

4-2018

# EnTagRec(++): An enhanced tag recommendation system for software information sites

Shawei WANG

David LO

Singapore Management University, davidlo@smu.edu.sg

Bogdan VASILESCU

Alexander SEREBRENIK

**DOI:** <https://doi.org/10.1007/s10664-017-9533-1>

Follow this and additional works at: [https://ink.library.smu.edu.sg/sis\\_research](https://ink.library.smu.edu.sg/sis_research)



Part of the [Computer and Systems Architecture Commons](#), and the [Software Engineering Commons](#)

---

## Citation

WANG, Shawei; LO, David; VASILESCU, Bogdan; and SEREBRENIK, Alexander. EnTagRec(++): An enhanced tag recommendation system for software information sites. (2018). *Empirical Software Engineering*. 23, (2), 800-832. Research Collection School Of Information Systems.

**Available at:** [https://ink.library.smu.edu.sg/sis\\_research/4127](https://ink.library.smu.edu.sg/sis_research/4127)

# ENTAGREC<sup>++</sup>: An enhanced tag recommendation system for software information sites

Shaowei Wang<sup>1</sup>  · David Lo<sup>2</sup> · Bogdan Vasilescu<sup>3</sup> · Alexander Serebrenik<sup>4</sup>

Published online: 21 July 2017  
© Springer Science+Business Media, LLC 2017

**Abstract** Software engineers share experiences with modern technologies using software information sites, such as Stack Overflow. These sites allow developers to label posted content, referred to as software objects, with short descriptions, known as tags. Tags help to improve the organization of questions and simplify the browsing of questions for users. However, tags assigned to objects tend to be noisy and some objects are not well tagged. For instance, 14.7% of the questions that were posted in 2015 on Stack Overflow needed tag re-editing after the initial assignment. To improve the quality of tags in software information sites, we propose ENTAGREC<sup>++</sup>, which is an advanced version of our prior work ENTAGREC. Different from ENTAGREC, ENTAGREC<sup>++</sup> does not only integrate the historical tag assignments to software objects, but also leverages the information of users, and an initial set of tags that a user may provide for tag recommendation. We evaluate its performance on five software information sites, STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT,

---

Communicated by: Romain Robbes

---

✉ Shaowei Wang  
shaowei@cs.queensu.ca

David Lo  
davidlo@smu.edu.sg

Bogdan Vasilescu  
vasilescu@cmu.edu

Alexander Serebrenik  
a.serebrenik@tue.nl

<sup>1</sup> SAIL, Queen's University, Kingston, Canada

<sup>2</sup> School of Information Systems, Singapore Management University, Singapore, Singapore

<sup>3</sup> School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA

<sup>4</sup> Department of Mathematics and Computer Science, Eindhoven University of Technology, Eindhoven, The Netherlands

SUPER USER, and FREECODE. We observe that even without considering an initial set of tags that a user provides, it achieves *Recall@5* scores of 0.821, 0.822, 0.891, 0.818 and 0.651, and *Recall@10* scores of 0.873, 0.886, 0.956, 0.887 and 0.761, on STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, SUPER USER, and FREECODE, respectively. In terms of *Recall@5* and *Recall@10*, averaging across the 5 datasets, it improves upon Tag-Combine, which is the prior state-of-the-art approach, by 29.3% and 14.5% respectively. Moreover, the performance of our approach is further boosted if users provide some initial tags that our approach can leverage to infer additional tags: when an initial set of tags is given, *Recall@5* is improved by 10%.

**Keywords** Software information sites · Recommendation systems · Tagging

## 1 Introduction

The growing online media has significantly changed the way people communicate, collaborate, and share information with one another (Vasilescu et al. 2014). This is also true for software developers, who create and maintain software by standing on the shoulders of others (Storey et al. 2010), reuse components and libraries originating from Open Source repositories (e.g., GITHUB, FREECODE, SOURCEFORGE), and forage online for information that will help them in their tasks (Brandt et al. 2009). When foraging for information, developers often turn to programming question and answer (Q&A) communities such as STACK OVERFLOW, ASK UBUNTU, and ASK DIFFERENT. Such sites supporting communication, collaboration, and information sharing among developers are known as *software information sites*, while their contents (e.g., questions and answers, project descriptions)—as *software objects* (Xia et al. 2013).

Typically, tags are short labels not more than a few words long, provided as metadata to software objects in software information sites. Users can attach tags to various software objects, effectively linking them and creating topic-related structure. Tags are therefore useful for providing a soft categorization of the software objects and facilitating search for relevant information. To accommodate new content, most software information sites allow users to create tags freely. However, this freedom comes at a cost, as tags can be idiosyncratic due to users’ personal terminology (Golder and Huberman 2006). As tagging is inherently a distributed and uncoordinated process, often similar objects are tagged differently (Xia et al. 2013). Idiosyncrasy reduces the usefulness of tags, since related objects are not linked together by a common tag and relevant information becomes more difficult to retrieve. Furthermore, some software information sites (e.g., STACK OVERFLOW) require users to add tags at the time of posting a question, even if they are unfamiliar with the tags in circulation at that time. Due to differences in personal terminology and tagging purpose, it is often difficult for users to select appropriate tags for their content. Having a tag recommendation system that can suggest tags to a new object (e.g., based on how other similar objects have been tagged in the past) could (i) help users select appropriate tags easily and quickly, and (ii) in time help homogenize the entire collection of tags such that similar objects are linked together by common tags more frequently.

To illustrate the importance of tags for the well functioning of a software information site, we note the considerable amount of discussion related to tags on META STACK OVERFLOW, a Q&A site with the same user interface as STACK OVERFLOW, that focuses STACK OVERFLOW’s functioning and administration: e.g., at the time of writing there were more

than 4,587 questions related to tags,<sup>1</sup> as opposed to only 1,312 related to user interface. Furthermore, tags on STACK EXCHANGE sites receive considerable attention from the user community, with between 14.4% and 22.5% of questions in our experiments involving tag re-editing (see Table 1). Finally, we also note that since the earlier, conference version of our work (Wang et al. 2014), STACK OVERFLOW has been experimenting with, and gradually phasing in across the STACK EXCHANGE network, a tag recommendation system of their own.<sup>2</sup>

In this work, we introduce an automatic tag recommendation system called ENTAGREC<sup>++</sup>, an enhanced version of our previous approach ENTAGREC (Wang et al. 2014). ENTAGREC learns from historical software objects and their tags, and recommends appropriate tags for new objects based on words that appear in the software objects. ENTAGREC consists of two inference components, Bayesian and frequentist, and tries to combine the advantages of the two opposite yet complementary lines of thought in the statistics community (Samaniego 2010).

To improve ENTAGREC, in ENTAGREC<sup>++</sup>, we integrate two additional components into ENTAGREC: User Information Component (UIC) and Additional Tag Component (ATC). We refer to ENTAGREC integrated with UIC alone as ENTAGREC<sup>+</sup>. The intuition behind these two components is as follows:

- Users in software information sites tend to exhibit particular interests, thus software objects posted by them are likely to focus on specific domains. In UIC, we leverage this intuition to improve tag recommendation. We first link historical software objects posted by the same user together. Next, for new software objects posted by the same user, we make use of software objects that the user has posted before, to help identify tags that are associated with the new object.
- We believe that it may be easier for a user to assign one or a few initial tags to a question he/she posts, but more difficult for her to provide a comprehensive set of tags. In ATC, we make use of an initial set of tags provided by a user to help identify additional relevant tags.

We evaluate ENTAGREC<sup>+</sup> on datasets from five popular software information sites, STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, SUPER USER, and FREECODE, by comparing it to TAGCOMBINE (Xia et al. 2013; Wang et al. 2015).<sup>3</sup> Our experimental results show that even without considering an initial set of tags that a user provides, our approach achieves *Recall@5* scores of 0.821, 0.822, 0.891, 0.818, and 0.651, and *Recall@10* scores of 0.873, 0.886, 0.956, 0.887, 0.761 on STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, SUPER USER, and FREECODE, respectively. Compared with TAGCOMBINE, ENTAGREC<sup>+</sup> improves TAGCOMBINE by 29.3% and 14.5% in terms of *Recall@5* and *Recall@10*, respectively. Furthermore, to evaluate the effectiveness of ATC, we compare ENTAGREC<sup>+</sup> with ENTAGREC<sup>++</sup>. We find that when an initial set of tags is given, on average ENTAGREC<sup>++</sup> improves ENTAGREC<sup>+</sup> by 10.0% in terms of *Recall@5*.

Our main contributions are:

---

<sup>1</sup><http://meta.stackexchange.com/questions/tagged/tags>

<sup>2</sup><http://meta.stackexchange.com/questions/206907/how-are-suggested-tags-chosen>

<sup>3</sup>Since the implementation of STACK OVERFLOW's proprietary system is, to the best of our knowledge, not documented publicly, a meaningful comparison was not possible.

- We propose ENTAGREC<sup>++</sup>, a novel automatic tag recommendation system for software information sites. ENTAGREC<sup>++</sup> composes a state-of-the-art Bayesian inference technique (labeled LDA), an enhanced frequentist inference technique that leverages a POS tagger and the spreading activation algorithm, and two other components that analyze the user who posts a software object and the initial set of tags that the user provides, to further boost the recommendation performance.
- We evaluate our proposed approach on datasets from five popular software information sites. Our study shows that our approach can achieve high recall, especially for STACK OVERFLOW, ASK UBUNTU, SUPER USER, and ASK DIFFERENT, and outperforms a prior state-of-the-art approach.

The rest of this article is organized as follows. We provide more background on tags in several software information sites and approaches to tag recommendation in Section 2. We present the high-level architecture of ENTAGREC<sup>++</sup> in Section 3, followed by detailed descriptions of the Bayesian, frequentist, user information, and additional tag inference components in Sections 4, 5, 6 and 7 respectively, and the specifics of how to integrate the four components in Section 8. We present our evaluation results in Section 9. Finally, we highlight related work in Section 10 and conclude in Section 11.

## 2 Preliminaries and Examples

In this section, we first describe some preliminary information on tags in software information sites. Then, we present some recent works on tag recommendation on software information sites. Finally, we show some motivating examples to illustrate why it is useful to consider incorporating user information and preliminary tags in the recommendation.

### 2.1 Tags in Software Information Sites

To facilitate navigation, search, and filtering, contents are marked with descriptive terms (Golder and Huberman 2006), known as tags; e.g., libraries associate books with authors' names and keywords, while scientific publishers require the authors to choose keywords themselves. In the digital world, tags can be used, e.g., to annotate weblog posts and shared links. Numerous software information sites employ tags, e.g., SOURCEFORGE<sup>4</sup> for code projects, Eclipse Marketplace<sup>5</sup> for plugins, Snipplr<sup>6</sup> for code fragments, and STACK OVERFLOW for questions.

An example STACK OVERFLOW question is presented in Fig. 1. The question pertains to the creation of an Eclipse plugin and it has two tags, representing the technical context of the question (`eclipse`) and a specific subject area (`eclipse-plugin`). Figure 2 shows the FREECODE description of Apache Ant: in addition to the textual description, two general tags are present, `Software Development` (describing the general domain of Apache Ant) and `Build Tools` (indicating a more specific functionality of Apache Ant, namely building Java programs).

Comparing Figs. 1 and 2 we observe that while the basic purpose of tagging—to facilitate navigation, search, and content filtering through the association of related contents via

---

<sup>4</sup><http://sourceforge.net/>

<sup>5</sup><http://marketplace.eclipse.org/>

<sup>6</sup><http://snipplr.com/>

# How to create an Eclipse plugin

---

i need a complete tutorial about Eclipse plugin. My plugin has not a graphical interface, but i need to use his function inside another plugin or java app.  
I use eclipse ONLY to load this plugin, but must work in eclipse.  
It should be easy, but i don't know how to do this.

eclipse eclipse-plugin

**Fig. 1** An example question in STACK OVERFLOW and its tags

linked descriptive terms—is common to both, specific policies how the tags should be used differ from site to site. For instance, FREECODE has no restriction on the number of tags per project, while STACK EXCHANGE sites restrict the total number of tags given to a question to five.

Most software information sites allow users to provide “free text tags”. Not being subject to the formal requirements of the sites, such tags can be expected to represent user intent in a more flexible way. However, tagging becomes a distributed and uncoordinated process, introducing different tags for similar objects, which might persist despite moderation or the ongoing correction efforts. For example, questions on STACK OVERFLOW entitled “SIFT and SURF feature extraction implementation using MATLAB”<sup>7</sup> and “Matlab implementation of Haar feature extraction”<sup>8</sup> are both related to image feature extraction but only the second one is labeled with the corresponding tag, i.e., `feature-extraction`.

