



COURSE-SKILL ATLAS: A NATIONAL LONGITUDINAL DATASET OF SKILLS TAUGHT IN U.S. HIGHER EDUCATION CURRICULA

A PREPRINT

 **Alireza Javadian Sabet**¹
alj112@pitt.edu

 **Sarah H. Bana**^{2,3}
sarah.bana@gmail.com

 **Renzhe Yu**^{4,5}
renzheyu@tc.columbia.edu

 **Morgan R. Frank**^{1,3,6*}
mrfrank@pitt.edu

¹Department of Informatics and Networked Systems, University of Pittsburgh, Pittsburgh, PA 15216, USA

²Argyros School of Business and Economics, Chapman University, Orange, CA, USA

³Digital Economy Lab, Institute for Human-Centered Artificial Intelligence,
Stanford University, Stanford, CA 94305, USA

⁴Teachers College, Columbia University, New York, NY 10027, USA

⁵Data Science Institute, Columbia University, New York, NY 10027, USA

⁶Media Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

September 17, 2024

ABSTRACT

Higher education plays a critical role in driving an innovative economy by equipping students with knowledge and skills demanded by the workforce. While researchers and practitioners have developed data systems to track detailed occupational skills, such as those established by the U.S. Department of Labor (DOL), much less effort has been made to document which of these skills are being developed in higher education at a similar granularity. Here, we fill this gap by presenting Course-Skill Atlas – a longitudinal dataset of skills inferred from over three million course syllabi taught at nearly three thousand U.S. higher education institutions. To construct Course-Skill Atlas, we apply natural language processing to quantify the alignment between course syllabi and detailed workplace activities (DWAs) used by the DOL to describe occupations. We then aggregate these alignment scores to create skill profiles for institutions and academic majors. Our dataset offers a large-scale representation of college education’s role in preparing students for the labor market. Overall, Course-Skill Atlas can enable new research on the source of skills in the context of workforce development and provide actionable insights for shaping the future of higher education to meet evolving labor demands, especially in the face of new technologies.

Keywords Skill · Higher Education · Future of Work · Labor Economics · Complexity · O*NET · Workplace Activity

Background & Summary

Skills are essential components of jobs and shape the career outcomes of workers in the labor market. Therefore, systematically studying skills and their sources is essential for predicting workers’ career trajectories and macro-level workforce dynamics [1, 2, 3, 4]. For example, recent research finds increasing demand for social skills for modern, flexible team-based work environments based on required skills in job postings [5]. In response to shifts in skills, employers need to consider skill profiles and skill development in their hiring and training. For instance, employers subjectively perceive the skill content of college majors when determining the requirements to include in online job

*corresponding author: Morgan R. Frank (mrfrank@pitt.edu)

advertisements [6]. Combined, the focus on skills in the labor market warrants a similar perspective on the sources of skills during workforce development and talent acquisition.

Higher education is arguably the most important source of skill development, which facilitates both economic and social mobility [7]. In the past few decades, empirical studies have consistently demonstrated that college-educated individuals earn higher wages, achieve more extensive professional networks, and collectively experience greater inter-generational upward mobility [8, 9]. Non-college educated workers now engage more in less skilled tasks than their counterparts compared to previous eras [10] and tend towards low-wage occupations. On the other hand, the economic returns of higher education vary across fields of study due to differing skill sets imparted by college majors [11, 12]. They also vary because of institutional selectivity [13]. Moreover, students from different demographic and socioeconomic backgrounds are sorted into different educational trajectories due to existing structural inequalities which may hinder the social mobility that higher education is intended to foster [14]. In recent years, as elevated dropout rates [15] and rising unemployment or underemployment rates of college graduates fuel concerns around the efficacy of higher education [16, 17], it is important to better understand *how* higher education imparts skills and prepares students for the labor market. This will require moving beyond using degree or credit attainment as proxies for skills, as these proxies fail to capture subtle variations in educational experiences between students in the same major or institution.

Recently, large-scale data about curricular and job *content* has become available in digital formats, which provides a new possibility of examining the mechanism of skill development. Some researchers have started to leverage these new data sources to connect higher education and jobs using natural language processing techniques [18, 19, 20, 21, 22, 23, 24, 25]. For example, Börner and colleagues [20] provided one of the earliest large-scale analysis of the alignment between college courses, job vacancies, and academic research. They leveraged an established skill taxonomy and connected the three pieces via skill mentions in their content. More recently, Light [19] measures changes in university course offerings over time by quantifying the semantic overlap between course descriptions and job postings. This novel line of research has provided emerging evidence of mismatches between labor market demands and skills taught in courses, as well as the uneven distribution of these mismatches along such dimensions as major areas, institutions, and geographical locations.

Despite the important findings, almost all the empirical studies to date acquire educational and job content information through proprietary data contracts with private vendors, and few have released their derived data about education-occupation alignment, making replication and extension efforts challenging for other researchers. In this context, we provide a new algorithmic pipeline and a public dataset to help analyze the skill development in American higher education [26]. We first introduce Syllabus20*NET, a natural language processing (NLP) framework designed to identify and interpret skills from curricular content, in line with the O*NET taxonomy [27] used by the U.S. Department of Labor (DOL) (see Figure 1). Applying Syllabus20*NET to the most extensive dataset of university course syllabi, we then present Course-Skill Atlas — a longitudinal, national dataset of inferred skill profiles across different institutions, academic majors, and student populations in the United States. To validate this dataset, we perform qualitative and quantitative explorations of the identified skills in reference to existing studies. We further discuss a handful of potential use cases of Course-Skill Atlas, including quantifying skill-salary correlations, analyzing temporal trends in curriculum design, and revealing gender skills gaps based on major and institution.

Overall, our provision makes three intellectual and practical contributions. First, we present a computational framework to describe the content alignment between education and workforce. While we developed the methodology based on two specific document types, it is applicable to other types of documents as well. Second, we provide an essential, public data source on the granular skill profiles of institutions, which can facilitate future research in such fields as higher education, labor economics, and future of work, especially in an era marked by rapid technological advancements and shifting economic landscapes [28, 29]. Third, our validation analyses illustrate some macro-level patterns of skills taught in higher education that warrant more in-depth research in the future.

Methods

Materials

Open Syllabus Project Dataset

Open Syllabus Project (OSP) (<https://opensyllabus.org/>) is a non-profit organization that curates a vast archive of over 20.9 million course syllabi from higher education institutions worldwide. The organization aims to map and analyze the curriculum across thousands of institutions, providing insights into the most commonly taught texts and subjects. OSP’s syllabi data comes from (i) scraped content from universities’ syllabi repositories, (ii) a broad web crawler with seeds from CommonCrawl (<https://commoncrawl.org/>) and manual curation, (iii) Internet Archive for 2021 Open Syllabus crawl (<https://archive.org/details/OPENSYLLABUS-20210506220126-crawl804>),

and (iv) syllabi donation from institutions and individuals. Through a research contract, we analyze one version 2.1 of the OSP data, which encompasses nearly 8 million course syllabi worldwide among which 3,162,747 syllabi across 62 fields of study (FOS) belong to 2,761 colleges and universities in the United States. Each course syllabus contains features such as course description, language, year, field of study, and information about the institution. In this paper, ‘major’ and ‘FOS’ are used interchangeably.

O*NET

O*NET (Occupational Information Network) (<https://www.onetonline.org/>) stands as a comprehensive database detailing worker attributes and job characteristics. Developed under the sponsorship of the U.S. Department of Labor/Employment and Training Administration, O*NET is essential for in-depth labor market and workforce analyses [30, 31, 32, 33, 34, 35, 36, 37, 38]. Educators, career counselors, and workforce development professionals leverage O*NET for evaluating job requirements against worker qualifications, aiding in curriculum development, career guidance, and labor market analysis [39]. Occupations form the core of the O*NET system, around which a standardized hierarchical taxonomy is organized, allowing for detailed analysis and comparison across diverse professional roles including:

- **Worker Characteristics:** These are essential in understanding the potential and capacity of the workforce. We specifically focus on:
 - *Ability* (<https://www.onetonline.org/find/descriptor/browse/1.A>): The performance of individual workers is influenced by their enduring abilities, the most granular component of worker characteristics. There are 52 abilities categorized into four key areas: cognitive, physical, psychomotor, and sensory.
- **Occupational Requirements:** These requirements define the specific demands of jobs and are integral to occupational definitions within the taxonomy. We analyze:
 - *Detailed Work Activity (DWA)* (https://www.onetcenter.org/dictionary/20.1/excel/dwa_reference.html): DWAs are precise descriptions of tasks and responsibilities of specific jobs. There are more than 2,000 DWAs across different occupations, which help understand the day-to-day activities and skills required for a particular role.
 - *Task Statement* (https://www.onetcenter.org/dictionary/20.1/excel/task_statements.html): Task is the basic unit of work. There are nearly 18,000 tasks in total, which provide the most detailed overview of job responsibilities.

Through its multidimensional taxonomy centered around occupations, O*NET not only facilitates a detailed understanding of job roles but also significantly aids in bridging educational preparation with labor market demands. This structure makes it an invaluable tool for policymakers, educators, and employment specialists. By focusing on occupations, O*NET enables a targeted analysis of the workforce, enhancing the relevance and applicability of labor market data in various professional settings.

Skill Inference Framework (Syllabus2O*NET)

Figure 1 provides an overview of our Syllabus2O*NET skill inference framework. This framework leverages natural language processing to estimate skill coverage in curricular content. Syllabus2O*NET begins by taking the raw texts of a course syllabus, which typically include course logistics (e.g., scheduling and grading rubrics) and learning content (e.g., learning objectives).

We then use Stanza [40] to partition the raw texts into individual sentences. Stanza leverages pre-trained neural network models to syntactically parse documents into sentences across diverse contexts and languages and has been shown to work even in the presence of complex punctuation and formatting [40]. This tool helps extract 322,473,524 sentences from the 3,162,747 course syllabi in the OSP dataset. On average, each syllabus contains 101.96 sentences (median 83).

Because not all sections of a syllabus reflect the subject matter and covered skills of the course, a human-in-the-loop approach is deployed to remove sentences pertaining to course logistics while keeping sentences about learning content. To do so, we compiled two distinct lists of keywords for labeling sentences, one for course logistics that includes 356 phrases (e.g., “plagiarism,” “attendance,” and “office hours”), and another for learning content including 51 phrases (e.g., “analyze,” “versus,” and “outcome”). The complete lists can be found on “course_logistics_terms” and “learning_content_terms” files on Figshare [26] and on the project’s GitHub page. We removed sentences from each syllabus that contained phrases related to course logistics or that lacked language related to learning content, resulting in

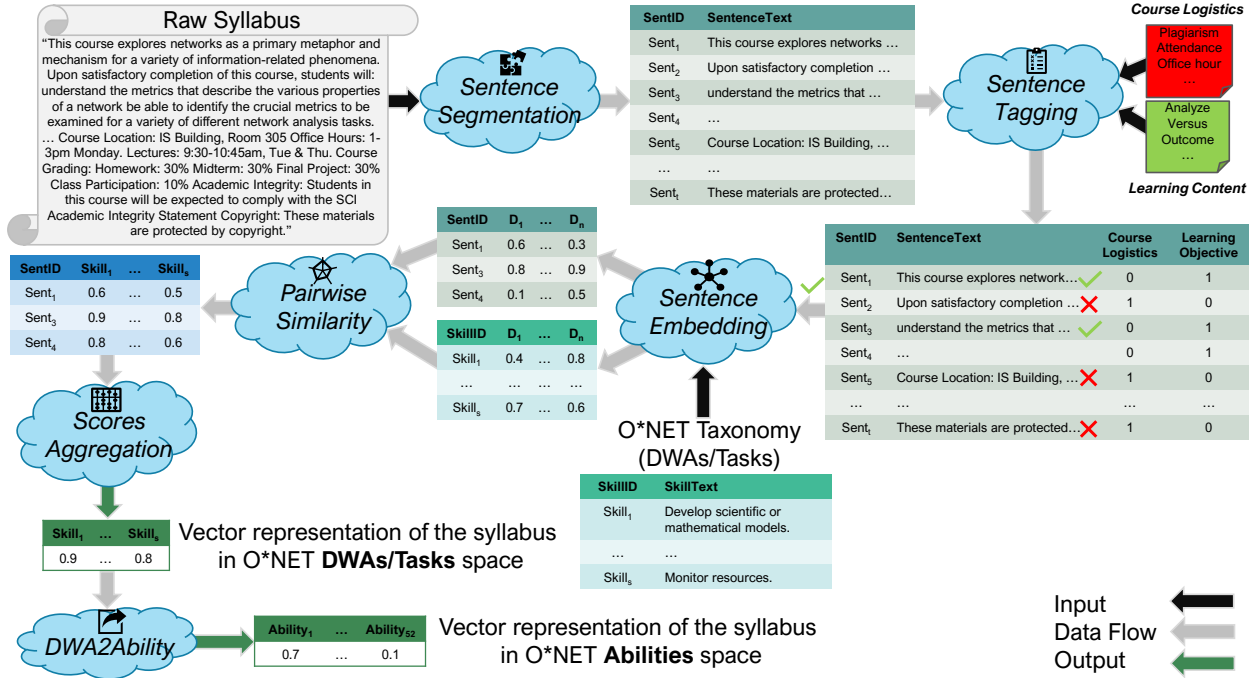


Figure 1: **The Syllabus2O*NET skill inference framework.** This natural language processing framework converts a course syllabus into a vector representing its coverage of individual “skills” defined by O*NET Detailed Work Activity (DWA) or Task. The pipeline receives a course syllabus as input and segments the raw texts into individual sentences. Then, using a curated dictionary, it identifies and keeps sentences related to learning content and transforms each sentence into a high-dimensional vector with sentence embedding (SBERT). Meanwhile, each skill is vectorized with the same approach. Next, pairwise cosine similarities between the embeddings of skills and learning content sentences are computed. Then, each skill’s maximum similarity score across the learning content sentences is used to indicate the skill coverage of the syllabus. Finally, DWA2Ability, which is composed of 52 Random Forest Regressors, maps the inferred skill coverage to related worker abilities.

the removal of 85.82% of the sentences in the raw data. After this cleaning process, each syllabus on average contains 17.61 learning content sentences (median = 12) (see Table 1 for details on the statistics of “learning content” sentence counts by FOS).

