



Finding associations between natural and computer languages: A case-study of bilingual LDA applied to the bleeping computer forum posts[☆]

Kundi Yao^{a,*}, Gustavo A. Oliva^a, Ahmed E. Hassan^a, Muhammad Asaduzzaman^b, Andrew J. Malton^c, Andrew Walenstein^c

^a Software Analysis and Intelligence Lab (SAIL) at Queen's University, Kingston, Canada

^b Lakehead University, Thunderbay, Canada

^c BlackBerry Ltd, Canada

ARTICLE INFO

Article history:

Received 1 February 2022

Received in revised form 16 January 2023

Accepted 20 February 2023

Available online 24 February 2023

Keywords:

Technical support

Logs

LDA

Multilingual LDA

Topic models

ABSTRACT

In the context of technical support, trails of technical discussions often contain a mixture of natural language (e.g., English) and software log excerpts. Uncovering latent links between certain problems and log excerpts that are often requested during the discussions of those problems enables the construction of a valuable knowledge base. Nevertheless, uncovering such latent links is challenging because English and software logs are two fundamentally different languages. In this paper, we investigate the suitability of multilingual LDA models to address the problem at hand. We study three models, namely: *enriched LDA* (M_+), *two-layer LDA* (M_{2L}), and *off-the-shelf bilingual LDA* (M_{bi}). We use approximately 8K discussion threads from a Bleeping Computer forum as our dataset. We observe that M_{2L} performs the best overall, although it yields a substantially coarser-grained view of the discussed themes in the threads (20 topics, 0.3% of the documents). We also note that M_+ outperforms M_{bi} achieving higher coherence, lower perplexity, and higher cross-lingual coverage ratio. We invite future studies to qualitatively assess the quality of the topics produced by the LDA models, such that the feasibility of employing such models in practice can be better determined.

© 2023 Elsevier Inc. All rights reserved.

1. Introduction

Software logs are a valuable resource in the diagnosis of failures in software systems. When end-users describe software failures, they typically combine natural language (e.g., English) and logs. Examples of channels in which end-users describe technical problems with a combination of natural language and logs include (i) tickets in IT ticket systems (e.g., Freshdesk¹), (ii) online chat with customer technical support (e.g., LiveChat²), and (iii) posts in technical support forums (e.g., Bleeping Computer³). We refer to discussion threads happening through these channels as *support threads*.

There might be latent associations between pieces of text written in natural language and certain log excerpts in support thread datasets. For instance, there might be multiple reports of a specific operating system crash that is associated with the log excerpt VIDEO_TDR_FAILURE (igdkmd64.sys). Discovering such latent associations enables the construction of a knowledge base that links problem discussions to log excerpts that are typically associated (requested/provided) with those discussions. Ultimately, such a knowledge base speeds up failure debugging and technical support service quality. The central question that we address in this paper is thus **how can one find associations between natural language and logs?**

Similar to previous research where source code and natural language are treated as separate channels (Casalnuovo et al., 2020), the problem of establishing connections between text and logs can be seen as the problem of establishing connections between *pieces of text that are written in different languages*. More specifically, system logs produced by a given tool (e.g., Mini-ToolBox⁴) can be seen as a semi-formal language that has its

[☆] Editor: Dr Martin Pinzger.

* Corresponding author.

E-mail addresses: kundi@cs.queensu.ca (K. Yao), gustavo@cs.queensu.ca (G.A. Oliva), ahmed@cs.queensu.ca (A.E. Hassan), md.asaduzzaman@lakeheadu.ca (M. Asaduzzaman), amalton@blackberry.com (A.J. Malton), awalenstein@blackberry.com (A. Walenstein).

¹ <http://freshdesk.com>

² <https://www.livechat.com>

³ <https://www.bleepingcomputer.com>

⁴ <https://www.bleepingcomputer.com/download/minitoolbox>

own syntax, vocabulary, and semantics. Recently, bilingual LDA models (Yang et al., 2019; Hao et al., 2018; Balikas, 2017) have been proposed to address this very problem. Although monolingual LDA has been extensively used in Software Engineering (SE) research (Chen et al., 2016), the extent to which bilingual LDA works for SE data is unknown. In particular, in the context of support threads, natural language and logs are often loosely related (as opposed to being direct translations of one another).

In this work, we conduct an exploratory study where we investigate whether bilingual LDA models are suitable for uncovering associations between English text and system logs. We evaluate three candidate bilingual models, namely: (i) enriched LDA (M_+), (ii) two-layer LDA (M_{2L}), and (iii) Mallet Bilingual LDA (M_{bi}). The first two models were designed by us and the third one is an off-the-shelf solution provided by the Mallet topic modeling toolkit.⁵ More specifically, we evaluate and compare these models along with three perspectives: internal quality, inferential power, and language alignment. The internal quality measures the semantic similarities of top words within the topic. The inferential power refers to the competency of a model in estimating unseen documents. And the language alignment measures the proportion of topics with top words from different languages. To evaluate the candidate models, we use data from the Bleeping Computer forum. Bleeping Computer is an information security and technology news website that offers a technical support forum where end-users can post their problems and get help from senior Bleeping Computer community members. Our research questions and associated key results are as follows:

