# EDUKG: a Heterogeneous Sustainable K-12 Educational Knowledge Graph

Bowen Zhao[1*], Jiuding Sun[2,3*], Bin Xu[1,3(✉)], Xingyu Lu[3], Yuchen Li[3], Jifan Yu[3], Minghui Liu[3], Tingjian Zhang[3], Qiuyang Chen[3], Hanming Li[3], Lei Hou[3], and Juanzi Li[3]

[1] Global Innovation Exchange, Tsinghua University, Beijing, China
`{zhaobw21,sjd22}@mails.tsinghua.edu.cn`
[2] Khoury College of Computer Sciences, Northeastern University, Boston MA, USA
[3] Department of Computer Science and Technology,
Tsinghua University, Beijing, China
`xubin@tsinghua.edu.cn`

**Abstract.** Web and artificial intelligence technologies, especially semantic web and knowledge graph (KG), have recently raised significant attention in educational scenarios. Nevertheless, subject-specific KGs for K-12 education still lack sufficiency and sustainability from knowledge and data perspectives. To tackle these issues, we propose EDUKG, a heterogeneous sustainable K-12 **Edu**cational **K**nowledge **G**raph. We first design an interdisciplinary and fine-grained ontology for uniformly modeling knowledge and resource in K-12 education, where we define 635 classes, 445 object properties, and 1314 datatype properties in total. Guided by this ontology, we propose a flexible methodology for interactively extracting factual knowledge from textbooks. Furthermore, we establish a general mechanism based on our proposed generalized entity linking system for EDUKG's sustainable maintenance, which can dynamically index numerous heterogeneous resources and data with knowledge topics in EDUKG. We further evaluate EDUKG to illustrate its sufficiency, richness, and variability. We publish EDUKG with more than 252 million entities and 3.86 billion triplets. Our code and data repository is now available at https://github.com/THU-KEG/EDUKG.

**Keywords:** Ontology · Knowledge Graph · K-12 Education

## 1 Introduction

*The object of education is to prepare the young to educate themselves throughout their lives*, as said by Robert M. Hutchins. Education, especially for K-12 children, plays a significant role in everyone's life. Intelligent education, which aims to leverage the Web and artificial intelligence (AI) technologies to improve students' learning efficiency [13,19], has always been an essential topic for researchers. In addition, the construction of educational knowledge graphs (KGs)

---

[*] Equal Contributions.

is a fundamental research with variable downstream applications, including educational data mining [22], learning management systems [1], question answering platforms [30], dialogue systems [18], etc.

Plentiful educational KGs have been proposed to help the development of computer-aided educational technologies. KnowEdu [4] and K12Edukg [3] are constructed by extracting concepts and prerequisite rules from subject-specific textbooks. However, entities in these KGs are only course concepts without other essential educational resources for students. CKGG [23] is proposed based on Chinese high-school-level geography education, yet they only integrate data for location entities. Meanwhile, several educational KGs are proposed based on massive online open courses (MOOCs). For instance, MOOC-KG [6] and HEKG [34] are built upon open course data, yet their ontology can only represent subject-specific knowledge at a shallow level. In particular, there are only 4 and 6 defined classes in MOOC-KG and HEKG, respectively. Furthermore, although KGs built upon open courses consist of heterogeneous data, they cannot dynamically develop with growing resources. Additionally, most KGs based on MOOCs are designed for higher education instead of K-12 education.

Despite that several KGs have been proposed for educational usage, they suffer from the following limitations:

• **Insufficient Knowledge Modeling.** Prior research pointed out that interdisciplinary teaching is beneficial for developing students' critical thinking, creativity, communication, and essential academia [12]. In the meantime, fine knowledge granularity is also beneficial for students' learning process [25]. Nevertheless, existing educational KGs only represent subject-specific knowledge on a coarse-grained level, lacking interdisciplinary entity relations.

• **Sophisticated Data Curation.** Education aims to teach students with broad ability instead of just knowledge in textbooks [28]. Educational resources, including examination questions and beyond, are proved to be beneficial for fostering students' abilities through *learning by doing* [2,24]. Also, existing data repositories for education, such as MOOCCUBEX [32] leverages a concept graph to organize heterogeneous data altogether. However, existing educational KGs still lack adequate resources due to data heterogeneity.

• **Neglected Information Growth.** Information for education is ever-growing from both knowledge and data perspectives. For knowledge, the educational reform in China is consistently changing the essential knowledge of education through time. For data, there are increasing online materials for students to learn. Nonetheless, prior educational KGs lack maintenance sustainability, i.e., the ability to capture and infuse new knowledge and resources incrementally.

To address these issues, we conclude that educational KGs should be built with an interdisciplinary schema that can represent not only knowledge but also resources. Meanwhile, towards maintaining sustainability, an educational KG should be able to grow and adapt incrementally to the change of real-world knowledge. Therefore, we propose EDUKG, a heterogeneous sustainable K-12 **Edu**cational **K**nowledge **G**raph for Chinese high-school-level education. We design an interdisciplinary fine-grained ontology that uniformly models knowledge,
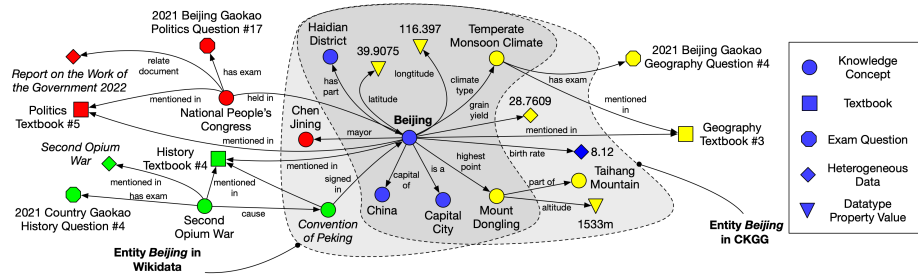
Fig. 1: EDUKG's data sufficiency compared with other KGs. Blue, yellow, red, and green points refer to knowledge and resources in interdisciplinary, geography, politics, and history subjects, respectively.

resources, and heterogeneous data. In total, we define 635 classes and 1759 properties in EDUKG ontology without subject boundaries. Guided by this ontology, we propose a semi-automated method for interactively acquiring knowledge from textbooks. Fig. 1 compares data in EDUKG with Wikidata[1] and CKGG, indicating that EDUKG consists of most sufficient information from both knowledge and data perspectives. Additionally, for sustainably maintaining EDUKG with growing data, we propose a general mechanism to index heterogeneous online data incrementally based on our proposed entity linking technique.

