

# Evaluating Professor Value-added: Evidence from Professor and Student Matching in Physics \*

YUTA KIKUCHI<sup>†</sup> RYO NAKAJIMA<sup>‡</sup>

March 6, 2017

## Abstract

This paper estimates a professor's value added to a postgraduate student's research achievement growth using unique panel data on matched advisor-advisee pairs in a world-leading physics graduate program. To address an identification problem related to the endogenous selection of advisors and advisees, we use professor turnover and estimate a semi-parametric lower bound of the variance in advisor quality affecting advisee research performance. We find that advisor vary greatly in quality to enhance advisee's research achievement. A one-standard-deviation increase in professor quality results in a 0.54 standard deviation increase in a doctoral student's research achievement growth, increasing the number of first-authored papers that are published in top journals by 0.64 at the doctoral level.

**Keywords:** knowledge creation; postgraduate education; faculty quality; research apprenticeship

**JEL codes:** D83; I23; J24

---

\*This research was supported by Grants-in-Aid for Scientific Research Nos. 14J08184 and 26380313 from the Japan Society for the Promotion of Science. The authors thank Donna Ginther and Daiji Kawaguchi for supportive comments and helpful suggestions for improving the manuscript. We are grateful to Kenichiro Aoki and Kensuke Kobayashi for providing helpful comments on postgraduate education and academic culture in Japanese physics communities, especially regarding the the University of Tokyo physics department.

<sup>†</sup>Graduate School of Economics and Business Administration, Hokkaido University Kita 9 Nishi 7, Kita-ku, Sapporo, 060-0809 JAPAN. E-mail: yuta.kikuchi@econ.hokudai.ac.jp.

<sup>‡</sup>Corresponding Author: Department of Economics, Keio University, 2-15-45 Mita, Minato Tokyo, 108-8345, JAPAN. E-mail: nakajima@econ.keio.ac.jp.

# 1 Introduction

## 1.1 Overview

How is knowledge created? Economists have had a particular and long-standing interest in knowledge production. Indeed, the new economic growth theory literature regards the way in which knowledge is created and accumulated as crucial for a nation to grow (Romer, 1990; Lucas, 1988; Grossman and Helpman, 1991).

Recently, there has been an increase in the number of empirical studies on knowledge creation in the field of science and technology. They have investigated how an individual's knowledge creation is affected by knowledge created by others, with a particular emphasis on knowledge spillovers between individuals within and across institutions. The evidence obtained thus far, however, has been mixed. Some studies (e.g., Azoulay, Zivin and Wang, 2010; Moser, Voena and Waldinger, 2014; Borjas and Doran, 2014) provide evidence in favor of positive knowledge spillovers, while others (Waldinger, 2012; Borjas and Doran, 2012) do not.

Surprisingly little attention has been devoted to knowledge *reproduction* processes across generations. As is often argued, scientific and technological knowledge is tacit (Polanyi, 1958, 1966). It is not easily translated and thus needs to be intentionally articulated, codified and diffused. Therefore, knowledge has long been reproduced through a deliberate process of education and learning whereby those with knowledge take voluntary action to pass it on to those who do not. Thus, to investigate the creation and diffusion of scientific and technological knowledge, it seems natural to distinguish *vertical* knowledge flow (i.e., the knowledge flow from an individual with high expertise to one with low expertise) from *horizontal* knowledge flow (i.e., the knowledge flow among individuals with the same level of expertise). If the two lines of knowledge flow differ in the efficiency of the transmission of know-how, the mixed results obtained by prior studies on the extent of spillovers might be explained with regard to such differences.

This paper focuses on the knowledge reproduction process whereby knowledge is conveyed through vertical relationships, including master-apprentice, teacher-student and senior-junior-collaborator relationships. We specifically focus on the *advisor-advisee* relationship in post-graduate education to examine its effectiveness in expanding scientific frontiers. Under the hypothesis that a professor's "quality" has a consequential impact on the growth of a student's research achievement, we estimate the professor's (advisor's) *value added* as the contribution to student's (advisee's) progress on research outcomes.

Empirically estimating a graduate school professor's value added is complicated by the endogenous selection process involving students, professors and schools. We anticipate that

students of promise will apply and be admitted to highly ranked schools. Moreover, professors with good academic standing are likely to have faculty positions at prestigious schools. These types of selective recruitment will lead to nonrandom sorting of students and professors across graduate programs. Furthermore, the existence and extent of sorting can be reinforced by the advisor-advisee matching process within a school, whereby students will choose and be chosen by faculty members.<sup>1</sup>

This paper estimates *within-school* professor value added at a world-leading postgraduate program in physics in Japan. To disentangle the influence of professors on students from the sorting and matching effects, we use an identification strategy that exploits professor *turnover* from events such as retirement, relocation, or death. We borrow this idea from Rivkin, Hanushek and Kain (2005, henceforth, RHK), who estimate a lower bound of the teacher quality effect on student achievement gains. Japanese graduate schools provide an ideal setting for applying RHK’s strategy of turnover-based value-added estimation. When an advisor exits a Japanese graduate school due to turnover, the advisees usually remain in the same program and continue their research projects under the supervision of new advisor.<sup>2</sup> Therefore, the advisees who experience advisor turnover are influenced by two advisors of different quality. Thus, the student’s research achievement growth, under the varied influences of different advisors, would be more volatile than that of advisees who did not suffer advisor turnover.

Figure 1 shows the distribution of students across initial advisors (set in the vertical axis) and across cohorts (set in the horizontal axis) for the physics graduate program that we use as an empirical testing ground. We base our analysis on a *lab*, defined by a cohort of students who were assigned to the same advisor. The red square and blue circle markers represent a lab in the treatment group where the advisor was replaced due to turnover and a lab in the control group where the advisor was not replaced, respectively. We demonstrate that cohort-to-cohort variation in the lab average of the student research achievement gains is larger for the treatment group than for the control group and is driven primarily by the change in advisor quality following turnover.

### Insert Figure 1

Certainly, factors other than advisor quality might affect an advisee’s research performance. This paper employs a semi-parametric education production function, which is widely

---

<sup>1</sup>These choices give rise to “assortative matching” between students and professors within a school with respect to research ability.

<sup>2</sup>There might be a concern that the original advisor who had left the program continued to provide research guidance to the students who remained in the program. In other words, the advisor who were recorded as a thesis advisor might be a “surrogate” of the true advisor. We address such concern in Section 5.2, and show that our lower bound estimate for the variance of advisor quality is still valid under such a case of “surrogate” advisor.

used in the economics of education literature, that attributes the student’s achievement gains to various fixed effects. Repeated observations of an individual student’s research outcomes, which are measured by publication records, in master’s and doctoral degree programs enable us to eliminate student fixed effects by taking the difference of the research outcomes of a given student from one degree program to another.

The estimation results provide strong evidence for the existence of professor value added. Indeed, the results consistently demonstrate that advisor vary greatly in quality to enhance advisee’s research achievement at the doctoral level, which is consistent with the expectation that knowledge and ideas are transmitted vertically from advisor to advisee. Specifically, our estimates indicate that a one-standard-deviation increase in advisor quality will increase an doctoral advisee’s research achievement by 0.54 standard deviations. We also find that a one-standard-deviation increase in advisor quality entails an increase in the number of articles published by a doctorate student in top journals as a first author by 0.64.

The findings of this paper are robust to different definitions of student research outcomes and are also insensitive to many different model specifications. Notably, we continue to find evidence of substantial professor value added even in a severe case for students in which their research contribution to the published papers is forced to zero whenever he or she collaborated with a advisor. We also find that the results are robust to a falsification exercise that examines whether the timing of the increased variability in the student research achievement gain agrees with that of advisor turnover, as predicted by the empirical model.

We also examine the robustness of the estimation results under a weaker conditional independence assumption on advisor switch. Our identification strategy for professor value added relies on the assumption that advisor switch due to turnover is *incidental*, that is, it is orthogonal to advisee’s unobservables, conditional on advisor’s observable characteristics. Yet, the assumption may not hold. There is a concern that the advisor switch (and non-switch) is *intentional* in a sense that students, having prior information on “scheduled” turnover of faculty members, self-select themselves into or out of specific labs. To address the concern, we restrict turnover events to those that seem incidental rather than intentional. Since students would hardly predict their advisors’ relocation or decease *ex ante*, advisor switch due to such non-retirement reasons deemed closer to be idiosyncratic. It thus does not seem unrealistic to assume that it is orthogonal to unobserved student characteristics that affect the research performance. We estimate the regression model using the only turnover due to non-retirement based reasons, and still find strong advisor quality effect.

We finally investigate alternative mechanisms for knowledge transmission other than that based on learning through the advisor-advisee relationship. The data indicate that advisor turnover does not have a significant unidirectional, positive or negative, impact on an advisee’s research achievement gain, *per se*, as is consistent with the mechanism that our value-added

model postulates, and is thus not fully explained by the other mechanisms such as that emphasizing (i) a disruption effect of advisor turnover, or (ii) a recombination role of various extant pieces of knowledge, Further analysis reveals that the effect of knowledge transmission from advisor to student *within* a lab tends to outweigh that from non-advisor to student *across* labs.

While we find a significant effect of advisor quality on advisee’s research performance growth, we concede that the finding may be specific to local contexts of physics discipline or the postgraduate education system in Japan. But, it can be at the least asserted that, in a process of scientific inquiry where prominent researchers, including several Nobel laureates, have been involved, advisors matters, and their quality of supervision varies substantially. Our findings on advisor effect speaks to a broad range of literature that evaluate the heterogeneous impact of single individuals on organizational performance.<sup>3</sup> For instance, the recent study by Lazear et al. (2015) analyzes a micro-data of one large firm and provides supportive evidence for the effect of a supervisor on the productivities of the workers in the team. It seems parallel with our finding that an advisor influences significantly the performance of advisee in the lab whom he or she supervises. Further research need to be carried out and accumulated to explore whether similar supervisor effect can be observed in other institutions and organizations in both academia and industry.

The remainder of the paper proceeds as follows. A brief literature review is provided in the remainder of this section. Section 2 describes the institutional background of postgraduate physics education in Japan. Section 3 presents the empirical model and describes a regression-based approach to estimate the lower bound of professor value added. Section 4 explains the data set used for the analysis. Section 5 discusses some empirical issues concerning value-added estimation. Section 6 presents the estimation results and provides robustness checks. Section 7 concludes.

## 1.2 Related Literature

This paper contributes to the literature by measuring the effectiveness of professors in promoting students’ research productivity growth at a postgraduate institution. The most closely related work to this paper is Waldinger (2010), who estimates the causal effect of prominent professors on the research outcomes of Ph.D. students in mathematics at German universities during the Nazi era. Although we share his view that “university quality is believed to be one of the key drivers for a successful professional career of university graduates ”(Waldinger, 2010, p.787), we highlight the importance of direct interactions between advisor and advisee

---

<sup>3</sup>See, for example, Jones and Olken (2005) for political leaders; Malmendier and Tate (2009) for chief executive officers; Branch et al. (2012) for school principals; Lacetera et al. (2016) for auctioneers.

as a medium whereby knowledge is memorized, transferred and accumulated. Indeed, anecdotal evidence (e.g., Zuckerman, 1977) suggests the importance of vertical social ties in scientific enterprises at academic institutions. However, to the best of our knowledge, no systematic quantitative study, especially one that carefully controls for endogenous matching between master (teacher, advisor or senior collaborator) and apprentice (student, advisee or junior collaborator), has been conducted to date.

Our findings validate the view of earlier studies (e.g., Azoulay et al., 2010; Moser et al., 2014; Borjas and Doran, 2014) that vertical social interactions among scientists are enduring and consequential for scientific and technological knowledge to be created and diffused. For example, a recent study by Moser et al. (2014), who estimate the effect of German Jewish émigrés on U.S. innovation, suggests that knowledge externalities occurred and were amplified through educational and collaborative ties in scientist networks such that U.S. junior scientists were trained by and collaborated with prominent Jewish senior scientists who emigrated. Borjas and Doran (2014) study the impact of the influx of Soviet mathematicians into the United States after the collapse of the Soviet Union and conclude that positive knowledge spillovers are generated through the relationships among collaborating mathematicians who regularly interact when at least one of them is an outstanding knowledge producer.

This study is also related to a voluminous education economics literature that evaluates teacher value added (e.g., Hanushek and Rivkin, 2006, 2010). We base our empirical analysis on the value-added model approach that is widely employed in the literature. Specifically, as mentioned above, we adopt a semi-parametric value-added model and employ the turnover estimator proposed by RHK. However, we depart from the previous literature on teacher value added in that we focus on value added at a level higher than secondary education. Although numerous studies estimate value added at the primary and secondary education levels (e.g., Hanushek and Rivkin, 2012, for a recent review), few studies (e.g., Hoffmann and Oreopoulos, 2009; Carrell and West, 2010) estimate a professor’s value added in the context of post-secondary institutions. While these studies on professor value added attempt to estimate the effectiveness of professors in improving students’ grade gains at the *undergraduate* level, we turn to professors’ value added to students’ research achievement gains at the *postgraduate* level and thus evaluate the effectiveness of professors in terms of their “quality” in advising or mentoring graduate students’ research projects.

To the best of our knowledge, no studies assess the impact of professor quality on graduate student research productivity growth by shedding light on the value-added contribution. A partial exception is the study by Hilmer and Hilmer (2009), who find a positive effect of an advisor’s research prominence on advisees’ early career publication success in U.S. economics Ph.D. programs. While they are successful in disentangling the effect of advisor quality from that of program quality on Ph.D. students’ publication outcomes, they do not address endoge-

nous advisor-advisee matching between professors and students within and across institutions. Thus, it seems questionable to interpret their finding of a positive correlation between the research productivity of advisors and advisees as causal.

## 2 Institutional Background

### 2.1 Postgraduate Physics Education in Japan

Postgraduate education in Japan, including in physics, has a two-tiered structure, that is, a two-year master's degree program followed by a doctoral program that typically lasts three or four years.<sup>4</sup> Leading Japanese research universities typically offer both master's and doctoral courses. In most cases, students enrolled in a doctoral degree program graduate with a master's degree from the same school. However, they are institutionally separated. Thus, a master's student seeking to pursue a doctorate must take an entrance examination, which is largely based on a master's thesis, to be admitted to a doctoral course even if it is offered by the same institution. In a sense, the master's degree program implicitly serves as a screening device for doctoral programs in Japan.

Three features are notable for graduate education in physics for master's programs in Japan. First, Japanese physics master's students are closely linked to their faculty advisors immediately after enrollment in a program. Indeed, applicants to a master's degree program must declare their desired field of specialization and submit a short list of faculty advisors from whom mentorship is sought upon admission. Only those who are approved for support by designated advisors are admitted to a graduate school.<sup>5</sup>

Second, physics education in Japan at the master's level is best characterized by research-based apprentice training, which is often contrasted by coursework-based training in the U.S. (Abe and Watanabe, 2012). Although Japanese master's students in physics are required to take some "coursework" credits toward their degrees, they can earn most of their credits through learning-by-doing style research "seminars" taught by a faculty advisor.<sup>6</sup>

Finally, for Japanese physics graduate students, a thesis is required to complete the master's program. It is expected to be original, as a doctoral thesis should be, although they are evaluated according to different criteria of scholarly maturity. Students are encouraged to begin original research in their chosen fields at an early stage of the master's degree program

---

<sup>4</sup>The basic structure has remained unchanged since World War II, although the organizational structure of universities has been reformed (see Ushiogi, 1993; Ogawa, 2002)

<sup>5</sup>This contrasts with U.S. graduate students, who are matched with their supervisors through the rotation of faculty labs after they complete their coursework and become Ph.D. candidates (see Gumpert, 1993).

<sup>6</sup>For example, for the master's degree program in physics at University of Tokyo, students must take at least thirty credits of coursework at the master's degree level. However, lab-based research "seminars" offered by thesis advisors constitute two thirds of their total credits.

under the instruction and guidance of a faculty advisor. Because the master's thesis is a critical factor for admittance to doctoral programs, Japanese students and professors attach great importance to a master's thesis as a pathway to doctoral study.

In contrast, the doctoral programs in physics at Japanese universities are more similar to their counterparts in Western countries than are the master's programs. Specifically, Japanese doctoral students and American Ph.D candidates are considered comparable in that there is no coursework requirement. Japanese students at the doctoral level, similar to Ph.D. candidates in the U.S., begin the research for their doctoral dissertations under the supervision of their research advisors. In general terms, Japanese physics students are required to write several articles published in refereed journals as a prerequisite for a doctoral degree. These publications are usually included in a doctoral thesis.

## 2.2 Physics Labs in Japanese Universities

Interaction between a graduate student and a faculty advisor is lab-oriented in Japanese physics graduate programs. Upon enrollment in the master's program, Japanese physics students are assigned individually to a lab, and the lab's leader (or sometimes sub-leader) becomes their thesis research advisor. Students acquire the knowledge necessary to conduct their research through frequent interaction with their advisors in a lab setting. The content of this lab-based teaching and learning includes basic research skills, such as how to read scientific articles, how to select research topics, how to present results at seminars and conferences, and how to write publishable papers, as well as the culture of physics such as the style of work, mode of thought, and a taste for "good" physics (Abe and Watanabe, 2012).

Japanese physics labs are generally democratic in tone. The "laboratory democracy" in Japanese physics communities can be traced back to the end of World War II, the period when there were immediate and insistent calls for the creation of a new "scientific Japan" under the control of the allied occupation (Low, 2005). To place this in perspective, it is broadly understood that Japanese physics labs are less prescriptive and less hierarchical than their U.S. counterparts (Kawashima and Maruyama, 1993; Gumport, 1993). There is also no strict division of labor among lab members, even between faculty members and graduate students, in Japanese physics labs (Traweek, 1988).<sup>7</sup> Hence, although it is not uncommon for the research topics of master's and doctoral theses to be suggested by advisors as a part of a large, ongoing project in a given lab, Japanese physics graduate students are, generally, given some autonomy to pursue their own research based on their original ideas.

---

<sup>7</sup>Interested readers should consult Appendix A that explains in detail the "laboratory democracy" among Japanese physicists.



### 3 Empirical Model

In this section, we introduce a simple value-added model that associates growth in student research achievement with the “quality” of the professor supervising the student. Then, we present a regression-based approach to estimate a lower bound of the variance in professor quality, which can be interpreted as the extent to which any professor differences matter in determining student research outcome growth.

#### 3.1 Value-added Specification

Following the standard value-added modeling approach (e.g., Hanushek and Rivkin, 2010), we employ a semi-parametric specification of a professor’s contribution to a student’s achievement growth.

Consider graduate student  $i$  who entered the master’s program of a graduate school in year  $c$ . Below, we treat year  $c$  as the student’s cohort. We denote the research outcome growth of a graduate student in the master’s degree program by  $g = m$  and in the doctoral degree program by  $g = d$ . The growth is measured by the *gains* in research output from the previous degree program to the current degree program.<sup>8</sup> Let  $\Delta outcome_{iag}^c$  be the research outcome growth of student  $i$  under the supervision of professor  $a \in \mathcal{A}$  in degree program  $g \in \{m, d\}$  in cohort  $c \in \mathcal{C}$ . We assume that it is given by the following function:

$$\Delta outcome_{iag}^c = \gamma_i + \theta_{ag} + \nu_{iag}^c, \tag{1}$$

where  $\gamma_i$  is student  $i$ ’s individual fixed effect,  $\theta_{ag}$  is professor  $a$ ’s quality that influences the student research outcome growth in degree program  $g$ , and  $\nu_{iag}^c$  is an idiosyncratic random shock.<sup>9</sup>

We assume that matching between student and professor is many-to-one, that is, multiple students are assigned to one advisor. Let us define a *lab* as a group of students (advisees) in the same cohort who were assigned to the same professor (advisor). Specifically, we use  $\ell(a, c)$  to denote a lab in which students are in cohort  $c$  and assigned to professor  $a$  as an advisor. Let  $L$  be the number of all labs in a school, and let students in lab  $\ell(a, c)$  be indexed by  $i = 1, \dots, I^{\ell(a, c)}$ , where  $I^{\ell(a, c)}$  is the number of students in lab  $\ell(a, c)$ . We use  $\mathcal{I}^{\ell(a, c)} \equiv \{1, \dots, I^{\ell(a, c)}\}$  to denote the set of students in the lab.

---

<sup>8</sup>We assume that the research output of students at the bachelor level is zero. We compute a publication-based research proficiency score, which is explained in detail in Section 4, for students in the sample when they are undergraduate students and find that it is negligible.

<sup>9</sup>Note that other effects, such as school fixed effects and research field fixed effects, are not included in the value-added model. We opt not to include these fixed effects because they are subtracted out of the estimation model in the process of “differencing”, as presented below.

We take the average of Equation (1) over all students in the same lab  $\ell(a, c)$ . Because the students in the same lab have the same advisor quality, we have the following equation for the lab-level average of the student research outcome growth:

$$\overline{\Delta outcome}_{ag}^{\ell(a,c)} = \bar{\gamma}^{\ell(a,c)} + \theta_{ag} + \bar{\nu}_{ag}^{\ell(a,c)}, \quad (2)$$

where the overbar notation indicates the group average.

Note that the *superscript*  $a$  denotes the *initial* advisor to whom the students in lab  $\ell(a, c)$  were assigned, while the *subscript*  $a$  denotes the advisor who supervised the students in degree program  $g$ . Thus, the advisors represented by the superscript and subscript could be different. For example, suppose that a turnover incident causes the students in lab  $\ell(a, c)$  to switch their research advisor from professor  $a$  in the master's degree program to professor  $b$  in the doctoral degree program. Here, the average student research outcome gain at the doctoral level, which is the left-hand side of Equation (2), is given by  $\overline{\Delta outcome}_{bd}^{\ell(a,c)}$ , where the index  $a$  in the superscript differs from the index  $b$  in the subscript.

