

The Education-Innovation Gap

Barbara Biasi, Song Ma

Impressum:

CESifo Working Papers

ISSN 2364-1428 (electronic version)

Publisher and distributor: Munich Society for the Promotion of Economic Research - CESifo GmbH

The international platform of Ludwigs-Maximilians University's Center for Economic Studies and the ifo Institute

Poschingerstr. 5, 81679 Munich, Germany

Telephone +49 (0)89 2180-2740, Telefax +49 (0)89 2180-17845, email office@cesifo.de

Editor: Clemens Fuest

<https://www.cesifo.org/en/wp>

An electronic version of the paper may be downloaded

- from the SSRN website: www.SSRN.com
- from the RePEc website: www.RePEc.org
- from the CESifo website: <https://www.cesifo.org/en/wp>

The Education-Innovation Gap

Abstract

This paper documents differences across higher-education courses in the coverage of frontier knowledge. Comparing the text of 1.7M syllabi and 20M academic articles, we construct the “education-innovation gap,” a syllabus’s relative proximity to old and new knowledge. We show that courses differ greatly in the extent to which they cover frontier knowledge. More selective and better funded schools, and those enrolling socio-economically advantaged students, teach more frontier knowledge. Instructors play a big role in shaping course content; research-active instructors teach more frontier knowledge. Students from schools teaching more frontier knowledge are more likely to complete a PhD, produce more patents, and earn more after graduation.

JEL-Codes: I230, I240, I260, J240, O330.

Keywords: education, innovation, syllabi, instructors, text analysis, inequality.

Barbara Biasi
EIEF, Rome / Italy &
Yale School of Management
barbara.biasi@yale.edu

Song Ma
Yale School of Management
Yale University, New Haven / CT / USA
song.ma@yale.edu

Please click here for the most updated version:

https://www.barbarabiasi.com/uploads/1/0/1/2/101280322/biasi_ma_2022.pdf

March 14, 2022

We thank Jaime Arellano-Bover, David Deming, Richard Freeman, David Robinson, Kevin Stange, Sarah Turner, and seminar and conference participants at Yale, Erasmus, Maastricht, Harvard (HBS; HGSE), Ohio State, HKU, Stanford (GSB; Hoover), UCL, Queens, Stockholm School of Economics, Duke, IIES Stockholm, NBER (Education; Entrepreneurship; Innovation), HEC (Paris), AEA, CEPR/Bank of Italy, MIT (Sloan), Boston University (Wheelock), UConn, Junior Entrepreneurial Finance and Innovation Workshop, Baruch, SOLE, IZA TOM and Economics of Education Conferences, and CESifo Economics of Education Conference for helpful comments. Xugan Chen provided outstanding research assistance. We thank the Yale Tobin Center for Economic Policy, Yale Center for Research Computing, Yale University Library, and Yale International Center for Finance for research support. All errors are our own.

1 Introduction

The dissemination of up-to-date knowledge is key for innovation and economic growth (Goldin and Katz, 2010; Jones, 2009). Higher education (HE) plays a central role in this process. Through the teaching of their curricula, HE programs facilitate human capital accumulation and nurture future innovators (Biasi, Deming, and Moser, 2020). These programs might differ, however, in their ability to equip students with up-to-date knowledge.¹ These differences can have important implications for labor market outcomes, education choices, and technological progress. Yet, they have so far remained unexplored; very little is known on how the content of HE varies across and within schools, how it is shaped, and how it relates to students' outcomes.

This paper brings together new data and a novel methodology to measure the extent to which HE courses cover frontier, i.e., recently produced, knowledge. Applying natural language processing (NLP) techniques to textual information on course syllabi (the content of HE courses) and academic publications (the frontier of knowledge), we build a novel metric: the education-innovation gap, designed to capture the distance between education content and the knowledge frontier. Specifically, we define the gap as a ratio of similarities of a course's content with knowledge from older vintages (covered by articles published decades ago) and new, frontier knowledge (covered by the most recent articles). For example, a Computer Science course that teaches *Visual Basic* (a relatively obsolete programming language) in 2020 would have a larger gap compared with a course that teaches *Julia* (a more recent programming language), because *Visual Basic* is mostly covered by old articles and *Julia* is mostly covered by recent articles.²

Using the education-innovation gap, we study the content of HE courses and provide four findings. First, HE courses differ greatly in their coverage of frontier knowledge, even conditioning on discipline and course level. Second, more selective and better funded institutions offer courses with lower gaps. These schools also enroll fewer disadvantaged students (Chetty et al., 2020), which implies that access to frontier knowledge is highly unequal. Third, instructors play a big role in shaping the content of their courses and research-active instructors teach more frontier knowledge, suggesting complementarities between teaching and research activities. Lastly, the dissemination of

¹Differences in HE programs attended have been associated with differences in earnings (Hoxby, 2020; Mountjoy and Hickman, 2020) and rates of invention (Bell et al., 2019).

²First released in 1991, *Visual Basic* is still supported by Microsoft in recent software frameworks, but the company announced in 2020 that it would not be further evolved (<https://visualstudiomagazine.com/articles/2020/03/12/vb-in-net-5.aspx>, retrieved 9/30/2020). *Julia* is a general-purpose language initially developed in 2009. Constantly updated, it is among the best for numerical analyses and computational science. As of July 2021 it was used at 1,500 universities, with over 29 million downloads and an 87 percent increase in a single year (<https://juliacomputing.com/blog/2021/08/newsletter-august/>, retrieved 9/30/2021).

frontier knowledge through HE courses is strongly and positively related to students' labor market outcomes and their ability to innovate in the future.

Our empirical analysis uses a novel source of information: the text of a sample of 1.7 million college and university syllabi, including about 540,000 courses taught at 800 four-year US institutions between 1998 and 2018. This sample represents about 5% of all courses taught in this time window, and it covers nearly all fields. While the sample over-represents courses from very selective schools, it is representative of the population in terms of fields, course levels (basic, advanced undergraduate, and graduate), and a broad set of school characteristics.

To construct the education-innovation gap, we start by calculating measures of textual similarity between each syllabus and the title, abstract, and keywords of over 20 million academic articles published in top academic journals since the journal's creation.³ Calculating pair-wise textual similarities involves three steps. First, we represent each document (a syllabus or an article) as a term frequency vector, projecting the text of the document on a comprehensive list of terms that refer to knowledge items. Each vector element is the frequency of a given term in the document, divided by the length of the document. Second, we use the "term-frequency-backward-inverse-document-frequency" (*TFBIDF*) approach (Kelly et al., 2021) to increase the importance of terms that are more informative of a document's content. This approach gives more weight to terms that are more "unique" to a document and de-emphasizes terms more commonly used across all documents. Third, we use these reweighted term frequency vectors to compute the cosine similarity between each syllabus and each article.

Armed with these cosine similarities, we construct the education-innovation gap of a given syllabus as the *ratio* of its average similarities with (a) older knowledge vintages, i.e., all articles published 13-15 years prior to the syllabus's date and (b) frontier knowledge, i.e., all articles published 1-3 years prior. Naturally, the gap is higher for syllabi that cover more older (rather than newer) knowledge. By virtue of being constructed as a ratio of similarities, the gap is not affected by idiosyncratic attributes of a syllabus such as length, structure, or writing style, which could introduce noise in the measurement of content. Moreover, the *TFBIDF*-adjustment implies that our method does not penalize syllabi for covering "classic" or "fundamental" knowledge. Terms pertaining to classic topics may belong to older knowledge vintages, but they are still widely taught; since they commonly appear across many documents, they receive a low weight.

A few empirical regularities confirm the ability of our measure to capture the distance between

³Previous works have used academic publications to capture the research frontier (for example, see Angrist et al., 2017, for economics research).

course content and the knowledge frontier. First, the gap is lower for syllabi that reference more recent articles and books in their lists of recommended readings. Second, the gap varies reasonably across course levels: It is the largest for basic undergraduate courses (taught in the first two years of a bachelor's degree and more likely to cover the fundamentals of a discipline) and smallest for graduate-level courses (master's and PhD). Third, using a simulation exercise, we show that gradually replacing "older" knowledge in a syllabus with "newer" knowledge (i.e., words most frequently appearing in old and new articles, respectively) progressively reduces the gap.

We begin by documenting significant differences in the gap across syllabi. To move a syllabus from the 25th to the 75th percentile of the gap distribution, approximately half of its content would have to be replaced with newer knowledge. Most of this variation occurs across courses and instructors, within field and course levels; a smaller share can be attributed to differences across fields and course levels. To account for these differences, the rest of our analysis compares syllabi within each field, course level, and year. The average syllabus in our data is more similar to newer than to older knowledge: Multiplying the gap by 100 for simplicity, its average equals 95.

Differences in the education-innovation gap across schools are useful to understand how the content of higher education is shaped. The gap is smaller in schools with a stronger focus on research (ranked as R1 in the Carnegie classification) and with more resources (higher endowment and spending on instruction and research). The gap is also smaller in more selective schools (for example the "Ivy-Plus," including the eight Ivy League colleges plus Stanford, MIT, Duke, and the University of Chicago) compared to non-selective schools. The magnitude of this difference is such that, in order to make the average syllabus in non-selective schools comparable to the average syllabus in an Ivy-Plus school, 8 percent of its content would have to be replaced with newer knowledge.

Importantly, differences across schools translate into disparities in access to up-to-date knowledge across students with different backgrounds. The education-innovation gap is significantly higher in schools enrolling students with lower median parental income and those with a higher share of Black or Hispanic students. This occurs because wealthier and more selective schools enroll more socio-economically advantaged students (Chetty et al., 2020).

In principle, part of these differences could be due to a "vertical differentiation" of educational content across schools. If students with greater ability enroll in more selective or better funded schools and are more capable of absorbing up-to-date content, cross-school differences in the gap might simply reflect schools' efforts to provide students with better tailored educational content.

We do not find evidence supporting this hypothesis: The negative correlation between the gap and parental income remains when we control for student ability, using the SAT and ACS scores of admitted students.

While the education-innovation gap varies significantly across schools with different characteristics, most of its variation (about a quarter) occurs within schools, across courses taught by different instructors. This can also be seen from the fact that the gap of the typical course remains stable over time, but it declines significantly when the instructor of a course changes.

Most higher-education instructors allocate their time and effort between teaching and research. As time is scarce, these tasks are often seen as competing (Hattie and Marsh, 1996; Courant and Turner, 2020). The nature of higher education, though, could also create complementarities between the two (Becker and Kennedy, 2005; Arnold, 2008). Our findings support the latter hypothesis. The education-innovation gap is significantly lower for courses taught by instructors who are more active in producing research (i.e., they publish more, are cited more, and receive more and larger grants). The gap is instead higher for non-ladder faculty, who specialize in teaching. The gap is also lower when the instructor's own research is closer to the topics of the course. These findings highlight that a proper deployment of faculty across courses can have important impacts on the content of education. They also suggest that investments in faculty research (both public, in the form of government grants, and institution-specific) can generate additional returns in the form of more updated instruction.

Our results so far unveil differences in the coverage of frontier knowledge across HE courses. Do these differences matter for the production of innovation and for students' outcomes? To answer this question, the ideal experiment would randomly allocate students to courses with different gaps. In the absence of this random variation, we settle on the more modest goal of characterizing the empirical relationship between the education-innovation gap and students' graduation rates, incomes, and measures of innovation, measured at the school level. In an attempt to account for students' selection into each school and other determinants of student outcomes related to instruction, we control for a large set of school observables such as institutional characteristics, expenditures, instructional characteristics, enrollment by demographic groups and major, selectivity, and parental background. We find that students in schools that offer courses with a lower gap are more likely to complete a PhD, produce more patents, and earn more after graduation. They are also more likely to graduate from college; a possible explanation is that taking more up-to-date courses makes students more motivated, and thus more likely to complete a program.

The education-innovation gap measures the academic content of each course. The richness of the information included in the syllabi allows us to go beyond knowledge and explore the skills students develop in each course. Recent works have highlighted the increasing importance of soft skills—non-cognitive attributes that shape the way people interact with others—for students’ success (Deming, 2017; Deming and Kahn, 2018). We measure the “soft-skills intensity” of each course as the extent to which evaluations are based on activities such as group projects, presentations, and surveys, which train soft skills. We find that courses with a lower education-innovation gap also tend to have a higher soft-skills intensity. More selective schools, those with more resources, and those serving more socio-economically advantaged students teach more soft-skills intensive courses. Within schools, research-active instructors are most likely to teach soft-skills intensive courses. Lastly, soft-skills intensity is strongly positively associated with student outcomes.

In the final part of the paper, we probe the robustness of our results to the use of alternative measures of knowledge composition. We consider three of them: The share of all “new” knowledge contained in a syllabus, designed not to penalize a syllabus that contains old and new knowledge compared with one that only contains new knowledge; a measure of “tail” knowledge, aimed at capturing the presence of the most recent content; and the education-innovation gap obtained using patent filings as a measure of frontier knowledge. All these alternative measures are significantly correlated with the gap, and our main results are qualitatively unchanged when we use them in lieu of the gap.

The main contribution of our paper is to document differences in the coverage of frontier knowledge across HE programs, a new and important dimension of heterogeneity. Analyzing the education-innovation gap, we shed new light on some of the most central questions related to innovation and higher education.

Several studies have characterized heterogeneity in the production of human capital, focusing on differences in the returns to educational attainment (Hanushek and Woessmann, 2012), majors and curricula (Altonji et al., 2012), college selectivity (Hoxby, 1998; Dale and Krueger, 2011), and the skill content of college majors (Hemelt et al., 2021; Li et al., 2021). Here, we take a novel approach: We directly examine curricula and educational content, among the most central components of higher education. With this approach, we document significant differences in the knowledge covered by each course, which could have important implications for students.

Our study is also related to the literature on education and the production of frontier knowledge and innovation. Earlier works (Nelson and Phelps, 1966; Benhabib and Spiegel, 2005) and

more recent ones (Akçigit et al., 2020; Bloom et al., 2021) have highlighted an important role for human capital and education—and education programs in particular—for the diffusion of ideas and technological advancements. Other studies have emphasized the importance of specific fields, such as STEM (Baumol, 2005; Toivanen and Väänänen, 2016; Bianchi and Giorcelli, 2019).⁴ Our findings highlight differences in the ability of HE programs to equip students with the knowledge necessary to innovate, which originate from heterogeneous course content. Importantly, these differences confirm a “lack of democratization” in access to valuable knowledge. US inventors have been shown to come from a small set of schools, enrolling very few low-income students (Bell et al., 2019). We find that these schools provide the most up-to-date educational content, which in turn suggests that access to frontier knowledge is highly unequal.

Lastly, we provide direct evidence on the importance of instructors in shaping the content of higher education. While some studies have found important effects on student outcomes (Hoffman and Oreopoulos, 2009; Carrell and West, 2010; Braga et al., 2016; Feld et al., 2020), much less is known on why and how instructors impact students (De Vlieger et al., 2020). We study instructors’ contribution to the production of educational content and carefully characterize differences across instructor types. Our findings also highlight complementarities between teaching and research activities.

2 Data

Our empirical analysis combines data from multiple sources. These include the text of course syllabi; the abstract of academic publications; job titles, publications, and grants of each instructor; characteristics of US higher education institutions; and labor market outcomes for the students at these institutions. More detail on the construction of our final data set can be found in [Appendix B](#).

2.1 College and University Course Syllabi

We obtained the raw text of a large sample of college and university syllabi from Open Syllabus (OSP), a non-profit organization that collects these data by crawling publicly accessible university and faculty websites to support educational research and applications. The initial sample contains more than seven million English-language syllabi of courses taught in over 80 countries between 1990 and 2018.

⁴The literature on the effects of education on innovation encompasses studies of the effects of the land grant college system (Kantor and Whalley, 2019; Andrews, 2017) and, more generally, of the establishment of research universities (Valero and Van Reenen, 2019) on patenting and economic activity. Education institutions also play a crucial role in fostering entrepreneurship (Tartari and Stern, 2021).

Most syllabi share a standard structure. The standard syllabus begins with basic details of the course (such as title, code, and the name of the instructor). It proceeds with a short description of its content, followed by a more detailed list of topics and required or recommended readings for each class session. Most syllabi contain information on evaluation criteria, such as assignments and exams; some also include general policies regarding grading, absences, lateness, and misconduct. Following this general structure, we parse each syllabus and extract four pieces of information: (i) basic course details, (ii) the course's content, (iii) the list of required and recommended readings, and (iv) a description of evaluation methods.⁵

Basic course details These include the name of the institution, the title and code of the course, the name of the instructor, as well as the quarter or semester and the academic year in which the course is taught. Course titles and codes allow us to classify each syllabus into one of three course levels: basic undergraduate, advanced undergraduate, or graduate. OSP assigns each syllabus to one of 69 detailed fields.⁶ We use this classification throughout the paper. For some tests, we further aggregate fields into four macro-fields: STEM, Humanities, Social Sciences, and Business.⁷

Course content We identify the portion of a syllabus that contains a description of the course's content by searching for section titles such as "Summary," "Description," and "Content."⁸ Typically, this portion describes the basic structure of the course, the key concepts that are covered, and (in many cases) a timeline of the content and the materials for each lecture.

Reference list We compile a list of bibliographic information for the required and recommended readings of each course by combining the list provided to us by OSP with all other in-text citations that we could find, such as "Biasi and Ma (2022)." We were able to compile a list of references for 71 percent of all syllabi. We then collect bibliographic information on each reference from Elsevier's SCOPUS database (described in more detail in Section 2.2); this includes title, abstract, journal, keywords (where available), and textbook edition (for textbooks).

Methods of evaluation To gather information on the methods used to evaluate students and the set of skills trained in the course, we use information on exams and other assignments. We identify and extract the relevant portion of the syllabus by searching for sections titled "Exam," "Assignment," "Homework," "Evaluation," and "Group." Using the text of these sections, we distinguish

⁵Angrist and Pischke (2017) use hand-coded syllabi from 38 universities to study the evolution of undergraduate econometrics classes.

⁶The field taxonomy used by OSP draws extensively from the 2010 Classification of Instructional Programs of the Integrated Postsecondary Education Data System, available at <https://nces.ed.gov/ipeds/cipcode/default.aspx?y=55>.

⁷Appendix Table B VII lists all 69 fields and shows the correspondence between fields and macro-fields.

⁸The full list of section titles used to identify each section is shown in Appendix Table B VI.

between hard skills (assessed through exams, homework, assignments, and problem sets) and soft skills (assessed through presentations, group projects, and teamwork). We were able to identify this information for 99.9 percent of all syllabi.

Sample restrictions and description To maximize consistency over time, we focus our attention on syllabi taught between 1998 and 2018 in four-year US institutions with at least one hundred syllabi in our sample.⁹ We exclude 35,917 syllabi (1.9 percent) with fewer than 20 words or more than 10,000 words (the top and bottom 1 percent of the length distribution).

Our final sample, described in panel (a) of Table 1, contains about 1.7 million syllabi of 542,251 courses at 767 institutions. Thirty-three percent of all syllabi cover STEM courses, ten percent cover Business, 30 percent cover Humanities, and 24 percent cover Social Science. Basic courses represent 39 percent of all syllabi and graduate courses represent 33 percent. A syllabus contains an average of 2,226 words in total, with a median of 1,778. Our textual analysis focuses on “knowledge” words, i.e., words that belong to a dictionary (see Section 3 for details). The average syllabus contains 420 unique knowledge words, with a median of 327.

2.2 Academic Publications

We use information from Elsevier’s SCOPUS database and compile the list of all peer-reviewed articles that appeared in the top academic journals of each field since the journal’s foundation.¹⁰ Top journals are defined as those ranked among the top 10 by Impact Factor (IF) in each field at least once since 1975 (or the journal’s creation, if it occurred after 1975).¹¹ Our final list of publications includes 20 million articles, corresponding to approximately 100,000 articles per year.¹²

Alternative measure of knowledge: Patents An alternative way to measure the knowledge frontier is to use the text of patents, rather than academic publications. To this purpose, we collected the text of more than six million patents issued since 1976 from the US Patents and Trading Office (USPTO) website.¹³ We capture the content of each patent with its abstract.

⁹For consistency and comparability, we removed 129,429 syllabi from one online-only university, the University of Maryland Global Campus.

¹⁰We accessed the SCOPUS data through the official API in April-August 2019.

¹¹Even if a journal appeared only once in the top 10, we collect all articles published since its foundation.

¹²SCOPUS classifies articles into 191 fields. To map each of these to the 69 syllabi fields, we calculate the cosine similarity (see Section 3) between each syllabus and each article. We then map each syllabi field to the SCOPUS field with the highest average similarity.

¹³Our web crawler collected the text content of all patents (in HTML format) from <http://patft.uspto.gov/netahtml/PTO/srchnum.htm>, with patent numbers ranging from 3850000 to 10279999.

2.3 Instructors: Research Productivity, Funding, and Job Titles

Nearly all course syllabi report the name of the course instructor. Using this information, we collected data on instructors' research productivity (publications and citations) and the receipt of public research funding. For a subset of instructors, we also collected information on job titles and annual salary.

Research Productivity Individual-level publications and citations data come from Microsoft Academic (MA). Discontinued at the end of 2021, MA was a search engine listing publications, working papers, other manuscripts, and patents for each researcher, together with citation counts for these documents. We linked MA records to syllabi instructors via fuzzy matching based on name and institution (details on this procedure are in [Appendix B](#)). We were able to successfully find 41 percent of all instructors, and we assume that the instructors we could not find never published an article (Table 1, panel (b)).

Using data from MA, we measure each instructor's research quantity and quality with the number of publications and received citations in the previous five years.¹⁴ On average, instructors published 6 articles in the previous five years, with a total of 172 citations (Table 1, panel (b)). The distributions of citation and publication counts are highly skewed: The median instructor in our sample only published one article in the previous five years and received no citations.

Funding We also collected information on government grants received by each researcher. Beyond research productivity, this information allows us to measure public investment in academic research. We focus on two of the main funding agencies of the U.S. government: the National Science Foundation (NSF) and the National Institute of Health (NIH).¹⁵ Our grant data include 480,633 NSF grants active between 1960 and 2022 (with an average size of \$582K in 2019 dollars) and 2,566,358 NIH grants active between 1978 and 2021 (with an average size of \$504K). We link grants to instructors via fuzzy matching between the name and institution of the investigator and those of the instructor (more details can be found in [Appendix B](#)). Eighteen percent of all syllabi instructors are linked to at least one grant; among these, the average instructor receives 10 grants, with a combined size of \$4,023K (Table 1, panel (b)).

¹⁴Using citations and publications in the previous five years helps address issues related to the life cycle of publications and citations, with older instructors having a higher number of citations and publications per year even if their productivity declines with time.

¹⁵These data are published by each agency, at <https://www.nsf.gov/awardsearch/download.jsp> and https://exporter.nih.gov/ExPORTER_Catalog.aspx. We accessed these data on May 25, 2021.

Job Titles In many US states, information on public college and university employees are disclosed online, to comply with state regulations on transparency and accountability. These records usually contain each employee’s name and job title. We were able to collect information on job titles for 35,178 instructors in our syllabi sample (10.6 percent of all instructors and 14.3 percent of public-sector instructors), employed in 490 public institutions in 16 states. On average, we observe instructors for two years (the modal year is 2017; we detail the coverage of these data in the [Appendix B](#)). Among all syllabi instructors for which we have job title information, 42 percent are ladder faculty (including 11 percent who are assistant professors, 13 percent who are associate professors, and 18 percent who are full professors; Appendix Figure [AI](#)).

