean scores occurred in Poland from 2000 to 2003.  The change between 2003 and 2006 is less clear.  After reweighting, the initial difference of 11.0 (or 11.6 in modal grade) almost disappears.  Nevertheless, we clearly observe substantial overall improvement after 2000.

Table 2: PISA 2000, 2003 and 2006 results for Poland in Reading
Factual (with survey weights), reweighted to the reference year
(with survey and propensity score weights), and modal for modal grade

| | Factual | Factual Modal grade | Factual | Reweighted | Factual Modal grade | Reweighted Modal grade |
|---|---|---|---|---|---|---|
| *Reweighting to 2000* | *2000* | | *2003* | | | |
| Mean score | 479.1 | 479.1 | 496.6 | 485.2 | 501.9 | 492.6 |
| Change from 2000 | - | - | 17.5 | 6.1 | 22.8 | 13.5 |
| *Reweighting to 2000* | *2000* | | *2006* | | | |
| Mean score | 479.1 | 479.1 | 507.6 | 502.8 | 513.5 | 509.7 |
| Change from 2000 | - | - | 28.5 | 23.7 | 34.4 | 30.6 |
| *Reweighting to 2003* | *2003* | | *2006* | | | |
| Mean score | 496.6 | 501.9 | 507.6 | 499.5 | 513.5 | 506.9 |
| Change from 2003 | - | - | 11.0 | 2.9 | 11.6 | 5.0 |

The change in mean scores is obviously interesting but looking at the change in whole distributions gives a more detailed picture.  Figures 4 and 5 show estimated factual distributions of scores in 2000, 2003 and 2006, together with reweighted scores for 2003 or 2006.  The figures clearly show that the whole score distributions are "shifted" to the right in 2003 or 2006 compared to 2000.  This means that the difference in achievement across PISA cycles is not only for low-achievers but also for high-achievers.  This way Poland closes the gap at all levels of performance.  While in PISA 2000 the percentage of students in the top two reading proficiency levels (4[th] and 5[th]) was 24.5 percent (compared to OECD average of 31.8 percent), this number in Poland increased in 2006 to 34.7 percent (compared to OECD average of 29.3 Percent).  At the bottom of the distribution, the percentage of Polish students at the 1[st] or below proficiency level was 23.3 percent in 2000 (17.9 percent in OECD) and 16.2 percent in 2006 (20.1 percent in

OECD) (OECD 2003: Table 2.1a; OECD 2007: Table 6.1a). That poses interesting questions about what caused the "shift" of the student score distribution. While extension of compulsory comprehensive education can explain higher performance for low-achievers who were mostly in vocational tracks, it is more complicated to explain higher performance of top-achievers. The question is whether the introduction of lower secondary schools could have an impact on students in former general secondary schools and what was in the reform that increased scores so significantly.

Figure 4: Change in reading literacy distribution between PISA 2000 and 2006
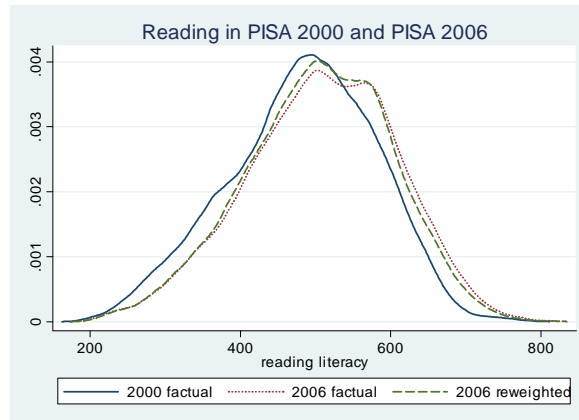


Figure 5: Change in reading literacy distribution between PISA 2003 and 2006



Estimates of score change for students in different tracks

Results for differences-in-differences propensity score matching estimates of the effect of abolishing tracking system for 15-year-olds are presented in Tables 3, 4 and 5. Table 3 contains estimates of factual and counterfactual mean scores for all students in PISA 2000, 2003 and 2006. In addition, results for students in vocational and non-vocational tracks are presented. Factual scores were weighted by survey weights provided in the official PISA datasets. Counterfactual scores were constructed using matching methods with survey weights taken into account in the manner described above.

