

# Mathematics Content Understanding for Cyberlearning via Formula Evolution Map

Zhuoren Jiang

<sup>1</sup>Sun Yat-sen University  
Guangzhou, China

<sup>2</sup>Peking University, Beijing, China  
jiangzhr3@mail.sysu.edu.cn

Zheng Gao

Indiana University Bloomington  
Bloomington, IN, USA  
gao27@indiana.edu

Liangcai Gao\*

Peking University  
Beijing, China  
glc@pku.edu.cn

Zhi Tang

Peking University  
Beijing, China  
tangzhi@pku.edu.cn

Ke Yuan

Peking University  
Beijing, China  
yuanke@pku.edu.cn

Xiaozhong Liu\*

<sup>1</sup>Alibaba Group  
Seattle & Hangzhou, China  
<sup>2</sup>Indiana University Bloomington  
Bloomington, IN, USA  
liu237@indiana.edu

## ABSTRACT

Although the scientific digital library is growing at a rapid pace, scholars/students often find reading Science, Technology, Engineering, and Mathematics (STEM) literature daunting, especially for the math-content/formula. In this paper, we propose a novel problem, “mathematics content understanding”, for cyberlearning and cyberreading. To address this problem, we create a Formula Evolution Map (FEM) offline and implement a novel online learning/reading environment, PDF Reader with Math-Assistant (PRMA), which incorporates innovative math-scaffolding methods. The proposed algorithm/system can auto-characterize student emerging math-information need while reading a paper and enable students to readily explore the formula evolution trajectory in FEM. Based on a math-information need, PRMA utilizes innovative joint embedding, formula evolution mining, and heterogeneous graph mining algorithms to recommend high quality Open Educational Resources (OERs), e.g., video, Wikipedia page, or slides, to help students better understand the math-content in the paper. Evaluation and exit surveys show that the PRMA system and the proposed formula understanding algorithm can effectively assist master and PhD students better understand the complex math-content in the class readings.

## KEYWORDS

Cyberlearning; Education; Formula Understanding; Formula Layout; Formula Evolution

\*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

CIKM '18, October 22–26, 2018, Torino, Italy

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-6014-2/18/10...\$15.00

<https://doi.org/10.1145/3269206.3271694>

## ACM Reference Format:

Zhuoren Jiang, Liangcai Gao, Ke Yuan, Zheng Gao, Zhi Tang, and Xiaozhong Liu. 2018. Mathematics Content Understanding for Cyberlearning via Formula Evolution Map. In *The 27th ACM International Conference on Information and Knowledge Management (CIKM '18)*, October 22–26, 2018, Torino, Italy. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3269206.3271694>

## 1 INTRODUCTION

Over the past decade, while the volume of Science, Technology, Engineering, and Mathematics (STEM) publications has increased dramatically in university and digital libraries, few efforts have been made to help readers, especially junior scholars and graduate students, understand them. From a learning and reading viewpoint, understanding the content (especially the math-content) of scientific publications in STEM remains daunting [16]. In a survey conducted involving computer science program (35 Master and Ph.D. students), participants rated readings (textbook/publications) in data mining to be difficult (45.71%) or very difficult (14.29%). Furthermore, students claimed that math content in the readings was too difficult and inscrutable to understand because of the readers' limited knowledge in statistics and mathematics, i.e., participants believed the mathematical content in the readings to be difficult (51.43%) or very difficult (22.86%). Meanwhile, all the participants believe these papers are important, and they hope they could get additional help to better understand the math content in these papers. This survey motivated our thinking about this new problem - **Mathematics Content Understanding (MCU)**, a.k.a. how can we propose a useful method to assist readers to better understand the math-content in an academic publication. Junior students who struggle with math problem or scholars who want to conduct interdisciplinary research can especially get benefit from this study. To the best of our knowledge, this is the first study investigates the MCU problem. However, MCU faces the following challenges:

**Math content representation.** There are significant differences between math content (formula in most cases) and natural language. First, the mathematical symbols of a formula are ambiguous. For instance, the variable “ $\alpha$ ” or “ $x$ ” could represent the same meaning if defined by different scholars. Second, formulae can carry recursive structures while natural language is usually linear in structure.

Third, formulae are highly structured and usually presented in a layout form, e.g.,  $\text{\LaTeX}$  or *MathML*. Existing text mining method can be hardly used to address MCU problem.

**Information need shifting.** Student’s math information need could potentially shift when facing a complicated formula. For instance, to understand the formula of “*Latent Dirichlet Allocation (LDA)*”, one may need to understand the formulae of “*Dirichlet/Beta distribution*” (component) or even “*Conditional Probability*” (foundation). From the evolutionary viewpoint, these formulae can be considered as the “ancestors” of the original formula, and have important auxiliary effects for MCU. However, such information can not be fully extracted from the formula context or citations.

**Information access  $\neq$  information understanding.** Though traditional formula retrieval models could help user to access the math information. But, understanding information is fundamentally different from accessing information [15]. User need more supportive information to understand the math-content in the publications. For instance, [14] showed that cyberlearning resources (i.e., slides or video) can be more helpful (than scholarly publications) for scientific understanding.

In order to address these challenges, this paper proposes a novel solution, **MCU via Formula Evolution Map**, in a broad area of information retrieval and education. Although STEM publications generally do not place a premium on writing for readability, in this study, we hypothesize that the formula evolution information can be important to assist readers to better understand the math-content in a paper. As Figure 1 shows, we investigate the following two processes:

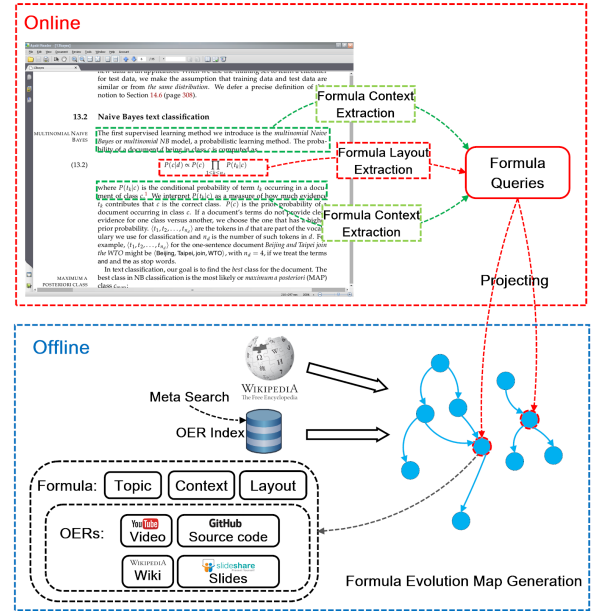
- In an *offline process*: By mining a large number of scientific knowledge-base documents and the associated formulae, from heterogeneous graph and joint embedding perspectives, we generate the **Formula Evolution Map (FEM)**, which encapsulates the mathematical evolutionary information over time.
- In an *online process*: By leveraging formula layout and context information extraction, we will **project the user information need (while reading the math-content in a publication) to the FEM**, as well as recommend useful resources to help readers better understand and consume the target publication and associated math content.

For FEM generation, the proposed algorithm explores more than 4 million documents and the associated math/formula/algorithmic information to characterize comprehensive and fundamental formula evolutionary relations in STEM (there are a total of 21,292,157 evolutionary relations on the FEM). To help readers/students better consume the target paper, in the FEM graph index, each formula vertex also associates the formula layout information, target topic (the formula belongs to), and a number of Open Educational Resources (OERs), i.e., video lectures, presentation slides, source codes, and Wikipedia pages, that may help readers to better consume the mathematical content in the online environment.

Meanwhile, in order to verify the proposed algorithm and cyberlearning hypothesis, we design a novel reading environment (PDF Reader with Math-Assistant, PRMA). When using PRMA system, a reader can easily highlight a formula with the mouse, and the system can automatically project the target formula in the paper to the vertexes on the backend FEM as well as recommend OERs for formula understanding. Evaluation results show that the proposed

method is important to help students understand the math-content in a paper, and its potential in cyberlearning is promising.