**Contributions.** In general, our contributions are summarized as follows:

1. An interdisciplinary, fine-grained ontology uniformly represents K-12 educational knowledge, resources, and heterogeneous data with 635 classes, 445 object properties, and 1314 datatype properties;
2. A large-scale, heterogeneous K-12 educational KG with more than 252 million entities and 3.86 billion triplets based on the data from massive educational and external resources;
3. A flexible and sustainable construction and maintenance mechanism empowers EDUKG to evolve dynamically, where we design guiding schema of the construction methodology as *hot-swappable*, and we simultaneously monitor 32 different data sources for incrementally infusing heterogeneous data.

**Outline.** In the following sections, we first illustrate the ontology for EDUKG in Sec. 2, and we present EDUKG construction and maintenance mechanisms in Sec. 3. Afterward, in Sec. 4, we introduce essential characteristics of EDUKG to prove its sufficient qualities. In Sec. 5, we present the impact and availability of EDUKG with its data, code, and applications. Finally, the related works are introduced in Sec. 6, and we conclude our paper in Sec. 7

## 2    Schema of EDUKG

In this section, we introduce the ontology of EDUKG, which uniformly represents knowledge, resource, and heterogeneous data.
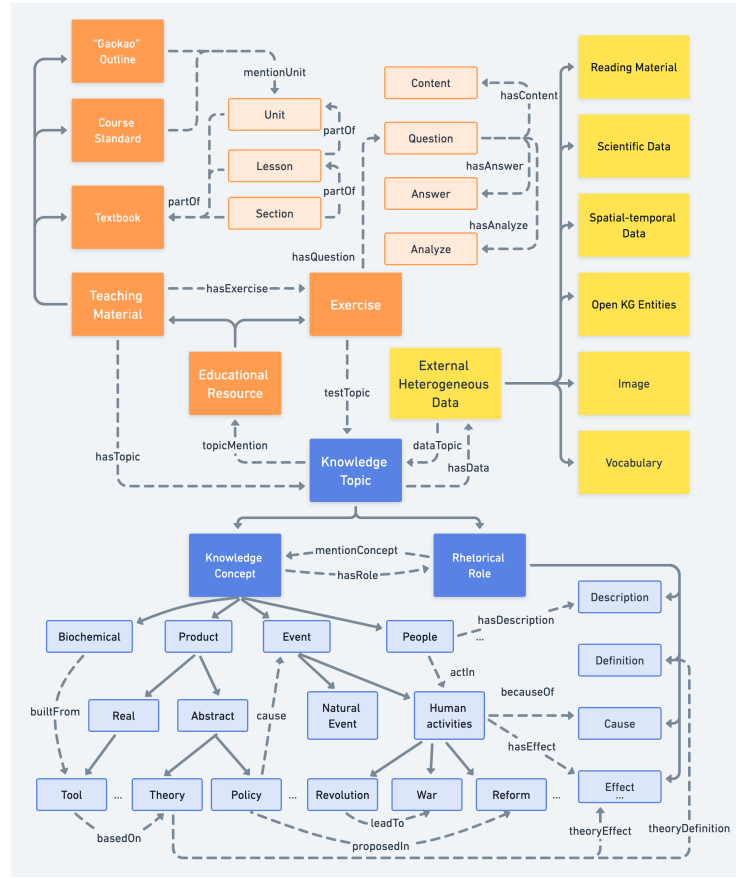
---

[1] https://www.wikidata.org

Fig. 2: An overview of EDUKG top-level ontology.

## 2.1   Overview of EDUKG Ontology

As shown in Fig. 2, we divide EDUKG ontology into three main sections, which are "Knowledge Topic", "Educational Resource", and "External Heterogeneous Data". Here we define the three top-level classes in EDUKG as follows:

– "Knowledge Topic": essential themes in some specific subjects [10] and their essential rhetorical roles.
– "Educational Resource": intra-curricular teaching and testing resources in K-12 education, for example, textbooks and examination exercises.
– "External Heterogeneous Data": extra-curricular resources and data give students vast approaches to learning more comprehensive knowledge.

Since EDUKG contains educational knowledge and resources, we investigate and follow multiple published knowledge and resource modeling standards. For knowledge, we reuse vocabularies from the widely-adopted RDF and RDFS

schema. As for resources, we adopt the LRMI Standard[2]. Furthermore, we use OWL for ontology representation and reuse the schema in existing ontologies and KGs, such as schema.org[3], YAGO [17], and Wikidata.

## 2.2 Intra- and Extra-curricular Resources

We divide resources in EDUKG into two sub-classes, i.e., intra-curricular "Educational Resource" and extra-curricular "External Heterogeneous Data". Meanwhile, "Educational Resource" consists of two main sub-classes, "Teaching Material" and "Exercise", where the former focuses on *learning to do* and the latter focuses on *learning by doing*. Their detailed composition is shown in Fig. 3.
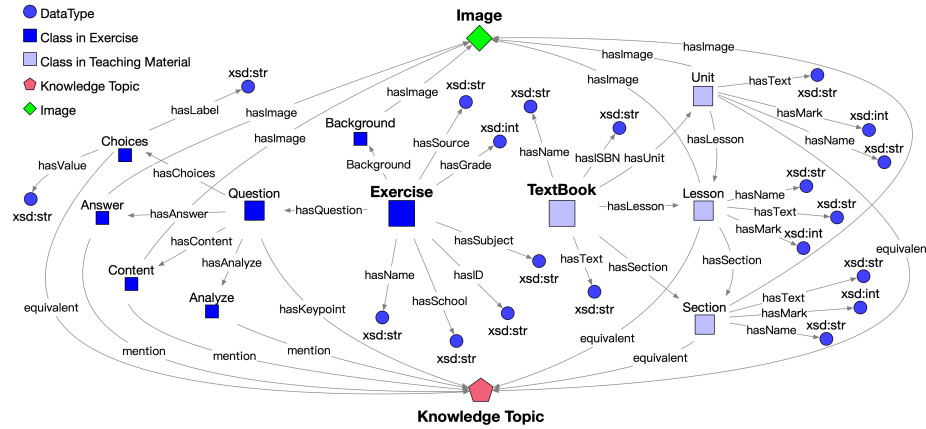


Fig. 3: Schema under "Educational Resource" class in EDUKG.