Table 3: Factual and counterfactual scores of students in different upper secondary tracks

| Reading achievement | PISA 2000 factual weighted mean score (no of obs) | PISA 2003 factual weighted mean score (no of obs) | PISA 2003 matched counterfactual score (no of matched obs) | | PISA 2006 factual weighted mean score (no of obs) | PISA 2006 matched counterfactual score (no of matched obs) | |
|---|---|---|---|---|---|---|---|
| | | | *Kernel matching* | *1-1 matching* | | *Kernel matching* | *1-1 matching* |
| | *(1)* | *(2)* | *(3)* | *(4)* | *(5)* | *(6)* | *(7)* |
| All schools | 479.1 (3654) | 496.6 (4196) | 497.9 (4151) | 495.2 (2528) | 507.6 (5233) | 514.9 (5229) | 514.1 (3056) |
| ISCED 3C schools | 357.6 (983) | - | 466.7 (4010) | 460.5 (926) | - | 484.3 (5141) | 474.4 (1090) |
| ISCED 3B schools | 478.4 (1491) | - | 491.4 (4150) | 487.7 (1527) | - | 507.3 (5163) | 501.8 (1823) |
| ISCED 3A schools | 543.4 (1180) | - | 525.6 (4064) | 524.9 (1233) | - | 543.0 (5221) | 547.0 (1376) |
| ISCED 3A and 3B schools | 513.6 (2671) | - | 507.3 (4157) | 507.0 (2206) | - | 524.8 (5233) | 520.5 (2609) |

Note: Standard errors are given in parentheses and were obtained from bootstrapping (kernel matching) or analytically (1-1 matching). * $p<0.05$, ** $p<0.01$, *** $p<0.001$

Not surprisingly, the counterfactual mean scores for all schools are similar to those reported earlier (see Table 2). Moreover, results for kernel matching and 1-to-1 matching are also similar. They differ slightly because of different matching methods and various number of matched control observations (provided in parentheses) but give qualitatively similar conclusions. This shows that the choice between reweighting or different matching methods has no crucial impact on final estimates.

Results are summarized in Table 4, where the estimates of score improvement are presented.[2] These estimates assess trends in performance for all students and across groups of students who without the reform would be in different secondary tracks. We see again overall improvement of average performance among 15-year-olds in Poland. Score improvement for all students is remarkable, at 16 to 18 points from 2000 to 2003 and around 35 to 36 points from 2000 to 2006. Crucial estimates are for the hypothetical performance improvement of 2000 in different tracks. Performance improvement for potential students of former vocational schools is simulated to be higher than 100 points from 2000 to 2003 and 120 points from 2000 to 2006. This is more than one standard deviation of PISA scores in OECD countries, which is a dramatic improvement, hardly comparable to effects of any known educational policy. Obviously, these estimates are statistically significant, supporting the hypothesis that 15-year-old students who without the reform would be placed in vocational tracks heavily benefited from the reform. However, the benefits for students in other tracks are not that visible. Students in mixed-general

---

[2] The numbers presented in the third row, after the name of the comparison and matching method, are showing how these differences were calculated from the results presented in Table 3. In each case, the difference was calculated by taking a counterfactual performance score of matched student from the 2003 or 2006 samples and subtracting from it the factual score of students tested in 2000. Standard errors for these differences were calculated by employing the BRR method, which properly accounts for complex survey design (stratification, clustering, and response adjustments).

schools improved their scores only slightly in 2003 and noticeably in 2006 only. Students in the general track would potentially have lower scores in 2003 and similar performance in 2006.