**The contribution of this paper is fourfold.** First, we propose an innovative MCU problem in an education context to help students and junior scholars better consume the STEM publications. Second, a new reading environment is employed to capture student information need while enabling them to highlight the confusing formula of the reading. Third, novel algorithms are proposed to characterize *formula evolution information* in a map by mining a massive scientific knowledge base. Last but not least, an extensive experiment (with 52 participants) is employed to qualitatively and quantitatively validate the proposed formula understanding hypothesis as well as the usefulness of the system and to evaluate the FEM generation quality. As MCU is a newly proposed but important problem, we share the algorithm generated FEM plus massive formula and math-topic information to motivate further investigation.



**Figure 1: The whole framework for this study**

Based on the experiment participants’ feedback, 72.73% of participants believed the proposed method can provide very useful information for math-understanding and 75.75% of participants believed the system recommended OERs (especially for videos and slides), comparing with the text content, are much more helpful for math-understanding. Algorithm evaluation also shows that FEM and the associated formula evolution relations are very important for math-understanding (can enhance precision and NDCG significantly).

## 2 PROBLEM FORMULATION

As aforementioned, simply return academic papers may not be enough to help junior scholars [14]. In this study, we hypothesize that formula evolution information along with OERs can be important to address MCU problem.

**DEFINITION 1. Formula Evolution Map (FEM).** FEM is defined as a weighted directed graph  $G = (F, R, \tau)$ , where  $F$  denotes the

formula vertex set, and  $R \subseteq F \times F$  denotes the directed evolution relation set.  $\omega$  is the relation weight set, denotes formula evolution probabilities.

FEM encapsulates the fundamental and enlightened formula evolution information, which could be especially useful for exploring the development of a formula as well as the details of its components. For instance, formula of “Bayes’ theorem” is the foundation of “Naive Bayes classifier” formula, and formulae of “Gaussian naive Bayes” and “Multinomial naive Bayes” are both the specific forms of general “Naive Bayes”. There are clear evolutionary paths among these formulae.

**DEFINITION 2. Mathematics Content Understanding via OER.** From OER recommendation viewpoint, the MCU problem can be defined as a conditional probability  $P(\text{OER}|\text{info-need})$ , i.e., the probability of an OER given a particular math information need, which can be formalized as:

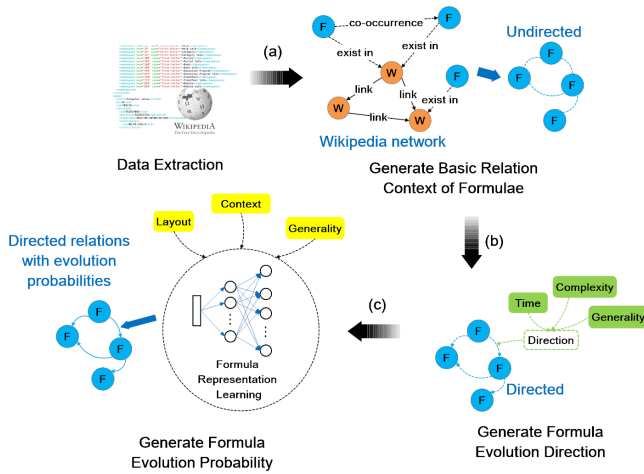
- **Input:** A mathematical content (a formula with its context).
- **Output:** A list of ranked OERs that could be potentially useful for understanding the target math-content.

**DEFINITION 3. Mathematics Content Understanding via FEM + OER.** Based on prior definitions, we can further integrate FEM into consideration, i.e.,  $P(\text{OER}|\text{info-need}) = P(\text{OER}|\text{formula}) * P(\text{formula}|\text{info-need})$ , where user’s math information need can be projected to a formula with its ancestors (vertexes) on FEM for MCU. Then, the projected formulae (on FEM) and their related OERs can help to address MCU problem.

### 3 METHODOLOGY

In this section, we discuss the research method in detail including: generating the Formula Evolution Map (FEM) offline (3.1), designing the novel cyberlearning environment to characterize readers’ information need when consuming mathematics content in a paper (3.2), and designing OER recommendation for formula understanding (3.3).

#### 3.1 Formula Evolution Map Generation



**Figure 2: Formula evolution map generation**

We generate a Formula Evolution Map (FEM), offline, to interconnect important/fundamental formulae in scientific publications.

Note that the formulae in the FEM cannot cover all the formulae in the readings. Instead, it provides the potential to 1) project any formula in any paper to the vertex(es) in FEM, and 2) trace the formula evolution information on FEM for math-understanding.

To achieve this goal, we employ a large knowledge base, Wikipedia, to generate FEM. There are two reasons we use Wikipedia: first, Wikipedia contains a wealth of mathematical information, including tens of thousands of fundamental formulae (with math-topic and formula layout information); second, Wikipedia provides links between pages, which can be important to generate formula evolution information. In this paper, we use the text information, formula layout information, and link topology information extracted from Wikipedia dump to generate FEM that enables an algorithm to estimate formulae evolution over time. More importantly, the proposed FEM can minimize the noisy formulae information.

The FEM generation process is illustrated in Figure 2, which involves three steps: (a) formula evolution relation generation, (b) formula evolution direction determination, (c) formula evolution probability calculation.

**(a) Formula Evolution Relation Generation:** By using the hyperlinks between Wikipedia pages and the co-occurrence relationships (in the same Wikipedia page) of formulae, we generate the basic relation context of formulae. The formula evolution relation generation can be modeled as:

$$R(f_a, f_b|W) = \text{Sgn}_r(w_a, w_b) = \begin{cases} 1 & \text{Co}(w_a, w_b) = 1 \text{ and } w_a = w_b \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Here,  $R$  is a signum function  $\text{Sgn}_r$ , that indicates the evolution relation existence between formula  $f_a$  and formula  $f_b$  based on the Wikipedia page network  $W$ .  $w_a$  is the home page of  $f_a$ , and  $w_b$  is the home page of  $f_b$  (extracted from the Wikipedia dump).  $\text{Co}(w_a, w_b) = 1$  indicates that there is a hyperlink between  $w_a$  and  $w_b$ , and  $w_a = w_b$  means  $f_a$  and  $f_b$  are hosted in the same page. When  $R(f_a, f_b) = 1$ , there could be a candidate evolution relation between  $f_a$  and  $f_b$ .

In the generated undirected formula relation network, each formula is characterized as a vertex with multiple attributes: (1) Wikipedia page title, (2) formula context information (250 characters), and (3) formula layout information.

**(b) Formula Evolution Direction Determination:** In this study, the evolution direction between formulae are determined by three indicators (assumptions): (1)  $\lambda_t(f)$ , formula birth time (formulae could evolve from past to present); (2)  $\lambda_g(f)$ , formula generality (formulae could evolve from fundamental to contextualized); (3)  $\lambda_c(f)$  formula layout complexity (formulae could evolve from simple to complex).

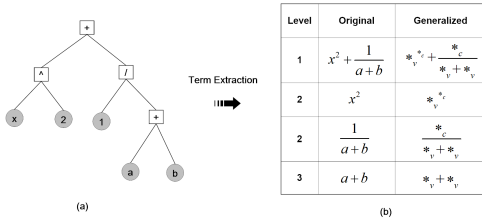
For a formula pair  $\{f_a, f_b\}$ ,  $R(f_a, f_b|W) = 1$ , the evolution direction is first decided by  $\lambda_t$ . However, not every formula has  $\lambda_t$  attribute. If the “birth time” of a formula is missing, we use  $\lambda_c$  to determine the direction. To avoid the uncertainty caused by the layout comparison, if the ratio of the complexity difference  $< 0.1$ , we use  $\lambda_g$  as the final direction indicator.