We use the event of professor turnover (e.g., retirement, relocation and decease) to identify the variance in professor quality. We implicitly assume that, when a professor exits a graduate program due to turnover, the students in the lab whom he or she initially supervised are re-assigned to a new advisor and continue their research projects in the same program.<sup>10</sup> In what follows, we therefore assume that an event of professor turnover on the faculty side leads to an event of advisor *switch* on the student side. In other words, we treat these two events, advisor turnover and advisor switch, identically. When advisor turnover occurs in a lab, two faculty members, whose quality levels are generally different, advised students in the lab.<sup>11</sup>

It should be noted that the professor, say  $b$ , who was assigned to the students in the lab of a professor, say  $a$ , after the latter exited due to turnover was not necessarily drawn at random from a pool of professors available at the school at that time. Indeed, the newly assigned professor might select the students that he or she is willing to take over. We thus allow the student fixed effect,  $\gamma_i$ , to be correlated with the quality of the re-assigned professor,  $\theta_{bd}$ , in the same way as we assume it to be correlated with the quality of the original advisor,  $\theta_{ad}$ .

---

<sup>10</sup>A joint transfer of faculty and students is quite rare in Japanese universities, and hence, even if a faculty member changes affiliation, the students usually remain in the same program.

<sup>11</sup>Based on the observed pattern of advisor replacement in our data, when advisor turnover occurred, the students were usually either assigned to a junior faculty member or the sub-leader of the same lab or they were moved to a different lab in closely related research fields within the same institution and were supervised by the faculty member who managed that lab.

### 3.2 A Lower-bound Estimation of the Variance in Advisor Quality

We are interested in estimating the *variance* for advisor’s quality,  $\sigma_g^2$ , rather than the the advisor’s *individual* quality,  $\theta_{ag}$ . Although the individual value added would be useful to answer important questions concerning who or which types of professors tend to have higher or lower effectiveness in research supervision, such estimation requires more extensive data and stricter identification assumptions than we adopted in this study.<sup>12</sup> We therefore set more modest goal, that is, to measure a *minimum degree* of speculum, if any, with which each professor contributes to the performance growth of his or her lab students.

For the purpose, we decompose the total variation in student outcome gains into the variation that can be attributed to professor quality,  $\theta_{ag}$ . First, take the difference of Equation (2) between the master’s degree and doctoral degree programs. Doing so eliminates the student fixed effect,  $\gamma_i$ , because it is constant across degree programs for a given student. If advisor turnover did not occur in lab  $\ell(a, c)$ , it is given by the following between-degree difference equation:

$$\overline{\Delta outcome}_{ad}^{\ell(a,c)} - \overline{\Delta outcome}_{am}^{\ell(a,c)} = (\theta_{ad} - \theta_{am}) + (\bar{v}_{ad}^{\ell(a,c)} - \bar{v}_{am}^{\ell(a,c)}). \quad (3)$$

In contrast, assume that there was advisor turnover in lab  $\ell(a, c)$ . As the students switched their advisors from advisor  $a$  in the master’s program to advisor  $b$  in the doctoral program, the between-degree difference equation, corresponding to Equation (3), is given by:

$$\overline{\Delta outcome}_{bd}^{\ell(a,c)} - \overline{\Delta outcome}_{am}^{\ell(a,c)} = (\theta_{bd} - \theta_{am}) + (\bar{v}_{bd}^{\ell(a,c)} - \bar{v}_{am}^{\ell(a,c)}). \quad (4)$$

Comparing Equations (3) and (4) shows that advisor turnover influences the development of student research achievement in different ways. There is a clear difference in student research outcome growth, which appears on the left-hand side of each equation, that responds differently to a change in advisors due to the difference in degree-level advisor effects,  $(\theta_{ad} - \theta_{am})$  and  $(\theta_{bd} - \theta_{am})$ , which are generally not equal. This plays a key role in the identification of the effect of advisor quality on student research outcome growth at each degree level.

The point is illustrated by Figure2, which depicts three labs with different cohorts,  $c_0$ ,  $c_1$  and  $c_2$ , whose initial advisor is professor  $a$ . In the figure, each lab is portrayed by a connected line segment, which represents the two-year master’s degree program (the first half of the segment) and the three-year doctoral degree program (the last half of the segment).<sup>13</sup> Here, advisor turnover did not occur in labs  $l(a, c_0)$  or  $l(a, c_1)$  before cohort  $c_2$ , and hence, the

<sup>12</sup>See Chetty et al. (2014), for example, about techniques to estimate a reliable individual teacher value added.

<sup>13</sup>For the ease of exposition, the labs’ cohorts are not overlapped in the figure, although this is not necessarily the case in the actual sample.

students in these labs were supervised by the same professor,  $a$ , throughout both the master's and doctoral programs. However, in lab  $l(a, c_2)$ , professor  $a$  exited the school due to turnover, and professor  $b$  took charge of the doctoral students.

Note that, on average, the research outcome gains of lab  $l(a, c_0)$  and  $l(a, c_1)$  students are the same, which is given by  $(\theta_{ad} - \theta_{am})$ , whereas, following advisor turnover, the average student research outcome gain of lab  $l(a, c_2)$ , which is given by  $(\theta_{bd} - \theta_{am})$ , could be better or worse than those of the previous cohorts, depending on whether the supervising quality of the newly assigned professor,  $b$ , is higher than that of the departing professor,  $a$ . In either case, irrespective of whether the achievement growth is positive or negative, an instance of turnover triggers a change in professor quality at the doctoral level and could thus result in a disparity in the between-degree research achievement gains between cohorts. We will use the induced divergence in research outcome growth as evidence of an advisor's impact on an advisee.

### Insert Figure 2

To improve the identification, we use the *double-differencing* approach as proposed by RHK to estimate a lower bound of the variance in unknown teacher quality. We take the difference of Equations (3) and (4) with respect to cohort year. Let  $c'$  denote the cohort before  $c$ , and let  $\tau$  be the years between  $c$  and  $c'$ . For professor  $a$ , consider two labs,  $\ell(a, c)$  and  $\ell(a, c')$ . Let  $W^{\ell(a, c, c')}$  denote a dummy variable indicating a change in advisor due to turnover: it takes value one if professor  $a$  is replaced in lab  $\ell(a, c)$  due to turnover and zero otherwise. Without loss of generality, we assume that supervisor replacement is from professor  $a$  to professor  $b$  such that, if there were advisor turnover, the students would have been supervised by two different professors,  $a$  and  $b$ , in the master's and doctoral degree programs, respectively. Then, we have the following double-differenced (*DD*) average student research outcome growth:

$$\begin{aligned}
& DD \overline{\Delta outcome}^{\ell(a, c, c')} \\
= & \begin{cases} [\overline{\Delta outcome}_{bd}^{\ell(a, c)} - \overline{\Delta outcome}_{am}^{\ell(a, c)}] - [\overline{\Delta outcome}_{ad}^{\ell(a, c')} - \overline{\Delta outcome}_{pm}^{\ell(a, c')}] & \text{if } W^{\ell(a, c, c')} = 1 \\ [\overline{\Delta outcome}_{ad}^{\ell(a, c)} - \overline{\Delta outcome}_{am}^{\ell(a, c)}] - [\overline{\Delta outcome}_{pd}^{\ell(a, c')} - \overline{\Delta outcome}_{pm}^{\ell(a, c')}] & \text{if } W^{\ell(a, c, c')} = 0 \end{cases} \\
= & \begin{cases} (\theta_{bd} - \theta_{ad}) + \text{error term} & \text{if } W^{\ell(a, c, c')} = 1 \\ \text{error term} & \text{if } W^{\ell(a, c, c')} = 0, \end{cases} \tag{5}
\end{aligned}$$

where the *error term* is a catchall random noise term that combines the average idiosyncratic errors.

Equation (5) shows that all of the fixed effects, except for doctoral-level advisor quality, are eliminated after the double difference is taken with respect to degree programs and cohorts. The *DD* measure is more variable, on average, for the pair of labs with and without a change

in advisor ( $W^{\ell(a,c,c')} = 1$ ) than that for the pair of labs without such a change ( $W^{\ell(a,c,c')} = 0$ ). The gap is attributable to a discrete change in doctoral-level advisor quality from  $\theta_{ad}$  to  $\theta_{bd}$  due to advisor turnover. Note that advisors' quality levels can be correlated with the lab averages of student fixed effects,  $\bar{\gamma}^{\ell(a,c)}$  and  $\bar{\gamma}^{\ell(a,c')}$ , and they can also be correlated with one another, that is,  $\text{Corr}(\theta_{ad}, \theta_{bd}) \neq 0$ . In what follows, we ascribe the sample variation in the  $DD$  measure as a series of variance and covariance components of advisor quality and idiosyncratic shocks.

### The Assumption on Advisor Quality and Idiosyncratic Shocks

We make the following assumptions concerning the distribution of advisor quality.

assumption 1.1: The expectation and variance of advisor quality are given by  $E(\theta_{ag}) = \mu_g$  and  $\text{Var}(\theta_{ag}) = \sigma_g^2$ , for any  $a \in \mathcal{A}$ ,  $g \in \{m, d\}$ , and  $c, c' \in \mathcal{C}$ .

assumption 1.2: The correlation of advisor quality across professors,  $a \neq b \in \mathcal{A}$ , is given by  $\text{Corr}(\theta_{ag}, \theta_{bg}) = \rho_g$ , for any  $a, b \in \mathcal{A}$ ,  $a \neq b$ ,  $g \in \{m, d\}$  and  $c, c' \in \mathcal{C}$ .

The assumption concerns the stationarity of the advisor quality distribution, which characterizes the notion that the professors' advising quality levels are drawn from a common distribution for each degree type. It requires that the grade-program-specific mean and variance do not vary across cohorts and that the correlation with any given advisor is constant. Specifically, we interpret  $\mu_g$  and  $\sigma_g^2$  as the long-run mean and variance of the stationary distribution of advisor quality in degree program  $g$  within a school. The stationarity assumption simplifies the estimation of professor value added because it reduces the number of parameters to be considered.

Appendix B.1 present the assumptions on the moments of the idiosyncratic shock after *demeaning* by each cohort. Let  $\bar{\nu}_g$  be the average of the random shock  $\nu_{iag}^c$ , the average of which is taken over all cohorts in each degree program,  $g$ , such that the *demeaned* random shock is given by  $\tilde{\nu}_{iag}^c = \nu_{iag}^c - \bar{\nu}_g$ . We assume that the random shocks demeaned by cohort are independent of turnover incidents (assumption 2.1). They can be serially correlated between degree programs within a student (assumption 2.2) and between students in each degree program if they are supervised by the same advisor (assumption 2.3). Otherwise, they are neither cross- nor serially correlated (assumptions 2.4 and 2.5). Note that, even if the demeaned random shock,  $\tilde{\nu}_{iag}^c$ , is uncorrelated with others under assumptions 2.4 and 2.5, the original random shock,  $\nu_{iag}^c$ , is allowed to be correlated through the common mean factor,  $\bar{\nu}_g$ .

## The Regression Model

Finally, given the assumptions presented above, we square both sides of Equation (5) and take the expectation conditional on the occurrence of turnover. We have the following result:<sup>14</sup>

$$\text{E} \left[ \left( \overline{DD\Delta outcome}^{\ell(a,c,c')} \right)^2 \middle| W^{\ell(a,c,c')} \right] = \alpha \left( \frac{1}{I^{\ell(a,c)}} + \frac{1}{I^{\ell(a,c')}} \right) + \{2\sigma_d^2(1 - \rho_d)\} W^{\ell(a,c,c')}, \quad (6)$$

where  $\alpha$  is a composite term of variances and covariances of the demeaned random shocks, and

Equation (6) provides a basis for estimating the variance in advisor quality at the doctoral level. Using the cohort examples,  $c_0, c_1$ , and  $c_2$ , that are depicted by Figure 2 for illustration, the squared difference measure of student research outcome growth, which is the right-hand side of Equation (6), is greater for  $\ell(a, c_1, c_2)$  than that for  $\ell(a, c_0, c_1)$  by  $2\sigma_d^2(1 - \rho_d)$ . We can therefore ascribe the large sample variation of the right-hand side of Equation (6), if any, to the variance in doctoral-level advisor quality,  $\sigma_d^2$ , unless the correlation coefficient,  $\rho_d$ , is equal to one.

We now present a regression model to obtain a lower-bound estimate of the variance of  $\sigma_d^2$ . Consider the following:

$$\left( \overline{DD\Delta outcome}_n \right)^2 = \alpha X_n + \beta W_n + \varepsilon_n, \quad (7)$$

where  $n = 1, \dots, N$  is the index of observations. Here, the unit of observation is each element of  $(a, c, c')$  for any advisor  $a \in \mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$ , where  $\tau$  is the period over which the difference is taken.<sup>15</sup> Note that, analogous to Equation (6), the covariate  $X_n = 2(1/I^{\ell(a,c)} + 1/I^{\ell(a,c')})$  is introduced into the regression. The random term  $\varepsilon_n$  is interpreted as the prediction error between the expected and observed values of the divergence measures, that is:

$$\varepsilon_n \equiv \text{E} \left[ \left( \overline{DD\Delta outcome}_n \right)^2 \middle| W_n \right] - \left( \overline{DD\Delta outcome}_n \right)^2.$$

Assume for a moment that the advisor switch indicator,  $W_n$ , is independent of the prediction error,  $\varepsilon_n$ . If the value of  $\rho_d$  were known perfectly, the OLS estimate  $\hat{\beta}$  in Equation (7)

<sup>14</sup>See Appendix B.2 for the derivation.

<sup>15</sup> To obtain the double-differenced average of the research outcome gain, which is the left-hand side of Equation (7), we take the difference between all cohorts within a period of  $\tau$  years. As  $\binom{\tau+1}{2} = \frac{\tau!}{(\tau-2)!2!}$  samples are created for each lab, the total sample size of the regression is given by  $N = \frac{\tau!L}{(\tau-2)!2!}$ , where  $L$  is the total number of labs.

would provide a consistent estimate of  $\sigma_d^2$  through the following equation:

$$\hat{\beta} = \{2\hat{\sigma}_d^2(1 - \rho_d)\}. \quad (8)$$

As the correlation is imperfect ( $\rho_d < 1$ ), a lower-bound estimate of  $\sigma_d^2$  is given by the last term of the following equation:

$$\hat{\sigma}_d^2 = \frac{\hat{\beta}}{2(1 - \rho_d)} \geq \frac{\hat{\beta}}{4}. \quad (9)$$

In other words, a lower-bound estimate of the within-school variance of faculty quality at the doctoral level is equal to the estimated coefficient,  $\hat{\beta}$ , of the regression model (7) divided by four.

## 4 Data

We assemble data sets of professors and students in a graduate program in physics in Japan. Among the numerous Japanese research universities that offer both master's and doctoral programs in the field of physics, we focus on the graduate program at University of Tokyo (henceforth, UTokyo), which is the oldest institution of its kind in the country and has enjoyed high prestige in the global academic community.<sup>16</sup>

The graduate program in physics at UTokyo consists of the department of physics as its core and other physics-related research institutes on campus.<sup>17</sup> The average number of graduates in recent years is 105.6 for the master's program and 58.4 for the doctoral program<sup>18</sup>. At present, there are more than 130 full-time faculty members. Many subfields of physics are covered by laboratories in UTokyo's physics graduate programs, such as nuclear physics, particle physics, condensed matter physics, and biophysics.

### 4.1 Data on Advisor and Advisee Pairs

To extract the information on matched advisor-advisee pairs, we use the master's and doctoral thesis catalogs for graduate students in UTokyo's physics program.<sup>19</sup> For each thesis entry in the catalog, the available information includes the degree date, the title of the thesis, the

---

<sup>16</sup>According to several world university rankings, UTokyo has been in the top 10 in the discipline of physics. The alumni include five Nobel laureates in physics as of 2015.

<sup>17</sup>The institutes are the Institute of Cosmic Ray Research, (ICRR), the Institute of Solid State Physics (ISSP), and the International Center for Elementary Particle Physics (ICEPP).

<sup>18</sup>These are the average figures over the period from 2010 to 2014.

<sup>19</sup>The catalogs are available on the department's website at <http://www.phys.s.u-tokyo.ac.jp/TOSH0/ronbun.html>.

name of the student, and the name of the faculty advisor who supervised the student.

We compile the thesis data for the students who obtained their doctoral degrees in the cohorts between 1970 and 2004 (35 years). Among all of the graduate students who were listed in both the master’s and doctoral thesis catalogs, we restrict our attention to those who earned doctorates within six years of enrollment. In addition, we restrict the analysis to those who were supervised by faculty members with the ranks of full and associate professors in the physics department or on-campus physics-related research institutions. There are 119 advisors (professors) and 1484 advisees (students).<sup>20</sup> The average numbers of students is 1.5 for each lab and 12.5 for each advisor in the sample.<sup>21</sup>

## 4.2 Data on Advisor Turnover and Switch

We obtain information on faculty turnover from the University Personnel Directory Book (“*Zenkoku Daigaku Shokuin Roku*”) published by *Koujyun Sha*, which includes information on the full name, rank, department, school, specialized fields and year of birth of all staff members at every Japanese university, public or private, in a given year. By compiling the roster of faculty members at UTokyo, we can obtain their turnover information.

We classify the reasons for turnover into the following three categories: (1) retirement if the instance of turnover occurred at the mandatory retirement age predetermined by UTokyo;<sup>22</sup> (2) move if turnover occurred before the retirement age and the faculty name began to reappear on other universities’ rosters beginning in the year after the turnover instance; and (3) decease/quit otherwise.<sup>23</sup> Figure 3 presents the graphs that plot the number of turnover incidents in each year of the sample period, broken down by the reasons.<sup>24</sup>

### Insert Figure 3

---

<sup>20</sup>We restrict advisors in the sample to those who supervised at least two cohorts of students within three years because, as shown in later sections, the observation unit of the baseline regression is set as the difference of two consecutive cohorts within three years.

<sup>21</sup>It should be noted that the reported average sizes of students per lab and per professor are undersized comparing to the actual sizes that make up the graduate program because the estimation samples are restricted to the students who took both master’s and doctoral degrees at the UTokyo’s physics.

<sup>22</sup> Before fiscal year 2000, the mandatory retirement age at UTokyo was 60. After the 2001 fiscal year, it was increased by one year every three years until it reached 65. As of 2004, which is the end of the sample period, the retirement age was 61.

<sup>23</sup>Note that the reasons for faculty turnover are not perfectly distinguishable. Indeed, the majority of faculty members categorized as “retire” did not actually retire from academic life and were reemployed at other universities or research institutions. This is possible because of the gap in retirement ages between universities: UTokyo set its faculty retirement age at 60 during the most of the sample period, while other Japanese universities, public and private, adopted retirement ages that were several years older.

<sup>24</sup>There is a considerable number of incidents in 1997, when the Institute for Nuclear Study (INS) at UTokyo, which was one of the on-campus research institutes affiliated with the physics department, was closed and merged with the High Energy Accelerator Research Organization (also known as the KEK (Kō Enerugi Kasokuki Kenkyū Kikō)), and some of the faculty members at the INS chose to leave UTokyo for the KEK.



The matched advisor-advisee data reveal that approximately 14.4 percent of graduate students switched advisors between the master’s program and the doctoral program. Instances of professor turnover are responsible for some, although not all, of the students’ observed changes in advisors. As mentioned previously, in Japanese universities, a joint transfer of faculty member and student is quite rare. If a faculty member exits a graduate program, another other faculty member – usually a sub-leader of the same lab or, sometimes, a faculty member from a different lab in the same institution whose research area is closely related to the professor who exited – becomes the new advisor of the students who are left behind. In either case, the student remains in the same program.

We identify an advisor switch due to turnover if a student’s master’s thesis advisor exited UTokyo before the student earned a doctoral degree. Such cases account for 53.2 percent of all advisor switches in the sample. We exclude students who switched advisors on their own initiative from the sample observations, as such student-side advisor switches are likely to be caused by a mismatch between advisor and advisee and could be correlated with student research outcomes.

From Figure 1 where the red square marker represents a lab where an advisor switch due to turnover occurred, while the blue circle marker represents a lab with no advisor turnover (the darker color marker means more students in a lab), it appears that faculty turnover incidents are relatively frequent. It thus seems to exist sufficient variation in the dependent variable of the regression equation (7) to identify the variance of advisor effectiveness on advisee’s research achievement gain within the same institution.

### 4.3 Data on Student Research Achievement

To measure a graduate student’s research achievement, we use the number of journal articles that he or she published. To obtain this information, we employ the Thomson Reuters Web of Science (WoS) archive. We collect physics articles with author names that match the name of the graduate student under consideration. We further restrict our attention to those articles published around the period when the target student was enrolled.

The articles selected by author name matching may contain false positives: these articles could have an author who coincidentally has the same name as the graduate student in the sample but is in fact a different person. To minimize such identification errors, we add a further restriction; that is, for an article to be identified as written by the student in question, we impose a restriction that the words in the article title should overlap to some extent with those in the title of the master’s or doctoral thesis.<sup>25</sup>

Based on a student’s publication records, we define the *research proficiency score* as the

---

<sup>25</sup>See Appendix C.1 for details on the score of word overlap in titles.

number of publication counts during *a given year*. Here, we employ two quality adjustment methods. First, we limit the publications to those published in twelve high-quality peer-reviewed journals, including three high-reputation general-interest science journals and nine highly ranked physics journals.<sup>26</sup> Second, we consider a student’s share of credit for an article if there are multiple authors. In physics, as in other scientific disciplines, papers are usually written by a group of authors whose contributions are not necessarily equal. We follow a standard bibliometric method (e.g., Liu and Fang, 2012; Waltman, 2012) based on the byline hierarchy rule to quantify an coauthor’s share of credit for an article with multiple authors.<sup>27</sup>

Figure 4 plots the average research proficiency score for our sample graduate students in each year. Note that, in the figure, we begin the graduate school year index at one in the year when a student entered the master’s program and increase it throughout the duration of the graduate program. For the sake of expedience, the graduate school year is also defined for the postdoctoral period after the student obtained a doctorate degree. In the figure, it corresponds to the period after the 6th year.

