

# Educational Forums at a Glance: Topic Extraction and Selection

Bernardo Pereira Nunes<sup>1</sup>, Ricardo Kawase<sup>2</sup>, Besnik Fetahu<sup>2</sup>,  
Marco A. Casanova<sup>1</sup>, and Gilda Helena B. de Campos<sup>3</sup>

<sup>1</sup> Department of Informatics – PUC-Rio – Rio de Janeiro, Brazil  
{bnunes, casanova}@inf.puc-rio.br

<sup>2</sup> L3S Research Center – Leibniz University Hannover – Hannover, Germany  
{kawase, fetahu}@L3S.de

<sup>3</sup> Department of Education – PUC-Rio – Rio de Janeiro, Brazil  
gilda@ccead.puc-rio.br

**Abstract.** Web forums play a key role in the process of knowledge creation, providing means for users to exchange ideas and to collaborate. However, educational forums, along several others online educational environments, often suffer from topic disruption. Since the contents are mainly produced by participants (in our case learners), one or few individuals might change the course of the discussions. Thus, realigning the discussed topics of a forum thread is a task often conducted by a tutor or moderator. In order to support learners and tutors to harmonically align forum discussions that are pertinent to a given lecture or course, in this paper, we present a method that combines semantic technologies and a statistical method to find and expose relevant topics to be discussed in online discussion forums. We surveyed the outcomes of our topic extraction and selection method with students, professors and university staff members. Results suggest the potential usability of the method and the potential applicability in real learning scenarios.

## 1 Introduction

Over the past decade, the World Wide Web became an important source of information and knowledge. The diversity and engagement of independent users and communities contributed to the creation and proliferation of a rich set of content available in different communication channels (such as social media, real-time channels, blogs, forums, etc.) as well as in formats (such as text, audio and video).

In particular, online discussion forums have played a key role in the process of knowledge creation [13], providing means for its users to exchange ideas, form opinions, position themselves and collaborate. As an outcome of the importance of online discussion forums is Wikipedia<sup>1</sup>, where for each Wikipedia article there is a forum-based page<sup>2</sup> that relies on the collaboration, discussion, consensus and collective effort of its users to keep Wikipedia constantly updated and curated.

---

<sup>1</sup> <http://www.wikipedia.org>

<sup>2</sup> [http://en.wikipedia.org/wiki/Help:Using\\_talk\\_pages](http://en.wikipedia.org/wiki/Help:Using_talk_pages)

Due to the benefits generated by users' participation in forums, most online courses combine educational materials and message boards. However, even though forums clearly leverage the creation of collective intelligence [19], the assessment of users' participation is still rather difficult [12,18]. Depending on the number of students and posts, manual assessment becomes impractical. Previous work addressed the problem of assessing the quality of students' participation [15,16]. However, they do not take into account whether a particular set of topics were addressed in a thread of a specific discipline.

Furthermore, different backgrounds in online discussion forums may lead a discussion to unforeseen directions, needing external support to realign the discussed topics of a thread. This task is often conducted by a tutor or moderator. But, as we will show in this paper, on average, 50% of forums discussing a specific subject with different audience or tutor/moderator cover distinct topics. This means even though online discussions are often different, a set of specific topics must be addressed to achieve the course goals. Therefore, if a given forum does not cover a set of expected topics, the assessment of the students might be hampered, since the acquired knowledge depends on the topics discussed in the forum.

In this paper, we combine semantic technologies and a statistical method to find, expose and recommend relevant topics as guidance to conduct debate forums. Briefly, with the help of semantic tools, the proposed method first performs Named-Entity Recognition (NER) and topic extraction, followed by a statistical approach that selects and ranks the most relevant topics of a forum thread. Finally, the method outputs the top-most representative topics discussed in a specific forum as well as a set of suggested topics to be discussed. We used 97 online forums from a Brazilian university to validate and assess our method.

Our main contribution in this work is the development of a well perceived semantic-based topic enrichment model for educational forums, in combination with its evaluation. Subsequently, this contribution accounts for positive effects in high-level assessment of tutor/moderator progress, topic recommendation and parity of knowledge acquisition by students in online forums.

The rest of this paper is organised as follows. Section 2 reviews related literature. Section 3 describes the use of forums in our context. Section 4 introduces the topic extraction and selection method. Moreover, we also extended Vygotsky's zone of proximal development to serve as a recommendation method. Section 5 presents the evaluation setup. Section 6 discusses the results obtained in the evaluation along with a brief analysis of the topics extracted from the forum threads. Finally, Section 7 discusses our outcomes and future work.