To generate  $\lambda_t(f)$ , we use the title of home Wikipedia page to represent the formula, then we use the greedy match algorithm in a large academic paper corpus to find the earliest appearance of the formula. Note that one formula may exist in multiple Wikipedia

pages. We use the first appearance time among the Wikipedia concepts as this formula's creation time. The smaller  $\lambda_t(f)$  is, the earlier the formula  $f$  appears.

Meanwhile, PageRank [23] is utilized to calculate the formula generality  $\lambda_g(f)$ , for measuring the fundamental level of a formula. The underlying assumption is similar as PageRank: more fundamental formulae are likely to receive more links from other formulae.  $\lambda_g$  is calculated via a formula-formula graph (generated from the page-page wiki-graph), and each vertex in the graph is a formula. The generality of a formula is voted by the links among formulae. We hypothesize that formulae could evolve from past to the present.

In this study, a formula semantic tree based approach is used for calculating  $\lambda_c(f)$ . We first parse the  $\text{\LaTeX}$  expressions of formulae and convert them into Presentation MathML expressions. Then we construct a formula tree using a semantic tree-constructed algorithm proposed in [12] (Figure 3 (a) shows an example of semantic tree presentation of the formula  $x^2 + \frac{1}{a+b}$ ). After that, we extract formula terms hierarchically from the constructed semantic tree. The extraction algorithm is described in Algorithm 1. In the proposed algorithm, there are two kinds of formula terms: original terms and generalized terms. The original terms are generated directly from the original substructures of the semantic tree presentation of the formula. The generalized terms are proposed by changing the variables and constants of the original terms into wildcards (describe the sketch of the formula structure, for fuzz representation). Variables are replaced by  $*_v$ , and constants are represented as  $*_c$ . There is a "level" attribute extracted for each term which denotes the level of the term in the semantic tree whose root's level is 1. Figure 3 (b) shows the sub-tree levels and terms of formula " $x^2 + \frac{1}{a+b}$ ". The similar formula tree layout presentation has been proven an effective method for formula retrieval task [5].



**Figure 3: (a) Semantic tree presentation and (b) term extraction of  $x^2 + \frac{1}{a+b}$**

Based on the formula semantic tree and its extracted terms, the  $\lambda_c(f)$  can be calculated as:  $\lambda_c(f) = \sum_{i=1}^N (L(t_i))$ , in which,  $N$  is the total formula term number,  $t_i$  is a formula term, and  $1 \leq i \leq N$ ,  $L(t_i)$  is the level of  $t_i$  in the formula tree. The greater  $\lambda_c(f)$  is, the more complicated formula  $f$  is.

**(c) Formula Evolution Probability Calculation:** Although Wikipedia provides multitudinous evidence that can be used to infer candidate evolution relations, some of them can be noisy. For instance, page "Artificial neural network" has a link to page "Algorithm", but not every formulae in "Algorithm" page can be evolutionary to the "Activation function" formula, some of them could be noisy. This kind of noisy formulae will not be useful for math-content understanding. To address this problem, we propose

the formula evolution probability to characterize the reliability of a formula evolution relation. The evolution probability from  $f_a$  to  $f_b$  can be modeled as:

$$P(f_a \xrightarrow{e} f_b) = \pi(\phi(f_a), \phi(f_b)) \quad (2)$$

where,  $P(f_a \xrightarrow{e} f_b)$  is the formula evolution probability from  $f_a$  to  $f_b$ ,  $\phi$  is a representation function, which can project each formulae to a low-dimensional joint embedding space from context, layout and generality viewpoints.  $\pi$  is probability scoring function based on the learned formula embeddings.

---

**Algorithm 1** Formula Semantic Term Extraction

---

**Input:** Formula Semantic Tree,  $ST$

**Output:** Set of formula terms,  $T$

```

1: Let  $O(ST)$  be the original semantic tree
2: Let  $G(ST)$  be the generalization of the semantic tree
3: Let  $L(ST)$  be the level of the semantic tree
4: procedure EXTRACTOR( $ST, L(ST)$ )
5:   if  $ST$  is not a leaf then
6:      $T+ = (O(ST), L(ST))$                                  $\triangleright$  original term
7:      $T+ = (G(ST), L(ST))$                                  $\triangleright$  generalized term
8:     for  $ST_i \leftarrow$  each child of  $ST$  do
9:       EXTRACTOR( $ST_i, L(ST) + 1$ )
10:    end for
11:  end if
12: end procedure

```

---

In this work, we construct the formula representation via semi-supervised graphical learning. Following the evolution relations in FEM, we can simulate a random walk of fixed length  $l$  with a set of parameters  $\theta$  to guide the walker on the graph. Let  $f_i$  denote the  $i$ th formula in the walk, which can be generated by the following distribution:

$$P(f_i | f_{i-1}) = \tanh[\tau(P_t, P_l, P_g | \theta)] \\ = \tanh[\theta_t P_t(f_i | f_{i-1}) + \theta_l P_l(f_i | f_{i-1}) + \theta_g P_g(f_i | f_{i-1})] \quad (3)$$

where  $P(f_i | f_{i-1})$  denotes the normalized transition probability between  $f_i$  and  $f_{i-1}$ .  $\tau(\cdot)$  is a trivariate function that can fusion three different kinds of transition probabilities: context  $P_t$ , layout  $P_l$ , and generality  $P_g$ .  $\theta = \{\theta_t, \theta_l, \theta_g\}$  is the non-negative fusion parameters that control the contribution of each transition probability. For this study, we set  $\theta_t = \theta_l = \theta_g = 1$ , and more sophisticated parameter tuning will be saved for future.  $\tanh(\cdot)$  is the hyperbolic tangent function for normalization.

For context transition probability estimation, a 250-word text window around the formulae was employed along with language model with Dirichlet smoothing [35].

$$P_t(f_i^t | f_{i-1}^t) \propto P_t(f_{i-1}^t | f_i^t) P_t(f_i^t) \quad (4)$$

$f_i^t$  represents the context of a formula,  $P_t(f_i^t | f_{i-1}^t)$  is the posterior probability,  $P_t(f_{i-1}^t | f_i^t)$  is the  $f_{i-1}^t$  likelihood given  $f_i^t$ ,  $p(f_i^t)$  is assumed to be uniform. We hypothesize that if two formulae share the similar context, they may have a high evolution probability.

For formula layout, we calculate formula transition probability by leveraging formula semantic layout tree and its extracted terms:

$$P_l(f_i^t | f_{i-1}^t) = \frac{\omega_{cov}(f_i^t, f_{i-1}^t) \sum_{t_n \in f_{i-1}^t} [\omega_{gen}(t_n) \omega_{lel}(t_n, f_i^t, f_{i-1}^t)]}{\sum_{t_n \in f_{i-1}^t} [\omega_{gen}(t_n)]} \quad (5)$$



where,  $f_i^l$  is the formula term set generated from the semantic layout tree;  $\omega_{cov}(f_i^l, f_{i-1}^l) = \frac{|f_i^l \cap f_{i-1}^l|}{|f_{i-1}^l|}$ , denotes the ratio between matched term number and total term number of  $f_{i-1}$ ;  $\omega_{gen}(t_n)$  is the penalty parameter for the generalized terms, if  $t_n$  is a generalized term,  $\omega_{gen}(t_n)$  is empirically set to 0.5 [12], otherwise,  $\omega_{gen}(t_n) = 1$ ;  $\omega_{lev}(t_n, f_i^l, f_{i-1}^l)$  is the term level weight, affected by the minimum level distance of the matched formula term  $t_n$  in  $f_i$  and  $f_{i-1}$ ,  $\omega_{lev}(t_n, f_i^l, f_{i-1}^l) = \frac{1}{1 + \min_j \{ |level(t, f_{i-1}^l) - level_j(t, f_i^l)| \}}$ . We hypothesize that, if two formulae have similar layout trees, they may have a high evolution probability.