#### Insert Figure 4

The figure illustrates the time pattern of how physics graduate students at UTokyo develop their research outcomes: the achievement curve rises and reaches its peak in the years near the completion of the doctoral degree (D1 and P1). Then, the research outcomes begin to decline during the postdoctoral periods (P1-P5). We suspect that this reflects two types of lag structure: the first relates to a publication lag, that is, the time lag from the submission to publication of articles in journals. The second concerns a gestation lag, that is, the time lag between project inception and completion.

## 5 Empirical Issues

In this section, we describe the empirical issues involved in estimating a lower bound of professor quality based on the regression model in Equation (7). We first address how to construct the squared difference measure of the student outcome growth variable, which is used as the dependent variable in the regression model. We next point out the possibility that supervisory period is overestimated, and discuss the bias in the lower bound estimate of advisor’s value added. We finally take up the non-randomness of professor turnover, which could cause an endogeneity problem and thus threaten the validity of the estimates. We then

---

<sup>26</sup>*Nature*, *Science* and *Proceedings of the National Academy of Sciences of the United States of America (PNAS)* are included as the general-interest science journals, and *Physical Reviews A, B, C, D, and E*; *Physical Reviews Letters*; and *Physics Letters A and B* are included as the top physics journals. We received advice from physicists regarding the selection of the top journals.

<sup>27</sup>See Appendix C.2 for details on computation of coauthor’s credit share.

propose a method to address this endogeneity concern.

## 5.1 Student Research Outcome Variable

Two issues arise regarding the double-differenced student research outcome measure: (i) the choice of years over which the student research outcomes are aggregated at the program level and (ii) the choice of interval years between the pair of cohorts that are differenced.<sup>28</sup>

To address the first issue, we examine the distribution of the student average research proficiency scores over the years at the graduate program. Figure 5 decomposes the student average research proficiency scores into those related to the master’s thesis and those related to the doctoral thesis. It is shown that the proficiency score associated with the master’s thesis peaks in the second year of the doctoral program (D2) and decreases thereafter, while the score related to the doctoral thesis continues to increase. We thus opt to aggregate the research proficiency scores over the period from the first year of the master’s program (M1) to the second year of the doctoral program (D2) to compute the research outcome at the master’s level. As for the research outcome at the doctoral level, we assemble the research proficiency scores from the first year of the doctoral program (D1) up to the fourth year of the postdoctoral period (P4). We choose a rather long aggregation period at the doctoral level in light of the lag between the time of article publication and the time the degree is awarded.

### Insert Figure 5

In sum, our benchmark student research outcomes are aggregated over the period from M1 to D2 and the period from D1 to P4 for the master’s degree and doctoral degree programs, respectively. Table 1 presents the descriptive statistics.<sup>29</sup>

### Insert Table 1

We turn to the second issue concerning the interval in years between cohorts. In Section 3.2,  $\tau$  denotes the number of years between two cohorts,  $c$  and  $c'$ , such that  $c - c' \leq \tau$  when determining the  $DD$  measure of the student research achievement growth. Note that there is no theoretical rule for which year should be used as  $\tau$ . On the one hand, the longer the interval is, the more efficient the estimator because it yields more samples for the regression analysis.<sup>30</sup> On the other hand, the shorter interval is, the better because it requires a weaker assumption on the covariance stationarity of the distribution of the demeaned random shocks

---

<sup>28</sup>Because the value-added model focuses on the student research achievement gain while in school, the magnitude may be minute and unnoticeable if it is measured by the annual gain. We thus select the unit of measure as each *degree program* period.

<sup>29</sup>The box plots of the research outcome distributions at the master’s and doctoral levels are presented in Table D.2 of Appendix D.1.

<sup>30</sup>As presented in footnote 15, the total sample size is given by  $N = \frac{\tau!L}{(\tau-2)!2!}$ , which is an increasing function of the adjacent period,  $\tau$ , *ceteris paribus*.

(assumption 2.3).<sup>31</sup> In light of balance, we adopt the adjacent cohort period of  $\tau = 3, 4$  and 5 years as the benchmark when implementing the regression.

## 5.2 Overestimation of Supervisory Period

Another empirical issue concerns imprecise measurement of supervisory period. We have thus far implicitly assumed that, when the incident of advisor turnover happened, the professor to which a student were reassigned supervised him or her for the entire doctoral course period. Yet, the assumed period of supervisory may be overstated. For instance, suppose that a student experienced advisor turnover in the middle of his or her doctoral course. Then the advisor who were newly assigned may supervise the student for *shorter* period than the entire period of the doctoral course.

The issue of overestimated supervisory period also arises when the professor who is recorded as a research advisor in the thesis catalog is a “surrogate” of the true advisor. In some cases of turnover, especially in cases where advisor relocates to other school, the original advisor who left the program continues to provide substantial research influence on the doctoral project of the former student. In this case, even if the student who experienced advisor turnover was assigned a new advisor, it is nominal. The substantive research guidance is not provided by the “surrogate” advisor but by the “hidden” advisor. The length of supervision of the “surrogate” advisor should be considered zero.

In Appendix D.1, we discuss how overstated supervisory period affects the the lower bound estimate of the variance of advisor quality, and show that the estimate is biased *downward*. Nevertheless, since the estimated value is considered as a conservative lower limit of the advisor effectiveness on advisee’s research achievement gain, it turns out that the interpretation of the estimate will not be changed.

## 5.3 Non-Random Turnover

Thus far, we have assumed that professor turnover is independent of various factors in the value-added model and thus does not affect student research performance except through the change in advisor quality. However, the assumption may be untenable. Arguably, a professor’s decision of whether to retire, move, or remain at a graduate program may be endogenous to the student’s performance. Table D.1 in Appendix D presents evidence suggesting that the sample is not balanced, that is, there are systematic differences between the groups with and without professor turnover on some characteristics.

---

<sup>31</sup>To be more precise, assumption 2.3 states that the covariance of the demeaned error terms is constant between any two students,  $i$  and  $j$ , in different cohorts,  $c$  and  $c'$ . This assumption may be reasonable only for adjacent cohorts.

Consequently, the regression model in Equation (7) may suffer from the standard endogenous variable problem, as the catch-all error term,  $\varepsilon_n$ , which influences student research outcome growth, will be confounded by the advisor switch dummy variable,  $W_n$ , through the heterogeneity of advisors, who systematically differ between those with and without turnover. In this case, we may not be able to obtain an unbiased estimate of  $\beta$  from the regression and thus be unable to obtain a reliable estimate of the lower bound of advisor quality.

To make the sample balanced and comparable, we thus employ a propensity score matching method. The basic idea is to match a turnover case with a case of no turnover that has approximately the same conditional likelihood, typically called the propensity score, that an incident of advisor turnover would have occurred. After constructing a new balanced sample based on the propensity score matching procedure, we estimate the regression model in Equation (7) using the balanced sample, as if advisor changes due to turnover occurred at random.

Note that, to account for the endogeneity of the advisor switch dummy variable in the regression model, we only control for advisor characteristics. It is potentially justifiable not to balance the sample on student characteristics because we exclude all cases in which a change in advisor occurs for a student's own reasons, as described in Section 4.2. The sample restriction can eliminate the possibility that student factors are confounded with the occurrence of an advisor switch, and therefore, it is deemed to occur exclusively for reasons on the faculty side. Hence, we control for the professor's characteristics in the propensity score analysis.

Appendix D.2 presents a detail procedure for the propensity score matching method that we employ to address the issue of non-random advisor turnover. Following standard practice in the literature, we estimate the propensity scores using a logit model. We determine a baseline specification of the model by a stepwise likelihood-test-based procedure, suggested by Imbens (2014) and Imbens and Rubin (2015). To address the problem caused by the limited common support of the propensity score distribution, we employ a systematic approach proposed by Crump et al. (2009).

## 6 Estimation Results

### 6.1 Benchmark Results

This section presents the estimation results for professors' value added to the students' research achievement gains.

## Variance Comparison

To make it easier to eyeball the results, we give a graphical representation of how the  $DD$  measure of the student research achievement growth, given by Equation (5), fluctuates when the incident of advisor turnover happened. As explained, we ascribe the observed change of the variance in the  $DD$  measure to disconnected change of advisor qualities due to turnover.

Figure 6 (left) presents the sample variances of the  $DD$  measure for (a) labs in the cohort where turnover occurred and (b) labs of the same advisor but their cohort is the latest one *before* turnover. Although the variances become larger once turnover happened for all the three cases of adjacent cohort periods, the changes are not substantial. Figure 6 (right) repeats the same comparison of the sample variances. But, this time, with the non-randomness of turnover being taken into account, it compare the one for (c) labs in the treatment group where turnover occurred with the one for (d) labs in the corresponding control group that are matched through the propensity score method. In contrast to the previous result, it is shown that the increase of the variances are considerable for all three cases of adjacent periods.

**Insert Figure 6**

## Baseline Regression Estimates

To quantify the magnitude of a professor's value added, we estimate the econometric model (7). The main estimate of interest is the lower bound of the variance in advisor quality at the doctoral level, which is given by one-fourth of the coefficient of the advisor switch indicator variable in the regression model.

Table 2 presents the baseline results. We report the regression estimates in rows (1) and (2). Columns (1), (2) and (3) are used to report the estimation results for the three cases of adjacent periods between cohorts,  $\tau = 3, 4$  and 5 years, respectively. As the estimated propensity scores are used for the true values, we compute resampling-based standard errors to correct for the additional sampling variability arising from estimation.<sup>32</sup> All estimates of  $\beta$ s are positive and statistically significant from zero at the 10 percent level except for one case.

**Insert Table 2**

Row (3) of Table 2 presents the estimated lower bound of advisor quality variance at the doctoral level. As the variance must be non-negative, we perform one-sided tests such that  $\overbrace{\text{Lower bound of } \sigma_d^2 = 0}^{\text{Lower bound of } \sigma_d^2 = 0}$  against the alternative  $\overbrace{\text{Lower bound of } \sigma_d^2 > 0}^{\text{Lower bound of } \sigma_d^2 > 0}$ . The results indicate

---

<sup>32</sup>Abadie and Imbens (2008) demonstrate that the bootstrap method generates biased estimates of the standard errors for a nearest-neighbor matching estimator and suggest the subsampling method developed by Politis and Romano (1994). We therefore use the subsampling method whereby we draw fewer observations than the same size at each iteration without replacement.

that the null hypothesis is rejected at least at the 5 percent level for all cases, indicating that a professor’s quality has a measurable effect on the research performance growth of the student to whom he or she is assigned.

For the results that we have presented thus far, we base the student research outcome on the research proficiency scores that are adjusted for the share of credit of each author. Alternatively, we can quantify the research outcome of a student *without* credit share adjustment. To this end, we count the number of *first-authored* articles that the student published as a lead author in the selected top general and field journals in physics. While the alternative research outcome measure might be crude and subject to a certain amount of noise — it might underrate the research achievement of a student because it ignores the articles for which he or she is not a lead author, or it might overrate the student’s attainment because it accords him or her all of the credit, even for multi-authored articles, irrespective of how many coauthors are involved — it nonetheless serves as a simple and easily interpreted yardstick.

The estimation results using the alternative research outcome measure are presented in columns (4) to (6) of Table 2. The regression estimates are larger than previous results that adjusted the author’s credit share. This is unsurprising because the first-author-based measure is greater than the original measure to the extent that the credit share is not weighted.<sup>33</sup> The estimated values of the lower bound of  $\sigma_a^2$ , reported in row (3), are correspondingly larger than those previously reported. Reassuringly, the null hypothesis that the variance in advisor quality is zero cannot be rejected at least at the 5 percent level. We therefore obtain qualitatively similar evidence on the professor’s value added as previously.

The results presented above indicate the effectiveness of professors in improving doctoral students’ research productivity growth. Indeed, better advisor quality causally affects advisees’ research achievement gains in graduate school. If we use 0.0489 as the most conservative estimate of the lower bound of the advisor quality variance among those reported in columns (1) to (3) of Table 2, we find that a one-standard-deviation increase in professor quality raises the average student research achievement gain at the doctoral level by at least 0.221, which corresponds to approximately 0.54 standard deviations of the total doctoral program research outcome distribution.

If we base the estimation results on the first-author-based research outcome measure reported in columns (4) to (6) of Table 2, we find that, if professor quality increases by one standard deviation, the average student publishes 0.64 more first-authored articles in top journals at the doctoral level.<sup>34</sup> We are thus able to conclude that professor’s value added to graduate student research outcomes is substantial.

---

<sup>33</sup>The mean and standard deviation of the first-author-based research outcome at the doctoral level are 0.39 and 0.96.

<sup>34</sup>When computing the standard deviation increase, we use 0.410 as the estimated value of the lower bound of advisor quality variance.

Our estimates of value added provide an interesting comparison with the professor value-added estimates at the undergraduate level reported by previous studies. For example, Hoffmann and Oreopoulos (2009) estimate professor value added to student’s achievement gains, measured by undergraduate course grades in a large Canadian university. They report that a one-standard-deviation increase in professor quality yields an approximately 0.05 standard deviation increase in a student’s grade. Carrell and West (2010) obtain a similar value-added estimate for professors at the U.S. Air Force Academy who teach introductory courses at the undergraduate level. They report that the standard deviation of value added is approximately 0.05. Therefore, our estimates of professor-value added at the postgraduate level are substantially larger than those standard-deviation estimates at the undergraduate level.

The observed difference in the estimates might not be too surprising considering factors that make our study distinct from other studies. First, the professor quality that we measure is different. We evaluate the dimension of professor quality that promotes a student’s *research* capability, whereas those previous studies assess the aspect of quality that enhances a student’s *academic* capability. Second, closely related to the first point, the student outcome is different. We focus on the research achievement gains of postgraduate students, while previous studies investigate the academic achievement gains of undergraduate students.

## Robustness against Various Specifications

We perform a series of checks on the estimation results regarding different specifications of student research outcomes. Firstly, we consider alternative configurations in terms of the period over which the research proficiency scores are aggregated for each degree program. Specifically, in addition to the benchmark case (M1-D2 for the master’s program and D1-P4 for the doctoral program), we examine alternative cases that change the aggregation period at the master’s and doctoral levels. Secondly, we examine whether the results are driven by a specific value of the threshold that is used to compute students’ research proficiency scores. As explained above, we consider research articles that are actually published by a target student if the author’s name matches the student’s name and, in addition, the degree of word overlap in the titles between the article and the student’s thesis exceeds some predetermined threshold value. While the default value is set to minimize both type 1 and type 2 errors, we employ both over-matching and under-matching criteria in the robustness exercise. Finally, we check the sensitivity of our estimates to the particular choice of top journals, we replicate the baseline analysis by narrowing the coverage to nine journals (two general-interest science and seven field journals) instead of twelve journals.<sup>35</sup>

---

<sup>35</sup>The three of the original twelve journals excluded here are *PNAS* in the general-interest science journal category and *Physics Letters A and B* in the field journal category. This is based on suggestions that we received from several physics researchers.



The estimation results are presented in Table E.1 and Table E.2 of Appendix E.2. It is found that all of the results are qualitatively similar to the previously reported findings. The null hypothesis that the variance in doctoral-level advisor quality is zero is rejected at the 10 percent level in all cases that change the aggregation period for the default value of the title matching threshold, and in most cases for both over-matching and under-matching criteria.

Considering the results, we can conclude that the findings from the regression model are not merely artifacts of the particular specifications of student research outcomes, and endorse the conclusion that professor quality plays a distinct role in enhancing a student’s research capacity in the doctoral program.

## 6.2 Robustness Tests

This section provides various robustness checks for the benchmark results. First, we implement a falsification test that investigates whether a false instance of an advisor switch predicts an increase in the volatility of student research outcomes between programs and cohorts. Second, we discuss the possibility that the lower bound of the estimate of the advisor quality variance might be overestimated. Third, we examine the robustness of the estimation results under a weaker conditional independence assumption on advisor switch due to turnover.

### Falsification Test

In our estimation framework, the variance in advisor quality is identified by an increase in the squared difference of the student research outcome gain at the time of advisor turnover. We thus implement a falsification exercise that examines whether the timing agrees with what is predicted by the empirical model.

To do so, we construct a *false* advisor switch dummy variable,  $\tilde{W}_n$ , that takes value one for the lab in one cohort before the actual incident and zero otherwise. Specifically, given lab  $\ell(a, c)$ , where advisor turnover occurred, the variable  $\tilde{W}_n$  is one in the *latest* cohort,  $c'$ , in which advisor  $a$  supervised at least one student before cohort  $c$ . We estimate a regression similar to regression model (7) using the dummy variable  $\tilde{W}_n$  as the regressor instead of using the true advisor switch dummy variable,  $W_n$ , with  $\tilde{\beta}$  being the coefficient of the variable  $\tilde{W}$ .

We report the estimated value of  $\tilde{\beta}$  in the panel (A) of Table 3, where we adopt the same definition of the student research outcome measures as the baseline case, and replicate the regression results except that we use the false advisor switch dummy variable.<sup>36</sup> It is shown that the false advisor switch dummy variable is sometimes negative and has no systematic impact on the the squared difference of the student research outcome gain. Indeed, in all cases

---

<sup>36</sup>The full estimation results are reported in Table E.3 of Appendix E.2.

except one, the false advisor switch dummy variable is not statistically significant, suggesting that the results survive the falsification test.

### Factors that Leads to Upward Bias

As we are interested in estimating a lower bound of the variance in advisor quality, downward bias would not be problematic. There is, however, a set of potential sources of upward bias.

The first possibility that might introduce upward bias concerns the allocation of the research *credit share* between advisor and advisee. One would consider that students are merely given a part of a larger research project, or subtopic, that the advisor has pursued, and thus, their contribution to the project in collaboration with their advisors is marginal.<sup>37</sup> If this is true, our turnover estimator for the lower bound of the variance in advisor quality might suffer from systematic upward bias, as we would then mistakenly ascribe the advisor’s research contribution to the student’s research achievement.

Because the actual collaboration process is not observed for joint research activities, it is impossible for us to allocate the true share of credit to each member of an advisor-advisee pair that engaged in a joint research project. We therefore consider an extreme case in which the student’s contribution is *zero* whenever he or she collaborated with a research advisor to highlight the sensitivity of the previous estimation results to the assumption on the allocation of research credit.

The panel (B) of Table 3 presents the estimated lower bound of advisor quality variance, assuming that the research proficiency score of student publication is equal to zero if it is coauthored with the advisor.<sup>38</sup> Looking across the columns of the table, the size of the estimated coefficients and the lower bound of advisor quality variance tend to be lower. Nonetheless, the one-sided test of the null hypothesis that doctoral-level advisor quality has no effect on an advisee’s research achievement growth is rejected at the 10 percent level. Because we consider a severe restriction on the allocation of the credit share to the side of advisees, which is overly severe for the advisees in terms of their research contributions, the reported evidence of positive professor value added reassuringly supports the conclusion that professors enhance their students’ research achievement gains by advising and mentoring their research projects at the postgraduate level.

Another potential sources of upward bias is that the assumption on the time-invariance

---

<sup>37</sup>Note that our empirical study relies on the assumption that the student made an original and substantial contribution to his or her thesis research projects and that the articles with titles that are closely associated with the master’s thesis or doctoral dissertation can be used as an unbiased yardstick to gauge the student’s in-school research achievement. The assumption appears somewhat reasonable for physics departments in Japanese universities, where, as described in Section 2.2, graduate students are typically accorded a fair amount of autonomy when choosing a research topic and approach.

<sup>38</sup>The full estimation results are reported in Table E.4 of Appendix E.2.

of advisor quality, given by assumption 1.1, might be violated. Suppose, contrary to the assumption, that it varies across cohorts *within* a professor. In particular, if it fluctuates as the end of a professor’s research career approaches, the squared difference measure, the dependent variable in regression model (7), becomes more volatile in the last cohorts before a professor’s turnover. In this case, the regression coefficient of the advisor switch dummy variable might overstate the lower bound of advisor quality variance.

To shed some light on this concern, we augment the regression model in Equation (7) by including a set of dummy variables that capture the possible change in advisor quality variance in the period near turnover.<sup>39</sup> The panel (C) of Table 3 presents the estimated lower bound of the variance in advisor quality from the augmented specification.<sup>40</sup> As shown, the estimates are not substantially affected by the inclusion of the cohort-specific dummy variables. Encouragingly, the null hypothesis that advisor quality in the doctoral program has no effect on student research outcome growth is rejected at least at the 10 percent level in all cases.<sup>41</sup>

### Insert Table 3

#### Non-Retirement Turnover

Our identification strategy for a lower limit of advisor quality relies on the assumption that advisor switch is *incidental*, that is, it is independent on advisee’s unobservables conditional on advisor’s observable characteristics. The assumption leads us to consider the baseline regression model where the advisor switch indicator,  $W_n$ , is assumed to be orthogonal to the student specific catch-all error term  $\varepsilon_n$ , once the observables of professors are controlled by the propensity score method.

There might be a concern about the validity on the conditional randomness of advisor switch, however. Indeed, advisor switch (and non-switch) will be *intentional* for students. It might be possible for students, having prior information on faculty members, to predict the “scheduled” retirement events, and choose their advisors and labs. If this kind of self-selection occurs, students in the lab with turnover might be systematically heterogeneous from those in the lab without turnover. This may cause the advisor switch indicator to be correlated with error term in the regression model through confounding factors related to student’s research performance. The problem of unobserved heterogeneity in the regression cannot be eliminated by the propensity score method since it takes into account only selection on observables on

---

<sup>39</sup>Specifically, we consider a regression that includes dummy variable  $D_k^{(a,c,c')}$  for  $k = 1, 2$ , and 3. The dichotomous variable takes value one if cohort  $c$  is within  $k$  years before professor  $a$  exited and zero otherwise.