2.4 Information on US Higher Education Institutions

The last component of our dataset includes information on all US colleges and universities of the syllabi in our data. Our primary source is the the Integrated Postsecondary Education Data System (IPEDS), maintained by the National Center for Education Statistics (NCES).¹⁶ For each school, IPEDS reports a set of institutional characteristics (such as name and address, sector, affiliation, and Carnegie classification); the types of degrees and programs offered; expenditure and endowment; characteristics of the student population, such as the distribution of SAT and ACT scores of all admitted students, enrollment figures for different demographic groups, completion rates, and graduation rates; and faculty composition (ladder and non-ladder). We linked each syllabus to the corresponding IPEDS record via a fuzzy matching algorithm based on school names. We were able to successfully link all syllabi in our sample.

We complement data from IPEDS with information on schools and students from three additional sources. The first one is the school-level dataset assembled and used by [Chetty et al. \(2020\)](#), which includes a school’s selectivity tier (defined using Barron’s scale), the incomes of students and parents, the number of patents obtained by all students, and a measure of intergenerational mobility (the share of students with parental income in the bottom quintile who reach the top income quintile as adults). These data are calculated using data on US tax records for a cross-section of cohorts who graduated between 2002 and 2004. The second is the Survey of Earned Doctorates, conducted by the NSF, which reports characteristics of all PhD receivers in US institutions each year. We use information on students’ graduating cohort and bachelor’s institution to construct the share of undergraduate students in each school and graduation year who eventually complete a PhD, for the

¹⁶IPEDS includes responses to surveys from all postsecondary institutions since 1993. Completing these surveys is mandatory for all institutions that participate, or apply to participate, in any federal financial assistance programs.

years 1998-2018.¹⁷ The third is the College Scorecard Database of the US Department of Education, an online tool designed to help users compare costs and returns of attending various colleges and universities in the US. This database reports the incomes of graduates ten years after the start of the program. We use these variables, available for the academic years 1997-98 to 2007-08, to measure student outcomes for each school.

Panel (c) of Table 1 summarizes the sample of colleges and universities for which we have syllabi data. On average, the median parental income of all students at each school is \$97,917. Across all schools, 3 percent of all students have parents with incomes in the top percentile. The share of minority students equals 0.22. Graduation rates average 61.4 percent in 2018, whereas students' incomes ten years after school entry, for the 2003-04 and 2004-05 cohorts, are equal to \$45,035. Students' average intergenerational mobility is equal to 0.29.

2.5 Data Coverage and Sample Selection

Our syllabi sample only covers a small fraction of all courses taught in US schools between 1998 and 2018. The number of syllabi increases over time, from 17,479 in 2000 to 68,792 in 2010 and 190,874 in 2018 (Appendix Figure AII).

To more accurately interpret our empirical results, it is crucial to establish patterns of sample selection. To do so, we compiled the full list of courses offered between 2010 and 2019 in a subsample of 161 US institutions (representative of all institutions included in IPEDS) using the course catalogs in the archives of each school.¹⁸ This allows us to compare our sample to the population of all courses for these schools and years.

This exercise does not reveal stark patterns of selection based on observables. The share of catalog courses covered by the syllabi sample remained stable over time, at 5 percent (Appendix Figure AIII). This suggests that, at least among the schools with catalog information, the increase in the number of syllabi over time is driven by an increase in the number of courses that are offered, rather than an increase in sample coverage. Our syllabi sample is also similar to the population in terms of field and course level composition. Between 2010 and 2018, STEM courses represent 33 percent of syllabi in our sample and 24 percent of courses in the catalog; Humanities represent

¹⁷The Survey of Earned Doctorates has been conducted since 1957. To assign a PhD recipient to their undergraduate institution, we use information on the institution where they obtained their bachelor's degree; to assign the recipient to a bachelor's degree cohort, we subtract 6 from their year of PhD receipt.

¹⁸We begin by randomly selecting 200 schools among all 4-year IPEDS institutions. Among these, we were able to compile course catalogs for 161 institutions, listed in Appendix Table AII. These look very similar in terms of observables to all schools in our sample (Appendix Table AIII). We focus our attention on years from 2010 onwards to maximize our coverage. For an example of a course catalogue, see <https://registrar.yale.edu/course-catalogs>.

30 and 32 percent, and Social Sciences represent 24 and 20 percent, respectively (Appendix Figure [AIV](#)). Similarly, basic undergraduate courses represent 39 percent of syllabi in our sample and 31 percent of courses in the catalog; advanced undergraduate courses represent 28 and 30 percent, and graduate courses represent 33 and 38 percent (Appendix Figure [AV](#)). These shares are fairly stable over time.

In addition, a school's portion of the catalog that is included in our sample and the change in this portion over time are unrelated to school observables. We show this in panel (a) of Table 2 (column 1), where we regress a school's share of courses included in our sample in 2018 on the following variables, one at the time and also measured in 2018: financial attributes (such as expenditure on instruction, endowment per capita, sticker price, and average salary of all faculty), enrollment, the share of students in different demographic categories (Black, Hispanic, alien), and the share of students graduating in Arts and Humanities, STEM, and the Social Sciences. We also test for the joint significance of all these variables. We find that these variables are individually and jointly uncorrelated with the share of courses in the syllabi sample, with an F-statistic close to one. In column 2 we repeat the same exercise, using the 2015-2018 change in the share of courses included in the syllabi as the dependent variable. Our conclusions are unchanged.

The only dimension in which our syllabi sample appears selected is school selectivity. Relative to non-selective institutions (for whom the share of courses in the sample is less than 0.1 percent), Ivy-Plus and Elite schools have a 2.4 percentage point higher share of courses included in the syllabi sample, and selective public schools have a 4.0 percentage point higher share. Taken together, these tests indicate that our syllabi sample does not appear to be selected on the basis of observable characteristics of schools and fields, although it does over-represent Ivy-Plus and Elite and selective public schools. By construction, though, we cannot test for selection based on unobservables. Our results should therefore be interpreted with this caveat in mind.

3 Measuring the Education-Innovation Gap

This section describes the construction of the education-innovation gap. We first explain how we measure similarities between course syllabi and academic publications. Then, we define and construct the gap using measures of similarity, implementing a series of adjustments to better describe each syllabus's content. Lastly, we validate our measure and describe its variation.

3.1 Measuring The Similarity Between Syllabi and Academic Publications

3.1.1 Constructing Term Frequency Vectors

We start by representing each document d (a syllabus or an article) as a term-frequency vector TF_d . Each element TF_{dw} of TF_d represents the frequency of term w in d :

$$TF_{dw} \equiv \frac{c_{dw}}{\sum_{k \in W} c_{dk}},$$

where, in the numerator, c_{dw} counts the number of times term w appears in d and the denominator is the total number of terms in d . To maximize our ability to capture the knowledge content of each document, we construct TF vectors focusing exclusively on terms related to knowledge concepts and skills, belonging to a dictionary W with $|W|$ terms (as a result, each term vector contains $|W|$ elements). Our primary dictionary is the list of all unique terms ever used as keywords in academic publications from the beginning of our publication sample until 2019.¹⁹ Appendix C contains details on the construction of term vectors and the use of a dictionary.

3.1.2 Adjusting for Term Relevance

When constructing similarity metrics, it is crucial to ensure that each term receives a weight proportional to its importance in capturing a document’s content. TF vectors give more weight to terms with a higher document frequency. However, terms that are very common across *all* documents receive more weight regardless of their ability to capture the content of a given document. For example, holding term frequency fixed, terms such as “Programming” or “Animals” – very common among Computer Science and Biology syllabi, respectively – are usually less informative of content than terms such as “Natural Language Processing” or “CRISPR.”²⁰

To this purpose, we use a leading approach in the text analysis literature called “term-frequency-inverse-document-frequency” (TFIDF, Kelly et al., 2021). This approach assigns each term a weight inversely proportional to the frequency of the term across all documents, underweighting terms that are not diagnostic of a document’s content.

We implement this approach by constructing an inverse-document frequency vector IDF (of

¹⁹We have also used the list of all terms that have an English Wikipedia webpage as of 2019. Our results are robust to this choice.

²⁰Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) is a family of DNA sequences found in the genomes of prokaryotic organisms such as bacteria and archaea. The term also refers to a recent technology that can be used to edit genes.

length $|W|$) with elements defined as

$$IDF_w \equiv \ln \left(\frac{|D|}{\sum_{n \in D} \mathbb{1}(c_{nw} > 0)} \right),$$

where D is the set of all documents (syllabi *and* articles). The denominator in parentheses is the total number of documents that contain word w . IDF_w is thus the inverse of the share of all documents containing word w . Using IDF , we can then transform TF_d into a term-frequency-inverse-document-frequency vector $TFIDF_d$, with elements equal to

$$TFIDF_{dw} = TF_{dw} \times IDF_w. \quad (1)$$

Accounting for Changes in Term Relevance Over Time The weighting approach described so far calculates the relative importance of each term for a given document pooling together documents published in different years. This is not ideal for our analysis, because we are interested in the novelty of the content of a syllabus d relative to research published in the years *prior* to d . Consider, for example, course CS229 at Stanford University, taught by Andrew Ng in the early 2000s and one of the first that entirely focused on *Machine Learning*. This term has become very popular in later years, so its frequency across all documents is very high and its IDF_w very low. Pooling together documents from different years would thus result in a very low $TFIDF_{dw}$ for the term “machine learning” in the course’s syllabus. Not accounting for changes in term frequency over time would then lead us to severely mischaracterize the course’s path-breaking content.

To overcome this issue, we modify the traditional $TFIDF$ and construct a retrospective or “point-in-time” version of IDF , meant to capture the inverse frequency of a term among all documents published *prior to* d . We call this vector “backward- IDF ,” or $BIDF_t$. It is indexed by t because it varies over time. We define the set of documents published prior to t as D_t ; the elements of $BIDF_t$ can be defined as

$$BIDF_{tw} \equiv \log \left(\frac{|D_t|}{\sum_{n \in D_t} \mathbb{1}(c_{nw} > 0)} \right).$$

The use of this weighting approach allows us to give a temporally appropriate weight to each term in a document. Using $BIDF_t$, we can then calculate a “backward” version of $TFIDF_d$ —called $TFBIDF_d$ —whose elements are

$$TFBIDF_{dw} = TF_{dw} \times BIDF_{t(d)w}, \quad (2)$$

where $t(d)$ is the publication year of document d .

3.1.3 Building Textual Similarities Between Syllabi and Articles

Armed with weighted term vectors, we can now construct measures of textual similarities between syllabi and articles. For simplicity, we denote $TFBIDF_d$ as V_d for each d . The measure of similarity we use is the cosine similarity, defined for two documents d and d' as

$$\rho_{d,d'} = \frac{V_d \cdot V_{d'}}{\|V_d\| \|V_{d'}\|} \quad (3)$$

where $\|V_d\|$ is the Euclidean norm of V_d . Since each element of V_d is non-negative, ρ lies in the interval $[0, 1]$. If d and d' use the exact same set of terms with the same frequency, $\rho_{d,d'} = 1$; if they have no terms in common, $\rho_{d,d'} = 0$.

3.2 Calculating the Education-Innovation Gap

We capture the similarity between each syllabus d and different vintages of knowledge using the average similarity of d with all the articles published in a three-year time period ending τ years before $t(d)$:

$$S_d^\tau = \frac{\sum_{n \in \Omega_\tau(d)} \rho_{dn}}{|\Omega_\tau(d)|}$$

where ρ_{dk} is the cosine similarity between syllabus d and an article k , $\Omega_\tau(d)$ is the set of all articles published in the three-year time interval $[t(d) - \tau - 2, t(d) - \tau]$, and $|\Omega_\tau(d)|$ is the total number of these articles.²¹

We construct the education-innovation gap as the ratio between the average similarity of a syllabus with older technologies (published in τ) and the similarity with more recent ones ($\tau' < \tau$):

$$Gap_d \equiv \left(\frac{S_d^\tau}{S_d^{\tau'}} \right) \quad (4)$$

Given this definition, the syllabus of a course taught in t has a lower education-innovation gap if its text is more similar to more recent research (published in $t - \tau'$) than to older research (published in $t - \tau$). For our analysis, we set $\tau = 13$ ($[t - 15, t - 13]$ vintage) and $\tau' = 1$ ($[t - 3, t - 1]$ vintage). We multiply the gap by 100 for readability.

Our measure features two attractive properties. First, being constructed as a ratio, the gap is not affected by syllabus-specific attributes such as style, format, or length, which could introduce

²¹Our main analysis uses three-year intervals; our results are robust to the use of one-year or two-year intervals.

noise in the ability of a simple measure of similarity to capture content. For example, a course could have a higher similarity with existing research compared with another course covering the same material, if the syllabus of the former is longer or uses more academic terms. We illustrate this point with a simulation exercise in [Appendix C](#).²²

Second, our measure does not heavily penalize syllabi for covering “classic” topics in the literature, as long as these are widespread across courses. This is guaranteed by the use of a *TFBIDF* approach, which reduces the impact on the gap of terms—such as classics—frequently used across all documents. For example, the term “Ordinary Least Squares” (“OLS”) refers to a relatively old but very common concept taught in most econometrics and statistics courses. As such, it will receive a low weight and syllabi will not be penalized much by covering it.

3.3 Validating The Measure and Interpreting Its Magnitude

We perform a series of tests to validate our measure’s ability to capture the distance between the content of a course and the research frontier. First, we show that the relationship between the gap and the average age of its reference list (defined as the average difference between the year of each syllabus and the publication year of each reference) is positive and significant (Figure 1, panel (a)). While the average reference age is easy to calculate, our text-based measure is available for all syllabi (including those for which the reference list is unavailable) and is more accurate in capturing the content of courses that only rely on very few bibliographic sources (for example, a textbook).

Second, we show that the gap varies reasonably across course levels. More advanced undergraduate courses and graduate-level courses have lower gaps compared with basic undergraduate courses. Controlling for field-by-year effects, basic undergraduate courses have a gap of 95.7; advanced undergraduate courses have a gap of 95.3, and graduate courses have a gap of 94.7 (Figure 1, panel (b)). This confirms the intuition that more advanced courses cover content that is closer to the knowledge frontier.

Third, we use a simulation exercise to confirm that our measure is able to pick up changes in a syllabus’s distance to different knowledge vintages as we change its textual content. Specifically, we randomly draw a subsample of 100,000 syllabi. Then, we progressively replace terms that are more frequent in older knowledge vintages (“old words”) with terms more frequent in newer vintage (“new words”), and we re-calculate the gap as we replace more words. Old words are those in the

²²We manually create a sample of 1.7 million syllabi as sets of dictionary terms, for which we know ex ante the ratio between “old” knowledge terms (more popular among old publications) and “new” knowledge terms (most popular among recent publications). The education-innovation gap performs much better as a measure of this ratio than a simple measure of similarity with new terms ([Appendix C](#)).

top 5 percent in terms of frequency in the old publication corpus between $t - 15$ and $t - 13$ or in the old publication corpus between $t - 15$ and $t - 13$ but not in the new publication corpus between $t - 3$ and $t - 1$; new words as those in the top 5 percent in terms of frequency in the new publication corpus between $t - 3$ and $t - 1$ or in the new publication corpus between $t - 3$ and $t - 1$ but not in the old publication corpus between $t - 15$ and $t - 13$. The gap monotonically decreases as we replace more old words with new ones (Figure 1, panel (c)). This simulation is also useful for gauging the economic magnitude of changes in the gap. In particular, a unit change in the gap is equivalent to the replacement of 10 percent of a syllabus's old words (or 34 old words, out of 330 words for the median syllabus).

3.4 The Education-Innovation Gap: Variation and Variance Decomposition

The average course has a gap of 95.3, with a standard deviation of 5.8, a 25th percentile of 91.6, and a 75th percentile of 98.8 (Table 1, panel (a) and Appendix Figure AVI). To give an economic meaning to this variation, we use the relationship illustrated in panel (c) of Figure 1. In order to move a syllabus from the 75th to the 25th percentile of the distribution (a 7.2 change in the gap) we would have to replace approximately 200 of its words, or 60 percent of the content of the median syllabus.

To better understand what drives variations in the gap, we perform a Shapley-Owen decomposition (Israeli, 2007) of its variance into five sets of factors: year, field, school, course, and instructor. For each factor j , we calculate the partial R^2 as

$$R_j^2 = \sum_{k \neq j} \frac{R^2 - R^2(-j)}{K!/j!(K-j-1)!}$$

where $R^2(-j)$ is the adjusted R^2 of a regression that excludes fixed effects for all factors except j ; this quantity captures the share of the variation captured by factor j .²³

This exercise indicates that differences across fields explain 4 percent of the total variation in the gap, while differences across schools explain 2 percent (Table 3, column 1). Courses explain a large 33 percent, indicating a great deal of persistence in the content of a course over time. Importantly, differences across instructors explain a large 25 percent. Results are similar when we replace courses with course levels; the latter explain less than 1 percent of the total variation (column 2).

²³We use adjusted R^2 throughout to account for the large number of fixed effects in the model.

4 The Education-Innovation Gap Across Schools

Our analysis starts by examining cross-school differences in educational content.

4.1 School Characteristics

We begin by testing how the education-innovation-gap relates to three sets of school attributes: (i) institutional, such as sector (public or private), research intensity (distinguishing between schools classified as R1 – “Very High Research Intensity” – according to the Carnegie classification, and all other schools) and emphasis on liberal arts and sciences relative to other subjects (distinguishing between Liberal Arts Colleges (LAC) and all other schools); (ii) financial, such as endowment and spending on instruction, faculty salaries, and research; and (iii) faculty composition and productivity, such as the share of non-ladder faculty, the share of tenure-track (non-tenured) faculty, and the number of academic publications per faculty.

We estimate pairwise correlations (captured by β in the following equation) between the gap and these attributes controlling for field, course level, and year of the syllabus:

$$\text{Gap}_i = \beta X_{s(i)} + \phi_{f(i)l(i)t(i)} + \varepsilon_i \quad (5)$$

where Gap_i measures the education-innovation gap of syllabus i , taught in school $s(i)$ and year $t(i)$; the variable X_s is the institutional characteristic of interest in school s ; and field-by-level-by-year fixed effects ϕ_{flt} control for systematic differences in the gap, common to all syllabi in the same field (f) and course level (l), that vary over time (t). We cluster standard errors at the institution level.

Institutional and financial characteristics Estimates of β for each school characteristic are shown in Figure 2. Public schools have a slightly larger gap compared with non-public schools, but this difference is indistinguishable from zero. No differences emerge between LACs and other schools. R1 schools have a 0.2 smaller gap compared with schools with a lower research intensity.

In order to quantify the economic magnitude of these differences, we can use the simulation results in Figure 1 (panel (c)). In order to close the difference in the gap between R1 and other schools, we would have to replace approximately 2 percent of the knowledge content of the median syllabus (7 terms). The difference between R1 and other schools, although significant, is thus quite small.

A statistically and economically significant relationship exists between the gap and financial characteristics, such as endowment and spending on instruction, faculty salary, and research. For

example, a 10-percent increase in instructional spending is associated with a 3.5 lower gap, or a 35 percent change in the syllabus; a 10-percent increase in research spending is associated with a unit lower gap or a 10 percent change in the syllabus.

Selectivity Next, we test whether the gap differs across schools with different selectivity. Following [Chetty et al. \(2020\)](#), we bin schools in four “tiers” according to their selectivity in admissions, measured with Barron’s 2009 ranking. “Ivy Plus” include Ivy League universities and the University of Chicago, Stanford, MIT, and Duke. “Elite” schools are all the other schools classified as tier 1 in Barron’s ranking. “Highly selective” schools include those in tiers 2 and 3, while “Selective” schools are those in tiers 4 and 5. Lastly, “Non-selective” schools include those in Barron’s tier 9 and all four-year institutions not included in Barron’s classification.

To compare the gap across different school tiers, we use the following equation:

$$\text{Gap}_i = \mathbf{S}'_i \boldsymbol{\beta} + \phi_{f(i)l(i)t(i)} + \varepsilon_i$$

where the vector \mathbf{S}'_i contains indicators for selectivity tiers (we omit non-selective schools), and everything is as before.

Point estimates of the coefficients vector $\boldsymbol{\beta}$ in equation (6), shown as diamonds in [Figure 2](#), indicate that more selective schools offer content that is closer to the research frontier. Ivy Plus and Elite schools have the smallest gap, 0.84 smaller than non-selective schools (corresponding to an 8 percent difference in the median syllabus). Highly selective schools have a 0.67 smaller gap and selective schools have a 0.51 smaller gap (5 percent). A possible interpretation for these differences is that more selective schools offer higher-quality education. However, if higher-ability students are better able to absorb frontier knowledge, another possibility is that schools tailor instruction to the abilities of their students. We attempt to test this hypothesis in the next section and in [Section 6](#), where we relate the education-innovation gap to student outcomes.

4.2 Students’ Characteristics

Schools with different characteristics serve different populations of students; for example, Ivy-Plus and Elite schools are disproportionately more likely to enroll students from wealthier backgrounds ([Chetty et al., 2020](#)). Cross-school differences might therefore translate into significant disparities in access to up-to-date knowledge among students with different backgrounds. Here, we focus on two dimensions of socio-economic background: parental income and race and ethnicity.

Parental income We re-estimate equation (5) using two measures of parental income as the explanatory variable: median parental income and the share of parents with incomes in the top percentile of the national distribution, constructed using tax returns for the years 1996 to 2004 (Chetty et al., 2020). These estimates, shown as the full triangles in Figure 2, indicate that schools serving more economically disadvantaged students offer courses with a higher gap. Specifically, a one-percent higher median parental income is associated with a 0.56 lower gap, which corresponds to a 5 percent difference in the median syllabus. Similarly, a 10-percentage point higher share of students with parental income in the top percentile is associated with a 0.42 lower gap (4 percent).

In principle, part of these differences could be due to a “vertical differentiation” of educational content across schools. If students with greater ability are better able to absorb more up-to-date content, cross-school differences in the gap might reflect schools’ efforts to provide students with appropriate educational content. Our data, however, do not support this hypothesis. Controlling for the average SAT score of students admitted at each school as a proxy for their ability yields only slightly smaller estimates compared with the baseline (Figure 2, hollow triangles). This rules out vertical differentiation as an explanation for cross-school differences in the gap.

Students’ race and ethnicity Schools that enroll a higher share of minority students (Black or Hispanic) also offer courses with a higher gap. Using the share of minority students as the explanatory variable in equation (5) reveals that a one-percentage point higher share is associated with a 0.58 higher gap, equivalent to a 6 percent change in the average syllabus. As before, this relationship holds (but is less precise) if we control for average student ability.

In line with existing evidence on disparities in access to selective schools among more and less advantaged students, our results document a new dimension of inequality: That in access to educational content that is close to the research frontier. Importantly, this inequality cannot be explained by differences in student ability.

5 The Role of Instructors

Instructors are considered one of the most important inputs for the production of student learning, and one of the most costly (De Vlieger, Jacob, and Stange, 2020). In line with this, our data show that most of the variation in the gap occurs within schools and across courses taught by different people. We now investigate in depth the role of instructors and their characteristics in shaping the content of higher education.

5.1 Persistency In A Course’s Content Over Time and Changes in Instructors

To understand how instructors shape the content of the courses they teach, we start by studying how the education-innovation gap of a course varies when the course instructor changes. We estimate an event study of the gap in a $[-4, 4]$ year window around the time of an instructor change:

$$\text{Gap}_i = \sum_{k=-4}^4 \delta_k \mathbb{1}(t(i) - T_{c(i)} = k) + \gamma_{c(i)} + \phi_{f(i)t(i)} + \varepsilon_i, \quad (6)$$

where i , f , and t denote a syllabus, field, and year respectively. The subscript c denotes a specific course within each school (for example, Econ 101 at Yale University); the variable T_c represents the first year in our sample in which the instructor of course c changes.²⁴ To more precisely capture the impact of an instructor change, we restrict our attention to courses taught by a maximum of two instructors in each year and set the indicator function to zero for all courses without an instructor change, which serve as the comparison group. We cluster standard errors at the course level. Assuming $\delta_0 = 0$, the parameters δ_k capture the differences between the gap k years after an instructor change relative to the year preceding the change.

