

Semantic-Awareness Recommendation with Linked Open Data in Web-Based Investigative Learning

Kang Ting and Shinobu Hasegawa

Abstract—Web-based investigative learning provides a platform for learners to create their own learning scenarios by organizing knowledge over the web in a self-directed way. This kind of knowledge management activity helps learners to achieve a proper cognitive load on the investigation. However, it is difficult for learners to discover related concepts among a vast number of unstructured web resources concurrently with a better knowledge construction process. Therefore, this research aims to propose a method to recommend semantic-related concepts with Linked Open Data for learners during the investigation of the web-based investigative learning process. We proposed a Semantic-awareness Recommendation System that extracts the semantic related concepts from DBpedia by sending the regulated SPARQL query. In this work, generating a regulated concept map based on the initial question for the recommendation, three significant elements would be considered: Semantic relations, Concept Importance Estimation and Filtering.

Index Terms—Web, investigative learning, recommendation, linked open data, self-directed learning.

I. INTRODUCTION

With the rapid development of information and communication technology in the 21st century, human life has also undergone tremendous improvements. In the field of education, we have considerable expectations and emphasis on the development of the web resources such as Linked Open Data (LOD) [1] which is combined a blend of Linked Data and Open Data. LOD breaks down barriers between different data formats and sources. Web-based investigative learning [2] is one of the learning approaches benefited from. It allows learners to investigate any topics to learn in a self-directed way.

The web-based investigative learning model included three processes [3]: Search for web resources, Navigational learning, and Question decomposition. Learners need to select suitable and reliable resources against an initial keyword for knowledge construction from a vast number of web resources by themselves. This learning means searching the meaning of the initial keyword and exhaustively investigating many concepts related to the initial question and construct broader and deeper knowledge. By repeating these processes cyclically, learners are expected to create a learning scenario that means turning those unstructured web resources into structured resources to make their knowledge construction process strengthen. However, it is difficult for

learners to discover related concepts among a vast number of unstructured web resources concurrently with a better knowledge construction process.

Therefore, this thesis aims to propose a method to recommend semantic-related concepts with LOD for learners during the investigation of the web-based investigative learning process. LOD is a set of structured data interlinking with related ones on the Web. In this work, we use DBpedia [4].

II. PROBLEM STATEMENT AND RELATED WORKS

This work is inspired by the previous research of web-based investigative learning providing high-quality recommendation and awareness of the relevance between concepts for learners to strengthen their knowledge construction process. Furthermore, this work also inspired by several research which providing recommendation list by means of LOD and semantic relations over the web.

A. Adaptive Recommendation for Question Decomposition in Web-Based Investigative Learning

According to previous works of Web-based investigative learning, Hagiwara [2] pointed out that learners often suffer from question decomposition during Web-based investigative learning. It is difficult for learners to make a sufficient investigation in concurrence with navigation and knowledge construction. Therefore, an adaptive recommending strategy for providing a related sub-question keyword against an initial Q-keyword by extracting the data from DBpedia was proposed. Owing to the finding of this work, the adaptive recommendation makes an effort to help learners who have difficulty constructing a learning scenario, sufficiency decomposing the learning scenario, and observing the relation between Q-keyword during web-based investigative learning. However, when we focus on learners' self-initiative, the recommendations should follow their learning process from the decomposition and the comprehensive concepts in which learners are newly interested.

B. Relevance between Q-Keywords Corresponding to Transition of Interest in Web-Based Investigative Learning

Regarding the transition of interest in web-based investigative learning, Yamauchi [5] pointed out that we should focus on the initial Q-keyword and those concepts learners are newly interested in. He defined three parameters to calculate the relevance between two questions by LOD as follows:

Manuscript received September 23, 2021; revised December 17, 2021.

Kang Ting and Shinobu Hasegawa are with Japan Advanced Institute of Science and Technology, Japan (e-mail: s1910155@jaist.ac.jp, hasegawa@jaist.ac.jp).

- Question distance: The number of nodes that appear in the shortest path to connect two questions on DBpedia.
- Question similarity: Simpson's coefficient between two sets consisted of related words of each question.
- Question coupling: The number of found elements connecting with question keyword in both directions on DBpedia.

By means of DBpedia, the relevance between a pair of Q-keywords could be calculated. His work breaks down the barriers of evaluating the relevance between different learning scenarios in Web-based investigative learning. His work does provide a good chance for learners to recognize the relation between learning scenarios partially by exploiting LOD. However, the capabilities of LOD were underutilized. It is not representative enough to express the relation between concepts comprehensively. Therefore, by exploiting the capabilities of LOD, we could express the relations between concepts comprehensively over the semantic web.

C. Research Related to Semantic-Awareness Recommendation

Several research projects [6], [7] focus on the semantic path-based ranking using LOD such as DBpedia to generate a ranked recommendation list and tuning the weights of features gathered from DBpedia to increase recommendation accuracy. However, according to the criteria of web-based investigative learning, every learner will create specific learning scenarios in a self-directed way against different Q-keywords.

In order to tackle those issues, we proposed a method to generate the regulated concept map against a selected keyword for recommending the relevant concepts at different levels without preventing learners from their self-directed investigation with LOD and Semantic relations between concepts. Therefore, for our proposed method, those recommended concepts related to the initial Q-keyword will be defined by three elements which included **Semantic Relations, Concept Importance Estimation, and Filtering.**

III. PRELIMINARIES

This section will introduce significant background knowledge for this research.

A. Connectivism

According to the definition of connectivism [8], learning is no longer just a process of personal acquisition of materialized knowledge, but a process of establishing connections to build an individual's internal cognitive network and external social network. Web-based investigative learning [3] provides a platform for learners to investigate the knowledge that resides in the Web to be connected. This kind of knowledge management activity helps learners to address those issues of organizational knowledge and transference. It could be seen from previous work [3] that web-based investigative learning improves the efficiency of the knowledge construction process for learners.

B. Web-Based Investigative Learning

Since the proposed recommending approach is for Web-based investigative learning, the basis of such learning and the cognitive tool named interactive Learning Scenario Builder (iLSB) [9] conducted for promoting question decomposition will be described.

In the previous work of Web-based investigative learning [3], this learning model included three stages (Fig. 1): Searching for web resources/pages, Navigational learning, and Question decomposition. Firstly, learners would search for web resources with a keyword representing an initial keyword also named Q-keyword. This stage aims to find out appropriate web resources for question investigation. Then, learners could navigate those resources selected in the previous stage for the knowledge construction process during the navigational learning stage. Meanwhile, they would also extract keywords from navigated resources. Finally, learners could build their own learning scenario by reviewing the knowledge constructed in the navigational learning stage to decompose the Q-keyword into sub-question to be further investigated.