**Educational Resource.** We first illustrate the two main sub-classes of "Educational Resource". "Teaching Material" represents materials delivering essential knowledge to students, i.e., *learning to do*, including textbooks, learning guidance, curriculum standards, etc. Meanwhile, "Exercise" refers to exam questions from "Gaokao"[4] and textbooks, which can be considered as tools to both examine and promote students' level of understanding and applying the knowledge they have learned, i.e., their ability of *learning by doing*.

**External Heterogeneous Data.** To increase the coverage of extra-curricular knowledge topics mentioned or even specifically tested in "Gaokao", we design the "External Heterogeneous Data" class to represent heterogeneous data from multiple sources to enhance EDUKG. To support downstream tasks across different subjects, we gather data from multiple sources on the Web, including webpages, images, tabular data, etc. Hence, We define the sub-classes of "External Heterogeneous Data" according to their format, as shown in Fig. 2.

---

[2] https://www.dublincore.org/specifications/dublin-core/dces

[3] https://schema.org/

[4] also known as the National College Entrance Examination (NCEE) in China

### 2.3   Interdisciplinary Knowledge Topic

For knowledge modeling, we propose a fine-grained and interdisciplinary sub-ontology for EDUKG's "Knowledge Topic" class with two sub-classes, namely "Knowledge Concept" and "Rhetorical Role".

**Knowledge Concept.** "Knowledge Concept" is a vital sub-class of "Knowledge Topic" for it refers to the main subjects taught during the course [31], for instance, *Industrial Revolution*, *Bourgeoisie*, *Britain*, etc. Since there are numerous existing ontologies, we do not need to build this sub-ontology from scratch. We adopt and modify ontologies from schema.org and YAGO as the top-level sub-ontology. For fine-grained classes on the bottom side of the hierarchy, we infuse classes and properties from several subject-specific ontologies [9,15,29] proposed for K-12 education. We further enrich essential properties for these classes according to a series of popular study guides for Chinese high-school education. In all, we create 585 classes, 336 object properties, and 1177 datatype properties for EDUKG "Knowledge Concept" sub-ontology.

Table 1: Essential Rhetorical Roles in EDUKG.

| Name | Description | Example |
|------|-------------|---------|
| Definition | Specifically defining as truth in the context of K-12 education | Equation **is defined as** the mathematical statement consisting of an equal symbol. |
| Process | Processes, developments and operations | **Step 1.** Formulating a hypothesis. |
| Mechanism | Describing mechanism and theory. | Fire extinguishers **work by** separating the fuel from the oxygen. |
| Reason | Expressing reasons and explanations | The emergence of capitalism is one of the **cause** of industrial revolution. |
| Effect | Expressing cause and effect | An increase of wealth is one of the **effect** of industrial revolution. |
| Significance | Describing Significance | Carbon dioxide is an **important** greenhouse gas that helps to trap heat. |
| Condition | Stating the condition of a proposition | The domain of the equation **must be** the subset of all real numbers. |

**Rhetorical Role.** Unlike open-domain KGs, for further improving knowledge granularity, knowledge topics in educational KGs should contain more than basic course concepts, i.e., named entities. The reason is that some rhetorical roles of concepts are also crucial for students during their study procedure. Rhetorical roles are defined as semantic units that segment documents into coherent units of information [16]. For example, the concept *Industrial Revolution* has the rhetorical role *the influence of Industrial Revolution toward Chinese society*, which is also a vital knowledge topic. Moreover, if only regarding rhetorical roles as datatype properties of knowledge concepts, some core concepts mentioned in a rhetorical role's name and content cannot be fully expressed using relations between entities in the KG. Meanwhile, using more sophisticated representations such as first-order logic to represent rhetorical roles is too difficult for users to understand. Thus, we define the "Rhetorical Role" class as sentences or phrases that illustrate a concept's essential properties, such as description, reason, re-

sult, significance, etc. The detailed definition and examples of some essential rhetorical roles are shown in Table 1.
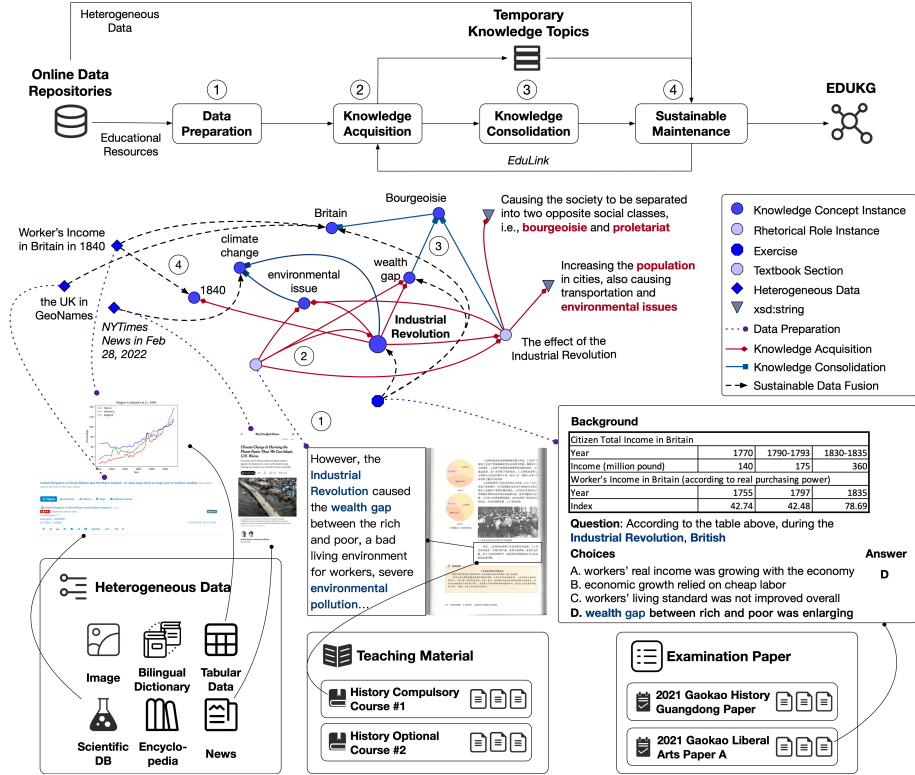
## 3   Construction and Maintenance of EDUKG



Fig. 4: The construction framework of EDUKG.