<sup>40</sup>The full estimation results are reported in Table E.5 of Appendix E.2.

<sup>41</sup>Furthermore, the estimation results shown in Table E.5 indicate that none of the coefficients concerning the added dummy variables are statistically significant.

the side of advisors.

To mitigate the concern, we examine how the lower bound value estimates of the variance in advisor quality are robust under a more plausible assumption on the conditional independence of advisor turnover. We perform regression analyses under an assumption that turnover due to *non-retirement* reasons is accidental and thus the event of advisor switch caused by the turnover is incidentally orthogonal to student unobserved characteristics.<sup>42</sup> It may be true that students can predict the timing of the scheduled mandatory retirement a faculty member with some degree of accuracy.<sup>43</sup> Nonetheless, it does not seem unreasonable to assume that students do not have sufficient information to foresee the future relocation, not to mention decease, of a faculty member *ax ante* at the timing of graduate enrollment. Given the limited ability of students to predict the timing of faculty’s move and death, it seems less controversial to assume that advisor switch caused by the non-retirement based turnover is more incidental than the one caused by retirement based turnover, and thus is likely to be uncorrelated to student specific unobserved factors.

Table 4 reports the estimated lower bounds for the variance of the advisor quality when the incidents of turnover are restricted to those due to non-retirement reasons.<sup>44</sup> The panel (A) presents the estimates when we replicates the baseline specification, while the panel (B) presents the estimates for the case where the research proficiency score of student publication is set to zero if it is a joint work with the advisor.

Overall we find that the estimates are positive and statistically significant, which bolsters the previous finding that advisor quality matters for advisee’s research performance growth. Indeed, the estimated lower bounds for the variance of the advisor quality are larger than those found in the corresponding baseline case. Yet, considering the results only pertain to advisor switches that occurred due to nonretirement reasons, we do not use them to update the information on the magnitudes of the variance for advisor’s quality. It seems nonetheless safe to say that the previous finding on the effectiveness of advisor is not driven by self selection behavior of advisees.

#### Insert Table 4

---

<sup>42</sup>Since the available information on students background characteristics is rather limited comparing to those on professors, we do not opt for including as many characteristics as possible to absorb all the potential student-side heterogeneity across labs.

<sup>43</sup>It looks that all the retirement events cannot be perfectly predictable, however. As explained in footnote 22, the mandatory retirement policy changed at UTokyo in the year 2001. The policy was made an official decision in the year 2000, but there had been considerable uncertainty as to whether the policy was enacted or not at the time. So, it seems highly likely that faculty members, not to mention graduate students, who were at UTokyo’s physics program in the late 90s can predict the change of retirement age that happened in the early 00s.

<sup>44</sup>The full estimation results are reported in Table E.6 and Table E.7 of Appendix E.2.

## 6.3 Additional Evidence for Professors' Influence on Students

### Other Mechanisms

The estimation results have shown that advisor turnover generates significant variations in an advisee's research achievement gains at UTokyo's department of physics. According to a standard value-added model, we ascribe the increased diversity of student research achievement gains to the discrete change in advisor quality at the time of turnover. Admittedly, however, there may remain other mechanisms that create such a pattern.

One possibility is that professor turnover always has a *positive* effect on students' research capacity and thus increases the variability in student achievement gains between cohorts with and without turnover. The positive advisor turnover effect could be caused by a mechanism that reflects a well-known understanding that innovation (and thus economic growth) is due to the recombination of existing ideas (e.g., Weitzman, 1998). It follows from this view that, as new innovation is likely to arise from recombining old knowledge elements, students who are supervised by different professors would have access to a wider variety of knowledge and ideas and can thus enhance their research capabilities.

Another mechanism is the one that yields a *negative* effect of professor turnover on students' research achievement gains. Turnover may lead to a disruption that impacts the students who are forced to change their research advisors. As is often noted in the education literature (e.g., Wisker and Robinson, 2013), if an advisor is lost due to turnover, an advisee who becomes an "orphan" occasionally perceives this as a traumatic event and suffers from psychological problems that might occasionally result in under-development of academic achievement. If this understanding is correct, advisor turnover would retard the advisee's research progress, irrespective of how high the quality of the newly assigned advisor is, and thus generate a noticeable gap in student research outcome gains between cohorts with and without turnover.

Recall that, according to the mechanism captured by the value-added model, the advisee's research outcome growth can be positive or negative after turnover – indeed, as explained in Section 3.1, the direction of growth depends decisively on the relative levels of advisor quality that were switched when turnover occurred and will thus not be predicted *a priori* unless the information on the exact quality levels is available. In the analysis that follows, we investigate which mechanism is more likely by estimating a regression similar to the regression model in Equation (7), except with the dependent variable being *in levels*,  $(DD\overline{\Delta outcome}_n)$ , not *in squares*  $(DD\overline{\Delta outcome}_n)^2$ . To identify the mechanism in place, we focus on the sign of the turnover effect on the advisor's research achievement gain.

The panel (A) of Table 5 presents the regression results for which all estimated coefficients of the advisor switch indicator are shown to be positive but are not statistically significant

in all cases.<sup>45</sup> We can interpret the results as indicating that, contrary to the predictions of the alternative mechanisms, advisor turnover can have a positive or negative impact on an advisee’s research productivity growth. As the individual impacts cancel one another out, the aggregate effect, as reflected by the integration, is not significantly different from the null in levels. It thus appears to confirm that the mechanism that the value-added model postulates should be a main driver of the empirical findings obtained thus far.

## Indirect Influence

Our analysis thus far has concentrated on the advisor-advisee relationship within a lab and intended to measure the effectiveness of knowledge transmission through a direct interaction channel within a lab. However, knowledge might be transmitted beyond master-apprenticeship-style contact.

We modify the baseline specification in Equation (2) by incorporating an “indirect” effect of professor. Consider the following augmented model of student research outcome gains:

$$\overline{\Delta outcome}_{eg}^{\ell(e,c)} = \bar{\gamma}^{\ell(e,c)} + \theta_{eg} + \sum_{f \in \mathcal{A}} \pi_{ef} \theta_{fg} + \bar{v}_{eg}^{\ell(e,c)}, \quad (10)$$

where we consider lab  $\ell(e, c)$  of professor  $e \in \mathcal{A}$  in cohort  $c \in \mathcal{C}$ . The parameter  $\pi_{ef}$  captures the magnitude of the indirect influence from non-advisor faculty member  $f$  on the average research outcome gain in program  $g$  for students in lab  $\ell(e, c)$ . In what follows, for the purpose of simplicity, we assume that  $\pi_{ef} = \pi$  if the research field or subfield of professor  $f$  is the same as or closely related to that of the direct advisor,  $e$ , and  $\pi_{ef} = 0$  otherwise.<sup>46</sup>

Analogous to Equation (6), we compute the conditional expectation of the squared double-differenced average student research outcome growth and construct a regression model based on the comparison of the conditional expectations between labs with and without a “treatment” assignment. To achieve this aim, let us use  $V^{\ell(e,c,c')}$  to denote an assignment indicator of an “indirect” turnover incident. Specifically, define  $V^{\ell(e,c,c')} = 1$  if a professor whose research subfield is the same as that of professor  $e$  is replaced due to turnover in cohort  $c$  and  $V^{\ell(e,c,c')} = 0$  otherwise.

We obtain the following result under the same assumptions as above on the distributions

---

<sup>45</sup>The full estimation results are reported in Table E.8 of Appendix E.2.

<sup>46</sup>In other words, we restrict the scope of indirect influence to that between professors and students within the *same* research field (or subfields).

of advisor quality and the idiosyncratic error terms:

$$\begin{aligned} & \mathbb{E} \left[ \left( \overline{DD\Delta outcome}^{\ell(e,c,c')} \right)^2 \mid W^{\ell(e,c,c')} = 0, V^{\ell(e,c,c')} \right] \\ &= \alpha \left( \frac{1}{\overline{I}^{\ell(e,c)}} + \frac{1}{\overline{I}^{\ell(e,c')}} \right) + \pi^2 \{2\sigma_d^2(1 - \rho_d)\} V^{\ell(e,c,c')}, \end{aligned} \quad (11)$$

where  $\alpha$  is the same as that given in Equation (6). This, in turn, leads to the following regression model using the subsamples that consist of labs in which advisor turnover did *not* occur ( $W^{\ell(e,c,c')} = 0$ ):

$$\left( \overline{DD\Delta outcome}_m \right)^2 = \alpha_{ind} X_m + \beta_{ind} V_m + \varepsilon_m, \quad (12)$$

where  $m = 1, \dots, M$  is the index of observations.<sup>47</sup>

Comparing Equations (11) and (12) leads to the parameter relationship that  $\beta_{ind} = \pi^2 \{2\sigma_d^2(1 - \rho_d)\}$ . Let us use  $\hat{\beta}_{dir}$  to denote an estimate of the coefficient of  $W_n$  from the baseline regression model given by Equation (7), and let  $\hat{\beta}_{ind}$  be an estimate of the coefficient of  $V_m$  from Equation (12) presented above. We therefore obtain  $\hat{\pi} = \sqrt{\hat{\beta}_{ind}/\hat{\beta}_{dir}}$ , which can be used as a measure of indirect knowledge transfer from a non-advisor professor in the same research field.

An empirical challenge is to identify a group of professors whose research subjects were close enough to that of the professor experiencing turnover. As the type of data necessary to judge the similarity between research subjects is absent or rarely present, we adopt a simple and heuristic method for identifying the same research subject groups, which exploits the information revealed by the actual turnover events. It is conceivable that, when an instance of professor turnover occurred, the students in the lab of the professor who exited were highly likely to be re-assigned to a professor whose research area was closely related to that of the original professor. In the empirical analysis that follows, we therefore assume that the original professor who exited and the re-assigned professor were working in the same research area. The detailed situation is illustrated in Appendix E.1.

The results shown in panel (B) of Table 5 presents the regression estimates.<sup>48</sup> We adopt the default setting for the student research outcome and use the same estimation method as before.<sup>49</sup> As shown, the estimates are ambiguous for  $\beta_{ind}$ . One of the estimates is negative, and in the case in which the estimates are positive, they are not statistically significant at

<sup>47</sup>Here, the unit of observation is each element of  $(e, c, c')$  such that  $W^{\ell(e,c,c')} = 0$  for any advisor  $e \in \mathcal{A}$  and cohorts  $c, c'$  such that  $0 < c - c' \leq \tau$ .

<sup>48</sup>The full estimation results are reported in Table E.9 of Appendix E.2.

<sup>49</sup>The research outcomes in the master's degree and doctoral degree programs are aggregated over the period from M1 to D2 and the period from D1 to P4, respectively. Furthermore, the set of "top journals" here consists of twelve journals.

the 10 percent level in any specification except one. The estimates of the squared indirect influence parameter,  $\pi^2$ , are reported in row (3).<sup>50</sup> Again, the signs of the estimates differ. The maximum estimate is 0.269, while the value where  $\hat{\beta}_{ind}$  is statistically significant is as low as approximately 0.1, as reported in column (6). This result implies that the indirect knowledge transfer effect from non-advisor faculty is  $\hat{\pi} = 0.33$ , suggesting that it is, at most, less than one-third of the direct effect from the advisor.

On balance, therefore, there appears to be little or no indirect influence from non-advisor faculty members across labs on doctoral student research productivity growth.

## 7 Conclusion

In this paper, we investigated the extent to which professors can affect the development of the research performance of the graduate students whom they supervise. By using detailed data on professors and students at UTokyo’s department of physics, we estimated a lower bound of the professor value added to student research achievement growth while in school. The estimation results consistently show that postgraduate research education based on an advisor-advisee relationship is quite effective — professors have a substantial impact on the students’ achievement gains in terms of the number of publications in top journals in physics. This corroborates the view of earlier studies (e.g., Azoulay et al., 2010; Moser et al., 2014; Borjas and Doran, 2014) that research interactions among scientists in vertically aligned relationships, including senior-and-junior-collaborator, teacher-student, and advisor-advisee relationships, matter for the creation and diffusion of scientific ideas and knowledge.

Our findings also suggest that the accumulation of prominent scientists in a comparatively small number of universities is explained, at least partially, by the results of successful education at the postgraduate level. For example, in Japan, five out of ten Nobel Prize winners in physics completed their doctoral degrees at UTokyo, and four earned their doctorate degrees at Nagoya University. Given our results on the effectiveness of professors in enhancing students’ research capability growth, we can speculate that the relatively high concentration of physics Nobel laureates in these two universities in Japan might be caused not only by the processes of students’ self-selection or schools’ selective recruitment but also by the beneficial reproduction of elite physicists, which was enabled by a deliberate process of teaching and learning in a lab. While previous studies (e.g., Waldinger, 2010) suggest that high-quality universities can facilitate human capital accumulation among graduate students, our paper specifically adds that this outcome is based on advisor-advisee-based education.

We need to highlight some limitations of this paper. First, our analysis of the professor’s

---

<sup>50</sup>We compute  $\hat{\pi}^2$  as the quotient of the estimate  $\hat{\beta}_{ind}$  over the estimate  $\hat{\beta}_{dir}$  that appears in the corresponding number of columns in Table 2.



value added is essentially short run. Although the estimation results reveal that research advisors can influence the research development of their students, the impact might be limited to the short span of time while the student is in graduate school or several years after the completion of graduate school. It is left to future research to examine whether a professor's supervision during a graduate program has a long-term impact on student research performance during their postgraduation careers.

Second, the analysis in the paper is limited to a small, albeit prominent, group of physicists. Thus, our conclusion regarding a professor's value added might not be generalizable to groups of other scientists from different disciplines or other graduate schools. We hope that the findings of this paper regarding the efficacy of professors in promoting student progress in research performance will be helpful to stimulate further research in related areas including the economics of higher education and the economics of science and technology.

## References

- Abadie, Alberto and Guido W Imbens**, “Notes and Comments on the Failure of the Bootstrap,” *Econometrica*, 2008, 76 (6), 1537–1557.
- Abe, Yasumi and Satoshi P Watanabe**, “Some Thoughts on Implementing US Physics Doctoral Education in Japanese Universities,” *Asia Pacific Education Review*, 2012, 13 (3), 403–415.
- Azoulay, Pierre, Joshua S. Graff Zivin, and Jialan Wang**, “Superstar Extinction,” *Quarterly Journal of Economics*, 2010, 125 (2), 549–589.
- Borjas, George J. and Kirk B. Doran**, “The Collapse of the Soviet Union and the Productivity of American Mathematicians,” *Quarterly Journal of Economics*, 2012, 127 (3), 1143–1203.
- and –, “Which Peers Matter? The Relative Impacts of Collaborators, Colleagues, and Competitors,” *Review of Economics and Statistics*, 2014. *forthcoming*.
- Branch, Gregory F, Eric A Hanushek, and Steven G Rivkin**, “Estimating the Effect of Leaders on Public Sector Productivity: The Case of School Principals,” Technical Report, National Bureau of Economic Research 2012. NBER Working Paper #17803.
- Carrell, Scott E. and James E. West**, “Does Professor Quality Matter? Evidence from Random Assignment of Students to Professors,” *Journal of Political Economy*, 2010, 118 (3), 409–432.
- Chetty, Raj, John N Friedman, and Jonah E Rockoff**, “Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-added Estimates,” *American Economic Review*, 2014, 104 (9), 2593–2632.
- Crump, Richard K., Joseph V. Hotz, Guido W. Imbens, and Oscar A. Mitnik**, “Dealing with Limited Overlap in Estimation of Average Treatment Effects,” *Biometrika*, 2009, 96 (1), 187–199.
- Grossman, Gene M. and Elhanan Helpman**, “Quality Ladders in the Theory of Growth,” *Review of Economic Studies*, 1991, 58 (1), 43–61.
- Gumport, Patricia J.**, “Graduate Education and Organized Research in the United States,” in “The Research Foundations of Graduate Education: Germany, Britain, France, United States, Japan,” University of California Press, 1993, pp. 225–260.

- Hanushek, Eric A. and Steven G. Rivkin**, “Teacher Quality,” in “Handbook of the Economics of Education,” Vol. 2, edited by Eric A. Hanushek and Finis Welch, Elsevier, 2006, chapter 18, pp. 1051–1078.
- Hanushek, Eric A and Steven G Rivkin**, “Generalizations about Using Value-added Measures of Teacher Quality,” *The American Economic Review*, 2010, *100* (2), 267–271.
- Hanushek, Eric A. and Steven G. Rivkin**, “The Distribution of Teacher Quality and Implications for Policy,” *Annual Review of Economics*, 2012, *4*, 131–157.
- Hilmer, Michael J. and Christiana E. Hilmer**, “Fishes, Ponds, and Productivity: Student-Advisor Matching and Early Career Publishing Success for Economics Phds,” *Economic Inquiry*, 2009, *47* (2), 290–303.
- Hoffmann, Florian and Philip Oreopoulos**, “Professor Qualities and Student Achievement,” *Review of Economics and Statistics*, 2009, *91* (1), 83–92.
- Imbens, Guido W.**, “Matching methods in practice: Three examples,” Technical Report, National Bureau of Economic Research 2014. NBER Working Paper #19959.
- Imbens, Guido W and Donald B Rubin**, *Causal Inference in Statistics, Social, and Biomedical Sciences*, Cambridge University Press, 2015.
- Jones, Benjamin F and Benjamin A Olken**, “Do Leaders Matter? National Leadership and Growth since World War II,” *Quarterly Journal of Economics*, 2005, *120* (3), 835–864.
- Kawashima, Tatsuo and Fumihiko Maruyama**, “The education of advanced students in Japan: engineering, physics, economics and history,” in “The Research Foundations of Graduate Education: Germany, Britain, France, United States, Japan,” University of California Press, 1993, pp. 326–353.
- Lacetera, Nicola, Bradley J Larsen, Devin G Pope, and Justin R Sydnor**, “Bid Takers or Market Makers? The Effect of Auctioneers on Auction Outcome,” *American Economic Journal: Microeconomics*, 2016, *8* (4), 195–229.
- Lazear, Edward P, Kathryn L Shaw, and Christopher T Stanton**, “The value of Bosses,” *Journal of Labor Economics*, 2015, *33* (4), 823–861.
- Liu, Xuan Zhen and Hui Fang**, “Fairly Sharing the Credit of Multi-authored Papers and its Application in the Modification of h-index and g-index,” *Scientometrics*, 2012, *91* (1), 37–49.
- Low, Morris**, *Science and the Building of a New Japan*, Palgrave Macmillan, 2005.

- Lucas, Robert E.**, “On the Mechanics of Economic Development,” *Journal of Monetary Economics*, 1988, *22*, 3–42.
- Malmendier, Ulrike and Geoffrey Tate**, “Superstar CEOs,” *Quarterly Journal of Economics*, 2009, *124* (4), 1593–1638.
- Moser, Petra, Alessandra Voena, and Fabian Waldinger**, “German-Jewish Émigrés and US Invention,” *American Economic Review*, 2014, *104* (10), 3222–3255.
- Ogawa, Yoshikazu**, “Challenging the Traditional Organization of Japanese Universities,” *Higher Education*, 2002, *43* (1), 85–108.
- Polanyi, Michael**, *Personal Knowledge: Towards a Post-Critical Philosophy*, University of Chicago Press, 1958.
- , *The Tacit Dimension*, Routledge, 1966.
- Politis, Dimitris N and Joseph P Romano**, “Large Sample Confidence Regions Based on Subsamples under Minimal Assumptions,” *The Annals of Statistics*, 1994, *22* (4), 2031–2050.
- Rivkin, Steven G., Eric A. Hanushek, and John F. Kain**, “Teachers, Schools, and Academic Achievement,” *Econometrica*, 2005, *73* (2), 417–458.
- Romer, Paul M.**, “Endogenous Technological Change,” *Journal of Political Economy*, 1990, *98* (5), S71–S102.
- Rosenbaum, Paul R and Donald B Rubin**, “Constructing a Control Group Using Multivariate Matched Sampling Methods that Incorporate the Propensity Score,” *The American Statistician*, 1985, *39* (1), 33–38.
- Tanabashi, Masaharu**, “Shoichi Sakata: His Life and Physics, Collections of Materials in Sakata Memorial Archival Library,” *Progress of Theoretical Physics Supplement*, 2012, *197*, 1–14.
- Tang, Li and John P. Walsh**, “Bibliometric Fingerprints: Name Disambiguation Based on Approximate Structure Equivalence of Cognitive Maps,” *Scientometrics*, 2010, *84* (3), 763–784.
- Trajtenberg, Manuel, Gil Shiff, and Ran Melamed**, “The ”Names Game”: Harnessing inventors’ Patent Data for Economic Research,” Technical Report, National Bureau of Economic Research 2006. NBER Working Paper #12479.

- Traweek, Sharon**, *Beamtimes and Lifetimes: The World of High Energy Physicists*, Harvard University Press, 1988.
- Ushioji, Morikazu**, “Graduate education and research organization in Japan,” in “The Research Foundations of Graduate Education: Germany, Britain, France, United States, Japan,” University of California Press, 1993, pp. 299–325.
- Waldinger, Fabian**, “Quality Matters: The Expulsion of Professors and the Consequences for PhD Student Outcomes in Nazi Germany,” *Journal of Political Economy*, 2010, 118 (4), 787–831.
- , “Peer Effects in Science: Evidence from the Dismissal of Scientists in Nazi Germany,” *Review of Economic Studies*, 2012, 79 (2), 838–861.
- Waltman, Ludo**, “An Empirical Analysis of the Use of Alphabetical Authorship in Scientific Publishing,” *Journal of Informetrics*, 2012, 6 (4), 700–711.
- Weitzman, Martin L**, “Recombinant Growth,” *Quarterly Journal of Economics*, 1998, 113 (2), 331–360.
- Wisker, Gina and Gillian Robinson**, “Doctoral “orphans”: Nurturing and Supporting the Success of Postgraduates Who Have Lost Their Supervisors,” *Higher Education Research & Development*, 2013, 32 (2), 300–313.
- Zuckerman, Harriet**, *Scientific elite: Nobel laureates in the United States*, Transaction Publishers, 1977.