## 2 Related Work

Li and Wu [11] combine approaches involving sentiment analysis and text mining to detect hotspot forums within a certain time span. Their method assists users to make decisions and predictions over polarised groups of messages in online forums. Despite not performing topic extraction in the hotspot forums, the emotional polarity information for each topic extracted would help users on understanding how a given topic is addressed in a discussion.

Cong et al. [3] present an approach for finding question-answer pairs in online forums based on Labeled Sequential Patterns (LSPs) and graph-based propagation model. While the creation of patterns for interrogative sentences is made using part-of-speech tags, the answers are detected and ranked using KL-divergence language model. Again, our approach is complementary to their approach, since our approach would serve as a filter for finding question and answers based on topics. Conversely, our method would benefit of this approach by identifying key posts in an online discussion.

Online forums play a key role in the student skills development as shown by Scaffidi et al. [17]. Their study focuses on the types of posts that facilitate discussions and collaboration amongst novice developers. The study of user behavior in online forums help to promote active interaction amongst users and therefore the construction of collective knowledge. We believe that the introduction of new topics to be discussed by such community of users could trigger new discussions and hence new knowledge.

Desanctis et al. [6] provide an interesting discussion about e-venues for learning such as video-conferenced classrooms, online communities and group discussion spaces. Although each venue influences the learning process of a particular group, they all have in common the need to bring new discussions that promote the development of knowledge of the participants. For instance, online communities usually last more than private group discussion spaces, since new participants with fresh questions can drop in at any-time. Thus, in order to maintain the group discussion, the recommendation of new topics for discussion would foster longer interactions amongst participants and knowledge refreshment.

Evidently, online discussions can also be fueled by tutors responsible for bringing new topics and questionings for the discussion. Previous studies [4] have shown that tutored venues can improve both retention and performance of the participants. In this paper, we use the tool for assisting tutors on addressing new topics relevant to the discussion.

A highly relevant direction of work goes onto the topic extraction from forums' text. Hulpus et al. [10] extract a set of topics from a given textual resource. The topics correspond to a DBpedia sub-graph category. Furthermore, the relationship between the topics and textual resources are quantified using graph importance measures. In our case, we simply aim at providing students with discussions from the forums specific to a topic. Hence, simple tf-idf techniques offer us the efficiency on distinguishing the topic specific discussions. This simple, yet efficient technique offers the scalability over large corpora, and at the same we avoid exhaustive computations that rely on graph centrality measures like the ones in [10].

Finally, research on topic modelling and extraction has been addressed on various ways and for different purposes. A well known approach LDA [2] has seen a wide applicability on modelling and extracting topics from textual documents. Other approaches like [7] extract and rank topics with respect to their relevance to specific datasets, which are extracted through a named entity disambiguation process. In contrast to the previous approaches, in our case we aim at suggesting forum pages for discussion specific to online learning scenarios, hence, the problem becomes simpler with respect to filtering specific forum pages rather than exhaustive rankings of forum pages and their corresponding topics.

### 3 Motivation

To illustrate the motivation of our research, we describe two scenarios where participants of online discussion forums would benefit from our method. Both scenarios result from the need of the staff from a Brazilian university to assess the participations in forums and topics discussed.

Online discussion forums are fundamental in the learning process and most of the online courses take advantage of their use to meet specific goals. Assessing student participation in forums is not a simple task, and due to the high number of posts, it can become impracticable. Hence, in order to maintain the quality of teaching and student experience, the university staff members required a tool to track the discussion progress.

The first scenario described by the university staff members is that tutors constantly overlook the discussion of relevant topics in favor of a better flow. Although the discussion flow is of utmost importance, tutors must conduct the forum in such a way that specific topics must be addressed and, at the same time, preserve the discussion flow. Hence, the university staff members are interested in the analysis of forums to check if particular topics were covered in a thread. By doing this, they can ensure that all participants had similar experience and learning situations that can contribute to the next activities. In the case that a set of topics are not covered, they would like to intervene and extend the forum closure or create a new forum thread to discuss the missing topics.

The second scenario aims at fostering the discussion with suggestions that may assist students in the discussion. For many reasons, some forums lack interaction and students must be encouraged to participate. In this manner, university staff members believe that a recommendation tool would promote the discussion and help to reach the forum's goal.