Formula generality transition probability can be calculated as:

$$P_g(f_i | f_{i-1}) = \lambda_g(f_{i-1}) \quad (6)$$

which is the formula generality of  $f_{i-1}$ . We hypothesize that, a fundamental formula can have a high evolution probability to its variants with more detailed contextual constraints.

We then use  $\phi: F \rightarrow \mathbb{R}^d$  as the mapping function (from formula vertexes) for representation learning, where  $d$  specifies the number of dimensions.  $\phi$  is a matrix of size  $|F| \times d$  parameters. Feature learning methods are based on the Skip-gram architecture [18]. For every starting formula vertex  $f \in F$ , we define  $N_S(f) \subset F$  as a network neighborhood (“context”) of vertex  $f$  generated through the proposed neighborhood sampling strategy  $S$  (semi-supervised random walk guided by  $\theta$ ). The objective function can be formalized as:

$$\max_{\phi} \sum_{f \in F} \log P(N_S(f) | \phi(\vec{f})) \quad (7)$$

Stochastic gradient ascent is used to optimize the joint embedding model parameters of  $\phi(\vec{\cdot})$ . Negative sampling [18] is applied for optimization efficiency. In this study, we use cosine similarity (with a ReLU function) of the optimized formula representations as  $\pi$  for scoring the formula evolution probability.

By calculating the evolution probability for formula relation, we are able to rule out the noisy relations (with low evolution probabilities) and better explore the formula evolution trajectory to help users understand the essence of the target formula.

### 3.2 Math-Information Need Characterization

To help students better understand the math-content of a scientific publication in a course environment, we design a novel system, PDF Reader with Math-Assistant (PRMA). As Figure 4 shows, the new system has three main functions:

- Capture evidence and characterize students’ emerging implicit/explicit information needs when encounter a formula understanding problem. For instance, students can ask a specific question given a formula (explicit information need), or easily highlight a formula with mouse in the paper, as evidence of an implicit information need. In either case, the PRMA is able to extract the formula layout presentation and formula context from the target PDF paper.
- Automatically project the target formula onto the formula evolution map, which allows students to navigate the formula evolution trajectory in FEM while helping them understand the target formula.
- Automatically recommend high quality OERs for the target formula, such as video lectures, slides, source code, or Wikipedia pages, to resolve students’ information needs while helping them

understand the formula. (We crawl and pre-index massive OERs by using meta-search algorithms provided in [15].)

At the front-end, the mathematical expressions are obtained by a symbol dominance based formulae recognition algorithm proposed in [36], and the formula context is extracted by a PDF parser. At the backend, the PRMA has access to the generated FEM, the list of assigned class readings (title, abstract, full content, associated topics and citation information), and the formula-based OER recommendations algorithm.

**Table 1: Formula projecting features for math-information need characterization\***

No.	Projecting feature	Mathematical Definition
1	Formula Context Feature	$p(f_i^t   f_c^t) p(f_c^t)$
2	Formula Context Keyword Feature	$\sum_{k_i \in f_i^t} p(f_i^{k_i}   f_c^t) p(f_c^t)$
3	Question Text Feature	$p(f_i^{qt}   f_c^t) p(f_c^t)$
4	Question Text Keyword Feature	$\sum_{q_{k,i} \in f_i^{qt}} p(f_i^{q_{k,i}}   f_c^t) p(f_c^t)$
5	Formula Layout Feature	$\frac{\sum_{t_i \in f_i^l} [\omega_{gen}(t_i) \omega_{lev}(t_i, f_c^l, f_i^l)]}{[\omega_{cov}(f_c^l, f_i^l)]^{-1} \sum_{t_i \in f_i^l} [\omega_{gen}(t_i)]}$
6	Paper Idea Feature	$p(f_i^{pabs}   f_c^t) p(f_c^t)$
7	Paper Keywords Feature	$\sum_{k_i \in f_i^p} p_K p(f_i^{k_i}   f_c^t) p(f_c^t)$
8	Weekly Topic Feature	$\sum_{w_i \in f_i^p} p_W p(f_i^{w_i}   f_c^t) p(f_c^t)$
9	Context Evolution Feature	$p(f_m^t   f_c^t) p(f_c^t)$
10	Layout Evolution Feature	$\frac{\sum_{t_i \in f_m^l} [\omega_{gen}(t_i) \omega_{lev}(t_i, f_c^l, f_m^l)]}{[\omega_{cov}(f_c^l, f_m^l)]^{-1} \sum_{t_i \in f_m^l} [\omega_{gen}(t_i)]}$
11	Generality Evolution Feature	$\lambda_g(f_c)$
12	Evolution Distance Feature	$ f_m \rightsquigarrow f_c $

Formula Text Feature Group
Formula Layout Feature Group

Paper Content Feature Group
Formula Evolution Feature Group

\*Because of the space limitation, the detailed feature description(hypothesis) will be available in <https://github.com/GraphEmbedding/FEM>

The algorithms presented in the next section can recommend the optimized OERs given the math-information need, which will be able to help readers better understand the essence of the targeted formula. Meanwhile, readers can also provide usefulness feedback for system recommended OERs. For instance, as Figure 4 shows, readers can click “Good”, “OK”, or “Bad” for each recommended OER given their information needs. The judgments and student click information will be saved as system logs. The formula evolution explore behavior (click the formula evolution map) will also be recorded, which will be important for formula-based OER recommendation algorithms, i.e., training learning to rank model, and algorithm evaluation.

Although FEM can be potentially helpful for MCU, it’s still challenging to characterize the student’s emerging information needs while facing a formula in the target paper. The critical problem is how to “project” a puzzling formula to one or a number of formula vertex(es) in the FEM.

To address this problem, by using PRMA, we employ multiple Formula Projecting Features (FPF) to characterize the math-information need. The proposed FPF mainly focuses on four aspects: (1) the math-information need can be related to formula layout, (2)

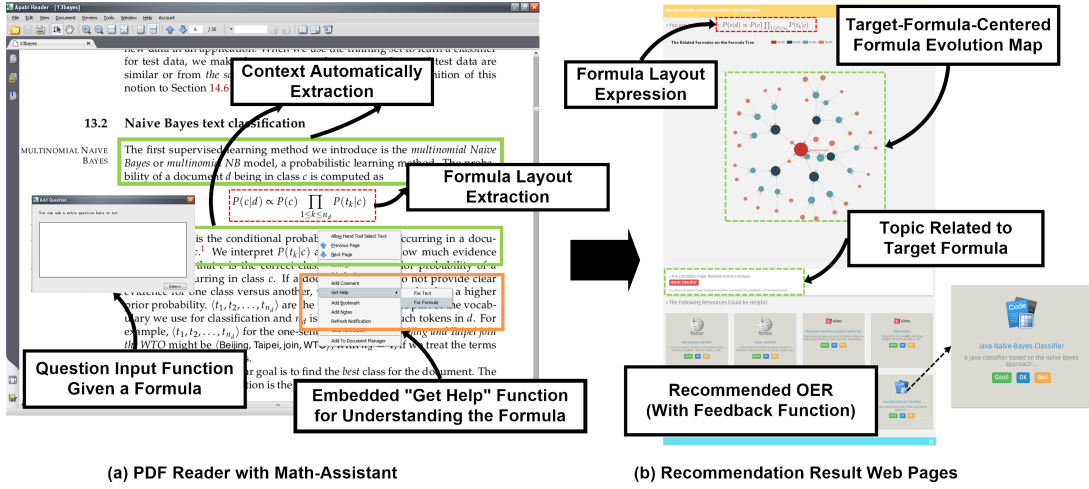


Figure 4: PDF Reader with Math-Assistant (PRMA) System

related to the text information of formula, e.g., students' questions about the formula, context and their associated topics, (3) related to the content of paper where formula exists, i.e., abstract, keywords and weekly topics in the syllabus (in a course environment), (4) related to the formulae located in the evolution trajectories in FEM given a matched formula. Detailed math definition is provided in Table 1.