OLS estimates of δ_k , shown in Figure 3, indicate that a change in a course’s instructor is associated with a sudden decline in the education-innovation gap. Estimates are indistinguishable from zero and on a flat trend in the years leading to an instructor change; the year of the change, the gap declines by 0.1. This decline is equivalent to replacing 2 percent of the content of a syllabus.

In Table 4 we re-estimate equation (6) for different subsamples of syllabi, pooling together years preceding and following an instructor change. After a change, the gap declines for all fields and course levels by about 0.1 on average (2 percent of a course’s content, column 1, significant at 1 percent). The decline is largest for Humanities and STEM courses (-0.14 and -0.11, columns 3 and 4, respectively), as well as for graduate courses (-0.12, column 8).

These results indicate that course updating is not a gradual process taking place over time. Instructors who teach the same course for many years tend to leave content unchanged. Instead, those who take over a course from someone else significantly update its content, bringing it closer to the knowledge frontier. Our findings also confirm that instructors play a crucial role in shaping the content of the courses they teach, particularly for advanced courses.

²⁴Our results are robust to using the median or the last year with an instructor change.

5.2 The Education-Innovation Gap and Instructors' Characteristics

The decline in the gap that follows an instructor change, though, could mask substantial differences across instructors. For example, the decline could differ for more research-active instructors, who spend less time teaching but are better informed on the frontier of knowledge. Similarly, the decline could depend on whether the new instructor is an expert on the topics covered by the course. We explore these possibilities next.

Ladder vs non-ladder faculty Ladder (i.e., tenure-track or tenured) faculty are generally more focused on research compared with non-ladder faculty, whose primary job is to teach. In recent years, universities have started to increasingly rely on non-ladder faculty to meet a rapid rise in enrollment (Goolsbee and Syverson, 2019).²⁵ Ex ante, whether one type of faculty or the other would be better at teaching up-to-date content is ambiguous. On the one hand, being specialized on teaching, non-ladder faculty might be better at keeping educational content updated. On the other hand, being better informed on frontier knowledge, ladder faculty might be more likely to include this knowledge in the courses they teach.

Comparing the education-innovation gap across job titles and controlling by field-by-course level-by-year effects, we find that non-ladder faculty (adjunct professors) have the largest gap, at 95.8 (Figure 4). Tenure-track assistant professors, on the other hand, have the lowest gap at 95. The difference between assistant and adjunct professors is equivalent to 7 percent of a syllabus's content.

Notably, the gap is almost as high for full (tenured) professors as it is for adjuncts, at 95.6. Associate professors have a slightly smaller gap than full at 95.5, but still significantly higher than assistant professors. Junior faculty on the tenure track thus appear to teach the courses with the most updated content.

Research productivity One possible explanation for these results is that assistant professors are more recently trained and face stronger incentives to be active in research. This might make them more informed about the knowledge frontier. We test this hypothesis directly by exploring the relationship between a course's gap and the research productivity of the instructor, measured using individual counts of citations and publications in the previous five years. We estimate the following

²⁵Employing non-ladder faculty makes it easier (and cheaper) for schools to face increases in enrollment. Colleges have monopsony power on tenure-track, but not ladder, faculty; the latter earn substantially lower wages and have a much higher elasticity of labor supply. This implies that, when enrollment increases, schools are better off hiring more non-ladder faculty to avoid increasing wages for tenure-track faculty (Goolsbee and Syverson, 2019)

equation:

$$\text{Gap}_i = \sum_{n=1}^4 \beta^n q_k^n(it(i)) + \gamma_{c(i)} + \psi_{f(i)t(i)} + \varepsilon_i \quad (7)$$

where q_k^n equals one if instructor k 's measure of research productivity (publications or citations) is in the n th quartile of the distribution (the omitted category is courses with instructors whose measure k equals zero). Course fixed effects $\gamma_{c(i)}$ and field-by-year fixed effects $\psi_{f(i)t(i)}$ control for unobserved determinants of the gap that are specific to a course in a given field and year. Estimates of β^n capture the difference in the gap between courses taught by faculty with productivity in the n th quartile and those taught by faculty with no citations or publications and are identified out of changes in instructors for the same course over time.

Estimates of β^n , shown in Table 5, indicate that the gap progressively declines as the research productivity of the instructor grows. In particular, a switch from an instructor without publications to one with a number of publications in the top quartile of the field distribution is associated with a 0.11 decline in gap (equivalent to updating 2 percent of a course's syllabus; Table 5, panel (a), column 1, significant at 1 percent). Similarly, a switch from an instructor without citations to one with citations in the top quartile is associated with a 0.06 lower gap (panel (b), column 1, significant at 5 percent). These relationships are stronger for Social Sciences courses (column 5) and for courses at the graduate level (column 8).²⁶

Fit with the course A natural explanation for this finding is that research-active instructors are better informed about the research frontier. If this is the case, we should expect the relationship between productivity and the gap be stronger for courses whose topics are more similar to the instructor's own research. To test for this possibility, we construct a measure of "fit" between the course and the instructor's research. This measure is defined as the cosine similarity between the instructor's research in the previous 5 years and the most updated course on the same topic across *all* schools (for example, Introductory Econometrics).²⁷ We then correlate this measure with the education-innovation gap, controlling for course and field-by-year fixed effects (as in equation 7). Estimates of this relationship indicate that a one-standard deviation higher instructor-course fit is associated with a 0.09 lower gap (Table 6, significant at 5 percent). This relationship is particularly

²⁶Panels (a) and (b) of Appendix Figure AVII show a binned scatterplot of the gap and either citations (panel (a)) or publications (panel (c)) in the prior 5 years, controlling for field effects. In this figure, the horizontal axis corresponds to quantiles of each productivity measure; the vertical axis shows the average gap in each quantile.

²⁷One attractive property of this measure is that it does not uniquely reflect the content of the syllabus itself, which is of course directly shaped by the instructor; rather, it aims at capturing the content of all courses on the same topic. Constructing this measure requires obtaining a unique identifier for courses on the same field or topic (e.g. Machine Learning) across schools. We describe the procedure we use to do this in Appendix B.

strong for STEM and Social Sciences courses (column 4) and for courses at the advanced undergraduate level (column 7).

Research funding In Table 7, we use data on the number of NSF and NIH grants received by each instructor and test whether the same relationship holds for research inputs, such as government grants; as before, we control for course and field-by-year effects. A switch from an instructor who never received a grant to one with at least one grant is associated with a 0.05 reduction in the gap (column 1, significant at 5 percent).²⁸ This suggests that public investments in academic research can yield additional private and social returns in the form of more up-to-date instruction.²⁹

Taken together, these findings indicate that instructors play a crucial role in shaping the content of the courses they teach. We also document some complementarities between research and teaching: Research-active instructors are more likely to cover frontier knowledge in their courses, especially when teaching advanced courses and courses closest in topic to their own research agenda. Our results suggest that a proper deployment of faculty across courses can have important impacts on the content of education, and that investments in faculty research (both public, in the form of government grants, and institution-specific) can generate additional returns in the form of more updated instruction.

6 The Education-Innovation Gap and Students' Outcomes

Significant differences in access to up-to-date knowledge exist both across and within schools, and across courses taught by different people. Do these differences matter for students' outcomes and for the production of innovation? To begin answering this question, we now explore the relationship between the gap and a) innovation measures, such as a school's share of undergraduate students who complete a PhD and the number of patents produced by all students; and b) labor-market measures, such as graduation rates, income, and intergenerational mobility.

All these outcomes are measured at the school level or at the school-by-cohort level (with the exception of the share of students who attend graduate school, also available by macro-field). The education-innovation gap is measured at the syllabus level. To construct a school-level measure, we follow the school value-added literature (Deming, 2014) and estimate the school component of

²⁸A binned scatterplot reveals a negative relationship between the gap and the number of NSF and NIH grants (Appendix Figure AVII, panel d).

²⁹For a review of the role of grant funding as a tool to promote innovation, see Azoulay and Li (2020).

the gap using the following model:

$$\text{Gap}_i = \theta_{s(i)} + \phi_{f(i)l(i)t(i)} + \varepsilon_i. \quad (8)$$

In this equation, the quantity θ_s captures the school component of the education-innovation gap for school s , accounting for flexible time trends that are specific to the level l and field f of the course. Because outcome measures refer to students who complete undergraduate programs at each school, we construct θ_s using only undergraduate syllabi; our results are robust to the use of all syllabi. Appendix Figure AX shows the distribution of θ_s ; its standard deviation is 0.85, corresponding to a 5 percent change in the average syllabus.

In the remainder of this section, we present estimates of the parameter δ in the following equation:

$$Y_{st} = \delta \hat{\theta}_s + X_{st}\gamma + \tau_t + \varepsilon_{st} \quad (9)$$

where Y_{st} is the outcome for students who graduated from school s in year t ; $\hat{\theta}_s$ is the school-level component of the gap (estimated from equation (8) and standardized to have mean zero and variance one); X_{st} is a vector of school observables; and τ_t are year fixed effects. We calculate bootstrapped standard errors, clustered at the level of the school, to account for the fact that $\hat{\theta}_s$ is an estimated quantity.

It should be stressed that the parameter δ does not necessarily capture the causal effect of the gap on outcomes. There might be school and student attributes related to both the content of a school's courses and student outcomes. To account for as many of these attributes as is possible, we control for a rich set of school observables from IPEDS and show how baseline estimates change when we implement this strategy. We include seven groups of controls, including institutional characteristics (private-public, selectivity tiers, and an interaction between selectivity tiers and an indicator for R1 institutions according to the Carnegie classification); instructional characteristics (student-to-faculty ratio and the share of ladder faculty); financials (total expenditure, research expenditure, instructional expenditure, and salary instructional expenditure per student); enrollment (share of undergraduate and graduate enrollment, share of white and minority students); selectivity (indicator for institutions with admission share equal to 100, median SAT and ACT scores of admitted students in 2006, indicators for schools not using either SAT or ACT in admission); major composition (share of students with majors in Arts and Humanities, Business, Health, Public and Social Service, Social Sciences, STEM, and multi-disciplinary fields); and family background, measured as

the natural logarithm of median parental income.

6.1 Innovation Measures

Obtaining a PhD We begin by studying the relationship between the gap and the share of students who obtain a PhD. We construct this variable using data from the NSF Survey of Earned Doctorates, separately for five macro-fields: STEM, Health, Business, Social Science, and Humanities. To match the level of aggregation of this variable, we aggregate the education-innovation gap at the school-by-macro field level (rather than just at the school level) and modify equation (9) slightly so that one observation in our data is a school-by-macro field in a year. In column 1 of Table 8 (panel (a)) we pool data across macro-fields. The unconditional correlation between the gap and the share of students who obtain a PhD is negative and statistically significant: A one-standard deviation lower gap is associated with a 0.4 percentage point higher share, or 17 percent compared with an average of 0.0265 percent. The correlation is particularly strong for Social Science (-0.0124) and Health (-0.0074), while it is small and indistinguishable from zero for STEM, Business, and Humanities. These correlations remain remarkably robust when we control for school characteristics (Table 8, panel (b)).

Invention Next, we test whether students at schools which offer courses with a lower gap produce more inventions later in their life, in the form of patents. We do so by substituting the total number of patents received after graduation by students at each school as the outcome in equation (9). Unconditionally, a one-standard deviation decline in the gap is associated with 27 additional patents at a given school, or 20 percent compared with an average of 131 patents (Table 8, panel (b), column 8, p-value equal to 0.11). The relationship remains robust and even becomes more precise controlling for school observables (Table 8, panel (b)).

6.2 Labor Market Outcomes

Graduation rates Next, we examine the relationship between the education-innovation gap and labor market outcomes. We begin with graduation rates, an outcome that immediately precedes entry in the labor market; graduation is in part also a function of choices made by the students, which could be impacted by the content of the courses they took.

Column 1 of Table 9 shows the relationship between the gap (measured in standard deviations) and graduation rates. An estimate of -0.05 in panel (a), significant at 1 percent, indicates that a one-standard deviation decline in the gap (or a 10 percent change in the content of a syllabus) is

associated with a 5 percentage point higher graduation rate. Compared with an average of 0.61, this corresponds to a 8 percent increase in graduation rate.

The estimate of δ declines as we control for observable school characteristics, indicating that part of this correlation can be explained by other differences across schools. However, it remains negative and significant at -0.007, indicating that that a one-standard deviation reduction in the gap is associated to a 1.1 percent increase in graduation rates (panel (b), column 1, significant at 5 percent).

Students' income and intergenerational mobility Graduation rates are a strictly academic measure of student success; however, they are also likely to affect students' long-run economic trajectories. To directly examine the relationship between the education-innovation gap and students' economic success after they leave college, in columns 2-8 of Table 9 we study the relationship between the gap and various income statistics.

Column 2 shows estimates on the natural logarithm of mean student income from the College Scorecard. While imprecise, this estimate indicates that a one-standard deviation in the gap is associated with a 0.7 percent increase in income controlling for the full set of observables (panel (b), p-value equal to 0.17). The College Scorecard also reports mean incomes for students with parental incomes in the bottom tercile of the distribution; for these students, the relationship is slightly larger at 0.8 percent (column 3, significant at 10 percent). Estimates are largely unchanged when we use median instead of mean income (column 4).

Information on mean student incomes at the school level is also reported by [Chetty et al. \(2020\)](#), calculated using tax records for a cross section of students. Unconditional estimates (which omit year effects due to the cross-sectional structure of the data) indicate that a one-standard deviation in the gap is associated with a 7 percent increase in students' mean income (panel (a), column 5, significant at 1 percent). This estimate is smaller, at 1.4 percent, when controlling for institutional characteristics (panel (b), column 5, significant at 1 percent).

In columns 6 through 8 of Table 9 we investigate the relationship between the gap and the probability that students' incomes reach the top echelons of the distribution. Estimates with the full set of controls indicate that a one-standard deviation decline in the gap is associated with a 0.84 percentage-point increase in the probability of reaching the top 20 percent (2.2 percent, panel (b), column 6, significant at 1 percent), a 0.53 percentage-point increase in the probability of reaching the top 10 percent (2.5 percent, column 7, significant at 5 percent), and a 0.31 percentage-point increase in the probability of reaching the top 5 percent (2.7 percent, column 8, significant at 10 percent).

Taken together, these results indicate a positive relationship between the school-level education-innovation gap and students' average and top incomes.

Lastly, in column 9 of Table 9 we study the association between the gap and intergenerational mobility. The unconditional correlation between these two variables is equal to -0.0293 , indicating that a one-standard deviation lower gap is associated with a 2.9 percentage-points increase in intergenerational mobility (9.9 percent, panel (a), column 9, significant at 1 percent). However, this correlation becomes smaller and indistinguishable from zero when we control for school observables, reaching -0.0047 when we include the full set of controls (column 9, panel (b), p-value equal to 0.15).

Summary Our findings indicate that a lower education-innovation gap at the school level is associated with more innovation and improved academic and economic student outcomes. The lack of experimental variation in the gap across schools prevents us from estimating a causal relationship. Yet, our results are robust to the inclusion of controls for a large set of school and student characteristics, indicating that these correlations are unlikely to be driven by cross-school differences in spending, selectivity, major composition, or parental background. These findings point to the potentially important role of up-to-date instruction in determining future innovation levels and the outcomes of students as they exit college and enter the labor market.

7 Soft Skills in Course Content

By definition, the education-innovation gap focuses on the novelty of a syllabus with respect to its *academic* content. Recent works have shown, though, how content might not be the only thing that matters for students; they have instead highlighted the importance of *skills* for students' later life outcomes. In particular, soft skills—defined as non-cognitive abilities that define how a person interacts with their colleagues and peers—are increasingly in high demand in the labor market and associated with more favorable outcomes (Deming, 2017).

Supported by this evidence, we now examine differences across syllabi in the extent to which they cover soft skills. We do so by focusing on each course's evaluation scheme. Specifically, we consider a course to be more soft-skills intensive if the assignments portion of the syllabus has a higher share of words such as “group”, “team”, “presentation”, “essay”, “proposal”, “report”, “drafting”, and “survey”. In the average syllabus, 33 percent of the words in the assignment portion of the syllabus refer to soft skills (Table 1, panel (a)).

The measure of soft-skills intensity is negatively correlated with the education-innovation gap

(with a correlation of -0.14 , Figure 5, panel (a)). Cross-school differences in the skill intensity of the courses display the same patterns we found for the education-innovation gap: The prevalence of soft skills is higher in schools with higher expenditure on instruction and salaries, increases with school selectivity, and is larger for schools with a higher median parental income and with a lower share of minority students (Figure AVIII, panel (a)). Soft skills are also more prevalent among courses taught by more research-productive instructors (Figure AIX, panel (a)).

We also study the relationship between soft-skills intensity and student outcomes. Controlling for the full set of school observables used in Tables 8 and 9, a one-standard deviation higher soft-skills intensity of a school's courses is associated with a 1.2 percentage-point higher graduation rate (2 percent, Table AIV, panel h, column 1, significant at 1 percent); a 1.7 percent higher mean income (column 2, significant at 1 percent); and a 1.2 percent higher chances of reaching the top income quintile for students with parental income in the bottom quintile (18 percent, column 9, significant at 1 percent).

Taken together, these findings indicate that differences across and within schools in course content are not limited to the extent to which content is up-to-date, but also extend to the skills that are trained. We interpret this as additional evidence for the importance of accounting for differences in content across courses when characterizing the heterogeneity of educational experiences for students at different schools.

8 Alternative Measures for The Education-Innovation Gap

In spite of its desirable properties, our measure of the education-innovation gap has some limitations. For example, the gap penalizes courses that include old *and* new content, relative to courses that include exactly the same new content but no old content. Being devised to measure the “average” age of content, the gap is also unable to distinguish courses with extremely novel content among those with the same gap. Lastly, the gap only captures the similarity of syllabi with academic publications. Especially in some fields, a course with relatively old academic content could still be novel in other dimensions, for example if it teaches recent technological innovations described in patents.

In this section, we probe the robustness of our results using alternative measures of a course's content designed to address these issues.

Presence of Old Content The education-innovation gap measures the presence, in a syllabus, of new content relative to older content. Consider two syllabi which both cover the same frontier

research in a given field; the first syllabus is shorter and only contains this new content, while the second one is longer and also contains older content. Our measure would assign a lower gap to the first syllabus compared to the second, even if both do an equal job in terms of covering frontier knowledge.

To address this limitation, we construct an alternative metric which measures the *share of old knowledge* of each syllabus, defined as one minus the ratio between the number of “new words” in each syllabus (defined as knowledge words that are (a) in the top 5 percent of the word frequency among articles published between $t - 3$ and $t - 1$, or (b) used in articles published between $t - 3$ and $t - 1$ but not in those published between $t - 15$ and $t - 13$) and the number of all new words. The correlation between the share of old knowledge and the education-innovation gap is 0.22 (Figure 5, panel (b)), and our main results carry through if we use this alternative formulation as a measure of novelty of a syllabus’s content (see panel (b) of Figure AVIII for the correlation with school-level characteristics; panel (b) of Figure AIX for the correlation with instructors’ research productivity; and panels a and b of Table AIV for the relationship with student outcomes).

Right Tail of Academic Novelty The education-innovation gap captures the “average” novelty of a syllabus. It is possible for two syllabi to have the same gap when one of them only covers content from five years prior while the other covers mostly material from fifteen years prior, but also a small amount of material from the previous year. To construct a measure that captures the presence of “extremely” new material in a syllabus, we proceed as follows. First, we draw 100 “sub-syllabi” from each syllabus, defined as subsets of 20 percent of the syllabus’s words, and calculate the corresponding education-innovation gap. We then recalculate the average gap among all sub-syllabi in the bottom 5 percent of the gap distribution of a given syllabus.³⁰ We refer to this as a “tail measure” of novelty.

The tail measure is positively correlated with the education-innovation gap, with a correlation of 0.67 (Figure 5, panel (c)). All our results hold when using the tail measure as a metric for syllabus novelty (see panel (c) of Figure AVIII, for the correlation with school-level characteristics; panel (c) of Figure AIX for the correlation with instructors’ research productivity; and panels c and d of Table AIV for the relationship with student outcomes).

Gap with Patents The education-innovation gap is defined using new academic publications as the frontier of knowledge. For STEM fields, knowledge advancements are also documented in the form of patents. To incorporate this information in our analysis, we construct a version of the

³⁰Our results are robust to the use of the top 10 and one percent.

education-innovation gap for STEM courses that uses patents in lieu of academic publications. This measure is positively correlated with the standard education-innovation gap (Figure 5, panel (d)). In addition, our main results carry over when using the patent-based gap (see panel (d) of Figure AVIII, for the correlation with school-level characteristics; panel (d) of Figure AIX for the correlation with instructors' research productivity; and panels e and f of Table AIV for the relationship with student outcomes).

Taken together, these results indicate that our main conclusions regarding the content of higher-education courses across schools and its relationship with instructors' characteristics and student outcomes are not dependent on the specific way in which we measure up-to-date content.

9 Conclusion

This paper uses the text of HE course syllabi to quantify the distance between the content of each course and frontier knowledge. Our approach centers around a new measure, the "education-innovation gap," defined as the textual similarity between course syllabi and knowledge from older vintages, relative to newer ones. We constructed this measure by applying NLP techniques to a novel data set that contains the full text of 1.7 million syllabi and 20 million academic publications.

Using our measure, we document a set of new findings about the dissemination of frontier knowledge across HE programs. Across and within schools, significant differences exist in the extent to which frontier knowledge is offered to students. More selective schools and those with more resources offer courses with a smaller gap. Since these schools enroll a lower portion of socio-economically disadvantaged students, access to updated knowledge is highly unequal across students from different backgrounds. Instructors play the largest role in shaping the content of the courses they teach. Among all instructors, those who are more research-active are more likely to teach courses with lower gaps. The education-innovation gap is strongly correlated with students' innovation and labor-market outcomes. In schools offering courses with lower gaps, students are more likely to graduate, earn a PhD, and produce patents. They also earn more once they enter the labor market.

Taken together, our findings indicate that the education-innovation gap can be an important metric for quantifying how frontier knowledge is produced and disseminated and could help shed new light on the way in which schools and instructors impact students' lives. A careful analysis of the causal impacts of a low-gap education on students' later life outcomes represents a fruitful avenue for future research.

References

- Akcigit, Ufuk, Jeremy G Pearce, and Marta Prato, 2020, Tapping into talent: Coupling education and innovation policies for economic growth, Technical report, National Bureau of Economic Research.
- Altonji, Joseph G, Erica Blom, and Costas Meghir, 2012, Heterogeneity in human capital investments: High school curriculum, college major, and careers, *Annual Review of Economics* 4, 185–223.
- Andrews, Michael, 2017, The role of universities in local invention: evidence from the establishment of us colleges, *Job Market Paper* .
- Angrist, Joshua, Pierre Azoulay, Glenn Ellison, Ryan Hill, and Susan Feng Lu, 2017, Economic research evolves: Fields and styles, *American Economic Review* 107, 293–97.
- Angrist, Joshua D, and Jörn-Steffen Pischke, 2017, Undergraduate econometrics instruction: through our classes, darkly, *Journal of Economic Perspectives* 31, 125–44.
- Arnold, Ivo JM, 2008, Course level and the relationship between research productivity and teaching effectiveness, *The Journal of Economic Education* 39, 307–321.
- Azoulay, Pierre, and Danielle Li, 2020, Scientific grant funding, Technical report, National Bureau of Economic Research.
- Baumol, William J, 2005, Education for innovation: Entrepreneurial breakthroughs versus corporate incremental improvements, *Innovation Policy and the Economy* 5, 33–56.
- Becker, William E, and Peter E Kennedy, 2005, Does teaching enhance research in economics?, *American Economic Review* 95, 172–176.
- Bell, Alex, Raj Chetty, Xavier Jaravel, Neviana Petkova, and John Van Reenen, 2019, Who becomes an inventor in america? the importance of exposure to innovation, *Quarterly Journal of Economics* 134, 647–713.
- Benhabib, Jess, and Mark M Spiegel, 2005, Human capital and technology diffusion, *Handbook of economic growth* 1, 935–966.