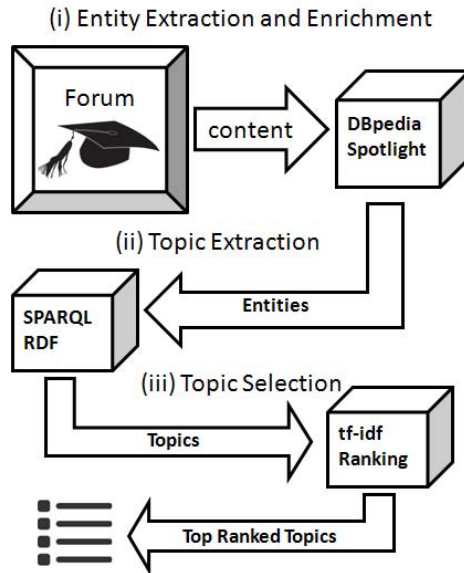
The current work assists university staff members and students to have a better overview of what is happening in the forum to take the right action and create situations/activities that can improve the learning experience of the students.

## 4 Topic Extraction and Selection

In this section we present the main steps for a coherent process chain that semantically and statistically selects the most relevant discussed topics in a given online discussion forum. The process chain, depicted in Figure 1 is composed of three steps described as follows: (i) Entity Extraction and Enrichment; (ii) Topic Extraction; and (iii) Topic Selection. We also present in this section a simple topic recommendation method used to assist learners and tutors in the teaching-learning process.

### 4.1 Entity Extraction and Enrichment

When dealing with online discussion forums, we are essentially working with unstructured data, which in turn hinders data manipulation and the identification of atomic elements in texts. To alleviate this problem, information extraction (IE) methods, such as Named-Entity Recognition (NER) and name resolution, are employed. These tools



**Fig. 1.** Topic extraction process workflow

automatically extract structured information from unstructured data and link to external knowledge bases in the Linked Open Data cloud (LOD), such as DBpedia<sup>3</sup>.

For instance, after processing the following sentence using an IE tool: “I agree with Barack Obama that the whole episode should be investigated.”, the entity “Barack Obama” is annotated and classified as *person* and linked to the DBpedia resource [http://dbpedia.org/resource/Barack\\_Obama](http://dbpedia.org/resource/Barack_Obama), where structured information about him is available.

We use the DBpedia Spotlight tool<sup>4</sup> to extract and enrich entities found in the posts within a forum thread. DBpedia Spotlight adds markups with semantic information surrounding atomic elements (entities) in the forum posts (as in [14]). These entities are the ones found in DBpedia dataset, and each one contains structured information extracted from Wikipedia[1].

Note that our method is language independent as long as we have a solid repository of entities (such as DBpedia or Freebase<sup>5</sup>) and a proper annotation tool (such as Spotlight, Alchemy<sup>6</sup> or WikipediaMiner<sup>7</sup>). However, the set of entities that can be identified by the annotation process is limited to the number known entities in the dataset, in our case, the Portuguese DBpedia dataset. This dataset currently contains 736,443 entities<sup>8</sup>.

<sup>3</sup> <http://www.dbpedia.org>

<sup>4</sup> <http://dbpedia-spotlight.github.io/demo/>

<sup>5</sup> <http://www.freebase.com>

<sup>6</sup> <http://www.alchemyapi.com>

<sup>7</sup> <http://wikipedia-miner.cms.waikato.ac.nz>

<sup>8</sup> <http://wiki.dbpedia.org/Datasets39/DatasetStatistics>

## 4.2 Topic Extraction

Given as starting point the entities that were found in the previous step, the topic extraction step begins by traversing the entity relationships to find a more general representation of the entity, i.e., the topics.

An entity is conventionally represented as a RDF (Resource Description Framework) triple in the form of (Subject, Predicate, Object), where each triple represents a fact, and the predicate names the relationships between the subject and the object. For example, a triple is (“Barack Obama”, “isPresidentOf”, “United States of America”). Furthermore, a set of RDF triples form a directed and labeled graph, where the nodes are a set of subjects and objects and the edges are represented by the predicate.