## 2.2 Tag Recommendation

Tags have been shown to aid users in navigating a site (Held et al. 2012; Cress et al. 2013; Bindelli et al. 2008; Zubiaga 2012). Thus, more complete tags can help in a number of scenarios, e.g.: (1) More complete tags may shorten the time it takes for a question to receive an answer. On sites such as STACK EXCHANGE, users are allowed to browse *unanswered questions* via tags. Giving more complete tags to a question may increase its chance to be discovered by a suitable user who can answer it well. (2) More complete tags also support developer learning. Developers can use the tags to browse through relevant questions and problems that others have encountered. This can help them avoid making similar mistakes and improve their programming and problem solving skills. (3) More complete tags may help reduce duplicated questions. Moderators can use the tags to identify related questions; by checking these related questions, moderators can decide whether a question is a duplicated one, and if so, it can be marked accordingly. Additionally, before posting new questions, users can use tags to browse for related ones, and avoid posting questions that were answered before.

Indeed, users of STACK EXCHANGE edited tags that were originally assigned to questions demonstrating that they appreciate more complete tags. Table 1 presents the ratio of questions involving tag re-editing on STACK OVERFLOW, ASK DIFFERENT, ASK UBUNTU, and SUPER USER. From the table, we can see that tag re-editing happens often. From the table, we can also notice that among the tag re-editing cases, 63.1–87.2% of the

---

<sup>7</sup><http://stackoverflow.com/q/5550896>

<sup>8</sup><http://stackoverflow.com/q/2058138>

# Apache Ant

Ant is a Java based build tool, similar to make, but with better support for the cross platform issues involved with developing Java applications. Ant is the build tool of choice for all Java projects at Apache and many other [Open](#) Source Java projects.

Tags

Software Development [Build Tools](#)

**Fig. 2** An example project in FREECODE and its tags

questions involve tag addition, which implies that tag addition is the major tag re-editing scenario.

A considerable number of studies have been done on tag recommendation for software information sites (Xia et al. 2013; Al-Kofahi et al. 2010). Among these studies, the approach TAGCOMBINE that was proposed by Xia et al. is shown to be the state-of-the-art (Xia et al. 2013) on software information sites. TAGCOMBINE combines three components: a multi-label ranking component, a similarity-based ranking component, and a tag-term based ranking component. The multi-label ranking component employs a multi-label classification algorithm (i.e., binary relevance method with naive Bayes as the underlying classifier) to predict the likelihood of a tag to be assigned to a software object. The similarity-based ranking component predicts the likelihood of a tag (to be assigned to a software object) by analyzing the tags that are given to the top-k most similar software objects that were tagged before. The tag-term based ranking component predicts the likelihood of a tag (to be assigned to a software object) by analyzing the number of times a tag has been used to tag a software object containing a term (i.e., a word) before. The multi-label ranking component of TAGCOMBINE constructs many one-versus-rest Naive Bayes classifiers, one for each tag. Each Naive Bayes classifier simply predicts the likelihood of a software object to be assigned a particular tag. However, mixture models have been shown to outperform one-versus-rest traditional multi-label classification approaches (Ramage et al. 2009; Ghamrawi and McCallum 2005; Puurula 2011). Thus, in our approach, we construct only one classifier which is a *mixture model* that considers all tags together to improve the effectiveness of the tag recommendation.

**Table 1** Objects with tag re-editing on STACK OVERFLOW, ASK DIFFERENT, ASK UBUNTU and SUPER USER

Dataset	Period	Questions involving tag re-editing	Total questions	Tag re-editing ratio	Tag addition ratio
STACK OVERFLOW	2015.1.1 – 2015.12.31	331,667	2,250,745	14.7%	72.1%
ASK DIFFERENT	Before 2016.1	11,243	63,276	17.7%	87.2%
ASK UBUNTU	Before 2016.1	48,942	191,191	25.6%	69.8%
SUPER USER	Before 2016.1	82,618	284,559	29.0%	63.1%

## 326 Tags

756	java	× 215	29	string	× 4
190	casting	× 2	27	php	× 14
186	assignment-operator	× 2	27	eclipse	× 13
186	operators		27	case-insensitive	

**Fig. 3** Highest rated tags associated to questions/answers posted by a user in STACK OVERFLOW

## 2.3 Motivating Examples

### 2.3.1 User Information

In addition to user-agnostic features (e.g., tag frequency), we also expect the information about the user posting the question to be useful when predicting tags. Indeed, one can conjecture that users have specific interests or expertise with certain technology and these interests or expertise are likely to manifest in the tags of their questions. To verify this conjecture we queried the users that asked more than one question on STACK OVERFLOW<sup>9</sup> and found that 51% of them have asked at least two questions labeled with the same tag. This suggests that users may post objects associated to some particular tags, rather than all tags, based on their personal background and interests.

For illustration, Fig. 3 presents the highest rated tags of a user in STACK OVERFLOW.<sup>10</sup> As of December 7, 2016, the user posted a number of questions/answers spanning 326 tags, and 215 of the user’s posts are tagged “java”, this user’s most commonly used tag. We can, therefore, conjecture that future software objects posted by the same user are more likely to be tagged with “java” tag, rather than other tags. Based on this observation, we can leverage user information to facilitate tag recommendation.

### 2.3.2 Additional Tag

In practice, it may be easy for users to label a question they post with a few tags. However, the set of initial tags may not be sufficient and for such cases, extra tags need to be added later - see Table 1. Intuitively, the initial tags can provide hints in the identification of missing tags.

We notice that some tags usually appear together. We could leverage tag co-occurrences to infer additional tags based on the initial tags that a user gave. Figure 4 presents an example of tag editing in Stack Overflow. The set of tags assigned to the question was refined – a tag “javascript” was added to the initial set of tags: “promise” and “pg-promise”. When we search questions that contain tags “promise”, and “pg-promise” and at least one other tag, we find 10 of such questions and 6 of them are also tagged with “javascript”. This suggests

<sup>9</sup><https://data.stackexchange.com/stackoverflow/queries>

<sup>10</sup><http://stackoverflow.com/users/137369/thirler?tab=tags>



2 edited tags link edited Jun 29 at 7:40  
yitaly-t 3,992 1 21 37

javascript promise pg-promise

1 source link asked Jun 29 at 3:25  
Parham 629 8 12

### Conditional task with pg-promise

I am trying to simply read a value from a table and based on the return value call for additional queries and return the combined results.

let's take a simple example: table `Users` has `id`, `name` and `emailid` and let's say if `emailid` is not null we want to call the email table and return a results like `{ id:[id], name:[name], email:[email]}`.

promise pg-promise

**Fig. 4** An example of adding a tag in STACK OVERFLOW

given an initial set of tags “promise” and “pg-promise”, it is likely that “javascript” should be included as well. This observation motivates us to recommend additional tags to users by analyzing the initial set of tags and leveraging tag co-occurrence.

## 3 General Architecture

In this section we describe the general architecture of our ENTAGREC<sup>++</sup> approach. ENTAGREC<sup>++</sup> contains six processing components: Preprocessing Component (PC), Bayesian Inference Component (BIC), Frequentist Inference Component (FIC), User Information Component (UIC), Additional Tag Component (ATC), and Composer Component (CC). Figure 5 presents the framework of ENTAGREC<sup>++</sup>.

Input software objects are processed by PC to generate a common representation. These textual documents are then input to the four main processing engines, namely BIC, FIC, UIC, and ATC. BIC and FIC infer tags based on words appearing in a software object. UIC infers tags based on the user who posts a software object; it works based on the assumption that a user tends to post similar software objects over time. ATC infers additional tags based on an initial set of tags given to a software object, by considering co-occurrences of tags. For some software objects, users who post them have provided some initial tags, and these tags can be used to better infer missing tags. CC combines the BIC, FIC, UIC, and ATC components.

ENTAGREC<sup>++</sup> works in two phases, a training phase and a deployment phase, as shown in Fig. 5a and b, respectively. In the training phase, ENTAGREC<sup>++</sup> trains several of its components using training software objects and corresponding tags. In the deployment phase, the trained ENTAGREC<sup>++</sup> is used to recommend tags for untagged software objects.

The common component in the training and deployment phase is PC, which converts each software object into a bag (or multiset) of words. The PC starts from the textual

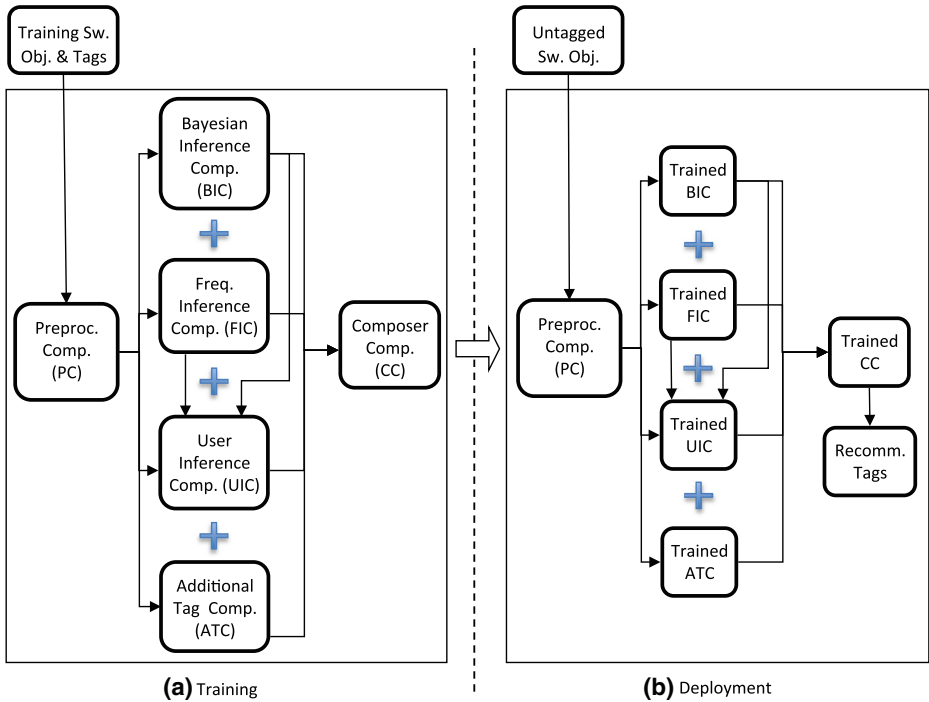


Fig. 5 ENTAGREC++ Architecture

description of a software object and performs tokenization, identifier splitting, number removal, stop word removal, and stemming. Tokenization breaks a document into word tokens. Identifier splitting breaks a source code identifier into multiple words. We split a token using two splitters: 1) Camel Casing splitter (Antoniol et al. 2002), e.g., the identifier “getMethodName” will be split into “get”, “method”, and “name”; 2) special sign splitter that splits tokens based on special signs (i.e.,  $\_$ ,  $-$ ), e.g., the identifier “get\_method\_name” will be split into “get”, “method”, and “name”. Number removal deletes numbers. Stop word removal<sup>11</sup> deletes words that are used in almost every document and, therefore, carry little document-specific meaning, e.g., “the”, “is”, etc. Finally, stemming reduces words to their root form. We use the Porter stemming algorithm (Porter 1997).

In the training phase, BIC, FIC, UIC, ATC, and CC are trained based on the training data. BIC uses the bag-of-words representation of the software objects and their corresponding tags to train itself. The result is a statistical model which takes as input a bag of words representing a software object, and produces a ranked list of tags along with their probabilities of being related to the input software object. FIC also processes the bag-of-words representations and the corresponding tags to train itself; it produces a statistical model (albeit in a different way than BIC) which also takes as input a bag of words, and outputs a ranked list of tags with their probabilities. UIC processes data about users who posted various software objects, to model the peculiar behaviors of the various users; it creates a statistical model which takes as input a user who posted a particular software object, and outputs a ranked

<sup>11</sup>Based on <http://www.textfixer.com/resources/common-english-words.txt>

list of tags with their probabilities. ATC takes as input a set of tags appearing together in the past. The result is a conditional statistical model, which takes an initial set of tags given by a user as input, and produces a ranked list of additional tags, with their probabilities. CC learns four weights for BIC, FIC, ATC, and UIC to generate a near-optimal combination of these four components from the training data.