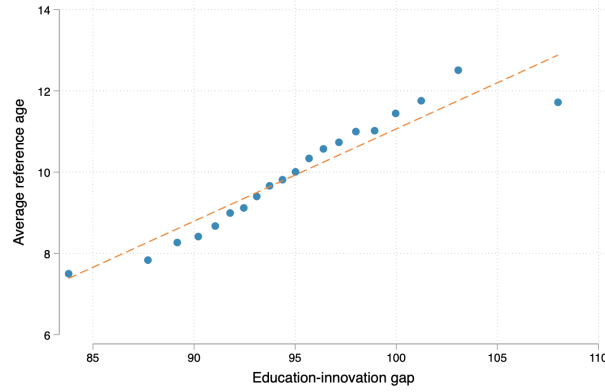
- Bianchi, Nicola, and Michela Giorcelli, 2019, Scientific education and innovation: from technical diplomas to university stem degrees, *Journal of the European Economic Association* .
- Biasi, Barbara, David J Deming, and Petra Moser, 2020, Education and innovation, in *The Role of Innovation and Entrepreneurship in Economic Growth* (University of Chicago Press).
- Bloom, Nicholas, Tarek Alexander Hassan, Aakash Kalyani, Josh Lerner, and Ahmed Tahoun, 2021, The diffusion of disruptive technologies, Technical report, National Bureau of Economic Research.
- Braga, Michela, Marco Paccagnella, and Michele Pellizzari, 2016, The impact of college teaching on students' academic and labor market outcomes, *Journal of Labor Economics* 34, 781–822.
- Carrell, Scott E, and James E West, 2010, Does professor quality matter? evidence from random assignment of students to professors, *Journal of Political Economy* 118, 409–432.
- Chetty, Raj, John N Friedman, Emmanuel Saez, Nicholas Turner, and Danny Yagan, 2020, Income segregation and intergenerational mobility across colleges in the united states, *Quarterly Journal of Economics* 135, 1567–1633.
- Courant, Paul N, and Sarah Turner, 2020, Faculty deployment in research universities, in *Productivity in Higher Education* (University of Chicago Press).
- Dale, Stacy, and Alan B Krueger, 2011, Estimating the return to college selectivity over the career using administrative earnings data, *NBER Working Paper* .
- De Vlieger, Pieter, Brian Jacob, and Kevin Stange, 2020, Measuring instructor effectiveness in higher education, in *Productivity in Higher Education* (University of Chicago Press).
- Deming, David, and Lisa B Kahn, 2018, Skill requirements across firms and labor markets: Evidence from job postings for professionals, *Journal of Labor Economics* 36, S337–S369.
- Deming, David J, 2014, Using school choice lotteries to test measures of school effectiveness, *American Economic Review* 104, 406–11.
- Deming, David J, 2017, The growing importance of social skills in the labor market, *Quarterly Journal of Economics* 132, 1593–1640.

- Feld, Jan, Nicolás Salamanca, and Ulf Zölitz, 2020, Are professors worth it? the value-added and costs of tutorial instructors, *Journal of Human Resources* 55, 836–863.
- Goldin, Claudia Dale, and Lawrence F Katz, 2010, *The Race Between Education and Technology* (Harvard University Press).
- Goolsbee, Austan, and Chad Syverson, 2019, Monopsony power in higher education: A tale of two tracks, Technical report, National Bureau of Economic Research.
- Hanushek, Eric A, and Ludger Woessmann, 2012, Do better schools lead to more growth? Cognitive skills, economic outcomes, and causation, *Journal of Economic Growth* 17, 267–321.
- Hattie, John, and Herbert W Marsh, 1996, The relationship between research and teaching: A meta-analysis, *Review of educational research* 66, 507–542.
- Hemelt, Steven W, Brad Hershbein, Shawn M Martin, and Kevin M Stange, 2021, College majors and skills: Evidence from the universe of online job ads, Technical report, National Bureau of Economic Research.
- Hoffman, F, and P Oreopoulos, 2009, Professor qualities and student performance, *Review of Economics and Statistics* 91, 83–92.
- Hoxby, Caroline M, 1998, The return to attending a more selective college: 1960 to the present, *Unpublished manuscript, Department of Economics, Harvard University, Cambridge, MA* .
- Hoxby, Caroline M, 2020, The productivity of us postsecondary institutions, in *Productivity in Higher Education*, 31–66 (University of Chicago Press).
- Israeli, Osnat, 2007, A shapley-based decomposition of the r-square of a linear regression, *The Journal of Economic Inequality* 5, 199–212.
- Jones, Benjamin F, 2009, The burden of knowledge and the death of the renaissance man: is innovation getting harder?, *Review of Economic Studies* 76, 283–317.
- Kantor, Shawn, and Alexander Whalley, 2019, Research proximity and productivity: long-term evidence from agriculture, *Journal of Political Economy* 127, 819–854.

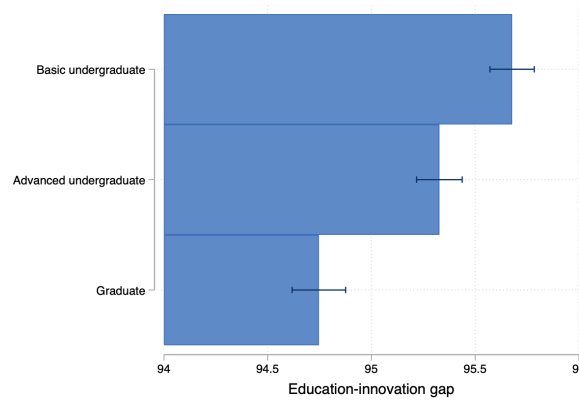
- Kelly, Bryan, Dimitris Papanikolaou, Amit Seru, and Matt Taddy, 2021, Measuring technological innovation over the long run, *American Economic Review: Insights* 3, 303–20.
- Li, Xiaoxiao, Sebastian Linde, and Hajime Shima, 2021, Major complexity index and college skill production, *Available at SSRN 3791651* .
- Ma, Xuezhe, and Eduard Hovy, 2016, End-to-end sequence labeling via bi-directional lstm-cnncrf, *arXiv preprint arXiv:1603.01354* .
- Mountjoy, Jack, and Brent Hickman, 2020, The returns to college (s): Estimating value-added and match effects in higher education, *University of Chicago, Becker Friedman Institute for Economics Working Paper* .
- Nelson, Richard R, and Edmund S Phelps, 1966, Investment in humans, technological diffusion, and economic growth, *American Economic Review* 56, 69–75.
- Sanh, Victor, Lysandre Debut, Julien Chaumond, and Thomas Wolf, 2019, Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, *arXiv preprint arXiv:1910.01108* .
- Tartari, Valentina, and Scott Stern, 2021, More than an ivory tower: The impact of research institutions on the quantity and quality of entrepreneurship, Technical report, National Bureau of Economic Research.
- Toivanen, Otto, and Lotta Väänänen, 2016, Education and invention, *Review of Economics and Statistics* 98, 382–396.
- Valero, Anna, and John Van Reenen, 2019, The economic impact of universities: Evidence from across the globe, *Economics of Education Review* 68, 53–67.

Figure 1: Validating The Education-Innovation Gap

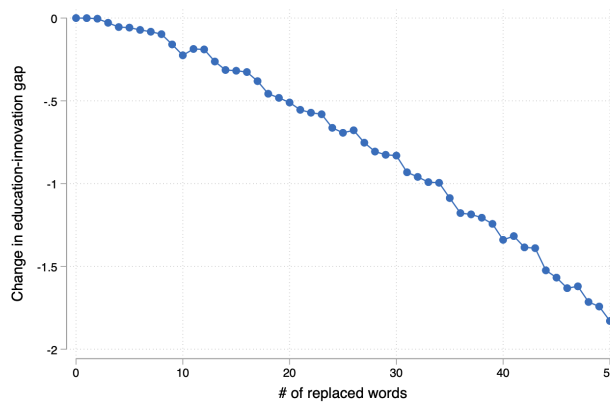
(a) Gap and Average Age of References Included in The Syllabi



(b) Gap by Course Level

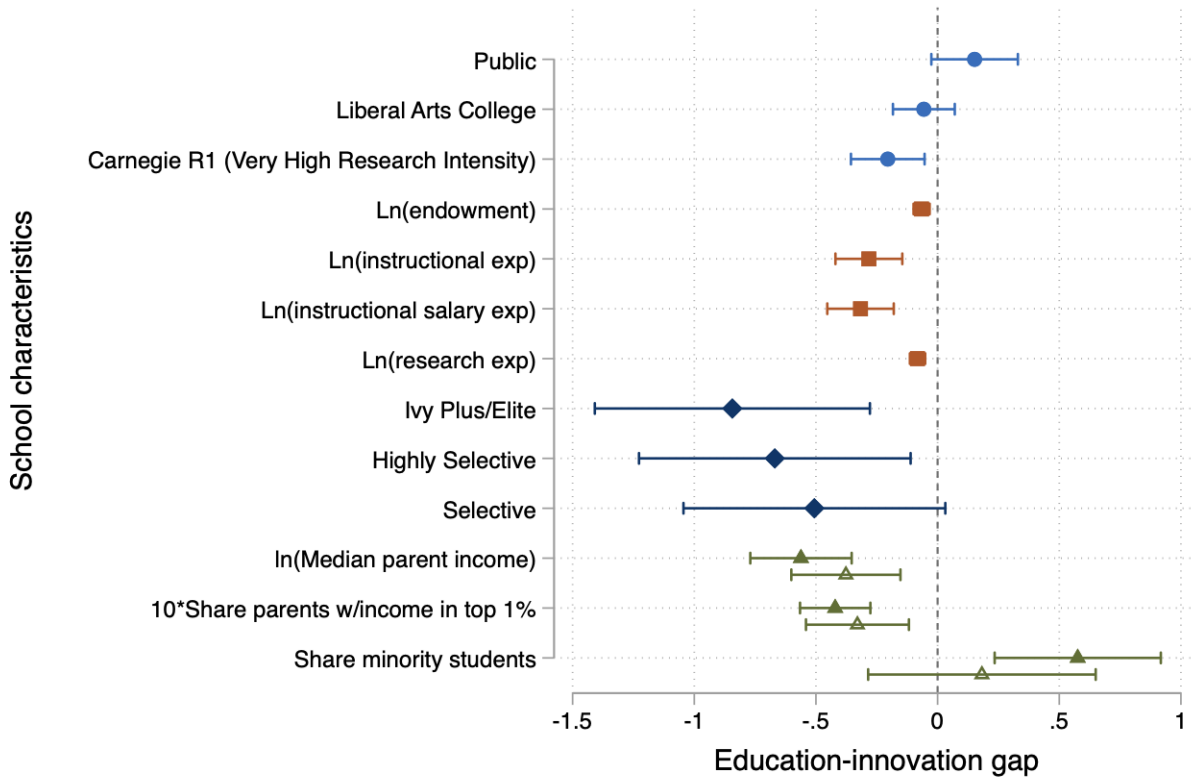


(c) Change in Gap as Old Words Are Replaced with New Words



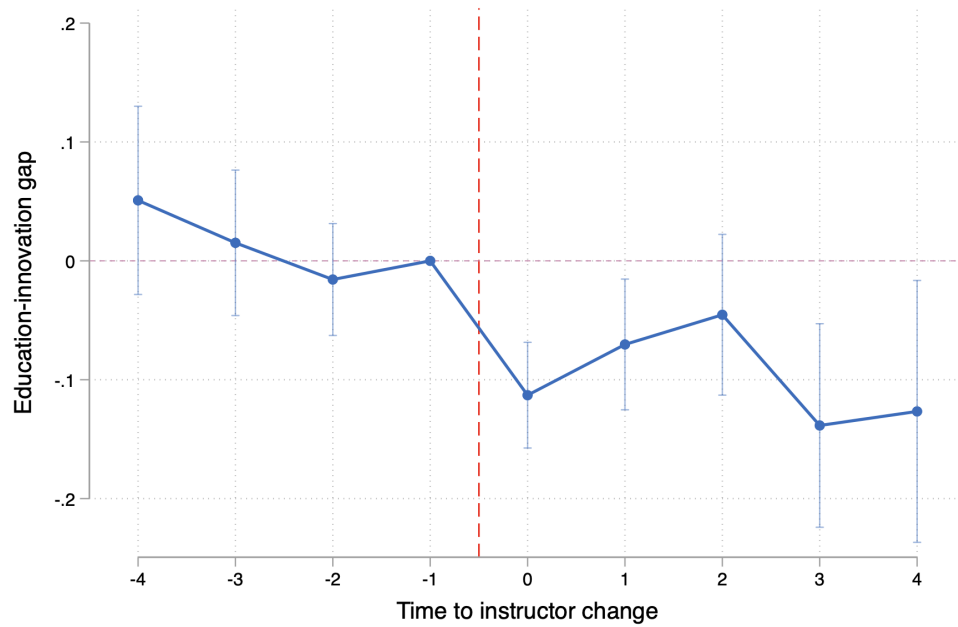
Note: Panel a) shows a binned scatterplot of the education-innovation gap and the average age of a syllabus’s references (required or recommended readings), where age is defined as the difference between the year of the syllabus and the year of publication of each reference. Panel b) shows the mean and 95-percent confidence intervals of the gap by course level, controlling for field-by-year effects. Panel c) shows the change in the gap for a subsample of 100,000 syllabi, in which we progressively replace “old” words with “new” words.

Figure 2: The Education-Innovation Gap and School Characteristics



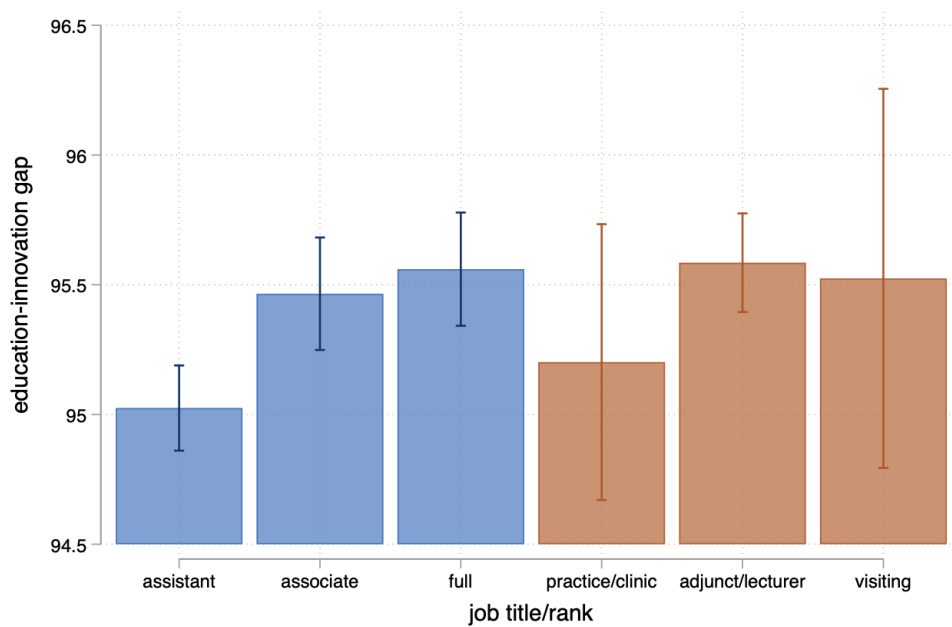
Notes: Point estimates and 95-percent confidence intervals of coefficient β in equation (5), i.e., the slope of the relationship between each reported variable and the education-innovation gap controlling for field-by-course level-by-year fixed effects. Each coefficient is estimated from a separate regression, with the exception of selectivity tiers (Ivy Plus/Elite, Highly Selective, Selective) which are jointly estimated. Endowment, expenditure, and share minority refer to the year 2018 and is taken from IPEDS. Estimates are obtained pooling syllabi data for the years 1998 to 2018. Standard errors are clustered at the school level.

Figure 3: Event Study: The Education-Innovation Gap Around An Instructor Change



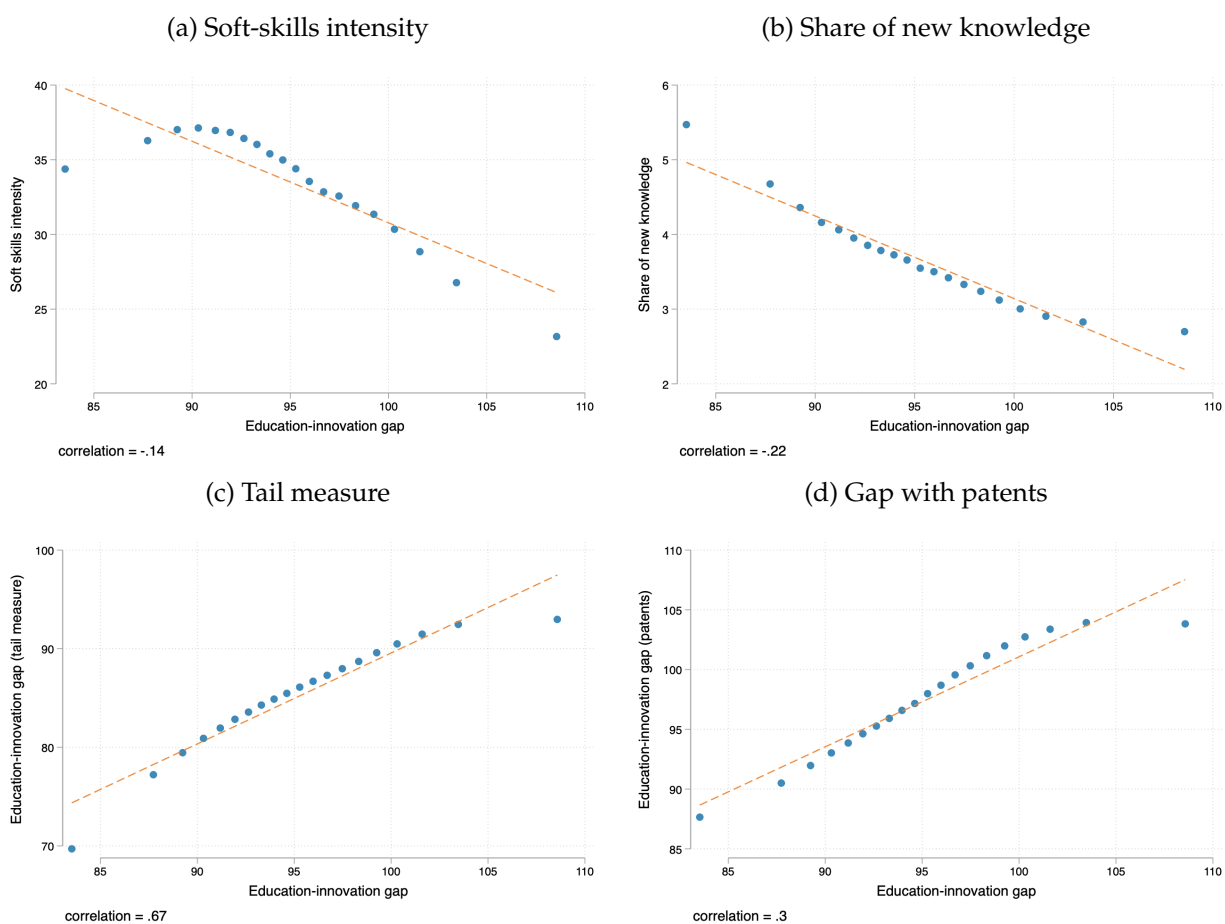
Notes: Estimates and confidence intervals of the parameters δ_k in equation (6), representing an event study of the education-innovation gap around an instructor change and controlling for course and field-by-year fixed effects. Observations are at the course-by-year level; we focus on courses with at most two episodes of instructor changes. Standard errors clustered at the course level.

Figure 4: Gap by Job Titles



Notes: Mean education-innovation gap by job title, along with 95-percent confidence intervals. Means are obtained as OLS coefficients from a regression of the gap on indicators for the job title of the instructor, as well as field-by-course level-by-year fixed effects. Estimates are obtained pooling data for multiple years. Standard errors are clustered at the school level.

Figure 5: The Education-Innovation Gap and Alternative Measures of Novelty: Binned Scatterplots



Notes: Binned scatterplots of the education-innovation gap and four alternative measures of novelty of each syllabus: a measure of soft-skills intensity, defined as the share of words in the assignment portion of a syllabus which refer to soft skills (panel (a)); a measure of new knowledge, defined as the share of all new words contained by each syllabus (where new words are knowledge words that are (a) in the top 5 percent of the word frequency among articles published between $t-3$ and $t-1$, or (b) used in articles published between $t-3$ and $t-1$ but not in those published between $t-15$ and $t-13$ (panel (b))); a “tail measure,” calculated for each syllabus by (a) randomly selecting 100 subsamples containing 20 percent of the syllabus’s words, (b) calculating the gap for each subsample, and (c) selecting the 5th percentile of the corresponding distribution (panel (c)); and the education-innovation gap calculated using the text of all patents as a benchmark, instead of academic articles (panel (d)).

Table 1: Summary Statistics: Courses, Instructors, and Schools

Panel (a): Syllabus (Course) Characteristics						
	count	mean	std	25%	50%	75%
Education-innovation gap	1,706,319	95.3	5.8	91.6	94.9	98.8
# Words	1,706,319	2226	1987	1068	1778	2796
# Knowledge words	1,706,319	1011	1112	349	656	1236
# Unique knowledge word	1,706,319	420	327	203	330	535
Soft skills	1,703,863	33.4	22.9	14.0	30.5	50.0
STEM	1,706,319	0.326	0.469	0	0	1
Business	1,706,319	0.103	0.304	0	0	0
Humanities	1,706,319	0.299	0.457	0	0	1
Social science	1,706,319	0.240	0.427	0	0	0
Basic	1,706,319	0.393	0.488	0	0	1
Advanced	1,706,319	0.275	0.446	0	0	1
Graduate	1,706,319	0.332	0.471	0	0	1

Panel (b): Instructor (Professor) Research Productivity						
	count	mean	std	25%	50%	75%
Ever Published?	332,064	0.41	0.49	0	0	1.00
# Publications per year	135,364	1.51	1.94	1.00	1.00	1.38
# Publications, last 5 years	111,404	6.01	14.89	0	1.00	5.42
# Citations per year	135,364	29.22	105.92	0	1.85	17.92
# Citations, last 5 years	111,404	172.46	887.99	0	0	54.32
Ever Grant?	332,064	0.18	0.38	0	0	0
# Grants	58,136	10.14	19.96	2.00	4.00	10.00
Grant amount (\$1,000)	54,462	4,023	19,501	236	912	3,201

Panel (c): Students' Characteristics and Outcomes at University Level						
	count	mean	std	25%	50%	75%
Median parental income (\$1,000)	767	97,917	31,054	78,000	93,500	109,900
Share parents w/income in top 1%	767	0.030	0.041	0.006	0.013	0.033
Share minority students	760	0.221	0.166	0.116	0.166	0.267
Graduation rates (2012–13 cohort)	758	0.614	0.188	0.473	0.616	0.765
Income (2003–04, 2004–05 cohorts)	762	45,035	10,235	38,200	43,300	49,800
Intergenerational mobility	767	0.294	0.138	0.182	0.280	0.375
Admission rate	715	0.642	0.218	0.533	0.683	0.800
SAT score	684	1104.4	130.5	1011.5	1079.5	1182.0

Note: Summary statistics of the variables used in the analysis.