Thus, for each extracted and enriched entity in the posts, we explore their relationships through the predicate *dcterms:subject*, which by definition<sup>9</sup> represents the topic of the entity. In that sense, to retrieve the topics, we use SPARQL query language for RDF over the DBpedia SPARQL endpoint<sup>10</sup>, where we navigate up in the DBpedia hierarchy to retrieve broader semantic relations between the entities and its topics. As it is shown in the following SPARQL query, we use the predicate *skos:broader*.

```
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
```

```
SELECT DISTINCT ?l1 ?l2 ?l3 ?l4
WHERE {
  <entity_uri> dcterms:subject ?l1 .
  ?l1 skos:broader ?l2 .
  ?l2 skos:broader ?l3 .
  ?l3 skos:broader ?l4
} LIMIT 1000;
```

The variable *entity\_uri* represents the entity in which we are interested in retrieving the topics extracted from the posts in a forum thread, while the variables *l1* to *l4* represent the topics that will be retrieved from the entity. Thus, given an entity, the topics of an entity are retrieved through the predicate *dcterms:subject* and *skos:broader*. The latter predicate is used to obtain a more general representation of the topic. This strategy will help us find the topics that best cover a forum thread.

Note that an entity/concept can be found in different levels of the hierarchical categories of DBpedia, and hence this approach would lead us to retrieve topics in different category levels. However, as in previous works [9,8], we take advantage of the co-occurrence of the topics in the different levels to find the most representative ones (see Section 4.3).

## 4.3 Topic Selection

Finally, in this last step, we select the most representative topics extracted from the posts that belong to a forum thread. For this, we rely on *tf-idf* (term frequency - inverse

<sup>9</sup> <http://dublincore.org/documents/2012/06/14/dcmi-terms/?v=terms#elements-subject>

<sup>10</sup> <http://pt.dbpedia.org/sparql> - DBpedia SPARQL endpoint in portuguese.

document frequency) score to statistically measure the importance of a topic in a forum thread.

Typically, *tf-idf* is used on information retrieval and text mining to measure the importance of a word to a document in a collection. However, in this paper, we adapted this metric to take into account entities and topics extracted from the posts instead of words.

To select the most representative topics, we compute *tf-idf* score twice, one for the entities extracted from the forum thread (i.e. the most representative entities in the collection) and another for the topics extracted from the entities (see Section 4.2).

Basically, to compute the term frequency (*tf*), we count the number of occurrences of an entity  $e$  in a post  $p \in P$ . As for the inverse document frequency (*idf*), we compute the (*idf*) score by dividing the total number of posts  $|P|$  by the number of posts containing the entity  $|P_e|$ , see Eq. 1.

$$tfidf(e, p, P) = tf(e, p) \times idf(e, P) \quad (1)$$

where *tf* is the raw frequency of a term in a post, and *idf* is the measure of commonness/rareness of an entity in a collection  $P$ . *tf* and *idf* can be computed by the equations 2 and 3, respectively.

$$tf(e, p) = frequency(e, p) \quad (2)$$

$$idf(e, P) = \log\left(\frac{|P|}{|P_e|}\right) \quad (3)$$

After computing the *tf-idf* score for each entity, the topmost representative entities are selected. From the selected entities, the topics are extracted according to the process described in Section 4.2.

With the topics in hands, we then compute the *tf-idf* score over the topics extracted from the entities and decreasingly rank them. Again, the topmost representative topics for a given forum thread are selected. Note that the number of topics that represent a forum is chosen by the user (in our case, the top 10 relevant topics). Finally, the top ranked topics are selected to represent the forum thread topics.

#### 4.4 Topic Recommendation

Another contribution of this paper lies in the recommendation of topics based on the Zone of Proximal Development (ZPD) introduced by Vygotsky [20]. Briefly, this concept of Vygotsky describes the distance between the independent performance of an individual to perform a certain task and the performance of the individual when assisted by more capable peers. Thus, as the ZPD concept suggests, the assistance of an external peer may improve the learners' skills.

In this paper we extended the ZPD concept to perform as a topic recommendation tool to learners participating in educational forums. As in our context forum threads

occur simultaneously, we consider a sibling-forum<sup>11</sup> as the more capable peer. Hence, the topic recommendation is based on the topics discussed in the sibling-forums.

Following the technique presented in the previous sections, the topmost representative topics discussed in a sibling-forum that are missing in the actual forum thread are recommended as topic seeds to foster the discussion. Although simple, the recommendation assists learners and tutors on addressing topics overlooked in the current thread and to broaden the discussion to topics that they would not address without the indirect assistance of their peers. We would like to emphasize that tutors and learners can always opt to accept or not the recommendation.