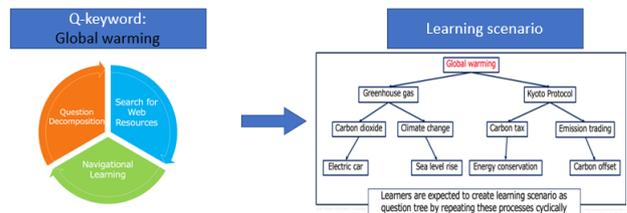


Fig. 1. The model of Web-based investigative learning [3].

The main feature of web-based investigative learning is to turn those unstructured web resources into structured resources by learners in their self-directed way. By comparing web resources with traditional text resources, text resources are well structured and provide learning scenarios that imply the questions to be investigated and their sequence, such as the table of contents. On the other hand, Web resources are unstructured and do not provide learning scenarios in advance. Therefore, learners need to decompose questions into related ones as sub-questions while constructing their knowledge. It implies that learners are expected to investigate questions in a self-directed way. Meanwhile, learners should create their learning scenarios and construct their knowledge concurrently. As a result, learners would have a high cognitive load on the investigation.

However, it is difficult for learners to discover concepts related to existing learning scenarios and estimate the relevance between a bunch of related concepts. If learners create a new learning scenario with weak relevance between previous scenarios they created, it is difficult for them to strengthen the knowledge construction process. Therefore, the necessity of recommendation should be regarded.

C. Interactive Learning Scenario Builder (iLSB)

In order to scaffold learners' investigative learning process as modeled, a cognitive tool named interactive Learning Scenario Builder (iLSB) (Fig. 2) [9] has been developed as an add-on for Firefox. iLSB provides scaffolding functions such as Searching engine for gathering learning resources,

Keyword repository for constructing their knowledge, and Question tree viewer for creating their learning scenario. Owing to previous work findings [3], it has ascertained that iLSB could promote question decomposition in Web-based investigative learning.

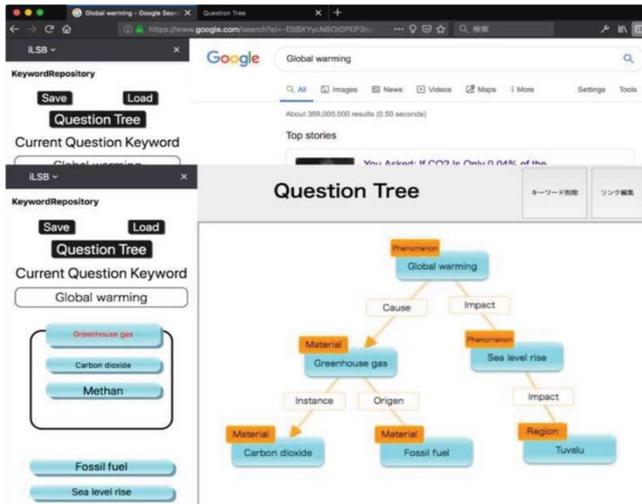


Fig. 2. The user interface of iLSB [9].

D. Linked Open Data

In 1989, Tim Berners-Lee invented the first proposal [1] for the World Wide Web. The proposal outlined the principal concepts, and it defined important terms behind the Web, such as describing the Internet as a system of an interlinked hypertext document. He founded the World Wide Web Consortium(W3C) to maintain the development of those open standards to ensure the long-term growth of the Web. In addition to the classic “Web of documents”, W3C also built a technology stack to support a “Web of data” named linked data. The ultimate goal of linked data is to enable computers to do more useful work and develop systems that can support trusted interactions. The term “Semantic Web” refers to W3C’s vision of the Web of linked data. Semantic Web technologies enable people to create data stores on the Web, build vocabularies, and write rules for handling data. Linked data are empowered by technologies such as RDF, SPARQL, OWL, and SKOS. In 2006, Berners-Lee released the principles of linked data [10]:

- Use URIs as names for things.
- Use HTTP URIs so that people can look up those names.
- When someone looks up a URI, provide useful information.
- Include links to other URIs. so that they can discover more things.

Therefore, LOD is a set of structured data interlinking with related ones on the Web that is linked and uses open resources. One remarkable example of a LOD set is DBpedia which extracts structured information from Wikipedia and makes it available on the Web. In this work, we will retrieve relevant concepts from DBpedia for the recommendation.

E. DBpedia

Since its establishment in 2007, the DBpedia project [4] has been sustainably releasing large and open data sets, which are extracted from Wikimedia projects(such as

Wikipedia and Wikidata [11]). The data has been extracted using a sophisticated software called DBpedia Information Extraction Framework (DIEF) and represented by using the Resource Description Framework (RDF) [12]. In the last few years, the system has received many extensions and fixes from the community, which leads to the creation of a stable release version. Furthermore, by the effort of the W3C Semantic Web Education and Outreach (SWEO) interest group, DBpedia interlinked to a lot of massive Linked Open Data sets.

The English version of DBpedia contains 6.0 million entities, of which 4.6 million have abstracts [13]. It means DBpedia has a huge range of subject coverage. Moreover, DBpedia also consists of 5.0 billion pieces of information (RDF triples) [13] extracted from the English edition of Wikipedia. Meanwhile, an information extraction framework that included extraction, clustering, uncertainty management, and query handling was developed by the DBpedia community [4]. Fig. 3 shows the overview of DBpedia components. It means that it is convenient for us to query those structured data represented by the Resource Description Framework, especially the relationships and properties of information. All in all, for retrieving concepts to provide a recommendation, DBpedia is a reliable data-set for this work.

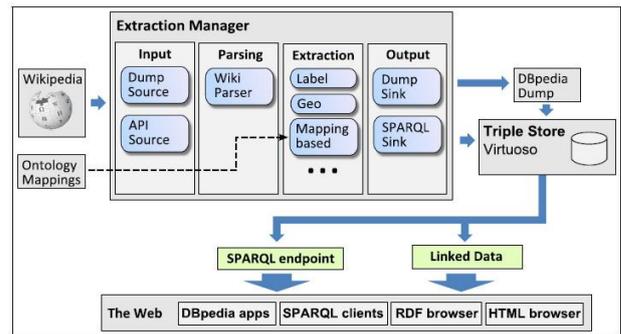


Fig. 3. The overview of DBpedia components [4].

F. Resource Description Framework

Resource description framework (RDF) [12] was conducted by the RDF Core Working Group under W3C. RDF data represent a data model for the information over the web as well as DBpedia. The model of RDF data expresses information in a triple, which included three elements such as subject, predicate, and object. Fig. 4 shows an instance on DBpedia that the relationship between subject and object is described by predicate which is directional and represented by a property.

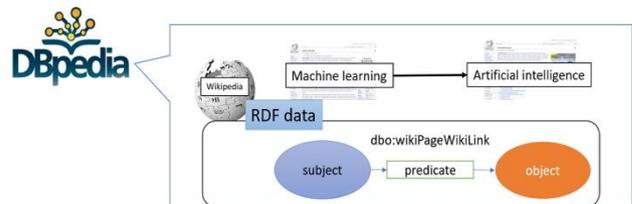


Fig. 4. An instance of RDF data on DBpedia.