Table 2: Selection Into The Sample: Share of Syllabi Included in the Sample and Institution-Level Characteristics

Panel (a): Share and Δ Share, Correlation w/ School Characteristics				
	Share in OSP, 2018		Δ Share in OSP, 2008-18	
	(1)	(2)	(3)	(4)
	Corr.	SE	Corr.	SE
ln Expenditure on instruction	0.002	(0.005)	0.015	(0.010)
ln Endowment per capita	-0.001	(0.002)	-0.001	(0.002)
ln Sticker price	0.003	(0.007)	0.007	(0.010)
ln Avg faculty salary	0.016	(0.020)	0.049	(0.024)
ln Enrollment	0.018	(0.009)	0.019	(0.011)
Share Black students	-0.030	(0.038)	0.035	(0.060)
Share Hispanic students	0.171	(0.145)	0.161	(0.115)
Share Asian students	0.186	(0.214)	0.324	(0.239)
Share grad in Arts & Humanities	0.159	(0.168)	0.189	(0.179)
Share grad in STEM	-0.001	(0.028)	0.064	(0.056)
Share grad in Social Sciences	0.014	(0.024)	0.104	(0.056)
Share grad in Business	0.037	(0.065)	0.116	(0.065)
F-stat	1.015		1.376	

Panel (b): Share and Δ Share, By School Tier				
	Share in OSP, 2018		Δ Share in OSP, 2008-18	
	(1)	(2)	(3)	(4)
	Mean	SE	Mean	SE
Ivy Plus/Elite	0.024	(0.008)	0.022	(0.009)
Highly Selective	0.003	(0.003)	0.006	(0.004)
Selective Private	0.029	(0.018)	0.001	(0.029)
Selective Public	0.040	(0.023)	0.009	(0.029)
F-stat	3.677		1.806	

Note: The top panel shows OLS coefficients (“means”) and robust standard errors (“SE”) of univariate regressions of each listed dependent variable on the corresponding independent variable. The bottom panel shows OLS coefficients (“means”) and syllabus-clustered standard errors (“SE”) of a regression of each dependent variable on indicators for school tiers. The dependent variables are the school-level share of syllabi contained in the OSP sample in 2018 (columns 1-2) and the change in this share between 2008 and 2018 columns (3-4). The F-statistics refer to multivariate regressions that include all the listed independent variables, and test for the joint significance of these variables.

Table 3: Decomposing the Variation In The Gap: Schools, Years, Fields, Courses, and Instructors

Variable	Partial R^2	
Year	0.169	0.180
Field	0.039	0.056
School	0.021	0.028
Course level	.	0.008
Course	0.330	.
Instructor	0.248	0.346
Total	0.161	0.124

Note: The table shows a decomposition of the adjusted R^2 of a regression of the education-innovation gap on all sets of listed fixed effects into the contribution of each set of fixed effects. This is done using a Shapley-Owen decomposition, which calculates the partial R^2 of each set of variables j as $R_j^2 = \sum_{k \neq j} \frac{R^2 - R^2(-j)}{K! / j!(K-j-1)!}$ where $R^2(-j)$ is the R^2 of a regression that excludes variables j . Column 1 includes course fixed effects; column 2 only includes course level fixed effects. We use adjusted R^2 in lieu of R^2 to account for the large number of fixed effects.

Table 4: The Education-Innovation Gap Around An Instructor Change

Instructor change	All Fields (1)	Business (2)	Humanities (3)	STEM (4)	Social Science (5)	Basic (6)	Advanced (7)	Graduate (8)
After change	-0.1021*** (0.0244)	-0.1009 (0.0683)	-0.1417*** (0.0464)	-0.1060*** (0.0399)	-0.0289 (0.0416)	-0.0897** (0.0456)	-0.0875* (0.0450)	-0.1152*** (0.0374)
N (Course x year)	379482	36325	105316	152974	95223	125493	112206	137721
# Courses	126343	11775	35947	45982	31805	43530	35395	46213
Course FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Field x Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

Note: OLS estimates; one observation is a course in a given year. The dependent variable is the education-innovation gap. The variable *After change* is an indicator for years following an instructor change, for courses with only one instructor and at most two instructor changes over the observed time period. All specifications control for course and field-by-year fixed effects. Standard errors in parentheses are clustered at the course level. * ≤ 0.1 , ** ≤ 0.05 , *** ≤ 0.01 .

Table 5: The Education-Innovation Gap and Instructors' Research Productivity: Publications and Citations

	All Fields (1)	Business (2)	Humanities (3)	STEM (4)	Social Science (5)	Basic (6)	Advanced (7)	Graduate (8)
Panel a): #publications								
1st quartile	-0.0219 (0.0178)	0.0485 (0.0472)	-0.0815** (0.0328)	0.0556 (0.0367)	-0.0641** (0.0291)	-0.0099 (0.0298)	-0.0277 (0.0324)	-0.0308 (0.0300)
2nd quartile	-0.0151 (0.0298)	0.0138 (0.0729)		0.0309 (0.0462)	-0.0418 (0.0425)	-0.0207 (0.0531)	0.0224 (0.0543)	-0.0366 (0.0471)
3rd quartile	-0.0045 (0.0302)	0.0596 (0.0712)	-0.0387 (0.0694)	0.0953* (0.0563)	-0.1057** (0.0459)	0.0374 (0.0574)	-0.0115 (0.0540)	-0.0356 (0.0462)
4th quartile	-0.1103*** (0.0376)	0.0220 (0.0894)	-0.1184 (0.0797)	-0.0638 (0.0699)	-0.1817*** (0.0621)	-0.0448 (0.0742)	-0.0927 (0.0698)	-0.1711*** (0.0551)
Panel b): #citations								
1st quartile	0.0288 (0.0248)	-0.0097 (0.0667)	0.0313 (0.0636)	0.1068** (0.0437)	-0.0469 (0.0361)	0.0438 (0.0427)	0.0543 (0.0450)	-0.0031 (0.0409)
2nd quartile	0.0194 (0.0282)	0.0050 (0.0675)	0.0106 (0.0682)	0.0827* (0.0499)	-0.0413 (0.0431)	0.0271 (0.0508)	0.0276 (0.0508)	0.0059 (0.0448)
3rd quartile	-0.0658** (0.0324)	-0.0464 (0.0775)	-0.0919 (0.0782)	-0.0355 (0.0584)	-0.1056** (0.0491)	0.0011 (0.0625)	-0.0965 (0.0600)	-0.0961** (0.0477)
4th quartile	-0.0713* (0.0412)	0.0667 (0.0946)	-0.1090 (0.1056)	-0.0151 (0.0722)	-0.1305** (0.0655)	-0.0200 (0.0799)	-0.0497 (0.0775)	-0.1385** (0.0601)
N (Course x year)	579622	60953	156970	195375	150731	209190	170946	199228
# Courses	153392	15156	43067	51873	39169	55444	43320	54557
Course FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Field x Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

Note: OLS estimates; one observation is a course in a given year. The dependent variable is the education-innovation gap; the independent variables are indicators for quartiles of the number of publications (panel (a)) and citations (panel (b)) of a course's instructors in the previous five years. The omitted category is courses with instructors with no publications or citations. For courses with more than one instructor, we consider the mean number of publications and citations across all instructors. All specifications control for course and field-by-year fixed effects. Standard errors in parentheses are clustered at the course level. * ≤ 0.1 , ** ≤ 0.05 , *** ≤ 0.01 .

Table 6: The Education-Innovation Gap and The Fit Between Instructors' Research and Course Content

	All Fields (1)	Business (2)	Humanities (3)	STEM (4)	Social Science (5)	Basic (6)	Advanced (7)	Graduate (8)
Fit w/top course (sd)	-0.0877** (0.0398)	0.1638 (0.0997)	0.0017 (0.1728)	-0.0756 (0.0559)	-0.0845 (0.0656)	-0.0637 (0.0832)	-0.1428* (0.0790)	-0.0611 (0.0558)
N (Course x year)	54591	3293	2270	35859	12626	16743	16224	21139
# Courses	17077	1040	781	11166	3923	5208	4833	6883
Course FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Field x Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

Note: OLS estimates; one observation is a course in a given year. The dependent variable is the education-innovation gap. The variable *Fit w/top course* is a measure of fit between the instructor's research and the content of the course, defined as the cosine similarity between the instructor's research in the previous 5 years and the content of the course with the smallest education-innovation gap among all courses in the same topic across all schools. All specifications control for course and field-by-year fixed effects. Standard errors in parentheses are clustered at the course level. * ≤ 0.1 , ** ≤ 0.05 , *** ≤ 0.01 .

Table 7: The Education-Innovation Gap and Instructors' Research Resources: NSF/NIH Grants

	All Fields (1)	Business (2)	Humanities (3)	STEM (4)	Social Science (5)	Basic (6)	Advanced (7)	Graduate (8)
At least one grant	-0.0453** (0.0199)	0.0045 (0.0571)	-0.0649 (0.0400)	-0.0316 (0.0367)	-0.0596* (0.0330)	-0.0476 (0.0327)	-0.0324 (0.0370)	-0.0567* (0.0336)
N (Course x year)	581995	60953	156970	195375	150731	210121	171867	199735
# Courses	153809	15156	43067	51873	39169	55594	43474	54663
Course FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Field x Year FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

Note: OLS estimates; one observation is a course in a given year. The dependent variable is the education-innovation gap. The variable *At least one grant* equals one if the course's instructor (or at least one of the course's instructors in case of multiple instructors) has received at least one NSF or NIH grant. All specifications control for course and field-by-year fixed effects. Standard errors in parentheses are clustered at the course level. * ≤ 0.1 , ** ≤ 0.05 , *** ≤ 0.01 .

Table 8: The Education-Innovation Gap and Innovation Measures: Share of Undergraduate Students Who Obtain a PhD and Total Nr of Patents

	Share of students who obtain a PhD, by field							Nr Patents
	All	STEM	Health	Business	Social Science	Humanities		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	
Panel (a): no controls								
Gap (sd)	-0.0044** (0.0018)	-0.0010 (0.0022)	-0.0074** (0.0033)	-0.0003 (0.0005)	-0.0124** (0.0051)	0.0054 (0.0076)	-26.8006 (16.6213)	
Mean dep. var.	0.0265	0.0452	0.0249	0.0021	0.0335	0.0228	129.7813	
N	65755	14714	9218	12698	14657	14468	1715	
Panel (b): w/ controls								
Gap (sd)	-0.0046** (0.0021)	0.0002 (0.0019)	-0.0066** (0.0030)	-0.0003 (0.0005)	-0.0101** (0.0046)	0.0061 (0.0074)	-21.8882* (12.5836)	
Mean dep. var.	0.0269	0.0461	0.0257	0.0021	0.0342	0.0228	131.0248	
N	47723	10673	6656	9243	10645	10506	1610	

Note: OLS estimates of the coefficient δ in equation (9). In columns 1-6, the variable *Gap* (*sd*) is a school-by-macro field-level education-innovation gap (estimated as $\theta_{s(i)}$ in equation (8), separately for each macro-field), standardized to have mean zero and variance one. In column 7, *Gap* (*sd*) is estimated at the school level pooling data from all fields. In columns 1-6, the dependent variable is the share of undergraduate students at each institution who eventually complete a PhD (from the NSF Survey of Doctorate Recipients, year 2000); in column 7, it is the total number of patents filed by students at each school, from Chetty et al. (2020). All columns in panel b control for sector (private or public), selectivity tiers, and an interaction between selectivity tiers and an indicator for R1 institutions according to the Carnegie classification; student-to-faculty ratio and the share of ladder faculty; total expenditure, research expenditure, instructional expenditure, and salary instructional expenditure per student; the share of undergraduate and graduate enrollment and the share of white and minority students; an indicator for institutions with admission share equal to 100, median SAT and ACT scores of admitted students in 2006, and indicators for schools not using either SAT or ACT in admission; the share of students with majors in Arts and Humanities, Business, Health, Public and Social Service, Social Sciences, STEM, and multi-disciplinary fields; and the natural logarithm of parental income. Columns 1-6 control for year effects. Column 1 also controls for macro field fixed effects. Bootstrapped standard errors in parentheses are clustered at the school level. * ≤ 0.1 , ** ≤ 0.05 , *** ≤ 0.01 .

Table 9: The Education-Innovation Gap and Student Outcomes

	Income (College Scorecard)			Income (Chetty et al., 2020)					
	Grad rate	Mean	$P_y \leq 33$ pctile	Median	Mean	P(top20%)	P(top10%)	P(top5%)	$P(\text{top20\%} P_y \text{ bottom 20\%})$
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Panel (a): no controls									
Gap (sd)	-0.0513*** (0.0068)	-0.0555*** (0.0104)	-0.0645*** (0.0106)	-0.0512*** (0.0088)	-0.0722*** (0.0124)	-0.0333*** (0.0057)	-0.0265*** (0.0046)	-0.0187*** (0.0036)	-0.0293*** (0.0053)
Mean dep. var.	0.5692					0.3694	0.2082	0.1143	0.2945
N	15683	3793	3566	3793	763	763	763	763	763
# schools	761	760	734	760					
Panel (b): w/ controls									
Gap (sd)	-0.0073** (0.0030)	-0.0067 (0.0041)	-0.0083* (0.0050)	-0.0090** (0.0045)	-0.0137*** (0.0048)	-0.0084*** (0.0025)	-0.0053** (0.0021)	-0.0031** (0.0015)	-0.0047* (0.0028)
Mean dep. var.	0.5816	10.8281	10.7605	10.7096		0.3710	0.2100	0.1159	0.2957
N	11471	1996	1843	1996	718	718	718	718	718
# schools	733	727	701	727					

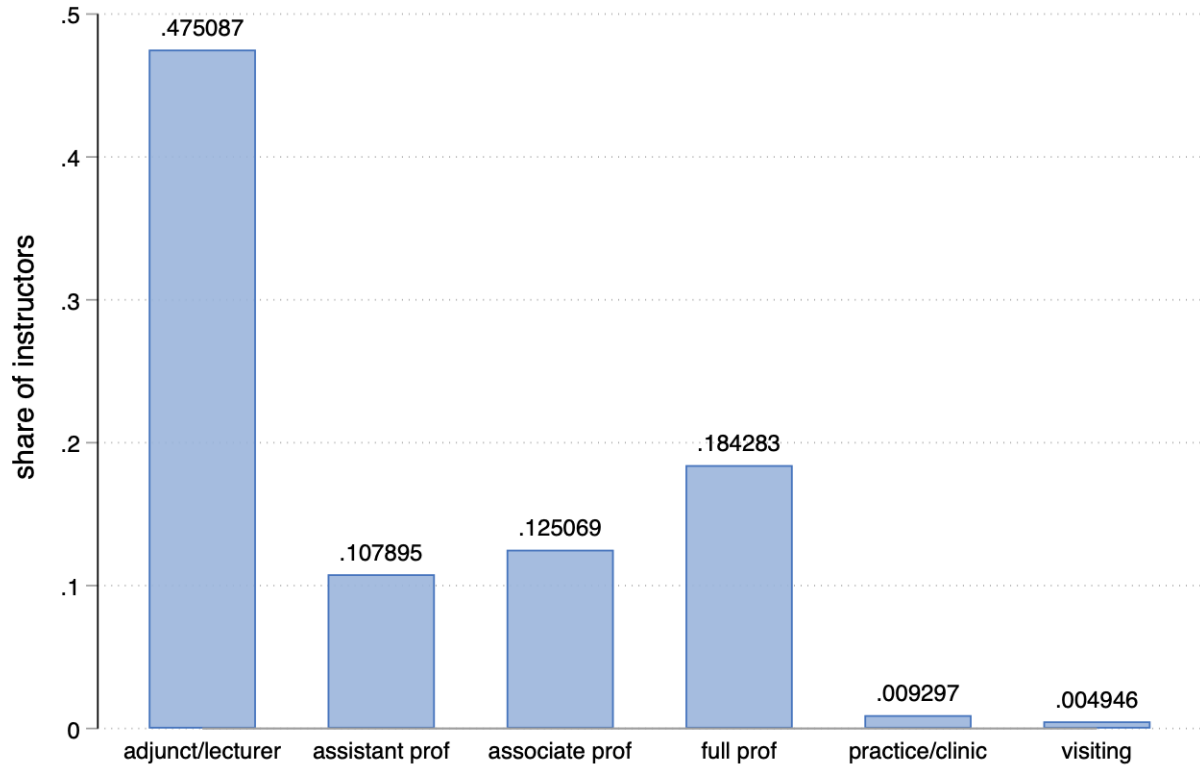
Note: OLS estimates of the coefficient δ in equation (9). The variable *Gap (sd)* is a school-level education-innovation gap (estimated as $\theta_{s(i)}$ in equation (8)), standardized to have mean zero and variance one. The dependent variable are graduation rates (from IPEDS, years 1998-2018, column 1); the log of mean student incomes from the College Scorecard, for all students (column 2) and for students with parental income in the bottom tercile (column 3); the log of median income from the College Scorecard (column 4); the log of mean income for students who graduated between 2002 and 2004 (from Chetty et al. (2020), column 5); the probability that students have incomes in the top 20, 10, and 5 percent of the national distribution (from Chetty et al. (2020), columns 6-8); and the probability that students with parental income in the bottom quintile reach the top quintile during adulthood (column 9). Columns 1-4 in panels a and b control for year effects. All columns in panel b control for sector (private or public), selectivity tiers, and an interaction between selectivity tiers and an indicator for R1 institutions according to the Carnegie classification; student-to-faculty ratio and the share of ladder faculty; total expenditure, research expenditure, instructional expenditure, and salary instructional expenditure per student; the share of undergraduate and graduate enrollment and the share of white and minority students; an indicator for institutions with admission share equal to 100, median SAT and ACT scores of admitted students in 2006, and indicators for schools not using either SAT or ACT in admission; the share of students with majors in Arts and Humanities, Business, Health, Public and Social Service, Social Sciences, STEM, and multi-disciplinary fields; and the natural logarithm of parental income. Bootstrapped standard errors in parentheses are clustered at the school level. * ≤ 0.1 , ** ≤ 0.05 , *** ≤ 0.01 .

Appendix

For online publication only

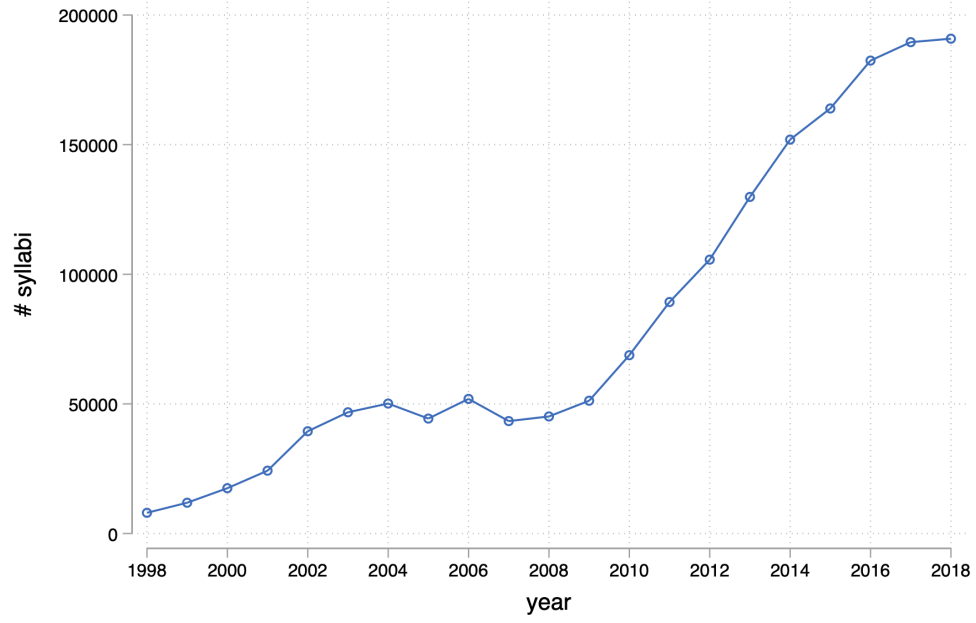
Appendix A Additional Tables and Figures

Figure AI: Distribution of Instructor Job Titles



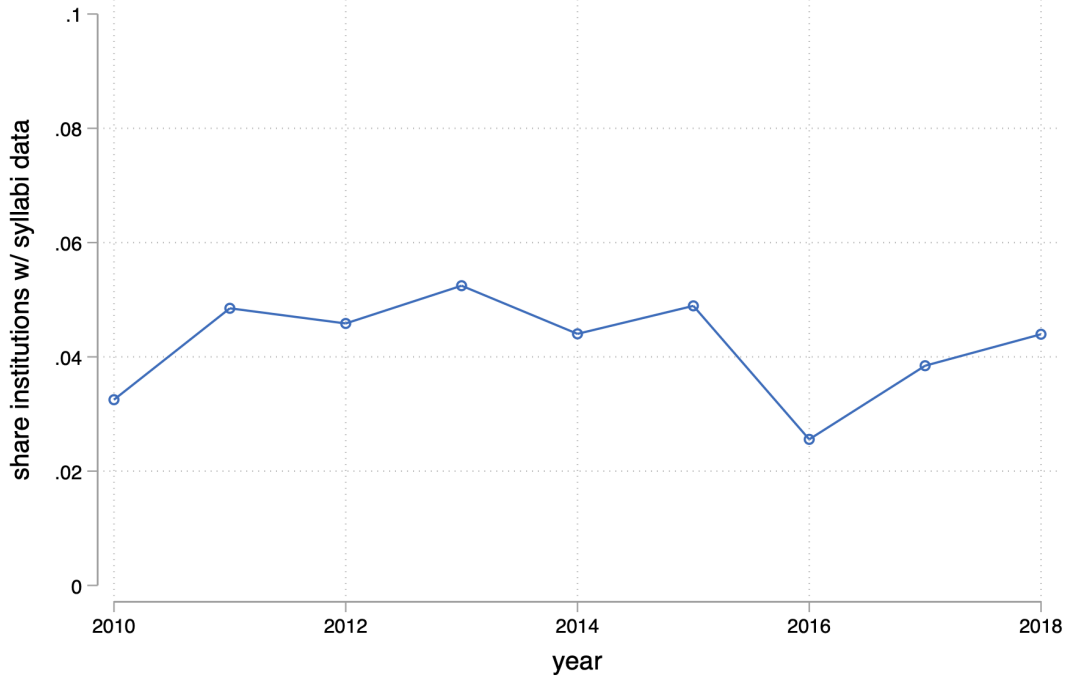
Note: Share of syllabi instructors by job title. The same is restricted to 35,178 instructors in public institutions for whom title information is available.

Figure AII: Number of Syllabi In The Sample, By Year



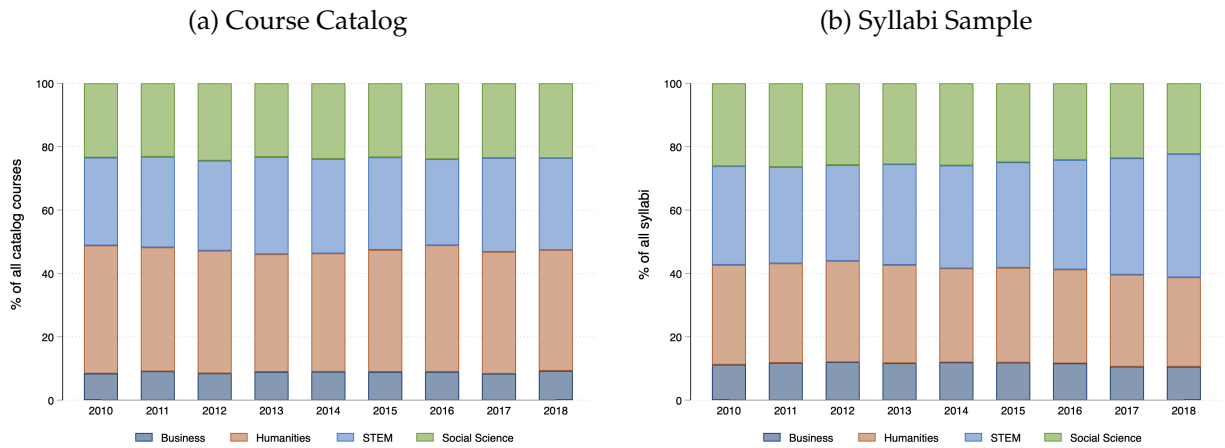
Note: Number of syllabi included in final sample, by year.