## 5 Evaluation Setup

In this section, we present the evaluations performed to validate the applicability of our method in real scenarios. The first evaluation consists of a questionnaire to university staff and participants of the forums. The second evaluation consists of expert manual assessment of the generated topics performed by two educators.

### 5.1 Technology Acceptance Model Evaluation

Over the course of our study, real data from online discussion forums were used to perform a comprehensive evaluation of our method. It was evaluated using 97 online discussion forums containing in total 10,785 anonymised posts provided by the distance education department of a Brazilian university. All selected threads occurred at least twice concurrently. Furthermore, each professor assessed the suggested topics from forums conducted by themselves.

Our main objectives included a thorough assessment of the recommendation of topics based on previous online discussion forums as well as the assessment of the selected topics that cover a forum discussion. For this, we submitted 3 questionnaires to 11 students, 4 professors and 3 coordinators of the distance education department to gather different perspectives and views of the proposed method.

The questionnaires were divided into three different categories of questions, namely *perceived usefulness*, *perceived ease-of-use* and additional suggestions. Basically, the questions followed the Technology Acceptance Model (TAM) proposed by Davis [5], arguably the most influential “Technology Acceptance Theory”.

Briefly, this theory states that there are two key aspects to measure users’ intention to adopt a new technology, the *perceived usefulness* and *perceived ease-of-use*. Perceived usefulness (PU) refers to “the degree to which a person believes that using a particular system would enhance his or her job performance”, while perceived ease of use (PEOU) refers to “the degree to which a person believes that using a particular system would be free of effort” [5].

Each questionnaire was divided in 6 PU questions, 6 PEOU questions and additional 3 opinion mining questions where we asked participants for further suggestions. Feedback for assertions such as ‘*The evaluation performed by teachers can be facilitated.*’,

---

<sup>11</sup> Two or more forums are considered as sibling-forums if they address the same subject, occur simultaneously and have different tutors and learners.



*'The tool can broaden the discussion.'*, and *'Suggested topics are relevant to the topics discussed.'* were collected in a 5-point Likert scale fashion.

Note that, in the case of university staff members, the topics were assessed over two randomly chosen forum threads, since they did not participate on the forum. Thus, a list of topics discussed in the forum and a list of suggested topics for each forum thread was available for their evaluation. As they are staff members of the university, they also have access to the forum discussions in case they would need additional information.

## 5.2 Expert Assessment

In parallel to the evaluation presented in Section 5.1, we recruited two experts in the distance education department. These experts are part of the senior university staff and are directly involved in research and in the management of online courses. The main objective of this evaluation is to have a first look on the performance of the proposed method in terms of precision. In practice, two distinct aspects were evaluated by the experts. First, (i) they evaluated the correctness of the assigned topics for a given forum thread. Second, (ii) they evaluated the topics that were recommended to a given forum thread based on previous forum threads.

For this evaluation, we randomly selected 22 forum threads from our corpus to be inspected by the two experts. In total, both experts read all 2070 posts corresponding to 19,1% of the total number of posts of our corpus.

**(i) First expert assessment: topic assignment.** After reading through all the posts of the randomly selected forum threads, the experts were presented with the top 10 topics that were automatically identified by our method. Next, the experts were asked to mark which topics were correct, and which were incorrect assigned to the forum (precision).

**(ii) Second expert assessment: topic recommendation.** Again, after reading through each randomly selected forum thread, the experts were presented with topic recommendations. We recall that the topic recommendation was based on the Zone of Proximal Development introduced by Vygotsky (see Section 4.4). We used the topics discussed in sibling-forums threads to recommend topics to the current thread. We considered the sibling-forum thread as the more capable peer and recommended the missing topics to the current discussion in order to assist the teaching-learning process. Similar to the previous assessment, the experts were asked to mark which recommended topics might be relevant for the discussion. Additionally, they were asked how important (in a 5-point Likert scale) the recommended topics were.

## 6 Results

Following the evaluation strategy presented in the previous section, we first present the results obtained by the TAM model and subsequently the results obtained by the evaluation performed with distance learning experts.