Next, we compute the semantic similarity between each O*NET DWA or Task (hereafter, “skill”) and each sentence in a syllabus. SBERT [41], a neural language model with a Siamese network structure, is used to convert a sentence into a fixed-size vector (also called “embedding”) that encodes its semantic meaning. We choose SBERT over alternative language models due to its diverse training corpora, faster computation, and superior performance on benchmark tasks. SBERT is trained on a diverse range of more than 1 billion sentences including S2ORC: The Semantic Scholar Open Research Corpus [42], WikiAnswers Corpus [43], PAQ: 65 Million Probably-Asked Questions [44], and GooAQ: Open Question Answering with Diverse Answer Types [45]. Specifically, we implement the “all-mpnet-base-v2” model [46] to embed each “learning content” sentence in course syllabi and each skill descriptor into a 768-dimension semantic space. Pairwise cosine similarity between these embeddings are calculated to measure semantic similarity between learning content and skills. For instance, “understand the metrics that describe the various properties of a network be able to identify the crucial metrics to be examined for a variety of different network analysis tasks” (in Figure 1) is semantically similar to the O*NET DWA “develop scientific or mathematical models,” with a cosine similarity of 0.9.

Finally, we create a vector for each syllabus to measure how much each skill is covered in the course, based on the semantic similarities. Specifically, for a given skill, we select the maximum similarity score across all the sentences within the syllabus (i.e., the score of the most similar sentence). This approach captures the most relevant skill information each course aims to develop and is particularly robust for handling the significant variations in syllabi’s length, detail, and structures. With this construction, the syllabus vectors have 2,070 dimensions for DWAs and 17,992 for Tasks.

Field of Study	#Sent.	# L. Sent.	% L. Sent.	Field of Study	# Sent.	# L. Sent.	% L. Sent.
Accounting	101	15	14.61%	Hebrew	65	7	11.32%
Agriculture	60	8	15.11%	History	88	9	10.62%
Anthropology	86	13	14.29%	Japanese	84	10	13.12%
Architecture	73	14	20.45%	Journalism	105	13	12.22%
Astronomy	76	10	12.50%	Law	60	7	12.12%
Atmospheric Sciences	66	11	14.29%	Library Science	107	13	12.04%
Basic Computer Skills	88	11	13.21%	Linguistics	72	9	11.43%
Basic Skills	76	11	14.29%	Marketing	106	16	15.79%
Biology	85	12	13.25%	Mathematics	75	10	12.73%
Business	97	15	15.22%	Mechanic/Repair Tech	62	10	17.95%
Chemistry	83	11	13.16%	Media/Communications	107	14	13.73%
Chinese	60	8	12.10%	Medicine	95	14	15.79%
Classics	47	4	9.30%	Military Science	71	11	14.35%
Computer Science	67	10	14.29%	Music	62	8	13.04%
Cosmetology	85	12	15.79%	Nursing	95	16	16.95%
Criminal Justice	79	12	15.38%	Nutrition	85	12	14.46%
Culinary Arts	61	14	22.22%	Philosophy	76	9	12.50%
Dance	59	11	17.54%	Physics	61	8	14.19%
Dentistry	69	16	25.93%	Political Science	96	12	12.23%
Earth Sciences	59	9	14.29%	Psychology	105	16	15.09%
Economics	74	11	13.92%	Public Safety	60	10	16.03%
Education	121	23	19.74%	Religion	84	11	12.61%
Engineering	51	9	16.28%	Sign Language	78	12	15.69%
English Literature	97	11	10.99%	Social Work	139	26	19.28%
Film and Photography	68	11	14.88%	Sociology	100	14	13.92%
Fine Arts	79	13	15.75%	Spanish	95	10	10.60%
Fitness and Leisure	54	10	17.19%	Theatre Arts	63	11	17.68%
French	69	7	10.45%	Theology	113	16	14.29%
Geography	70	10	14.29%	Transportation	94	12	13.46%
German	71	8	11.39%	Veterinary Medicine	70	10	11.39%
Health Technician	75	11	16.22%	Women’s Studies	81	13	15.19%

Table 1: **Sentence statistics per FOS.** The table presents the median of the number of sentences (# Sent.), number of identified Learning Objective related sentences (# L. Sent.), and the percentage of the identified Learning Objective related sentences (% L. Sent.) per FOS.

Some studies might aim at depicting the competencies of workers and do not have detailed information about jobs. In line with this use case, we further establish the connection between learning content and O*NET abilities by mapping inferred skills to abilities. Because O*NET does not provide a standardized crosswalk linking DWAs, tasks, and abilities, we create a subprocess `DWA2Ability` to achieve this. We start with the O*NET database profiles of DWAs for each occupation. Next, we extract importance scores of each O*NET ability within each occupation. We formulate a map between the two sets of occupation profiles as a regression using DWAs as independent variables and ability scores as dependent variables. We train a Random Forest Regressor [47] for each ability and fine-tune hyperparameters via Grid Search and 5-fold cross-validation. This approach yielded 52 models (i.e., one per O*NET ability), each achieving mean squared error of at most 0.025 (see Table 2 for details on model performance). Using the trained models, we map syllabi’s DWA scores to abilities. If a syllabus does not teach content that provides the students with a certain ability, the corresponding ability score is 0.

Negative values in cosine similarity To determine how much each skill is covered in a course, we calculate the cosine similarity between the skill and each sentence in a syllabus. Cosine similarity ranges from -1 to 1 , where a value of 1 indicates that the two vectors (or sentences in this context) are perfectly aligned and have the same direction, while a value of -1 indicates that the vectors are diametrically opposed. In general, a negative cosine similarity value suggests that the two sentences are not only dissimilar but also convey opposite or contrasting meanings. However, in the context of skill inference, the meaning of dissimilarity might not always be applicable. For example, the top five DWAs with the most negative values include “Trim trees or other vegetation,” “Adjust the tension of nuts or bolts,” “Install carpet or flooring,” “Install trim or paneling,” and “Apply sealants or other protective coatings.” In the output of our `Syllabus2O*NET` pipeline, fewer than 10 DWAs have negative values in more than 5% of the syllabi, with the highest being 7.17%. For transparency and flexibility, we have kept these negative values as they are in the released aggregated datasets, allowing users to take appropriate actions according to their specific needs (e.g., change them to zero or normalize them).

Ability	MSE	Ability	MSE	Ability	MSE	Ability	MSE
Arm-Hand Steadiness	0.025	Fluency of Ideas	0.014	Number Facility	0.012	Speech Clarity	0.008
Auditory Attention	0.014	Glare Sensitivity	0.017	Oral Comprehension	0.010	Speech Recognition	0.014
Category Flexibility	0.016	Gross Body Coordination	0.013	Oral Expression	0.011	Speed of Closure	0.010
Control Precision	0.019	Gross Body Equilibrium	0.018	Originality	0.014	Speed of Limb Movement	0.020
Deductive Reasoning	0.011	Hearing Sensitivity	0.012	Perceptual Speed	0.012	Stamina	0.019
Depth Perception	0.014	Inductive Reasoning	0.010	Peripheral Vision	0.013	Static Strength	0.024
Dynamic Flexibility	0.007	Information Ordering	0.017	Problem Sensitivity	0.010	Time Sharing	0.013
Dynamic Strength	0.016	Manual Dexterity	0.023	Rate Control	0.019	Trunk Strength	0.017
Explosive Strength	0.016	Mathematical Reasoning	0.011	Reaction Time	0.020	Visual Color Discrimination	0.019
Extent Flexibility	0.019	Memorization	0.013	Response Orientation	0.012	Visualization	0.017
Far Vision	0.011	Multilimb Coordination	0.022	Selective Attention	0.008	Wrist-Finger Speed	0.024
Finger Dexterity	0.015	Near Vision	0.018	Sound Localization	0.025	Written Comprehension	0.012
Flexibility of Closure	0.014	Night Vision	0.015	Spatial Orientation	0.017	Written Expression	0.014

Table 2: DWA2Ability models training performance for each ability. The mean squared error (MSE) obtained from the 5-fold cross-validation (CV) for finding the best model of ability.

Skill Normalization

Some skills, such as “Maintain student records” and “Document lesson plans”, are ubiquitous across fields of study (FOS) and, therefore, do not distinguish the learning content of one FOS from another. To address this issue and control for widespread skills we propose two approaches when using our data. Although we use them in some of the validation analysis, the published dataset remains intact.

The first approach is applying Revealed Comparative Advantage (RCA) (a.k.a., “location quotient” [48, 33, 49, 50, 51]). RCA is a concept used to identify areas where an entity has a relatively high presence compared to others. In our context, RCA helps to reveal which skill s most strongly distinguishes one FOS m from others. RCA is calculated as follows:

$$rca(m, s) = \frac{onet(m, s) / \sum_{s' \in S} onet(m, s')}{\sum_{m' \in S} onet(m', s) / \sum_{m' \in M, s' \in S} onet(m', s')}, \quad (1)$$

where $m \in M$ denotes a FOS (i.e., a college major) and $s \in S$ denotes a skill (i.e., an O*NET DWA). If $rca(m, s) > 1$, then s is more related to m than would be expected across all DWAs and all FOS; therefore, s is a relatively distinctive skill identifying m more strongly than other FOS.

The second approach, which we used to produce Fig. 11, is to mask frequent skills in empirical analyses that use the skill scores. What is considered “frequent” should depend on the specific application, so we do not mask any skills in the published dataset. However, as a useful resource, we create a table “top10_DWA_per_FOS” on Figshare [26] which lists the top 10 inferred workplace activities with the highest average DWA scores per FOS.

Data Records

The Course-Skill Atlas dataset is freely available at Figshare [26].

Schema and Variables

Our research contract with OSP requires that we do not release information at the individual syllabus level. As such, we create a dataset of inferred skills aggregated at the institution-year-FOS level. The Course-Skill Atlas dataset includes three key components: **DWAs**, **Tasks**, and **Abilities**. Each represents the corresponding vectors produced by Syllabus2O*NET and aggregated at our chosen granularity. Figure 2 provides the entity relationship diagram of these components.

Each record in these datasets includes the year, institution name, and UnitID (i.e., the unique identifier assigned by Integrated Postsecondary Education Data System (IPEDS) (<https://nces.ed.gov/collegenavigator/>) to each institution), the geographical location of the institution (i.e., the city and state), the FOS name along with its CIP code(s), as well as the institution’s sector — one of nine institutional categories created by combining an institution’s control and level (e.g., “public 4-year or above”), which we obtained from the Carnegie Classification of Institutions of Higher Education (CCIHE) (<https://carnegieclassifications.acenet.edu/>).