Figure AIII: Share of Catalog Courses in the Syllabi Sample



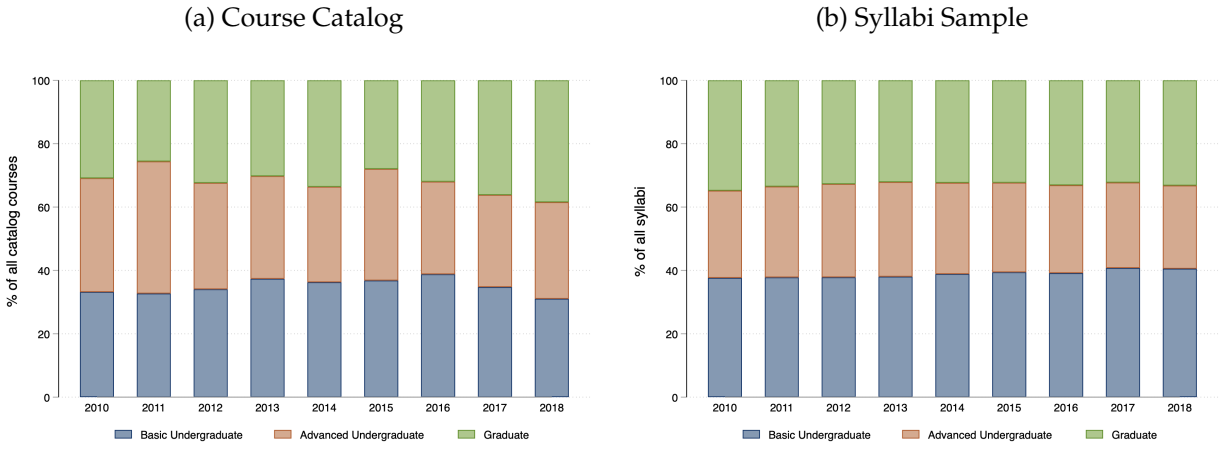
Note: Share of courses from full course catalogs whose syllabi are included in the syllabi sample.

Figure AIV: Macro-Field Coverage, Course Catalogs and Syllabi Sample



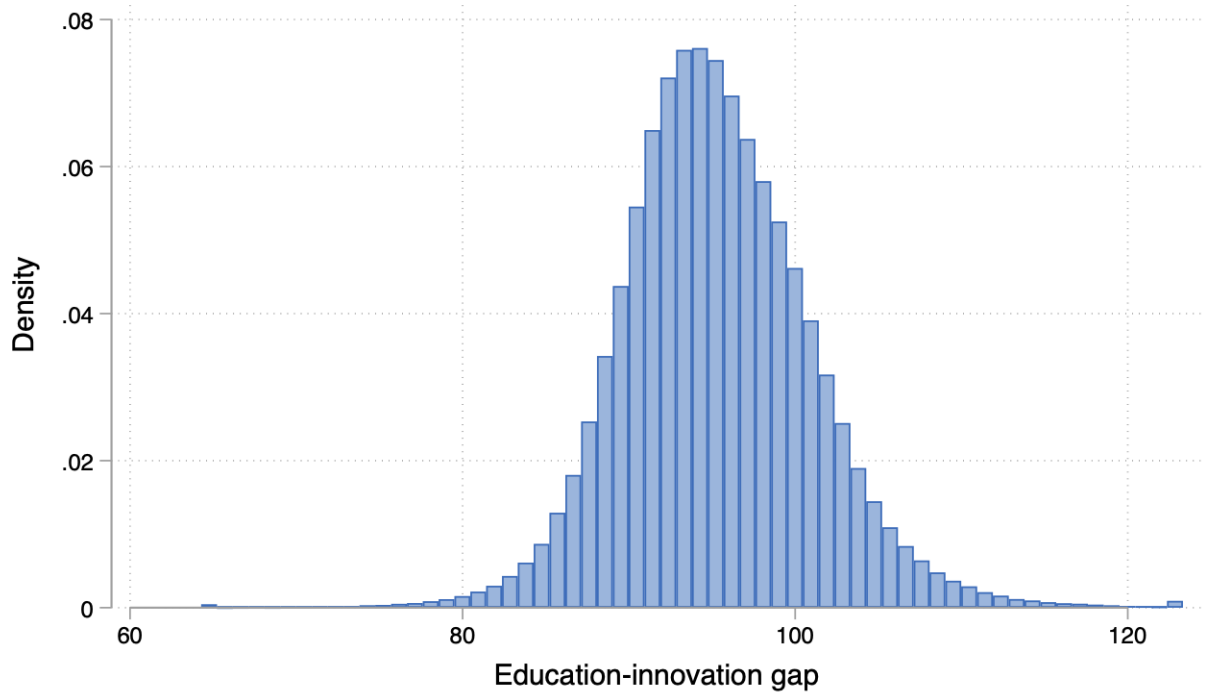
Note: Composition across macro fields, for all courses included in a sample of school catalogs (panel (a)) and for courses included in the syllabi sample (panel (b)).

Figure AV: Course Level Coverage, Course Catalogs and Syllabi Sample



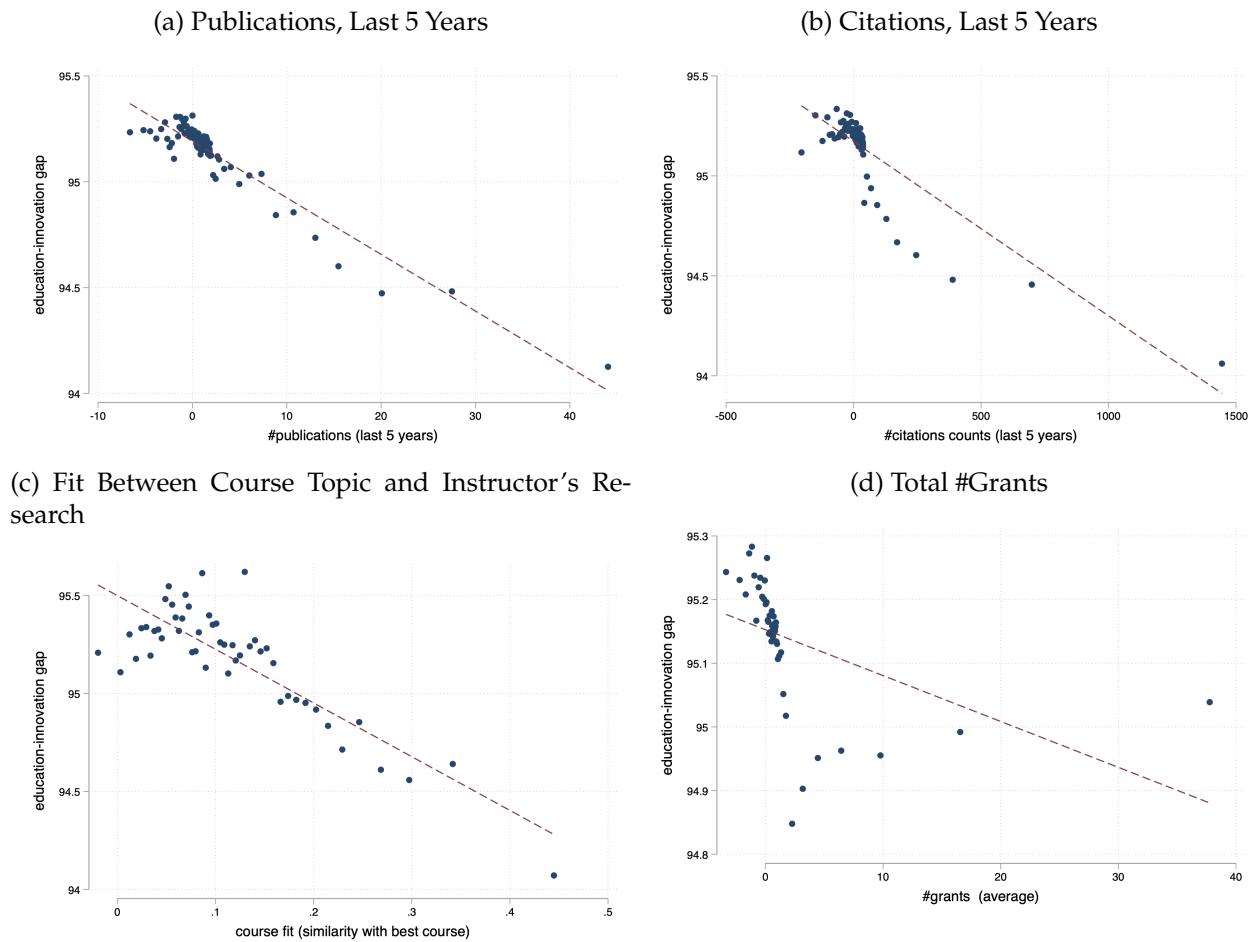
Note: Composition across course levels, for all courses included in a sample of school catalogs (panel (a)) and for courses included in the syllabi sample (panel (b)).

Figure AVI: Education-Innovation Gap: Distribution



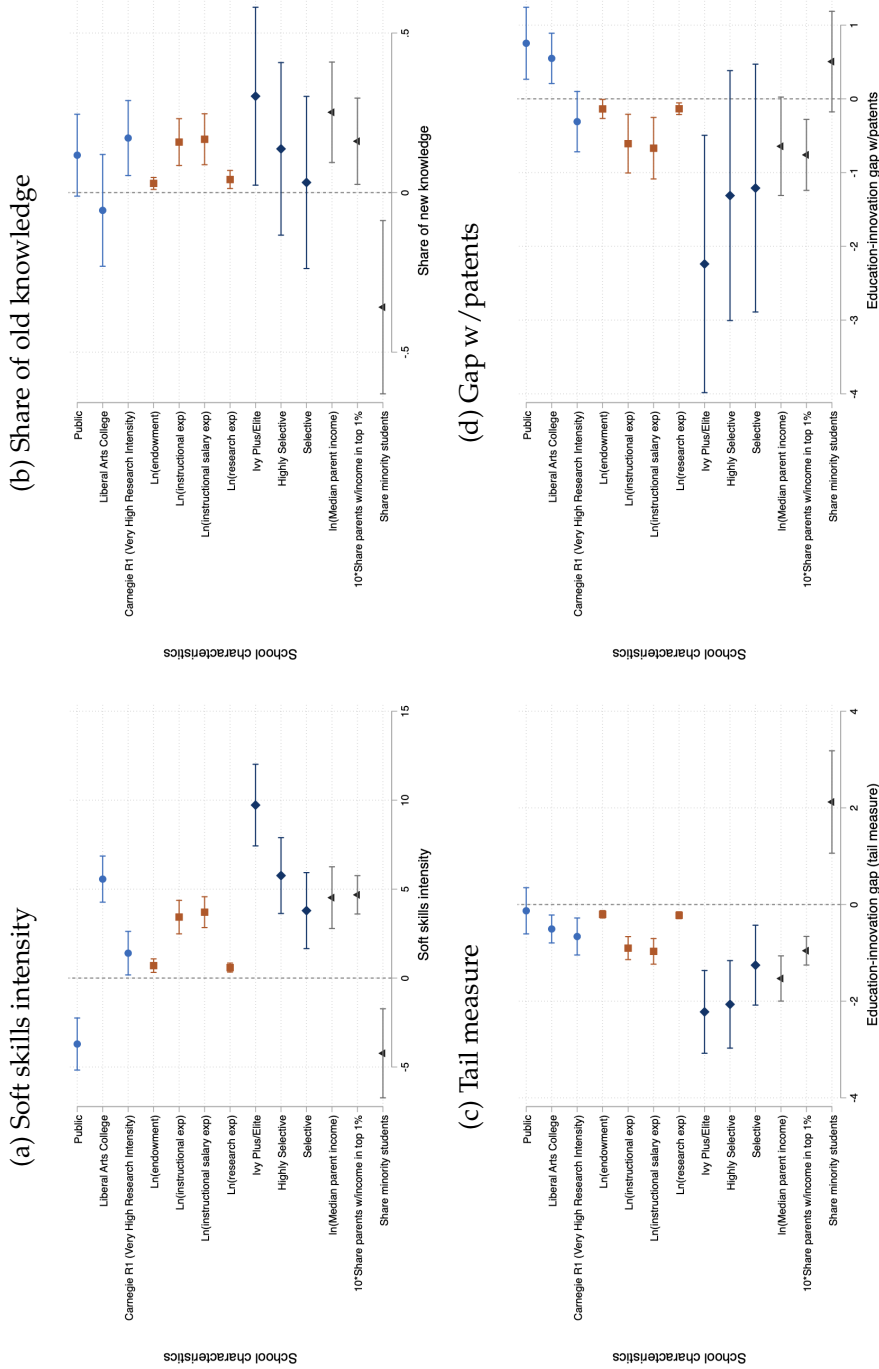
Notes: Histogram of the education-innovation gap.

Figure AVII: Instructors' Research Productivity, Funding, and Fit with The Course The Education-Innovation Gap



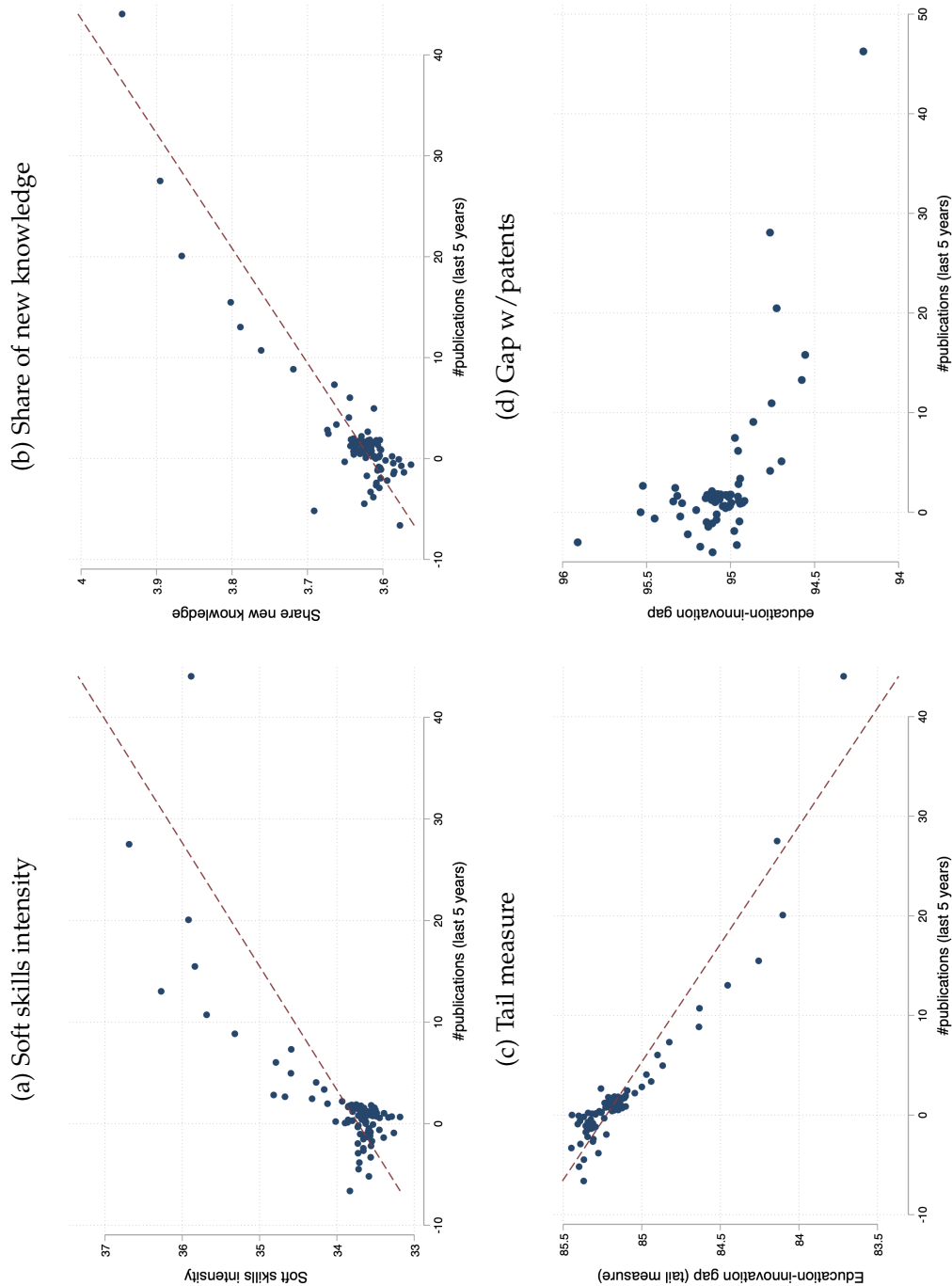
Notes: Binned scatterplot of the gap (vertical axis) and measures of research productivity, funding, and fit between the course topic and the research of the instructor. These measures are the number of publications in the last 5 years (panel a); the number of citations in the last 5 years (panel b); the fit between the instructor's research agenda and the course content, calculated as the cosine similarity between the instructor's publications and the syllabus of the course with the lowest gap among all courses on a given topic (for example, Advanced Microeconomics) across schools in each year (panel c); and the total number of NSF and NIH grants ever received (panel d). All graphs control for field fixed effects.

Figure A.VIII: School Characteristics and Alternative Measures of Course Novelty



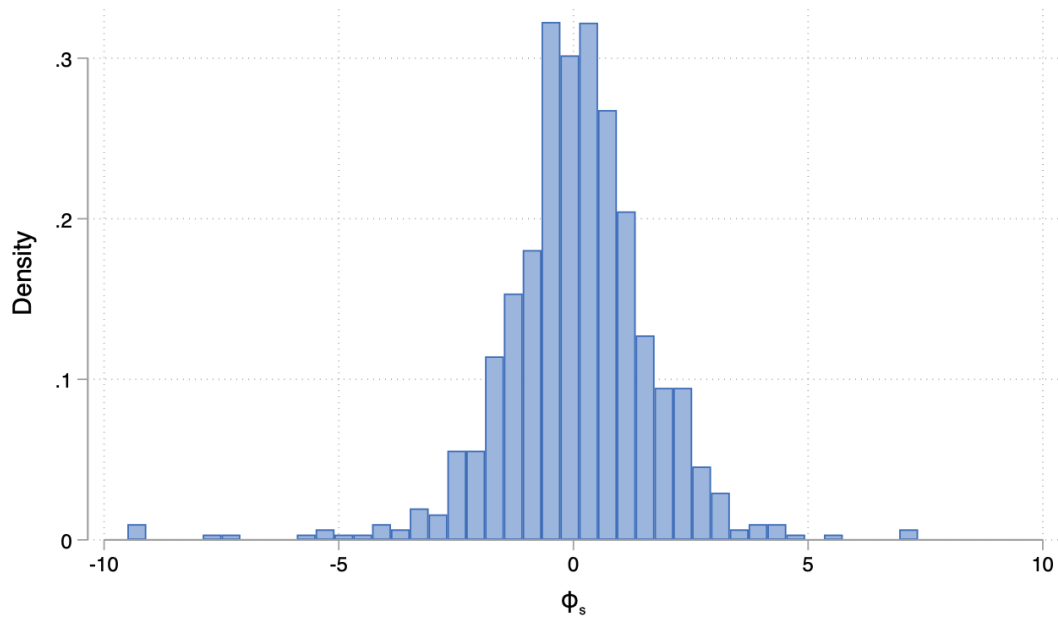
Notes: Point estimates and 95-percent confidence intervals of coefficient β in equation (5), using alternative measures of course novelty: a measure of soft skills intensity, defined as the share of words in the assignment portion of a syllabus which refer to soft skills (panel d), a measure of new knowledge, defined as the share of all new words contained by each syllabus (where new words are knowledge words that are (a) in the top 5 percent of the word frequency among articles published between $t - 3$ and $t - 1$, or (b) used in articles published between $t - 3$ and $t - 1$ but not in those published between $t - 15$ and $t - 13$, panel b); a “tail measure,” calculated for each syllabus by (a) randomly selecting 100 subsamples containing 20 percent of the syllabus’s words, (b) calculating the gap for each subsample, and (c) selecting the 5th percentile of the corresponding distribution (panel c); and the education-innovation gap calculated using the text of all patents as a benchmark, instead of academic articles (panel d). Each coefficient is estimated from a separate regression, with the exception of selectivity tiers (Ivy Plus+/Elite, Highly Selective, Selective) which are jointly estimated. Endowment, expenditure, and share minority information refers to the year 2018 and is taken from IPEDS. Estimates are obtained pooling syllabi data for the years 1998 to 2018. Standard errors are clustered at the school level.

Figure AIX: Instructor Productivity (# Publications) and Alternative Measures of Course Novelty



Notes: Binned scatterplots of a measure of instructor productivity (the number of citations in the prior 5 years) and four alternative measures of course novelty: a measure of soft skills intensity, defined as the share of words in the assignment portion of a syllabus which refer to soft skills (panel a), a measure of new knowledge, defined as the share of all new words contained by each syllabus (where new words are knowledge words that are (a) in the top 5 percent of the word frequency among articles published between $t - 3$ and $t - 1$, or (b) used in articles published between $t - 3$ and $t - 1$ but not in those published between $t - 15$ and $t - 13$, panel b); a "tail measure," calculated for each syllabus by (a) randomly selecting 100 subsamples containing 20 percent of the syllabus's words, (b) calculating the gap for each subsample, and (c) selecting the 5th percentile of the corresponding distribution (panel c); and the education-innovation gap calculated using the text of all patents as a benchmark, instead of academic articles (panel d). Relationships are plotted controlling for field effects.

Figure AX: Distribution of School-Level Gap



Note: Distribution of ϕ_s , the school-level component of the gap, corresponding to $\theta_{s(i)}$ in equation (8).

Table AI: Categorization of Course (Macro-)Fields

Macro-field	Fields
Business	Business, Accounting, Marketing, Public Administration
Humanities	English Literature, Media / Communications, Philosophy, Theology, Criminal Justice, Library Science, Classics, Women's Studies, Journalism, Religion, Sign Language, Liberal Arts, Music, Theatre Arts, Fine Arts, History, Film and Photography, Dance, Anthropology, Japanese, French, Chinese, German, Spanish, Hebrew
Science	Mathematics, Biology, Chemistry, Physics, Earth Sciences, Astronomy, Atmospheric Sciences, Dentistry, Medicine, Nutrition, Nursing, Veterinary Medicine, Natural Resource Management
Engineering	Computer Science, Engineering, Architecture, Agriculture, Basic Computer Skills, Engineering Technician, Transportation
Social Sciences	Psychology, Political Science, Economics, Law, Social Work, Geography, Education, Linguistics, Sociology Education, Criminology
Other	Fitness and Leisure, Basic Skills, Mechanic / Repair Tech, Cosmetology, Culinary Arts, Health Technician, Public Safety, Career Skills, Construction, Military Science

Note: Mapping between the “macro-fields” used in our analysis and syllabi’s “fields” as reported in the OSP dataset.