These findings are in line with economic intuition. The short-term effects of the reform could be harmful for general school students who were mixed with low-achievers in newly introduced lower secondary schools. In the longer term, this negative impact disappears. It could be that teachers adjusted their methods to more diversified classrooms or that segregation between and within lower secondary schools recreated the former stratification. Students in mixed-general schools obviously benefited from the reform when one considers general skills tested in PISA. Effects are again more visible in the longer term probably because of similar adjustments and mixing with high-achieving students. Positive effects for vocational school students were expected because after the reform they spend much more time learning non-vocational subjects. However, what is interesting is the enormous magnitude of the effect that is about one standard deviation of PISA international scores. Again, it is not surprising because of the change in classroom time allocated after the reform to the general instead of vocational subjects. However, what is surprising is that counterparts of vocational school students "adapted" so quickly to the new system. In other words, it is striking that just a few additional months of comprehensive instead of vocational education changes general skills for an important number of students so dramatically.

Table 4: Propensity score matching estimates of score change for students in different upper secondary school tracks

| Reading achievement | Score change: PISA 2003 – PISA 2000 | | Score change: PISA 2006 – PISA 2000 | |
|---|---|---|---|---|
| | *Kernel matching* (1) - (3) | *1-to-1matching* (1) - (4) | *Kernel matching* (1) - (6) | *1-to-1 matching* (1) - (7) |
| All schools | 18.8 (4.3) | 16.1 (4.5) | 35.8 (4.4) | 35.0 (4.5) |
| ISCED 3C schools | 109.2 (5.8) | 103.0 (5.8) | 126.8 (5.7) | 116.9 (6.3) |
| ISCED 3B schools | 13.0 (5.7) | 9.3 (6.5) | 28.9 (5.8) | 23.4 (7.2) |
| ISCED 3A schools | -17.8 (5.4) | -18.5 (4.3) | -0.4 (5.1) | 3.6 (5.0) |
| ISCED 3A and 3B schools | -6.3 (4.3) | -6.6 (4.3) | 11.2 (4.2) | 6.9 (4.4) |

Notes: Propensity score matching with common support restriction; Standard errors are given in parentheses and were obtained through BRR method accounting for complex survey design.

Finally, we turn to relative improvement or differences-in-differences estimator of performance change for vocational school students. Relevant estimates are presented in Table 5. They are based on simple calculations from the tables above but clearly show the improvement

of vocational school students versus score change for students in other tracks. The first row shows estimates of the relative performance change of vocational school students versus all students in other tracks. This is the most reliable comparison because of the highest possible sample size. As we noted above, the estimates show that the relative improvement in performance of vocational school students is higher than one standard deviation of international scores (100). Relative improvement in comparison to students in mixed general-vocational schools is slightly lower but still substantial.

| Table 5: Relative score change (differences-in-differences) for students in vocational schools | | | | |
|---|---|---|---|---|
| Relative score change | from PISA 2000 to PISA 2003 | | from PISA 2000 to PISA 2006 | |
| | *Kernel matching* | *1-1 matching* | *Kernel matching* | *1-1 matching* |
| ISCED 3C versus ISCED 3A+3B | 115.5 | 109.6 | 115.7 | 110.0 |
| ISCED 3C versus ISCED 3A | 127.1 | 121.5 | 127.2 | 113.3 |
| ISCED 3C versus ISCED 3B | 96.2 | 93.7 | 98.0 | 93.5 |

Summing up, there is no doubt that students who were in vocational tracks in 2000 would score much lower without the reform. The results show that the reform improved the overall mean performance of 15-year-olds in Poland, mainly by boosting the performance of students in former vocational and mixed general-vocational tracks. There are two remaining question which are of policy relevance. One is whether this large positive impact of the reform was long lasting. More precisely, the question is whether 15-year-old students in lower secondary schools still have higher achievement one or two years later, after they were again separated into tracks at the upper secondary school level. Secondly, it seems interesting to find out what particular changes in curriculum or organization of the school system boosted student scores. These two issues are investigated below by utilizing data from the PISA 2006 national option in Poland which provides performance scores for 16- and 17-year-olds and by employing decomposition analysis.