(RQ1) What is the achieved level of internal quality by each candidate model? LDA models need to produce a coherent set of topics (i.e., the set of words that describe a topic should be semantically connected). We evaluated the candidate models using the *topic coherence* metric (the higher, the better), which scores each topic by measuring the degree of semantic similarity between high scoring words in that topic. We observed that: *The M_{2L} model achieves the highest level of internal quality (highest coherence).*

(RQ2) What is the achieved level of inferential power by each candidate model? LDA models are useful when they can correctly guess the topic of a new, unseen document. We evaluated the candidate models using the *perplexity* metric (the lower the better), which measures how “perplexed” an LDA model gets when it encounters the words of a new, unseen document. We observed that:

The M_{2L} model achieves the highest inferential power (lowest perplexity).

(RQ3) What is the achieved level of language alignment by each candidate model? In the context of support threads, LDA models become more useful when they produce topics that contain both English words and log words. We evaluated our candidate LDA models using our proposed *language alignment* metric. This metric corresponds to the ratio of LDA topics that contain both English words and log words. We observed that:

The M_{2L} model achieves the highest language alignment (highest cross-lingual coverage ratio).

We conducted a follow-up study to determine which model yields the best balance between the three quality attributes of interest. To quantify such a balance, we define an *aggregated score* that ranges from zero to one. Our conclusions are as follows. Although M_{2L} unsurprisingly achieves the highest aggregated score (0.87), it provides a substantially coarser-grained view of the discussed themes in the documents. Comparison between M_+ (aggregated score of 0.77) and M_{bi} (aggregated score of 0.69) reveals that M_+ achieves a higher overall performance. When single

measure is preferred, M_+ produces topics with higher coherence, while M_{bi} produces topics with lower perplexity. Therefore, the choice between those two models depends on which measure one prioritizes.

Paper structure. Section 2 describes how we collected data from the Bleeping Computer forum. Section 3 describes our candidate LDA models. Section 4 shows the motivation, the approach, and the results associated with each of our research questions. Section 5 describes a follow-up study in which we determine what candidate model achieves the best balance between internal quality, inferential power, and language alignment. Section 6 describes related work revolving around multilingual LDA as well as LDA being applied to Software Engineering. Section 7 discusses the threats to the validity of our study. Finally, Section 8 states our concluding remarks.

2. Data collection

We collect data from the community question-answering forum of the Bleeping Computer website. In the following, we describe this data source as well as our approach for collecting the data.

2.1. Data source

The Bleeping Computer forum is organized into a number of *sections* and *subsections*. For instance, there exists a *Microsoft Windows Support* section, which centers discussion around that operating system. Its subsections include (but are not limited to) *Windows Crashes and Blue Screen of Death (BSOD) Help and Support*, *Windows 8 and 8.1*, and *Windows 10 Support*. A subsection might contain other subsections. For instance, the *Windows 10 Support* subsection contains two other subsections: *Windows 10 Discussion* and *Windows 10 Insider Preview Builds*.

Users can ask technical *questions* by starting a *thread* on a given forum subsection. Any user can respond to a question by providing one or more *replies*. In a reply, forum members typically ask for question clarifications, suggest the question-poster to run a diagnosis software tool, examine the produced logs by a diagnosis software tool, and offer potential solutions to the problem. We use the term *post* to refer to either the question or any reply to it. A thread thus comprises the set of posts corresponding to the question and its associated replies. A thread contains English text and may contain system logs. For the sake of simplicity, we henceforth refer to the former as *text* and to the latter as *logs*.

2.2. Approach

In the following, we describe our data collection approach. An overview is shown in Fig. 1.

Step 1: Determine which threads to collect. We collect all threads that (i) contain log snippets that are produced by the MiniToolBox tool and (ii) belong to the *am-i-infected-what-do-i-do* forum section. We explain our rationale below.

– *Step 1a: Choose a log tool to focus on.* The produced logs from each log tool differ in terms of structure, vocabulary, and grammar. We thus consider that each log tool produces logs in a different language. Therefore, we choose a log system to focus on. While inspecting a random sample of threads from Bleeping Computer, we notice that users typically submit logs that are produced by the following log tools: Trend Micro HijackThis, DDS, MiniToolBox, ComboFix, and OTL. To choose the tool, we take a step back and reflect on the goals of our study. We wish to uncover associations between text and logs. As such, it does not make sense to pose constraints on how the text or logs should look like. In the most general case, logs will contain

⁵ <http://mallet.cs.umass.edu>

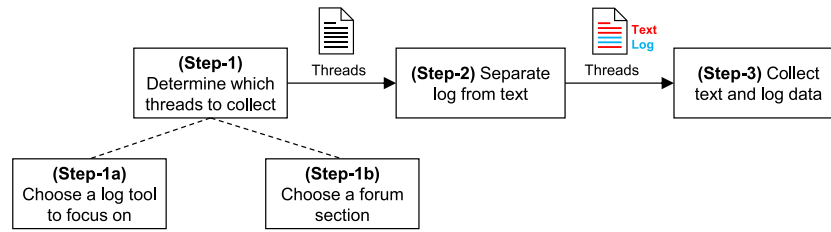


Fig. 1. An overview of our the data collection process.