After ENTAGREC<sup>++</sup> is trained, it is used in the deployment phase to recommend tags for untagged objects. For each such object, we first use PC to convert it to a bag of words. Next, we feed this bag of words, including the user information and additional tags (if available), to the trained BIC, FIC, UIC, and ATC. Each of them will produce a list of tags with their likelihood scores. CC will compute the final likelihood score for the tags based on the weights that it has learned in the training phase. The top few tags with the highest likelihood scores will be output as the predicted tags of an input untagged or partially tagged software object.

The following sections detail each of the five major components of ENTAGREC<sup>++</sup>, BIC, FIC, UIC, ATC, and CC.

## 4 Bayesian Inference

The goal of BIC is to compute the probabilities of various tags, given a bag of words representing a software object, using Bayesian inference. Given a tag  $t$  and a software object  $o$ , BIC computes the conditional probability of  $t$  being assigned to  $o$ , given the words  $\{w_1, \dots, w_n\}$  that appear in  $o$ . This is denoted as  $P(t|w_1 \dots w_n)$ . Using the Bayes theorem (Gelman et al. 2003), this probability can be computed as:

$$P(t|w_1 \dots w_n) = \frac{P(w_1 \dots w_n|t) \times P(t)}{P(w_1 \dots w_n)} \quad (1)$$

The probabilities on the right hand side of the above equation can be estimated based on training data.

A state-of-the-art Bayesian inference algorithm is Latent Dirichlet Allocation (LDA) (Blei et al. 2003). LDA has been shown effective to process various software engineering data for various tasks, e.g., Lukins et al. (2010), Baldi et al. (2008), Asuncion et al. (2010), Panichella et al. (2013), and Rebouças et al. (2016). LDA takes as input a set of documents and a number of topics  $K$ , and outputs the probability distribution of topics per document. Our problem can be readily mapped to LDA, where a document corresponds to a software object, and a topic corresponds to a tag. Using this setting, LDA outputs the probability distribution of tags for a software object.

However, LDA is an unsupervised learning algorithm. It does not take as input any training data and it is not possible to pre-define a set of tags as the target topics to be assigned to documents. Fortunately, recent advances in the natural language processing community introduced extensions to LDA, such as Labeled LDA (L-LDA) Ramage et al. (2009). For L-LDA, the labels can be predefined and a training set of documents can be used to train the LDA, such that it will compute the probability distribution of topics, coming from a pre-defined label set (tags, in our case), for a document (a software object, in our case), based on a set of labeled training data. In this work, we use L-LDA as the basis for the Bayesian inference component.

BIC works on two phases: training and deployment. In the training phase, BIC takes as input a set of bags of words representing software objects, and their associated tags. These are used to train an L-LDA model. In the deployment phase, given a bag of words

corresponding to a software object, the trained L-LDA model is used to infer the set of tags for the input software object along with their probabilities. In the end, the top  $K_{Bayesian}$  inferred tags for the object will be output and fed to the Composer Component (CC).

*Example 1* Consider an object has following words {install, eclipse} and a tag eclipse. In order to compute the probability  $P(\text{eclipse}|\text{install, eclipse})$ , we need to estimate the value of  $P(\text{install, eclipse}|\text{eclipse})$ ,  $P(\text{eclipse})$ , and  $P(\text{install, eclipse})$  first based on the (1). By using L-LDA, we could estimate the the value of  $P(\text{install, eclipse}|\text{eclipse})$ ,  $P(\text{install, eclipse})$  and  $P(\text{eclipse})$ . Suppose the estimated value of  $P(\text{install, eclipse}|\text{eclipse}) = 0.02$ ,  $P(\text{install, eclipse}) = 0.001$ , and  $P(\text{eclipse}) = 0.005$ , thus the value of  $P(\text{eclipse}|\text{install, eclipse})$  is 0.1.

## 5 Frequentist Inference

FIC computes the probability that a software object is assigned a particular tag based on the words that appear in the software object, while taking into account the *number* of words that appear along with the tag in software objects in a training set. Section 5.1 describes our basic approach and several extensions are presented in Section 5.2. Hereafter, unless stated otherwise, FIC refers to the extended approach.

### 5.1 Basic Approach

Consider software object  $o$  with  $n$  words:  $\{w_1, w_2, \dots, w_n\}$  and a tag  $t$ , the weight of tag  $t$  for object  $o$  can be computed as the proportion of the  $n$  words that co-appear with tag  $t$  in the training data. More formally, the weight is defined as

$$W(o, t) = \frac{\sum_{w_i \in o} I(t, w_i)}{|o|}, \quad (2)$$

where

$$I(t, w_i) = \begin{cases} 1, & \exists o \in TRAIN | o \text{ contains } w_i \&o \text{ tagged with } t \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

The higher the weight  $W(o, t)$ , the more representative FIC deems tag  $t$  to be for software object  $o$ .

### 5.2 Extended Approach

There are several problems with the basic approach. First, often not all words in a software object are related to the tags that are assigned to the software object. Although the pre-processing component (PC) has removed stop words, still many non-stop words are unrelated to software object tags, thus need to be removed. Second, we have a data sparsity problem, since many tags are not used frequently in the training set. Thus, often a tag is not characterized by sufficiently many words. To address this problem, we leverage the relationships among tags to recommend additional associated tags to an input untagged software object.

### 5.2.1 Removing Unrelated Words with POS Tagger

One problem in estimating the probabilities  $P(t|w_i)$  is that not all words that appear in a software object are related to the tags. We use the example in Fig. 1 to illustrate this. Words “need” and “work” are not stop words, but they are unrelated to the tags `eclipse` and `eclipse-plugin`. Thus, there is a need to filter out these unrelated words before we estimate the probabilities.

We observe that nouns and noun phrases are often more related to the tags than other kinds of words. Past studies have also found that nouns are often the most important words (Capobianco et al. 2013; Shokripour et al. 2013). Thus, in this extension, we remove all words except nouns and noun phrases. To identify these nouns and noun phrases, we use the Part-Of-Speech (POS) Tagger (Toutanova et al. 2003) to infer the POS of each word in the representative bag of words of a software object. In this paper, we use the Stanford Log-linear Part-Of-Speech Tagger.<sup>12</sup> To illustrate this extension, consider the words that appear in the software object shown in Fig. 1. After this step, only the words “tutorial”, “eclipse”, “plugin”, “interface”, “function”, “java”, and “app” remain.

Note that we only did this for FIC and not BIC as L-LDA assigns different probabilities to words that are associated to a topic (i.e., a tag). Unrelated words will receive low probabilities. In FIC, the words that appear in objects tagged with tag  $t$  are treated as equally important. Thus, we only perform this extended processing step for FIC.

We refer to the basic approach extended by this processing step as *FrePOS*. Given an untagged software object, *FrePOS* outputs the top  $K_{\text{Frequentist}}$  tags.

### 5.2.2 Finding Associated Tag with Spreading Activation

Due to the data sparseness problem, *FrePOS* might miss some important tags that are not adequately represented in the training data. To find additional tags, we leverage relationships among tags using a technique named spreading activation (Crestani 1997). Spreading activation takes as input a network containing weighted nodes that are connected with one another with weighted edges, and a set of starting nodes. Initially, all nodes except the starting nodes are assigned weight 0. Spreading activation then processes the starting nodes, one at a time. For each starting node, it spreads (or propagates) the node’s weight to its neighboring nodes which are at most  $MH$  hops away from it (where  $MH$  is a user-defined threshold). At the end of the process, we output all nodes with non zero weights and their associated weights. In our context, the network is a tag network, the starting nodes are the nodes corresponding to tags returned by *FrePOS*, and the weights of these starting nodes are the probabilities assigned to the corresponding tags by *FrePOS*.

To perform spreading activation, we first need to construct a network of tags. Each node in the network corresponds to a tag, and each edge connecting two nodes in the network corresponds to the relationship between the corresponding tags. The weight of each edge measures how similar two tags are. We measure this based on the co-occurrence of tags in software objects in the training set. Consider a set of tags where each of them is used to label at least one software object in the training set. We denote this set as:  $\text{Tags} = \{t_1, t_2, t_3, \dots, t_k\}$ , where  $k$  is the total number of unique tags. We denote an edge between two tags  $t_i$  and  $t_j$  as  $e_{t_i, t_j}$ . The weight of  $e_{t_i, t_j}$  depends on the number of

---

<sup>12</sup><http://nlp.stanford.edu/software/tagger.shtml>

software objects that are tagged by  $t_i$  and  $t_j$  in the training set. It can be calculated as follows:

$$weight(e_{t_i,t_j}) = \frac{|Doc(t_j) \cap Doc(t_i)|}{|Doc(t_i) \cup Doc(t_j)|} \quad (4)$$

where  $Doc(t_i)$  and  $Doc(t_j)$  are the sets of objects tagged with  $t_i$  and  $t_j$ , respectively, and  $|\cdot|$  denotes cardinality.

The edge connecting two tags is assigned a higher weight if the tags appear together more frequently, which means they are more associated with each other. We denote the set of edges connecting pairs of nodes as *Links*. The tag network is then a graph  $TN$  defined as (*Tags*, *Links*). Given a tag  $t$ , we denote the node in  $TN$  corresponding to  $t$  as  $TN[t]$ . Given a node  $n$  and an edge  $E(n_1, n_2)$ , we denote their weights as  $weight(n)$  and  $weight(E(n_1, n_2))$ , respectively.

The pseudocode of our approach to infer associated tags from the initial set of tags returned by *FrePOS* is shown in Algorithm 1. The algorithm takes as input a tag network  $TN$  constructed from all tags in the training data, a set of starting tags  $SST$  returned by *FrePOS*, and a threshold  $MH$  that restricts the weight propagation to a maximum number of hops. Then, it initializes the weights of nodes corresponding to tags in the set of starting tags with the probabilities returned by *FrePOS*, and it sets the weights of other nodes to 0 (Lines 8–11). For each starting tag, our algorithm then performs spreading activation starting from the corresponding node in the tag network by calling the procedure `SpreadingActivation` (Lines 12–14). Finally, the algorithm outputs all nodes in the set of starting tags, along with the associated tags, which correspond to nodes in  $TN$  whose weights are larger than zero (Line 15).

---

**Algorithm 1** Find associated tags

---

```

1: FindAssociatedTags
2: Input:
3:  $TN$ : Tag network
4:  $SST$ : Set of starting tags
5:  $MH$ : Maximum hop
6: Output: Set of candidate tags
7: Method:
8: Initialize the weight of each tag in  $TN$  with 0
9: for each tag  $t$  in  $SST$  do
10:   Set  $weight(TN[t]) =$  Probability of tag  $t$  inferred by FrePOS
11: end for
12: for each tag  $t$  in  $SST$  do
13:   Call SpreadingActivation( $TN, TN[t], 0, MH$ )
14: end for
15: return  $SST \cup \{t \mid weight(TN[t]) > 0\}$ 

```

---

The procedure `SpreadingActivation` spreads the weight of a node to its neighbors. It takes as input a tag network  $TN$ , a starting node  $N$ , the current hop  $CH$ , and the maximum hop  $MH$ . The procedure first checks if it needs to propagate the weight of node  $N$ —it only propagates if the current hop  $CH$  does not exceed the threshold  $MH$ , and the weight of the current node is larger than zero (Lines 8–10). It then iterates through nodes  $N'$  that are directly connected to  $N$  (Lines 11–17). For each such node, we compute a weight  $w$  which is a product of the weight of node  $N$  and the weight of the edge  $N-N'$  (Line 12). If the weight of node  $N'$  is less than  $w$ , we assign  $w$  as the weight of node  $N'$  (Lines 13–14). The

procedure then tries to propagate the weight of  $N'$  to its neighbors by a recursive call to itself (Line 15).