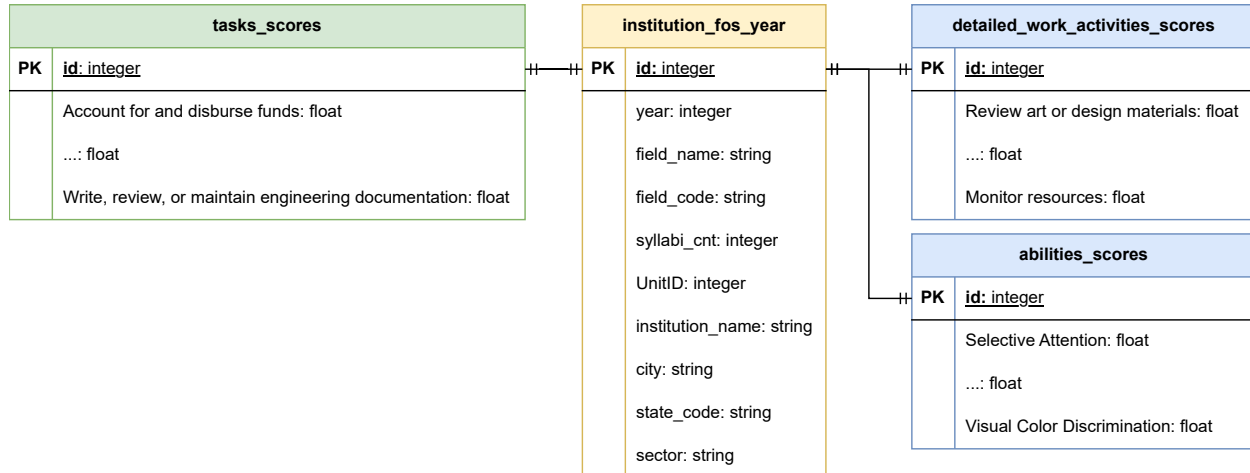


Figure 2: **The entity relationship diagram of the skills extracted from U.S. course syllabi.** In each table, PK represents the table’s primary key. “institution_fos_year” comprise the main data table encompassing 281, 153 records. For each corresponding id from “institution_fos_year” Table, “task_scores”, “detailed_work_activities_scores”, and “abilities_scores” tables contain the scores for 17, 992 tasks, 2, 070 DWAs, and 52 abilities respectively inferred using Syllabus20*NET). For brevity, we replaced the remaining tasks, DWAs, and abilities with “...”. Lines connecting tables indicate the presence of a relational table.

The field_code string contains one or more IPEDS CIP codes (<https://nces.ed.gov/ipeds/>), representing the field(s) of study most associated with the syllabus. OSP’s field classifier relies on the IPEDS 2010 CIP taxonomy (<https://nces.ed.gov/ipeds/cipcode/default.aspx?y=55>) to determine the most relevant field of study (FOS) for each syllabus. CIP codes are structured in lengths of two, four, and six digits, where two-digit codes represent a broad discipline, four-digit codes represent subdivisions of that discipline, and six-digit codes provide further subdivisions. For instance, the two-digit CIP code ‘01’ corresponds to “Agriculture, Agriculture Operations, and Related Sciences”; within this category, the four-digit code “01.01” denotes “Agricultural Business and Management,” and “01.0103” specifies “Agricultural Economics.” OSP’s FOS classifier is trained and tested on a curated subset of the CIP taxonomy that OSP has found most effective for describing syllabi. Occasionally, OSP combines codes, but only within the same two-digit branch of the taxonomy. In these instances, codes are separated by a forward slash (/). For example, the code “45.09/45.10” merges “International Relations and National Security Studies” and “Political Science and Government,” both of which fall under the two-digit code ‘45’ for “Social Sciences.” This combined FOS is labeled as “Political Science.” Note that the assigned CIP code(s) are consistent across all the syllabi; meaning that all “Political Science” syllabi are mapped to the same list of CIP codes, i.e., “45.09/45.10” (see Table field_name_and_code on Figshare [26] for the list of fields name and their CIP code(s) mappings). If OSP is unable to confidently assign an academic field to the syllabus, the value of this column is null. Moreover, we enriched each record by the sector of the institutions (i.e., control and level combined (<https://carnegieclassifications.acenet.edu/wp-content/uploads/2023/03/CCIHE2021-PublicData.xlsx>)). Further university characteristics can be added to the dataset by merging the syllabi Table and the target dataset on the UnitID variable.

The aggregated scores are the **average scores** across all the syllabi belonging to the corresponding year, university, and FOS. For example, let’s consider a given year, university, and FOS which has two syllabi. The score of DWA_1 in $Syllabus_a$ is 0.8 and the score of the same DWA in $Syllabus_b$ is 0.6. Taking the average of the two scores, the aggregated score of DWA_1 for the given triplet is 0.7.

For the remainder of this paper, we only use DWA scores for various descriptive and validation analyses, which can be easily applied to task and ability scores as well.

Descriptive Statistics

Figure 3 depicts the geographical, temporal, and institutional distribution of the syllabi data. Across each U.S. state, between 32% and 76% of the postsecondary institutions in each state provide at least eight course syllabi (i.e., corresponding to the 25th percentile. See Fig. 3a). The majority of the syllabi belong to the period post-2000, with sparse coverage between 1966 and 1999. (See Fig. 3b) Across the entire data set, the majority of universities (nearly 2000) contribute at least 10 syllabi to the data, but some contribute up to 10^5 syllabi across all FOS (see Fig. 3c).

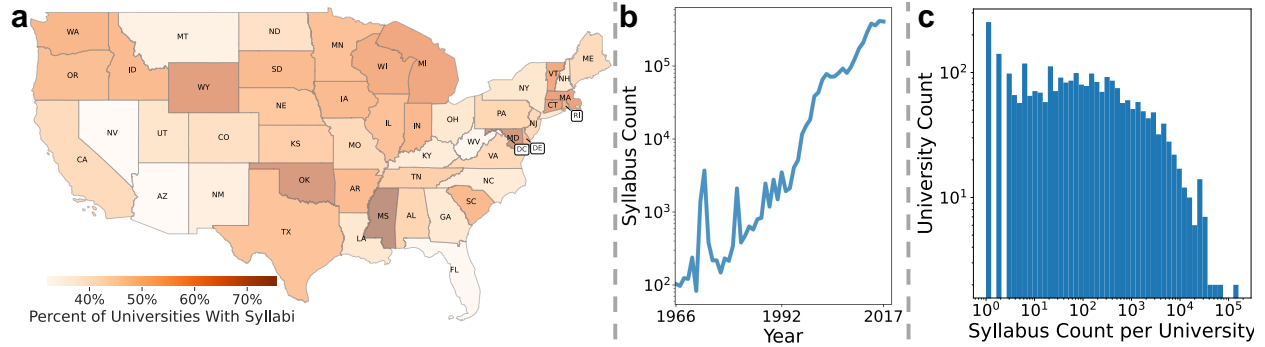


Figure 3: **Descriptive statistics of the Open Syllabus Project (OSP) dataset.** (a) Percentage of the universities with at least eight course syllabi (25th percentile) available per state. (b) Total number of syllabi per year. (c) The syllabus count per university across all years and all FOS.

Figure 4 shows the syllabi count distribution of the aggregated records. Table 3 lists the frequency of syllabi per FOS. Table 4 details the geographical coverage of the OSP dataset. Number of educational institutions within each state is obtained from CCIHE. For example, Texas with 865,973 syllabi has the largest number of syllabi (27.85%) in the dataset. According to CCIHE, there are 226 universities and educational institutions located in Texas, among which 54.42% have at least 8 syllabi (25th percentile) in the OSP dataset. Moreover, more than 80% of the syllabi belong to the public universities with a majority belonging to the 4-year universities. Nearly 15% are from private not-for-profit, 4-year or above universities (see Table 5 for the frequency and percentage of syllabi per university sector). In addition, Table 6 lists the syllabus count per top 100 universities across all years and all FOS. Finally, Figure 5 details the number of syllabi per FOS between 2000 and 2017 (see Figure 6 for FOS according to 2-digit Classification of Instructional Programs (CIP) 2010 taxonomy).

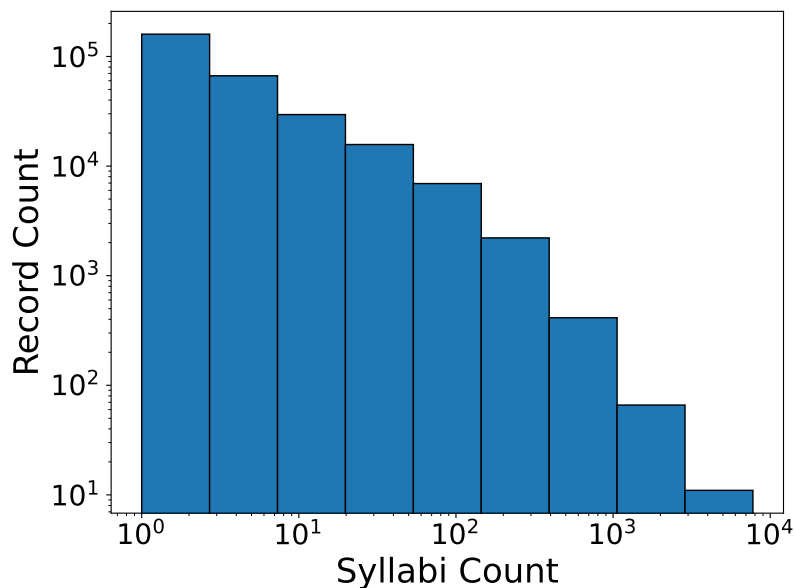


Figure 4: **Syllabi count distribution.** The syllabi count distribution of the records (a record is the average score of all syllabi belonging to a FOS, year, and institution triplet).

Field of Study	# Syllabi	Field of Study	# Syllabi	Field of Study	# Syllabi
Mathematics	258,160	Accounting	51,984	Religion	14,440
English Literature	232,065	Sociology	46,836	French	14,305
Business	201,100	Physics	44,802	Journalism	12,712
Computer Science	184,649	Film and Photography	42,690	Nutrition	11,883
Biology	140,187	Criminal Justice	39,805	Dentistry	10,367
Education	140,182	Spanish	39,650	Culinary Arts	9,430
Fitness and Leisure	131,262	Health Technician	38,268	Sign Language	8,665
Psychology	122,387	Social Work	36,745	German	8,385
History	107,676	Philosophy	35,583	Classics	7,813
Media / Communications	85,561	Agriculture	35,305	Cosmetology	7,291
Music	82,329	Marketing	31,430	Astronomy	7,286
Fine Arts	75,722	Law	31,421	Transportation	7,121
Basic Skills	73,362	Theatre Arts	29,087	Japanese	5,456
Engineering	70,084	Theology	24,584	Women's Studies	5,237
Political Science	69,111	Public Safety	23,931	Chinese	5,054
Basic Computer Skills	68,028	Earth Sciences	21,870	Linguistics	4,859
Nursing	63,603	Anthropology	21,509	Military Science	3,202
Mechanic / Repair Tech	62,423	Library Science	20,234	Atmospheric Sciences	2,231
Chemistry	61,280	Dance	19,694	Veterinary Medicine	2,105
Economics	56,157	Architecture	19,379	Hebrew	1,674
Medicine	55,161	Geography	17,935		

Table 3: Frequency of syllabi per FOS.

Name	# Syllabi	% Syllabi	# Inst.	% Covered Inst.	Name	# Syllabi	% Syllabi	# Inst.	% Covered Inst.
Alabama	66,548	2.14%	60	48.33%	Montana	6,929	0.22%	24	37.50%
Alaska	3,271	0.11%	8	50.00%	Nebraska	4,609	0.15%	34	52.94%
Arizona	16,476	0.53%	66	31.82%	Nevada	16,541	0.53%	22	31.82%
Arkansas	9,018	0.29%	53	56.60%	New Hampshire	3,066	0.10%	24	41.67%
California	553,589	17.80%	425	46.82%	New Jersey	38,225	1.23%	82	47.56%
Colorado	18,718	0.60%	62	45.16%	New Mexico	28,830	0.93%	36	38.89%
Connecticut	8,695	0.28%	38	65.79%	New York	77,699	2.50%	295	43.05%
Delaware	1,737	0.06%	7	57.14%	North Carolina	76,644	2.46%	134	41.04%
District of Columbia	13,356	0.43%	22	36.36%	North Dakota	4,120	0.13%	20	45.00%
Florida	78,441	2.52%	161	32.92%	Ohio	80,494	2.59%	160	42.50%
Georgia	72,955	2.35%	107	42.06%	Oklahoma	35,158	1.13%	46	69.57%
Hawaii	19,016	0.61%	17	52.94%	Oregon	24,273	0.78%	50	56.00%
Idaho	40,312	1.30%	14	57.14%	Pennsylvania	72,515	2.33%	193	48.19%
Illinois	65,455	2.10%	152	55.26%	Puerto Rico	40	0.00%	86	2.33%
Indiana	22,389	0.72%	66	57.58%	Rhode Island	5,322	0.17%	15	60.00%
Iowa	20,861	0.67%	56	57.14%	South Carolina	37,259	1.20%	66	57.58%
Kansas	14,598	0.47%	63	50.79%	South Dakota	2,880	0.09%	21	57.14%
Kentucky	49,872	1.60%	59	40.68%	Tennessee	23,144	0.74%	83	50.60%
Louisiana	24,547	0.79%	52	42.31%	Texas	865,973	27.85%	226	54.42%
Maine	8,252	0.27%	30	46.67%	Utah	18,612	0.60%	23	43.48%
Maryland	157,830	5.08%	54	70.37%	Vermont	5,722	0.18%	16	62.50%
Massachusetts	36,102	1.16%	106	65.09%	Virginia	41,186	1.32%	108	47.22%
Michigan	120,455	3.87%	87	63.22%	Washington	45,493	1.46%	72	58.33%
Minnesota	68,253	2.19%	85	55.29%	West Virginia	5,836	0.19%	41	34.15%
Mississippi	7,174	0.23%	33	75.76%	Wisconsin	23,880	0.77%	67	61.19%
Missouri	47,594	1.53%	93	47.31%	Wyoming	19,570	0.63%	9	66.67%

Table 4: **Geographical distribution of the OSP dataset.** For each U.S. state “# Syllabi” and “% Syllabi” detail the number and percentage of course syllabi from the universities and institutions located in that state (the sum of “% Syllabi” equals to 100%). “# Inst.” specifies the number of institutions located in the given state based on the CCIHE. “% Covered Inst.” specifies the percentage of the number of universities with at least 8 course syllabi (25th percentile) in the OSP dataset.