# Figures

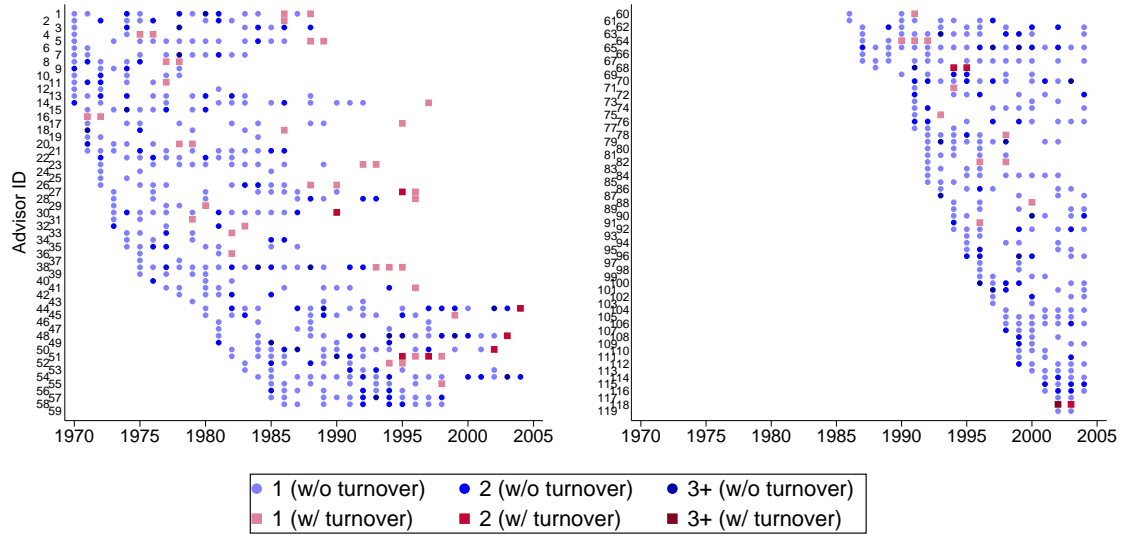


Figure 1: Distribution of Students across Initial Advisors and Cohorts

*Notes:* Initial advisors are set in the vertical axis and cohorts are set in the horizontal axis. The red square and blue circle markers represent a lab where the advisor was replaced due to turnover and a lab where the advisor was not replaced, respectively. The darker color marker means more students in a lab.

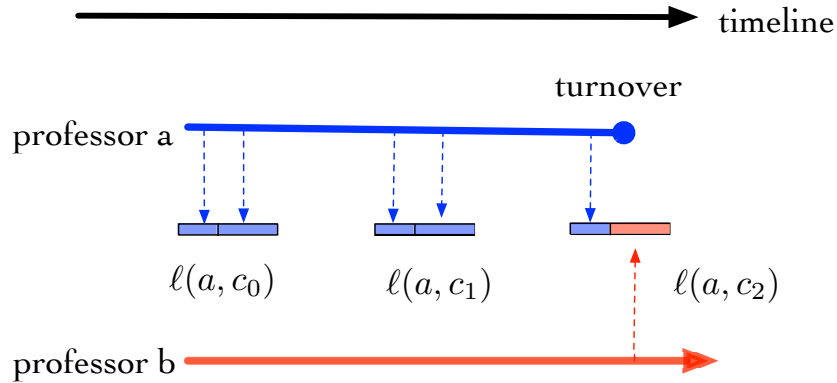


Figure 2: Example of Labs with and without Turnover

*Notes:* There are three labs with different cohorts,  $c_0$ ,  $c_1$  and  $c_2$ , whose initial advisor is professor  $a$ . Each lab is portrayed by a connected line segment, which represents the two-year master's degree program and the three-year doctoral degree program. Advisor turnover did not occur in labs  $\ell(a, c_0)$  or  $\ell(a, c_1)$  before cohort  $c_2$ . On the other hand, in lab  $\ell(a, c_2)$ , professor  $a$  exited the school due to turnover, and professor  $b$  took charge of the doctoral students.

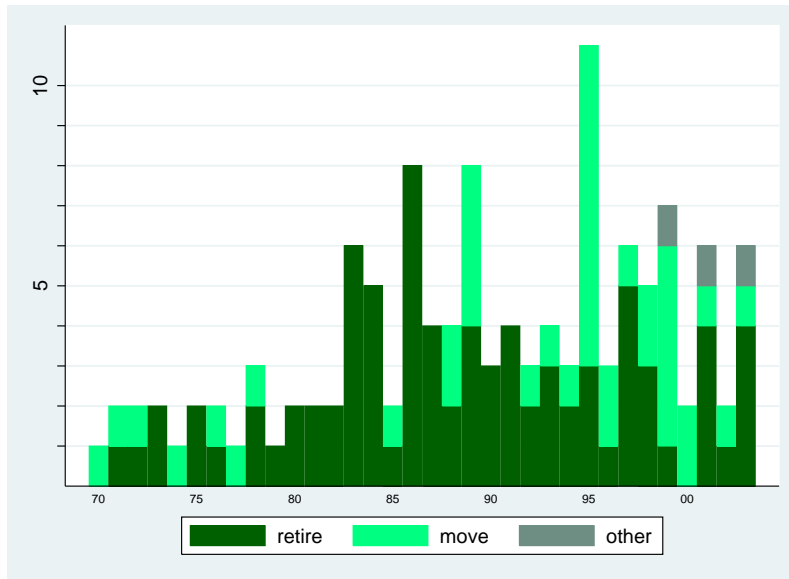


Figure 3: Number of Turnover Incidents in Each Year (1970-2004)

*Notes:* The reasons for turnover are classified into the following three categories: (1) retirement if the instance of turnover occurred at the mandatory retirement age; (2) move if turnover occurred before the retirement age and the faculty name began to reappear on other universities' rosters beginning in the year after the turnover instance; and (3) decease/quit otherwise.

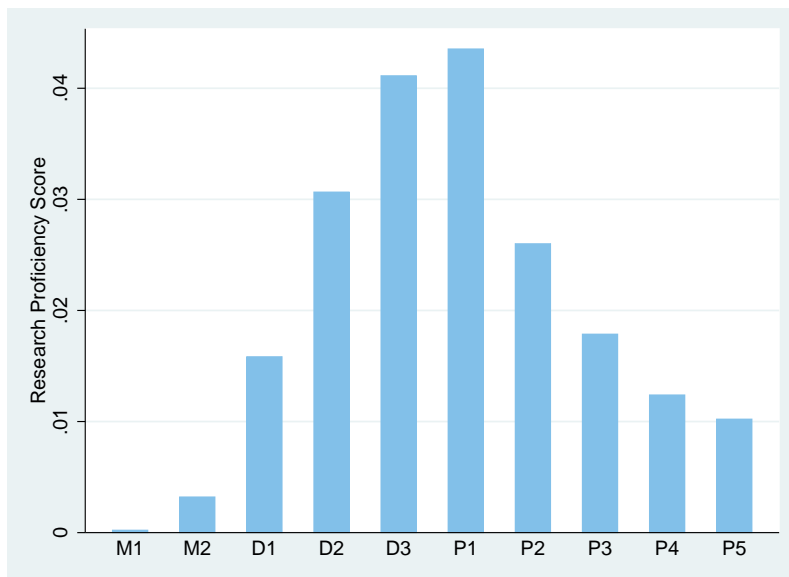


Figure 4: Average Student Research Proficiency Scores

*Notes:* The research proficiency score is defined as the number of publication counts of a student's publication records during a given year. Two quality adjustment methods are employed. First, we limit the publications to those published in twelve high-quality peer-reviewed journals. Second, we consider a student's share of credit for an article if there are multiple authors.

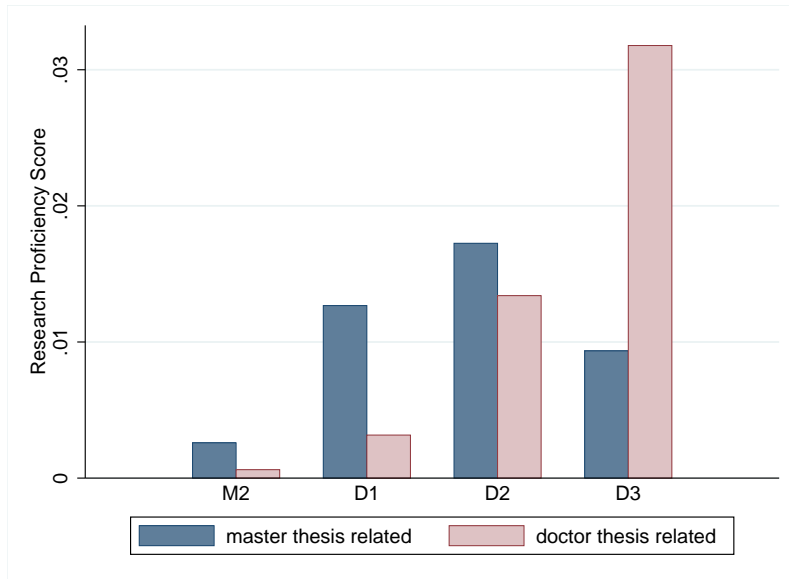


Figure 5: Decomposition of the Student Research Proficiency Scores

*Notes:* The student average research proficiency scores are decomposed into those related to the master's thesis and those related to the doctoral thesis. We determine whether research articles are actually published by a target student or not if the degree of word overlap in the titles between the article and the student's master's and doctoral theses exceeds some predetermined threshold value.

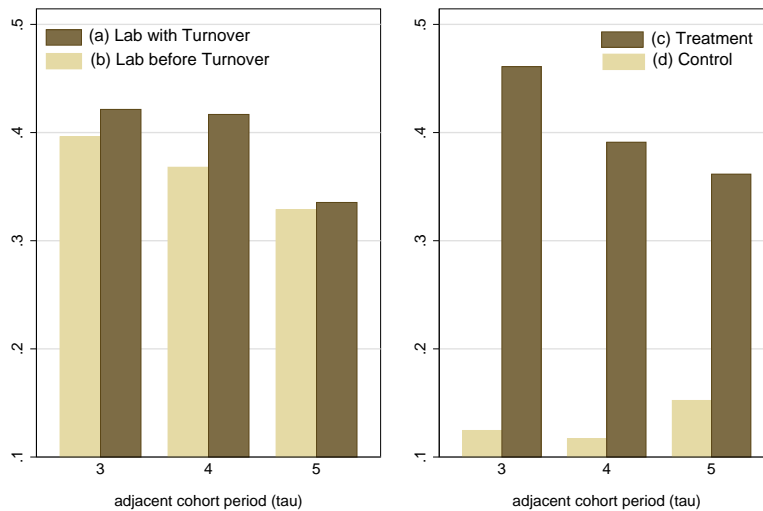


Figure 6: Comparison of the Double-differenced (DD) Average Student Research Outcome Growth between Labs with and without Turnover

*Notes:* The left figure presents the sample variances of the DD measure for (a) labs in the cohort where turnover occurred and (b) labs of the same advisor but their cohort is the latest one before turnover. The right figure provides the same comparison of the sample variances between (c) labs in the treatment group where turnover occurred and (d) labs in the corresponding control group that are matched through the propensity score method.



# Tables

Table 1: Descriptive Statistics for the Student Research Outcomes in Levels and in Differences

	Research Outcome at the Master's Level	Research Outcome at the Doctoral Level	Research Outcome Gain at the Doctoral Level
	$outcome_{iam}^c$	$outcome_{iad}^c$	$\Delta outcome_{iad}^c$
Mean	0.0677	0.2202	0.1481
S.D.	0.2184	0.5075	0.4068
Min	0.0000	0.0000	-0.4738
Max	2.3175	4.7303	4.7303
Sample Size	1019	1019	1019

*Notes:* The research outcome at each degree level is computed based on the research proficiency scores. The aggregation years are M1-D2 for the master's level and D1-P4 for the doctoral level, respectively. The research outcome gain at the doctoral level is given by the difference of the research outcome from the doctoral level to the master's level. Since the research outcome at the bachelor's level is normalized as zero, the research outcome gain at the master's level is equal to the research outcome at the master's level.

Table 2: Baseline Estimation Results: The Effect of Advisor Turnover on Student Research Outcome Growth at the Doctoral Level

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\overline{\Delta outcome}]^2$			$[DD\overline{\Delta outcome}]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\hat{\alpha}$	0.0667 *** (0.0237)	0.0742 *** (0.0225)	0.0960 *** (0.0134)	0.0570 (0.0682)	0.1267 ** (0.0545)	0.4025 *** (0.0720)
(2) $\hat{\beta}$	0.3371 * (0.1746)	0.2663 ** (0.1204)	0.1956 ** (0.0985)	2.3091 * (1.3249)	2.1322 ** (0.9220)	1.6401 ** (0.8251)
(3) Lower bound of $\sigma_d^2$	0.0843 ** [0.0268]	0.0666 ** [0.0135]	0.0489 ** [0.0236]	0.5773 ** [0.0407]	0.5331 ** [0.0104]	0.4100 ** [0.0234]
Sample Size						
Total	925	1202	1446	925	1202	1446
After matching	104	186	271	104	186	271

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The total sample size is given by the number of observations for each tuple of  $(a, c, c')$  for any advisor  $a$  in  $\mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$  where  $\tau$  is the period over which the difference is taken. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The after-matching sample size is the sum of the numbers of observations for the treatment group where turnover occurred and the corresponding control group that are matched through the propensity score method. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that Lower bound of  $\sigma_d^2 = 0$  against the alternative

Lower bound of  $\sigma_d^2 > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table 3: Robustness Test Results

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\overline{\Delta outcome}]^2$			$[DD\overline{\Delta outcome}]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(A) FALSIFICATION TEST						
$\hat{\beta}$	0.2979 (0.2123)	0.2039 (0.1725)	0.0814 (0.1270)	0.6315 ** (0.2616)	-0.0969 (0.3965)	-0.7784 (0.4350)
(B) THE STUDENT PROFICIENCY SCORE IS ZERO IF COAUTHORED WITH ADVISOR						
Lower bound of $\sigma_d^2$	0.0514 [0.1168]	0.0416 * [0.0803]	0.0593 ** [0.0101]	0.3444 * [0.0790]	0.3343 ** [0.0247]	0.3422 *** [0.0093]
(C) WHEN A CHANGE IN ADVISOR AUALITY VARIANCE IS ALLOWED						
Lower bound of $\sigma_d^2$	0.0848 ** [0.0326]	0.0608 ** [0.0301]	0.0536 ** [0.0220]	0.5830 * [0.0503]	0.5302 ** [0.0144]	0.4277 ** [0.0283]

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The estimated coefficient of the false advisor switch indicator,  $\tilde{W}_n$ , and the estimated lower bounds of the variance in advisor quality at the doctoral level are reported. The full estimation results are in Table E.1, Table E.2 and Table E.3 of Appendix E.2. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that Lower bound of  $\sigma_d^2 = 0$  against the alternative Lower bound of  $\sigma_d^2 > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table 4: Estimation Results when Non-Retirement Turnover Events Are Used: The Estimated Lower Bounds for the Variance of the Advisor Quality at the Doctoral Level

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\Delta outcome]^2$			$[DD\Delta outcome]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(A) BASELINE CASE						
Lower bound of $\sigma_d^2$	0.1863 ** [0.0496]	0.1872 *** [0.0035]	0.1240 *** [0.0087]	1.2319 * [0.0917]	1.3214 ** [0.0129]	1.0214 ** [0.0097]
(B) THE STUDENT PROFICIENCY SCORE IS ZERO IF COAUTHORED WITH ADVISOR						
Lower bound of $\sigma_d^2$	0.1482 * [0.0926]	0.1433 ** [0.0195]	0.1326 ** [0.0103]	0.9007 [0.1019]	0.8801 ** [0.0227]	0.7196 ** [0.0164]

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The estimated lower bounds of the variance in advisor quality at the doctoral level are reported. The full estimation results are in Table E.4 and Table E.5 of Appendix E.2. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that Lower bound of  $\sigma_d^2 = 0$  against the alternative Lower bound of  $\sigma_d^2 > 0$ . The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses.

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table 5: Additional Evidence for Professors' Influence on Students

Dependent Adjacent Period	Credit Share Weighted			First-authored-paper Based		
	$\tau = 3$ (1)	$\tau = 4$ (2)	$\tau = 5$ (3)	$\tau = 3$ (4)	$\tau = 4$ (5)	$\tau = 5$ (6)
(A) THE <i>DD</i> MEASURE IN LEVELS IS USED						
$\hat{\beta}$	0.1439 (0.1746)	0.1480 (0.1204)	0.0310 (0.0985)	0.2655 (2.1883)	0.3968 (1.5374)	0.2276 (1.3783)
(B) EFFECT OF NON-ADVISOR TURNOVER						
$\hat{\beta}_{ind}$	-0.0230 (0.0374)	0.0274 (0.0320)	0.0527 (0.0336)	0.0764 (0.1192)	0.1007 (0.0924)	0.1768 (0.0924)
$\hat{\pi}^2$	-0.0682	0.1030	0.2694	0.0331	0.0472	0.1078

*Notes:* The dependent variables are the double-differenced average student research outcome growth in level and in square for panel (A) and (B), respectively. The estimated coefficients of the advisor switch indicator due to turnover,  $W_n$ , and the assignment indicator of “indirect” turnover,  $V_m$ , are reported. The full estimation results are in Table E.6 and Table E.7 of Appendix E.2. We define  $\hat{\pi}^2 = \hat{\beta}_{ind}/\hat{\beta}_{dir}$  where  $\hat{\beta}_{dir}$  is taken from the estimates in Table 2. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses.

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

# Appendix

## A Laboratory Democracy

Low (2005) notes that professor Shouichi Sakata at the physics department of Nagoya Imperial University, an influential physicist at that time, played an important role in developing the new democratic lab system in the Japanese physics community. Sakata, who was under the philosophical influence of Marxism, introduced a charter for the physics department at Nagoya in 1946. The charter holds that democracy should serve as the guiding principle in department affairs; all faculty members and students should be treated equally concerning physics research (Department of Physics, Nagoya University, 2015). The idealism of Sakata's "laboratory democracy" then spread. Soon after the Nagoya Charter was announced, several physics departments at other universities introduced similar systems. See Tanabashi (2012) for details on Sakata's laboratory democracy.

Traweek (1988) reports that decision-making in Japanese physics labs was based on the consensus of the members. Traweek, an anthropologist who studied various research groups of elementary particle physicists in Japan and the U.S., offers a first-hand account of the democratic nature of labs in Japan by asking group leaders of a lab for the source of new ideas for experimental design or data analysis. She writes, (p.147) " [lab leaders] generally credited the graduate students ... they said the group then responds to their ideas, perhaps modifying or amplifying them".

## B The Assumption on the Random Shock

### B.1 Assumptions on the Moments of the Idiosyncratic Shocks

The following assumptions impose restrictions on the moments of the idiosyncratic shock after *demeaning* by each cohort.

assumption 2.1: The conditional expectation and variance of the demeaned random shock,  $\tilde{\nu}_c$ , for student  $i \in \mathcal{J}^{\ell(a,c)}$  are  $E(\tilde{\nu}_{iag}^c | W^{\ell(a,c,c')}) = 0$  and  $\text{Var}(\tilde{\nu}_{iag}^c | W^{\ell(a,c,c')}) = \phi_g^2$ , respectively, for any  $a \in \mathcal{A}$ ,  $g \in \{m, d\}$ , and  $c, c' \in \mathcal{C}$ .

assumption 2.2: The covariance of the demeaned random shocks *between* degree programs *within* the same student,  $i \in \mathcal{J}^{\ell(a,c)}$ , is given by  $\text{Cov}(\tilde{\nu}_{iam}^c, \tilde{\nu}_{iad}^c | W^{\ell(a,c,c')}) = \phi_{md}$  for any  $a \in \mathcal{A}$ , and  $c, c' \in \mathcal{C}$ .

assumption 2.3: The covariance of the demeaned random shocks between different students  $i \in \mathcal{J}^{\ell(a,c)}$  and  $j \in \mathcal{J}^{\ell(a,c)}$  who are advised by the *same* professor in degree program  $g$  is given by  $\text{Cov}(\tilde{\nu}_{iag}^c, \tilde{\nu}_{jag}^c | W^{\ell(a,c,c')}) = \text{Cov}(\tilde{\nu}_{iag}^c, \tilde{\nu}_{jag}^{c'} | W^{\ell(a,c,c')}) = \psi_g$ , for any  $a \in \mathcal{A}$ ,  $g \in \{m, d\}$ , and  $c, c' \in \mathcal{C}$ .

assumption 2.4: The covariance of the demeaned random shocks between different students  $i \in \mathcal{J}^{\ell(a,c)}$  and  $j \in \mathcal{J}^{\ell(a',c')}$  who are advised by *different* professors in degree program  $g$  is zero, that is,

$$\text{Cov}(\tilde{\nu}_{iag}^c, \tilde{\nu}_{ja'g}^{c'} | W^{\ell(a,c,c')}) = \text{Cov}(\tilde{\nu}_{iag}^c, \tilde{\nu}_{ja'g}^{c'} | W^{\ell(a,c,c')}) = 0,$$

for any  $a, a' \in \mathcal{A}$ ,  $a \neq a'$ ,  $g \in \{m, d\}$ , and  $c, c' \in \mathcal{C}$ .

assumption 2.5: The covariance of the demeaned random shocks *between* different students  $i \in \mathcal{J}^{\ell(a,c)}$  and  $j \in \mathcal{J}^{\ell(a',c')}$  *between* degree programs is zero, that is,

$$\text{Cov}(\tilde{\nu}_{iam}^c, \tilde{\nu}_{ja'd}^{c'} | W^{\ell(a,c,c')}) = \text{Cov}(\tilde{\nu}_{iad}^c, \tilde{\nu}_{ja'm}^{c'} | W^{\ell(a,c,c')}) = 0$$

for any  $a, a' \in \mathcal{A}$ , and  $c, c' \in \mathcal{C}$ .