---

**Algorithm 2** Spreading activation for a node

---

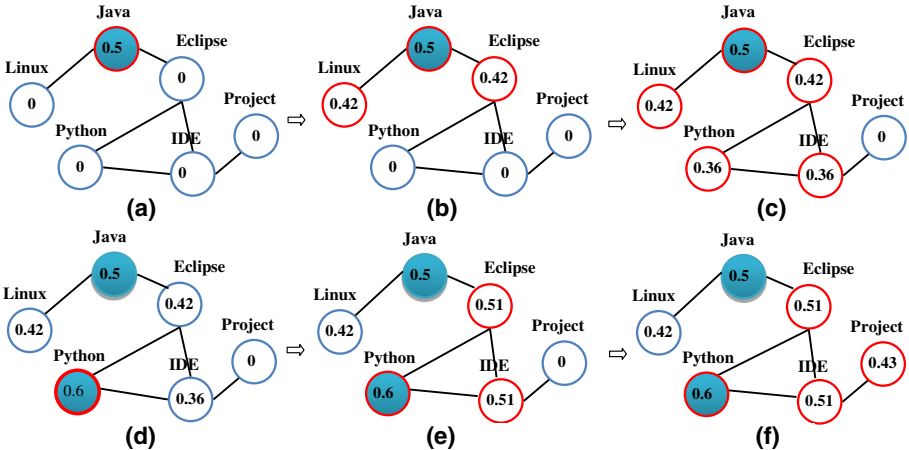
```

1: SpreadingActivation
2: Input:
3:  $TN$ : Tag network
4:  $N$ : Current node
5:  $CH$ : Current hop
6:  $MH$  Maximum hop
7: Method:
8: if  $CH > MH$  or  $weight(N) = 0$  then
9:   return
10: end if
11: for each node  $N'$  that is directly connected to  $N$  do
12:   Set  $w = weight(N) \times weight(E(N', N))$ 
13:   if  $weight(N') < w$  then
14:      $weight(N') = w$ 
15:     SpreadingActivation( $TN, N', CH+1, MH$ )
16:   end if
17: end for

```

---

*Example 2* Consider a set of starting tags  $SST = \{JAVA = 0.5, PYTHON = 0.6\}$  output by *FrePOS*, a tag network  $TN$  shown in Fig. 6 and a threshold  $MH = 2$ . Let us assume the weights of all edges in the tag network are 0.85. At the beginning, our approach initializes the weight of the node corresponding to tag *JAVA* in  $TN$  with 0.5 (Fig. 6a). Then, the weight of node *JAVA* is propagated to its neighbors *LINUX* and *ECLIPSE* and their weights are both updated to 0.42 (Fig. 6b). The weight is recursively propagated to all neighbors of node *JAVA* of distance  $MH$  hops or less (Fig. 6c). Then, our approach processes tag *PYTHON*, node *PYTHON*'s weight is updated to 0.6, which is the weight of tag *PYTHON* output by *FrePOS*.



**Fig. 6** Finding associated tags using spreading activation: an example

(Fig. 6d). Our approach then propagates the weight of node PYTHON to its neighbors. If a neighbor's weight is lower than that which is propagated from PYTHON, the original weight is replaced with the new weight. Otherwise, the original weight remains unchanged. Thus, the weights of ECLIPSE and IDE are updated to 0.51 (0.51 exceeds 0.425, the current weights of these tags, Fig. 6e). The weight of node PROJECT is updated to 0.43 (Fig. 6f). Finally, the tags JAVA = 0.5, PYTHON = 0.6, ECLIPSE = 0.51, IDE = 0.51, LINUX = 0.42, and PROJECT = 0.43 will be output.

The spreading activation process requires a parameter  $MH$  (maximum hop); by default, we set the parameter  $MH$  to 1, as the complexity of spreading activation is exponential to the value of  $MH$ . At the end, our FIC component outputs candidate tags that are output by *FrePOS* and the associated tags that are output by the spreading activation procedure described above. These tags are input to the composer component (CC).

Note that we only apply this spreading activation step to FIC and not BIC. L-LDA used in BIC is more robust than *FrePOS* to the data sparsity problem. We find that the application of this step to BIC does not improve its effectiveness.

## 6 User Information Component

In this component, we make use of tags attached to software objects that a user has posted before, in order to infer tags for a new software object posted by the same user. As we show in the Section 2.3.2, the user usually posts questions that associated to certain specific tags. Based on this intuition, we compute the weight of a tag  $t$  given a new software object  $o$  posted by a user  $u$  as:

$$w(t, o, u) = \begin{cases} \frac{|\{o \in Doc(u) | o \text{ tagged with } t\}|}{|Doc(u)|}, & \text{if } t \in T_{BIC} \cup FIC, \\ 0, & \text{otherwise} \end{cases}, \quad (5)$$

where  $Doc(u)$  is the set of past objects posted by  $u$ ,  $T_{BIC} \cup FIC$  is the set of tags with non-zero weights from BIC and FIC, and  $|\cdot|$  denotes cardinality. Note that we define  $w(t, o, u)$  as 0 for tags not in  $T_{BIC} \cup FIC$  to avoid noise due to irrelevant tags.<sup>13</sup>

*Example 3* To illustrate how UIC works, consider the user in Fig. 3. Suppose she posts a new question  $o$ . We first find questions and answers posted by her in the past. Second, we use BIC and FIC to get a list of candidate tags for the question  $o$ , say  $T_{BIC} \cup FIC$  is {java, netbeans, string, algorithm}. Third, using (5), we compute the tags' weights: java receives weight  $\frac{203}{317}$ , string  $\frac{4}{317}$ , algorithm  $\frac{2}{317}$ , and netbeans  $\frac{2}{317}$  (she happens to have used all four tags recommended by BIC&FIC).

## 7 Additional Tag Component

In this component, we make use of an initial set of tags provided by a user to infer additional tags based on historical tag co-occurrences. Consider a software object  $o$  and a set of initial

<sup>13</sup>Our experiments show that the effectiveness of UIC substantially degrades if it takes into consideration all tags.



tags  $\{t_1, t_2, \dots, t_k\}$  provided by a user. The probability of the object  $o$  to be assigned a tag  $t$  is:

$$P(o, t|t_1, t_2, \dots, t_k) = \prod_{i=1}^k P(o, t|t_i) \quad (6)$$

The above probabilities ( $P(o, t|t_i)$  for  $1 \leq i \leq k$ ) can be estimated from the training data as follows:

$$P(o, t|t_i) = \frac{\text{Number of objects labeled with } t \text{ and } t_i}{\text{Number of objects labeled with } t_i} \quad (7)$$

*Example 4* Suppose a user posts a question  $o$  and provides an initial tag string. In the training data, let us say there are 100 software objects labeled with tag `string`, and among them 40, 30, 10, 20, and 10 posts are also labelled with tags `java`, `c#`, `io`, `javascript`, and `python`, respectively. Using (7), we estimate the probabilities:  $P(o, \text{java}|\text{string}) = 0.4$ ,  $P(o, \text{c\#}|\text{string}) = 0.3$ ,  $P(o, \text{io}|\text{string}) = 0.1$ ,  $P(o, \text{javascript}|\text{string}) = 0.2$  and  $P(o, \text{python}|\text{string}) = 0.1$ .

## 8 Composer Component

Given a target software object  $o$  and a tag  $t$ , BIC, FIC, UIC, and ATC each produces a probability for the tag to be relevant. We need to combine these probabilities to estimate the overall relevance of the tag. A simple solution is to take an average of the four probabilities. However, this assumes that BIC, FIC, UIC, and ATC are equally accurate in predicting tag relevance, which may not be the case. To accommodate for differences in the effectiveness of the four components, we can assign weights to them. More accurate components can be given higher weights, and these weights can be learned from a training data. After these weights are learned, for every tag, we can compute a weighted average of its probabilities and use it as its overall relevance score. This score can then be used to produce a final ranked list of tags. This strategy is commonly referred to as fusion via a linear combination of scores which is a classical information retrieval technique (Vogt and Cottrell 1999).

More formally, we define  $\text{ENTAGREC}^{++}$  ranking score as  $\text{ENTAGREC}^{++}_o(t)$  as follows:

$$\text{ENTAGREC}^{++}_o(t) = \alpha \times B_o(t) + \beta \times F_o(t) + \gamma \times U_o(t) + \delta \times A_o(t), \quad (8)$$

where  $B_o(t)$ ,  $F_o(t)$ ,  $U_o(t)$ , and  $A_o(t)$  are the probabilities of tag  $t$  computed by BIC, FIC, UIC, and ATC, respectively, and  $\alpha$ ,  $\beta$ ,  $\gamma$ ,<sup>14</sup> and  $\delta \in [0, 1]$  are the weights the composer component assigns to BIC, FIC, UIC, and ATC, respectively.<sup>15</sup> Note that if there is no additional tag provided by users, ATC will be deactivated and, correspondingly,  $\delta$  will be set to 0.

To automatically tune  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ , we use a set of training software objects and employ grid search (Bergstra and Bengio 2012). The pseudocode of our weight tuning procedure is shown in Algorithm 3. The weight tuning procedure takes as input the set of training software objects  $TO$ , an evaluation criterion  $EC$ , and the four sets of tags returned by BIC, FIC, UIC, and ATC (along with their probabilities). Our tuning procedure initializes  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  to 0 (Line 12). Then, it incrementally increases the value of  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  by

<sup>14</sup>By construction,  $\gamma$  is an extra weight given to some of the tags in  $T_{\text{BIC} \cup \text{FIC}}$ .

<sup>15</sup>Since  $\text{ENTAGREC}^{++}_o(t)$  is itself a probability score, it could also be expressed as a function of only three coefficients  $\alpha'$ ,  $\beta'$ , and  $\gamma'$ , with the fourth being automatically  $1 - \alpha' - \beta' - \gamma'$ . We chose the four-coefficient expression to better reflect the four components of  $\text{ENTAGREC}^{++}$ .

0.1 until they reach 1.0 (Lines 13–16). For each combination of four parameters and each software object  $o$  in  $TO$ , our tuning procedure computes the ENTAGREC<sup>++</sup> scores for each tag returned by BIC, FIC, UIC, and ATC (Lines 17–19). Then tags are ordered based on their ENTAGREC<sup>++</sup> scores (Line 21). This is the ranked list of tags that are recommended for  $o$ . Next, our tuning procedure evaluates the quality of the resulting ranking based on particular  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  values using  $EC$  (Line 22). The process is repeated for all objects in  $TO$  and again the quality of the resulting ranking is evaluated using  $EC$  (Line 24). The process continues until all combinations of  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  have been exhausted and our tuning procedure finally outputs the best combination of  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  based on  $EC$  (Line 29).

---

**Algorithm 3** Weight Tuning Algorithm

---

```

1: TuneWeights
2: Input:
3:  $TO$ : Training Tagged Software Objects
4:  $EC$ : Evaluation Criterion
5:  $Tags^B$ : Set of tags inferred by BIC
6:  $Tags^F$ : Set of tags inferred by FIC
7:  $Tags^U$ : Set of tags inferred by UIC
8:  $Tags^A$ : Set of tags inferred by ATC
9: Output:
10:  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ 
11: Method:
12: Set  $\alpha = 0$ ,  $\beta = 0$ ,  $\gamma = 0$ ,  $\delta = 0$ 
13: for each  $\alpha$  from 0 to 1, each step increases  $\alpha$  by 0.1 do
14:   for each  $\beta$  from 0 to 1, each step increases  $\beta$  by 0.1 do
15:     for each  $\gamma$  from 0 to 1, each step increases  $\gamma$  by 0.1 do
16:       for each  $\delta$  from 0 to 1, each step increases  $\delta$  by 0.1 do
17:         for each object  $o$  in  $TO$  do
18:           for each tag  $t$  in  $Tags^B \cup Tags^F \cup Tags^U \cup Tags^A$  do
19:             Compute ENTAGREC++ $o$ ( $t$ ) according to Eq. 8
20:           end for
21:           Sort tags based on their ENTAGREC++ scores (desc. order)
22:           Evaluate the effectiveness of  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  on  $o$  based on  $EC$ 
23:         end for
24:       Evaluate the effectiveness of  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  on  $TO$  based on  $EC$ 
25:     end for
26:   end for
27: end for
28: end for
29: return the best  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  based on  $EC$ 

```

---

Various evaluation criteria can be used in our weight tuning procedure. In this paper, we make use of  $Recall@k$ , which has been used as the evaluation criterion in many past tag recommendation studies, e.g., Al-Kofahi et al. (2010) and Zangerle et al. (2011).  $Recall@k$  was also used in the previous state-of-the-art study on tag inference for software information sites (Xia et al. 2013).

**Definition 1** Consider a set of  $n$  software objects. For each object  $o_i$ , let the set of its correct (i.e., ground truth) tags be  $Tags_i^{correct}$ . Also, let  $Tags_i^{topK}$  be the top- $k$  ranked tags that are recommended by a tag recommendation approach for  $o_i$ .  $Recall@k$  for  $n$  is given by:

$$Recall@k = \frac{1}{n} \sum_{i=1}^n \frac{|Tags_i^{topK} \cap Tags_i^{correct}|}{|Tags_i^{correct}|} \quad (9)$$

In the deployment phase, the composer component combines the recommendations made by the inference components by computing the ENTAGREC<sup>++</sup> scores using (8) for each recommended tag. It then sorts the tags based on their ENTAGREC<sup>++</sup> scores (in descending order) and outputs the top-k ranked tags.