Technical Validation

Duplicate Analysis

How many duplicate syllabi exist in our dataset? The syllabi data may have “duplicates” because an instructor might teach a course across multiple years with minimal syllabus change, or some introductory courses may have some standard design adopted across institutions. To assess the prevalence of duplicate syllabi in our dataset, we conducted

Institution Name	Count	Institution Name	Count
Alamo Colleges	160,041	South Texas College	12,009
University of Maryland University College	137,257	Dallas County Community College District	11,641
Amarillo College	75,198	Monterey Peninsula College	11,258
Lansing Community College	63,945	University of Minnesota System	11,220
University of Alabama, Tuscaloosa	54,278	South Plains College	11,184
Texas State Technical College	48,291	Wilkes University	11,036
Clark State Community College	46,094	Bellevue College	10,935
Houston Community College System	45,401	Reedley College	10,785
Santa Rosa Junior College	35,621	Modesto Junior College	10,780
Texas A&M University	33,579	University of Texas Rio Grande Valley	10,653
Rowan-Cabarrus Community College	33,403	University of Southern California	10,480
North Idaho College	31,292	Kentucky Community and Technical College System	10,264
University of Texas at Dallas	30,922	Santa Barbara City College	10,232
University of Georgia	30,349	Fullerton College	9,993
Texas State University–San Marcos	30,051	Pennsylvania State University	9,303
University of Texas at Arlington	28,651	Lewis and Clark Community College	9,103
San Diego Community College District	27,768	Southwestern Community College	9,004
Western Kentucky University	27,684	Chaminade University of Honolulu	8,866
Park University	27,096	Great Basin College	8,837
Sam Houston State University	26,523	University of Akron	8,787
Stephen F. Austin State University	26,170	University of California, San Diego	8,746
University of Michigan–Ann Arbor	25,883	University of Washington	8,707
Midwestern State University	25,460	Nova Southeastern University	8,555
Fresno City College	24,786	University of Colorado Boulder	8,553
George Mason University	24,046	Rutgers University	8,494
Oral Roberts University	23,950	University of South Florida	8,193
San Jose State University	23,505	San Mateo County Community College District	8,124
Minnesota State Colleges and Universities System	23,112	Palomar College	8,112
Texas Tech University	22,812	Westmont College	8,068
McLennan Community College	22,561	Mt. San Jacinto College	8,015
University of California, Irvine	22,505	Stony Brook University	7,982
Galveston College	22,077	New York University	7,959
Tyler Junior College	20,986	Victoria College	7,919
Clemson University	20,453	Iowa State University	7,885
University of Texas at Austin	20,387	University of Maryland, College Park	7,842
Collin College	17,608	Butte College	7,728
New Mexico Junior College	16,635	Merced College	7,388
Hartnell College	16,522	Ventura County Community College District	7,368
University of Texas at El Paso	15,842	Alvin Community College	7,347
Loyola University New Orleans	15,830	Ohlone College	7,240
University of North Texas	15,680	Imperial Valley College	7,174
University of Florida	14,919	Chaffey College	7,076
Casper College	14,665	Santa Monica College	6,870
University of Texas at San Antonio	14,595	California State University, San Marcos	6,865
University of Houston–Clear Lake	14,386	Palm Beach State College	6,742
Angelo State University	13,663	Carnegie Mellon University	6,702
Excelsior College	13,517	El Paso Community College	6,624
Texas A&M University–Commerce	13,264	Massachusetts Institute of Technology	6,329
Princeton University	13,182	Florida International University	6,287
Santa Ana College	12,969	Napa Valley College	6,224

Table 6: **Most frequent universities.** The syllabus count per top 100 universities across all years and all FOS.

an analysis on the original, disaggregated data. Specifically, we compared the DWA skill vectors from syllabi within the same field of study and university across and within various academic years. In this context, “duplicate” syllabi are defined as those with either identical textual content or learning materials that yield the same similarity score on our NLP framework’s assessment. 25.20% of the total syllabi within a university-major-year are duplicates. This number grows to 31.60% when measuring the total duplicates within a university-major pair. By taking the differences between these two, we observe that 6.40% of duplicates are across years. This relatively low value suggests that instructors are updating their syllabi over time, which leads to differences in skills critical for our analyses. These results indicate that majority of the duplicates come from multiple courses in the same major, university and academic year teaching the same content. Table 7 reports these values by field of study. Additionally, Table duplicate_counts on Figshare [26], details the total counts of both original and duplicate syllabi for each university and field of study pairing.

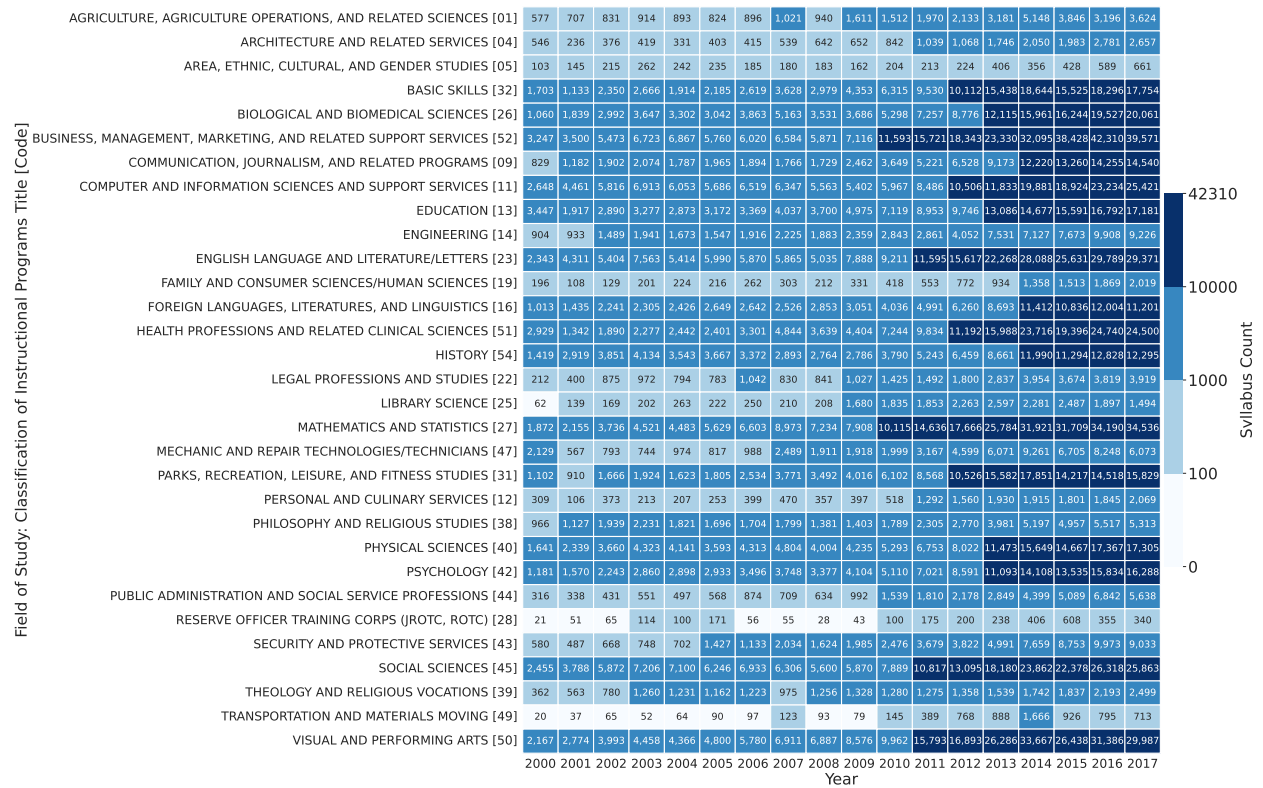


Figure 6: Frequency of syllabi per FOS according to 2-digit CIP 2010 taxonomy between 2000 and 2017.

Representativeness Analysis

How representative is our data of US higher education? A recent study using the OSP dataset [52] shows that the syllabi sample consistently represents about 5% of all courses taught in US institutions from 1998 to 2018. In their analysis of courses from 161 representative US institutions, they find that the sample slightly overrepresents Ivy-Plus schools but is broadly representative in terms of field and course level distribution. Despite this bias, the dataset adequately reflects U.S. higher education offerings during the period, with no significant bias in financial resources or student demographics.

Moreover, we calculate the coverage of institutions across different FOS according to the number of Bachelor’s degrees awarded in a FOS across US institutions between 2000 and 2017 obtained from IPEDS completions data (<https://nces.ed.gov/ipeds/datacenter/DataFiles.aspx?year=-1&surveyNumber=3&sid=6758f146-5ae5-44a0-8982-9821e70a8757&rtid=7>). We calculate the coverage as the percentage of the number of institutions with at least one syllabus from a institution-FOS-year combination divided by the total number of institutions with a positive number of bachelor’s degrees from the same institution-FOS-year combination obtained from IPEDS. Fields such as “Engineering” and “Social Sciences” show a consistent coverage over 40%. On the other hand, several fields like “Military Technologies” and “Natural Resources and Conservation” have consistently zero coverage (see Fig. 7). Table `ipeds_2digit_grad_2000_2017` on Figshare [26] provides the number of graduates per university and 2-digit CIP code between 2000 and 2017 downloaded from IPEDS.

Sufficient Number of Syllabi

How stable are our estimates of the skill content associated with each FOS and institution in each year? We analyze DWA institution-FOS-year combinations with at least two syllabi post-2000 (i.e., 2, 686, 066 syllabi). For each institution-FOS-year combination, we calculate the Manhattan and Euclidean distances between the aggregated published skill profile (i.e., the average scores of the complete set of syllabi) and the skill profile of a randomly selected subsample of syllabi (from one syllabus up to the maximum number of syllabi for the given institution-FOS-year combination minus one). We perform ten trials for each subset size. Finally, we average the distances of each subset size for all syllabi and