In this section, we introduce the EDUKG construction and maintenance framework as shown in Fig. 4. We first introduce our data preparation procedure. We then use a semi-automated method to acquire and consolidate essential knowledge from teaching resources. Finally, we propose a sustainable EDUKG maintenance mechanism generalized from our proposed entity linking system that can consistently index massive, heterogeneous online data.

### 3.1   Data Preparation

**Teaching Material.** We first gather a complete series of textbooks for Chinese K-12 education from an online repository, including 208 books that cover 16 different subjects. Towards data adequacy, we only select 46 books from the eight main subjects (literature, mathematics, physics, chemistry, biology, history,

geography, and politics) of Chinese high-school-level education since they are highly correlated with "Gaokao" in most regions of China.

To increase the granularity of teaching materials, as shown in the bottom middle part of Fig. 4, we develop an HTML-based parsing algorithm for segmenting the selected 46 textbooks into 256 units, 779 lessons, and 2371 sections in total. Afterward, we manually identify 2608 topics from *Chinese national curriculum standards for high-school education* and *"Gaokao" outline* as the key topic candidates for those sections. We identify key concepts for each section according to their semantic similarity and prerequisite rule. In general, we summarize the similarity score between sections and topics as follows:

$$\mathrm{S}(d_i, c) = \mathcal{D}(d_i, c) \cdot \prod_{j=1}^{i-1} \mathcal{O}(d_j, c) \tag{1}$$

where $d_i$ refers to the i-th section in textbooks, while $c$ is the identified concept. $\mathcal{O}(d_j, c)$ the truth function denotes whether concept $c$ is mentioned in section $d_j$, and $\mathcal{D}$ is the cosine-similarity function of TF-IDF embedding.

**Examination Papers.** For timeliness, we only gather "Gaokao" exam papers within the past five years. We collect data from two online examination databases and gather 302 examination papers with 6518 exercises. We categorize the exercises into five sub-classes, namely *choice question*, *text-filling question*, *number-filling question*, *question and answer*, and *writing question*. Afterward, we split each exercise into four main sections, i.e., "Background", "Question", "Answer", and "Analysis". More specifically, for *choice question*, there is another special section "Choice" defined for it as illustrated in Fig. 4. To resolve exercise data according to the schema mentioned in Sec. 2.2, we develop a template-based method to parse and classify those 6518 exercises. Then we manually correct some wrongly parsed exercises. Additionally, we ask human annotators to manually identify the key concepts of each exercise based on its content.

### 3.2   Semi-automated Knowledge Acquisition

To construct EDUKG with high precision and broad coverage, we apply a semi-automated method for acquiring knowledge from teaching materials. We further adopt several NLP techniques to consolidate EDUKG to improve its fact coverage and completeness.

**Knowledge Acquisition from Textbooks.** To acquire key topics from teaching materials, we leverage both named entity recognition (NER) and entity linking methods for detecting in-text entities because NER can provide broad coverage while entity linking can ensure precision. For NER, We fine-tune the Chinese RoBERTa [5] model on CLUENER [26] dataset for extracting fine-grained entities from textbooks. As for EL, we use XLink [33] system for discovering open-domain knowledge concepts. We also use the entity linking system proposed in EDUKG (EduLink), which will be introduced in detail in Sec. 3.3. As shown in the middle part of Fig. 4, the knowledge concepts "Industrial Revolution" and "wealth gap" are identified via the XLink system. It is worth noting

that the backbone KG for EduLink is interactively changing during the human annotation. Each identified entity's confidence level is calculated as:

$$P(c) = S(c) + \alpha \cdot (f_{pos}(c) + f_{neg}(c)) \qquad (2)$$

where $S(c)$ is the sum of scores calculated by NER and EL algorithms, and $f_{pos}, f_{neg}$ correspond to the frequency of positive and negative labels of the given entity, while $\alpha$ is a pre-set parameter indicates the weights of human feedback.

Moreover, for easily mining relations from given documents, we further align entities' infobox data with their potential mentions in the document as candidate knowledge triplets. Meanwhile, the entities' co-occurrence pairs are also regarded as candidates. Subsequently, we use the OpenIE API provided by News-Miner [8] system for jointly extracting open entities and relations from textbook documents. Finally, We merge the entities and relations with span information, and we map the predicates in extracted relations to the properties defined in EDUKG ontology by calculating similarity scores via Sentence-BERT [21] model in Text2vec Toolkit [27].

We invite ten annotators with outstanding scores in "Gaokao" to label the extracted knowledge triplet candidates from teaching materials based on our extracted knowledge triplets. We divide the labeling into two stages to tackle the error-propagation issue and reduce cost labor. The first stage focuses on entity (knowledge topic) recognition, and the second focus on triplets, i.e., relations and properties extraction. Furthermore, since there is no existing entity in EduLink at the start of the annotation procedure, we infuse several pre-built KGs based on a series of study guides in China to solve this cold-start issue. Overall, we design this ontology-guided knowledge acquisition method as *hot-swappable*, where the guiding ontology is decoupled with the knowledge acquisition methods, so we can flexibly change its schema based on the latest curriculum standards in China. Our detailed knowledge acquisition task design is illustrated in the Sec. A of appendix in our supplementary materials.

**Self-supervised Knowledge Consolidation.** KGs built through handcrafted or semi-automated methods may suffer from several issues from both entity and relation perspectives, including **(1) Insufficient Concept Coverage** and **(2) Incomplete Entity Relation**. Due to these issues, we leverage the following techniques to consolidate EDUKG and improve its scalability and completeness.