## 9 Experiments and Results

In this section, we first present our experiment settings in Section 9.1. Our experiment results are then presented in Sections 9.2. We discuss some interesting points in Section 9.3.

### 9.1 Experimental Setting

We evaluate ENTAGREC<sup>++</sup> on five datasets: STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, SUPER USER (all four part of the STACK EXCHANGE network), and FREECODE, which were used to evaluate ENTAGREC (Wang et al. 2014). STACK OVERFLOW is a Q&A site for software developers to post general programming questions. ASK DIFFERENT is a Q&A site related to Apple devices, e.g., iPhone, iPad, mac. ASK UBUNTU is a Q&A site about Ubuntu. SUPER USER is a Q&A site for systems administrators and power users. FREECODE is a site containing descriptions of many software projects.

Table 2 presents descriptive statistics of the four datasets, including period of the data, number of objects, number of tags, maximum and average number of objects for each per tag, and average elapsed time since registration of all studied users. The STACK OVERFLOW and FREECODE datasets are obtained from Xia et al. and they have been used to evaluate TAGCOMBINE (Xia et al. 2013). The ASK UBUNTU, ASK DIFFERENT, and SUPER USER datasets are new. We collect all questions in ASK UBUNTU, ASK DIFFERENT, SUPER USER that are posted before April 2012. Following (Xia et al. 2013), to remove noise corresponding to tags that are assigned idiosyncratically, we filter out tags that are associated with less than 50 objects. These tags are less interesting since not many people use them, and thus they are less useful to be used as *representative tags* and recommending them does not help much in addressing the tag synonym problem addressed by tag recommendation studies. The numbers summarized in Table 2 are *after filtering*.

We perform ten-fold cross validation (Han et al. 2011) for evaluation. We randomly split the dataset into ten subsamples. Nine of them are used as training data, to train ENTAGREC<sup>++</sup>, and one subsample is used for testing. We repeat the process ten times and use *Recall@k* as the evaluation metric. Note that we conduct ten-fold cross validation 100 times

**Table 2** Basic statistics of the four datasets

Dataset	Period	Objects	Tags	Objects per tag		Avg. age of users (days)
				Max	Avg	
STACK OVERFLOW	2008.6 – 2008.12	47,668	437	6,113	234.93	54
FREECODE	2001.1 – 2012.6	39,231	243	9,615	545.08	NA
ASK UBUNTU	Before 2012.4	37,354	346	6,169	234.03	237
ASK DIFFERENT	Before 2012.4	13,351	153	2,019	180.88	253
SUPER USER	Before 2012.4	47,996	460	7,009	245.7	745

and take averages as results. Unless otherwise stated, we set the values of  $K_{Bayesian}$  and  $K_{Frequentist}$  at 70 as the setting in ENTAGREC. We conduct all our experiments on a Windows 2008 server with 8 Intel®2.53GHz cores and 24GB RAM.

## 9.2 Evaluation Results

The goal of our evaluation is to compare the effectiveness of ENTAGREC<sup>++</sup> with those of ENTAGREC and TAGCOMBINE. TAGCOMBINE is the tag recommendation approach proposed by Xia et al. (2013). ENTAGREC is our earlier version of ENTAGREC<sup>++</sup> (Wang et al. 2014), which did not include the user information component and the additional tag component. Our goal can be refined into the following research questions:

RQ1. How effective is ENTAGREC<sup>+</sup> compared to ENTAGREC and TAGCOMBINE in terms of *Recall@k*?

To answer this research question, we perform ten-fold cross validation, and compare ENTAGREC<sup>++</sup>, ENTAGREC and TAGCOMBINE in terms of *Recall@5* and *Recall@10*. To make the comparison fair, we do not provide an initial set of tags to ENTAGREC<sup>++</sup> because ENTAGREC and TAGCOMBINE cannot accept an initial set of tags as input. We refer to the version of ENTAGREC<sup>++</sup> without additional tags as ENTAGREC<sup>+</sup>.

RQ2. Does the additional tag component improve the effectiveness of ENTAGREC<sup>+</sup>?

The additional tag component makes use of the additional tags provided by a user to recommend associated tags. We want to investigate whether this component improves ENTAGREC<sup>+</sup>. To answer this research question, we collect the questions whose tags have been updated in the past from STACK OVERFLOW, ASK UBUNTU, and ASK DIFFERENT. In the experiment, we take the initial tags labeled by users when they created the questions, and use the recent tags of the question as the ground truth. We also remove the tags that are associated with less than 50 objects, as before in RQ1. We use the datasets in Table 1 for this experiment; the number of tags and objects after filtering is shown at Table 3. For ASK UBUNTU, ASK DIFFERENT and SUPER USER, we collect all questions involving tag re-editing before December 2015. After filtering, we are left with 483 tags and 31,881 objects associated with the tags from ASK UBUNTU, with 157 tags and 7,762 objects from ASK DIFFERENT, and with 196 tags and 13,796 objects from SUPER USER. For STACK OVERFLOW, because the number of questions involving tag re-editing is too large, we randomly sample 50,000 questions; after filtering, we are left with 649 tags and 42,493 objects. We do not consider FREECODE for this experiment because we cannot obtain historical tag data from FREECODE.

**Table 3** Basic statistics of questions involving tag re-editing on STACK OVERFLOW, ASK UBUNTU, SUPER USER, and ASK DIFFERENT

Dataset	Tags	Objects	Avg. initial tag set size	Avg. edited tag set size
STACK OVERFLOW	649	42,493	2.9	3.4
ASK UBUNTU	483	31,881	2.4	2.8
ASK DIFFERENT	157	7,762	2.3	3.1
SUPER USER	196	13,796	2.0	2.7

### 9.2.1 RQ1: Overall Effectiveness of ENTAGREC<sup>+</sup>

We compare ENTAGREC<sup>+</sup> with competing approaches: TAGCOMBINE proposed by Xia et al. (2013) and ENTAGREC by Wang et al. (2014).

Table 4 summarizes the comparison between ENTAGREC<sup>+</sup>, ENTAGREC, and TAGCOMBINE. We also show the beanplots of the comparison between those three approaches in Fig. 7. We performed ten-fold cross-validation 100 times and evaluated the approaches in terms of the average *Recall@5* and *Recall@10*. ENTAGREC<sup>+</sup> achieves sizeable improvements over TAGCOMBINE for the Stack Exchange datasets (more than 34.7% for *Recall@5* and more than 18.3% for *Recall@10*), and performs comparably to TAGCOMBINE on FREECODE. Averaging across the 5 datasets, ENTAGREC<sup>+</sup> improves TAGCOMBINE in terms of *Recall@5* and *Recall@10* by 27.8% and 14.1% respectively. We perform a Wilcoxon signed-rank test (Wilcoxon 1945) to test the significance of the differences in the performance of TAGCOMBINE and ENTAGREC<sup>+</sup> measured in terms of *Recall@5* and *Recall@10*. We also perform Benjamini Yekutieli procedure (Benjamini and Yekutieli 2001) to adjust the *p*-value obtained from Wilcoxon signed-rank test to deal with the impact of multiple comparisons. For the Stack Exchange datasets (STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, SUPER USER), the results show that ENTAGREC<sup>+</sup> outperforms TAGCOMBINE in terms of *Recall@5* and *Recall@10* significantly. For FREECODE, ENTAGREC<sup>+</sup> outperforms TAGCOMBINE in terms of *Recall@5* significantly. However, TAGCOMBINE significantly outperforms ENTAGREC<sup>+</sup> in terms of *Recall@10* but the absolute difference is small (0.016).

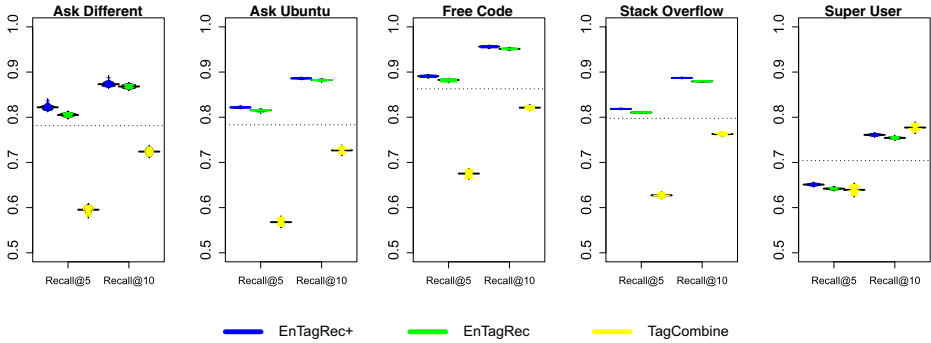
We also compare the effectiveness of ENTAGREC<sup>+</sup>, which employs user information, to that of ENTAGREC in terms of *Recall@5* and *Recall@10*. Analyzing Table 4, we observe that ENTAGREC<sup>+</sup> outperforms ENTAGREC on all datasets. In terms of *Recall@5*, ENTAGREC<sup>+</sup> improves ENTAGREC by 1.3% on average. In terms of *Recall@10*, ENTAGREC<sup>+</sup> achieves a 0.3% improvement over ENTAGREC. We perform a Wilcoxon signed-rank test (Wilcoxon 1945) and Benjamini Yekutieli procedure (Benjamini and Yekutieli 2001) on each dataset. The results indicates the improvement achieved by ENTAGREC<sup>+</sup> is statistically significant (adjusted *p*-value < 0.05).

To investigate if the differences in the *Recall@5* and *Recall@10* values are *substantial*, we also compute Cliff’s Delta (Grissom and Kim 2005) which measures effect size. The

**Table 4** The comparison of ENTAGREC<sup>+</sup>, ENTAGREC and TAGCOMBINE in terms of *Recall@5* and *Recall@10*

Dataset	ENTAGREC	ENTAGREC <sup>+</sup>	TAGCOMBINE
<i>Recall@5</i>			
STACK OVERFLOW	0.805	<b>0.821</b>	0.595
ASK UBUNTU	0.815	<b>0.822</b>	0.568
ASK DIFFERENT	0.882	<b>0.891</b>	0.675
SUPER USER	0.810	<b>0.818</b>	0.627
FREECODE	0.642	<b>0.651</b>	0.639
<i>Recall@10</i>			
STACK OVERFLOW	0.868	<b>0.873</b>	0.724
ASK UBUNTU	0.882	<b>0.886</b>	0.727
ASK DIFFERENT	0.951	<b>0.956</b>	0.821
SUPER USER	0.879	<b>0.887</b>	0.763
FREECODE	0.754	0.761	<b>0.777</b>

The highest value is typeset in boldface



**Fig. 7** Beanplots of ENTAGREC<sup>+</sup>, ENTAGREC and TAGCOMBINE in terms of *Recall@5* and *Recall@10*

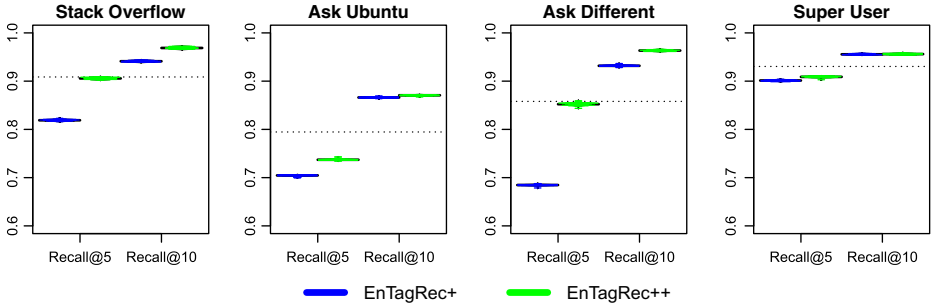
results are shown in Table 5. It interprets the effect size values as small for  $0.147 < |d| < 0.33$ , medium for  $0.33 < |d| < 0.474$ , and large for  $|d| > 0.474$  (Grissom and Kim 2005). If the effect size is close to 0, it means that all values of one group are larger than those of another group. From the results we can conclude that ENTAGREC<sup>+</sup> substantially outperforms TAGCOMBINE and ENTAGREC on STACK OVERFLOW, ASK UBUNTU, SUPER USER, and ASK DIFFERENT datasets (with large effect sizes). For the FREECODE dataset, ENTAGREC<sup>+</sup> also substantially outperforms ENTAGREC (with large effect sizes). In terms of *Recall@5*, ENTAGREC<sup>+</sup> outperforms TAGCOMBINE substantially. However, in terms of *Recall@10*, TAGCOMBINE outperforms ENTAGREC<sup>+</sup> substantially, which means that all values in one group are larger or smaller than those in another group when comparing two groups.