### 3.3 OER Recommendation via FEM Mining

Figure 5 illustrates the graphical OER recommendation via FEM mining towards the mathematical query from PRMA, the vertexes and edges are depicted in Table 2. There are totally 48 offline OER ranking features (ORF) constructed<sup>1</sup> for this study, which can be divided into two groups:

(1) text-mining-based feature group, for instance, we calculate the  $p(OER|formula)$  based on the OER's text description and formula context using the language model. Note that we utilize the reading paper's content information (abstract, keywords, and weekly topics) for constructing text-mining-based features, which means even for a same formula, if the reading paper is changed, the recommended OERs will change correspondingly.

(2) Heterogeneous-graph-ranking feature group, a formula vertex in FEM can random walk to the OERs in the scholarly heterogeneous graph, i.e., OER ranking given a formula vertex in FEM. In this study, we employ meta-path plus random walk from the formula to candidate OER as graph-ranking features. For instance,  $F^* \xrightarrow{m} K \xrightarrow{p} R^2$  is a graph-ranking feature, which denotes that if an OER has a high probability relation from the keywords (topics) vertexes extracted from question formula's context, this OER should be recommended. The random walk probability can be estimated by:

$$r(v_i^{(1)}, v_j^{(l+1)}) = \sum_{t=v_i^{(1)} \rightsquigarrow v_j^{(l+1)}} RW(t), RW(t) = \prod_j w(v_{ij}^{(j)}, v_{i,j+1}^{(j+1)}) \quad (8)$$

<sup>1</sup>Because of the space limitation, we cannot provide more detailed OER ranking features. The detailed feature list will be available in <https://github.com/GraphEmbedding/FEM>

where  $t$  is a tour from  $v_i^{(1)}$  to  $v_j^{(l+1)}$  following the meta-path[30], and  $RW(t)$  is the simulated random walk probability of the tour  $t$ . Suppose  $t = (v_{i1}^{(1)}, v_{i2}^{(2)}, \dots, v_{il+1}^{(l+1)})$ ,  $w(v_{ij}^{(j)}, v_{i,j+1}^{(j+1)})$  is the weight of edge  $v_{ij}^{(j)} \rightarrow v_{i,j+1}^{(j+1)}$ . More detailed algorithm can be found in [15].

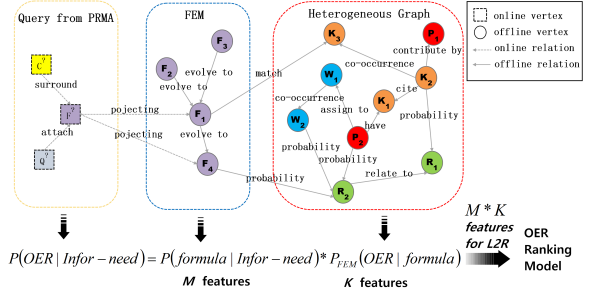


Figure 5: Graphical representations of OER recommendation via FEM mining

In section 3.2, we proposed a number of FPF (i.e.,  $M$  formula projecting features, see Table 1) for formula projecting, and the ORF (i.e.,  $K$  OER ranking features) are presented in this section. By using the formulae in FEM as the transition, we integrate FPF and ORF to recommend useful OER for math content understanding. The OER recommendation probability  $P(v_r|v_f)$ , i.e., an OER  $r$  based on the query formula  $f$ , can be calculated as:

$$P(v_r|v_f) = \sum_{m=1}^M \sum_{k=1}^K \omega_{m,k} \cdot \mathbb{E}_{F_m \in FPF}(v_f, v_f) \cdot \mathbb{E}_{F_k \in ORF}(v_f, v_r) \quad (9)$$

Here,  $\omega_{m,k}$  is the feature weight for  $m_{th}$  FPF and  $k_{th}$  ORF,  $\mathbb{E}_{F_m \in FPF}(v_f, v_f)$  is the online query formula projecting probability based on different FPF,  $\mathbb{E}_{F_k \in ORF}(v_f, v_r)$  is the OER ranking probability based on different ORF. There are totally  $M * K$  features for OER recommendation. In order to avoid laborious parameter tuning, we utilize learning to rank algorithm [11] to jointly optimize the

feature weights for FPF and ORF. The training data (OER usefulness judgments) are collected via PRMA system.

**Table 2: Vertexes and edges of OER recommendation via FEM mining**

Vertex	Description
$R$	Open Education Resource
$P$	Paper
$K$	Keyword
$W$	Weekly Topic (from Syllabus)
$F$	Formula
$F^?$	Query Formula
$C^?$	Context of Query Formula
$Q^?$	User Additional Question
Edge	Description
$P \xrightarrow{h} K$	Paper is related to keyword (using Labeled LDA [26])
$P \xrightarrow{a} W$	Paper is assigned to weekly topic (probability)
$P \xrightarrow{p} R$	Paper-resource relationship based on $p(R P)$ (language model)
$K \xrightarrow{cite} K$	Keyword cites keyword (probability)
$K \xrightarrow{co} K$	Keyword-keyword co-occurrence (probability)
$K \xrightarrow{cont} P$	Keyword is contributed by paper (using PageRank with prior [32])
$K \xrightarrow{p} R$	Keyword-resource relationship based on $p(R K)$ (language model)
$W \xrightarrow{co} W$	Weekly topic-weekly topic co-occurrence (probability)
$W \xrightarrow{p} R$	Weekly topic-resource relationship based on $p(R W)$ (language model)
$R \xrightarrow{r} R$	OER is related to OER (collected from service sites)
$F \xrightarrow{p} R$	Formula context-OER content relation based on $p(R F)$ (language model)
$F \xrightarrow{m} K$	Formula context is related to keyword (greedy match algorithm)
$F \xrightarrow{e} F$	Formula evolves to formula (probability)
$C^? \xrightarrow{s} F^?$	Context surrounds formula
$Q^? \xrightarrow{at} F^?$	Additional question is attached to formula
$F^? \xrightarrow{p} F$	Query formula is online projected to the formula on FEM

## 4 EXPERIMENT

### 4.1 Dataset and Experiment Setting

We tested this reading system and the associated mathematics content understanding algorithms in a real learning environment. Two graduate-level information retrieval courses at Indiana University were used for this experiment. A total of 52 students (Masters and Ph.D.s) voluntarily participated this experiment, and they were required to use the PRMA system for eight weeks (with 15 required readings, 10 chapters of an IR book, and 5 ACM journal papers). They could use PRMA with university account, and PRMA enables formula understanding function (e.g., highlight a formula in the reading and access formula evolution trajectories in FEM as well as recommended OERs). Meanwhile, we asked each participant to provide OER relevance judgments for the system-recommended OERs. There were a total of 7,099 valid judgments collected (for 622 student requests), and we used those judgments to train the learning to rank model and to evaluate the algorithm performance. Among the 622 student requests, only 29 (4.7%) requests contained an explicit question. This phenomenon indicates that most students don't want to input a specific question when facing a formula in PRMA.

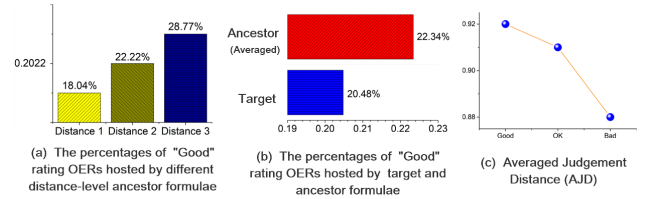
At the backend of PRMA, we created a heterogeneous graph for OER recommendation. For paper vertexes, we used 248,893

publications from 1,553 venues (in ACM digital library). The paper vertexes were connected to 7,190 keyword labeled topics, and the publication data were also used to generate formula birth time  $\lambda_t(f)$ . According to the syllabus, there were a total of 60 weekly topics. By using meta-search, we collected a total of 1,112,718 OERs.