Field of Study	Duplicates per University, Major Combination (Total)		Duplicates per University, Major, Year Combination (Within Year)		Duplicates Across Years (Total - Within Year)	
	Number	Percentage	Number	Percentage	Number	Percentage
Library Science	12,626	64.82%	11,985	61.53%	641	3.29%
Transportation	3,544	56.11%	3,050	48.29%	494	7.82%
Public Safety	10,987	48.91%	9,525	42.40%	1,462	6.51%
Mechanic / Repair Tech	25,092	46.58%	20,504	38.06%	4,588	8.52%
Culinary Arts	3,665	41.95%	2,929	33.53%	736	8.42%
Health Technician	13,711	41.41%	11,019	33.28%	2,692	8.13%
Cosmetology	2,870	40.80%	2,243	31.89%	627	8.91%
Basic Computer Skills	24,240	40.65%	20,633	34.61%	3,607	6.04%
Basic Skills	25,887	39.81%	22,071	33.94%	3,816	5.87%
Japanese	1,971	39.55%	1,601	32.12%	370	7.43%
Dentistry	3,440	39.40%	2,897	33.18%	543	6.22%
Criminal Justice	14,208	38.03%	11,534	30.87%	2,674	7.16%
Fitness and Leisure	45,192	37.65%	34,921	29.09%	10,271	8.56%
Music	26,393	37.13%	19,150	26.94%	7,243	10.19%
Sign Language	2,454	36.74%	1,920	28.75%	534	7.99%
German	2,899	36.56%	2,283	28.79%	616	7.77%
Atmospheric Sciences	760	35.70%	566	26.59%	194	9.11%
Military Science	1,070	35.64%	842	28.05%	228	7.59%
Mathematics	79,952	35.40%	64,279	28.46%	15,673	6.94%
Spanish	12,352	34.92%	10,086	28.51%	2,266	6.41%
Nursing	19,529	34.52%	16,520	29.20%	3,009	5.32%
Dance	5,824	34.48%	4,313	25.54%	1,511	8.94%
English Literature	71,645	34.41%	63,061	30.28%	8,584	4.13%
Nutrition	3,590	33.39%	2,868	26.68%	722	6.71%
Accounting	15,836	32.59%	12,447	25.61%	3,389	6.98%
Business	61,178	32.22%	49,171	25.90%	12,007	6.32%
Computer Science	53,616	32.20%	42,614	25.59%	11,002	6.61%
Media / Communications	24,793	31.92%	20,681	26.63%	4,112	5.29%
Agriculture	9,459	31.74%	6,387	21.43%	3,072	10.31%
Astronomy	1,841	31.49%	1,374	23.50%	467	7.99%
Biology	39,526	31.36%	30,387	24.11%	9,139	7.25%
Fine Arts	21,049	30.81%	16,321	23.89%	4,728	6.92%
Medicine	15,591	30.63%	12,420	24.40%	3,171	6.23%
French	3,931	29.97%	2,962	22.58%	969	7.39%
Earth Sciences	5,946	29.74%	4,226	21.14%	1,720	8.60%
Chemistry	16,349	29.72%	11,917	21.67%	4,432	8.05%
Physics	11,255	29.14%	8,148	21.10%	3,107	8.04%
Theatre Arts	7,853	28.94%	6,142	22.63%	1,711	6.31%
Chinese	1,337	28.83%	1,003	21.63%	334	7.20%
Film and Photography	11,182	28.48%	8,784	22.37%	2,398	6.11%
History	26,505	27.77%	21,129	22.14%	5,376	5.63%
Law	7,545	27.41%	5,542	20.14%	2,003	7.27%
Veterinary Medicine	370	26.93%	270	19.65%	100	7.28%
Engineering	17,067	26.91%	12,318	19.42%	4,749	7.49%
Marketing	7,965	26.64%	6,150	20.57%	1,815	6.07%
Sociology	11,411	26.57%	9,199	21.42%	2,212	5.15%
Religion	3,509	25.03%	2,792	19.91%	717	5.12%
Psychology	28,009	24.83%	21,883	19.40%	6,126	5.43%
Classics	1,717	24.72%	1,110	15.98%	607	8.74%
Economics	12,356	24.14%	8,904	17.40%	3,452	6.74%
Social Work	8,340	23.01%	6,806	18.78%	1,534	4.23%
Geography	3,681	22.64%	2,550	15.68%	1,131	6.96%
Philosophy	6,881	22.04%	5,391	17.27%	1,490	4.77%
Political Science	13,088	21.37%	10,097	16.49%	2,991	4.88%
Hebrew	334	20.99%	185	11.63%	149	9.36%
Journalism	2,552	20.51%	1,871	15.04%	681	5.47%
Anthropology	3,993	19.94%	3,175	15.86%	818	4.08%
Education	26,628	19.70%	21,342	15.79%	5,286	3.91%
Architecture	3,115	17.76%	2,303	13.13%	812	4.63%
Women's Studies	872	17.38%	726	14.47%	146	2.91%
Theology	3,511	14.75%	2,098	8.82%	1,413	5.93%
Linguistics	597	12.99%	374	8.14%	223	4.85%

Table 7: **Duplicate syllabi per field of study.** A duplicate refers to a set of DWA vector with the same scores. The first column shows the total number and percentage of the duplicates obtained from each institution-FOS combinations. The second column shows the number and percentage of the duplicates within a year calculated based on institution-FOS-year combinations. And the last column is the total duplicates minus the the within year duplicates. The duplicate column shows the average of duplicate syllabi (i.e., having the same DWA skill vector) within the same field of study and university across different years.

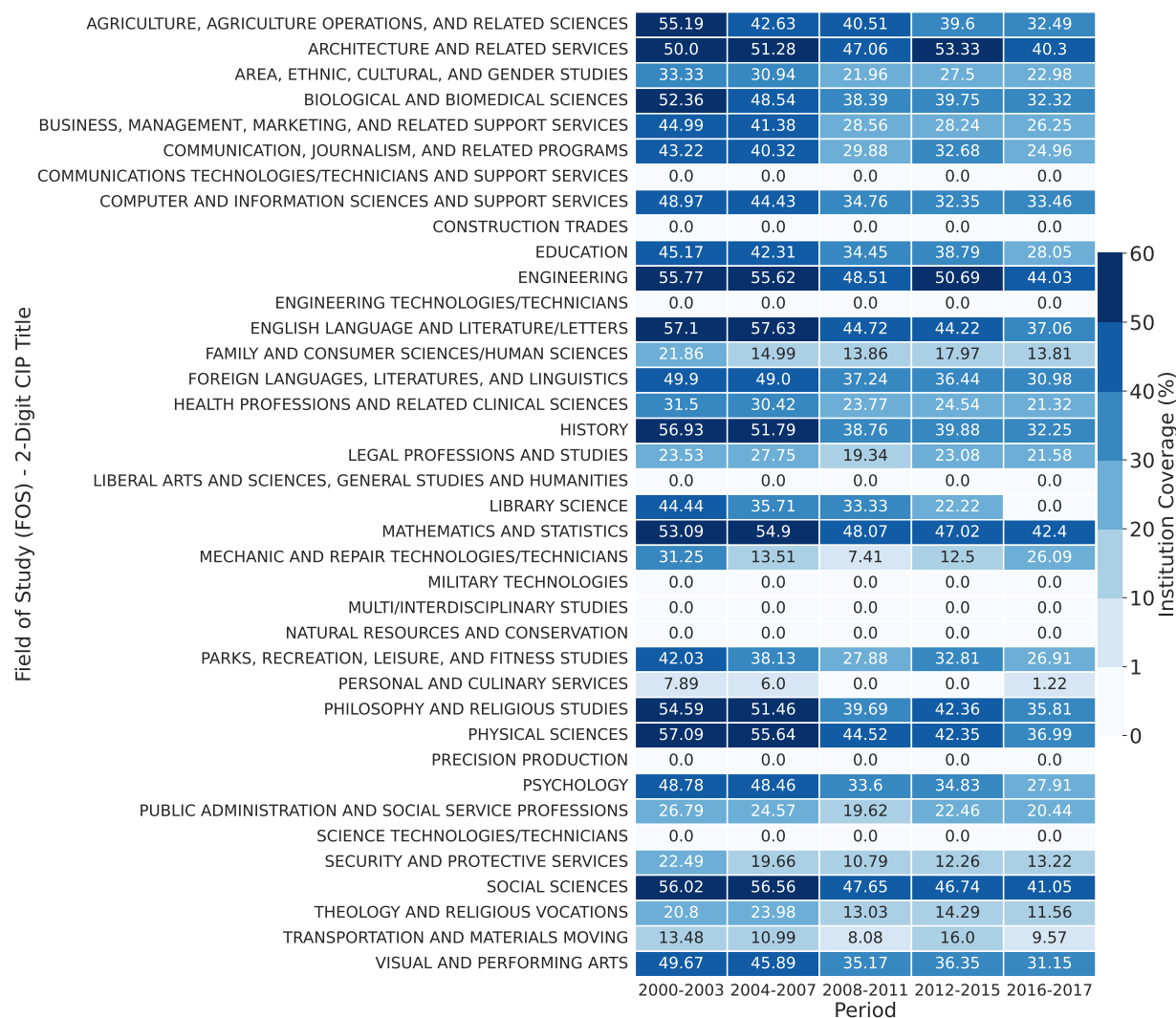


Figure 7: **Trends in the coverage of institutions across different fields of study.** The heatmap visualizes the coverage trends for various fields of study (2-digit CIP titles) across six time periods (2000-2003, 2004-2007, 2008-2011, 2012-2015, and 2016-2017). The coverage is the percentage of the number of institutions with at least one syllabus from a institution-FOS-year combination divided by the total number of institutions with a positive number of bachelor's degrees from the same institution-FOS-year combination obtained from IPEDS.

calculate their distances. Figure 8 illustrates this analysis using Manhattan (Fig. 8a) and Euclidean (Fig. 8b) distances between the aggregated and sampled profiles. Each point in the plots represents the mean distance for a given number of syllabi and the error bars represent a 95% confidence interval. The elbow points [53, 54], annotated with “X”, indicate where the rate of decrease in distance significantly slows down thus identifying the sufficient number of syllabi is equal to 9. Out of 281, 153 published institution-FOS-year combinations, 49, 750 have at least 9 syllabi (17.69%). Note that the minimum number of syllabi (i.e., 9) is obtained using nearly all the available syllabi and the y-axes of the figures are limited to 50 for the visualization purpose. Moreover, to see how the sufficient number might vary within each FOS, we redo the mentioned procedure for all the institution-FOS-year combinations within a given FOS. The minimum number of syllabi for a majority of the fields of study is between 8 and 10 with some exceptions for Transportation and Veterinary Medicine (see Fig. 9 for examples and Figure euclidean_n_syllabi_elbow_per_major on Figshare [26] for the complete list.). We publish all the calculated distances in Table manhattan_euclidean_distances on Figshare [26].

Lastly, how does the minimum number of syllabi per cohort relate to the total number of graduates with available syllabi in our dataset? We explore by counting the number of graduates per institution-FOS combination between 2003 and 2017 while varying the minimum number of syllabi from different cohorts between 2000 and 2017 (see Fig. 10). The

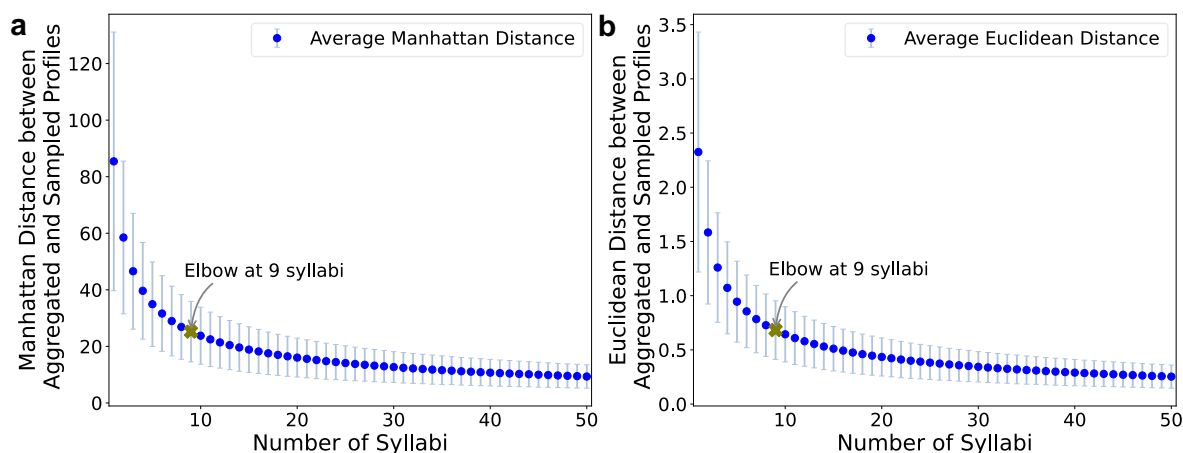


Figure 8: **Identification of the sufficient number of syllabi.** The average (a) Manhattan Distance and (b) Euclidean Distance between the aggregated and a given number of randomly selected syllabi. Each point represents the mean distance for a given number of syllabi, with their corresponding error bar. The elbow points, marked with an olive ‘X’ and annotated, indicate the sufficient number of syllabi (here, 9 syllabi) where the rate of decrease in distance significantly slows down. These points help identify the threshold beyond which additional syllabi contribute minimally to reducing the distance, providing a practical cutoff for data aggregation. Note that the x-axis is limited to 50 for the visualization purpose.