### 9.2.2 RQ2: Effectiveness of Additional Tag Component

To answer the research question, we compare the effectiveness of two versions of ENTAGREC<sup>++</sup>: one with additional tag component (ENTAGREC<sup>++</sup>) and another without additional tag component (ENTAGREC<sup>+</sup>). Figure 8 presents the beanplots of ENTAGREC<sup>+</sup> and

**Table 5** Effect sizes

Dataset	ENTAGREC <sup>+</sup> vs. ENTAGREC	ENTAGREC <sup>+</sup> vs. TAGCOMBINE
<i>Recall@5</i>		
STACK OVERFLOW	1	1
ASK UBUNTU	1	1
ASK DIFFERENT	1	1
SUPER USER	0.99	0.99
FREECODE	1	0.92
<i>Recall@10</i>		
STACK OVERFLOW	0.68	1
ASK UBUNTU	1	1
ASK DIFFERENT	1	1
SUPER USER	0.99	0.99
FREECODE	1	-1



**Fig. 8** Beanplots of ENTAGREC<sup>+</sup> and ENTAGREC<sup>++</sup> in terms of *Recall@5* and *Recall@10*

ENTAGREC<sup>++</sup> in terms of *Recall@5* and *Recall@10*. Table 6 summarizes the comparison between ENTAGREC<sup>++</sup> and ENTAGREC<sup>+</sup> in terms of *Recall@5* and *Recall@10*. From the results, we notice that ENTAGREC<sup>++</sup> outperforms ENTAGREC<sup>+</sup> on all datasets. On average, ENTAGREC<sup>++</sup> achieves 10.0% and 4.8% improvements over ENTAGREC<sup>+</sup> in terms of *Recall@5* and *Recall@10*, respectively. We also perform a Wilcoxon signed-rank test (Wilcoxon 1945) and Benjamini Yekutieli procedure (Benjamini and Yekutieli 2001) on each dataset, which indicates the improvement achieved by ENTAGREC<sup>++</sup> is statistically significant (adjusted  $p$ -value < 0.05). Thus, we demonstrate that additional tag component helps to improve ENTAGREC when additional tag provided.

We also compute Cliff’s Delta (Grissom and Kim 2005) which measures effect size to test if the differences in the recall values are *substantial*. The results are shown in Table 7.

### 9.3 Discussion

**Illustrative examples** Figure 9 shows a software object from STACK OVERFLOW with the ruby and rdoc tags. TAGCOMBINE cannot infer any of the tags. On the other hand ENTAGREC can infer all tags. This is one of the many examples where the performance of ENTAGREC is better than TAGCOMBINE.

Figure 10 presents a software object from STACK OVERFLOW with tags python, apache-spark, apache-spark-sql, and pyspark. The object is initially tagged with python, apache-spark, and apache-spark-sql. Later, the tag pyspark

**Table 6** The comparison of ENTAGREC<sup>++</sup> and ENTAGREC<sup>+</sup> in terms of *Recall@5* and *Recall@10*

Dataset	ENTAGREC <sup>++</sup>	ENTAGREC <sup>+</sup>
<i>Recall@5</i>		
STACK OVERFLOW	<b>0.905</b>	0.819
ASK UBUNTU	<b>0.737</b>	0.705
ASK DIFFERENT	0.852	0.685
SUPER USER	<b>0.908</b>	0.901
<i>Recall@10</i>		
STACK OVERFLOW	<b>0.968</b>	0.941
ASK UBUNTU	<b>0.87</b>	0.866
ASK DIFFERENT	<b>0.963</b>	0.932
SUPER USER	<b>0.956</b>	0.955

**Table 7** Effect sizes

Dataset	ENTAGREC <sup>++</sup> vs. ENTAGREC <sup>+</sup>
<i>Recall@5</i>	
STACK OVERFLOW	1
ASK UBUNTU	1
ASK DIFFERENT	1
SUPER USER	0.99
<i>Recall@10</i>	
STACK OVERFLOW	1
ASK UBUNTU	0.99
ASK DIFFERENT	1
SUPER USER	0.28

is added. ENTAGREC<sup>++</sup> can infer the tag `pyspark` given the initial set of tags, while ENTAGREC fails to do so.

**The impact of different MH on EnTagRec<sup>+</sup>** To understand the impact of parameter *MH* in Algorithm 1 on our approach, we test different values of *MH* of ENTAGREC<sup>+</sup> on the five datasets and see how the effectiveness of ENTAGREC<sup>+</sup> varies. Figure 11 presents the results, which show that the effectiveness of ENTAGREC<sup>+</sup> remains stable when we increase *MH* from 1 to 5. Since the difference in effectiveness is negligible for different *MH* values, we choose to set *MH* to 1 to reduce the computing cost.

**Stack exchange sites vs. FreeCode.** From the experimental results, we note that ENTAGREC<sup>+</sup> performs much better than TagCombine on STACK EXCHANGE Sites (i.e., STACK OVERFLOW, ASK UBUNTU, SUPER USER, ASK DIFFERENT), while it performs similarly to TagCombine on FREECODE. To understand why the performance of ENTAGREC<sup>+</sup> varies on different sites, we check the length (in words) of objects in the five datasets; the summary is shown in Table 8. We see that the length of objects in FREECODE is much shorter than that of STACK EXCHANGE sites. This may explain why the performance of ENTAGREC<sup>+</sup> on FREECODE is not as good as on the other sites. BIC is based on L-LDA, which usually requires training documents to be relatively long in order to achieve good results – c.f. Hong and Davison (2010). Unfortunately, objects in FREECODE are short, which results in poor results from BIC. To further verify our conjecture, we divide the objects into two groups. We sort the objects of each dataset by their length (in words) in ascending order. We take

## How do I add existing comments to RDoc in Ruby?

I've got all these comments that I want to make into 'RDoc comments', so they can be formatted appropriately and viewed using `rdoc`. Can anyone get me started on understanding how to use RDoc?

`ruby` `rdoc`

edited May 5 '13 at 0:24

 Taryn East

9,133 ●3 ●29 ●61

[add comment](#)

asked Aug 1 '08 at 13:38

 CodingWithoutComments

9,106 ●12 ●52 ●73

**Fig. 9** ENTAGREC correctly suggests tags `ruby` and `rdoc` for this STACK OVERFLOW question, while TAGCOMBINE does not



## Adding a new column in Data Frame derived from other columns (Spark)

▲ I'm using Spark 1.3.0 and Python. I have a dataframe and I wish to add an additional column which is derived from other columns. Like this,

2

```
>>old_df.columns
[col_1, col_2, ..., col_m]

>>new_df.columns
[col_1, col_2, ..., col_m, col_n]
```

where

```
col_n = col_3 - col_4
```

How do I do this in PySpark?

python apache-spark apache-spark-sql pyspark

share improve this question

edited Jul 10 '15 at 10:03



zero323

44.3k ● 11 ● 48 ● 76

asked Jul 10 '15 at 5:55



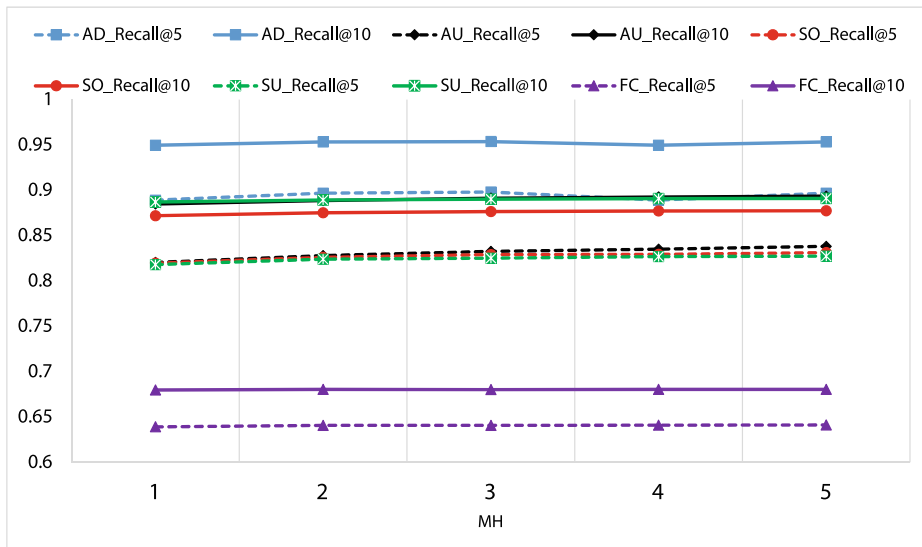
menorah84

738 ● 5 ● 20

add a comment

**Fig. 10** ENTAGREC<sup>++</sup> correctly suggests tags `pyspark` for this STACK OVERFLOW question given the initial tags `python`, `apache-spark`, and `apache-spark-sql`, while ENTAGREC does not

the top 50% of the objects as one group (i.e.  $Group_{short}$ ) and the rest as the another group (i.e.,  $Group_{long}$ ). We evaluate the effectiveness of ENTAGREC<sup>+</sup> on each group of the five datasets in terms of  $Recall@5$  and  $Recall@10$ . The results are shown at Table 9. We could see that ENTAGREC<sup>+</sup> consistently achieves better  $Recall@k$  scores on  $Group_{long}$ , which suggests that our approach is more effective on long objects rather than short ones.



**Fig. 11** The results of different values of  $MH$  on STACK OVERFLOW (SO), ASK UBUNTU (AU), ASK DIFFERENT (AD), SUPER USER (SU), and FREECODE (FC). Y axis is truncated (i.e., 0.6–1.0) and differences are smaller than they appear

**Table 8** The summary of the length (in words) of objects in the five dataset

Dataset	Entire dataset	$Group_{short}$	$Group_{long}$
STACK OVERFLOW	75.8	30.9	118.8
ASK UBUNTU	77.9	28.4	125.5
ASK DIFFERENT	60.9	26.4	93.6
SUPER USER	62.7	27.9	96.5
FREECODE	19.5	11.6	27.3

**Precision@k results** Aside from  $Recall@k$ ,  $Precision@k$  (Definition 2) has also been used to evaluate information retrieval techniques. In this paper, we focus on  $Recall@k$  as the evaluation metric. A similar decision was made by past tag recommendation studies (Zangerle et al. 2011; Xia et al. 2013). This is the case as the number of tags that are attached to an object is often small (much less than  $K$ ). Thus, the value of  $Precision@k$  is often very low and is not meaningful.