collected troubleshooting data from several sources (e.g., different hardware components and different parts of the operating system). Hence, we opt to choose the log tool that produces the most complete type of log. Upon an analysis of each tool, we selected MiniToolBox. This tool focuses on troubleshooting network-related problems. It is also capable of producing logs that pack information from several parts of the Microsoft Windows operating system, including: content of the host files, IP configurations, Winsock entries, last 10 Event Viewer log errors, installed programs, installed devices, registered user accounts, disk partitions, memory size, and Windows restore points. In contrast, the logs of other tools tend to be simpler and contain less heterogeneous types of information. For instance, Trend Micro HijackThis aids in the detection of Windows malware by scanning a computer and producing a log of the running processes and suspicious Windows registry keys.

– *Step 1b: Choose a forum section.* We collect data from the *am-i-infected-what-do-i-do* forum section, which contains 36,477 threads. The rationale is that this section has the largest ratio of threads containing log excerpts (21.8%) among all sections in the Bleeping Computer forum. Next, since MiniToolBox logs typically start with either of the two following sentences: “MiniToolBox by Farbar Version:dd-mm-yyyy” or “Scan result of Farbar Recovery Scan Tool (FRST) (x64) Version:dd-mm-yyyy”, we select the threads that contain MiniToolBox logs based on such patterns. In total, we collect 7,969 threads in this step.

Step-2: Separate log from text. Since we only consider MiniToolBox logs and text, logs generated from other diagnosis software tools are considered noise data which shall be removed. In this step, we discuss our approaches to identifying both *complete* and *incomplete* log excerpts.

– *Identify complete log excerpt.* Diagnosis software diagnosis tools usually include a prefix pattern with the tool’s name in its generate logs, as we discussed in Step 1b. In this step, we summarize patterns from different diagnosis software tools and leverage such patterns to separate logs from text. First, we study the diagnosis software tools that are both hosted on the Bleeping Computer website⁶ (*internal tools*) and referred from external sources (e.g., ESET Online Scanner⁷) (*external tools*). From all *internal tools*, we choose those that occur in studied threads through keyword matching. We then extract and examine external links in all posts to summarize all *external tools*. Finally, we create regular expression rules for each tool based on the characteristics and availability of their corresponding logs, and such rules are leveraged to identify logs from texts.

– *Identify incomplete log excerpt.* We observe incomplete log excerpts (e.g., part of a diagnosis log) which are not identified by regular expression rules in certain posts. Inspired by the work of Bettenburg et al. (2011), we leverage a spellchecker to determine the proportion of unknown words (i.e., words not included in the English dictionary) in a data excerpt. Data excerpts with n

$> t$ unknown terms are considered log data and vice versa, where t is some predefined threshold. In order to determine a suitable threshold t , we first manually build a dataset from our studied threads containing 500 text excerpts and 500 log excerpts. Next, we (i) select 10 random samples with 50 text excerpts and 50 log excerpts, (ii) run *pyspellchecker*,⁸ (iii) and search for the threshold that maximizes the accuracy of log and text identification on each sample (from 0.3 to 1.0, with steps of 0.01). We use the median of the 10 threshold values as our optimal threshold (0.685).

Step 3: Collect text and log data. For each selected thread, we collect text and log data. Text data comprises the English text found in all posts of a given thread, including the question being asked (thread name). Log data is obtained as per Step-2. We collect all MiniToolBox log snippets that appear in a given thread.

Summary

- Data source: Bleeping Computer.
- Dataset: All data in the *am-i-infected-what-do-i-do* forum as of December 15th, 2020.
- Pieces of collected data: Discussion threads containing both English text and log snippets produced by MiniToolBox.
- Total number of collected threads: 7,969.

3. Computation of the LDA models

In this section, we first introduce the data preprocessing procedures that we performed according to the different natures of log and text. Next, we describe how we compute our candidate LDA models. Finally, we demonstrate how we implement these models.

3.1. LDA-oriented data preprocessing

Each input document corresponds to a thread that we collected from the Bleeping Computer forum. We remove line breaks and any HTML tags from both log and text parts of each input document. We also replace sequences of whitespaces with a single space and remove special characters (e.g., '). Next, we apply tokenization and separate the words.

Text. We perform *stop word removal* and *lemmatization*, which are common LDA-oriented data preprocessing techniques for natural languages, to further process text words.

Logs. Logs are computer-generated and thus have their own structure, vocabulary, and grammar. Compared to natural language, log data tends to be more repetitive and with fewer unique words (Yao et al., 2020). Logs also comprise system-specific words (e.g., *block_9*) that consist of more than alphabetical characters (which form most natural language words). Furthermore, log words (e.g., *getProductID*) do not follow the same composition and syntax as natural language. Therefore, the common text preprocessing rules such as *lemmatization* are not applicable to logs. Hence, we choose *not* to preprocess logs.