We could immediately realize that a collection of triples can be represented as a graph data model named RDF graph [12] labelled and directed. This kind of structured data

benefited the construction of content that we are going to recommend in our proposed method.

G. SPARQL Query

For querying the RDF data, SPARQL Protocol and RDF Query Language (SPARQL) was developed by the W3C RDF Data Access Working Group (DAWG) [14]. It is a standard query language and protocol for RDF graph data. SPARQL query is used for querying required and optional RDF graph patterns with specifying conjunctions and disjunctions. Generally, in the SPARQL query, whether subject, predicate, or object could be the target variable of the RDF graph data. That is to say, sending a SPARQL query is a process to search the RDF graph data, which matches with required graph patterns. According to DAWG [14], the SPARQL query form included **SELECT**, **CONSTRUCT**, **DESCRIBE** and **ASK**. By combining the modifiers such as **LIMIT**, **ORDER BY**, **FILTER**, and so on, we could easily query all of the required graph patterns as we need against LOD. By sending the SPARQL query to Public DBpedia SPARQL endpoint, We could extract all of the RDF graphs in DBpedia.

H. Simple Knowledge Organization System

Simple knowledge Organization System (SKOS) [15] is a W3C recommendation document that defined a standard data model for sharing and linking knowledge organization systems via the semantic web. The principal element categories of SKOS are concepts, labels, notations, documentation, semantic relations, mapping properties, and collections. It is useful for us to describe the relationship between concepts, such as in a semantic-awareness way. The associated elements are listed in the Fig. 5.

SKOS Vocabulary Organized by Theme					
Concepts	Labels & Notation	Documentation	Semantic Relations	Mapping Properties	Collections
Concept	pref_Label	note	broader	broadMatch	Collection
ConceptScheme	alt_Label	changeNote	narrower	narrowMatch	orderedCollection
inScheme	hidden_Label	definition	related	relatedMatch	member
hasTopConcept	notation	editorialNote	broaderTransitive	closeMatch	memberList
topConceptOf		example	narrowerTransitive	exactMatch	
		historyNote	semanticRelation	mappingRelation	
		scopeNote			

Fig. 5. Properties of Simple Knowledge Organization System [15].

I. PageRank Algorithm

PageRank algorithm is a major algorithm that Google uses to evaluate the relevance or importance of a web page. Two different versions of the PageRank algorithm were published by Lawrence Page and Sergey Brin in several publications [16], [17].

Quoting the description of the Original PageRank algorithm published by Page and Brin [16] is given by:

$$PR(A) = (1 - d) + d \left(\frac{PR(T_1)}{C(T_1)} + \dots + \frac{PR(T_n)}{C(T_n)} \right) \quad (1)$$

where:

- PR(A) is the Original Google PageRank of page A.
- PR(T_n) is the Original Google PageRank of page T_n which is linked to page A.
- C(T_n) is a number of the outbound links on page T_n.
- d is a damping factor that could be set between 0 and 1.

- n is the total number of pages linked to page A.

In the Original Google's PageRank algorithm, the importance of a page T is constantly weighted by the number of its outbound links C(T). It means that the more outbound links a page T have, the less page A would be benefited by a link from page T. The PageRank value of page A would be the sum of the inbound links multiplied by a damping factor d is generally set to 0.85 [18].

For the second version, the PageRank value of page A is as follows [17]:

$$PR(A) = \frac{(1-d)}{N} + d \left(\frac{PR(T_1)}{C(T_1)} + \dots + \frac{PR(T_n)}{C(T_n)} \right) \quad (2)$$

Obviously, these two versions of the PageRank algorithm have no fundamental difference between each other. However, in the second version, it adapts (1 - d)/N instead of (1 - d) where N is the sum of all web pages. It means the probability of a random user surfer a web page is weighted by the total number of web pages. It forms a probability distribution over web pages, and the sum of PageRank value of all pages on the web would be 1.

IV. SEMANTIC-AWARENESS RECOMMENDATION SYSTEM

A. System Design

The design of the Semantic-awareness recommendation system is illustrated in Fig. 6. Firstly, the system requests learners to input an initial question for extracting relevant concepts from DBpedia. During the Regulated Concept Map Generation process, the initial question could be the concepts in the learning scenarios used to create or other keywords in which learners are newly interested in. By means of a regulated SPARQL query strategies, we can retrieve related concepts at different levels. Secondly, for Concept Importance Estimation, we employ the PageRank algorithm to estimate the importance of those concepts we retrieved in the previous section. The PAGERANK algorithm would calculate the importance of nodes in generated Regulated Concept Map. Finally, it is significant for us to define the filtering condition before updating the recommendation list to the learner. The filtering strategy is based on the content containment of the concept, which means that if there is no definition on DBpedia for the concept, we have to filter it. As a result, the proposed system updates the recommendation list for learners to continue the navigational learning process.

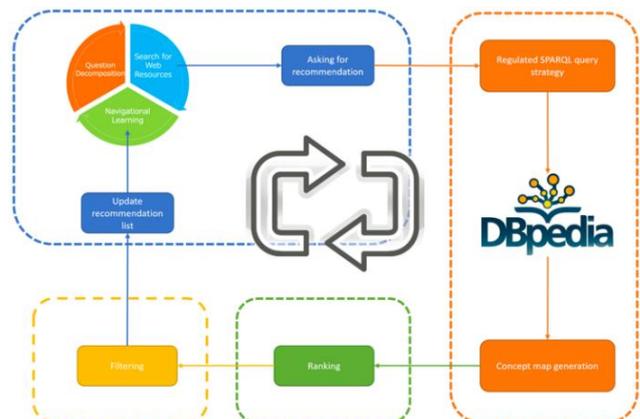


Fig. 6. Overview of Semantic-awareness recommendation system.

B. Regulated Concept Map Generation

DBpedia is a linked open data project which extracts structured content from Wikipedia [4]. Those structured content represented as RDF graph allowed the user to query the relationships and properties of Wikipedia resources semantically. We can either download the entire data set or access it by the public SPARQL endpoint. In this work, accessing DBpedia by public SPARQL endpoint is preferred since the dataset of DBpedia will be updated in the future.

As we mentioned previously, The important elements **Concepts** and **Semantic Relations** of the SKOS were employed in this work. SKOS concept is defined as RDF resources, and SKOS semantic relations are designed to declare the relationship between concepts within the scheme. The associated elements were employed in this work are listed in the Fig. 7.

SKOS Vocabulary Organized by Theme					
Concepts	Labels & Notation	Documentation	Semantic Relations	Mapping Properties	Collections
Concept	prefLabel	note	broader	broadMatch	Collection
ConceptScheme	allLabel	changeNote	narrower	narrowMatch	orderedCollection
inScheme	hiddenLabel	definition	related	relatedMatch	member
hasTopConcept	notation	editorialNote	broaderTransitive	closeMatch	memberList
topConceptOf		example	narrowerTransitive	exactMatch	
		historyNote	semanticRelation	mappingRelation	
		scopeNote			

Fig. 7. Properties of SKOS were employed in this work [15].