**7. Additional Analyses**

PISA gives an option to participating countries to conduct additional research using the framework and measurement tools of PISA. Poland utilized this option in 2006 for the first time to conduct additional surveys among 16- and 17-year-old students (see Federowicz 2007 for the report on PISA 2006 in Poland). These students were mostly in upper secondary school, which opens up a possibility for comparisons with students in lower secondary schools (from the international sample). After taking into account the difference in student age, one can compare the performance of 15-, 16- and 17-year-olds, also across educational tracks of upper secondary schools. In other words, knowing the achievement at the end of lower secondary schools we can compare how much students updated their skills during the time of education in different types of upper secondary schools.
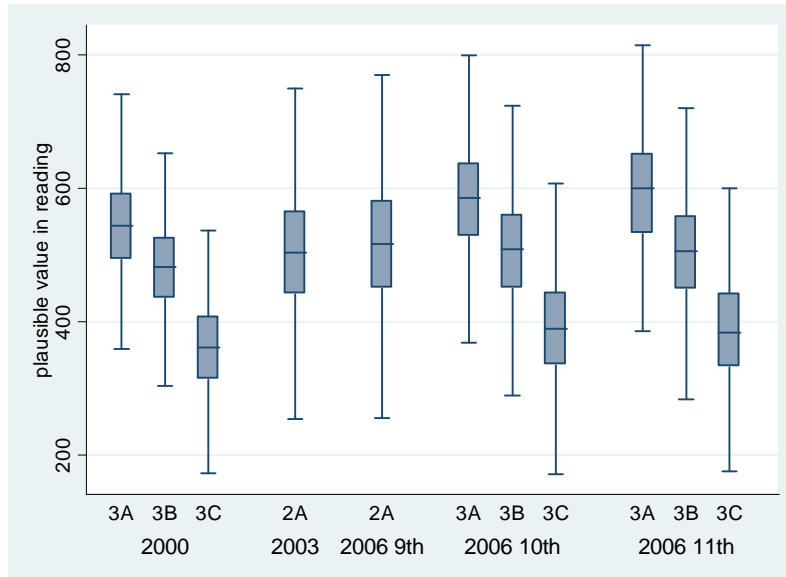
Analysis of PISA 2006 "national option" samples

Estimates of mean achievement by PISA cycle, grade and type of school program are presented in Table 6. First, 16-year-old students in the 10th grade score on average higher than do 15-year-olds in the 9th grade, and 17-year-old in the 11th grade score higher than 16-year-olds. This is in line with intuition that older students are more able to solve PISA tests. However, when we look at the type of school programs, it is clear that mainly students in ISCED 3A schools improved, whereas 17-year-old students in vocational schools have even lower scores. This seems to be counterintuitive but there are two possible, highly likely explanations. Firstly it has to be noted that students change tracks, mostly in the 10th grade, and these are mostly low performing students who are forced to move to the vocational or mixed general-vocational track. Because of such changes in the population, student achievement in mixed general-vocational or vocational upper secondary schools could be lower in higher grades. Second, students in ISCED 3C tracks devote more time for vocational training in higher grades. Therefore, their general skills tested in PISA could be diminished. Consequently, slightly lower achievement in ISCED 3C is not that surprising.

| Table 6: Mean achievement by PISA wave, grade and type of school program | | | | | |
|---|---|---|---|---|---|
| PISA wave: | 2000 | 2003 | 2006 | | |
| Type of school program: | 9th grade | 9th grade | international 9th grade | national 10th grade | national 11th grade |
| Mean achievement | 479.1 | 501.9 | 513.5 | 520.1 | 528.3 |
| ISCED 2A *lower secondary school* | - | 501.9 | 513.5 | - | - |
| ISCED 3A general secondary | 543.4 | - | - | 580.8 | 592.6 |
| ISCED 3A/B general, profiled secondary | - | - | - | 494.9 | 494.6 |
| ISCED 3B vocational secondary | 478.4 | - | - | 505.9 | 508.8 |
| ISCED 3C vocational (basic) | 357.6 | - | - | 388.8 | 384.1 |