⁶ <https://www.bleepingcomputer.com/download/windows/>

⁷ <https://www.eset.com/int/home/online-scanner/>

⁸ <https://pypi.org/project/pyspellchecker/>

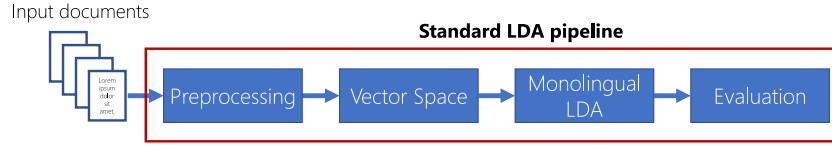


Fig. 2. Standard LDA pipeline.

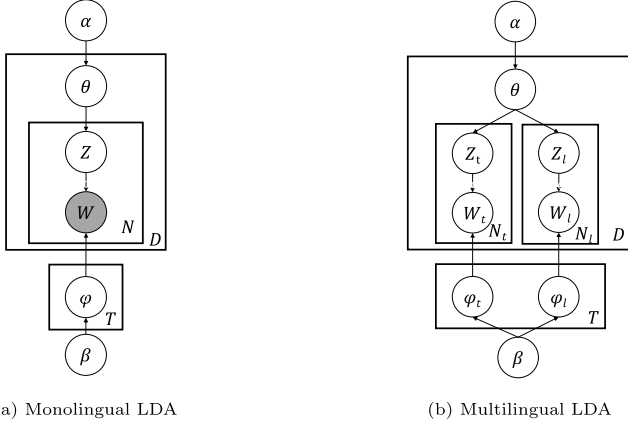


Fig. 3. Plate notation of adopted LDA models. α and β are the parameters of the Dirichlet prior over θ (per-document topic distribution) and φ (per-topic word distribution). W and Z denotes each word and the topic of each word. (e.g., $W_{d,n}$ is the n^{th} word of the d^{th} document and $Z_{i,j}$ is the topic of the j^{th} word of the i^{th} document). N , D , and T represents the total words in a document, total number of documents in a corpus, and the total number of topics, respectively.

3.2. Candidate LDA models

In the context of topic modeling, a topic model is meant to *discover* topics in a corpus (collection of documents), and a topic is a distribution over terms that is biased around those associated with a single theme (Wang and Blei, 2011). Specifically, Latent Dirichlet Allocation (LDA) (Blei et al., 2003) is a probabilistic topic model that has been recurrently studied and adopted in prior software engineering research (Chen et al., 2016). Fig. 2 shows a standard topic modeling pipeline with monolingual LDA as a reference model.

To handle documents containing both natural language (English) and logs, we design three new LDA model pipelines. Our first two pipelines use monolingual LDA models, while the third uses a multilingual LDA model. The plate notation of monolingual and multilingual LDA models is shown in Fig. 3.

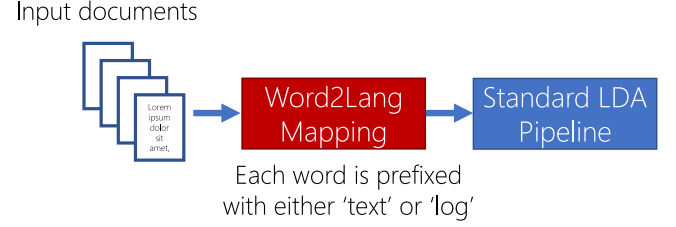
Our three candidate LDA-based topic model pipelines are shown below (for brevity purposes, we refer to them as our *candidate LDA models*):

- **Enriched LDA (M_+)**: adds a prefix to the words of the original documents, such that text and log words can be differentiated.
- **Two-layer LDA (M_{2L})**: uses internal LDA runs to summarize documents in an effort to counterbalance the prevalence of text in comparison to logs.
- **Off-the-shelf bilingual LDA (M_{bi})**: treats text and logs as different languages and uses Mallet's bilingual LDA implementation.

In the following, we describe our candidate models in detail. For each candidate model, we describe (i) the rationale behind its design and (ii) how we construct it.

3.2.1. Enriched LDA (M_+)

An overview of our Enriched LDA (M_+) candidate model is shown in Fig. 4.

Fig. 4. Enriched LDA (M_+) pipeline.

- **Rationale**: Standard LDA models are monolingual, which means that they cannot align topics across languages. Nevertheless, standard LDA models can still produce topics containing words from multiple languages if run on multilingual input documents. Such multilingual topics can uncover relationships between words from different languages (e.g., text and logs). The key drawback in this scenario is the lack of an explicit indication of the language in which a topic word is written in. Our Enriched LDA model addresses such a drawback.

- **Construction**: Our Enriched LDA model performs a preemptive data preprocessing step before pushing the input documents to the standard LDA pipeline. Each text word is prefixed with the term *text_* and each log word is prefixed with the term *log_*. We refer to this prefixing process as *Word2Lang* mapping. As a result, the standard LDA pipeline produces topics containing words that explicitly indicate the language in which they are written.