Extracting semantic related concepts from DBpedia, we firstly extract SKOS Concepts (RDF graph) from DBpedia using SPARQL query. Then, related concepts with semantic relations (Broader-Narrower) would be returned. The essential property: **SKOS: broader** would be used. This property represents a hierarchical relation between concepts. For example, *A SKOS: broader B* means B is broader and has more general meaning than A. Narrower follows in the same pattern. Take an initial Q-keyword **Machine learning** as an instance (Fig. 8). If we want to find out those concepts have broader relation with the initial Q-keyword, we first send the SPARQL query to extract keywords with semantic relations.

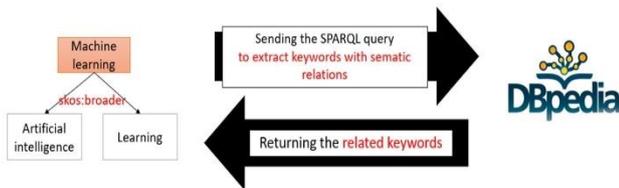


Fig. 8. Extracting SKOS Concepts from DBpedia using the SPARQL query.

Public SPARQL endpoint could return the query results in different formats such as JSON, CSV, and N-Triples. In this work, we select N-Triples for the result format since it is convenient for us to convert undirected graphs to a directed graph with two directed edges for each undirected edge. N-Triples is a line-based, plain text format for encoding an RDF graph. It is not only for the calculation of PageRank algorithm but also for the **A is SKOS:broader of B** relation between concepts narrower of **A SKOS:broader B**.

It is important for us to regulate the SPARQL query strategies if we aim to recommend related concepts at different levels without preventing learners from their self-directed investigation. The regulated concepts map is a

collection of entities called nodes, which are concepts that we are going to recommend to learners. Concepts are linked by edges with the properties **SKOS:broader** and **is SKOS:broader of**. The queries asking for broader nodes and narrower nodes used in this research are as follows (Fig. 9 and Fig. 10).

```
sparql.setQuery(f"""
CONSTRUCT {{ ?child skos:broader {concept} .}}
WHERE {{ ?child skos:broader {concept} .}}
""")
```

Fig. 9. The query for extracting broader concepts of the initial Q-keyword.

```
sparql.setQuery(f"""
CONSTRUCT {{ {concept} skos:broader ?parent .}}
WHERE {{ {concept} skos:broader ?parent .}}
""")
```

Fig. 10. The query for extracting narrower concepts of the initial Q-keyword.

By combining the queries above, Fig. 11 shows all of the semantic related concepts at different levels on DBpedia could be extracted. We need to pay attention to the definition of the PageRank algorithm. If the node has no outbound link, its importance would be 0. Therefore, it is necessary for us to adjust the range of the SPARQL query for the regulated concept map. Fig. 12 shows an instance that when we focus on the importance of Parent Nodes and Sibling Nodes from Parents Node, the range of SPARQL query should be more in-depth.

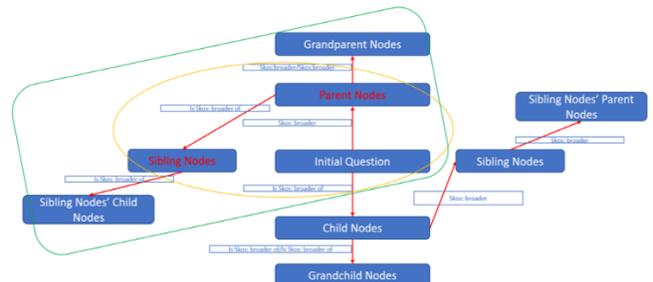


Fig. 11. Overview of the regulated concept map generation.

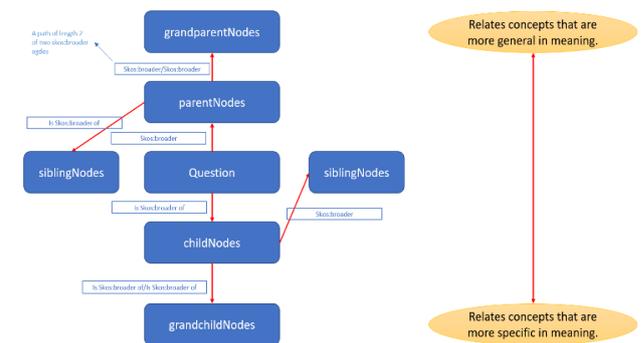


Fig. 12. An instance that how we arrange the range of the SPARQL query.

C. Concept Importance Estimation (Semantic Aware PageRank)

The original PageRank algorithm was introduced previously. In the original paper that Lawrence Page and Sergey Brin published, they consider the PageRank algorithm as a model of user behaviour who randomly suffers a web page with a certain probability, and the probability is given by the links on that web page. Moreover, due to the

previous work finding [19], the PageRank algorithm is also generally used as an index to decide the importance of nodes in a directed graph such as the RDF graph. Therefore, the PageRank algorithm is suitable for the concept importance estimation of this work, and we named it as Semantic-aware PageRank. In this research, a regulated concept map is a set of interlinked nodes, and we defined:

- The number of all nodes in Regulated Concept Map as $|R|$.
- A set of nodes x with the links $\{x, r\} \in E$ where $r \in R$ as B_r .
- The number of links from node x as C_x .

Eventually, we can calculate the PageRank value for all nodes in Regulated Concept Map PR_r based on the equation below:

$$PR_r = \frac{1-d}{N} + d \sum_{x \in B_r} \frac{PR_x}{C_x} \quad (3)$$

where d is a damping factor, which is set as 0.85 [18].

We assume that the importance of a concept node is determined by the number of outbound links on that concept. The probability of random surfer a node is weighted by the total number of nodes in the Regulated Concept Map. Equation (3) forms a probability distribution only over all the nodes in the Regulated Concept Map, and the sum of them would be 1.

In fact, there existed the calculated PageRank value for DBpedia. There are the reasons that we did not generally use the calculated PageRank value on DBpedia for concept importance estimation. The first reason is that it is significant for us to regard learners' knowledge construction process during Web-based investigative learning. Concepts were recommended to guide learners to navigate related concepts with strong relevance between initial Q-keyword. The second reason is that it is essential for us to regard the content containment of those recommended concepts. Concepts such as Time period have certain importance in general cases. However, for the knowledge construction process, the utility of the concept is not important for the recommendation. These reasons concluded the importance of the directional SPARQL query strategies.

D. Filtering

Before updating the recommendation list to learners, we have to filter those concepts which are not important. In this work, we would filter concepts based on the concept's utility. Since not every concept has a definition on DBpedia, the hypothesis that the concept has no definition on DBpedia is not significant for the recommendation. We would explain it through a practical case.