Box plots presented below summarizes score distribution for the categories presented in Table 6 (Figure 6). This time data for vocational upper-secondary schools and general profiled (mixed) upper-secondary schools were collapsed into one category, ISCED 3B. A slight improvement is visible from 2000 to 2006 and for the 10th and 11th grades. However, it is also evident that mean scores increased because of the improvement at the top of the achievement distribution. Looking at the vocational ISCED 3C schools it is clear that while some students caught up with their colleagues in other tracks, students performing at the lowest proficiency levels are still numerous.

Figure 6: PISA scores compared over time and with 16 and 17 year-olds



An interesting comparison is between PISA 2000 results and the PISA 2006 additional "national option" sample. Table 7 gives estimates of the relative difference between achievement of students in vocational and other tracks in 2000 and in 2006 separately for the 10th and 11th grades. The results are striking. While the overall mean performance of Polish students improved significantly, the difference between students in vocational and other tracks remained almost the same, and even increased for 17-year-olds. Thus, the stratification of Polish students in the old secondary school system still exists under the new name of upper secondary schools. It seems that the reform helped to update the skills of the average student, but the negative effect of the tracking system was simply postponed by one year. The achievement gap noted in PISA 2000 is still visible and almost of the same magnitude. From one point of view this is not surprising, since the reform focused on primary and lower secondary education. On the other hand, it is now evident that the overall effect of the reform is not so positive. Intuitive claims that upper secondary education did not improve that much seems to be supported by the results presented here. Clearly, while there are visible positive effects of the reform, there are also doubts whether these positive effects are long lasting or affect all students similarly. Still, students in vocational tracks lack knowledge and skills needed to fully benefit from the modern society and economy and the reform did not change that.

Table 7: Estimates of Relative Differences in Achievement in Vocational and Other Tracks in 2000 and 2006, and for the 10th and 11th grade special sample

| | 2000 9th grade | 2006 10th grade | 2006 11th grade |
|---|---|---|---|
| ISCED 3A + 3B | 513.6 | 544.4 | 552.7 |
| ISCED 3C | 357.5 | 388.8 | 384.1 |
| Difference | 156.0 | 155.6 | 168.6 |
| (standard error) | (7.5) | (10.2) | (10.3) |

Decomposition results

      We present the results of a simple decomposition in an attempt to explain one of the possible pathways by which the reform may have led to improved student achievement.  Table 8a and 8b presents the results of production function estimates along with the decomposition results in reading.  Overall, two-thirds of the observed test score differential between PISA 2000 and 2006 is explained by the variables in the model.  Thus, only one-third is unexplained, or in this case, due to the returns to characteristics.  In terms of unexplained, most of the difference is due to returns to student characteristics, and of that, all is due to age.  That is, the returns to being older increased immensely over time.  In terms of the unexplained, most is due to school characteristics.  Moreover, of that, almost all is due to hours of instruction.  That is, receiving more than four hours per week of reading classes is associated with a higher score.  The returns to hours of reading class increased over time, but much more of an increase was seen in the proportion of students that received more than four hours of language class, from 1 percent in 2000 to 76 percent in 2006.