Table AII: List of Institutions in the Catalog Data

Institution	Institution
Aiken Technical College	Minnesota State University Moorhead
Alabama A & M University	Mississippi College
Alabama State University	Mississippi Community College Board
Alexandria Technical & Community College	Missouri State University
Arkansas Tech University	Mitchell Technical Institute
Asnuntuck Community College	Montgomery College
Bay Path University	Morehead State University
Benedictine University	Mountain Empire Community College
Bentley University	Mountwest Community and Technical College
Bluegrass Community and Technical College	Mt San Antonio College
Briar Cliff University	New Mexico State University-Alamogordo
Brown University	Niagara University
Bryan College	Nichols College
California Baptist University	North Carolina State University
California Lutheran University	North Florida College
California Polytechnic State University	NorthWest Arkansas Community College
Camden County College	Oakwood University
Campbell University	Oral Roberts University
Cardinal Stritch University	Orangeburg Calhoun Technical College
Carlow University	Oregon State University
Catawba College	Oxnard College
Cecil College	Penn State New Kensington
Cedarville University	Plymouth State University
Cerritos College	Princeton University
Coe College	Richland Community College
College for Creative Studies	Robeson Community College
College of Alameda	Rocky Mountain College
College of Southern Nevada	SUNY College at Old Westbury
College of the Siskiyous	SUNY College at Potsdam
Columbia University	SUNY Oneonta
Concordia University Texas	SUNY Orange
Copiah-Lincoln Community College	San Diego Mesa College
County College of Morris	San Diego Miramar College
Dartmouth College	San Diego State University
Daytona State College	Schenectady County Community College
Dominican University	South Arkansas Community College
Drury University	Southern University at New Orleans
Duke University	Spring Arbor University
ENMU-Ruidoso Branch Community College	Spring Hill College
Eastern Nazarene College	Stanford University
Elmhurst College	Suffolk County Community College
Florida Gulf Coast University	Texas Lutheran University
Florida Institute of Technology	The University of Montana
Fresno Pacific University	The University of Texas Rio Grande Valley

(Continued)

Table AII. Continued

Institution	Institution
Frostburg State University	Three Rivers Community College
George Mason University	Trevecca Nazarene University
Georgia State University	Trocaire College
Glendale Community College	University of Akron
Grays Harbor College	University of Central Oklahoma
Green River College	University of Chicago
Grossmont College	University of Colorado Denver
Helena College University of Montana	University of Evansville
Herkimer County Community College	University of Louisville
Hibbing Community College	University of Maine at Presque Isle
Hood College	University of Missouri-St Louis
Hudson County Community College	University of North Carolina at Chapel Hill
Indiana University-Northwest	University of North Dakota
Iowa Central Community College	University of North Texas
Jackson State Community College	University of Notre Dame
Jefferson State Community College	University of Pennsylvania
Kankakee Community College	University of Pittsburgh
Kellogg Community College	University of South Carolina Aiken
Kettering University	University of South Florida-Sarasota-Manatee
Keystone College	University of Wisconsin-River Falls
King's College	Upper Iowa University
Kutztown University of Pennsylvania	Vanderbilt University
Lake Forest College	Virginia Highlands Community College
Las Positas College	Wayne State College
Lassen Community College	Weber State University
Leeward Community College	Webster University
Lincoln University	Wenatchee Valley College
Long Beach City College	Wentworth Institute of Technology
Los Medanos College	Wesleyan University
Louisiana State University-Shreveport	Western Colorado University
MacMurray College	Western Dakota Technical Institute
Marian University	William Jewell College
Marian University	William Woods University
Marietta College	Yale University
Martin Luther College	Youngstown State University
Martin Methodist College	Yuba College
Millsaps College	

Note: List of schools for which we collected course catalog data.

Table AIII: Characteristics of Schools In and Out of Catalog Data

	Mean for Institutions In the Sample # Institutions = 158	Mean for Institutions Out of the Sample # Institutions = 1,956	<i>t</i> -statistics	<i>p</i> -values
ln Expenditure on instruction (2013)	8.693	8.601	-1.725	0.085
ln Endowment per capita (2000)	6.857	6.483	-1.304	0.193
ln Sticker price (2013)	9.197	9.153	-0.520	0.603
ln Avg faculty salary (2013)	8.890	8.850	-1.897	0.058
ln Enrollment (2013)	8.708	8.634	-0.685	0.494
Share Black students (2000)	0.109	0.112	0.153	0.879
Share Hispanic students (2000)	0.063	0.065	0.183	0.855
Share alien students (2000)	0.025	0.022	-1.030	0.303
Share grad in Arts & Humanities (2000)	7.581	7.958	0.382	0.703
Share grad in STEM (2000)	14.861	14.050	-0.772	0.440
Share grad in Social Sciences (2000)	21.068	19.202	-1.342	0.180

Note: Balance test of universities in and out of the catalog sample.

Table AIV: Alternative Measures of Novelty and Student Outcomes

	Income (College Scorecard)			Income (Chetty et al., 2020)					
	Grad rate (1)	Mean (2)	$P_y \leq 33$ pctile (3)	Median (4)	Mean (5)	$P(\text{top } 20\%)$ (6)	$P(\text{top } 10\%)$ (7)	$P(\text{top } 5\%)$ (8)	$P(\text{top } 20\% P_y \leq 20 \text{ pctile})$ (9)
Panel (a): Share new knowledge, no controls									
Gap (sd)	0.0424*** (0.0081)	0.0594*** (0.0103)	0.0678*** (0.0121)	0.0499*** (0.0101)	0.0755*** (0.0137)	0.0338*** (0.0066)	0.0303*** (0.0053)	0.0226*** (0.0040)	0.0310*** (0.0062)
Mean dep. var.	0.5692				763	0.3694 763	0.2082 763	0.1143 763	0.2945 763
N	15683	3793	3566	3793	763	763	763	763	763
# schools	761	760	734	760					
Panel (b): Share new knowledge, with controls									
Gap (sd)	0.0040 (0.0036)	0.0034 (0.0049)	0.0027 (0.0057)	0.0018 (0.0051)	0.0109** (0.0045)	0.0048 (0.0032)	0.0041* (0.0021)	0.0032* (0.0017)	0.0004 (0.0032)
Mean dep. var.	0.5816	10.8281	10.7605	10.7096	718	0.3710 718	0.2100 718	0.1159 718	0.2957 718
N	11471	1996	1843	1996	718	718	718	718	718
# schools	733	727	701	727					
Panel (c): Tail measure, no controls									
Gap (sd)	-0.0503*** (0.0090)	-0.0643*** (0.0105)	-0.0714*** (0.0119)	-0.0580*** (0.0101)	-0.0882*** (0.0125)	-0.0393*** (0.0056)	-0.0336*** (0.0050)	-0.0245*** (0.0036)	-0.0385*** (0.0056)
Mean dep. var.	0.5692				763	0.3694 763	0.2082 763	0.1143 763	0.2945 763
N	15683	3793	3566	3793	763	763	763	763	763
# schools	761	760	734	760					
Panel (d): Tail measure, with controls									
Gap (sd)	-0.0023 (0.0034)	-0.0123*** (0.0043)	-0.0166*** (0.0047)	-0.0137*** (0.0049)	-0.0194*** (0.0048)	-0.0113*** (0.0027)	-0.0089*** (0.0023)	-0.0057*** (0.0016)	-0.0121*** (0.0030)
Mean dep. var.	0.5816	10.8281	10.7605	10.7096	718	0.3710 718	0.2100 718	0.1159 718	0.2957 718
N	11471	1996	1843	1996	718	718	718	718	718
# schools	733	727	701	727					
Panel (e): Gap w/patents, no controls									
Gap (sd)	-0.0232*** (0.0068)	-0.0323*** (0.0116)	-0.0434*** (0.0122)	-0.0282*** (0.0099)	-0.0404*** (0.0138)	-0.0144** (0.0067)	-0.0140** (0.0059)	-0.0120*** (0.0042)	-0.0146** (0.0064)
Mean dep. var.	0.5692				763	0.3694 763	0.2082 763	0.1143 763	0.2945 763
N	15683	3793	3566	3793	763	763	763	763	763

(Continued)

Table AIV. Continued

	Grad rate	Mean	$P_y \leq 33$ pctile	Median	Mean	P(top 20%)	P(top 10%)	P(top 5%)	$P(\text{top } 20\% P_y \leq 20 \text{ pctile})$
# schools	761	760	734	760					
Panel (f): Gap w/patents, with controls									
Gap (sd)	-0.0049 (0.0032)	-0.0003 (0.0038)	-0.0023 (0.0044)	-0.0007 (0.0042)	-0.0039 (0.0046)	0.0004 (0.0025)	-0.0015 (0.0020)	-0.0023* (0.0012)	-0.0014 (0.0029)
Mean dep. var.	0.5816	10.8281	10.7605	10.7096	718	0.3710	0.2100	0.1159	0.2957
N	11471	1996	1843	1996	718	718	718	718	718
# schools	733	727	701	727					
Panel (g): Soft skills intensity, no controls									
Gap (sd)	0.0982*** (0.0065)	0.0935*** (0.0091)	0.0966*** (0.0113)	0.0818*** (0.0085)	0.1125*** (0.0115)	0.0497*** (0.0052)	0.0394*** (0.0044)	0.0293*** (0.0035)	0.0521*** (0.0053)
Mean dep. var.	0.5692	3793	3566	3793	763	0.3694	0.2082	0.1143	0.2945
N	15683	760	734	760	763	763	763	763	763
# schools	761	760	734	760					
Panel (h): Soft skills intensity, with controls									
Gap (sd)	0.0116*** (0.0034)	0.0172*** (0.0052)	0.0096 (0.0068)	0.0209*** (0.0058)	0.0125** (0.0057)	0.0103*** (0.0031)	0.0028 (0.0027)	0.0007 (0.0020)	0.0119*** (0.0038)
Mean dep. var.	0.5816	10.8281	10.7605	10.7096	718	0.3710	0.2100	0.1159	0.2957
N	11471	1996	1843	1996	718	718	718	718	718
# schools	733	727	701	727					

Note: OLS estimates of the coefficient δ in equation (9). The variable Gap (sd) is a school-level education-innovation gap (estimated as $\theta_{s(i)}$ in equation (8)), standardized to have mean zero and variance one. The dependent variable are graduation rates (from IPEDS, years 1998-2018, column 1); the log of mean student incomes from the College Scorecard, for all students (column 2) and for students with parental income in the bottom tercile (column 3); the log of median income from the College Scorecard (column 4); the log of mean income for students who graduated between 2002 and 2004 (from Chetty et al. (2020), column 5); the probability that students have incomes in the top 20, 10, and 5 percent of the national distribution (from Chetty et al. (2020), columns 6-8); and the probability that students with parental income in the bottom quintile reach the top quintile during adulthood (column 9). Columns 1-4 in panels a and b control for year effects. All columns in panel b control for control (private or public), selectivity tiers, and an interaction between selectivity tiers and an indicator for R1 institutions according to the Carnegie classification; student-to-faculty ratio and the share of ladder faculty; total expenditure, research expenditure, instructional expenditure, and salary instructional expenditure per student; the share of undergraduate and graduate enrollment and the share of white and minority students; an indicator for institutions with admission share equal to 100, median SAT and ACT scores of admitted students in 2006, and indicators for schools not using either SAT or ACT in admission; the share of students with majors in Arts and Humanities, Business, Health, Public and Social Service, Social Sciences, STEM, and multi-disciplinary fields; and the natural logarithm of parental income. Bootstrapped standard errors in parentheses are clustered at the school level. * ≤ 0.1 , ** ≤ 0.05 , *** ≤ 0.01 .

Appendix B Dataset Construction

B.1 Syllabi

We obtained data on of university and college syllabi from the Open Syllabus Project (OSP).³¹ The dataset includes nearly 7 million syllabi, collected from 7,365 institutions across the world. OSP provided us with basic information on each syllabus, the full text, and the list of references (papers, textbooks, articles, etc.) included in each syllabus, for a total of 1.8 million unique titles.

We use the following variables from the OSP database:

- `id`: The unique identifier assigned to each syllabus.
- `text`: The text of the syllabus.
- `textmd5`: The md5sum of the text, which can also be used as unique identifier.
- `language`: The language of the document.
- `year`: The academic year when the syllabus was taught.
- `fieldname`: The name of the academic field most associated with the syllabus.
- `institutionid`: The unique identifier for the institution of the course.
- `unitid`: The IPEDS identifier for the institution.
- `countrycode`: The ISO 3166-1 alpha-2 code of the country the syllabus was taught in.
- `institutionname`: The name of the institution of the course.

In the paper, we focus on syllabi that satisfy the following criteria.

- (i) Taught in a four-year, non-online university based in the US (`countrycode` equal to "US") with at least 100 syllabi in the data;
- (ii) Taught in English;
- (iii) Taught between 1998 and 2018;
- (iv) With a word length between 20 and 10,000.

The number of syllabi we keep in each step, and the associated syllabi characteristics, are shown in Table BV.

³¹<https://opensyllabus.org>

Table BV: Summary Statistics of Open Syllabus Project

	# of records	Syllabus word length (raw)	Syllabus word length ("knowledge content")
Original data	6,852,971		
Keep syllabus based in the United States (Syllabus language is English)	3,995,483		
Keep syllabus from four-year university	1,951,933	2,725.41	1,435.09
Year from 1998 to 2018	1,937,284	2,732.09	1,436.77
Extracted syllabus length must be in [20, 10000]	1,901,367	2,279.66	1,057.35
Number of syllabi per institution larger than 100	1,882,224	2,274.55	1,056.77
Remove syllabus from online-only univer- sities	1,752,795	2,218.08	1,010.82

Note: Counts of syllabi, raw word length, and knowledge content (number of words remaining after the cleaning process is complete).

Course catalog data To complement the syllabi data and determine selection patterns into it, we also obtained the entire list of course offerings from university catalogs, for a sample of US institutions. We begin by randomly selecting 10% of all universities in our sample (212 universities). Then, we manually search and download electronic copies (usually in the PDF format) of university catalogs for those universities for all years available, which list all courses offered in that institution and year. Out of the 212 universities selected, 161 have at least one catalog available. We downloaded and processed a total of 2,348 catalogs for these 161 universities (14.5 catalogs per university). Due to random selection, these schools are representative of the full sample on the basis of standard school-level characteristics. A balance test of characteristics between the full sample and the catalog sample is shown in Table [AIII](#).

University catalog data provide the following information: course code, course name, and course level (classified into Basic, Advanced, and Graduate). Some course catalogs also provide a brief course description.

B.1.1 Extracting A Course's Content From Its Syllabus

The full text of a syllabus is contained in the variable `text` of the OSP database. To transform text into usable content, we (i) clean it by removing html language left over from web scraping or correcting obvious errors from OCR procedures; (ii) identify the various sections of the syllabus in it; and (iii) remove text unrelated to content (e.g., course policy, absence policy, accommodation rules, etc.). We now explain these steps in more detail.

B.1.2 Cleaning The Raw Text

To clean the text of each syllabus, we proceed as follows.

- (i) We use the Unidecode Python Package³² to convert Unicode text into ASCII text. This includes legacy code that does not support Unicode, non-Roman names on a US keyboard, and ASCII approximations for symbols and non-Latin alphabets.
- (ii) We remove browser information, often present in the header of a syllabus, by searching for keywords such as "Internet Explorer", "Newer Browser", "JavaScript Enabled", "Cookies Are", "Download Info", "Login", "Log In", "Print", and "Search".

B.1.3 Identifying Syllabi Sections

The average syllabus contains a set of sections, some of which are relevant for our analysis. The relevant sections include: instructor and course information (such as code, course level, and title); course description, requirements, and objectives; an outline; homework, exams, and other evaluation methods; and other policies. A syllabus often also includes other information that we do not use in the analysis and, as such, we want to remove. These include the honor code, policies related to disability, classroom laptop and cellphone policies, and many others.

To parse among sections, we developed a supervised algorithm based on a set of section title keywords. The algorithm identifies a section type by searching through a set of keywords belonging to each category. Table **BVI** provides section types along with the corresponding keywords.

Using these keywords, the algorithm separates the text into different sections of the syllabus by combining keywords with the formatting rules of each syllabus. In Figure **BXI**, we use part of a syllabus as an example to present our process step-by-step.

³²<https://pypi.org/project/Unidecode/>

Table BVI: Section Title Keywords List

Section type	Keywords
<i>Course Description</i>	Syllabi, Syllabus, Title, Description, Method, Instruction, Content, Characteristics, Overview, Tutorial, Intro, Abstract, Methodologies, Summary, Conclusion, Appendix, Guide, Document, Module, Introduction, Approach, Lab, Background
<i>Requirements</i>	Requirement, Applicability, Required
<i>Objectives</i>	Objectives, Achievement, Outcome, Motivation, Purpose, Statement, Skill, Competency, Performance, Goal
<i>Outline</i>	Outline, Schedule, Timeline, Guideline
<i>Materials</i>	Text, Material, Resource, Recommend, Reference, Book, Calendar, Textbook, Guidebook
<i>Instructor information</i>	Instructor, About, Email, Phone, Contact, Professor, Staff, Faculty, Information
<i>Projects, homework, papers, and exams</i>	Personal, Total, Individual, Exercise, Essay, Submission, Assign, Homework, Paper, Final, Examing, Midterm, Term, Semester, Proposal, Application, Demonstration, Program, Task, Report, Pracical, Drafting, Project, Plan, Deadline, Makeup, Advising, Advisor, Survey, Assignment, Planning, Practice, Group, Participation, Team, Research, Activity, Complaint, Design, Analysis, Strategy, Procedure, Working, Work, Exam, Examination, Training, Professional, Test, Case, Discussion, Grade, Presentation, Quiz, Essay, Layout, Sample, Rewrite
<i>Grades</i>	Assessment, Point, Scope, Evaluation, Record, Grading, Composition, Review
<i>Other Policies</i>	Academic, Justice, Administration, Rule, Discipline, Disclaimer, Regulation, Standard, Affair, Dishonesty, Plagiarism, Misconduct, Offence, Medical, Absent, Absence, Trip, Religious, Observance, Ttendance, Honesty, Origination, Originator, Help, Technology, Attendance, Accessing, Service, Oppotunity, Administrative, Accommodation, Support, Policy, Right, Responsibility, Disability, Weather, Integrity, Copyright
<i>Notes</i>	Remark, Notice, Additional, Acknowledgement, Absolutely, Absolute, Important, Note, Cannot, Can, Must, Should, Will, Please, No
<i>Other Words</i>	Course, Lecture, Catalog, Campus, Commuity, Class, Classroom, College, Univerity, Discussion, Seminar

Note: Keywords used to identify the corresponding section types in a syllabus of a syllabus. In the implementation, we use both the singular and plural versions of each term.

1. For each syllabus, we identify the section titles based on the word list described above and the formatting features. We mark all cases when the section title phrases appear as all uppercases or consecutive initial capital letters using regular expressions.
 - In Figure **BXI**, underlined sentences satisfy the features of a section title, such as “Course Description”.
2. We divide the syllabus into parts and we use Arabic numerals to mark them out. Finally, we select sections with relevant titles and extract the cleaned text.
 - In Figure **BXI**, we focus on highlighted sections, such as “Course Objective,” “Prerequisites,” and “Text”.

B.1.4 Extracting Additional Information

Instructor Names To extract the name of the instructor from each syllabus, we build a neural network model based on the BiLSTM-CNNs-CRF model for named entity recognition (NER).³³ The training/test dataset is built via the following three steps:

- (i) We select syllabi that contain at least one keyword such as “Doctor”, “Doctors”, “Dr”, “Professor”, “Prof”, “Instructor”, “Instructors”, “Tutor”, “Tutors” in the first 3,500 characters.
- (ii) We use the Spacy³⁴ package to identify whether the words following those keywords are names of people (entity label is “PERSON”).
- (iii) We process the syllabus text sentence by sentence as the training and test data of the model.

We also apply a few additional filters: (a) we remove single letter names; (2) all the words in the name are required to appear in the Python Library *English First and Last Names Data Set*³⁵; (c) after the first two filters, we only keep the first instructor name. With this algorithm, we are able to assign an instructor name to 86.23% of all syllabi. The out-of-sample precision of this algorithm is 85.18%.

Course Level: Basic, Advanced, Graduate To assign a course level (basic undergraduate, advanced undergraduate, and graduate) to each syllabus, we trained a Natural Language Processing (NLP) algorithm. Our training sample consists of 56,831 syllabi taught in universities for which we

³³BiLSTM-CNNs-CRF model for named entity recognition (NER), Ma and Hovy (2016).

³⁴<https://spacy.io/>

³⁵<https://github.com/philipperemy/name-dataset>

have catalog information, for whom we can manually code the course levels. Specifically, in the catalog data, we label a course as basic undergraduate if the course belongs to the undergraduate catalog of a university and the course code starts with 1 or 2; we label the course as advanced undergraduate if the course belongs to the undergraduate catalog and the course code starts with 3 or 4; finally, we label the course as graduate if the course belongs to the graduate catalog or the first digit of the course code is larger than 4. We link syllabi to catalog information using institution and course code. Once we have obtained course levels for these syllabi, we use course levels as labels and the text of each syllabus as input in the training model. The model we use is Distilled BERT³⁶ (Sanh et al., 2019), accessed via the transformers library.³⁷ The out-of-sample prediction precision is 85.04%.

Course code Our data extraction process allows us to obtain the course code corresponding to each syllabus. However, these courses are institution-specific and often vary over time. To be able to identify courses of the same level (e.g., basic undergraduate) covering the same topic (e.g., Principles of Microeconomics), both within and across schools, we proceed as follows. First, we construct a unified within-school course code using the raw course code and the course name. We do so as follows: (a) we remove the punctuations and multiple whitespaces from codes and names; (b) for course names, we further remove stop-words and isolate the course stem name (the common base form of the words). We then consider two courses as sharing a course code if (a) they share the same name and code (b) they share the same name, even if the course code changes over time. This procedure accounts for the possibility that the course code system might have changed within a school over time.

Once we have a disambiguated identifier for courses within the same school, we assign courses a cross-school identifier. Specifically, we assign two courses the same cross-school identifier if they share the same standardized course name.

B.2 References and Recommended Readings in Each Syllabus

In addition to syllabi text and metadata, OSP provided us with two additional datasets: “Matches” and “Catalog.” “Matches” allows us to link syllabi to records in “Catalog;” “Catalog” is the set of 1.8 million bibliographic records assigned to at least one syllabus. We use the following variables from the “Matched” dataset:

- `MatchID`: The unique identifier of the match;

³⁶<https://arxiv.org/abs/1910.01108>

³⁷<https://huggingface.co/transformers/index.html>

- ID: The id of the syllabus;
- WorkID: The id of the catalog record.

We use the following variables from the “Catalog” dataset:

- WorkID: The id of the catalog record.
- Publicationtype: The type of publication (“journal” or “book”).
- Publicationyear: The year of publication.

B.2.1 Syllabi Field

The OSP database classifies syllabi into one of 69 fields. For some of our analyses, we group these into macro-fields. The grouping is illustrated in Table [BVII](#).

Table BVII: Categorization of Course (Macro-)Fields

Macro-field	Fields
Business	Business, Accounting, Marketing, Public Administration
Humanities	English Literature, Media / Communications, Philosophy, Theology, Criminal Justice, Library Science, Classics, Women's Studies, Journalism, Religion, Sign Language, Liberal Arts, Music, Theatre Arts, Fine Arts, History, Film and Photography, Dance, Anthropology, Japanese, French, Chinese, German, Spanish, Hebrew
Science	Mathematics, Biology, Chemistry, Physics, Earth Sciences, Astronomy, Atmospheric Sciences, Dentistry, Medicine, Nutrition, Nursing, Veterinary Medicine, Natural Resource Management
Engineering	Computer Science, Engineering, Architecture, Agriculture, Basic Computer Skills, Engineering Technician, Transportation
Social Sciences	Psychology, Political Science, Economics, Law, Social Work, Geography, Education, Linguistics, Sociology Education, Criminology
Other	Fitness and Leisure, Basic Skills, Mechanic / Repair Tech, Cosmetology, Culinary Arts, Health Technician, Public Safety, Career Skills, Construction, Military Science

Note: Mapping between the “macro-fields” used in our analysis and syllabi “fields” as reported in the OSP database.