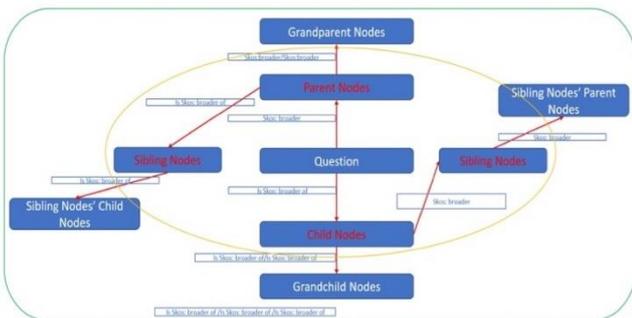


Fig. 13. The range and nodes will be applied for the practical case.

Setting **Natural Language Processing** as an initial Q-keyword, Fig. 13 shows the range of the directional SPARQL query and nodes that its concept importance would be calculated.

By means of the proposed Concept importance estimation approach, we got the ranking list below (Table I):

TABLE I: TOP 10 CONCEPTS RELATED TO NATURAL LANGUAGE PROCESSING SORT BY SEMANTIC-AWARE PAGERANK

From	Concept	PageRank value
Parent	Artificial intelligence applications	0.04195
Parent	Natural language and computing	0.04156
Siblingparent	Character encoding	0.03461
Siblingparent	Computing by natural language	0.03171
Siblingparent	Computational linguistics	0.03020
Siblingparent	Language software	0.02873
Grandparent	Linguistics	0.02483
Siblingparent	Internationalization and localization	0.02327
Child	Corpus linguistics	0.02227
Siblingparent	Language-specific Linux distributions	0.02150

By sending the SPARQL query to the public SPARQL endpoint as follow (Fig. 14), we could easily discover that a concept of **Language-specific Linux distributions** has no definition on DBpedia, and it would be filtered.

```
SELECT DISTINCT ?definition
WHERE{
  (<http://dbpedia.org/resource/concept> db:abstract ?definition
  FILTER (LANG(?definition) = 'en')
}
```

Fig. 14. The SPARQL query for filtering.

V. EVALUATION

In this section, the ranking results would be analyzed by Spearman's correlation coefficient. Spearman's correlation coefficient measures the strength and direction of the association between two ranked variables. Furthermore, a case study would be conducted in order to test the hypothesis that using the Semantic-awareness recommendation with linked open data could help learners strengthen the knowledge construction process by discovering semantic related concepts during Web-based investigative learning.

A. Comparison between Semantic-Aware PageRank, DBpagerank, and User expectation.

This experiment would be conducted to measure the strength of the association between Semantic-aware PageRank, DBpagerank and User expectation against the concepts extracted by the regulated concept map. DBpagerank [4] is the PageRank value for all the resources in DBpedia calculated by the DBpedia community. If the values are same, the rank of the average value is returned. For the ranking list arranged by Professor, we consider it as the user expectation. In total, three pairs of variables and the direction of the relationship would be analyzed. We would employ two initial Q-keywords, **Machine learning** and **Smoking**, and two regulated query strategies (More general and More specific) would be applied. Three pairs of variables would be analyzed:

Pair A

- Ranking list sort by Semantic-aware PageRank.
- Ranking list sort by DBpagerank.

Pair B

- Ranking list sort by DBpagerank.
- Ranking list arranged by Professor.

Pair C

- Ranking list sort by Semantic-aware PageRank.
- Ranking list arranged by Professor.

Spearman's correlation coefficient [20] is a statistical measure of the strength of a *monotonic* relationship between paired data. In a population, it is denoted by r_s and is by design constrained as follows:

$$-1 \leq r_s \leq 1 \quad (4)$$

According to the definition, the closer r_s is to ± 1 , the stronger the monotonic relationship. Since the correlation is effect size, we could verbally describe the strength of the correlation using the following guide for the absolute value of r_s [20]:

- 0.19: Very weak
- 0.20 - 0.39: Weak
- 0.40 - 0.59: Moderate
- 0.60 - 0.79: Strong
- 0.80 - 1.0: Very strong

For determining the significance of this test, we have to test the null hypothesis H_0 where there is no monotonic correlation to the population against the alternative hypothesis H_1 , where there is a monotonic correlation. Let ρ_s as the Spearman's population correlation coefficient then we can thus express this test as follow:

$$\begin{aligned} H_0: \rho_s &= 0 \\ H_1: \rho_s &\neq 0 \\ \alpha &= 0.05 \end{aligned}$$

Table II shows the analyzing results of the Spearman's Correlation Coefficient.

TABLE II: THE ANALYZING RESULTS OF THE SPEARMAN'S CORRELATION COEFFICIENT

Machine Learning (More general)	Pair A	Pair B	Pair C
Coefficient	0.0587	0.4304	0.5714
N	20	20	20
T Statistic	0.2494	2.0230	2.9542
Degree of Freedom	18	18	18
P-value	0.8058	0.0582	0.0085*

Machine Learning (More specific)	Pair A	Pair B	Pair C
Coefficient	-0.3560	-0.0607	-0.2497
N	20	20	20
T Statistic	1.6161	0.2581	1.0941
Degree of Freedom	18	18	18
P-value	0.1235	0.7993	0.2883

Smoking (More general)	Pair A	Pair B	Pair C
Coefficient	-0.0942	-0.1325	0.4938
N	20	20	20
T Statistic	0.4016	0.5673	2.4092
Degree of Freedom	18	18	18
P-value	0.6927	0.5775	0.0269*

Smoking (More specific)	Pair A	Pair B	Pair C
Coefficient	0.1743	0.0682	0.4573
N	20	20	20
T Statistic	0.7510	0.2900	2.1817
Degree of Freedom	18	18	18
P-value	0.4623	0.7752	0.0426*

*P<0.05 Two tailed

By the observation of the analyzing results above, for the recommendation list of **Machine learning (More general)** ($r_s = 0.5714, n = 20, P < 0.05$), **Smoking (More general)** ($r_s = 0.4938, n = 20, P < 0.05$) and **Smoking (More specific)** ($r_s = 0.4573, n = 20, P < 0.05$), Pair C maintained the highest value of the Spearman's Correlation Coefficient. It shows that, in most cases, there is a **moderate, positive monotonic** correlation between Ranking lists sort by Semantic-aware PageRank value and ranking lists arranged by Professor (User expectation).

B. Case Study

The criteria of this case study are as follows: By applying two directional SPARQL query strategies (Fig. 15) to 4 question keywords which are **Machine learning, Nuclear power, Governance, and Smoking**, we test the efficiency of the proposed recommendation system in Web-based investigative learning. By asking participants to pick out those concepts that are more general or more specific in meaning to the initial question keyword in the limited time. Meanwhile, the link to the concepts' definition page on DBpedia will be provided as a reference during the tasks. It simulates the navigational learning of Web-based investigative learning that learners could navigate those resources recommended by the proposed approach for knowledge construction process during the navigational learning stage.