Pros and cons:

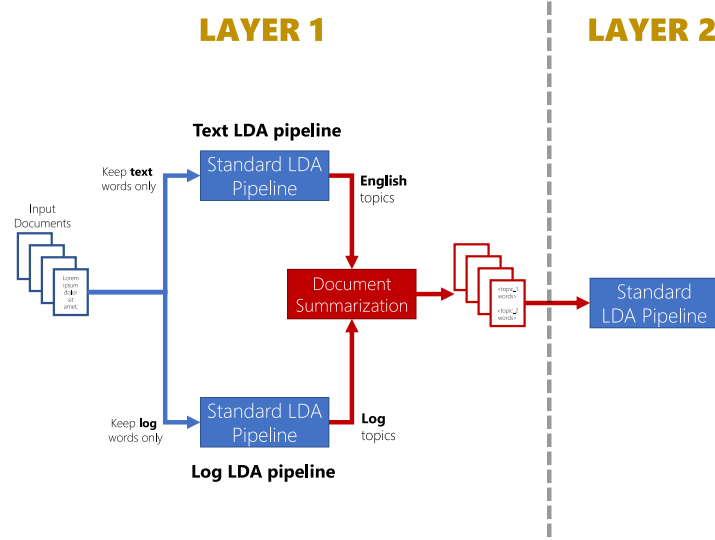
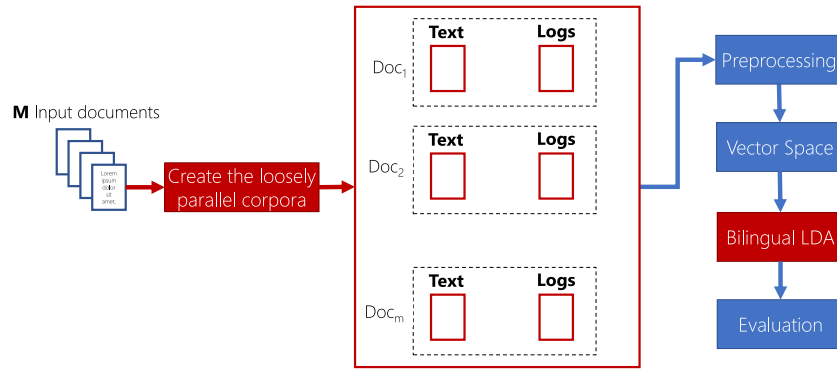
- + Simple and straightforward to interpret
- + Topics tend to be coherent (i.e., represent themes).
- If text is more prevalent than logs, then topics are more likely to contain text words only.

3.2.2. Two-layer LDA (M_{2L})

An overview of our two-layer LDA (M_{2L}) candidate model is shown in Fig. 5.

- **Rationale**: Our Enriched LDA candidate model (M_+) is biased towards producing topics that frequently include text words only when text is more prevalent than logs. Indeed, logs tend to be repetitive and contain a much lower number of unique words compared to text (Yao et al., 2020). Our two-layer LDA candidate model aims to tackle the imbalance between text and logs.

- **Construction**: We design a two-layer LDA pipeline. In the first layer, we prefix words with either text or log (analogously to M_+). Next, we generate two topic models: one using text words only (*LDA-LY1-text*) and another using log words only (*LDA-LY1-log*). Once these models are built using the standard LDA pipeline, we proceed to the *document summarization* step. In this step, we replace the contents of the original documents with their text topics and log topics. More specifically, for each topic model, we first identify the document-specific topic according to the topic probabilities, then we extract the top 10 associated terms by term probability. To identify the origin of each

Fig. 5. Two-layer LDA (M_{2L}) pipeline.Fig. 6. Off-the-shelf bilingual LDA (M_{bi}) pipeline.

term, we also append the topic indexes at the end of each term (e.g., LOG:stamp;TOPIC:399). The top terms extracted from both models are then merged by documents. Finally, the newly created documents now act as input to a third LDA topic model (layer 2).

Pros and cons:

- + Designed to counter-balance the prevalence of English words.
- More complicated.
- Higher computational cost (three LDA models across two layers).

3.2.3. Off-the-shelf bilingual LDA (M_{bi})

An overview of the off-the-shelf bilingual LDA (M_{bi}) candidate model is shown in Fig. 6.

- **Rationale:** Text and logs can be seen as belonging to two inherently different languages. Yet, we expect that the text and logs in a given forum thread will have some relationship. Multilingual LDA models are designed to find topic alignments across documents written in different languages. However, practitioners

should be cautious when choosing a proper multilingual LDA model since it usually requires comparable corpora. The weak alignment between different languages may introduce noise and thus impair the effectiveness of the topic modeling process.

- **Construction:** We use Mallet's multilingual LDA implementation, which is one of the most popular, readily available multilingual LDA implementations. Most importantly, Mallet's multilingual LDA implements the pLDA (a.k.a., pLTM) algorithm (Mimno et al., 2009), which is arguably suitable for *loosely parallel corpora* (a.k.a., *loosely comparable corpora*). Loosely parallel corpora can be formalized as a set of document pairs $\langle d_1, d_2 \rangle$ where d_1 and d_2 are each written in a different language, yet they share *similar* content (i.e., similar themes) (Zhu et al., 2013). In other words, Mallet's implementation does *not* require d_1 and d_2 to be direct translations of one another (i.e., a parallel corpora). The authors of pLDA observe that (i) “pLDA is appropriate for aligning topics in corpora that have a small subset of comparable documents” and (ii) “relatively small numbers of topically comparable document tuples are sufficient to align topics between languages in non-comparable corpora” (Mimno et al., 2009). Hence, we choose to build our model using Mallet's multilingual LDA implementation. The input to pLDA is the loosely parallel

corpora, which we build by extracting the text and log parts of each document (Fig. 6).