• **Knowledge Concept Expansion.** We adopt the entities from Xlore [11] to expand the concepts in EDUKG for enriching its knowledge coverage. For each entity $e$ in EDUKG, we first manually identify its equivalent entity $\hat{e}$ in Xlore. We calculate the scores of all neighbouring nodes $c$ of $\hat{e}$ as follows

$$\text{score}(c) = \text{Sim}(\hat{e}, c) \cdot \frac{1}{|\mathcal{N}(e)|} \sum_{n \in \mathcal{N}(e)} w_n \text{Sim}(n, c) \qquad (3)$$

where $\mathcal{N}(e)$ denotes the neighbours of entity $e$ in EDUKG, while $w_n$ indicates the weight assigned to neighbour $n$ according to its relation with $\hat{e}$. We use the cosine similarity of Sentence-BERT embedding as the similarity function $\text{Sim}(\cdot)$.

Fig. 4 indicates that by leveraging external KG, we can infuse more knowledge concepts such as entity "climate change" from Wikidata.

• **Rhetorical Roles Linking.** For the convenience of data labeling, we only ask the annotators to identify rhetorical roles as concepts' property values at first. To identify rhetorical roles, we first design templates for each sub-class of "Rhetorical Role" and use a template-based method to recognize them from the entities' datatype properties. After recognizing them, we leverage an entity mention detection method for linking rhetorical roles with its mentioned knowledge concepts. As shown in the middle right part of Fig. 4, we can link the rhetorical role "The effect of the Industrial Revolution" with the concepts "Bourgeoisie" and "environmental issues" based on its content. In total, we recognize 18080 rhetorical roles and detect 20596 concept mentions in all rhetorical roles.
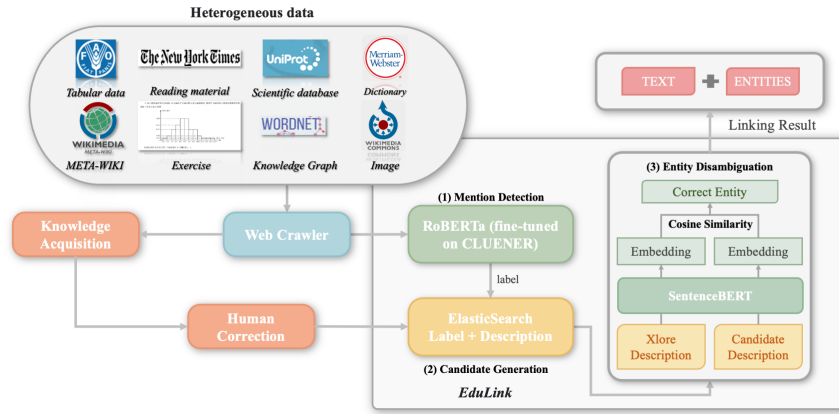
### 3.3   Sustainable Maintenance



Fig. 5: EDUKG Sustainable Maintenance Mechanism.

Unlike most traditional KGs built once and for all, we focus on the sustainable maintenance of our proposed EDUKG. Specifically, we regard the sustainability of EDUKG from two different aspects: (1) **Knowledge aspect**, as we have mentioned in Sec. 3, we design the ontology-guided knowledge acquisition method as *hot-swappable*, where the schema of the knowledge acquisition can be flexibly modified based on the latest curriculum standards; (2) **Data aspect**, as shown in Fig. 5, we propose and generalize our educational entity linking system EduLink, which leverages knowledge topics in EDUKG to dynamically index massive online data gathered through Web crawlers.

**Heterogeneous Data Compilation.** We consistently gather massive, heterogeneous online data, including news, government policies, statistics index, lexical databases, domain-specific KGs, etc., from multiple sources via Web crawlers, and the sources of these heterogeneous data are shown in Table 2.

Table 2: Heterogeneous Data Sources for EDUKG.

| Related Subject | Source | Format | Description |
|---|---|---|---|
| Biochemics | Basechem | CSV | Properties of chemicals |
| | PubChem | TTL | Chemical KG |
| | BioGRID | TSV | Protein and protein reaction |
| | UniProt | DAT | Protein and gene |
| Comprehensive Liberal Arts | FAOSTAT | CSV | Countries' data provided by FAO |
| | IHS | XLS | International historical data since 1846 |
| | NBS | CSV | National and provincial data in China |
| | DataBank | CSV | Historical data in the world |
| | GovNews | JSON | News published by Chinese government |
| | GeoNames | RDF-XML | World geographical KG |
| Language | HowNet | TXT | Chinese sememe knowledge resource |
| | NYTimes | JSON | New York Times news |
| | WordNet | NT | English electronic lexical database |
| | Framster | NQ | Integrated semantic knowledge base |
| | MW | TXT | Merriam-Webster english dictionary |
| All | CSKG | TSV | Commonsense KG |

We categorize the data above into three general types, which are: (1) Unstructured data, i.e., data without any pre-defined schema, including news, governmental policies, images, etc.; (2) Semi-structured data, i.e., data with limited or weak schema, such as JSON and XML files. In particular, we also regard tabular data (such as spatio-temporal geography data and statistics index) as semi-structured data, for they can only be modeled in relational databases but not in graph manners; (3) Structured data, i.e., data with specific schema or ontology, such as KGs with RDF format.

**Educational Knowledge Indexing.** As illustrated in the middle left part of Fig. 4, to use knowledge topics in EDUKG for effectively organizing gathered heterogeneous data, we propose our entity linking system, *EduLink*, and generalize it as a tool for heterogeneous data indexing. Following previous research [20], we separate entity linking to *mention detection*, *candidate generation*, and *entity disambiguation*. Meanwhile, we adjust several implementation details of each component for linking "Rhetorical Role" instances in EDUKG. We also adjust these components to support the linking for unstructured, semi-structured, and structured data. The evaluation of EduLink is presented in Sec. B of the appendix in supplementary materials. The detailed implementation of each module is shown as follows.