We used a Wikipedia dump of July 30, 2014. There were 358,116 raw formulae, 34,683 formula home pages, 198,336 page-formula ownership relations, and 74,947,670 page hyperlinks. For the experiment, we kept the formulae featuring at least two variables and three operators. There were 194,150 formulae left, and the generated FEM had 21,292,157 potential formula evolution relations.

As Figure 4 (b) shows, for each query formula, we visualized a three-level (distance) target-formula-centered evolution sub-graph. Meanwhile, the visualized formula vertex had at least 0.5 probability ( $P(f_a \xrightarrow{e} f_b) \geq 0.5$ ) to connect to its neighbor.

### 4.2 Experiment Results









**Figure 6: Statistics of OER judgements**

Among the 622 student requests, 610 of them (98.1%) contained at least one OER rated as "Good" or "OK" for their "mathematics content understanding". For all the OER judgments, participants rated 19.72% of the recommended OERs as "Good", 37.61% as "OK", and 42.67% as "Bad". Note that, the students were asked to rate at least top 5 OERs for each request, because we need to collect a amount of "Bad" judgments for model training and evaluation. Based on the target-centered formula evolution map provided by PRMA, users could not only consume the recommended OERs for target formula but also freely explore the formula evolution map to check the recommended OERs for ancestor formulae. For any formula request, PRMA could record the user OER judgement, the ancestor formula (on FEM) that hosted the OERs, and the distance from ancestor formula to target formula on FEM. As Figure 6 (a) shows, while the (evolutionary) distance between target and ancestor formula increasing, the percentage of "Good" judgements for OERs (ancestor formula hosted) raises (from 18.04% to 28.77%). Meanwhile, as Figure 6 (b) indicates, there are more percentages of "Good" rating OERs from the ancestor formulae. This finding demonstrates that the ancestor formulae on FEM can be especially important to assist students better consume the math-content in a paper, and, when the target formula is complex, the background information (e.g., the ancestor formulae on FEM) can be more useful.

In order to further explore the relationship between OER judgements and evolutionary distance, we calculate the Average Judgement Distance (AJD) from ancestor to question (target) formula on FEM for each type of judgements:  $AJD_{type} = \frac{\sum_{i=1}^{N_{type}} d_i}{N_{type}}$ .  $N_{type}$  is the total number of a specific type of judgement (e.g., "Good", "OK", etc.);  $d_i$  should be one of  $\{0, 1, 2, 3\}$ , representing the evolutionary distance between OER hosted formula and target formula, for  $i_{th}$

judgement of this specific type. The result,  $AJD_{\text{Good}} (0.92) > AJD_{\text{OK}} (0.91) > AJD_{\text{Bad}} (0.88)$ , indicates the remoter ancestor formulae (with recommended OERs) can be more helpful for math information understanding. This finding also proves FEM can be effective to address the student’s math-information need for math-understanding (without FEM, it is not feasible to trace the ancestor formulae and math-evolution for information understanding).

**Table 3: Baseline groups and comparison groups for OER ranking experiment**

Baselines	
 $Rank_{abs}$	Use the reading’s abstract to represent the user’s math-information needs, then generate the OER ranking based on $p(OER abstract)$ (using Language Model with Dirichlet prior smoothing [35], the same below).
$Rank_{keyword}$	Use the reading’s keywords to represent the user’s math-information needs: $p(OER keywords)$ .
$Rank_{context}$	Use the formula’s context to represent the user’s math-information needs: $p(OER context)$ .
Comparison Groups	
 $L2R_{layout}$	L2R (learning to rank) model using formula layout feature[12] plus all OER ranking features.
$L2R_{text}$	L2R model using formula context feature, formula context keyword feature, question text feature and question text keyword feature plus all OER ranking features.
$L2R_{content}$	L2R model using paper idea feature, paper keywords feature and weekly topic feature plus all OER ranking features.
$L2R_{multiple}$	L2R model using formula layout feature group, formula text feature group and paper content feature group plus all OER ranking features.
 $L2R_{all}$	L2R model using all formula projecting features plus all OER ranking features.
 Without FEM  With Partial FEM  With Complete FEM	

As this study is not focusing on learning to rank (L2R), we used a relative simple list-wise algorithm, Coordinate Ascent[17], which iteratively optimizes a multivariate objective ranking function, for formula and OER features integration and algorithm evaluation. Meanwhile, we needed to employ the baseline groups for comparison. However, as mathematics content understanding is a newly proposed problem, and few existing algorithms addressed this problem. We chose three classic methods as baseline groups (without FEM assisted) and five L2R models for different feature groups (with partial or complete FEM assisted) as comparison groups. The baselines (e.g., text and formula retrieval models) and comparison groups are listed in Table 3.

From an NDCG viewpoint, we scored  $Good = 2$ ,  $OK = 1$ , and  $Bad = 0$ . The OER recommendation performance can be found in Table 4. As the OER recommendation was more like a QA problem and students were more interested to find the first useful resource, we used MRR (Mean Reciprocal Rank) as the metric to train the learning to rank model. For evaluation, 10-fold cross-validation was used.

From a performance viewpoint, evaluation results show that, first, FEM can provide important information for OER recommendation and math-understanding. For instance,  $L2R_{all}$  (the best performed method empowered with complete FEM information), compared with the baseline groups (without FEM assisted),  $P@3$  has an average increase of 15.2%,  $NDCG@3$  improves with 20.8%,  $MAP$  increases 9.7% and  $MRR$  enhances 10.9%. Meanwhile,  $MRR$  score of  $L2R_{all}$  is higher than 0.88 (which means students are finding

the useful OERs (“Good” or “OK”) in the 1st position in almost result ranking list). This finding is also confirmed in the exit survey where 72.73% of participants believe the PRMA system along with recommended OERs can provide precise and useful information for math-understanding. It is clear that FEM plays a critical role in the proposed framework, and FEM can provide very helpful information for math content understanding.

Second, based on the student judgments, we found that a number of formula projecting features, can be potentially useful. From a ranking viewpoint, all L2R models outperform the baseline groups of  $Rank_{abs}$  and  $Rank_{keyword}$  (L2R model using only formula layout feature group or formula text feature group can not outperform the baseline of  $Rank_{context}$ ). By combining multiple formula projecting feature groups, the ranking model  $L2R_{multiple}$  outperforms all baseline groups (including  $Rank_{context}$ ), which also proves L2R approach is an effective method for integrating features.

Third, while evolution relations among massive formulae are very useful for math-understanding, various kinds of information, i.e., formula layout, context, and generality, can be all useful for evolution relation discovery. For instance, though the comparison groups (based on partial FEM features) perform decently in the experiment,  $L2R_{all}$  (with comprehensive FEM mining features) is significantly superior ( $p < 0.0001$ ) than all other groups for almost all the evaluation metrics. This finding supports our initial hypothesis that evolution relation is a latent variable hiding behind formula context and layout, and we can hardly explore it by using a single kind of evidence.

### 4.3 Exit Survey

The goal of this proposed study is to design a novel algorithm/system to assist students to better understand the mathematics content in a paper. Although from an OER ranking viewpoint the experiment results are positive, we designed another exit survey to further proof the usefulness of the new reading environment and the effectiveness of the new formula-understanding method. In the survey, we asked each participant seven questions (at the end of the experiment), including “precision” (Q1), “satisfaction” (Q2), “usefulness” (Q3), “relevance” (Q4), “user-friendliness” (Q5), “usability” (Q6) and “effectivity” (Q7)<sup>2</sup>. Based on the students’ feedback, 72.73% of participants believed the new system along with the formula-understanding method can provide precise and useful information, and 43.75% find the proposed method can help them better understand the math-content in a paper “most of the time” (another 31.25% reported it is helpful “about half of the time”). From a system usability perspective, 78.85% of participants found the formula highlight function in PDF and OER recommendation functions were easy or very easy to use.