Figure BXI: Dividing A Syllabus Into Sections: An Example

Econ 561a	Yale University	Fall 2005	
Prof. Tony Smith (Part I)		Prof. Michael Keane (Part II)	
Syllabus for	COMPUTATIONAL METHODS FOR ECONOMIC DYNAMICS		ECON 561a

Course Objectives:
Most of the dynamic economic models used in modern quantitative research in economics do not have analytical (closed-form) solutions. For this reason, the computer has become an indispensable tool for conducting research in dynamic economics. The goal of this two-part course is precisely to teach students computational tools for conducting numerical analysis of dynamic economic models. The focus of the first half of the course, taught by Prof. Tony Smith, is on solving dynamic programming problems and on computing competitive equilibria of dynamic economic models. The first half of the course also provides an introduction to some of the basic tools of numerical analysis, including minimization, root-finding, interpolation, function approximation, and integration. The focus of the second half course, taught by Prof. Michael Keane, is on solving and estimating discrete-choice dynamic programming models of economic behavior. Taken together, the two halves of the course provide students with a thorough introduction to the numerical analysis of dynamic economic models in both microeconomics and macroeconomics.

Contact Information (Prof. Tony Smith)
Office: 28 Hillhouse, Room 306 Office phone: (203) 432-3583
Email address: tony.smith@yale.edu Course Web site: www.econ.yale.edu/smith/econ561a
Office hours: Thursdays from 10AM–noon, or by appointment

Class Meetings:
The course meets on Mondays and Wednesdays from 2:30PM to 3:50PM in a room to be determined.

Prerequisites:
This course is designed for graduate students in economics who have taken first-year graduate courses in microeconomics, macroeconomics, and econometrics. No prior knowledge of either numerical methods or computer programming is assumed, but some familiarity with a programming language would prove helpful.

Texts:
The required textbook for this course is:
Numerical Recipes in Fortran 77: The Art of Scientific Computing, Second Edition (Volume 1 of Fortran Numerical Recipes) by William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery (Cambridge University Press, 1992). This book, as well as its (optional) companion Numerical Recipes in Fortran 90: The Art of Parallel Scientific Computing, Second Edition (Volume 2 of Fortran Numerical Recipes), is available online at: www.library.cornell.edu/nr/.
Other (optional) books that students might find useful are:

- Numerical Methods in Economics by Kenneth L. Judd (MIT Press, 1998).
- Handbook of Computational Economics (Volume 1), edited by Hans M. Amman, David A. Kendrick, and John Rust (North-Holland, 1996).
- Computational Methods for the Study of Dynamic Economies, edited by Ramon Marimon and Andrew Scott (Oxford University Press, 1999).
- Dynamic Economics: Quantitative Methods and Applications by Jérôme Adda and Russell Cooper (MIT Press, 2003).
- Applied Computational Economics and Finance by Mario J. Miranda and Paul L. Fackler (MIT Press, 2002).

Grading:
The course grade will be based on two (equally-weighted) projects, one for the first part of the course and one for the second part of the course. Each project consists of writing a program in Fortran to solve an assigned problem. Students must submit their code as well as a brief (roughly five pages) description of their numerical findings. The first project will involve solving for the competitive equilibrium of a dynamic macroeconomic model; the second project will involve solving and estimating a discrete-choice dynamic programming model. Fortran is the language of choice for most researchers in computational economics; requiring that the code for the projects be written in Fortran will help students to become proficient in this powerful and useful language. The first project is due on Monday, November 14 and the second project is due at the end of the semester. Occasional short programming problems may also be assigned as the course proceeds. The purpose of these assignments is to help students develop the skills they need to complete the projects; these assignments will not be graded.

Approximate Schedule of Lectures (Part I)
I. INTRODUCTION
Lecture 1 Introduction to numerical dynamic programming (built around the stochastic growth model and the Aiyagari (1994) model). General considerations in numerical analysis: convergence, roundoff error, truncation error. Numerical differentiation.
Readings:

- Aiyagari, S.R. (1994), “Uninsured Idiosyncratic Risk and Aggregate Saving,” Quarterly Journal of Economics 109, 659–684.
- Numerical Recipes: Chapters 1 and 5.7
- Judd: Chapters 1, 2, and 7.7

II. BASIC NUMERICAL METHODS
Lecture 2 Root-finding in one or more dimensions: bisection, secant method, Newton’s method, fixed-point iteration, Gauss-Jacobi, Gauss-Seidel, Brent’s method.
Readings:

- Numerical Recipes: Chapter 9

.....

Note: Example of a syllabus from OSP, in its original format. Subsections are identified using the algorithm described in this appendix.

B.3 Academic Publications

To construct the education-innovation gap, we collected a large sample of academic articles from top journals. We describe here how this sample is defined, constructed, and collected.

B.3.1 List of Top Journals

We begin by compiling a list of top academic journals within each discipline. Our starting point is the Journal Citation Reports (JCR), an annual report published by Thomson Reuters (formerly ISI) to provide citation and publication data of academic journals in the science and social science fields by means the impact factor.³⁸ We consider as top journals those that were ranked within the top ten of their respective field at least once since their establishment. This leaves us with 3,962 journals in 223 fields.

B.3.2 Collecting Academic Articles

Having compiled a list of top journals, we collect information on all the articles ever published in these journals. These data come from Scopus, an Elsevier-owned database containing abstracts and citations of academic articles.³⁹ To extract the metadata of journal articles, we accessed Scopus's API and searched for the ISSN of each journal ("ISSN(0022-1082)"). We then extracted all the metadata of all articles of the relative journal, for all available years. We focus our attention on the following variables:⁴⁰:

- `EID`: electronic ID, used as the unique identifier of each article;
- `title`: title of the article;
- `ISSN`: ISSN of publisher;
- `coverdate`: publication date;
- `description`: abstract;
- `authkeywords`: keywords.

Our initial search yielded 20,779,713 articles, of which we discarded those without an abstract.

³⁸<https://jcr.clarivate.com/>

³⁹<https://www.scopus.com>

⁴⁰The full list of variables available through Scopus is available at <https://dev.elsevier.com/guides/ScopusSearchViews.htm>

B.3.3 Data Cleaning

The main information from academic articles that we use in our analysis is the abstract, contained in the variable `description` of the SCOPUS database. We further clean the content of this variable to remove copyright disclaimers, usually present at the beginning or at the end of each abstract and unrelated to content. We do this using keyword recognition techniques. Starting from the first sentence of an abstract, we remove it if it contain at least one among the words “copyright”, “©”, “published”, “publisher”, “all right”, “all rights reserved”; we repeat this procedure until the first sentence does not contain any of these words. We then repeat the same procedure starting from the last sentence.

B.4 Research Productivity

We use information from Microsoft Academic (MA) to measure the research productivity of all people listed as instructors in the syllabi. We download these data from Microsoft Academic Knowledge Graph (MAKG).⁴¹ MAKG is a large resource-description framework (RDF) knowledge graph with over eight billion triples containing information about scientific publications and related entities, including authors, institutions, journals, and fields of study. The data set is based on the Microsoft Academic Graph and licensed under the Open Data Attributions license. For each researcher, Microsoft Academic lists publications, working papers, other manuscripts, and patents, together with the counts of citations to each of these documents. Due to differences in counting citations, Microsoft Academic citations do not necessarily match those from similar services such as Web of Science or Google Scholar; the correlations between all these services’ citations numbers, though, are very high.

We link instructor records from the text of the syllabi to Microsoft Academic records using names, a person’s history of institutions, and research fields. After restricting the sample of syllabi to those whose instructor can be matched to a unique MA researcher, we are able to successfully match 38.93% of all instructors in the syllabi data. Out of all syllabi with an instructor name (2,433,683), 34.77% (846,287) of them can be matched to MA, while 25.86% (629,365) of them have duplicated matching. In our selected syllabi sample, the matching ratio is 38.93% (= 682,286 / 1,752,795).

⁴¹We downloaded the data based on the Microsoft Academic Graph data as of 2020-05-29 from <https://zenodo.org/record/3936556#.YFndr2Qza3J>

B.5 Patents

We obtain data on patents from the publicly available Patent Full-Text Database (PatFT)⁴² of the US Patent and Trademark Office (USPTO). This database provides records for all patents ever issued since 1976. We used a web crawler to collect the text content of patents over this period, which includes patents with numbers ranging from 3,850,000 to 10,279,999. We use the following variables for each patent record:

- `PatentNumber`: The unique identifier assigned to each patent record;
- `Abstract`: The abstract in each patent filings;
- `Year`: The year that the patent was issued;
- `Class`: The International Patent Classification (IPC) assigned to each patent.

B.6 National Science Foundation and National Institute of Health Grants

We collected information on grants awarded by the National Science Foundation (NSF)⁴³ and the National Institutes of Health (NIH)⁴⁴ to construct measures of research investment and productivity. These data are provided directly by the respective organizations; the versions used in the paper were accessed on May 25, 2021.

The NSF grant data include 480,633 grants with effective starting year ranging from 1960 to 2022. The NIH grant data include 2,566,358 grants with effective year ranging from 1978 to 2021. Both NSF and NIH grant data contain information on the host institution (institution name, country, state, and city) and the investigator (investigator name and role). In the NSF data, investigators can be listed under four figures: principal investigator (PI), co-PI, former PI, and former co-PI. In the NIH data, they can be listed under two figures: contact and non-contact.

B.6.1 Linking NSF/NIH Institutions to Syllabi Institutions

To link grants to institutions in the syllabi data and IPEDS, we use information on the institution's name and location (country, state, and city). To do so, we first perform an exact match using institution names as listed in the NSF/NIH data and in IPEDS, stripped of punctuations and stop words (including "and" and "the"). Then, for the remaining unmatched NSF/NIH institutions, we conduct a fuzzy matching based on name and location. We require the matching algorithm to meet the

⁴²<http://patft.uspto.gov/netathtml/PTO/index.html>

⁴³<https://www.nsf.gov/awardsearch/download.jsp>

⁴⁴https://exporter.nih.gov/ExPORTER_Catalog.aspx

following two conditions: (1) the two institutions must be in the same city; (2) the fuzzy matching ratio must be larger than a certain threshold (specifically, we use partial ratio and token set ratio in the FuzzyWuzzy Package).⁴⁵ This method sometimes leads us to match a NSF/NIH institution to multiple IPEDS institutions; in this case, we consider the IPEDS institution with the largest average matching ratio .

We are able to match 11.30% (2,402) NSF institutions to IPEDS, covering 82.05% (= 394,383 / 480,633) of all NSF grants. Similarly, we are able to match 6.73% (1,573) NIH schools to IPEDS, covering 66.53% (= 1,707,498/2,566,358) of all NIH grants. The unmatched NSF and NIH institutions are mostly non-academic, private or not-for-profit research institutes.

B.6.2 Linking NSF/NIH Investigators to Instructors

Next, we match grant investigators to course instructors in the syllabus data. We do this via a fuzzy matching algorithm using names. The NSF and NIH data provide different investigator information to be used in the fuzzy matching, so the matching methods differ slightly between the two datasets.

NSF To match NSF investigators to instructors, we first remove duplicates within NSF based on first name, last name, email, and institutions since NSF does not provide investigator unique identifiers. We consider two investigators to be the same person if (1) they share the same email, or (2) they have exactly the same first name and last name in the same school in a certain year. Next, we perform a many-to-one fuzzy matching between NSF investigators and syllabi instructors based on the names and history of institutions the researcher was employed at. We proceed in three steps:

- (i) After removing any punctuations from name strings, we fuzzy-match each NSF investigator name with syllabus instructor names. We calculate matching scores using the Whoswho Package⁴⁶, a Python library for determining whether two names belong to the same person.
- (ii) If a match has a score of 100, we consider it successful. For matches with scores larger than 95 who ever worked at the same school, assign an investigator to one and only one instructor as follows.
 - (a) If an NSF investigator and a set of syllabi instructors have spent some common period of time at the same institution as we can observe it, we link the investigator to the instructor with the maximum matching score.

⁴⁵The package uses Levenshtein Distance to calculate the differences between sequences, and its homepage is <https://github.com/seatgeek/fuzzywuzzy>, and we use a threshold of 80.

⁴⁶<https://github.com/rlieb/whoswho>

- (b) If they have not spent any common period of time at the same institution, we link the investigator to the instructor with the maximum matching score and minimum temporal distance between the time spent at each institution.
- (iii) For matches with matching score larger than 95 but in different schools,
 - (a) If an instructor and an investigator are observed for the same period of time in the two datasets, we choose the match with the maximum matching score.
 - (b) Otherwise, we choose the matching with the maximum matching score and shorter time distance between observed periods between the two datasets.

This procedure leaves us with 232,206 unique investigators, 23.31% ($= 54,118 / 232,206$) of whom can be matched to one syllabus instructor, and corresponding to 44.28% ($= 208,857 / 471,646$) of all grants.

NIH Data from NIH contain investigator unique identifiers, which implies we do not have to remove duplicates. We use these to perform a one-to-one matching between each NIH investigator and a syllabus instructor. We follow the same process as with NSF grant data.

Of all syllabus instructors, 15.46% ($= 68,236 / 441,452$) of them have at least one NSF or NIH grant; these instructors correspond to 2,237 IPEDS schools and 190,738 syllabi, accounting for about 12.82% ($= 190,738 / 1,487,820$) of all syllabi with a listed instructor. This procedure leaves us with 298,687 unique investigators, 10.07% ($= 30,087 / 298,687$) of whom can be matched to one syllabus instructor, corresponding to 17.69% ($= 450,339 / 2,546,123$) of all grants.

B.7 Instructors' Job Titles and Salaries

We were able to collect the salaries of instructors employed at 490 public college and universities in 16 states. As the regulations on the disclosure of public-sector employees' salaries vary across states and over time, the temporal coverage of our data differs across states. Table **BVIII** describes the coverage and source of the salary data data.

Together with the salary data, the job title of each employee is also disclosed. We were able to identify the following titles: assistant professor, associate professor, full professor, lecturer, adjunct professor, clinical professor, professor of practice, and visiting professor. Table **BIX** describes how we assign job titles based on the information available in the data.

Table BVIII: Coverage and Source of Salary and Job Title Data

State	Data available for	Source
CA	2011-2018	https://transparentcalifornia.com/agencies/salaries/
CT	2010-2018	http://transparency.ct.gov/html/searchPayroll.asp
GA	2010-2018	https://open.ga.gov/openga/salaryTravel/index
IA	2009-2018	https://www.legis.iowa.gov/publications/fiscal/salaryBook
IL	2010-2018	https://salary.bettergov.org/
IN	2012-2018	https://gateway.ifionline.org/default.aspx
KS	2012-2018	http://kanview.ks.gov/DataDownload.aspx
MA	2010-2018	https://cthrupayroll.mass.gov/
MD	2012-2018	https://salaries.news.baltimoresun.com/
MI	2014-2018	https://www.mackinac.org/salaries
MN	2011-2018	https://mn.gov/mmb/transparency-mn/payrolldata.jsp
NV	2009-2018	https://transparentnevada.com/
NY	2008-2018	https://www.seethroughny.net/payrolls
OK	2010-2018	https://data.ok.gov/dataset
RI	2011-2018	http://www.transparency.ri.gov/payroll/
WA	2016-2018	http://fiscal.wa.gov/salaries.aspx

Note: States for which instructor salary and job title data is available, together with available year and source.

Table BIX: Assigning Job Titles

Job Title	Definition
Adjunct Professor	Any word of the job title starts with "adjunct", "adj", "temporary", "temporari", "temporar", or "part time".
Clinical Professor	Any word of the job title starts with "clinic" or "clin".
Professor of Practice	Any word of the job title starts with "practic" or "pract".
Visiting Professor	Any word of the job title starts with "visiting" or "visit".
Lecturer	(1) Any word of the job title starts with "lectur", "lect", "instructor", "instruct", "instr", "teacher", or "teach"; (2) AND any word of the job title does not end with "ship"; (3) AND job title is not identified as adjunct professor, clinical professor, professor of practice, and visiting professor.
Professor	(1) Any word of the job title starts with "professor", "prof", or "tenur"; (2) OR any word of the job title includes "tenr trk" or "tenur track"; (3) AND any word of the job title does not end with "profession"; (4) AND job title is not identified as adjunct professor, clinical professor, professor of practice, or visiting professor.
Assistant Professor	Job title is identified as professor; (2) AND any word of the job title starts with "assist", "asst", or "assi".
Associate Professor	Job title is identified as professor; (2) AND any word of the job title starts with "associ", "assoc", or "asso".
Full Professor	Job title is identified as professor; (2) AND detailed job title is not identified as assistant professor or associate professor.

Note: Procedure used to assign job titles to salary records.

Appendix C Calculating The Education-Innovation Gap: Additional Details and A Simulation Exercise

We now explain in detail the process we use to identify the knowledge terms used in our analysis, extract them from the text of syllabi and academic publications, and calculate the gap.

C.1 Extracting Knowledge Terms From Each Document

Dictionary The first step is to build a dictionary, i.e., a list of all such terms. We use the list of all unique words and expressions ever used as a keywords in academic publications. We extract these keywords from the data described in Section B.3.

Term Extraction Next, we convert the text content of each document (syllabi and academic papers) into numerical data for statistical analyses. To do so, our starting point is to clean the text. First, we convert the text of each document into ASCII text using the Unidecode Python Package.⁴⁷ This allows us to handle host legacy code that does not support Unicode, non-Roman names on a US keyboard, and ASCII approximations for symbols and non-Latin alphabets. Next, we convert all capitalized characters to lowercase and use the NLTK Python Toolkit to strip out all non-word text elements, such as punctuation, numbers, and HTML tags. We also remove all occurrences of 280 “stop word”, which include propositions, punctuation, pronouns, and other words that carry little semantic content.⁴⁸

Once we have cleaned the text, we convert it into numerical data using a term-extraction algorithm called NGramMatch. This algorithm performs exact string matching of the text of each document, consisting in N-grams with N ranging from 1 to 7, with the dictionary. To do so, the algorithm extracts N-grams from text, to form a basic term set. Then, it filters out all the terms which cannot be linked to any dictionary entry. In the final set, the algorithm assigns each document a frequency vector based on matched dictionary words.

C.2 A Simulation Exercise

To better understand how the education-innovation gap captures the academic novelty of a syllabus’s content and to illustrate his properties, we perform a simulation exercise. In this simulation, we manually construct a set of syllabi by combining dictionary words that can be found in academic publications. Each syllabus is characterized by a year (t , ranging from 1998 to 2018 to match

⁴⁷<https://pypi.org/project/Unidecode/>

⁴⁸We use the stopwords list using the union of all single letters and Stanford CoreNLP package: <https://github.com/stanfordnlp/CoreNLP/blob/master/data/edu/stanford/nlp/patterns/surface/stopwords.txt>.

our data), a known gap (gap , ranging between 0 and 1) and a parameter governing his style ($style$); we define the latter two parameters in what follows. For each of these syllabi, we calculate the education-innovation gap with the procedure described in the text; we then compare it with the known gap to assess its performance.

The three parameters characterizing each syllabus govern the way the terms in it are drawn from three different buckets of words: new knowledge terms, old knowledge terms, and style words.

- New knowledge terms are (i) in the top 5% of the word frequency distribution among articles published between $t - 3$ and $t - 1$ or (2) words that appear in articles published between $t - 3$ and $t - 1$ but not those published between $t - 15$ and $t - 13$;
- Old knowledge terms are (i) in the top 5% of the word frequency distribution among articles published between $t - 15$ and $t - 13$ or (2) words that appear in articles published between $t - 15$ and $t - 13$ but not those published between $t - 3$ and $t - 1$;
- Style words are those terms that appear in academic articles, but do not belong to the previous two groups;
- gap is the share of old to new knowledge words in a syllabus.

To generate each syllabus, we use the following algorithm:

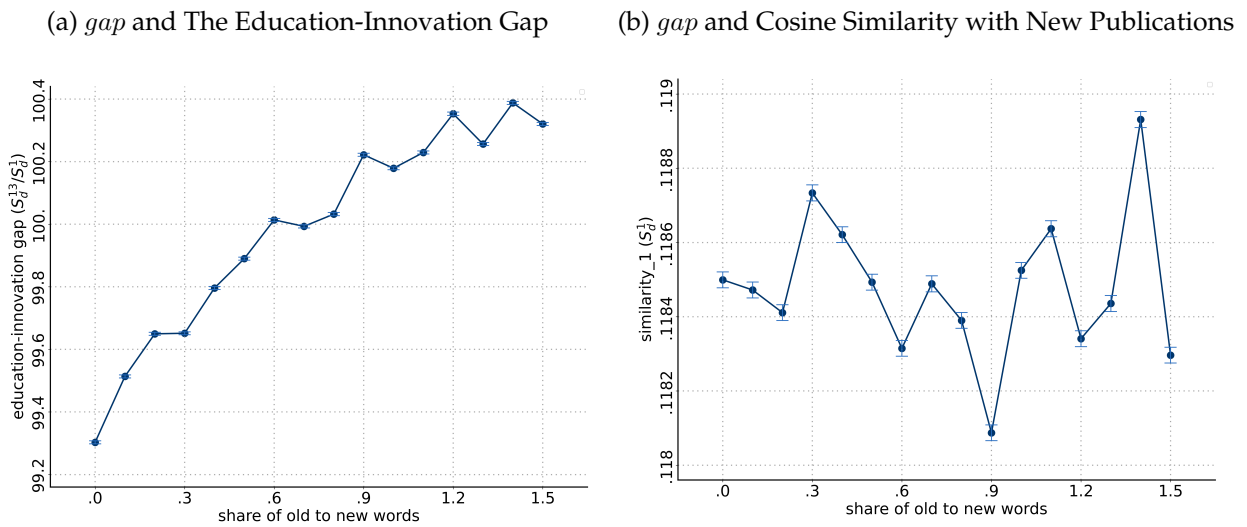
- We assign the syllabus a length of L , where $L = 10 * U$ and U is drawn from a discrete uniform distribution between 1 and 50 (so that L lies between 10 and 500, with increments of 10);
- We assign the syllabus a number $L_s = L \times style$ style words, where $style$ ranges between 0.01 and 0.1 in increments of 0.01;
- The remaining $L - L_s = L_k$ words in the syllabus are drawn from the new and old knowledge terms buckets; among these, $L_k \times (1 + gap)^{-1}$ are from the new knowledge terms bucket and $L_k \times gap \times (1 + gap)^{-1}$ are from the old knowledge terms bucket.

With this algorithm, we generate 10 syllabi for each set of parameters $\{t, L, style, gap\}$; the total number of generated syllabi is thus $= 10 \times 21 \times 46 \times 11 \times 16 = 1,700,160$, which is close to the sample size in our data.

Figure BXII (panel (a)) shows the relationship between gap and our estimated education-innovation gap. The correlation between these variables is strong and equal to 0.96. By contrast, in panel (b)

we show the relationship between *gap* and the cosine similarity between the syllabus and new publications (appeared in $t - 3$ to $t - 1$), i.e., the denominator of the education-innovation gap. The relationship is much noisier. This could happen because a simple cosine similarity is likely to be affected by the overall style of the syllabus, whereas the *gap* is not.

Figure BXII: Simulated Syllabi and Their “True” Gap Measure



Note: Panel (a) shows the relationship between *gap* and the education-innovation gap as defined and constructed in the paper. Panel (b) shows the relationship between *gap* and the cosine similarity between the syllabus and new publications (appeared in $t - 3$ to $t - 1$).