• **Mention Detection.** We leverage the aforementioned fine-tuned Chinese RoBERTa model for detecting entity mentions in plain texts. Meanwhile, since there are rhetorical roles defined as entities in EDUKG, we reuse the aforementioned template-based method in Sec. 3.2 for recognizing each sub-class of rhetorical roles from the given texts. Besides, only unstructured data need to be processed via the mention detection model because data with structural information are properly segmented.

• **Candidate Generation.** As for candidate generation, to support automatic fuzzy search of entities, we load the names of EDUKG knowledge topics into

*ElasticSearch*[5] along with its descriptions, and then use its searching function to generate entity candidates given the entity names detected in the step above.
• **Entity Disambiguation.** We consider the entity disambiguation task as a sentence similarity ranking task. We calculate the similarity score between its input context and its description stored in EDUKG by the Sentence-BERT model for each detected entity. For unstructured data indexing, the input context is its original sentence or image caption; for semi-structured tabular data, the context is other properties along with their column names within the same record; for structured data, the context is the entity's description in its original KG.

## 4   Evaluation of EDUKG

We illustrate the quality of EDUKG in this section by presenting specific task-related characteristics. Meanwhile, we compare EDUKG with existing educational KGs, including HEKG [34], KnowEdu [4], K12EduKG [3], CKGG [23], MEduKG [14], and MOOC-KG [6]. We present characteristics of EDUKG compared with existing educational KGs as shown in Table 3.

Table 3: Characteristics comparison between existing KGs, where Pred., Rhet., Mat., Exe., Ext., and Dom. are short for predicate, rhetorical role, teaching material, exercise, external data source, subject domain, respectively.

| KG | Class | Pred. | Entity | Triplet | Concept | Rhet. | Mat. | Exe. | Ext. | Dom. |
|---|---|---|---|---|---|---|---|---|---|---|
| HEKG | 6 | 7 | N/A | N/A | N/A | ✗ | 1.2k | ✗ | 4 | MOOC |
| KnowEdu | N/A | N/A | N/A | N/A | N/A | ✗ | ✗ | ✗ | ✗ | N/A |
| K12EduKG | N/A | N/A | N/A | N/A | N/A | ✗ | ✗ | ✗ | ✗ | Mathematics |
| CKGG | 754 | 389 | 412k | 1,500,000k | 412k | ✗ | ✗ | ✗ | 20+ | Geography |
| MEduKG | N/A | N/A | N/A | N/A | N/A | ✗ | ✗ | ✗ | ✗ | N/A |
| MOOC-KG | 4 | 14+ | 28k | N/A | N/A | ✗ | ✗ | ✗ | 4 | MOOC |
| Ours | 635 | 1759 | 252,328k | 3,860,446k | 36.79k | 18k | 3k | 6k | 32 | Interdisciplinary |

**Knowledge Sufficiency.** In Table 3, although CKGG has more classes and triplets than EDUKG, its original paper claimed that there are 655 classes and 353 properties in CKGG which are not populated, meaning that EDUKG still achieves better knowledge granularity. Besides, entities in CKGG are mostly locations, and triplets in CKGG are almost datatype properties of locations, lacking data variability. Thus, Compared with existing educational KGs, EDUKG achieves remarkably better knowledge sufficiency from the following aspects:
• **Interdisciplinary Knowledge Modeling.** EDUKG contains rich interdisciplinary knowledge topics and relations extracted from textbooks. On average, each extracted knowledge topic occurs in 1.70 different textbooks. For example, the concept *Leonardo Da Vinci*, who is a great artist, scientist, engineer, and theorist, is essential in subjects like mathematics, physics, and history. Furthermore, we show a case study in Fig. 6, where we illustrate the sufficiency and practicality of interdisciplinary knowledge toward education.
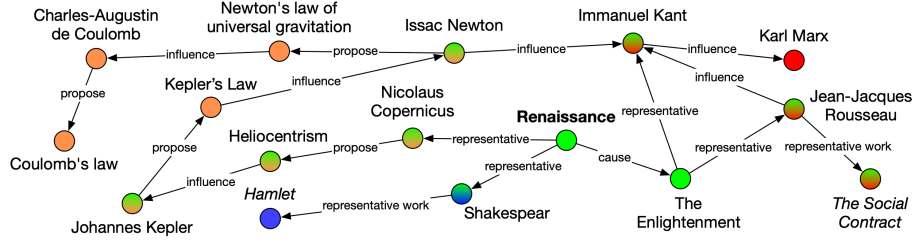
---

[5] https://www.elastic.co/

Fig. 6: A case study of EDUKG interdisciplinary concepts, where green, blue, orange, and red points refer to subjects of history, literature, physics, and politics, respectively. Points with multiple colors refer to interdisciplinary concepts.

• **Fine-grained Knowledge Topics.** We organize extracted knowledge concepts based on our proposed fine-grained sub-ontology under the "Knowledge Concept" class. The average instances' depth on the class hierarchy of "Knowledge Concept" sub-ontology is 3.87, while 48.39% of the instances of "Knowledge Concept" are classified under the leaf type of its class hierarchy, which addresses the fine granularity of extracted concepts. Also, as mentioned in Sec. 2.3, to our knowledge, we are the first to infuse rhetorical roles as entities in KG to improve the knowledge granularity. Each rhetorical role, on average, mentions 1.14 knowledge concepts in EDUKG, showing that infusing rhetorical roles in KG can densify our KG effectively. The distribution of each knowledge concept and each rhetorical role's in- and out-degree are shown in Fig. 7a, presenting the high density of EDUKG's knowledge topics.
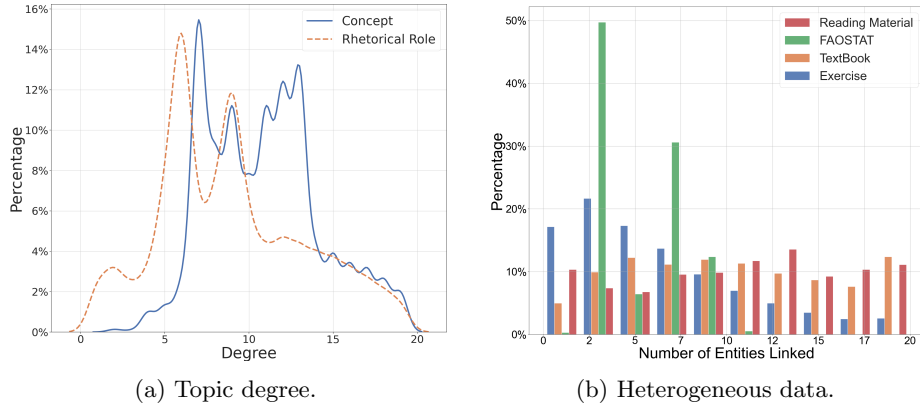


(a) Topic degree.
(b) Heterogeneous data.