Overall, 63.63% of participants reported that the system and formula understanding method can be helpful or very helpful to assist them to better understand the target paper (only 12.12% reported the new functions are not helpful, and others were neutral). For the OER usefulness, 75.75% participants are satisfied with the quality of the recommended OERs (especially for videos and slides). Participants reported, comparing narrative content, the OERs are more helpful for math-understanding. It is clear that the proposed system and math information understanding method achieves the

<sup>2</sup>More detailed information can be found at <https://github.com/GraphEmbedding/FEM>



**Table 4: Measures of different OER ranking algorithms (Significant test: L2R<sub>all</sub> vs. other groups;  $\dagger p < 0.01$ ,  $\dagger\dagger p < 0.001$ ,  $\dagger\dagger\dagger p < 0.0001$ )**

Ranking	NDCG@3	NDCG@5	NDCG@all	P@3	P@5	MAP	MRR
<span style="color: red;">■</span> <i>Rank<sub>abs</sub></i>	0.5536	0.5860	0.7156	0.5544	0.4932	0.7309	0.7764
<span style="color: red;">■</span> <i>Rank<sub>keyword</sub></i>	0.5744	0.6004	0.7268	0.5635	0.4973	0.7344	0.7833
<span style="color: red;">■</span> <i>Rank<sub>context</sub></i>	0.6483	0.6642	0.7589	0.6293	0.5422	0.7689	0.8328
<span style="color: yellow;">■</span> <i>L2R<sub>layout</sub></i>	0.5771	0.6031	0.7253	0.5680	0.5000	0.7393	0.7861
<span style="color: yellow;">■</span> <i>L2R<sub>text</sub></i>	0.6408	0.6531	0.7567	0.6202	0.5361	0.7671	0.8509
<span style="color: yellow;">■</span> <i>L2R<sub>content</sub></i>	0.6571	0.6731	0.7651	0.6440	0.5551	0.7815	0.8497
<span style="color: yellow;">■</span> <i>L2R<sub>multiple</sub></i>	0.6798	0.6908	0.7777	0.6497	0.5612	0.7961	0.8717
<span style="color: blue;">■</span> <i>L2R<sub>all</sub></i>	<b>0.7150<sup>†††</sup></b>	<b>0.7252<sup>†††</sup></b>	<b>0.7969<sup>†††</sup></b>	<b>0.6712<sup>††</sup></b>	<b>0.5776<sup>††</sup></b>	<b>0.8171<sup>†††</sup></b>	<b>0.8848<sup>†</sup></b>
<span style="color: red;">■</span> Without FEM Assisted	<span style="color: yellow;">■</span> With Partial FEM Assisted		<span style="color: blue;">■</span> With Complete FEM Assisted				

goals, and the algorithms developed in this paper can be promising for scaffolding in education domain.

## 5 RELATED WORK

**Scientific Information Understanding:** Help students better understand and consume publications is an essential task in education and cyberlearning domains. Since 1976 [34], the term “scaffolding” has been widely used in educational research [24, 25]. In particular, the concept of scaffolding is applied to the studies of computer-assisted learning environments, also known as computer-mediated scaffolding [25]. One of the most recent efforts utilizes existing social tagging and annotation tools. Social tagging/annotation has produced positive results in a number of tasks, including the promotion of learning [9, 29, 33]. Prior studies, however, also found those scaffolding approaches, by leveraging social tagging, can be quite limited [15, 22], and students cannot essentially benefit from such systems when reading a challenging text. Researchers only recently began to focus on the usefulness of Open Educational Resources (OERs). For instance, Dennis, et al. [3] found that an additional video presentation had significant positive impacts on students’ learning, and more recently, Liu [13, 14] found ODRs, e.g., presentation videos/slides and Wikipedia pages, can help scholars better understand scientific readings. Meanwhile, more recent studies [8, 15] found that text and graph mining methods can be used to automatically recommend high quality OER to help students understand the paper text content. However, this approach cannot be applied to address the formula understanding problem.

**Scientific Topic Evolution:** Topic dynamics and evolution has been recently investigated. Laura Dietz et al. [4], for example, devised a probabilistic topic model that explains the generation of documents. The model incorporated topical innovation and topical inheritance via citations. Blei and Lafferty [2] proposed a Dynamic Topic Model (DTM), which explicitly characterized the chronological nature of sequential corpora by utilizing a Markov chain of term distributions over time. Based on [2], Gerrish and Blei proposed the Document Influence Model (DIM) [6]. This model respected the ordering of the documents and not only tracked how underlying theme has changed over time, but also captured how past articles exhibit varying influence on future articles. More recently, Jiang et al. [7] investigated topic evolution problem by integrating both text and citation data. Unlike earlier efforts, [7] generated a heterogeneous graph with various relations between topics and paper, i.e., citation and topic evolution, and supervised random walk was used for citation recommendation.

However, all the existing methods cannot be used to address scientific formula evolution for two reasons. First, one complex formula (in a paper) may implicitly associate with different kinds of topics, and these topics may not appear in the formula context (text information is not complete). Second, scholarly publication citation information may not be sufficient for formula understanding, e.g., “LDA” studies do not necessarily cite “Beta distribution” or “Bayesian inference” foundations (that can be important to help readers to understand the formula).

**Formula Search and Layout Mining:** Formula search is an important area in information retrieval. Recently, National institute of informatics Testbeds and Community for Information access Research (NTCIR) developed an evaluation collection for mathematical formula search with the aim of facilitating and encouraging research in formula search and its related fields [1]. For formula search, one of the main challenges should be formula information extraction, namely how to convert formulae into terms which were utilized to build the index. The same with text tokenizer, a text-based category of tokenizers was employed for formula retrieval and formula layout mining [19–21]. Different from plain text, formulae were highly structured. They can be expressed and parsed as tree structures. Thus, tree-based methods were the most important tokenization approach in recently proposed formula search systems [10, 27, 28, 31]. However, formula understanding and formula evolution mining are novel problems, and existing formula search methods cannot be directly applied.

## 6 CONCLUSION

In this study, we propose a novel problem-Mathematics Content Understanding-to assist readers to better understand and consume the math-content in scientific publications by leveraging Formula Evolution Map (FEM) and high-quality OERs. By using the PRMA cyberreading system, students/scholars can easily highlight a target formula in a PDF reading, and the proposed algorithms can project the query formula to the formula(e) vertex(es) in FEM as well as recommend OERs to users. In the offline process, we extract formula evolution relations from a massive Wikipedia dump. Evaluation shows that formula relations on FEM are fundamental and enlightened for math-understanding. Most of the experiment participants find the proposed method/system can effectively help them better understand the math-content and readings in a cyberreading environment. Meanwhile, students reported that the formula highlight function was easy to use and the recommended OERs can be very useful for their understanding of the math-content.

For the algorithm evaluation, we found that: first, the proposed OER recommendation via FEM mining is effective. It achieves the best performance for all evaluation metrics. Compared with the baselines without FEM assistance, the new model is significantly superior. Second, formula evaluation information is especially important for math-content understanding tasks, i.e., compared with the other comparison groups with partial FEM assistance, the proposed method with complete FEM information has an average increase of 8.2% for P@3 and an average increase of 11.9% for NDCG@3.

The methodological limitations of this work is that the parameter  $\theta$  of evolution probability calculation is treated equally without tuning. This is caused by two reasons. First, there is no existing formula evolution relations for training. Second, the number of scientific knowledge base documents and associated formulae is huge, and the time cost of direct parameter tuning can be very high. In the future, we will explore the combination of knowledge base and academic corpus in depth, and try to use more features (e.g., selected paper citation relation information) in the formula evolution mining. Meanwhile, we will propose more sophisticated optimization methods for FEM evolution probability parameter tuning.