Table 8a: PISA Reading Scores Decomposition for Poland, PISA 2000-2006

| Test Scores | b2000 | b2006 | X2000 | X2006 | Determinants of Test scores Differentials | | as % of total test score diff | |
|---|---|---|---|---|---|---|---|---|
| | | | | | *Endowments* | *Unexplained* | *Endowments* | *Unexplained* |
| | | | | | $b_{2006}(X_{2006}-X_{2000})$ | $X_{2006}(b_{2006}-b_{2000})$ | | |
| Constant | 296.47 | 161.49 | 1.00 | 1.00 | 0.00 | -134.98 | 0.0 | -205.2 |
| | | | | | | | | |
| *Schools* | | | | | | | | |
| Student - teacher ratio | 2.08 | -0.14 | 12.01 | 11.33 | 0.09 | -26.61 | 0.1 | -40.5 |
| % of certified teachers | -23.92 | 18.85 | 0.90 | 0.97 | 1.21 | 38.57 | 1.8 | 58.6 |
| achievement data used to evaluated teachers and principal performance | 51.94 | 7.00 | 0.98 | 0.92 | -0.43 | -44.01 | -0.7 | -66.9 |
| More than 4 hours per week of language class | 3.27 | 42.77 | 0.01 | 0.76 | 32.09 | 0.41 | 48.8 | 0.6 |
| attend to public school | 13.89 | -22.18 | 0.98 | 0.98 | -0.13 | -35.25 | -0.2 | -53.6 |
| *Student characteristics* | | | | | | | | |
| Age | 0.28 | 12.85 | 15.73 | 15.71 | -0.23 | 197.76 | -0.4 | 300.6 |
| Female | 36.12 | 32.53 | 0.51 | 0.51 | -0.05 | -1.83 | -0.1 | -2.8 |
| *Family background* | | | | | | | | |
| Mother - Upper secondary | 4.68 | 27.11 | 0.74 | 0.77 | 0.70 | 16.65 | 1.1 | 25.3 |
| Mother -University | 41.49 | 63.09 | 0.17 | 0.15 | -1.52 | 3.65 | -2.3 | 5.6 |
| 11–100 books | 31.38 | 30.58 | 0.39 | 0.54 | 4.75 | -0.31 | 7.2 | -0.5 |
| 101-500 books | 52.90 | 67.39 | 0.47 | 0.35 | -8.03 | 6.87 | -12.2 | 10.4 |
| Computer at home | 22.74 | 33.89 | 0.47 | 0.80 | 11.22 | 5.19 | 17.1 | 7.9 |
| | | | | | | | | |
| Total | | | | | 39.7 | 26.1 | 60.3 | 39.7 |
| Overall | | | | | 65.8 | | 100.0 | |

Source: Program for International Student Assessment ( PISA) 2000 and 2006

Table 8b: Determinants of PISA Differentials, Reading 2000-2006

| | as % of total test score diff | |
|---|---|---|
| | *Endowments* | *Unexplained* |
| Constant | 0.0 | -205.2 |
| Schools | 49.9 | -101.7 |
| Family | 10.8 | 48.7 |
| Student | -0.4 | 297.9 |
| Total | 60.3 | 39.7 |
| Overall | 100 | |

Source: Program for International Student Assessment ( PISA) 2000 and 2006

We also present a modified decomposition (Table 9).  The results are very similar.  That is, most of the differential is explained, and most is due to school characteristics.  Most of the differences in school characteristics are the increased hours of math and language class instructions that students were subject to because of the reform.

Table 9: Modified Decomposition Results

|  | Explained (%) | Unexplained (%) |
|---|---|---|
| PISA  Math 2000-2006 |  |  |
| Overall | 68.7 | 31.3 |
| Schools | 71.0 |  |
| Family | 28.9 |  |
| Student | 0.2 |  |
| PISA  Reading 2000-2006 |  |  |
| Overall | 66.1 | 33.9 |
| Schools | 83.6 |  |
| Family | 16.6 |  |
| Student | 0.2 |  |

## 8.    Conclusions

The vocationalization of the secondary school curricula has been advocated for many decades.  The call for technical and vocational schooling used to be a standard recommendation promoted by international organizations and followed by several countries.  Unfortunately, the fervor for this approach has run ahead of substantial evidence on the impact of this policy.  This paper contributes towards modestly filling that void by studying the effects of a change in curricular structure on educational quality.