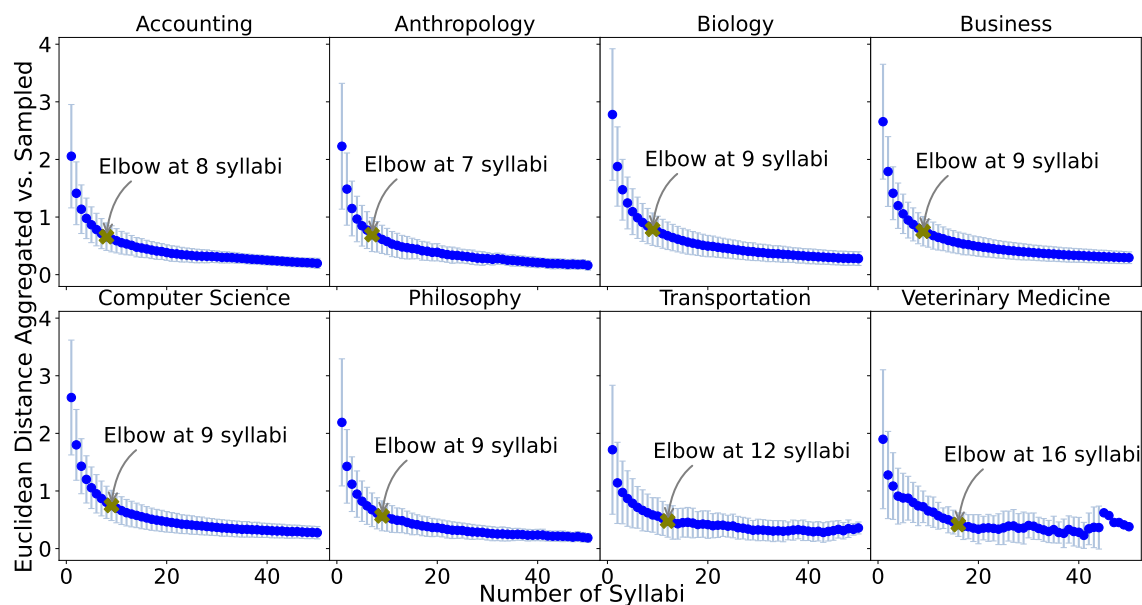


Figure 9: **Identification of the sufficient number of syllabi for sample of FOS.** The average Euclidean Distance between the the aggregated and a given number of randomly selected syllabi within each FOS. Each point represents the mean distance for a given number of syllabi, with their corresponding error bar. The elbow points, marked with an olive ‘X’ and annotated, indicate the sufficient number of syllabi (e.g., 8 syllabi for Accounting and 10 syllabi for Architecture) where the rate of decrease in distance significantly slows down. These points help identify the threshold beyond which additional syllabi contribute minimally to reducing the distance, providing a practical cutoff for data aggregation. Note that the x-axis is limited to 50 for the visualization purpose.

first cohort (2003) contains the syllabi between 2000 to 2003. For example, considering all the institution-FOS-year combinations with at least one syllabus allows us to analyze the course materials of nearly 1.1 million nationwide Bachelor’s degree graduates per year. However, restricting to institution-FOS-year combinations with at least 9 syllabi narrows the analysis to around 622, 000 graduates, i.e., around 35.46% of graduates per year.

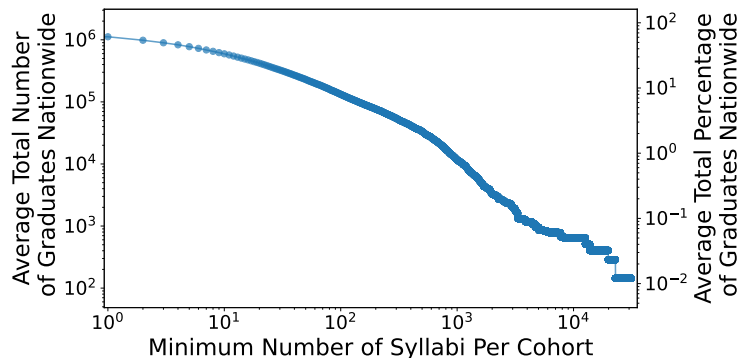


Figure 10: **The relationship between the minimum number of syllabi per cohort and the average total number of graduates nationwide between 2000 and 2017.** The x-axis shows the minimum number of syllabi per cohort (institution, FOS, period) and the y-axes show the average total number (left) and percentage (right) of graduates nationwide. The plot shows how the total number of graduates decreases as the minimum number of syllabi per cohort increases. The shaded area corresponds to the 95% confidence interval.

Qualitative Analysis of the Inferred Workplace Activities

As a face-validity check of Syllabus20*NET, we list the ten DWAs that are most strongly associated with three example fields of study (FOS): Agriculture, Biology, and Computer Science (see Fig. 11). Some DWAs (e.g., “Prepare informational or reference materials”) are common across many FOS and, therefore, obscure the DWAs that most distinguish one FOS from the others. To normalize for ubiquitous DWAs, we also present the DWAs with the greatest revealed comparative advantage (RCA) in each field (see Section “Skill Normalization”). For instance, “Plant crops, trees, or other plants” emerges as the foremost skill in *Agriculture*, “Research diseases or parasites” is predominant in *Biology*, and “Coordinate software or hardware installation” is leading in *Computer Science* according to RCA scores. Top 10 DWA per FOS contains similar results for each FOS (see Table top10_DWA_per_FOS on Figshare [26]).

Relating Fields of Study from Skill Similarity

How similar are fields of study based on their skills? Following existing work [32], we employ agglomerative hierarchical clustering technique [55] on the DWA vector representations of academic majors, aiming to elucidate their hierarchical relationships. Hierarchical clustering generates a nested sequence of clusters, allowing for an in-depth exploration of clusters at varying levels of granularity without predefining a specific number of categories (in this context, groups of majors). In this framework, FOS are deemed similar if they share the same work activities (see Figure 12 for DWAs and Figure 13 for Tasks).

The resulting dendrogram offers another face-validity check as similar FOS (e.g., STEM majors) tend to require similar DWAs. For instance, Marketing and Economics are closely related, as are Linguistics and History. Notably, just before the final clustering step, which amalgamates all majors (indicated in blue), two predominant clusters are discernible: one representing technical majors including those in STEM (in green) and the other humanities-based majors (in orange). For example, although Film and Photography is not a STEM-designated program, it requires skills that are common in STEM fields, such as “Draw detailed or technical illustrations” and “Design video game features or details” (see Table top10_DWA_per_FOS on Figshare [26]).

Dynamic Differences Between Inferred Workplace Activities and Labor Market Workplace Activities

How responsive are the skills taught in higher education to the skills required in the U.S. labor market? A “skill mismatch” may occur if higher education fails to adapt to the demands of the labor market [56, 57, 58] (e.g., by teaching more theoretical skills than practical skills [56]). Our dataset naturally offers an avenue of examining this mismatch as the scores are computed via comparison between taught content in higher education and skills in the labor market defined by the federal government. To validate this utility, we perform a similar analysis to that done by Börner and colleagues [20]. In their study, skill mentions were identified in course syllabi and in job postings to compare skills taught and demanded in computer science related fields in the US, where skills came from a skill taxonomy established by the Burning Glass Technologies. Then, Kullback-Leibler (KL) divergence was used to quantify the difference between skill distributions in course syllabi and those in job postings, and between different time periods.

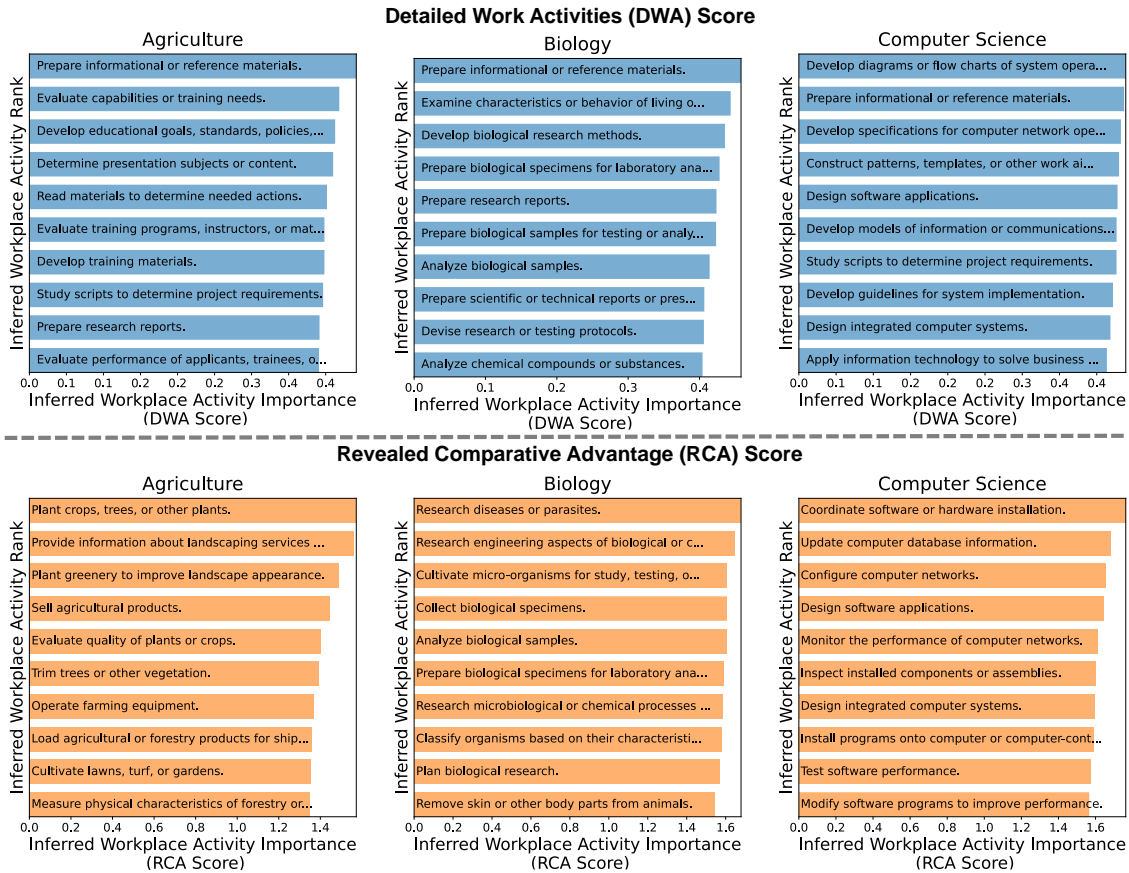


Figure 11: **The DWAs most strongly associated with Agriculture, Biology, and Computer Science.** (*top*) Top 10 inferred workplace activities with the highest DWA scores. (*bottom*) Top DWAs according to their RCA scores.

With the Course-Skill Atlas dataset, we have measures of skills taught in course syllabi. On the labor market side, each occupation in the economy can be characterized by the same measure of skill, in this case, DWAs. To generate a measure of the economy-wide skill demand, we weight by the national employment share associated with each occupation from the US Bureau of Labor Statistics (BLS) Occupation Employment and Wage Statistics (OEWS) (<https://www.bls.gov/emp/ind-occ-matrix/occupation.xlsx>). Across all FOS, the decreasing values of KL divergence over time in Fig. 14a indicate that the skills taught in course syllabi are becoming more similar to the skills required in the labor market. Specifically, comparisons between earlier syllabi (e.g., 00-03) and later labor force periods (e.g., 12-16) show a trend of decreasing divergence, reflecting an increasing alignment of skills. *This pattern indicates that educational entities are progressively aligning their course content more closely with the requirements of professional environments.* This alignment might result from a heightened recognition of occupational needs, enhancements in educational methodologies, or influences from regulatory agencies and corporate collaborations. Additionally, the reduction in KL divergence over successive periods underscores that recent syllabi not only integrate more pertinent skills but also likely eliminate outdated elements less relevant to contemporary professional demands. These gradual modifications suggest an encouraging evolution toward educational outputs that are directly advantageous for students transitioning into professional roles. These results suggest that taught skills are forward-looking. However, going beyond existing research, our dataset enables us to make direct comparisons between specific FOS and labor market dynamics for individual occupations. For example, motivated by earlier analysis of CS syllabi and CS-related job postings [20], we compare CS syllabi to Computer and Mathematical occupations (i.e., Standard Occupation Classification code: 15-0000. See Fig. 14b). Despite the trend across all FOS, over time, the labor force skill distribution becomes increasingly dissimilar to older course syllabi which confirms the rapidly changing nature of such domains. Comparing the KL Divergence scores of the syllabi among different periods (top left box of Fig. 14b), we observe that syllabi are staying stagnant, and as a result, they are moving away from the frontier of knowledge required in the labor force.