**Definition 2** Consider a set of  $n$  software objects. For each object  $o_i$ , let the set of its correct (i.e., ground truth) tags be  $Tags_i^{correct}$ . Also, let  $Tags_i^{topK}$  be the top-k ranked tags that are recommended by a tag recommendation approach for  $o_i$ . Average  $Precision@k$  for  $n$  is given by:

$$Precision@k = \frac{1}{n} \sum_{i=1}^n \frac{|Tags_i^{topK} \cap Tags_i^{correct}|}{|Tags_i^{topK}|} \quad (10)$$

Still, for the sake of completeness, we show the  $Precision@k$  results in Table 10. The results show that ENTAGREC<sup>+</sup> outperforms ENTAGREC and TAGCOMBINE on all datasets in terms of  $Precision@5$ . ENTAGREC<sup>+</sup> outperforms ENTAGREC on ASK UBUNTU, ASK DIFFERENT, SUPER USER, and FREECODE in terms of  $Precision@10$ . In terms of  $Precision@10$ , ENTAGREC<sup>+</sup> also outperforms TAGCOMBINE on four out of the five datasets. The difference between  $Precision@10$  of ENTAGREC<sup>+</sup> and TAGCOMBINE is small (i.e., 0.008). We also performed the Wilcoxon signed-rank test (Wilcoxon 1945) and Benjamini Yekutieli procedure (Benjamini and Yekutieli 2001). We found that in terms

**Table 9** The comparison between  $Group_{long}$  and  $Group_{short}$

Dataset	$Group_{long}$	$Group_{short}$
<i>Recall@5</i>		
STACK OVERFLOW	<b>0.912</b>	0.680
ASK UBUNTU	<b>0.882</b>	0.703
ASK DIFFERENT	0.954	0.807
SUPER USER	<b>0.908</b>	0.689
FREECODE	<b>0.635</b>	0.629
<i>Recall@10</i>		
STACK OVERFLOW	<b>0.956</b>	0.751
ASK UBUNTU	<b>0.940</b>	0.866
ASK DIFFERENT	<b>0.986</b>	0.781
SUPER USER	<b>0.965</b>	0.778
FREECODE	<b>0.680</b>	0.674

**Table 10** *Precision@5* and *Precision@10* for three approaches ENTAGREC<sup>+</sup>, ENTAGREC, and TAGCOMBINE.

Dataset	ENTAGREC	ENTAGREC <sup>+</sup>	TAGCOMBINE
<i>Precision@5</i>			
STACK OVERFLOW	0.346	<b>0.353</b>	0.221
ASK UBUNTU	0.358	0.361	0.251
ASK DIFFERENT	0.364	0.373	0.278
SUPER USER	0.376	0.380	0.285
FREECODE	0.382	0.396	0.381
<i>Precision@10</i>			
STACK OVERFLOW	0.187	0.187	0.151
ASK UBUNTU	0.196	0.197	0.158
ASK DIFFERENT	0.205	0.202	0.173
SUPER USER	0.201	0.207	0.177
FREECODE	0.240	0.241	<b>0.249</b>

The highest value is typeset in boldface

of *Precision@5*, ENTAGREC<sup>+</sup> significantly outperforms TAGCOMBINE. Also, in terms of *Precision@10*, ENTAGREC<sup>+</sup> significantly outperforms ENTAGREC on all datasets and TAGCOMBINE on STACK OVERFLOW, ASK UBUNTU, SUPER USER, and ASK DIFFERENT. For FREECODE, TAGCOMBINE significantly outperforms ENTAGREC<sup>+</sup> terms of *Precision@10*.

When the additional tags are given, the precision of ENTAGREC<sup>+</sup> and ENTAGREC<sup>++</sup> are presented at Table 11. ENTAGREC<sup>++</sup> outperforms ENTAGREC on all datasets in terms of *Precision@5* and *Precision@10*. We also performed the Wilcoxon signed-rank test (Wilcoxon 1945) and Benjamini Yekutieli procedure (Benjamini and Yekutieli 2001). We found that in terms of *Precision@5*, ENTAGREC<sup>++</sup> significantly outperforms ENTAGREC on all dataset.

To investigate if the differences in the precision values are substantial, we also compute Cliff’s Delta which measures effect size. The results are shown in Table 12. From the results we can conclude that ENTAGREC<sup>+</sup> substantially outperforms TAGCOMBINE and ENTAGREC on the STACK OVERFLOW, ASK UBUNTU, SUPER USER, and ASK DIFFERENT datasets (with at least medium effect sizes). For the FREECODE dataset, ENTAGREC<sup>+</sup> still substantially outperforms TAGCOMBINE in terms of *Precision@5*. However, TAGCOMBINE

**Table 11** *Precision@5* and *Precision@10* for two approaches ENTAGREC<sup>+</sup> and ENTAGREC<sup>++</sup>

Dataset	ENTAGREC <sup>++</sup>	ENTAGREC <sup>+</sup>
<i>Precision@5</i>		
STACK OVERFLOW	<b>0.225</b>	0.202
ASK UBUNTU	<b>0.202</b>	0.191
ASK DIFFERENT	0.225	0.181
SUPER USER	<b>0.249</b>	0.247
<i>Precision@10</i>		
STACK OVERFLOW	0.122	0.118
ASK UBUNTU	<b>0.122</b>	0.121
ASK DIFFERENT	<b>0.130</b>	0.125
SUPER USER	<b>0.133</b>	0.132

The highest value is typeset in boldface

**Table 12** Effect sizes (Precision)

Dataset	ENTAGREC <sup>+</sup> vs. ENTAGREC	ENTAGREC <sup>++</sup> vs. ENTAGREC <sup>+</sup>	ENTAGREC <sup>+</sup> vs. TAGCOMBINE
<i>Precision@5</i>			
STACK OVERFLOW	0.939	1	1
ASK UBUNTU	1	1	1
ASK DIFFERENT	1	1	1
SUPER USER	0.461	0.92	1
FREECODE	1	NA	1
<i>Precision@10</i>			
STACK OVERFLOW	-0.024	1	1
ASK UBUNTU	1	0.568	1
ASK DIFFERENT	1	1	1
SUPER USER	0.99	0.47	0.99
FREECODE	0.341	NA	-0.457

substantially outperforms ENTAGREC in terms of *Precision@10*. ENTAGREC<sup>++</sup> substantially outperforms ENTAGREC<sup>+</sup> on STACK OVERFLOW, ASK UBUNTU, SUPER USER, and ASK DIFFERENT datasets (with large and medium effect sizes). ENTAGREC<sup>+</sup> substantially outperforms ENTAGREC in terms of *Precision@5*. In terms of *Precision@10*, ENTAGREC<sup>+</sup> outperforms ENTAGREC on the ASK UBUNTU, SUPER USER, and ASK DIFFERENT datasets with large size and on the FREECODE with medium size.

**Efficiency** We find that ENTAGREC<sup>++</sup> runtimes for the training and deployment phases are reasonable. ENTAGREC<sup>++</sup>'s training time can mostly be attributed to training an L-LDA model in the Bayesian inference component of ENTAGREC<sup>++</sup>, which never exceeds 18 minutes (it is the maximum time across the ten iterations measured on the Stack Overflow dataset, the largest of the four). The Frequentist inference component is much faster; its runtime never exceeds 40 seconds (measured on the STACK OVERFLOW dataset). The training time of user information component and additional tag component never exceeds 1 minute. In the deployment phase, the average time ENTAGREC<sup>++</sup> takes to recommend a tag never exceeds 0.14 seconds.

**Retraining frequency** Since ENTAGREC<sup>++</sup> is efficient (i.e., model training can be completed in minutes), we can afford to retrain it daily. For example, a batch script can be run at a scheduled hour every day. Within a day, software objects and tagging behaviors are very likely to remain unchanged, and thus there is no need to retrain ENTAGREC<sup>++</sup> more frequently.

**ATC usage** Since more complete tags may shorten the time it takes for a question to be discovered and receive answer, we suggest to apply ATC just after a question is created with initial tags. Moreover, ATC can even be applied in real time, i.e., when users are entering tags.

**Threats to validity** Threats to external validity relate to the generalizability of our results. We have analyzed five popular software information sites (i.e., four STACK EXCHANGE

sites and FREECODE) and more than 160,000 software objects. In the future, we plan to reduce this threat further by analyzing even more software objects from more software information sites. As a threat to internal validity, we assume that the data in the software information sites are correct. To reduce the threat we only used older data—assuming people correct wrongly/poorly assigned tags. Also, two of our datasets (i.e., STACK OVERFLOW and FREECODE) were used in a past study (Xia et al. 2013). We use a lot of data and only consider tags that are used to label at least 50 objects to further reduce the impact of noise. Furthermore, manual inspection of a random sample of 100 STACK OVERFLOW objects (questions) revealed that only 1 had a clearly irrelevant tag (out of a total of 3 tags for that object).

Threats to construct validity relate to the suitability of our evaluation metrics. We have used *Recall@k* and *Precision@k* to evaluate our proposed approaches ENTAGREC and ENTAGREC<sup>++</sup> in comparison with other approaches. These measures are standard information retrieval measures used by prior tag recommendation studies, e.g., Al-Kofahi et al. (2010), Zangerle et al. (2011), and Xia et al. (2013). We have also performed statistical test and effect size test to check if the differences in *Recall@k* and *Precision@k* are significant and substantial. Thus, we believe there is little threat to construct validity.

## 10 Related Work

**Tag Recommendation:** Al-Kofahi et al. proposed TAGREC which recommends tags in work item systems (e.g., IBM Jazz) (Al-Kofahi et al. 2010). There are a number of studies from the data mining research community, that recommend tags for social media sites like Twitter, Delicious, and Flickr (Jäschke et al. 2007; Sigurbjörnsson and van Zwol 2008; Zangerle et al. 2011). Among these studies, the work by Zangerle et al. is the latest approach to recommend hashtags for short messages in Twitter (Zangerle et al. 2011). Xia et al. proposed TAGCOMBINE, which combines three components: a multi-label ranking component, a similarity-based ranking component, and a tag-term based ranking component (Xia et al. 2013). Xia et al. have shown that TAGCOMBINE outperforms TAGREC and Zangerle et al.’s approach in recommending tags in software information sites.

The closest work to ours is TAGCOMBINE proposed by Xia et al. which is also the prior state-of-the-art work (Xia et al. 2013). There are a number of technical differences between ENTAGREC, proposed in our preliminary work (Wang et al. 2014), and TAGCOMBINE. ENTAGREC combines two components: a Bayesian inference component that employs Labeled LDA (BIC), and an enhanced frequentist inference component that removes unrelated words with the help of a parts-of-speech (POS) tagger, and finds associated tags with a spreading activation algorithm (FIC). Our BIC is related to the multi-label ranking component of TAGCOMBINE since both of them employ Bayesian inference. The multi-label ranking component of TAGCOMBINE constructs many one-versus-rest Naive Bayes classifiers, one for each tag. Each Naive Bayes classifier simply predicts the likelihood of a software object to be assigned a particular tag. In ENTAGREC, we construct only one classifier which is a *mixture model* that considers all tags together. Mixture models have been shown to outperform one-versus-rest traditional multi-label classification approaches (Ramage et al. 2009; Ghamrawi and McCallum 2005; Puurula 2011). Also, our FIC removes unrelated words (using POS tagger) and finds associated tags (using spreading activation) while none of the three components of TAGCOMBINE perform these. We have compared our approach with TAGCOMBINE, on four datasets: STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, and FREECODE. We show that our approach outperforms TAGCOMBINE

on three datasets (i.e., STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT), and performs as well as TAGCOMBINE on one dataset (i.e., FREECODE). ENTAGREC<sup>++</sup> extends ENTAGREC by including two additional components, User Information Component (UIC) and Additional Tag Component (ATC), which boosts performance further.

**Tagging in software engineering** The need for automatic tag recommendation has been recognized both by practitioners (Warbox 2009; Her 2011; Jmac 2013) and by researchers. Aside from tag recommendation studies mentioned above, there are several software engineering studies that also analyze tagging and leverage tags for various purposes. Treude et al. performed an empirical study on the impact of tagging on a large project with 175 developers over a two years period (Treude and Storey 2009). Wang et al. analyzed tags of projects in FREECODE, inferred the semantic relationships among the tags, and expressed the relationships as a taxonomy (Wang et al. 2012). Thung et al. detected similar software applications using software tags (Thung et al. 2012). Storey et al. proposed an approach called TagSEA that allows one to create, edit, navigate, and manage annotations in source code (Storey et al. 2009). Treude et al. performed an empirical study on several professional projects that involved more than 1,000 developers, and found that tagging can play an important role in the development process (Treude and Storey 2012). They found that tags are helpful in articulation work, finding of tasks, and exchange of information. Cabot et al. conducted an empirical study on the labels that are used to classify issues on issue tracking system and they found that the use of such labels improves issue resolution process (Cabot et al. 2015). Wang et al. have demonstrated that the practice of tagging helps in assisted tracing (a process where analysts inspect results produce by automated traceability techniques) (Wang et al. 2015). Through a user study, they find that tagging is readily adopted by analysts and improve the quality of the trace matrices produced at the end of the study.

Furthermore, several studies of STACK OVERFLOW have used tags to focus on questions or answers pertaining to a certain technology (Bazelli et al. 2013; Vasilescu et al. 2013) or to enhance studies of related websites such as GITHUB (Pletea et al. 2014) or Wikipedia (Joorabchi et al. 2015).