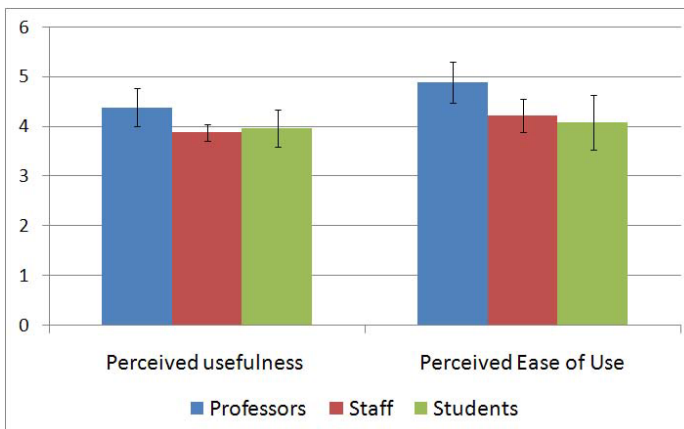
## 6.1 Technology Acceptance Model Results

The results of the questionnaires are summarised in Figure 2. The error bar charts show that all participants reported a high positive perception for the proposed topics, the implications and the applicability of the results. In particular, professors had a slightly better acceptance, when compared to the other tiers of participants. The coefficient of internal consistency Cronbach's  $\alpha$  of 0.65 for PU and 0.72 for PEOU indicated a good reliability of the results. These questionnaire results suggest the potential usability of our proposed topic extraction and selection method.

Regarding the suggestions included in the questionnaires, we observe that the most controversial question referred to whether or not the recommended topics should be available for professors, students or both. All professors suggested that the topics should be available only to them. All staff members suggested that topics should be available to both. Interestingly, students did not come to a common agreement. While the majority (64%) agreed that suggested topics should be available to professors and students, 36% opined that topic suggestions should be available only to professors.

We believe that the controversy is raised by the different backgrounds each group of participants had and the understanding they had of the topics. Staff members, who are not effectively involved in the online forums, assumed that the discussed topics should come out from an agreement between professors and students. On the other hand, the opinion of professors that a tool should present topics recommendation directly to them in fact reflects their need to control those around them. Finally, the split students' opinions lie in the fact that some students are still skeptical that online educational forums can smoothly evolve without proper moderation.

Unlike the questionnaires given to students and staff members, professors' questionnaire had an additional question regarding whether other professors can benefit from the suggested topics. The results reported that 75% of the professors strongly agree that *other* professors would take advantage of the suggested topics.



**Fig. 2.** Error Bars for survey questions regarding perceived usefulness and perceived ease-of-use

Finally, all staff members and professors (strongly) agree that the assessment of students would be facilitated if disparate forums addressed the same topics. Likewise, all staff members (strongly) agree that the proposed method would help in the assessment of the professor regarding the coverage of topics addressed in the forums. Nevertheless, 88% of all participants agree that the use of such method should be optional, and therefore, preserve the independence of tutors and learners.

## 6.2 Assessment Results by Distance Learning Experts

Table 1 depicts the results of the experts' assessment. The table expose individual (for each expert) and combined results. The combined results required the positive matching of both experts' opinion. Thus, for given a topic to be classified as *correct*, both experts must agree. If one of the experts marked as *incorrect*, the topic is automatically classified as incorrect.

In the (i) topic assignment assessment, each discussion forum was assigned with the 10 best ranked topics, thus the results are presented as precision (P@1, P@5 and P@10). In the (ii) topic recommendation assessment, since the recommendations were originated from sibling-forums, not all of them received the same number of recommendations. In average, each forum received 4.95 topic recommendations ( $\sigma = 2.68$ ). This means that the number of suggested topics is equivalent to, in average, 50% of the topics discussed in sibling-forums. This result demonstrates that sibling-forums being conducted by different tutors and having different learners can take different directions. Thus, the topic recommendation method has shown to be extremely important to assist tutors in the guidance of the forum thread and to align the topics being coverage in sibling-forums.

The results given by the experts' assessment show a high precision achieved by our proposed method. It reaches close to 100% precision for the top 1 topic assignment and impressive results above 82% for top 5 and top 10 topics. For the harder task of topic recommendation, we also observe quality results with average precision above 73%. These results reinforce the findings from the Technology Acceptance Model evaluation, reaffirming the benefits of our topic extraction and recommendation method.

**Table 1.** Expert assessment results

	Topic Assignment			Topic Recommendation
	(P@1)	(P@5)	(P@10)	(Avg. Precision)
Expert 1	100	94.63	90.97	83.90
Expert 2	97.56	90.24	86.34	79.51
Combined	97.56	87.80	82.19	73.17

## 6.3 Discussion Evolution

It is noteworthy that by applying our proposed method, we are able to post-identify the top topics of a given forum discussion. However, for effective topic guidance support,

it is important that the top topics are identified on the fly, during the progress of the discussion. To understand the convergence of the automatic identified topics, we incrementally generated topics for 88 forum discussions. The selected forums had at least 50 posts each.