## ACKNOWLEDGMENTS

The work is supported by the National Science Foundation of China (61472014), Guangdong Province Frontier and Key Technology Innovative Grant (2015B010110003, 2016B030307003) and the Opening Project of State Key Laboratory of Digital Publishing Technology.

## REFERENCES

- [1] Akiko Aizawa, Michael Kohlhase, Iadh Ounis, and Moritz Schubotz. 2014. NTCIR-11 Math-2 Task Overview. In *NTCIR*. Citeseer.
- [2] David M Blei and John D Lafferty. 2006. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*. ACM, 113–120.
- [3] Alan R Dennis, Kelly O McNamara, Stacy Morrone, and Joshua Plaskoff. 2015. Improving Learning with eTextbooks. In *Proceedings of the 48th Hawaii International Conference on System Sciences*. 5253–5259.
- [4] Laura Dietz, Steffen Bickel, and Tobias Scheffer. 2007. Unsupervised prediction of citation influences. In *Proceedings of the 24th international conference on Machine learning*. ACM, 233–240.
- [5] Liangcai Gao, Ke Yuan, Yuehan Wang, Zhuoren Jiang, and Zhi Tang. 2016. The math retrieval system of ICST for NTCIR-12 MathIR task. *Proc. NTCIR-12* (2016).
- [6] Sean Gerrish and David M Blei. 2010. A language-based approach to measuring scholarly impact. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*. 375–382.
- [7] Zhuoren Jiang, Xiaozhong Liu, and Liangcai Gao. 2015. Chronological Citation Recommendation with Information-Need Shifting. In *Proceedings of the 24th ACM International Conference on Information and Knowledge Management*. ACM, 1291–1300.
- [8] Zhuoren Jiang, Xiaozhong Liu, Liangcai Gao, and Zhi Tang. 2016. Community-based Cyberreading for Information Understanding. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. ACM, 789–792.
- [9] Tristan E Johnson, Thomas N Archibald, and Gershon Tenenbaum. 2010. Individual and team annotation effects on students' reading comprehension, critical thinking, and meta-cognitive skills. *Computers in human behavior* 26, 6 (2010), 1496–1507.
- [10] Michael Kohlhase and Ioan Sucan. 2006. A search engine for mathematical formulae. In *International Conference on Artificial Intelligence and Symbolic Computation*. Springer, 241–253.
- [11] Hang Li. 2014. Learning to rank for information retrieval and natural language processing. *Synthesis Lectures on Human Language Technologies* 7, 3 (2014), 1–121.
- [12] Xiaoyan Lin, Liangcai Gao, Xuan Hu, Zhi Tang, Yingnan Xiao, and Xiaozhong Liu. 2014. A mathematics retrieval system for formulae in layout presentations. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. ACM, 697–706.
- [13] Xiaozhong Liu. 2013. Generating metadata for cyberlearning resources through information retrieval and meta-search. *Journal of the American Society for Information Science and Technology* 64, 4 (2013), 771–786.
- [14] Xiaozhong Liu and Han Jia. 2013. Answering academic questions for education by recommending cyberlearning resources. *Journal of the American Society for Information Science and Technology* 64, 8 (2013), 1707–1722.
- [15] Xiaozhong Liu, Zhuoren Jiang, and Liangcai Gao. 2015. Scientific information understanding via open educational resources (OER). In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 645–654.
- [16] Xiaozhong Liu and Jian Qin. 2014. An interactive metadata model for structural, descriptive, and referential representation of scholarly output. *Journal of the Association for Information Science and Technology* 65, 5 (2014), 964–983.
- [17] Donald Metzler and W Bruce Croft. 2007. Linear feature-based models for information retrieval. *Information Retrieval* 10, 3 (2007), 257–274.
- [18] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119.
- [19] Bruce R Miller and Abdou Youssef. 2003. Technical aspects of the digital library of mathematical functions. *Annals of Mathematics and Artificial Intelligence* 38, 1-3 (2003), 121–136.
- [20] Robert Miner and Rajesh Munavalli. 2007. An approach to mathematical search through query formulation and data normalization. In *Towards Mechanized Mathematical Assistants*. Springer, 342–355.
- [21] Jozef Mišutka and Leo Galamboš. 2008. Extending full text search engine for mathematical content. *Towards Digital Mathematics Library*. Birmingham, United Kingdom, July 27th, 2008 (2008), 55–67.
- [22] Elena Novak, Rim Razzouk, and Tristan E Johnson. 2012. The educational use of social annotation tools in higher education: A literature review. *The Internet and Higher Education* 15, 1 (2012), 39–49.
- [23] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1999. *The PageRank citation ranking: Bringing order to the web*. Technical Report. Stanford InfoLab.
- [24] Roy D Pea. 2004. The social and technological dimensions of scaffolding and related theoretical concepts for learning, education, and human activity. *The journal of the learning sciences* 13, 3 (2004), 423–451.
- [25] Sadhana Puntambekar and Roland Hubscher. 2005. Tools for scaffolding students in a complex learning environment: What have we gained and what have we missed? *Educational psychologist* 40, 1 (2005), 1–12.
- [26] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D Manning. 2009. Labeled LDA: A supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1-Volume 1*. Association for Computational Linguistics, 248–256.
- [27] Thomas Schellenberg, Bo Yuan, and Richard Zanibbi. 2012. Layout-based substitution tree indexing and retrieval for mathematical expressions. In *IS&T/SPIE Electronic Imaging*. International Society for Optics and Photonics, 82970I–82970I.
- [28] Petr Sojka and Martin Liška. 2011. Indexing and searching mathematics in digital libraries. In *International Conference on Intelligent Computer Mathematics*. Springer, 228–243.
- [29] Addison Su, Stephen JH Yang, Wu-Yuin Hwang, and Jia Zhang. 2010. A Web 2.0-based collaborative annotation system for enhancing knowledge sharing in collaborative learning environments. *Computers & Education* 55, 2 (2010), 752–766.
- [30] Y. Sun, J. Han, X. Yan, P. S. Yu, and T. Wu. 2011. PathSim: Meta path-based top-k similarity search in heterogeneous information networks. In *Proc. 2011 Int. Conf. Very Large Data Bases (VLDB'11)*. Seattle, WA.
- [31] Yuehan Wang, Liangcai Gao, Simeng Wang, Zhi Tang, Xiaozhong Liu, and Ke Yuan. 2015. WikiMirs 3.0: a hybrid MIR system based on the context, structure and importance of formulae in a document. In *Proceedings of the 15th ACM/IEEE-CS Joint Conference on Digital Libraries*. ACM, 173–182.
- [32] Scott White and Padhraic Smyth. 2003. Algorithms for estimating relative importance in networks. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 266–275.
- [33] Joanna Wolfe. 2008. Annotations and the collaborative digital library: Effects of an aligned annotation interface on student argumentation and reading strategies. *International Journal of Computer-Supported Collaborative Learning* 3, 2 (2008), 141–164.
- [34] David Wood, Jerome S Bruner, and Gail Ross. 1976. The role of tutoring in problem solving\*. *Journal of child psychology and psychiatry* 17, 2 (1976), 89–100.
- [35] Chengxiang Zhai and John Lafferty. 2001. A study of smoothing methods for language models applied to ad hoc information retrieval. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 334–342.
- [36] Xiaode Zhang, Liangcai Gao, Ke Yuan, Runtao Liu, Zhuoren Jiang, and Zhi Tang. 2017. A Symbol Dominance Based Formulae Recognition Approach for PDF Documents. In *14th IAPR International Conference on Document Analysis and Recognition, ICDAR 2017, Kyoto, Japan, November 9-15, 2017*. 1144–1149.