### B.2 Derivation of Equation (6)

We compute the conditional expectation of the squared left-hand side of Equation (5). Under the assumption that the random shock,  $\nu_{iag}^c$ , is orthogonal to advisor quality,  $\theta_g$ , for any student  $i \in \mathcal{J}^{\ell(a,c)}$ , professor  $a \in \mathcal{A}$ , cohort  $c \in \mathcal{C}$  and program  $g \in \{m, d\}$ , the conditional

expectation is given as follows:

$$\begin{aligned}
& \mathbb{E} \left[ \left( \overline{DD\Delta Outcome}^{\ell(a,c,c')} \right)^2 \middle| W^{\ell(a,c,c')} \right] \\
&= \mathbb{E} \left[ (\theta_{bd} - \theta_{ad})^2 \middle| W^{\ell(a,c,c')} = 1 \right] \cdot W^{\ell(a,c,c')} \\
&+ \mathbb{E} \left[ \left\{ \left( \bar{\nu}_{bd}^{\ell(a,c)} - \bar{\nu}_{am}^{\ell(a,c)} \right) - \left( \bar{\nu}_{ad}^{\ell(a,c')} - \bar{\nu}_{am}^{\ell(a,c')} \right) \right\}^2 \middle| W^{\ell(a,c,c')} = 1 \right] \cdot W^{\ell(a,c,c')} \\
&+ \mathbb{E} \left[ \left\{ \left( \bar{\nu}_{ad}^{\ell(a,c)} - \bar{\nu}_{am}^{\ell(a,c)} \right) - \left( \bar{\nu}_{ad}^{\ell(a,c')} - \bar{\nu}_{am}^{\ell(a,c')} \right) \right\}^2 \middle| W^{\ell(a,c,c')} = 0 \right] \cdot (1 - W^{\ell(a,c,c')}). \quad (\text{A.1})
\end{aligned}$$

Under assumption 1.1-1.2, we can compute the first part of Equation (A.1) as follows:

$$\mathbb{E} \left[ (\theta_{bd} - \theta_{ad})^2 \middle| W^{\ell(a,c,c')} = 1 \right] = 2\sigma_d^2(1 - \rho_d). \quad (\text{A.2})$$

We turn to the second part of Equation(A.1), which is related to the conditional expectation of when turnover occurred,  $W^{\ell(a,c,c')} = 1$ . We have the following equality concerning the value within the expectation operator:

$$\begin{aligned}
& \left\{ \left( \bar{\nu}_{bd}^{\ell(a,c)} - \bar{\nu}_{am}^{\ell(a,c)} \right) - \left( \bar{\nu}_{ad}^{\ell(a,c')} - \bar{\nu}_{am}^{\ell(a,c')} \right) \right\}^2 \\
&= \left\{ \left( \bar{\tilde{\nu}}_{bd}^{\ell(a,c)} - \bar{\tilde{\nu}}_{am}^{\ell(a,c)} \right) - \left( \bar{\tilde{\nu}}_{ad}^{\ell(a,c')} - \bar{\tilde{\nu}}_{am}^{\ell(a,c')} \right) \right\}^2 \\
&= \left\{ \left( \bar{\tilde{\nu}}_{bd}^{\ell(a,c)} \right)^2 + \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c)} \right)^2 - 2 \left( \bar{\tilde{\nu}}_{bd}^{\ell(a,c)} \right) \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c)} \right) \right\} + \left\{ \left( \bar{\tilde{\nu}}_{ad}^{\ell(a,c')} \right)^2 + \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c')} \right)^2 - 2 \left( \bar{\tilde{\nu}}_{ad}^{\ell(a,c')} \right) \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c')} \right) \right\} \\
&- 2 \left\{ \left( \bar{\tilde{\nu}}_{bd}^{\ell(a,c)} \right) \left( \bar{\tilde{\nu}}_{ad}^{\ell(a,c')} \right) - \left( \bar{\tilde{\nu}}_{bd}^{\ell(a,c)} \right) \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c')} \right) - \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c)} \right) \left( \bar{\tilde{\nu}}_{ad}^{\ell(a,c')} \right) + \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c)} \right) \left( \bar{\tilde{\nu}}_{am}^{\ell(a,c')} \right) \right\}, \quad (\text{A.3})
\end{aligned}$$

where  $\bar{\tilde{\nu}}_{ag}^{\ell(a,c)}$  is the lab  $\ell(a,c)$  average of  $\tilde{\nu}_{iag}^c$ , and hence, we use  $\bar{\nu}_{ag}^{\ell(a,c)} = \bar{\tilde{\nu}}_{ag}^{\ell(a,c)} + \bar{\nu}_g$  in the computation above. We take the conditional expectation of each piece of the last term of Equation (A.3) under assumptions 2.1 - 2.5.

1. Consider the squared term of the average demeaned error,  $\bar{\tilde{\nu}}_{pg}^{\ell(a,t)}$ , where professor  $p \in \{a, b\}$ , program  $g \in \{d, m\}$  and cohort  $t \in \{c, c'\}$ . We have the following equation:

$$\begin{aligned}
\left( \bar{\tilde{\nu}}_{pg}^{\ell(a,t)} \right)^2 &= \left[ \frac{1}{I^{\ell(a,t)}} \sum_{i \in I^{\ell(a,t)}} (\tilde{\nu}_{pig}^t) \right]^2 \\
&= \left( \frac{1}{I^{\ell(a,t)}} \right)^2 \left\{ \sum_{i \in I^{\ell(a,t)}} (\tilde{\nu}_{pig}^t)^2 + 2 \sum_{j \in I^{\ell(a,t)}} \sum_{k \neq j \in I^{\ell(a,t)}} (\tilde{\nu}_{jpg}^t) (\tilde{\nu}_{kpg}^t) \right\}.
\end{aligned}$$



Assumptions 2.1 and 2.3 lead to the following conditional expectation:

$$\mathbb{E} \left[ \left( \bar{v}_{pg}^{\ell(a,t)} \right)^2 \middle| W^{\ell(p,c)} = 1 \right] = \frac{\phi_g^2 + 2\psi_g}{I^{\ell(a,t)}}. \quad (\text{A.4})$$

2. Consider the cross-term of the average demeaned errors,  $\bar{v}_{pd}^{\ell(a,t)}$  and  $\bar{v}_{am}^{\ell(a,t)}$ , between master's and doctoral programs *within* lab  $\ell(a,t)$  for cohort  $t \in \{c, c'\}$ , where professor  $p = b$  if the professor switched from  $a$  to  $b$  due to turnover and  $p = a$  if not. Then, we have:

$$\begin{aligned} \left( \bar{v}_{pd}^{\ell(a,t)} \right) \left( \bar{v}_{am}^{\ell(a,t)} \right) &= \left[ \frac{1}{I^{\ell(a,t)}} \sum_{i \in I^{\ell(a,t)}} \left( \tilde{v}_{ipd}^t \right) \right] \cdot \left[ \frac{1}{I^{\ell(a,t)}} \sum_{i \in I^{\ell(a,t)}} \left( \tilde{v}_{iam}^t \right) \right] \\ &= \left( \frac{1}{I^{\ell(a,t)}} \right)^2 \left\{ \sum_{i \in I^{\ell(a,t)}} \left( \tilde{v}_{ipd}^t \right) \left( \tilde{v}_{iam}^t \right) + \sum_{j \in I^{\ell(a,t)}} \sum_{k \neq j \in I^{\ell(a,t)}} \left( \tilde{v}_{jpg}^t \right) \left( \tilde{v}_{kam}^t \right) \right\}. \end{aligned}$$

Given Assumption 2.2, the conditional expectation is given by:

$$\mathbb{E} \left[ \left( \bar{v}_{pd}^{\ell(a,t)} \right) \left( \bar{v}_{pm}^{\ell(a,t)} \right) \middle| W^{\ell(p,c)} = 1 \right] = \frac{\phi_{md}}{I^{\ell(a,t)}}. \quad (\text{A.5})$$

3. Consider the cross-term of the average demeaned errors between  $\bar{v}_{pg}^{\ell(a,c)}$  and  $\bar{v}_{p'g'}^{\ell(a,c')}$  *across* cohorts  $c$  and  $c'$ , where professors  $p \in \{a, b\}$  and  $p' \in \{a, b\}$  and grad programs  $g \in \{d, m\}$  and  $g' \in \{d, m\}$ . It is equal to:

$$\begin{aligned} \left( \bar{v}_{pg}^{\ell(a,c)} \right) \left( \bar{v}_{p'g'}^{\ell(a,c')} \right) &= \left[ \frac{1}{I^{\ell(a,c)}} \sum_{i \in I^{\ell(a,c)}} \left( \tilde{v}_{ipg}^c \right) \right] \cdot \left[ \frac{1}{I^{\ell(a,c')}} \sum_{j \in I^{\ell(a,c')}} \left( \tilde{v}_{jp'g'}^{c'} \right) \right] \\ &= \left( \frac{1}{I^{\ell(a,c)}} \right) \left( \frac{1}{I^{\ell(a,c')}} \right) \left\{ \sum_{i \in I^{\ell(a,c)}} \sum_{j \neq i \in I^{\ell(a,c')}} \left( \tilde{v}_{ipg}^c \right) \left( \tilde{v}_{jp'g'}^{c'} \right) \right\}. \end{aligned}$$

The conditional expectation is zero under Assumption 2.4-2.5. That is:

$$\mathbb{E} \left[ \left( \bar{v}_{pg}^{\ell(a,c)} \right) \left( \bar{v}_{p'g'}^{\ell(a,c')} \right) \middle| W^{\ell(a,c,c')} = 1 \right] = 0. \quad (\text{A.6})$$

Using results (A.4), (A.5), and (A.6) presented above, the conditional expectation of

Equation (A.3 ), regardless of whether an advisor switch occurred, is equal to:

$$\begin{aligned} & \mathbb{E} \left[ \left\{ \left( \bar{v}_{bd}^{\ell(a,c)} - \bar{v}_{am}^{\ell(a,c)} \right) - \left( \bar{v}_{ad}^{\ell(a,c')} - \bar{v}_{am}^{\ell(a,c')} \right) \right\}^2 \middle| W^{\ell(a,c,c')} = 1 \right] \\ &= \left( \frac{1}{I^{\ell(a,c)}} + \frac{1}{I^{\ell(a,c')}} \right) \{ \phi_d^2 + \phi_m^2 + 4(\psi_d + \psi_m) - 2\phi_{md} \}. \end{aligned} \quad (\text{A.7})$$

Similar computation reveals that the third part of Equation (A.1 ), for the case in which turnover did not occur,  $W^{\ell(a,c,c')} = 0$ , is the right-hand side of Equation (A.7 ).

Based on Equations (A.2 ) and (A.7 ), we therefore have the following result:

$$\mathbb{E} \left[ \left( \overline{DD\Delta Outcome}^{\ell(a,c,c')} \right)^2 \middle| W^{\ell(a,c,c')} \right] = \alpha \left( \frac{1}{I^{\ell(a,c)}} + \frac{1}{I^{\ell(a,c')}} \right) + \beta \cdot W^{\ell(a,c,c')},$$

where  $\alpha = \phi_d^2 + \phi_m^2 + 4(\psi_d + \psi_m) - 2\phi_{md}$  and  $\beta = 2\sigma_d^2(1 - \rho_d)$ .

## C Supplementary Materials for Section 4

### C.1 On computation of coauthor's credit share

What follows illustrates how the coauthor's credit share is constructed. Suppose that the names of the authors are ordered alphabetically. Then, the contribution weight is fractional: each author receives equal credit. Suppose this alphabetical approach is not used. Then, each author receives a share of credit that decreases in the authorship ranking. Following Liu and Fang (2012), the credit formula is given by  $n^{-1/k}r^{-(1-1/k)}$  for the  $r$ -th author of a paper with  $n$  authors. The integral constant,  $k$ , controls the declining rate of credit allocated in proportion to that of the first author. According to the suggestion of Liu and Fang (2012), we set  $k = 3$  for our analysis. Waltman (2012) notes that authorship could unintentionally be alphabetical, especially when the number of authors is small, despite the authors' intention to list their names based on a non-alphabetical criterion. Therefore, we account for the probabilities of both such incidental and intentional alphabetical authorship and use the expected value as the final research outcome measure.

### C.2 A Score of Word Overlap in Titles

As described in Section 4, our measure of a graduate student's research achievement is based on the number of articles that he or she published in selected physics journals.

To identify the articles that were authored by each student in the sample, we compile physics papers from the Thomson Reuters WoS archive that satisfy the following three conditions: (1) the author names match the name of the student; (2) the publication dates are in the period from the year in which the student was enrolled in graduate school to four years after he or she received a doctoral degree; and (3) the words in the title overlap to some extent with those in the title of the student's master or doctorate thesis.

The first and second conditions can be easily verified because the authors' names and publication dates of articles are available from the WoS database, whereas the student names and the degree date of each student are found in the the master's and doctoral thesis catalogs of UTokyo's physics department.

To enforce the third condition, we define a score that assesses the degree of overlap in the words in titles. Let  $\tilde{R}_i$  be the set of all physics articles that are associated with student  $i \in \mathcal{S}$  after the first and second conditions presented above are satisfied. Note that, although all articles in the set  $\tilde{R}_i$  include authors whose names are the same as student  $i$ , the student may or may not actually be the author of these articles. Such misidentification arises because of false positives in author name matching.

We use  $t(r_{ij})$  to denote the *title* of article  $r_{ij} \in \tilde{R}_i$  and use  $t_i$  to denote the title of student

$i$ 's thesis (either master's or doctoral, depending on the context). Each title of an article or a thesis consists of *words*. For each article  $r_{ij} \in \tilde{R}_i$ , we compute the following score of word overlap in titles:

$$m_{ij} = \frac{\sum_{w \in \{t_i \cap t(r_{ij})\}} \phi(w)}{\max \left\{ \sum_{w \in t_i} \phi(w), \sum_{w \in t(r_{ij})} \phi(w) \right\}}, \quad (\text{A.8})$$

where  $\phi(w)$  is a weighting of word  $w$  that measures the rareness of the word.

Indeed, the frequency of words used in article titles varies substantially, some being common and others rare. Clearly, such information is potentially useful in deciding whether an article sharing the author name with a thesis is actually authored by the person who wrote the thesis. If the words included in both the titles of an article and thesis are relatively rare, there is a higher likelihood that the authors are the same, whereas the converse is true if the words are relatively common.

To utilize the intuition,  $\phi(w)$  assigns high weight to relatively rare words and low weight to relatively common words. Following a similar approach to that proposed by Tang and Walsh (2010), we determine the weight,  $\phi(w)$ , based on the relative frequency of word  $w$ , which is computed by dividing its count frequency by the total counts of all technical terms that appear in all titles of the master's and doctoral theses of UTokyo's students. Specifically, we sort all words used in titles into five categories or quintiles based on their relative frequencies. For word  $w_k$  that is in the  $k$ -th quintile, the weight is given by  $\phi(w_k) = (6 - k)^{-2/3}$  for  $k = 1, 2, 3, 4, 5$ .

One remaining issue concerns words referring to the same concept in physics that are rendered differently. For instance, words such as “energy”, “energies”, “energetics”, and “energetic” are considered to represent the same notion. We address this issue by “standardizing” the words. Specifically, we undertake the following actions. First, we transform all non-letter, non-Greek characters and symbols into spaces. Second, we convert all words into lower case. Third, we reduce inflected (or derived) words to their word stem using a stemming algorithm.<sup>51</sup> For instance, the stemming algorithm reduces the words “energy”, “energies”, “energetics”, and “energetic” to the unique root word, “energi”. Fourth, we eliminate all of the non-informative “stopwords”, that is, very high-frequency words such as *the*, *to*, *of*, and *study*. For example, consider an article with the title “*ENERGY-LEVEL STATISTICS OF METALLIC FINE PARTICLES*.” In this case, the title is decomposed into the set of standardized root words as “energi”, “level”, “statist”, “metal”, “fine” and “particl”.

We use the title word overlap score, given by Equation (A.8), when we identify that

---

<sup>51</sup>Specifically, we use Porter's stemming algorithm, which is the most commonly used algorithm for word stemming in English.

article  $r_{ij} \in \tilde{R}_i$  is authored by student  $i$ , depending on whether the score,  $m_{ij}$ , exceeds the predetermined threshold,  $\bar{m}$ . Let  $\hat{R}_i$  be the set of articles associated with student  $i$  by the word-overlapping-score method presented above such that  $\hat{R}_i \subseteq \tilde{R}_i$ .

### C.3 An Optimal Threshold

How can we determine the threshold,  $\bar{m}$ , for the title word overlap score when matching articles and theses? Two types of matching errors are possible. We refer the first as a type 1 error, which occurs if we under-match articles, i.e., if we miss articles that are indeed authored by a student by regarding them as being written by another author. However, the second error, referred to as a type 2 error, arises when we include articles that are not authored by a target student. A type 1 error is likely to occur when we impose a threshold value,  $\bar{m}$ , that is too high, whereas a type 2 error will be more likely when we impose a low threshold,  $\bar{m}$ , and end up with spurious matches that actually belong to different authors.

One fundamental problem regarding the problem of identifying students' publications is that the true set,  $R_i$ , is unknown for student  $i \in \mathcal{S}$ , and therefore, the degrees of type 1 and type 2 errors cannot be assessed.

However, we might be able to obtain a reasonably accurate approximation set of published articles for certain students, especially for those who became academic researchers and published their CVs on the web. Let  $\bar{\mathcal{S}} \subseteq \mathcal{S}$  be the set of such students/researchers. We acquired the CVs of 40 such researchers by a random web search and parsed the research publication information to create the benchmark set of articles. Our expectation is that the benchmark article set,  $\bar{R}_i$ , will contain reliable and comprehensive information on the true set,  $R_i$ , at least for student/researcher  $i \in \bar{\mathcal{S}}$ . Nevertheless, the set  $\bar{R}_i$  might include some articles that are not directly related to their thesis projects. In this regard, the benchmark set should be close to but somewhat larger than the true set.

We use the benchmark article set to evaluate the performance of the matching procedure based on the word overlap score in titles. Specifically, to gauge the performance at each threshold value, we use two goodness-of-fit indices, *GOFI2a* and *GOFI2b*, proposed by Trajtenberg et al. (2006). Let  $\bar{R}_i$  be the benchmark set of student  $i \in \bar{\mathcal{S}}$  and  $\hat{R}_i(m)$  the corresponding set estimated by the matching procedure based on the word overlap score in titles, with  $m$  being the threshold value.

Those measures are defined as:

$$\begin{aligned} GOFI2a(m) &\equiv \text{Average} \left[ \frac{|\bar{R}_i \cap \hat{R}_i(m)|}{|\bar{R}_i|} \right] \\ GOFI2b(m) &\equiv \text{Average} \left[ \frac{|\bar{R}_i \cap \hat{R}_i(m)|}{|\hat{R}_i(m)|} \right], \end{aligned}$$

where the average is taken over all persons in the selected set  $\bar{\mathcal{J}}$ . In essence, if our matching procedure tends to under-match or over-match,  $GOFI2a(m)$  or  $GOFI2b(m)$  decrease, respectively. Therefore, we should seek to increase these indices to avoid type 1 and type 2 errors to the greatest extent possible, but a trade-off exists between the two goals.

Figure C.1 presents those two indices for various values of  $m$  in increments of 0.05.  $GOFI2b(m)$ , which is presented as a solid blue line, increases in the range of a smaller threshold value,  $m$ , and reaches nearly 0.65 when  $m = 0.25$  with no improvement being observed if  $m > 0.25$ . This leads to the implication that type 2 error will no longer be reduced dramatically if we set  $m > 0.25$ . Turning to  $GOFI2a(m)$ , which is presented as a dashed red line, it decreases consistently as the threshold value,  $m$ , rises, implying that type 1 error will be alleviated as the value of  $m$  decreases.

Accordingly, we consider the optimal threshold to be  $\bar{m} = 0.25$ , as this is the value that balances the two goodness-of-fit measures —  $GOFI2a(m)$  is maximized (thus, type 1 error is minimized) on the *condition* that  $GOFI2b(m)$  remains at a high level (thus, an increase in type 2 error is reduced as much as possible).