## 11 Conclusion and Future Work

In this work, we propose a novel approach to recommend tags to software information sites. Our approach, named ENTAGREC<sup>++</sup>, an enhanced version of ENTAGREC, learns from tags of historical software objects to infer tags of new software objects. To recommend tags, ENTAGREC<sup>++</sup> enhances ENTAGREC by adding two more inference components. One, named user information component (UIC), makes use of historical tagging information peculiar to a user to infer tags for a current software object the user creates. Another one, named additional tag component (ATC), makes use of an initial set of tags given by a user to recommend additional tags better. ENTAGREC<sup>++</sup> composes the four components by finding the best weights that optimize the performance of ENTAGREC<sup>++</sup> on a training dataset. We evaluate the performance of ENTAGREC<sup>++</sup> on four datasets, STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, SUPER USER, and FREECODE, which contain 47,688, 39,231, 37,354, and 13,351 software objects, respectively. We find that that without leveraging ATC, our approach (named ENTAGREC<sup>+</sup>) achieves *Recall@5* scores of 0.821, 0.822, 0.891 and 0.651, and *Recall@10* scores of 0.873, 0.886, 0.956 and 0.761, on STACK OVERFLOW, ASK UBUNTU, ASK DIFFERENT, SUPER USER, and FREECODE, respectively. In terms of

*Recall@5* and *Recall@10*, averaging across the 4 datasets, ENTAGREC<sup>+</sup> improves TAG-COMBINE (Xia et al. 2013), which is the prior state-of-the-art approach, by 29.1% and 14.2% respectively. In addition, with ATC, ENTAGREC<sup>++</sup> achieves a 13.1% improvement over ENTAGREC<sup>+</sup> in terms of *Recall@5*. We have published the code and datasets that we used online.<sup>16</sup> Admittedly, we have only tested our approach on Stack Exchange sites and FreeCode.

As future work, we plan to reduce the threats to validity by experimenting with more software objects from more software information sites. In this paper, we only consider the major tag re-editing scenario – tag addition (see Table 1). In the future, we also plan to support tag deletion and tag correction. Furthermore, we plan to improve the *Recall@5* and *Recall@10* of ENTAGREC further by investigating cases where ENTAGREC<sup>++</sup> is inaccurate, and by building a more sophisticated machine learning solution.

## References

- Al-Kofahi JM, Tamrawi A, Nguyen TT, Nguyen HA, Nguyen TN (2010) Fuzzy set approach for automatic tagging in evolving software ICSM, pp 1–10
- Antoniol G, Canfora G, Casazza G, De Lucia A, Merlo E (2002) Recovering traceability links between code and documentation. IEEE Trans Softw Eng 28(10):970–983
- Asuncion HU, Asuncion AU, Taylor RN (2010) Software traceability with topic modeling ICSE, pp 95–104
- Baldi P, Lopes CV, Linstead E, Bajracharya SK (2008) A theory of aspects as latent topics OOPSLA, pp 543–562
- Bazelli B, Hindle A, Stroulia E (2013) On the personality traits of stackoverflow users. In: 2013 IEEE international conference on software maintenance, pp 460–463
- Benjamini Y, Yekutieli D (2001) The control of the false discovery rate in multiple testing under dependency. Ann Stat 29:1165–1188
- Bergstra J, Bengio Y (2012) Random search for hyper-parameter optimization. JMLR 13:281–305
- Bindelli S, Criscione C, Curino C, Drago ML, Eynard D, Orsi G (2008) Improving search and navigation by combining ontologies and social tags. In: On the move to meaningful internet systems, OTM 2008 Workshops, OTM confederated international workshops and posters, ADI, AWeSoMe, COMBEK, EI2N, IWSSA, MONET, OnToContent + QSI, ORM, PerSys, RDDS, SEMELS, and SWWS 2008, Monterrey, Mexico, November 9-14, 2008. Proceedings, pp 76–85
- Blei DM, Ng AY, Jordan MI (2003) Latent dirichlet allocation. JMLR, 993–1022
- Brandt J, Guo PJ, Lewenstein J, Dontcheva M, Klemmer SR (2009) Two studies of opportunistic programming: interleaving web foraging, learning, and writing code CHI. ACM, pp 1589–1598
- Cabot J, Izquierdo JLC, Cosentino V, Rolandi B (2015) Exploring the use of labels to categorize issues in open-source software projects. In: 22nd IEEE international conference on software analysis, evolution, and reengineering, SANER 2015. Montreal, QC, Canada, March 2-6, 2015, pp 550–554
- Capobianco G, Lucia AD, Oliveto R, Panichella A, Panichella S (2013) Improving IR-based traceability recovery via noun-based indexing of software artifacts. J Softw Evol Process 25(7):743–762
- Cress U, Held C, Kimmerle J (2013) The collective knowledge of social tags: direct and indirect influences on navigation, learning, and information processing. Comput Educ 60(1):59–73
- Crestani F (1997) Application of spreading activation techniques in information retrieval. Artif Intell Rev 11(6):453–482
- Gelman A, Carlin J, Stern H, Rubin D (2003) Bayesian data analysis. CRC Press
- Ghamrawi N, McCallum A (2005) Collective multi-label classification CIKM, pp 195–200
- Golder SA, Huberman BA (2006) Usage patterns of collaborative tagging systems. J Inf Sci 32(2):198–206
- Grissom RJ, Kim JJ (2005) Effect sizes for research. A broad practical approach
- Han J, Kamber M, Pei J (2011) Data mining: concepts and techniques. Morgan Kaufmann Publishers Inc
- Held C, Kimmerle J, Cress U (2012) Learning by foraging: the impact of individual knowledge and social tags on web navigation processes. Comput Hum Behav 28(1):34–40

- Hong L, Davison BD (2010) Empirical study of topic modeling in twitter. In: Proceedings of the first workshop on social media analytics, SOMA '10, pp 80–88
- Jäschke R, Marinho LB, Hotho A, Schmidt-Thieme L, Stumme G (2007) Tag recommendations in folksonomies PKDD
- Jmac (2013) Select and display 'suggested tags' for all posts based on related questions (or other logic). <http://meta.stackexchange.com/q/196702/182512>
- Joorabchi A, English M, Mahdi AE (2015) Automatic mapping of user tags to wikipedia concepts: the case of a q&a website. *Stackoverflow*. *J Inf Sci* 41(5):570–583
- Her J (2011) Tag recommendations for stack overflow. <http://meta.stackexchange.com/q/88611/182512>
- Lukins SK, Kraft NA, Etkorn LH (2010) Bug localization using latent dirichlet allocation. *Inf Softw Technol* 52(9):972–990
- Panichella A, Dit B, Oliveto R, Di Penta M, Poshyvanyk D, Lucia AD (2013) How to effectively use topic models for software engineering tasks? An approach based on genetic algorithms ICSE, pp 522–531
- Pletea D, Vasilescu B, Serebrenik A (2014) Security and emotion: Sentiment analysis of security discussions on github. In: Proceedings of the 11th working conference on mining software repositories, MSR 2014. ACM, New York, pp 348–351
- Porter MF (1997) An algorithm for suffix stripping Readings in information retrieval. Morgan Kaufmann, pp 313–316
- Puurula A (2011) Mixture models for multi-label text classification. In: 10th New Zealand computer science research student conference
- Ramage D, Hall D, Nallapati R, Manning CD (2009) Labeled lda: a supervised topic model for credit attribution in multi-labeled corpora. In: EMNLP '09, pp 248–256
- Rebouças M, Pinto G, Ebert F, Torres W, Serebrenik A, Castor F (2016) An empirical study on the usage of the swift programming language. In: 2016 IEEE 23rd international conference on software analysis, evolution, and reengineering (SANER), pp 634–638
- Samaniego FI (2010) A comparison of the bayesian and frequentist approaches to estimation. Series in Statistics, Springer
- Shokripour R, Anvik J, Kasirun ZM, Zamani S (2013) Why so complicated? Simple term filtering and weighting for location-based bug report assignment recommendation MSR
- Sigurbjörnsson B, van Zwol R (2008) Flickr tag recommendation based on collective knowledge WWW '08, pp 327–336
- Storey M-A, Ryall J, Singer J, Myers D, Cheng L-T, Muller M (2009) How software developers use tagging to support reminding and refinding. *IEEE Trans Softw Eng* 35(undefined):470–483
- Storey M-A, Treude C, van Deursen A, Cheng L-T (2010) The impact of social media on software engineering practices and tools. In: FoSER '10, pp 359–364
- Thung F, Lo D, Jiang L (2012) Detecting similar applications with collaborative tagging. In: ICSM, pp 600–603
- Toutanova K, Klein D, Manning CD, Singer Y (2003) Feature-rich part-of-speech tagging with a cyclic dependency network. In: HLT-NAACL
- Treude C, Storey M-A (2009) How tagging helps bridge the gap between social and technical aspects in software development. In: ICSE '09, pp 12–22
- Treude C, Storey M-A (2012) Work item tagging: communicating concerns in collaborative software development. *IEEE Trans Softw Eng* 38(1):19–34
- Vasilescu B, Serebrenik A, Devanbu PT, Filkov V (2014) How social Q&A sites are changing knowledge sharing in open source software communities. In: CSCW, pp 342–354
- Vasilescu B, Serebrenik A, van den Brand MGJ (2013) The babel of software development: linguistic diversity in open source. In: Jatowt A, Lim E-P, Ding Y, Miura A, Tezuka T, Dias G, Tanaka K, Flanagan A, Dai BT (eds) Proceedings of the social informatics: 5th international conference, SocInfo 2013, Kyoto, Japan, November 25-27, 2013. Springer International Publishing, pp 391–404
- Vogt CC, Cottrell GW (1999) Fusion via a linear combination of scores. *Inf Retr* 1(3):151–173
- Wang S, Lo D, Jiang L (2012) Inferring semantically related software terms and their taxonomy by leveraging collaborative tagging. In: ICSM, pp 604–607
- Wang S, Lo D, Vasilescu B, Serebrenik A (2014) EnTagRec: an enhanced tag recommendation system for software information sites. In: 30th IEEE international conference on software maintenance and evolution, Victoria, BC, Canada, September 29 - October 3, 2014. IEEE Computer Society, pp 291–300
- Wang W, Niu N, Liu H, Wu Y (2015) Tagging in assisted tracing. In: 2015 IEEE/ACM 8th international symposium on software and systems traceability, pp 8–14
- Wang X-Y, Xia X, Lo D (2015) Tagcombine: recommending tags to contents in software information sites. *J Comput Sci Technol* 30(5):1017–1035



- Warbox D (2009) Auto-tagging. <http://meta.stackoverflow.com/questions/1377/auto-tagging>
- Wilcoxon F (1945) Individual comparisons by ranking methods. *Biom Bull* 1(4):80–83
- Xia X, Lo D, Wang X, Zhou B (2013) Tag recommendation in software information sites. In: *MSR '13*, pp 287–296
- Zangerle E, Gassler W, Specht G (2011) Using tag recommendations to homogenize folksonomies in microblogging environments. In: *SocInfo'11*, pp 113–126
- Zubiaga A (2012) Enhancing navigation on wikipedia with social tags. *CoRR*, arXiv:1202.5469



**Shaowei Wang** is a Postdoc in the Software Analysis and Intelligence (SAIL) Lab at Queens University, Canada. He obtained his PhD from Singapore Management University, and BSc from Zhejiang University. His research interests include code mining and recommendation, software maintenance, developer forum analysis, and mining software repositories. He has served as a reviewer of a number of high-quality journals (e.g., *IEEE Transaction on Software Engineer*, *Empirical Software Engineering*). More information at: <https://sites.google.com/site/wswshaoweiwang/>.



**David Lo** his PhD degree from the School of Computing, National University of Singapore in 2008. He is currently an Associate Professor in the School of Information Systems, Singapore Management University. He has close to 10 years of experience in software engineering and data mining research and has more than 200 publications in these areas. He received the Lee Foundation Fellow for Research Excellence from the Singapore Management University in 2009, and a number of international research awards including several ACM distinguished paper awards for his work on software analytics. He has served as general and program co-chair of several well-known international conferences (e.g., *IEEE/ACM International Conference on Automated Software Engineering*), and editorial board member of a number of high-quality journals (e.g., *Empirical Software Engineering*).



**Bogdan Vasilescu** is an assistant professor at Carnegie Mellon University's School of Computer Science, where he is engaged in interdisciplinary research at the intersection of software engineering and social computing. Bogdan explores large-scale software-related data using a mixture of quantitative and qualitative methods, to develop and validate theories about the processes involved in software engineering and computer-supported collaborative work. Prior to joining CMU, Bogdan was a postdoctoral researcher at University of California, Davis. He received his PhD and MSc in Computer Science at Eindhoven University of Technology, both with cum laude distinction.



**Alexander Serebrenik** (PhD KULeuven, Belgium 2003) is an associate professor of software evolution at Eindhoven University of Technology, The Netherlands. He studies how software systems and their developers' communities change with time, and focuses both on the social and on the technical aspects of this process. He has co-edited a book on Software evolution, co-authored 30 journal articles, more than 90 conference papers and acted as the steering committee chair, general chair and program co-chair of international conferences on software maintenance and evolution.