Fig. 7: The detailed characteristics for knowledge and resources in EDUKG.

**Resource Richness.** As shown in Table 3, there are only a few existing educational KGs consisting of teaching materials or other educational-related resources. EDUKG includes 256 units, 779 lessons, and 2371 sections segmented from 46 textbooks of eight main subjects in Chinese high-school-level education. Moreover, each section has 22.17 related knowledge topics on average, indicat-

ing the high cohesion between educational resources and knowledge in EDUKG. Furthermore, there are 6518 exercises in EDUKG, which are further parsed into 10602 questions in total. Each exercise, on average, is linked with 1.66 knowledge topics, where 65.29% among them are not explicitly mentioned yet essential for the problems' solving.

**Data Variability.** Since EDUKG is a heterogeneous educational KG, data variability plays an essential role in EDUKG's data quality. As mentioned in Sec. 3.3, we leverage our entity linking system for indexing online heterogeneous data. In total, we convert data from 32 sources into more than 250 million entities, where on average, each entity corresponds to 4.80 knowledge topics in EDUKG. The detailed distribution of related knowledge topics for some sub-classes of "External Heterogeneous Data" is shown in Fig. 7b. In addition, we also infuse image data in EDUKG. On average, each textbook section is linked with 3.11 images, and each exercise is linked with 1.36 images, illustrating our data heterogeneity.

**Supported Tasks.** Prior KGs for education suffer from the lack of sufficient knowledge and heterogeneous resources and data. Thus, they could not fully support variable downstream educational tasks. However, EDUKG can be regarded not only as a KG but also as a data cube for numerous educational-related usages, including learning management system development, intelligent tutoring system research, educational data mining exploration, etc. We show a question answering demo platform in Sec. C of the appendix in supplementary materials.

## 5    Impact and Availability

In this section, we first present the availability of EDUKG, including its code, data, and applications. Afterward, we highlight the impact of EDUKG on research and society, along with its beneficial groups.

EDUKG is published under the terms of the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. EDUKG ontology is available under persistent URI in w3id[6], and the RDF dump of EDUKG is available at our Github homepage[7]. All data correspond to knowledge, resources, and heterogeneous data is free to download. The ontology and sampled data of EDUKG can be found in supplementary materials. Additionally, our toolkits for constructing and sustainably maintaining EDUKG are also open-sourced. Furthermore, we build an open platform[8] upon EDUKG, where users can efficiently browse massive K-12 educational knowledge and resources. We also provide vast kinds of open API, including knowledge searching, entity linking, question answering, etc., on our open platform for the usage of developers and researchers.

**Educational Research Development.** EDUKG can provide essential knowledge extracted from textbooks in a fine-grained manner, which is essential for researchers to explore cutting-edge educational technologies. Meanwhile, these data are also valid for developers to build downstream educational applications.

---

[6] https://w3id.org/edukg/ontology/
[7] https://github.com/THU-KEG/EDUKG
[8] http://open.edukg.cn/home

**Educational Equity Promotion.** Students nowadays from different regions still suffer from regional education inequality problem [7]. With the help of the EDUKG-based online learning platform, students can access comprehensive learning resources and conduct self-learning efficiently, which effectively ceases the education inequality issue.

**Beneficial Groups.** EDUKG aims to support the development of variable beneficial groups, including (1) students, parents, and teachers with limited learning resources, (2) educational application developers with insufficient data, and (3) researchers of educational-AI fields to access abundant knowledge and resources from our platform efficiently.

## 6   Related Work

KnowEdu [4] is a KG construction system extracting entities and mining prerequisite rules from K-12 teaching materials. Meanwhile, K12EduKG [3] leverages existing NLP techniques, such as NER, for extracting key educational concepts from K-12 mathematical textbooks. It also proposes an association rule mining algorithm to identify prerequisite rules between entities. Besides, several educational KGs are proposed based on MOOCs. For example, MOOC-KG [6] is constructed based on MOOC platforms, where numerous learning materials are integrated. Besides, HEKG [34] uses entity detection and relation extraction techniques based on a simple schema with the crawled information on MOOC websites. In addition, CKGG [23] is claimed as a KG for Chinese high-school-level geography education. However, it should be regarded as a KG for geography rather than educational usage because it integrates massive professional geography data beyond high-school-level education.

## 7   Conclusion and Future Work

This paper proposes EDUKG, a heterogeneous sustainable K-12 educational KG based on Chinese high-school education with more than 252 million entities and 3.86 billion triplets. To the best of our knowledge, EDUKG is the first large-scale interdisciplinary KG for K-12 education. Meanwhile, we design and propose a series of toolkits for sustainable data maintenance to dynamically collect and extract knowledge and resources from massive online data.

As for future work, there are several approaches for improving the data sufficiency of EDUKG. A substantial improvement is to add multilingual and multimodal data in EDUKG. Meanwhile, for adaptive learning tasks, the student behavioral data is essential for downstream tasks. We would further conduct user experimentation based on EDUKG for broader task support. Moreover, our future work will also focus on developing variable downstream applications. For instance, we are working on several applications for real high-school education scenarios based on educational data mining. EDUKG can also provide rich knowledge and data for question answering platform development. We hope that EDUKG can be beneficial to the development of educational technologies.