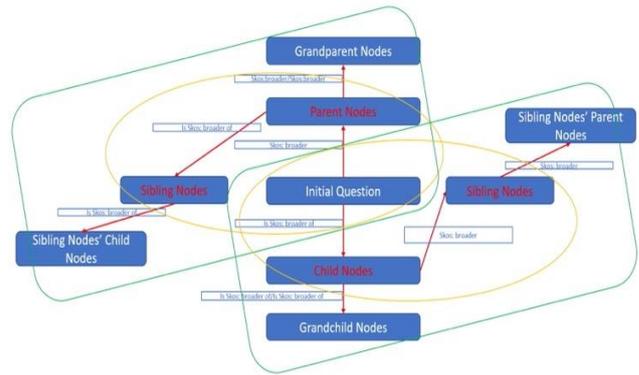


Fig. 15. Two directional query strategies for the case study.

There are 4 sections in total. For each section, participants would process 3 tasks as follows:

- Picking out those concepts that are more general or more specific in meaning to the initial question keyword.
- Writing down the three most important concepts that they selected in the previous task.
- Post-test questionnaire.

1) Analyzing results of first task

For measuring the difference in the number of concepts that participants pick out during the limited time between two groups, the Mann-Whitney U Test [21] is employed. The purpose of this non-parametric measurement is to compare the difference between the two populations. The basis on which we make inferences is also based on the sampling distribution composed of all the possible sample characteristics. Therefore, we test the hypothesis below:

- H_0 : There is no difference in the number of concepts selected between the two groups.

- H_1 : There is a difference in the number of concepts selected between the two groups.
- $\alpha = 0.05$
- Sampling distribution $Z_{critical} = \pm 1.96$

Instead of calculating the difference in the average, we calculate the verification statistical value U which is based on the grade of the variable score in the sample. For calculating the U value, we first merge all the observations from both groups to one set the two populations which are the number of concepts that participants pick out during the limited time, and then give the grades according to the value of the variable items. The higher value would maintain a higher grade, and then they would sort by order (from high to low). Then add up the grades assigned to each sample. Finally, compare the difference in the sum of levels between the two populations. Here we have:

$$U_1 = N_1N_2 + \frac{N_1(N_1+1)}{2} - \sum R_1 \quad (5)$$

where:

- N_1 and N_2 are the sample size of each group
- $\sum R_1$ is the sum of the ranks in the controlled group which is equal to 1305.

An equally valid formular U_2 is as follow:

$$U_2 = N_1N_2 + \frac{N_2(N_2+1)}{2} - \sum R_2 \quad (6)$$

where:

- $\sum R_2$ is the sum of the ranks in experimental group which is equal to 1935.

and

$$U = \min(U_1, U_2) = 485$$

For large samples, U is approximately normally distributed. Therefore, the value of the standardized value $Z_{statistic}$ would be:

$$Z_{statistic} = \frac{U - \mu_u}{\sigma_u} \quad (7)$$

where μ_u and σ_u are the mean and standard deviation. μ_u and σ_u are given by:

$$\mu_u = \frac{N_1N_2}{2} \quad (8)$$

$$\sigma_u = \sqrt{\frac{N_1N_2(N_1+N_2+1)}{12}} \quad (9)$$

As a result,

$$Z_{statistic} = \frac{485 - 800}{73.94} \approx -4.26 < -1.960$$

Since $Z_{statistic} < -1.960$ with $\alpha = 0.05$ (Two tailed), we have to reject H_0 , and state that there is a difference in a number of concepts selected between the controlled group and experimental group.

Skewness is a measure of symmetry, and Kurtosis describes the tail shape of the data's distribution. By observing Descriptive statistics (Table III) for two groups in this case study, the data of the experimental group forms a negative skewness. The data distribution of the experimental group is left-skewed which means during the task of the experimental group. Learners tend to select more concepts than the controlled group. Moreover, the data distribution of the experimental forms a positive kurtosis which indicates a

fat-tailed distribution. It refers to an increase in the probability of concepts being selected for an extreme number in the experimental group.

TABLE III: DESCRIPTIVE STATISTICS FOR TWO GROUPS IN THE FIRST TASK

	Concepts selected in Controlled group	Concepts selected in Experimental group
Mean	7.025	8.25
Standard Error	0.4244	0.3257
Median	7	8
Mode	4	8
Standard Deviation	2.6841	2.0600
Sample Variance	7.2045	4.2436
Kurtosis	-0.4968	1.7845
Skewness	0.3070	-0.9471
Range	11	10
Minimum	3	3
Maximum	14	13
Sum	281	330
Count	40	40

2) Analyzing results of second task

During the second task, participants were asked to write down the three most important concepts they selected in the previous task. In the previous task, there are 20 concepts in the recommendation list (10 for more general and 10 for more specific). The ranking list would be adjusted according to the requirement. For example, in the controlled group, if the task is asking the participants to pick out more general concepts to the initial Q-keyword, those top 10 concepts recommended by the proposed method sort by DBpagerank would be moved to the top. Similarly, in the experimental group, the top 10 concepts recommended by the proposed method sort by Semantic-aware PageRank against the Regulated Concept Map would be moved to the top. Therefore, we slipped concepts in the recommendation list to the following three levels:

- Level 1: Rank 1-3(sort by DBpagerank or Semantic-aware PageRank)
- Level 2: Rank 4-10(sort by DBpagerank or Semantic-aware PageRank)
- Level 3: Rank 11-20(sort by DBpagerank or Semantic-aware PageRank)

The observed results of the concepts selected by the participants corresponding to the levels are summarized in Table IV. The Chi-square test [22] of association evaluates relationships between those three levels above. We test the hypothesis as follow:

- H_0 : There is no relationship that exists on the three levels in the population.
- H_1 : There is a relationship that exists on the three levels in the population.
- $\alpha = 0.05$

TABLE IV: CONCEPTS SELECTED BY THE PARTICIPANTS CORRESPONDING TO THE LEVEL

	Level 1	Level 2	Level 3	Row Totals
Controlled group	44	53	23	120
Experimental group	53	51	16	120
Column Totals	97	104	39	240

The calculation of the Chi-Square statistic is as follow:

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e} \quad (10)$$

where f_0 is the observed counts in the cells and f_e is the expected frequency if NO relationship existed between three levels. By calculating, the chi-square statistic is 2.1299. The p-value is 0.344741. The result is not significant at $p < 0.05$. We have to reject H_1 , and state that there is no relationship among the three levels in the population. It means whatever the recommended concepts are sort by DBpagerank or Semantic-aware PageRank. Participants still could pick out those concepts which are expected as important concepts by the proposed method under a certain probability.

3) Analyzing results of post-test questionnaire

In the third task of each section, the post-test questionnaire was conducted. The purpose is to investigate the participants' perception of the satisfaction and effectiveness of the system. Likert Scale was employed as a measure. Regarding the satisfaction of the system, three questions will be asked:

- I am satisfied with the recommendations (0-5).
- The recommendations are useful to me (0-5).
- The recommendations are unanticipated (0-5).