The Polish education reform program generated an exogenous variation in vocational school attendance at the secondary school level across time that provided us the opportunity to assess the impact on test scores.  Our identification strategy used the fact that likely vocational graduates did not have that option in PISA 2003, thus providing a comparison group for our empirical approach, propensity score matching and differences-in-differences estimation.

Our results suggest that vocationalization, on average, reduces test scores by a full standard deviation.  Clearly, given the importance of the reform program, other factors played a role in the increase in Poland's scores in PISA.  Nevertheless, the delayed vocationalization played a major role.  The pathway, we argue, is through increased hours of math instruction, possibly more exposure to testing, and increased motivation on the part of students and teachers.

We substantiated our findings by taking advantage of the application of PISA to 16 and 17 year olds in PISA.  We find that once vocational school options are available again, by the time students are 16, then test scores decline for those students who enter the vocational track.  While this goes a long way towards proving our initial findings, it also serves as a caution to policymakers about the effectiveness of vocational schooling – when that schooling is not designed to improve math and reading skills, which we show such students can learn when they are given the opportunity, and which have become the real vocational skills in the world of work today.  To increase test scores, it helps to have students study the subjects of the tests.  We conclude that much of the test score increase in Poland in recent years has to do with the delayed vocationalization of the secondary school curriculum.

**References**

Abadie, A. 2005. "Semiparametric Difference-in-Differences Estimators." *Review of Economic Studies* 72(1): 1-19.

Attanasio, O., A. Kugler and C. Meghir. 2009. "Subsidizing Vocational Training for Disadvantaged Youth in Developing Countries: Evidence from a Randomized Trial." IZA Discussion Papers 4251, Institute for the Study of Labor (IZA).

Becker, G. 1964. *Human Capital*. Chicago: The University of Chicago Press.

Becker, S. and A. Ichino. 2002. "Estimation of average treatment effects based on propensity scores." *Stata Journal* (StataCorp LP) 2(4): 358-377.

Betcherman, G., K. Olivas and A. Dar. 2004. "Impacts of active labor market programs: new evidence from evaluations with particular attention to developing and transition countries." World Bank Social Protection Discussion Paper 0402.

Blinder, A. 1973. "Wage discrimination: Reduced form and structural estimates." *Journal of Human Resources 8*(4): 436–455.

Card, D.E., P. Ibarraran, F. Regalia, D. Rosas and Y. Soares. 2007. "The Labor Market Impacts of Youth Training in the Dominican Republic: Evidence from a Randomized Evaluation." NBER Working Paper No. W12883.

Currie, J. and D. Thomas. 1999. "Early test scores, socioeconomic status, and future outcomes." NBER Working Paper no. 6943.

DiNardo, J N., M. Fortin and T. Lemieux. 1996. "Labor Market Institutions and the Distribution of Wages, 1973-1992." *Econometrica* 64(5): 1001- 1044.

Ekström, E. 2002. "The value of a third year in upper secondary vocational education: Evidence from a piloting scheme." Institute for Labour Market Policy Evaluation Working Paper No. 2002:23. Uppsala, Sweden.

Federowicz, M. 2007. Badanie PISA. Umiejętności polskich gimnazjalistów. IFiS PAN, Warszawa 2007.

Fertig, M. 2003. "Who is to blame? The determinants of German students' Achievement in the PISA 2000 study." IZA Discussion Paper 739.

Fertig, M. and C.M. Schmidt. 2002. "The role of background factors for reading literacy: Straight national scores in the PISA 2000 study." IZA Discussion Paper no. 545.

Foster, P.J. 1965. "The vocational education fallacy in development planning," in C. Anderson and M.J. Bowman, eds., *Education and Economic Development*, Chicago: Aldine.

Fryer, R. and S. Levitt. 2002. "Understanding the Black-White test-score gap in the first two years of school." NBER Working Paper no. 8975.

Glewwe, P. 2002. "Schools and skills in developing countries: Education policies and socioeconomic outcomes." *Journal of Economic Literature* 40(2): 436-82.