Pros and cons:

- + Prioritizes language alignment.
- + Easy to use and readily available (off-the-shelf solution)
- + Lower computation cost compared to our *two-layer LDA* candidate model.
- Text and logs are very loosely related. Off-the-shelf implementations such as Mallet's may thus struggle to align text and logs, which can result in fewer discovered topics and/or less coherent topics.

3.3. Implementation of candidate models

We modify gensim's Mallet LDA wrapper⁹ to implement the candidate LDA models. We added multiple features based on the original wrapper to accommodate our requirements, including (i) supporting multilingual LDA modeling and analysis, (ii) calculating different evaluation metrics, and (iii) extending configurations for Mallet modeling. The detailed implementation is available in our replication package.¹⁰

4. Results

In this section, we present our evaluation results to three proposed research questions (RQs). In each RQ, we explain the motivation and introduce corresponding approaches prior to discussing the experimental results.

4.1. RQ1: What is the achieved level of internal quality by each candidate model?

Motivation. We measure the quality of the generated topic models in an effort to understand the effectiveness of those models in capturing the underlying topics from a corpus. High-quality LDA models generate coherent and interpretable topics that can provide practitioners with valuable insights into the data for downstream tasks such as pattern discovery and text summarization – ultimately supporting practitioners in identifying the root underlying cause of failures.

Approach. We calculate the topic coherence achieved by each candidate LDA model after hyper tuning their parameters. Subsequently, we compare the coherence scores obtained. In the following, we explain our approach in more detail:

• *Step (1) Coherence calculation.* A set of statements or facts is said to be coherent if they support each other (Röder et al., 2015). Several metrics have been proposed to quantify the coherence of a fact set (e.g., a topic in LDA) (Bouma, 2009; Newman et al., 2010; Mimno et al., 2011; Aletras and Stevenson, 2013). We use the C_V measure to calculate topic coherence. In simple terms, C_V measures the relative distance between words within a topic and outputs an average computed over all topics. The C_V measure ranges from 0.0 to 1.0. We use *gensim* to calculate C_V .

It is challenging to interpret the absolute values of C_V , since there is no established nor commonly used guideline. In our practical experience using LDA and C_V (Chen et al., 2016, 2012; Barua et al., 2014; Li et al., 2018; Wan et al., 2021), we observe that extreme values for C_V are unlikely. Values too close to zero (0.0

to 0.2) often indicate a data collection problem (e.g., corrupt data). Similarly, values too close to one (0.8 to 1.0) often indicate that the top words in a topic always appear in the input documents as bigrams or trigrams, which is uncommon. The values in between tend to occur much more frequently, especially those between 0.4 and 0.6. Some LDA practitioners seem to share a similar interpretation of the C_V values.¹¹ Hence, we use the following thresholds to interpret the C_V measure.

$$\text{Coherence}(C_V) = \begin{cases} \text{Abnormally low,} & \text{if } 0.0 \leq |C_V| < 0.2 \\ \text{Very low,} & \text{if } 0.2 \leq |C_V| < 0.4 \\ \text{Low,} & \text{if } 0.4 \leq |C_V| < 0.5 \\ \text{Medium,} & \text{if } 0.5 \leq |C_V| < 0.55 \\ \text{High,} & \text{if } 0.55 \leq |C_V| < 0.6 \\ \text{Very high,} & \text{if } 0.6 \leq |C_V| < 0.8 \\ \text{Abnormally high,} & \text{if } 0.8 \leq |C_V| \leq 1.0 \end{cases}$$

• *Step (2) Finding the LDA parameter set that optimizes coherence.* LDA models have three key parameters: *alpha*, *beta*, and the *number of topics*. LDA uses Bayesian inference, which relies on the notion of *prior distributions*. Alpha is a parameter of the prior distribution of topics over documents. A low alpha value puts more weight on having each document composed of only a few dominant topics. Beta is a parameter of the prior distribution of words over topics. A low beta value puts more weight on having each topic composed of only a few dominant words. The number of topics indicates how many topics the LDA model should produce.

In this step, we aim to find the parameter set that optimizes the coherence of each candidate model (this set is likely different for each model). There are multiple hyperparameter tuning algorithms that can be used to achieve this goal. For instance, Agrawal et al. (2018) employed Differential Evolution (DE), Panichella et al. (2013) used Genetic Algorithms (GA), and popular LDA tools such as Mallet use Grid Search. Hence, there is no clear guideline indicating which algorithm should be used in each situation. Therefore, we choose to delegate this decision to OpenTuner (Ansel et al., 2014). OpenTuner uses ensembles of disparate search techniques simultaneously. Techniques that perform well receive larger testing budgets and techniques that perform poorly are disabled. As such, OpenTuner adapts to the problem at hand. Additionally, OpenTuner is extensible by design and supports flexible specification of hyperparameters and objective functions. We run OpenTuner on a Canada Compute cluster called Cedar.¹² In terms of resource allocation, we requested 32 CPU and 64Gb of RAM in total. We optimize one model at a time with a separate job for each. The execution time budget that we allocated for each OpenTuner job is LDA-dependent: 48 h for M_+ , 48 h for M_{bi} , and $48 \times 3 = 144$ h for M_{2L} (since this model entails training three LDAs). We use the following search space (parameter range) in all cases: $0.001 \leq \alpha \leq 50.0$; $0.001 \leq \beta \leq 50.0$; $10 \leq \text{numTopics} \leq 800$.