Similarly, regarding the effectiveness of the system, three questions will be asked:

- The recommendations are relevant to the initial keyword (0-5).
- The recommendations enable me to strengthen the knowledge construction process (0-5).
- The recommendations make the investigation more efficient (0-5).

To measure the difference in the participants' perception of the satisfaction and effectiveness of the system between the controlled group and experimental group, the Mann-Whitney U Test [21] is employed. The purpose of this non-parametric measurement is to compare the difference between the two populations. The basis on which we make inferences is also based on the sampling distribution composed of all the possible sample characteristics. Therefore, we test the hypothesis below:

- H_0 : There is no difference in the participants' perception of the satisfaction of the system between the two groups.
- H_1 : There is a difference in the participants' perception of the satisfaction of the system between the two groups.
- $\alpha = 0.05$
- Sampling distribution $Z_{critical} = \pm 1.96$

We will skip the calculation of the Mann-Whitney U Test as mentioned previously. Since $Z_{statistic} \approx -5.478 < -1.960$ with $\alpha = 0.05$ (Two tailed), we have to reject H_0 , and state that there is a difference in the participants' perception of the satisfaction of the system between two groups.

By the observation of Descriptive statistics (Table V) for the participants' perception of the satisfaction of the system between two groups, the data of both groups forms a negative skewness. The data distribution of both groups is left-skewed which means whatever the concepts recommended by the proposed method are sort by DBpagerank or PageRank value against the regulated concept map, learners tend to be satisfied. However, only the data distribution of experimental forms a positive kurtosis which indicates a fat-tailed distribution. It refers to an increase in the probability of satisfaction scores being selected for an extreme number in

the experimental group.

TABLE V: DESCRIPTIVE STATISTICS OF THE PARTICIPANTS' PERCEPTION OF THE SATISFACTION

	Concepts selected in Controlled group	Concepts selected in Experimental group
Mean	11.675	12.6
Standard Error	0.2491	0.2449
Median	12	13
Mode	13	13
Standard Deviation	1.5752	1.5492
Sample Variance	2.4814	2.4000
Kurtosis	-0.2680	1.2748
Skewness	-0.7131	-0.7195
Range	6	7
Minimum	8	8
Maximum	14	15
Sum	467	504
Count	40	40

Similarly, we also test the hypothesis below regarding the effectiveness of the system:

- H_0 : There is no difference in the participants' perception of the effectiveness of the system between the two group.
- H_1 : There is a difference in the participants' perception of the effectiveness of the system between the two group.
- $\alpha = 0.05$
- Sampling distribution $Z_{critical} = \pm 1.96$

Since $Z_{statistic} \approx -5.748 < -1.960$ with $\alpha = 0.05$ (Two tailed), we have to reject H_0 , and state that there is a difference in the participants' perception of the effectiveness of the system between the two groups. By the observation of Descriptive statistics (Table VI) for the participants' perception of the effectiveness of the system between the two groups, the data of both groups form a negative skewness. The data of the experimental group forms more serious negative skewness. The data distribution of both groups is left-skewed which means whatever the concepts recommended by the proposed method are sort by DBpagerank or PageRank value against the regulated concept map, learners tend to be satisfied with the efficiency of the system. Moreover, the data distribution of the experimental forms a positive kurtosis which indicates a fat-tailed distribution. It refers to an increase in the probability of the participants' perception of the effectiveness scores being selected for an extreme number in the experimental group.

TABLE VI: DESCRIPTIVE STATISTICS OF THE PARTICIPANTS' PERCEPTION OF THE EFFECTIVENESS

	Concepts selected in Controlled group	Concepts selected in Experimental group
Mean	12.5	13.35
Standard Error	0.2375	0.2280
Median	12.5	14
Mode	12	14
Standard Deviation	1.5021	1.4420
Sample Variance	2.2564	2.0795
Kurtosis	-0.2524	1.5395
Skewness	-0.2628	-1.1959
Range	6	6
Minimum	9	9
Maximum	15	15
Sum	500	534
Count	40	40

VI. DISCUSSION

By observing the analysis results above, three major issues affect the performance of the proposed method.

Firstly, there existed unexpected recommendations by the proposed method. Take the more general concepts of the initial question **Machine learning** recommended by the proposed method as an instance (Table VII). We could realize that those concepts such as **Neuropsychological assessment** and **Euthenics** are quite difficult for learner to construct the knowledge for the initial question even though those concepts do exist semantic relationships between initial question.

TABLE VII: MORE GENERAL CONCEPTS OF MACHINE LEARNING RECOMMENDED BY PROPOSED METHOD AND ITS RANKING

Concepts	Sort by Semantic-aware PageRank
Artificial_intelligence	1
Learning	2
Personhood	3
Cognition	4
Memory	5
Computational_neuroscience	6
Cybernetics	7
Unsolved_problems_in_computer_science	8
Futurology	9
Emerging_technologies	10
Education	11
Cognitive_science	12
Neuroscience	13
Formal_sciences	14
Behavior	15
Computer_science	16
Intelligence	17
Neuropsychological_assessment	18
Cognitive_neuroscience	19
Euthenics	20

Secondly, we employed the semantic relations over the web to represent the interlinking between concepts extracted by the proposed method. However, as we mentioned the definition of the Connectivism, making a decision for interlinking between concepts is one of the knowledge management activities that learners would build an individual's cognitive network. During this process, learners may not consider the semantic relations over the web. The first task of the case study simulates the navigational learning of Web-based investigative learning that learners could navigate those resources recommended by the proposed approach for the knowledge construction process. Take Fig. 16 as an instance, it could not be wrong when we consider this decision making as a cognitive process.

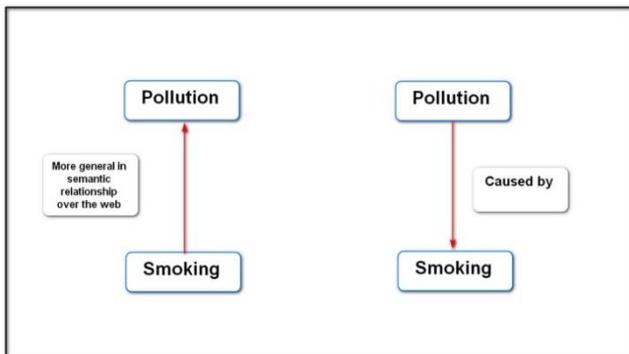


Fig. 16. An example of interlinking between two concepts.

Thirdly, according to the concept importance estimation we proposed above, we assume that the importance (Semantic-aware PageRank) of a concept node in Regulated Concept Map is determined by the number of outbound links on that concept. It means the evaluation of the relevance [5] between the recommended concepts and initial Q-keyword was regardless. The analyzed results of the second task in the case study confirmed that whatever the recommended concepts are sort by DBpagerank or Semantic-aware PageRank, participants still could pick out those concepts which are expected as important concepts by the proposed method under a certain probability. However, there is a gap between the priority of the relevance and importance when participants decide to pick out the related concept. For example, there existed a concept contain a serious Semantic-aware PageRank value over the Regulated Concept Map, but it is far away from the initial Q-keyword in DBpedia. Therefore, participants may tend to mark this concept as the related but not important one.