Hanushek, E. 1986. "The economics of schooling: Production and efficiency in public schools." *Journal of Economic Literature* 24(3): 1141-1177.

Hanushek, E. 2002. "Publicly provided education," in A.J. Auerbach and M. Feldstein (eds.), *Handbook of public economics* (vol. 4). Amsterdam: Elsevier.

Heckman, J.J., H. Ichimura and P. Todd. 1998. "Matching as an Econometric Evaluation Estimator." *Review of Economic Studies* 65(2): 261-94.

Heckman, J., R. Lalonde and J. Smith. 1999. "The economics and econometrics of active labor market programs," in O. Ashenfelter and D. Card, eds., *Handbook of Labor Economics* (vol. 3A): 1865-2097.

Gruber, J. 1994. "The Incidence of Mandated Maternity Benefits." *American Economic Review* 84(3): 622-41.

Karlan, D. and M. Valdivia. 2006. "Teaching entrepreneurship: Impact of business training on microfinance clients and institutions." Yale University.

Konarzewski, K. 2008. *Przygotowanie uczniów do egzaminu: pokusa łatwego zysku*. Raport badawczy. ISP, Warszawa.

Konarzewski, K. 2004. *Reforma oświaty, podstawa programowa i warunki kształcenia*. ISP, Warszawa.

Lee, M.-J. 2005. *Micro-econometrics For Treatment-effect Analysis*. Oxford: Oxford University Press.

Meyer, B.D. 1995. "Natural and Quasi-Experiments in Economics." *Journal of Business and Economic Statistics* 13: 151-162.

Ministry of National Education. 1998. *Reform of the education system: proposal*. WSiP, Warsaw.

Murnane, R.J., J.B. Willett and F. Levy. 1995. "The Growing Importance of Cognitive Skills in Wage Determination." *Review of Economics and Statistics* 77(2): 251-66.

Oaxaca, R. 1973. "Male-female wages differentials in urban labor markets." *International Economic Review 14*(3): 693–709.

OECD. 2007. *PISA 2006 Science Competencies for Tomorrow's World*. Volume 1 and 2. OECD, Paris.

OECD. 2005. *PISA 2003 Technical Report*. OECD, Paris.

OECD. 2003. *Literacy Skills for the World of Tomorrow - Further Results from PISA 2000*. OECD, Paris.

OECD. 2002. *PISA 2000 Technical Report*. OECD, Paris.

OECD. 2001. *Knowledge and Skills for Life: First Results from PISA 2000*.

Psacharopoulos, G. 1997. "Vocational Education and Training Today: challenges and responses." *Journal of Vocational Education and Training* 49(3): 385-393.

Psacharopoulos, G. 1987. "To vocationalize or not to vocationalize: that is the curriculum question." *International Review of Education* 33: 187-211.

Psacharopoulos, G. and W. Loxley. 1985. *Diversified Secondary Education and Development: Evidence from Colombia and Tanzania*. Baltimore: Johns Hopkins University Press.

Psacharopoulos, G. and H.A. Patrinos. 2004. "Returns to investment in education: a further update." *Education Economics* 12(2): 111-134.

Rosenbaum, P.R. and D.B. Rubin. 1983. "The Central Role of the Propensity Score in Observational Studies for Causal Effects." *Biometrika* 70: 41–55.

Royston P. 2004. Multiple imputation of missing values. Stata Journal 4(3):227-241.

Tarozzi, A. 2007. "Calculating Comparable Statistics from Incomparable Surveys, with an Application to Poverty in India." *Journal of Business and Economic Statistics* 25(3): 314-336.

Todd, P.E. and I. Wolpin. 2003. "On the specification and estimation of the production function for cognitive achievement." *Economic Journal* 113: F3-F33.

Zanutto, E. 2006. "A Comparison of Propensity Score and Linear Regression Analysis of Complex Survey Data." *Journal of Data Science* 4: 67-91.