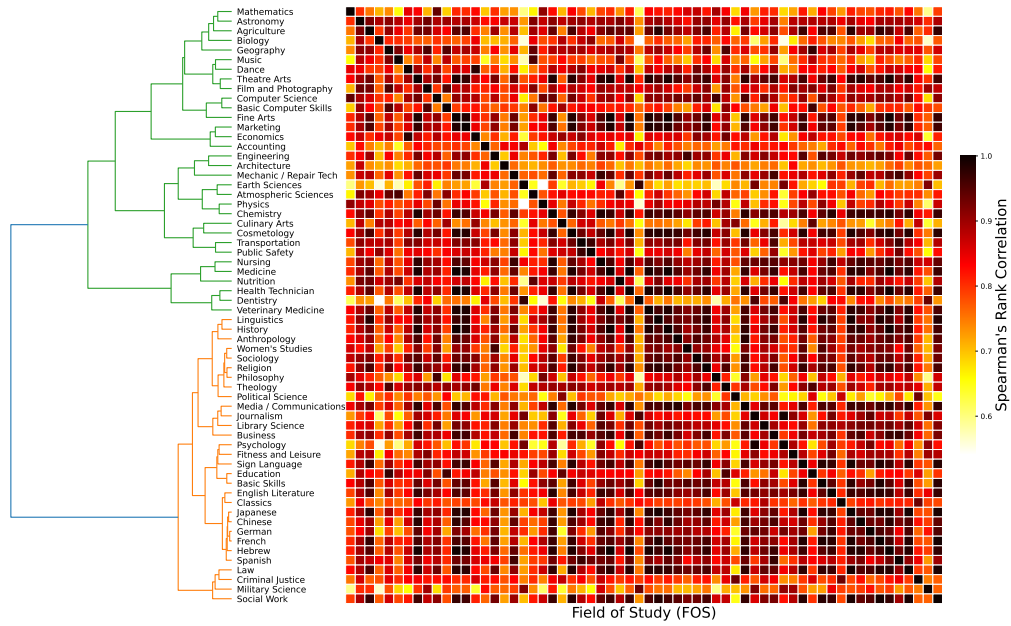


Figure 12: **The similarity of FOS according to their DWA profiles.** The heatmap shows the Spearman’s rank correlation between the fields of study and the dendrogram represents the results of hierarchically clustering similar FOS. The dendrogram on the left side organizes the FOS into clusters based on their similarity. Each branch point (node) indicates a point where two branches merge, showing the hierarchical relationship between the fields of study. Fields of study that cluster together (i.e., merge at lower levels) are more similar to each other in terms of their task profiles. For instance, closely related fields like “Physics” and “Chemistry” are grouped together.

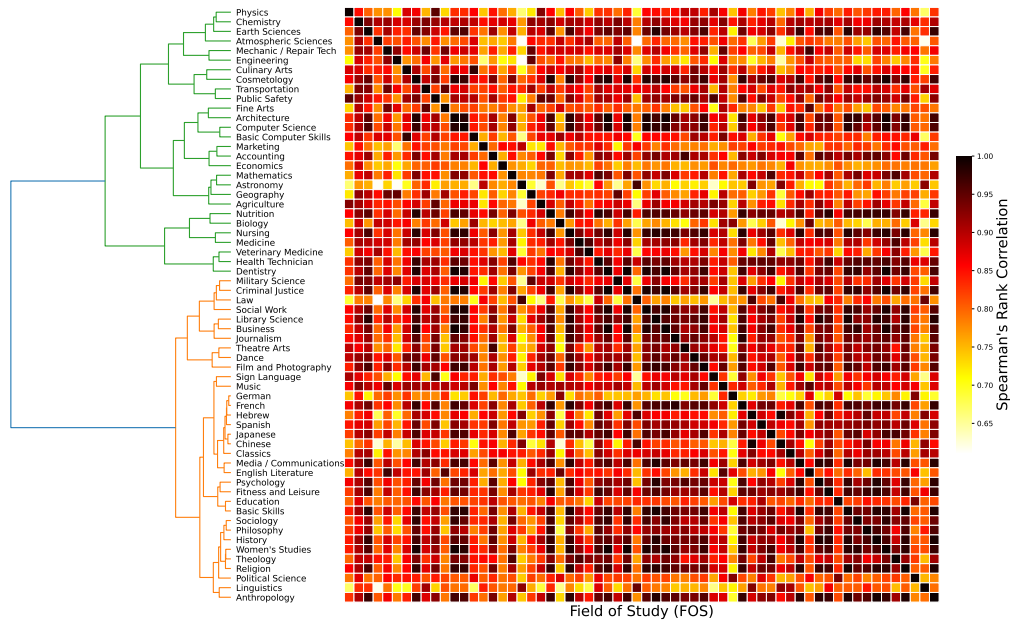


Figure 13: **The similarity of FOS according to their Task profiles.** The heatmap shows the Spearman’s rank correlation between the fields of study and the dendrogram represents the results of hierarchically clustering similar FOS. The dendrogram on the left side organizes the FOS into clusters based on their similarity. Each branch point (node) indicates a point where two branches merge, showing the hierarchical relationship between the fields of study. Fields of study that cluster together (i.e., merge at lower levels) are more similar to each other in terms of their task profiles. For instance, closely related fields like “Physics” and “Chemistry” are grouped together.

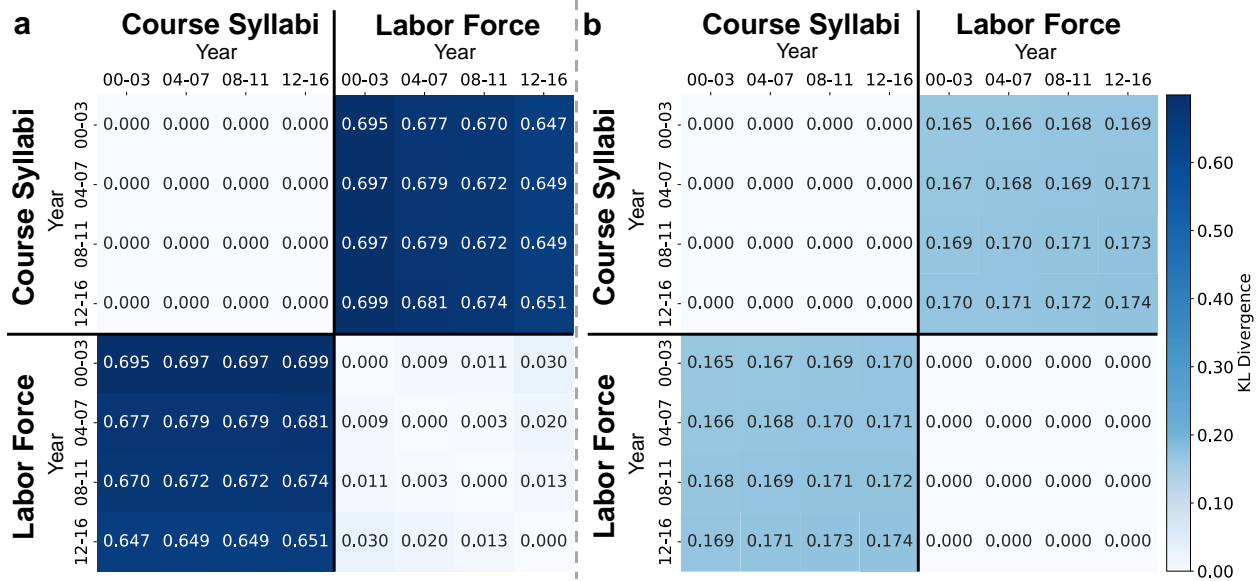


Figure 14: **Dynamic differences between skill (DWA) distributions in course syllabi and the labor force.** The matrix reports pairwise Kullback-Leibler (KL) divergence between course syllabi and the labor force for (a) syllabi from all FOS and employment-weighted O*NET DWA profiles for occupations requiring a university degree, and (b) Computer Science and Mathematics course syllabi and employment-weighted O*NET DWA profiles for Computer and Mathematical Occupations (SOC 15-0000). The off-diagonal elements in the Course Syllabi and Labor Force cells are looking at the correlation between the DWAs in those time periods and the syllabi in those time periods. These are not required to be symmetric, meaning that $D_{KL}(P \parallel Q)$ is not necessarily equal to $D_{KL}(Q \parallel P)$.

Usage Notes

Course-Skill Atlas offers a versatile tool for addressing a variety of research questions pertinent to education and workforce development across multiple domains. In the following, we briefly touch on potential research questions utilizing this dataset, including exploring differences in skill sets across gender profiles in U.S. higher education, the trend of abilities in teaching activities, and utilizing skill scores for salary estimation. Lastly, we discuss our data’s limitations.

- How does the specificity of skills taught in different college majors affect labor market outcomes, such as wages, career adaptability, and the likelihood of obtaining managerial positions?

Research has consistently shown a strong relationship between college majors, skills, and wages [59, 60, 61]. For example, some majors may offer more diverse skill sets with more general skills profiles that lead to adaptable careers after a student graduates and enters the workforce [6]. The salary gap among majors is multifaceted, involving factors like labor market demands [62] and major distribution’s effect on gender wage disparities [60]. There is a growing literature in labor and education economics on how general versus specific majors affect occupational choice and wages [63, 64]. Majors with higher specificity, such as education and nursing, generally lead to higher earnings compared to more general majors like music and psychology, driven by higher hourly wages. However, graduates from specific majors are less likely to hold managerial positions, with those from majors of average specificity being most likely to become managers [63]. Our dataset provides the opportunity to investigate such differences.

- How have teaching strategies and curriculum design evolved over time across different majors and universities?

Our dataset enables the study of skill differences within and across majors and universities over time. Taking active learning in social sciences as an example, a recent critique of active learning and the employability agenda in higher education within the social sciences [65] identified an inadvertent neglect of key skills including reading, listening, and note-taking due to the lack of proper active learning strategies. The findings from such a direction of research could pave the way for further investigation into how educational strategies can be developed to effectively balance traditional academic skills with the competencies essential for active learning environments.

- What role do educational institutions play in shaping the differences in skills between genders, particularly in relation to course syllabi?

Existing research finds that males and females tend to possess different workplace skills on aggregate [66, 67], which may correspond to gender stereotypes shaping careers [68]. But are these differences the result of education or labor market outcomes? In general, these questions can only be studied through enrollment and graduation statistics from the US Department of Education without taking into account the granular differences in taught skills across different majors and institutions. However, our dataset enables the study of this heterogeneity and, thus, creates an opportunity to explain career outcomes from the differences in skills taught during higher education—even differences within a given FOS based on varied enrollment across educational institutions. Our dataset has the potential to explain the skill differences between majors and institutions based on course syllabi.

Limitations This study produces a novel large-scale data set reflecting the skills taught to US college and university students across majors. While useful for understanding one of the major pathways for workforce development in the US, there are some limitations to the current data set. First, the syllabus dataset is, to our knowledge, the largest collection of university syllabi available, but as reported in [52] it is slightly skewed by school selectivity, overrepresenting Ivy-Plus, Elite, and selective public schools by 2.4 to 4.0 percentage points, while including less than 0.1 percent from non-selective institutions. Although the sample isn't biased by observable characteristics and consistently represents about 5% of all courses taught in US institutions post 1998, the potential for unobservable selection remains a caveat in interpreting the results. Second, we propose a new approach for inferring taught skills (i.e., O*NET DWAs and Tasks) from syllabus text, but it is difficult to confirm the effectiveness of our approach without wide-scale comprehensive exams to test the skills that students actually obtain during a course. Such an effort would be extremely cumbersome because each student would ideally be assessed on over 2,000 DWAs; it's not clear how to empirically validate each of these DWAs and implementing such an examination across universities and majors throughout the US would be an immense undertaking. Effectively, it is crucial to acknowledge that teaching does not necessarily equate to learning. Third, another limitation of this work relates to handling potential prerequisites. Prerequisites might appear in a syllabus in two forms. If they appear as administrative details (e.g., the course code), since our data preparation pipeline removes such details, they will not affect the inferred skill vector. On the other hand, when a syllabus includes the content of prerequisites — similar to learning objectives, they are processed in the same manner as the skills for the course itself. However, due to the unstructured nature of each course description (i.e., presented in a single string format), we are unable to identify and exclude these prerequisites from the final skill scores. This limitation affects the accuracy of the skill assessment by conflating course skills with prerequisites. Fourth, due to the lack of enough metadata in the raw dataset, we are unable to distinguish between undergrad and graduate courses. Fifth, our approach relies on the O*NET database which is designed to describe workers in the US workforce, and not explicitly designed to describe learning outcomes. While O*NET serves as a standardized taxonomy for communicating results to policymakers, its coverage across all occupations, and by extension, academic majors, is not uniformly comprehensive. Sixth, existing research [69] show that the distribution of the course credits varies for college students even with the same field of study. Consequently, using field of study as a stand-in for an individual's complete set of skills is inadequate. Due to the lack of data on enrollment per major and coursework taken by the students, we based our coverage analysis solely on the number of graduates per major.

Code availability

The source code for Syllabus2O*NET and DWA2Ability is available at <https://github.com/AlirezaJavadian/Syllabus-to-ONET>.

Acknowledgments

This work received funding from Russell Sage Foundation (G-2109-33808) and is supported by the University of Pittsburgh Center for Research Computing.

Author contributions statement

A.J.S. processed the data, performed all calculations, and produced all figures and statistics. All authors designed the research, analyzed the results, and reviewed the manuscript.