VII. CONCLUSION

In this work, we proposed a Semantic-awareness recommendation method of extracting and presenting related concepts at different levels for an initial question through LOD to promote the efficiency in the knowledge construction process in Web-based investigative learning. To prevent learners from the self-directed investigation, we proposed the regulated concept map generation to retrieve the relevant concepts at different levels with LOD and Semantic relations. For evaluating the relevance between initial Q-keyword and concepts in the regulated concept map, we defined the relativity as Semantic relations, node importance, and content containment for concepts we extracted from DBpedia. Owing to the finding of analyzing results measure by Spearman's Correlation Coefficient, we proposed the methodology of concept importance estimation maintained the most serious strength of the association between learner's expectation. We have also reported a case study whose purpose was to evaluate that using the Semantic-awareness recommendation with linked open data could help learners strengthen the knowledge construction process by discovering semantic related concepts during Web-based investigative learning. The results of the study suggest that Semantic-awareness recommendation with linked open data promotes the efficiency of the knowledge construction process.

VIII. FUTURE WORK

First of all, the proposed method in this work only supports recommending concepts based on one initial question. In order to support long-term learning scenario creation, a recommendation against multiple Q-keywords is needed. Secondly, as we mentioned above, the capabilities of LOD were underutilized. We could not tell that the semantic relations over the web are the best way for recommending concepts. Therefore, by combining different relations over the web, the results could be highly anticipated. Thirdly, there still is a room to improve the concept importance

estimation and filtering, such as applying certain techniques like machine learning and natural language processing.

CONFLICT OF INTEREST

The authors declare no conflict of interest.

AUTHOR CONTRIBUTIONS

Kang Ting conducted the research under the supervision of Shinobu Hasegawa. All authors had approved the final version.

REFERENCES

- [1] T. Berners-Lee and CERN, *Information Management: A Proposal*, 1989.
- [2] M. Hagiwara, A. Kashihara, S. Hasegawa, K. Ota, and R. Takaoka, "Adaptive recommendation for question decomposition in web-based in-vestigative learning," in *Proc. 2019 IEEE International Conference on Engi- neering, Technology and Education (TALE)*, 2019, pp. 1–9.
- [3] K. Kashihara and N. Akiyama, "A model of meta-learning for web-based navigational learning," *International J. of Advanced Technology for Learning*, vol. 2, no. 4, pp. 198–206, 2005.
- [4] J. Lehmann, R. Isele, M. Jentzsch *et al.*, "DBpedia - a large-scale, multilingual knowledge base extracted from Wikipedia," *Semantic Web*, vol. 6, no. 2, pp. 167–195, 2015.
- [5] T. Yamauchi, "Relevance between q-keywords corresponding to transition of interest in web-based investigative learning," Master's thesis, Japan Advanced Institute of Science and Technology, School of Information Science, Nomishi Ishigawa, 2020.
- [6] P. T. T. D. Noia, V. C. Ostuni, and E. D. Sciascio, "Sprank: Semantic path-based ranking for top-n recommendations using linked open data," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 8, no. 9, pp. 1–34, 2016.
- [7] M. G. C. Musto, G. Semeraro, and P. Lops, "Tuning personalized pagerank for semantics-aware recommendations based on linked open data," in *Proc. European Semantic Web Conference 2017*, 2017, pp. 169–183.
- [8] G. Siemens, "Connectivism: A learning theory for the digital age," *Elearnspace.org*, 2012.
- [9] A. Kashihara and N. Akiyama, "Learner-created scenario for investigative learning with web resources," in H. Chad Lane, Kalina Yacef, Jack Mostow, and Philip Pavlik, Eds. *Artificial Intelligence in Education*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 700–703.
- [10] Tim Berners-Lee design issues: Linked data. [Online]. Available: <https://www.w3.org/DesignIssues/LinkedData.html>
- [11] D. Vrande'ci c and M. Kroetzsch, "Wikidata: A free collaborative knowledgebase," *Communications of the ACM*, vol. 57, no. 10, pp. 78–85, 2014.
- [12] W3C Recommendation: Resource description framework (rdf) model and syntax specification. [Online]. Available: <https://www.w3.org/TR/1999/REC-rdf-syntax-19990222/>
- [13] DBpedia Blog: New 2016-04 DBpedia release. [Online]. Available: <https://blog.dbpedia.org/2016/10/19/yeah-we-did-it-again-new-2016-04-dbpedia-release/>
- [14] W3C Recommendation: SPARQL protocol for rdf. [Online]. Available: <https://www.w3.org/TR/rdf-sparql-protocol/>
- [15] W3C Recommendation SKOS simple knowledge organization system reference. [Online]. Available: <https://www.w3.org/TR/skos-reference>
- [16] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web search engine," *Computer Networks and ISDN Systems*, vol. 30, pp. 107–117, 1998.
- [17] L. Page, S. Brin, R. Motwani, and T. Winograd, "The pagerank citation ranking: Bringing order to the web," Stanford InfoLab, SIDL-WP-1999- 0120(1999-66), 1999.
- [18] H. H. Fu, K. J. Lin, and H. T. Tsai, "Damping factor in google page ranking," *Applied Stochastic Models in Business and Industry*, vol. 22, pp. 431–444, 2005.
- [19] S. Ichinose, I. Kobayashi, M. Iwazume, and K. Tanaka, "Ranking the results of dbpedia retrieval with sparql query," in *Revised Selected Papers of the Third Joint International Conference on Semantic Technology - Volume 8388*, JIST 2013, Berlin, Heidelberg: Springer-Verlag, 2013, pp. 306–319.
- [20] J. W. Gooch, *Spearman Rank Correlation Coefficient*, New York: Springer New York, 2008, pp. 502–505.
- [21] P. E. McKnight and J. Najab, *Mann-Whitney U Test*, American Cancer Society, 2010, p. 1.
- [22] K. L. Wuensch, *Chi-Square Tests*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 252–253.

Copyright © 2022 by the authors. This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).



Kang Ting received the B.Ed. degree in mathematics education from the Education University of Hong Kong in 2017.

He is a student at Japan Advanced Institute of Science and Technology, studying for the master's degree in information science under the supervision of Shinobu Hasegawa since 2019. His research interests are learning informatics and educational technology.



Shinobu Hasegawa received his B.S., M.S., and Ph.D. degrees in systems science from Osaka University in 1998, 2000, and 2002, respectively.

He is now an associate professor in the Research Center for Advanced Computing Infrastructure, Japan Advanced Institute of Science and Technology. His research areas include support for web-based learning, gamification, distance learning, and ICT in education.