Findings. Observation (1) The M_{2L} model achieves the highest level of internal quality (highest coherence). However, it produces 4 to 65 times more topics than the other models. Fig. 7 depicts the coherence scores that we obtained from the hyperparameter tuning of our three candidate models. In all three cases, hyperparameter tuning led to a higher coherence score compared to that produced with the parameters fixed at default values. M_{2L} ranks first with a final coherence of 0.73 (*very high*), M_+ ranks second with a final coherence of 0.65 (*very high*), and M_{bi} ranks last with a coherence of 0.51 (*medium*).

⁹ <https://radimrehurek.com/gensim-3.8.3/models/wrappers/ldamallet.html>

¹⁰ <https://bit.ly/3Tdlp81> Once the paper is accepted, this replication package will be made publicly available on GitHub.

¹¹ <https://stackoverflow.com/questions/54762690/what-is-the-meaning-of-coherence-score-0-4-is-it-good-or-bad>

¹² <https://docs.computecanada.ca/wiki/Cedar>

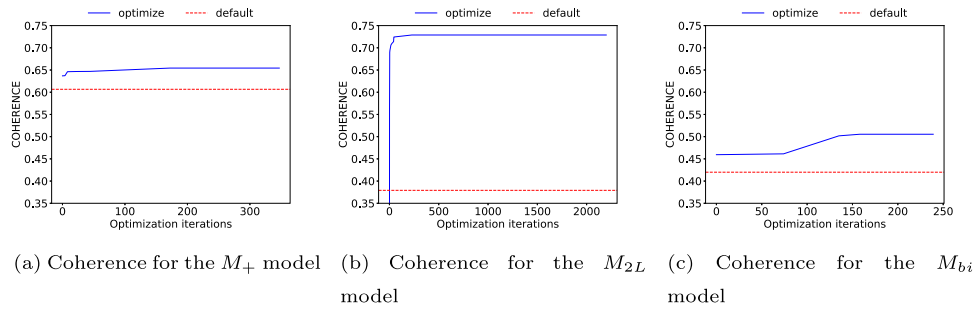


Fig. 7. Coherence scores of our candidate LDA models (higher is better). The red dotted line serves as a reference and indicates the coherence score achieved when gensim's default values for alpha and beta are used to build the LDA model. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 1

Optimal parameter set for each candidate model when maximizing coherence.

	Alpha	Beta	# Topics	Optimal coherence
M_+	24.02	6.25	157	0.65
M_{2L}	49.04	19.03	651	0.73
M_{bi}	6.07	11.59	10	0.51

The optimal parameter-set for each model is shown in Table 1. As the table indicates, M_{2L} produces 4-65x more topics than the other candidate models. Hence, if practitioners care about coherence only, they may select either M_{2L} or M_+ depending on how coarse-grained they want the topics to be.

4.2. RQ2: What is the achieved level of inferential power by each candidate model?

Motivation. LDA models are useful when they can correctly guess the topic of a new, unseen document. For instance, assume that an LDA found produced a topic T that contains both English words and log words. Now, assume that a forum member posts a problem P in a Bleeping Computer forum. If the LDA topic of P happens to be T and the LDA model in question has good inferential power, then log-specific information (i.e., the log words in T) can be requested to the poster of P . The provision of such log-specific information is likely to help to address P .

Approach. We computed the *perplexity* of each candidate model. Perplexity is the geometric mean of the inverse marginal probability of each word in a held-out set of documents. Perplexity is considered a suitable measure to determine the predictive power of a topic model. Low perplexity indicates that the model is good at predicting unseen data. To determine the LDA parameter-set that minimizes perplexity, we employed the same approach of RQ1 (i.e., use of OpenTuner in a Canada Compute cluster).

Findings. Observation (2) The M_{2L} model achieves the highest inferential power (lowest perplexity). Fig. 8 depicts the perplexity scores that we obtained from the hyperparameter tuning of our three candidate models. The hyperparameter tuning led to major perplexity drops for models M_+ and M_{bi} . M_{2L} 's perplexity remained fairly still, but it already started at a remarkably low level. In all cases, hyperparameter tuning led to a lower perplexity compared to that produced with the parameters fixed at their default values. M_{2L} ranks first with a final perplexity of 13.15, M_{bi} ranks second with a final perplexity of 100.88 (7.7 times higher than M_{2L}), and M_+ ranks last with a final perplexity score of 132.62 (10.1 times higher than M_{2L} , 1.3 times higher than M_{bi}).