We considered for this experiment a 10-post step granularity, i.e. after every 10 new posts in a discussion forum, we re-generated the list of top 10 topics and compare with the previous list. We used the overlap of topics in the lists (precision) as a metric for comparison. We defined a convergence of topics if the overlap between the lists is equal or greater than 90%.

Out of the 88 forum discussions that were used in this evaluation, in only 10 cases (11.4%) we observed that the identified topics diverge after converging above the threshold.

In average, the topics converge after 37.9% ( $\sigma = 26.7$ ) of the posts in a discussion forum. In practice, 52.3% of the discussion forums have the assigned topics converging after 20 posts, and 79.5% after 30 posts.

From this analysis we infer that, with 30 posts as input, the method can provide descriptive topics with descent performance. This result is important for the setup of the method and deployment in real scenarios. The topic recommendation also fosters new discussions and open new directions in the discussion. Once again, we recall to ZPD concept to show that with external assistance the discussion can become richer.

## 7 Discussion and Outlook

We presented a method for automatically generating topics that represent a forum thread in distance learning environments. We combined semantic and statistical techniques in a coherent process chain to extract, select and rank the most relevant topics of a forum. Moreover, we also introduced a simple topic recommendation method based on Vygotsky's educational theory.

Our experiments showed that most professors, university staff and students are willing to use our proposed approach in future forums. Moreover, 75% of the professors reported that other professors would benefit from the suggested topics.

Reviewing a sample of 97 forum threads, we verified that, on the average, 50% of the topics discussed in disparate forums addressing the same subject are different. This situation resulted in a concern with regard to the topics addressed in the forums and the post assessment of the students. A priori, students in disparate forums covering the same subject should have a similar experience and learn the same topics.

Thus, providing a method to overview the topics discussed in different forums will help university staff members, such as course coordinators, to rapidly intervene in forums that topics are being overlooked. The topic recommendation method has proven useful for the alignment and diversification of the discussions between sibling-forums. As reported in Section 6.3 52.3% of the topics discussed in forum threads converge after 20 posts, and 79.5% after 30 posts. Thus, if the suggested topics are taken into consideration by tutors and learners during the discussion, it may last longer, active and cover the expected topics for the discussion.

In theory, the use of the proposed method would bring more control of what is being taught in a forum and, therefore, ensure quality. In practice, this can be different and

some considerations arose out of the purpose of using the proposed method by a few interviewed respondent.

The first consideration lies in the freedom of the professors in guiding the forums. As every professor has its own teaching style and may also have a different point-of-view when they approach a subject, the concern of having to address specific topics in a forum might decrease the creativity and engagement of some professors. On the other hand, (assistant) professors may also take advantage of the suggested topics to guide the forum.

Another consideration with respect to the suggested topics is its availability to students. In the same time a topic suggestion may trigger an insight or make some students more confident, other students may stick only to the suggested topics. In the latter case, professors may take advantage of the students' participation and use it as a starting point to new discussions.

In general, the proposed method aims at assisting university staff members, professors and students to have a better overview of what is being discussed in the forum and, therefore, enable professors to take more informed actions to preserve discussion flow, improve students' experience and ensure topic coverage.

Our method also provides to the university staff members the possibility of assessing forum coverage, tracking what students are learning in different forums and, in some cases, detecting deviations in the forums. Adopting the method, depends on the instructional design of the course. The set-up of the course is crucial to determine which methods must be used and who will use it (professors, students or both).

As for future work, we plan to expand the method to accept external topic suggestions. For instance, professors involved in the course can also add topics to the discussion. Furthermore, we also plan to create a Moodle plugin.