## Figure

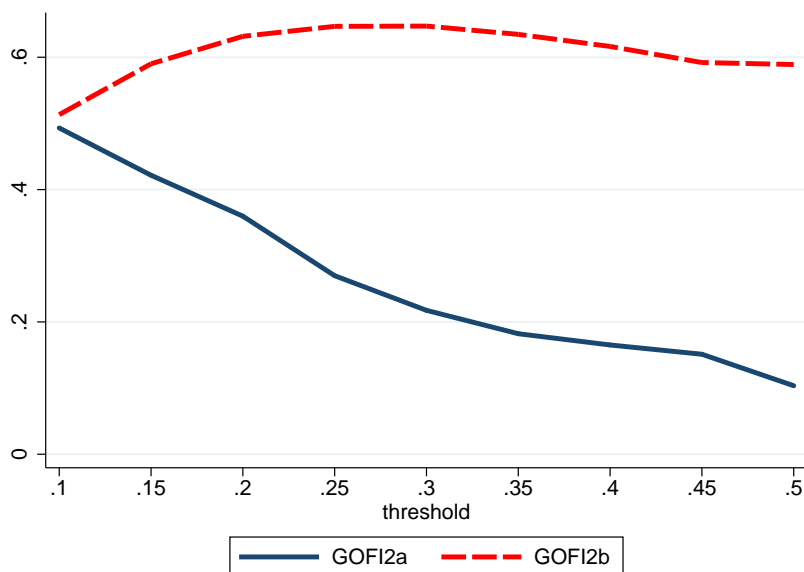


Figure C.1: Comparison of Two Goodness-of-Fit Indices over Various Thresholds for the Word Overlap Score in Titles

*Notes:* Two goodness-of-fit indices,  $GOFI2a$  and  $GOFI2b$  proposed by Trajtenberg et al. (2006), are plotted for various values of the threshold value of the matching procedure based on the word overlap score in titles in increments of 0.05.

## D Supplementary Materials for Section 5

### D.1 Discussion on Overstated Supervisory Period

In this subsection we discuss the effect of overstated supervisory period on the the lower bound estimate of the variance of advisor quality.

Consider a situation where advisor turnover happened in lab  $\ell(a, c)$  and professor  $a$  left the program and professor  $b$  took over the students who were left behind. Let  $i$  be a student in lab  $\ell(a, c)$  and let  $\lambda_i$  be a number between zero and one that captures the timing of student  $i$ ' advisor switch. Figure D.1 illustrates the situation. Student  $i$  was supervised by professor  $a$  at the master's degree course and the earlier part of the doctoral course. Since advisor switch happened in the middle of the doctoral course, the remaining proportion of  $(1 - \lambda_i)$  of the doctoral course period was supervised by the newly assigned professor  $b$ . Assume for simplicity that  $\lambda_i$  is random and independent of the factors that affect an advisee's research performance as well as the incident of advisor switch. Note that, in the baseline specification, we assume that  $\lambda_i = 0$ .

Given mixed influence on research performance from two different advisors, the research outcome growth of student  $i$  is given by

$$\Delta outcome_{iad}^c = \gamma_i + \lambda_i \theta_{ad} + (1 - \lambda_i) \theta_{bd} + \nu_{iad}^c. \quad (D.1)$$

It should be noted that the specification above encompasses the case of "hidden" supervision where the advisor who left the program due to turnover continued to provide research guidance on the former student, and the professor who was recorded as the research advisor on the thesis catalog was indeed a "surrogate" of the true advisor. In this case, we interpret the parameter  $\lambda_i$  as the magnitude of the research influence on student  $i$  from the true but "hidden" advisor  $a$ . In other words,  $(1 - \lambda_i)$  represents the fraction of the research achievement growth of student  $i$  that is attributable to the nominal "surrogate" professor  $b$ .

It follows that the  $DD$  average student research outcome growth between two labs,  $\ell(a, c)$  and  $\ell(a, c')$ , conditional on the advisor switch due to turnover,  $W^{\ell(a, c)} = 1$ , is given by

$$DD \overline{\Delta outcome}^{\ell(a, c, c')} = (1 - \bar{\lambda}^{\ell(a, c)}) (\theta_{bd} - \theta_{ad}) + error\ term \quad (D.2)$$

where  $\bar{\lambda}^{\ell(a, c)}$  is the average of  $\lambda_i$  over students in lab  $\ell(a, c)$  and the definitions of the other variables are the same as those shown in Equation (5) in Section 3.2. Then, we have the

following conditional expectation of the squared left-hand side of Equation (D.2 )

$$\begin{aligned}
& \mathbb{E} \left[ \left( \overline{DD\Delta Outcome}^{\ell(a,c,c')} \right)^2 \middle| W^{\ell(a,c,c')} \right] \\
&= \mathbb{E} \left[ (1 - \bar{\lambda}^{\ell(a,c)})^2 (\theta_{bd} - \theta_{ad})^2 \middle| W^{\ell(a,c,c')} = 1 \right] \cdot W^{\ell(a,c,c')} \\
&+ \mathbb{E} \left[ \left\{ \left( \bar{\nu}_{abd}^{\ell(a,c)} - \bar{\nu}_{am}^{\ell(a,c)} \right) - \left( \bar{\nu}_{ad}^{\ell(a,c')} - \bar{\nu}_{am}^{\ell(a,c')} \right) \right\}^2 \middle| W^{\ell(a,c,c')} = 1 \right] \cdot W^{\ell(a,c,c')} \\
&+ \mathbb{E} \left[ \left\{ \left( \bar{\nu}_{ad}^{\ell(a,c)} - \bar{\nu}_{am}^{\ell(a,c)} \right) - \left( \bar{\nu}_{ad}^{\ell(a,c')} - \bar{\nu}_{am}^{\ell(a,c')} \right) \right\}^2 \middle| W^{\ell(a,c,c')} = 0 \right] \cdot (1 - W^{\ell(a,c,c')}). \tag{D.3}
\end{aligned}$$

The first part of Equation (D.3 ) is the same as the one found in Equation (A.1 ) in the Appendix B.2, except that the squared difference of advisor qualities,  $(\theta_{bd} - \theta_{ad})^2$ , is multiplied by  $(1 - \bar{\lambda}^{\ell(a,c)})^2$ , while the second and third part of Equation (D.3 ) is exactly the same as the corresponding terms of Equation (A.1 ). Under the assumptions presented in Section 3.2, the first part of Equation (D.3 ) is equal to

$$(1 - \bar{\lambda}^{\ell(a,c)})^2 \times 2\sigma_d^2(1 - \rho_d). \tag{D.4}$$

Considering the agreement between Equation (A.1 ) and Equation (D.3 ) except the first term, we therefore have the following result:

$$\mathbb{E} \left[ \left( \overline{DD\Delta Outcome}^{\ell(a,c,c')} \right)^2 \middle| W^{\ell(a,c,c')} \right] = \alpha \left( \frac{1}{N^{\ell(a,c)}} + \frac{1}{N^{\ell(a,c')}} \right) + \check{\beta} \cdot W^{\ell(a,c,c')}, \tag{D.5}$$

where  $\alpha = \phi_d^2 + \phi_m^2 + 4(\psi_d + \psi_m) - 2\phi_{md}$  and  $\check{\beta} = 2\sigma_d^2(1 - \rho_d)(1 - \bar{\lambda}^{\ell(a,c)})^2$ .

If we turn to the regression equation where  $(\overline{DD\Delta Outcome}_n)^2$  is regressed on advisor switch indicator variable  $W_n$ , as presented in Equation (7) in Section 3.2, the coefficient  $\check{\beta}$  of  $W_n$  relates to the variance through the following equation

$$\check{\beta} = 2\sigma_d^2(1 - \rho_d)(1 - \bar{\lambda}^{\ell(a,c)})^2 \tag{D.6}$$

Since correlation is imperfect ( $|\rho_d| < 1$ ) and supervision is partial ( $0 < \bar{\lambda}_d < 1$ ), a lower bound estimate is still given by the last term of the following inequality

$$\hat{\sigma}_d^2 = \frac{\check{\beta}}{2(1 - \rho_d)(1 - \bar{\lambda}^{\ell(a,c)})^2} \geq \frac{\check{\beta}}{2(1 - \rho_d)} \geq \frac{\check{\beta}}{4}. \tag{D.7}$$

That is,  $\widehat{\text{Lower bound of } \sigma_d^2}$  is given by one-fourth of the estimated coefficient of the advisor switch indicator variable in the regression model.



## Figure

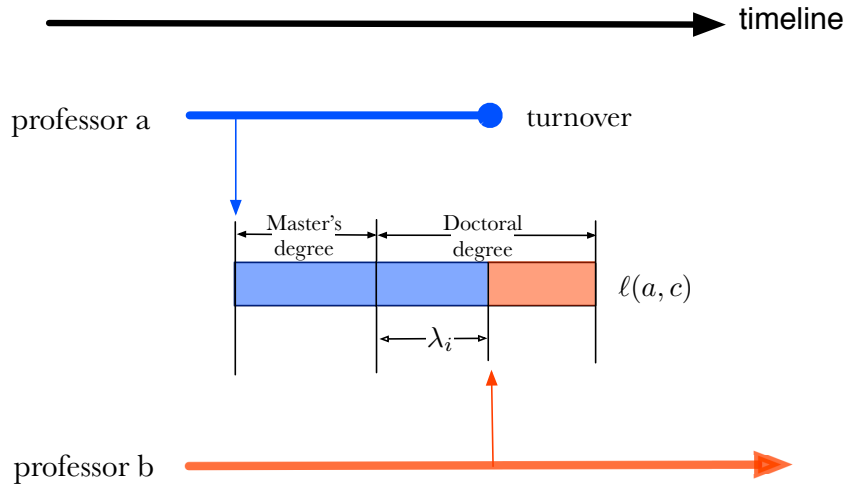


Figure D.1: Example of Partial Supervisory Period

*Notes:* Advisor turnover happened in lab  $\ell(a, c)$  and professor  $a$  left the program and professor  $b$  took over the students who were left behind. Let  $\lambda_i$  be a number between zero and one that captures the timing of student  $i$ ' advisor switch.

## D.2 Discussion on Non-random Turnover

### (1) Evidence suggesting that the sample is not balanced

Table D.1 reports the descriptive statistics for some characteristics of advisors and compares those of advisors when turnover occurred and the corresponding advisor characteristics when it did not. We find that, for some characteristics, the differences in means between the two groups, professors with turnover in column (1) and those without in column (2), are statistically significant at the 5 percent level. We also find that the absolute values of the standardized differences, reported in column (3), are large for some characteristics.<sup>52</sup>

### (2) A procedure for the propensity score matching method

To make the sample balanced and comparable, we thus employ a propensity score matching method. Following standard practice in the literature, we estimate the propensity scores using a logit model. - We include all of the characteristics presented in Table D.1 scores. Specifically, in the first step, we begin with a set of basic covariates and add an additional linear term based on a likelihood ratio test for the null hypothesis that the coefficient of the added variable is

<sup>52</sup>The standardized difference considers the size of the difference in means of a conditioning variable, scaled by the square root of the variances of the treatment and control groups in the original sample. According to the suggestion of Rosenbaum and Rubin (1985), an absolute value of the standardized difference greater than 0.2 should be considered “large”.

equal to zero. In the second step, we proceed to the choice of the quadratic and cross-product terms and apply the same type of likelihood test as that used in the first step. We follow the suggestion of Imbens and Rubin (2015) that the threshold values for the likelihood ratio test should be  $C_L = 1.0$  and  $C_Q = 2.71$  for the linear and quadratic terms, respectively. The results of the logit estimation of the propensity score can be found in Table XX. Given the estimated propensity scores, we match a case with  $W_n = 1$  (a lab with an advisor switch) to one with  $W_n = 0$  (a lab without an advisor switch) that share approximately identical estimated propensity scores. We employ a one-to-one nearest-neighbor matching method.

To assess the quality of the propensity score matching, we present Figure D.3 that depicts the absolute values in the standardized differences of the variables for the original and matched samples. The imbalance between the treatment and control cases is attenuated on many professor characteristics. For example, professor's age differs between the treatment and control labs by more than the average standard deviation (the absolute standard deviation is 1.129) before matching, whereas the difference is considerably reduced (the absolute standard deviation is 0.006) after matching.

Figure D.4 presents the distributions of the estimated propensity scores for the treatment labs (left) and control labs (right) in each case of the adjacent period,  $\tau = 3, 4,$  and  $5$ . The top and bottom groups in the graphs correspond to those before and after matching, respectively. Before matching, the shapes of distributions differ considerably between the treatment and control groups. Nevertheless, the propensity score distributions have some degree of overlap. Moreover, after matching, the dissimilarity of the distributions between the treatment and control groups is considerably reduced.

One might worry that the spread of the common support of the propensity score distributions should not be across the full range  $[0, 1]$  and hence that the observations of the treatment group, especially those with high propensity scores, are matched forcibly with those of the control group, the propensity scores of which are not sufficiently close. To address the problem caused by the limited common support of the propensity score distribution, we employ a systematic approach proposed by Crump et al. (2009) and discard all observations with estimated propensity scores outside the range of  $[0.1, 0.9]$ .

## Figures

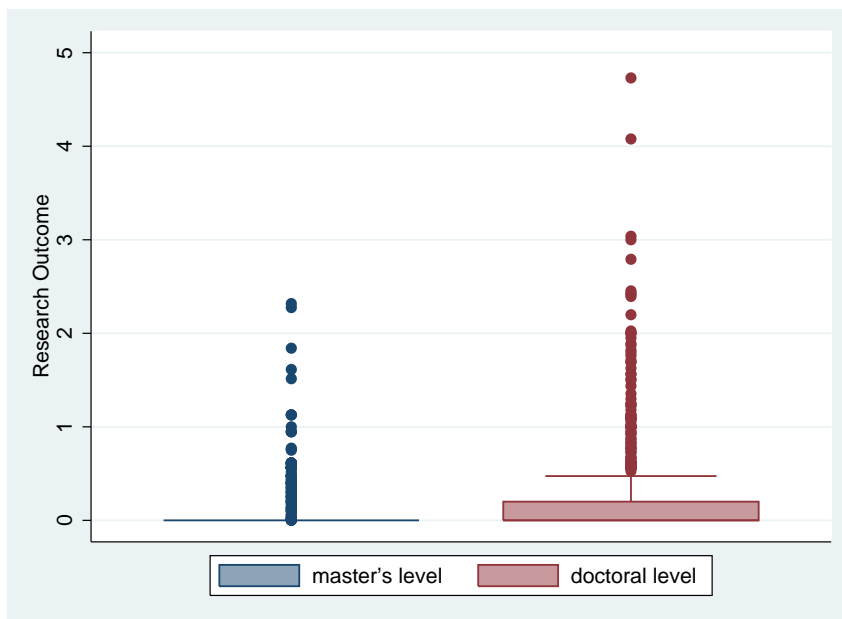


Figure D.2: Box Plots of the Student Research Outcome Distributions in the Master's and Doctoral Programs

*Notes:* Student research outcomes are aggregated over the period from M1 to D2 for the master's degree, and the period from D1 to P4 for doctoral degree program.

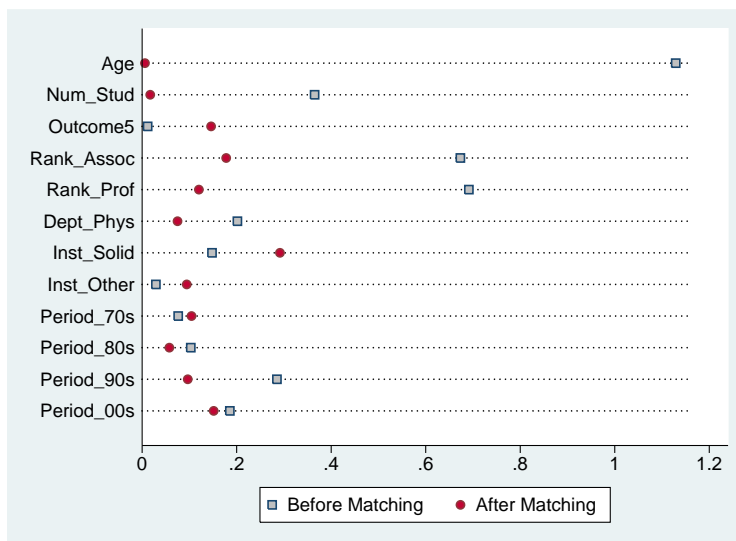


Figure D.3: Comparison of the Absolute Values of the Standardized Differences between Treatment and Control Groups

*Notes:* The standardized difference considers the size of the difference in means of a conditioning variable, scaled by the square root of the variances of the treatment and control groups in the original sample. According to the suggestion of Rosenbaum and Rubin (1985), an absolute value of the standardized difference greater than 0.2 should be considered large.

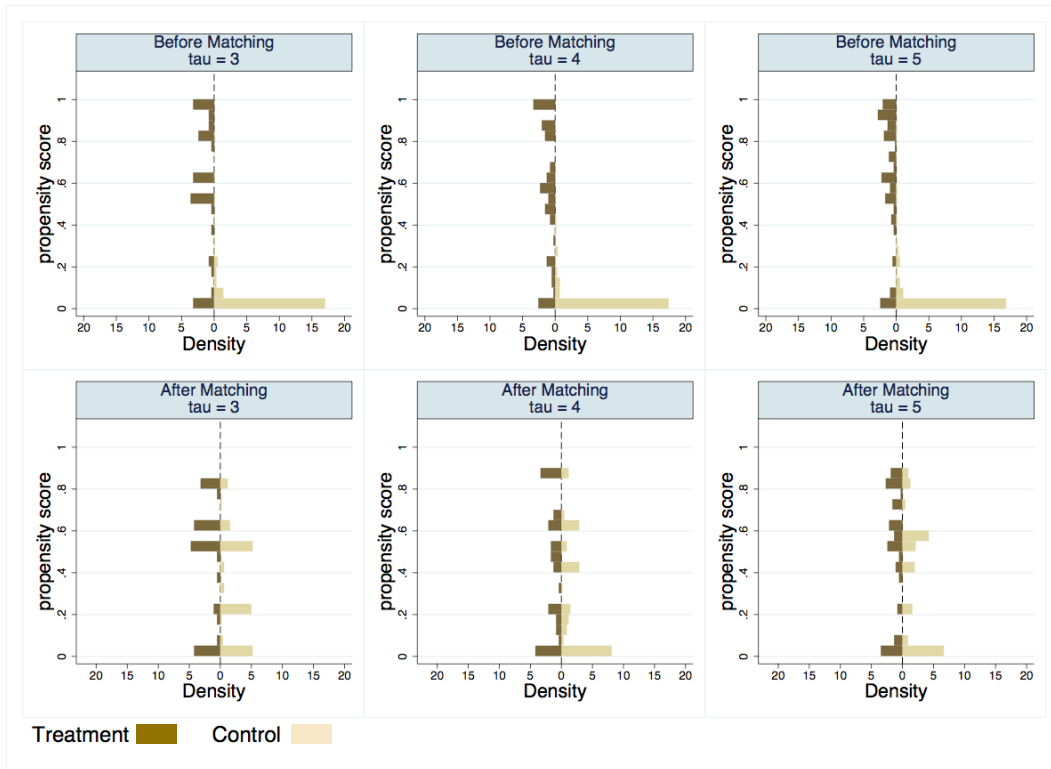


Figure D.4: Distribution of the Propensity Score for the Treatment and Control Groups: Before and After Matching

*Notes:* The distribution of the estimated propensity scores are presented for the treatment labs and control labs (right) in each case of the adjacent period,  $\tau = 3, 4,$  and  $5$ . The top and bottom groups in the graphs correspond to those before and after matching, respectively.

## Tables

Table D.1: Descriptive Statistics of Advisors: Comparison between Advisors When Turnover Occurred and When It Did Not

Variable	Description	With Turnover	Without Turnover	t-stat	Absolute Standardized Difference
<i>Age</i>	Professor Age	53.72 (6.37)	47.06 (5.77)	-9.00 ***	1.10
<i>Num_Stud</i>	Number of Students	1.16 (0.41)	1.27 (0.52)	1.63	0.22
<i>Outcome5</i>	Professor's Research Outcome (5 years average)	0.18 (0.29)	0.21 (0.40)	0.50	0.07
<i>Rank_Assoc</i>	Associate Professor Dummy	0.21 (0.41)	0.44 (0.50)	3.73 ***	0.51
<i>Rank_Prof</i>	Full Professor Dummy	0.79 (0.41)	0.53 (0.50)	-4.27 ***	0.59
<i>Dept_Phys</i>	Department of Physics Dummy	0.72 (0.45)	0.75 (0.43)	0.62	0.08
<i>Inst_Solid</i>	Institute of Solid State Dummy	0.72 (0.41)	0.22 (0.41)	0.27	0.03
<i>Inst_Other</i>	Other Institutes Dummy	0.34 (0.48)	0.29 (0.45)	-0.85	0.11
<i>Period_70s</i>	70's Dummy	0.15 (0.36)	0.21 (0.41)	1.19	0.16
<i>Period_80s</i>	80's Dummy	0.16 (0.37)	0.26 (0.44)	1.71 *	0.23
<i>Period_90s</i>	90's Dummy	0.57 (0.50)	0.35 (0.48)	-3.72 ***	0.46
<i>Period_00s</i>	00's Dummy	0.12	0.19	1.46	0.20

*Notes:* The absolute standardized difference is given by the size of the difference in means of a conditioning variable, scaled by the square root of the variances in the original samples (Rosenbaum and Rubin, 1985).