# References

1. Aliyu, I., Kana, A., Aliyu, S.: Development of knowledge graph for university courses management. International Journal of Education and Management Engineering **10**(2), 1 (2020)
2. Bailey, K.M.: Working for washback: A review of the washback concept in language testing. Language testing **13**(3), 257–279 (1996)
3. Chen, P., Lu, Y., Zheng, V.W., Chen, X., Li, X.: An automatic knowledge graph construction system for k-12 education. In: L@S. pp. 1–4 (2018)
4. Chen, P., Lu, Y., Zheng, V.W., Chen, X., Yang, B.: Knowedu: A system to construct knowledge graph for education. Ieee Access **6**, 31553–31563 (2018)
5. Cui, Y., Che, W., Liu, T., Qin, B., Yang, Z.: Pre-training with whole word masking for chinese bert. IEEE/ACM Transactions on Audio, Speech, and Language Processing **29**, 3504–3514 (2021)
6. Dang, F., Tang, J., Li, S.: Mooc-kg: a mooc knowledge graph for cross-platform online learning resources. In: ICEIEC. pp. 1–8 (2019)
7. Hannum, E., Meiyan, W.: Geography and educational inequality in china. China economic review **17**(3), 253–265 (2006)
8. Hou, L., Li, J., Wang, Z., Tang, J., Zhang, P., Yang, R., Zheng, Q.: Newsminer: Multifaceted news analysis for event search. Knowledge-Based Systems **76**, 17–29 (2015)
9. Hu, J., Li, Z., Xu, B.: An approach of ontology based knowledge base construction for chinese k12 education. In: ICMIP. pp. 83–88 (2016)
10. Ilkou, E., Abu-Rasheed, H., Tavakoli, M., Hakimov, S., Kismihók, G., Auer, S., Nejdl, W.: Educor: An educational and career-oriented recommendation ontology. In: ISWC. pp. 546–562 (2021)
11. Jin, H., Li, C., Zhang, J., Hou, L., Li, J., Zhang, P.: Xlore2: large-scale cross-lingual knowledge graph construction and application. Data Intelligence **1**(1), 77–98 (2019)
12. Jones, C.: Interdisciplinary approach-advantages, disadvantages, and the future benefits of interdisciplinary studies. Essai **7**(1), 26 (2010)
13. Kuiper, E., Volman, M., Terwel, J.: The web as an information resource in k–12 education: Strategies for supporting students in searching and processing information. Review of educational research **75**(3), 285–328 (2005)
14. Li, N., Shen, Q., Song, R., Chi, Y., Xu, H.: Medukg: A deep-learning-based approach for multi-modal educational knowledge graph construction. Information **13**(2), 91 (2022)
15. Li, S., Xu, B., Yang, Y.: Drte: A term extraction method for k12 education. Journal of Chinese Information Processing **32**(3), 101–109 (2018)
16. Malik, V., Sanjay, R., Guha, S.K., Nigam, S.K., Hazarika, A., Bhattacharya, A., Modi, A.: Semantic segmentation of legal documents via rhetorical roles. arXiv preprint arXiv:2112.01836 (2021)
17. Pellissier Tanon, T., Weikum, G., Suchanek, F.: Yago 4: A reason-able knowledge base. In: ESWC. pp. 583–596 (2020)
18. Peng, Y., Chen, P., Lu, Y., Meng, Q., Xu, Q., Yu, S.: A task-oriented dialogue system for moral education. In: AIED. pp. 392–397 (2019)
19. Piech, C., Bassen, J., Huang, J., Ganguli, S., Sahami, M., Guibas, L.J., Sohl-Dickstein, J.: Deep knowledge tracing. NEURIPS **28** (2015)
20. Ravi, M.P.K., Singh, K., Mulang, I.O., Shekarpour, S., Hoffart, J., Lehmann, J.: Cholan: A modular approach for neural entity linking on wikipedia and wikidata. In: EACL. p. 504–514 (2021)

21. Reimers, N., Gurevych, I.: Sentence-bert: Sentence embeddings using siamese bert-networks. In: EMNLP-IJCNLP. pp. 3982–3992 (2019)
22. Romero, C., Ventura, S.: Educational data mining and learning analytics: An updated survey. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery **10**(3), e1355 (2020)
23. Shen, Y., Chen, Z., Cheng, G., Qu, Y.: Ckgg: A chinese knowledge graph for high-school geography education and beyond. In: ISWC. pp. 429–445 (2021)
24. Shohamy, E.: The power of tests: The impact of language tests on teaching and learning. nflc occasional papers. (1993)
25. Tiberghien, A.: Learning and teaching: differentiation and relation. Research in Science Education **27**(3), 359–382 (1997)
26. Xu, L., Dong, Q., Liao, Y., Yu, C., Tian, Y., Liu, W., Li, L., Liu, C., Zhang, X., et al.: Cluener2020: fine-grained named entity recognition dataset and benchmark for chinese. arXiv preprint arXiv:2001.04351 (2020)
27. Xu, M.: Text2vec: Text to vector toolkit (2 2022), `https://github.com/shibing624/text2vec`
28. Xu, S., Yang, R., Hao, B.: China's National College Entrance Examination Report (2022). Xinhua Publishing House (2021)
29. Yang, Y., Xu, B., Hu, J., Tong, M., Zhang, P., Zheng, L.: Accurate and efficient method for constructing domain knowledge graph. Journal of Software **29**(10), 2931–2947 (2018)
30. Yang, Z., Wang, Y., Gan, J., Li, H., Lei, N.: Design and research of intelligent question-answering (q&a) system based on high school course knowledge graph. Mobile Networks and Applications **26**(5), 1884–1890 (2021)
31. Yu, J., Wang, C., Luo, G., Hou, L., Li, J., Tang, J., Liu, Z.: Course concept expansion in moocs with external knowledge and interactive game. In: ACL. pp. 4292–4302 (2019)
32. Yu, J., Wang, Y., Zhong, Q., Luo, G., Mao, Y., Sun, K., Feng, W., Xu, W., Cao, S., Zeng, K., et al.: Mooccubex: A large knowledge-centered repository for adaptive learning in moocs. In: CIKM. pp. 4643–4652 (2021)
33. Zhang, J., Cao, Y., Hou, L., Li, J., Zheng, H.T.: Xlink: An unsupervised bilingual entity linking system. In: CCL. pp. 172–183 (2017)
34. Zheng, Y., Liu, R., Hou, J.: The construction of high educational knowledge graph based on mooc. In: 2017 IEEE 2nd Information Technology, Networking, Electronic and Automation Control Conference (ITNEC). pp. 260–263 (2017)