## References

1. Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: Dbpedia: A nucleus for a web of open data. In: Aberer, K., Choi, K.-S., Noy, N., Allemang, D., Lee, K.-I., Nixon, L.J.B., Golbeck, J., Mika, P., Maynard, D., Mizoguchi, R., Schreiber, G., Cudré-Mauroux, P. (eds.) *ASWC 2007 and ISWC 2007*. LNCS, vol. 4825, pp. 722–735. Springer, Heidelberg (2007)
2. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. In: Dietterich, T.G., Becker, S., Ghahramani, Z. (eds.) *NIPS*, pp. 601–608. MIT Press (2001)
3. Cong, G., Wang, L., Lin, C.-Y., Song, Y.-I., Sun, Y.: Finding question-answer pairs from online forums. In: *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2008*, pp. 467–474. ACM, New York (2008)
4. Cottam, J.A., Menzel, S., Greenblatt, J.: Tutoring for retention. In: *Proceedings of the 42nd ACM Technical Symposium on Computer Science Education, SIGCSE 2011*, pp. 213–218. ACM, New York (2011)
5. Davis, F.D.: Perceived usefulness, perceived ease of use, and user acceptance of information technology. *MIS Quarterly*, 319–340 (1989)
6. DeSanctis, G., Fayard, A.-L., Roach, M., Jiang, L.: Learning in online forums. *European Management Journal* 21(5), 565–577 (2003)

7. Fetahu, B., Dietze, S., Pereira Nunes, B., Antonio Casanova, M., Taibi, D., Nejdil, W.: A scalable approach for efficiently generating structured dataset topic profiles. In: Presutti, V., d'Amato, C., Gandon, F., d'Aquin, M., Staab, S., Tordai, A. (eds.) *ESWC 2014. LNCS*, vol. 8465, pp. 519–534. Springer, Heidelberg (2014)
8. Fetahu, B., Dietze, S., Pereira Nunes, B., Antonio Casanova, M., Taibi, D., Nejdil, W.: A scalable approach for efficiently generating structured dataset topic profiles. In: Presutti, V., d'Amato, C., Gandon, F., d'Aquin, M., Staab, S., Tordai, A. (eds.) *ESWC 2014. LNCS*, vol. 8465, pp. 519–534. Springer, Heidelberg (2014)
9. Fetahu, B., Dietze, S., Nunes, B.P., Taibi, D., Casanova, M.A.: Generating structured profiles of linked data graphs. In: Blomqvist, E., Groza, T. (eds.) *International Semantic Web Conference. CEUR Workshop Proceedings*, vol. 1035, pp. 113–116. CEURWS.org (2013)
10. Hulpus, I., Hayes, C., Karnstedt, M., Greene, D.: Unsupervised graph-based topic labelling using dbpedia. In: Leonardi, S., Panconesi, A., Ferragina, P., Gionis, A. (eds.) *WSDM*, pp. 465–474. ACM (2013)
11. Li, N., Wu, D.D.: Using text mining and sentiment analysis for online forums hotspot detection and forecast. *Decision Support Systems* 48(2), 354–368 (2010)
12. Mazzolini, M., Maddison, S.: Sage, guide or ghost? the effect of instructor intervention on student participation in online discussion forums. *Computers Education* 40(3), 237–253 (2003)
13. Nonaka, I.: A dynamic theory of organizational knowledge creation. *Organization Science* 5(1), 14–37 (1994)
14. Pereira Nunes, B., Mera, A., Casanova, M.A., Kawase, R.: Boosting retrieval of digital spoken content. In: Graña, M., Toro, C., Howlett, R.J., Jain, L.C. (eds.) *KES 2012. LNCS*, vol. 7828, pp. 153–162. Springer, Heidelberg (2013)
15. Pendergast, M.: An analysis tool for the assessment of student participation and implementation dynamics in online discussion forums. *SIGITE Newsl.* 3(2), 10–17 (2006)
16. Romero, C., López, M.-I., Luna, J.-M., Ventura, S.: Predicting students' final performance from participation in on-line discussion forums. *Comput. Educ.* 68, 458–472 (2013)
17. Scaffidi, C., Dahotre, A., Zhang, Y.: How well do online forums facilitate discussion and collaboration among novice animation programmers? In: King, L.A.S., Musicant, D.R., Camp, T., Tymann, P.T. (eds.) *SIGCSE*, pp. 191–196. ACM (2012)
18. Shaul, M.: Assessing online discussion forum participation. *IJICTE* 3(3), 39–46 (2007)
19. Veerman, A.L., Andriessen, J.E.B., Kanselaar, G.: Collaborative learning through computer-mediated argumentation. In: *Proceedings of the 1999 Conference on Computer Support for Collaborative Learning, CSCL 1999. International Society of the Learning Sciences* (1999)
20. Vygotsky, L.: *Mind in society. The development of higher psychological processes*. Harvard University Press, Cambridge; edited by cole, michael et al. edition, 0 1978 / 1930