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table D.2: Estimation Results of Propensity Score

Adjacent Period	$\tau = 3$		$\tau = 4$		$\tau = 5$	
<i>Age</i>	-5.823	***	-6.494	***	-6.784	***
	[9.70]		[9.39]		[6.55]	
<i>Num_Stud</i>	21.900		-0.701	*	-0.212	
	[0.02]		[1.86]		[0.46]	
<i>Rank_Assoc</i>	0.832		1.109		-0.552	
	[0.76]		[1.00]		[0.26]	
<i>Inst_Other</i>	-0.208		-0.462		-0.308	
	[0.46]		[0.99]		[0.46]	
<i>Period_90s</i>	1.014	**	1.777	***	1.101	
	[1.99]		[4.11]		[1.49]	
<i>Period_00s</i>	-1.724	***	-1.375	**	-1.782	*
	[2.67]		[2.44]		[1.92]	
<i>Age</i> <sup>2</sup>	0.063	***	0.072	***	0.073	***
	[10.11]		[9.58]		[6.84]	
<i>Num_Stud</i> <sup>2</sup>	-7.263					
	[0.02]					
<i>Age</i> $\times$ <i>Period_80s</i>	1.465	**			1.157	
	[2.47]				[1.56]	
<i>Num_Stude</i> $\times$ <i>Period_80s</i>	-80.110	**			-63.27	
	[2.43]				[1.53]	
<i>Outcome5</i> $\times$ <i>Inst_Other</i>	-1.816	*	-2.325	**	-3.059	**
	[1.76]		[2.06]		[1.99]	
<i>Rank_Assoc</i> $\times$ <i>Inst_Solid</i>			2.346	*	2.941	*
			[1.90]		[1.83]	
<i>Inst_Solid</i> $\times$ <i>Period_00s</i>	3.487	***	3.168	**		
	[3.02]		[2.29]			
<i>Inst_Other</i> $\times$ <i>Period_80s</i>	-3.700	***			-2.757	*
	[2.95]				[1.83]	
Constant	111.7		138.900	***	148.800	***
	[0.16]		[8.98]		[6.05]	
Sample Size	1446		1202		925	

*Notes:* The dependent variable is the advisor switch indicator due to turnover,  $W_n$ . The definitions of the independent variables are given in Table D.1. The specification of the model is given by a stepwise likelihood-test-based procedure, suggested by Imbens (2014) and Imbens and Rubin (2015).

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

## E Supplementary Materials for Section 6

### E.1 Indirect Influence from Non-advisor

We explain an indirect effect from non-advisor faculty members on students across labs by Figure E.1 illustrate, which parallels that illustrated in Figure 2. As we have assumed previously, there are three cohorts,  $c_0$ ,  $c_1$  and  $c_2$  where an instance of turnover involving professor  $a$  occurred in cohort  $c_2$ , and the students in lab  $\ell(a, c_2)$  switched their research advisor from professor  $a$  to professor  $b$  in the doctoral program. Then, professor  $b$ , whose research area is the same as that of professor  $a$ , took over the students in lab  $\ell(a, c_2)$ , whereas he had supervised two labs,  $\ell(b, c_0)$  and  $\ell(b, c_1)$ , before the incident occurred, and oversaw another lab,  $\ell(b, c_2)$ , at the time of the incident. We assume that professor  $a$ 's turnover affects the doctoral research productivity of the students in lab  $\ell(b, c_2)$  because the indirect influence from the professor,  $\theta_{ad}$ , ceases to exist after turnover.

In the estimation that follows, we choose the lab of professor  $b$  that was influenced “indirectly” by professor  $a$ 's turnover if the turnover occurred while the students in the lab were in the doctoral program (i.e., from the first doctoral year to the final year of the doctoral program). We require this because the indirect influence from professor  $a$  at the doctoral degree level, not the master's degree level, needs to be changed. In this case, to identify the magnitude of the indirect impact, we essentially compare the gap in student research outcome growth between labs  $\ell(b, c_2)$  and  $\ell(b, c_1)$  (treatment group with  $V^{\ell(b, c_2, c_1)} = 1$ ) with the same gap between labs  $\ell(b, c_1)$  and  $\ell(b, c_0)$  (control group with  $V^{\ell(b, c_1, c_0)} = 0$ ).

# Figure

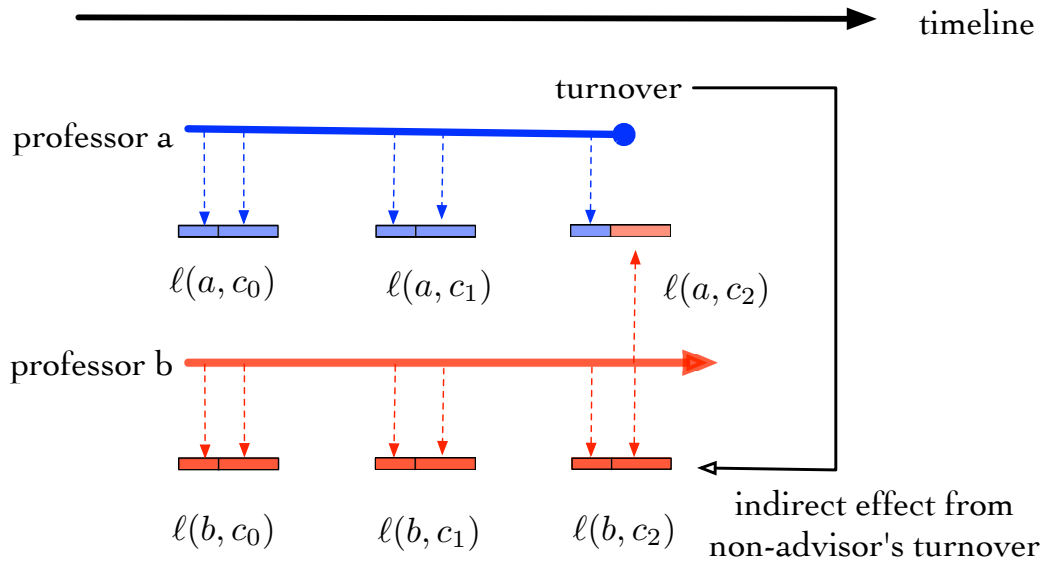


Figure E.1: Example of Labs with and without Turnover for Advisors and Non-Advisors

*Notes:* There are three cohorts,  $c_0$ ,  $c_1$  and  $c_2$  where an instance of turnover involving professor  $a$  occurred in cohort  $c_2$ , and the students in lab  $\ell(a, c_2)$  switched their research advisor from professor  $a$  to professor  $b$  in the doctoral program. Then, professor  $b$ , whose research area is the same as that of professor  $a$ , took over the students in lab  $\ell(a, c_2)$ , whereas he had supervised two labs,  $\ell(b, c_0)$  and  $\ell(b, c_1)$ , before the incident occurred, and oversaw another lab,  $\ell(b, c_2)$ , at the time of the incident.



## E.2 Supplementary Tables

Table E.1: Estimation Results Regarding Different Specifications of Student Research Outcomes: The Estimated Lower Bound of the Variance in Advisor Quality at the Doctoral Level Is Reported for the Case of Twelve Journals

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\Delta outcome]^2$			$[DD\Delta outcome]^2$		
	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
Adjacent Period	(1)	(2)	(3)	(4)	(5)	(6)
Threshold 0.25 (default)						
M1-D2/D1-P4 <sup>††</sup>	0.0843 ** [0.0268]	0.0666 ** [0.0135]	0.0489 ** [0.0236]	0.5773 ** [0.0407]	0.5331 ** [0.0104]	0.4100 ** [0.0234]
M1-D2/D1-P3 <sup>†</sup>	0.0689 ** [0.0232]	0.0498 ** [0.0181]	0.0257 * [0.0973]	0.4578 ** [0.0310]	0.4099 *** [0.0080]	0.2793 ** [0.0338]
M1-D1/D1-P4 <sup>†</sup>	0.1460 ** [0.0488]	0.1243 ** [0.0210]	0.0910 ** [0.0389]	0.9274 ** [0.0450]	0.8522 ** [0.0133]	0.7525 ** [0.0145]
M1-D1/D1-P3 <sup>†</sup>	0.1215 ** [0.0485]	0.0966 ** [0.0271]	0.0572 * [0.0913]	0.7711 ** [0.0384]	0.6911 ** [0.0119]	0.5745 ** [0.0175]
Threshold 0.20 (overmatch)						
M1-D2/D1-P4 <sup>†</sup>	0.0953 ** [0.0407]	0.0790 ** [0.0184]	0.0576 ** [0.0341]	0.7246 ** [0.0460]	0.6833 ** [0.0115]	0.5492 ** [0.0213]
M1-D2/D1-P3 <sup>†</sup>	0.0621 ** [0.0323]	0.0425 ** [0.0341]	0.0205 [0.1505]	0.4672 ** [0.0280]	0.4179 *** [0.0071]	0.2641 ** [0.0438]
M1-D1/D1-P4 <sup>†</sup>	0.1659 * [0.0580]	0.1466 ** [0.0231]	0.1214 ** [0.0247]	1.1114 ** [0.0490]	1.0467 ** [0.0134]	2.2262 *** [0.0085]
M1-D1/D1-P3 <sup>†</sup>	0.1159 * [0.0565]	0.0913 ** [0.0349]	0.0639 * [0.0665]	0.7805 ** [0.0364]	0.7055 ** [0.0105]	0.6066 ** [0.0127]
Threshold 0.30 (undermatch)						
M1-D2/D1-P4 <sup>†</sup>	0.0509 * [0.0712]	0.0378 ** [0.0476]	0.0191 [0.1037]	0.3490 ** [0.0222]	0.3135 *** [0.0042]	0.1795 * [0.0548]
M1-D2/D1-P3 <sup>†</sup>	0.0401 * [0.0719]	0.0282 * [0.0609]	0.0023 [0.4395]	0.2662 ** [0.0110]	0.2282 *** [0.0019]	0.0884 [0.1246]
M1-D1/D1-P4 <sup>†</sup>	0.1123 * [0.0599]	0.0988 ** [0.0237]	0.0683 * [0.0541]	0.6166 ** [0.0345]	0.5541 *** [0.0094]	0.4385 ** [0.0220]
M1-D1/D1-P3 <sup>†</sup>	0.0923 * [0.0582]	0.0783 ** [0.0261]	0.0360 [0.1507]	0.4971 ** [0.0265]	1.0168 *** [0.0094]	0.3002 ** [0.0331]

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The estimated lower bounds of the variance in advisor quality at the doctoral level,  $\overline{\text{Lower bound of } \sigma_d^2}$  are reported. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that  $\overline{\text{Lower bound of } \sigma_d^2} = 0$  against the alternative  $\overline{\text{Lower bound of } \sigma_d^2} > 0$ . The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses.

<sup>†</sup> (master's level aggregation period) / (doctoral level aggregation period)

<sup>††</sup> The baseline cases.

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table E.2: Estimation Results Regarding Different Specifications of Student Research Outcomes: The Estimated Lower Bound of the Variance in Advisor Quality at the Doctoral Level Is Reported for the Case of Nine Journals

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\Delta outcome]^2$			$[DD\Delta outcome]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
M1-D2/D1-P4 <sup>†</sup>	0.0449 *** [0.0005]	0.0363 *** [0.0000]	0.0211 *** [0.0079]	0.1862 *** [0.0039]	0.1591 *** [0.0003]	0.0241 [0.3100]
M1-D2/D1-P3 <sup>†</sup>	0.0377 *** [0.0027]	0.0262 *** [0.0000]	0.0070 [0.2448]	0.1403 *** [0.0010]	0.1117 *** [0.0000]	0.0753 ** [0.0449]
M1-D1/D1-P4 <sup>†</sup>	0.0829 *** [0.0080]	0.0675 *** [0.0017]	0.0355 ** [0.0431]	0.4216 ** [0.0130]	0.3441 [0.5151]	0.2079 ** [0.0395]
M1-D1/D1-P3 <sup>†</sup>	0.0680 *** [0.0085]	0.0495 *** [0.0026]	0.0115 [0.2533]	0.3389 *** [0.0084]	0.2588 *** [0.0025]	0.1168 * [0.0873]

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The estimated lower bounds of the variance in advisor quality at the doctoral level, Lower bound of  $\sigma_d^2$  are reported. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that Lower bound of  $\sigma_d^2 = 0$  against the alternative Lower bound of  $\sigma_d^2 > 0$ . The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses.

<sup>†</sup> (master's level aggregation period) / (doctoral level aggregation period)

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table E.3: Falsification Test Results (Full Results)

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\overline{\Delta outcome}]^2$			$[DD\overline{\Delta outcome}]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\tilde{\alpha}$	0.1827 *** (0.0418)	0.1902 *** (0.0409)	0.1736 *** (0.0443)	0.2819 *** (0.0458)	0.6911 *** (0.1828)	0.8996 *** (0.2287)
(2) $\tilde{\beta}$	0.2979 (0.2123)	0.2039 (0.1725)	0.0814 (0.1270)	0.6315 ** (0.2616)	-0.0969 (0.3965)	-0.7784 (0.4350)
Sample Size						
Total	887	1136	1351	887	1136	1351
After matching	422	763	603	422	763	603

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The total sample size is given by the number of observations for each tuple of  $(a, c, c')$  for any advisor  $a$  in  $\mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$  where  $\tau$  is the period over which the difference is taken. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The after-matching sample size is the sum of the numbers of observations for the treatment group where turnover occurred and the corresponding control group that are matched through the propensity score method. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that  $\overline{\text{Lower bound of } \sigma_d^2} = 0$  against the alternative  $\overline{\text{Lower bound of } \sigma_d^2} > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table E.4: Estimation Results: The Student Proficiency Score is Set to Zero If the Student Coauthored with the Advisor (Full Results)

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\overline{\Delta outcome}]^2$			$[DD\overline{\Delta outcome}]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\hat{\alpha}$	0.0544 *** (0.0230)	0.0670 *** (0.0225)	0.0439 *** (0.0091)	0.0215 (0.0448)	0.0713 * (0.0419)	0.0528 ** (0.0250)
(2) $\hat{\beta}$	0.2057 (0.1727)	0.1665 (0.1187)	0.2374 ** (0.1022)	1.3776 (0.9758)	1.3374 ** (0.6807)	1.3690 ** (0.5814)
(3) $\overline{\text{Lower bound of } \sigma_d^2}$	0.0514 [0.1168]	0.0416 * [0.0803]	0.0593 ** [0.0101]	0.3444 * [0.0790]	0.3343 ** [0.0247]	0.3422 *** [0.0093]
Sample Size						
Total	925	1202	1446	925	1202	1446
After matching	104	186	271	104	186	271

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The total sample size is given by the number of observations for each tuple of  $(a, c, c')$  for any advisor  $a$  in  $\mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$  where  $\tau$  is the period over which the difference is taken. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The after-matching sample size is the sum of the numbers of observations for the treatment group where turnover occurred and the corresponding control group that are matched through the propensity score method. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that  $\overline{\text{Lower bound of } \sigma_d^2} = 0$  against the alternative  $\overline{\text{Lower bound of } \sigma_d^2} > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table E.5: Estimation Results: When a Change in Advisor Quality Variance Is Allowed during the Period Near Turnover (Full Results)

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\overline{\Delta outcome}]^2$			$[DD\overline{\Delta outcome}]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\hat{\alpha}$	0.0717 *** (0.0234)	0.0604 *** (0.0156)	0.0982 *** (0.0144)	0.0622 (0.0674)	0.1395 ** (0.0614)	0.4297 *** (0.0752)
(2) $\hat{\beta}$	0.3391 * (0.1839)	0.2431 * (0.1294)	0.2145 ** (0.1065)	2.3320 (1.4204)	2.1209 ** (0.9704)	1.7108 * (0.8976)
(3) $\hat{\delta}_1$	—	-0.1207 (0.032)	-0.0085 (0.022)	—	-0.2789 (0.123)	-0.5772 (0.181)
(4) $\hat{\delta}_2$	-0.1512 (0.2259)	0.4291 (0.296)	-0.2504 (0.152)	-0.4563 (1.5777)	-0.1015 (0.938)	-1.4842 (1.093)
(5) $\hat{\delta}_3$	-0.1195 (0.0431)	—	—	-0.1036 (0.0843)	—	—
(6) Lower bound of $\sigma_d^2$	0.0848 ** [0.0326]	0.0608 ** [0.0301]	0.0536 ** [0.0220]	0.5830 * [0.0503]	0.5302 ** [0.0144]	0.4277 ** [0.0283]
Sample Size						
Total	925	1202	1446	925	1202	1446
After matching	104	186	271	104	186	271

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The total sample size is given by the number of observations for each tuple of  $(a, c, c')$  for any advisor  $a$  in  $\mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$  where  $\tau$  is the period over which the difference is taken. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The after-matching sample size is the sum of the numbers of observations for the treatment group where turnover occurred and the corresponding control group that are matched through the propensity score method. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that Lower bound of  $\sigma_d^2 = 0$  against the alternative

Lower bound of  $\sigma_d^2 > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table E.6: Estimation Results when Non-Retirement Turnover Events Are Used: Baseline Case (Full Results)

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\overline{\Delta outcome}]^2$			$[DD\overline{\Delta outcome}]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\hat{\alpha}$	-0.0047 (0.0369)	0.0104 (0.0174)	0.0659 ** (0.0264)	0.0953 (0.2117)	0.0996 (0.1453)	0.3749 *** (0.1269)
(2) $\hat{\beta}$	0.7453 * (0.4522)	0.7490 *** (0.2782)	0.4959 ** (0.2085)	4.9277 (3.7047)	5.2857 ** (2.3709)	4.0857 ** (1.7469)
(3) Lower bound of $\sigma_d^2$	0.1863 ** [0.0496]	0.1872 *** [0.0035]	0.1240 *** [0.0087]	0.0000 * [0.0000]	0.9448 ** [0.0907]	0.9980 ** [0.0110]
Sample Size						
Total	894	1154	1381	894	1154	1381
After matching	27	52	85	27	52	85

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The total sample size is given by the number of observations for each tuple of  $(a, c, c')$  for any advisor  $a$  in  $\mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$  where  $\tau$  is the period over which the difference is taken. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The after-matching sample size is the sum of the numbers of observations for the treatment group where turnover occurred and the corresponding control group that are matched through the propensity score method. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that Lower bound of  $\sigma_d^2 = 0$  against the alternative

Lower bound of  $\sigma_d^2 > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table E.7: Estimation Results when Non-Retirement Turnover Events Are Used: The Student Proficiency Score is Set to Zero If the Student Coauthored with the Advisor (Full Results)

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\overline{\Delta outcome}]^2$			$[DD\overline{\Delta outcome}]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\hat{\alpha}$	-0.0316 (0.0298)	0.0053 (0.0171)	0.0436 (0.0218)	-0.1398 (0.1937)	0.0785 (0.1079)	0.1862 (0.0760)
(2) $\hat{\beta}$	0.5929 (0.4475)	0.5731 (0.2777)	0.5303 (0.2292)	3.6030 (2.8352)	3.5205 (1.7587)	2.8782 (1.3485)
(3) Lower bound of $\sigma_d^2$	0.1482 * [0.0926]	0.1433 ** [0.0195]	0.1326 ** [0.0103]	0.9007 [0.1019]	0.8801 ** [0.0227]	0.7196 ** [0.0164]
Sample Size						
Total	894	1154	1381	894	1154	1381
After matching	27	52	85	27	52	85

Table E.8: Estimation Results: the Double-Difference Measure in Levels Is Used as the Dependent Variable (Full Results)

Dependent	Credit Share Weighted			First-authored-paper Based		
	[ $DD\overline{\Delta outcome}$ ]			[ $DD\overline{\Delta outcome}$ ]		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\hat{\alpha}$	-0.0462 (0.0237)	-0.0186 (0.0225)	0.0362 *** (0.0134)	-0.0655 (0.0961)	-0.0693 (0.0856)	0.0963 (0.0893)
(2) $\hat{\beta}$	0.1439 (0.1746)	0.1480 (0.1204)	0.0310 (0.0985)	0.2655 (2.1883)	0.3968 (1.5374)	0.2276 (1.3783)
Sample Size						
Total	925	1202	1446	925	1202	1446
After matching	104	186	271	104	186	271

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The total sample size is given by the number of observations for each tuple of  $(a, c, c')$  for any advisor  $a$  in  $\mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$  where  $\tau$  is the period over which the difference is taken. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The after-matching sample size is the sum of the numbers of observations for the treatment group where turnover occurred and the corresponding control group that are matched through the propensity score method. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that Lower bound of  $\sigma_d^2 = 0$  against the alternative

Lower bound of  $\sigma_d^2 > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.

Table E.9: Estimation Results: Effect of Non-Advisor Turnover on Student Research Outcome Growth at the Doctoral Level (Full Results)

Dependent	Credit Share Weighted			First-authored-paper Based		
	$[DD\Delta outcome]^2$			$[DD\Delta outcome]^2$		
Adjacent Period	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 3$	$\tau = 4$	$\tau = 5$
	(1)	(2)	(3)	(4)	(5)	(6)
(1) $\hat{\alpha}_{ind}$	0.0434 *** (0.0175)	0.0356 *** (0.0056)	0.0335 *** (0.0060)	0.0988 *** (0.0358)	0.0863 *** (0.0094)	0.0755 *** (0.0136)
(2) $\hat{\beta}_{ind}$	-0.0230 (0.0374)	0.0274 (0.0320)	0.0527 (0.0336)	0.0764 (0.1192)	0.1007 (0.0924)	0.1768 * (0.0924)
(3) $\hat{\pi}^2 = \hat{\beta}_{ind}/\hat{\beta}_{dir}$	-0.0682	0.1030	0.2694	0.0331	0.0472	0.1078
Sample Size						
Total	858	1105	1317	858	1105	1317
After matching	145	282	288	145	282	288

*Notes:* The dependent variable is the squared double-differenced average student research outcome growth. The total sample size is given by the number of observations for each tuple of  $(a, c, c')$  for any advisor  $a$  in  $\mathcal{A}$  and cohort  $c, c'$  such that  $0 < c - c' \leq \tau$  where  $\tau$  is the period over which the difference is taken. To make the sample balanced, a propensity score matching method is used. A logit model is used to estimate the propensity scores. The after-matching sample size is the sum of the numbers of observations for the treatment group where turnover occurred and the corresponding control group that are matched through the propensity score method. The standard errors that are computed by the subsampling method of Politis and Romano (1994) are in parentheses. The numbers in square brackets are p-values for the one-sided tests such that  $\text{Lower bound of } \sigma_a^2 = 0$  against the alternative Lower bound of  $\sigma_a^2 > 0$ .

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

\* Significant at the 10 percent level.