Competing interests

The authors declare no competing interests.

Acknowledgments

This work received funding from Russell Sage Foundation and is supported by the University of Pittsburgh Center for Research Computing.

Author contributions statement

A.J.S. processed the data, performed all calculations, and produced all figures and statistics. All authors designed the research, analyzed the results, and reviewed the manuscript.

Competing interests

The authors declare no competing interests.

References

- [1] Deming, D. & Kahn, L. B. Skill requirements across firms and labor markets: Evidence from job postings for professionals. *Journal of Labor Economics* **36**, S337–S369 (2018).
- [2] Spenner, K. I. Skill: Meanings, methods, and measures. *Work and Occupations* **17**, 399–421 (1990). <https://doi.org/10.1177/0730888490017004002>.
- [3] Attewell, P. What is skill? *Work and Occupations* **17**, 422–448 (1990). <https://doi.org/10.1177/0730888490017004003>.
- [4] Warhurst, C., Mayhew, K., Finegold, D. & Buchanan, J. *The Oxford handbook of skills and training* (Oxford University Press, 2017).
- [5] Deming, D. J. The value of soft skills in the labor market. *NBER Reporter* **4**, 7–11 (2017).
- [6] Hemelt, S. W., Hershbein, B., Martin, S. & Stange, K. M. College majors and skills: Evidence from the universe of online job ads. *Labour Economics* **85**, 102429, <https://doi.org/10.1016/j.labeco.2023.102429> (2023).
- [7] Haveman, R. & Smeeding, T. The role of higher education in social mobility. *The Future of children* 125–150 (2006).
- [8] Chetty, R., Friedman, J. N., Saez, E., Turner, N. & Yagan, D. Mobility report cards: The role of colleges in intergenerational mobility. Tech. Rep., national bureau of economic research (2017).
- [9] Kerckhoff, A. C. Education and social stratification processes in comparative perspective. *Sociology of education* 3–18 (2001).
- [10] Acemoglu, D. & Autor, D. Chapter 12 - skills, tasks and technologies: Implications for employment and earnings. In Card, D. & Ashenfelter, O. (eds.) *Handbook of Labor Economics*, vol. 4, 1043–1171, [https://doi.org/10.1016/S0169-7218\(11\)02410-5](https://doi.org/10.1016/S0169-7218(11)02410-5) (Elsevier, 2011).
- [11] Altonji, J., Arcidiacono, P. & Maurel, A. Chapter 7 - the analysis of field choice in college and graduate school: Determinants and wage effects. In Hanushek, E. A., Machin, S. & Woessmann, L. (eds.) *Handbook of the Economics of Education*, vol. 5, 305–396, <https://doi.org/10.1016/B978-0-444-63459-7.00007-5> (Elsevier, 2016).
- [12] Altonji, J. G., Blom, E. & Meghir, C. Heterogeneity in human capital investments: High school curriculum, college major, and careers. *Annual Review of Economics* **4**, 185–223 (2012). <https://doi.org/10.1146/annurev-economics-080511-110908>.
- [13] Triventi, M. The role of higher education stratification in the reproduction of social inequality in the labor market. *Research in social stratification and mobility* **32**, 45–63 (2013).
- [14] Lovenheim, M. & Smith, J. Chapter 4 - returns to different postsecondary investments: Institution type, academic programs, and credentials. In Hanushek, E. A., Machin, S. & Woessmann, L. (eds.) *Handbook of the Economics of Education*, vol. 6, 187–318, <https://doi.org/10.1016/bs.hesedu.2022.11.006> (Elsevier, 2023).

- [15] Shapiro, D. *et al.* Completing college: A national view of student attainment rates by race and ethnicity—fall 2010 cohort (signature 12 supplement). *National Student Clearinghouse* (2017).
- [16] Bonvillian, W. B. & Sarma, S. E. *Workforce education: a new roadmap* (MIT Press, 2021).
- [17] Taylor, P. *et al.* *Is college worth it?* (Pew Social and Demographic trends, 2011).
- [18] Yu, R. *et al.* A research framework for understanding education-occupation alignment with NLP techniques. In *Proceedings of the 1st Workshop on NLP for Positive Impact*, 100–106, 10.18653/V1/2021.NLP4POSIMPACT-1.11 (Association for Computational Linguistics (ACL), 2021).
- [19] Light, J. Student demand and the supply of college courses. *Available at SSRN 4856488* (2024).
- [20] Börner, K. *et al.* Skill discrepancies between research, education, and jobs reveal the critical need to supply soft skills for the data economy. *Proceedings of the National Academy of Sciences* **115**, 12630–12637 (2018).
- [21] del Pilar Garcia-Chitiva, M. & Correa, J. C. Soft skills centrality in graduate studies offerings. *Studies in Higher Education* **49**, 956–980, 10.1080/03075079.2023.2254799 (2024).
- [22] Chau, H., Bana, S. H., Bouvier, B. & Frank, M. R. Connecting higher education to workplace activities and earnings. *PLOS ONE* **18**, e0282323, 10.1371/JOURNAL.PONE.0282323 (2023).
- [23] Desikan, B. S. & Evans, J. Misalignment between skills discovered, disseminated, and deployed in the knowledge economy. *Journal of Social Computing* **3**, 191–205, 10.23919/JSC.2022.0013 (2022).
- [24] Chang, X., Wang, B. & Hui, B. Towards an automatic approach for assessing program competencies. In *Proceedings of the 12th International Learning Analytics and Knowledge Conference (LAK '22)*, 119–129, 10.1145/3506860.3506875 (ACM, 2022).
- [25] Lastra-Anadon, C. X., Das, S., Varshney, K. R., Raghavan, H. & Yu, R. How universities can mind the skills gap (higher education and the future of work) (2021).
- [26] Javadian Sabet, A., Bana, S. H., Yu, R. & Frank, M. Course-Skill Atlas: A national longitudinal dataset of skills taught in U.S. higher education curricula, 10.6084/m9.figshare.25632429.v7 (2024).
- [27] O*NET OnLine, (2024). Available at <https://www.onetonline.org/>.
- [28] Akour, M. & Alenezi, M. Higher education future in the era of digital transformation. *Education Sciences* **12**, 784 (2022).
- [29] Brasca, C., Krishnan, C., Marya, V., Owen, K. & Sirois, J. How technology is shaping learning in higher education. *McKinsey & Company. June* **15** (2022).
- [30] Lim, J., Aklın, M. & Frank, M. R. Location is a major barrier for transferring us fossil fuel employment to green jobs. *Nature Communications* **14**, 5711 (2023).
- [31] Frank, M. R., Sun, L., Cebrian, M., Youn, H. & Rahwan, I. Small cities face greater impact from automation. *Journal of the Royal Society Interface* **15**, 20170946 (2018).
- [32] Chau, H., Bana, S. H., Bouvier, B. & Frank, M. R. Connecting higher education to workplace activities and earnings. *Plos one* **18**, e0282323 (2023).
- [33] Alabdulkareem, A. *et al.* Unpacking the polarization of workplace skills. *Science advances* **4**, eaao6030 (2018).
- [34] Frank, M. R. *et al.* Toward understanding the impact of artificial intelligence on labor. *Proceedings of the National Academy of Sciences* **116**, 6531–6539 (2019).
- [35] Moro, E. *et al.* Universal resilience patterns in labor markets. *Nature communications* **12**, 1972 (2021).
- [36] Frank, M. R. *et al.* Network constraints on worker mobility. *Nature Cities* **1**, 94–104 (2024).
- [37] Agnihotri, A. & Misra, R. K. Managerial competencies: A comparative study of us-india employer’s needs. *Global Business and Organizational Excellence* **43**, 92–106 (2024).
- [38] Cabell, A. Supporting the career development of black adults during the covid-19 pandemic. *Journal of Employment Counseling* **60**, 62–71 (2023).
- [39] Council, N. R. *et al.* A database for a changing economy: Review of the occupational information network (o*net) (2010).
- [40] Qi, P., Zhang, Y., Zhang, Y., Bolton, J. & Manning, C. D. Stanza: A python natural language processing toolkit for many human languages. *arXiv preprint arXiv:2003.07082* (2020).
- [41] Reimers, N. & Gurevych, I. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084* (2019).

- [42] Lo, K., Wang, L. L., Neumann, M., Kinney, R. & Weld, D. S2ORC: The semantic scholar open research corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4969–4983 (Association for Computational Linguistics, Online, 2020).
- [43] Fader, A., Zettlemoyer, L. & Etzioni, O. Open Question Answering Over Curated and Extracted Knowledge Bases. In *KDD* (2014).
- [44] Lewis, P. *et al.* Paq: 65 million probably-asked questions and what you can do with them. *Transactions of the Association for Computational Linguistics* **9**, 1098–1115 (2021).
- [45] Khashabi, D. *et al.* Gooaq: Open question answering with diverse answer types. *arXiv preprint* (2021).
- [46] Reimers, N. & Inui, K. Sentence-transformers/all-mpnet-base-v2 · hugging face (2019).
- [47] Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011).
- [48] Hoen, A. R. & Oosterhaven, J. On the measurement of comparative advantage. *The Annals of Regional Science* **40**, 677–691 (2006).
- [49] Hidalgo, C. A. & Hausmann, R. The building blocks of economic complexity. *Proceedings of the national academy of sciences* **106**, 10570–10575 (2009).
- [50] Hausmann, R. & Hidalgo, C. A. The network structure of economic output. *Journal of economic growth* **16**, 309–342 (2011).
- [51] Hartmann, D., Guevara, M. R., Jara-Figueroa, C., Aristarán, M. & Hidalgo, C. A. Linking economic complexity, institutions, and income inequality. *World development* **93**, 75–93 (2017).
- [52] Biasi, B. & Ma, S. The education-innovation gap. Working Paper 29853, National Bureau of Economic Research (2022). 10.3386/w29853.
- [53] Ng, A. Clustering with the k-means algorithm. *Machine Learning* 1–2 (2012).
- [54] Kodinariya, T. M., Makwana, P. R. *et al.* Review on determining number of cluster in k-means clustering. *International Journal* **1**, 90–95 (2013).
- [55] Johnson, S. C. Hierarchical clustering schemes. *Psychometrika* **32**, 241–254 (1967).
- [56] Aloysius, O. I., Ismail, I. A., Suandi, T. & Arshad, M. Enhancing university’s and industry’s employability-collaboration among nigeria graduates in the labor market. *International Journal of Academic Research in Business and Social Sciences* **8**, 32–48 (2018).
- [57] Pujol-Jover, M., Riera-Prunera, C. & Abio, G. Competences acquisition of university students: Do they match job market’s needs? *Intangible Capital* **11**, 612–626 (2015).
- [58] Jackson, D. & Tomlinson, M. The relative importance of work experience, extra-curricular and university-based activities on student employability. *Higher Education Research & Development* **41**, 1119–1135 (2022).
- [59] Grogger, J. & Eide, E. Changes in college skills and the rise in the college wage premium. *Journal of Human Resources* 280–310 (1995).
- [60] Eide, E. College major choice and changes in the gender wage gap. *Contemporary Economic Policy* **12**, 55–64 (1994).
- [61] Long, M. C., Goldhaber, D. & Huntington-Klein, N. Do completed college majors respond to changes in wages? *Economics of Education Review* **49**, 1–14 (2015).
- [62] Altonji, J. G., Kahn, L. B. & Speer, J. D. Cashier or consultant? entry labor market conditions, field of study, and career success. *Journal of Labor Economics* **34**, S361–S401 (2016).
- [63] Leighton, M. & Speer, J. D. Labor market returns to college major specificity. *European Economic Review* **128**, 103489 (2020).
- [64] Martin, S. M. College major specificity, earnings growth, and job changing. *University of Michigan Doctoral Dissertation* (2022).
- [65] David, M. & Maurer, H. Reclaiming agency: skills, academics and students in the social sciences. *European Political Science* 1–17 (2022).
- [66] Christl, M. & Köppl-Turyna, M. Gender wage gap and the role of skills and tasks: evidence from the austrian piasac data set. *Applied Economics* **52**, 113–134 (2020).
- [67] Azmat, G. & Ferrer, R. Gender gaps in performance: Evidence from young lawyers. *Journal of Political Economy* **125**, 1306–1355 (2017).

-
- [68] Fana, M., Villani, D. & Bisello, M. Gender gaps in power and control within jobs. *Socio-Economic Review* **21**, 1343–1367 (2022). <https://academic.oup.com/ser/article-pdf/21/3/1343/50969983/mwac062.pdf>.
- [69] Light, A. & Schreiner, S. College major, college coursework, and post-college wages. *Economics of Education Review* **73**, 101935, <https://doi.org/10.1016/j.econedurev.2019.101935> (2019).