The optimal parameter-set for each model is shown in Table 2. As the table indicates, M_{2L} and M_{bi} produce remarkably fewer

Table 2

Optimal parameter set for each candidate model when minimizing perplexity.

	Alpha	Beta	# Topics	Optimal perplexity
M_+	9.25	0.06	184	132.62
M_{2L}	0.46	13.15	10	13.15
M_{bi}	5.12	0.00	502	100.88

Table 3

Optimal parameter set for each candidate model when maximizing cross-lingual coverage ratio.

	Alpha	Beta	# Topics	Optimal Crossl. ratio
M_+	16.87	47.67	490	0.99
M_{2L}	44.77	33.70	694	1.00
M_{bi}	37.95	48.78	800	0.95

topics than M_+ . Different from the coherence optimization (RQ1), M_{2L} now produces the least number of topics. Hence, although coherence and perplexity are both popular LDA evaluation measures, we conclude that optimizing one in lieu of the other can lead to very different models.

4.3. RQ3: What is the achieved level of language alignment by each candidate model?

Motivation. In this work, we are interested in discovering associations between text and logs. Therefore, we assess the ability of our candidate models in discovering topics that contain both text and log words.

Approach. We refer to topics containing both log and text words as *cross-lingual topics*. Our approach consists of determining the coverage ratio of cross-lingual topics produced by each candidate model. We proceed analogously to RQ1 and RQ2.

Findings. Observation (3) The M_{2L} model achieves the highest cross-lingual coverage ratio. Fig. 9 depicts the cross-lingual coverage ratios that we obtained from the hyperparameter tuning of our three candidate models. Hyperparameter tuning led to very high cross-lingual coverage ratios overall (at least 0.95 in all three models), clearly outperforming the default parameter configurations. More specifically, M_{2L} ranked first with a final cross-lingual coverage ratio of 1.0 which is very close to M_+ with a final cross-lingual coverage ratio of 0.99, and M_{bi} ranked last with a final cross-lingual coverage ratio of 0.95.

The optimal parameter-set for each model is shown in Table 3. The numbers of topics produced by the models are all generally high (at least 490). This result contrasts to those obtained in RQ1 and RQ2, where M_{bi} and M_{2L} produced 10 topics only.

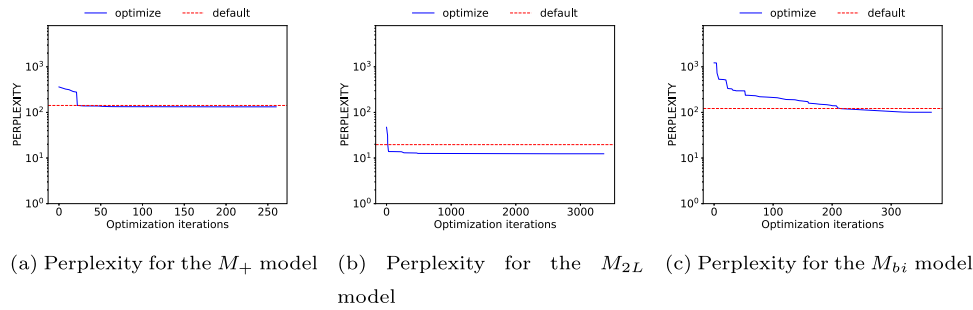


Fig. 8. Perplexity scores of our candidate LDA models (lower is better). The red dotted line serves as a reference and indicates the coherence score achieved when gensim's default values for alpha and beta are used to build the LDA model. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

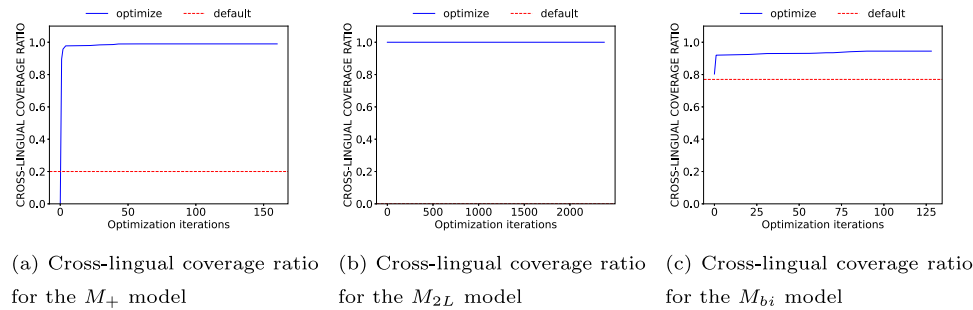


Fig. 9. Cross-lingual coverage ratio scores for our candidate LDA models.

4.4. Discussion

In prior sections, we present the experimental results of our proposed candidate models. In this section, we further study the extracted topics from the candidate models to discover the associations between log and text words. Such a relationship shall provide extra information to assist in failure debugging and improving the technical support service quality.

We follow a similar approach as Panichella et al. (2013) to understand the extracted topics using representative words. First, we choose the candidate models that achieve different optimal evaluation metric scores, respectively. Then for each model, we pick the most representative topics of each thread based on their probability in the corresponding topic distribution over documents. After mapping threads to topics, we then choose the top 20 words for each topic with the highest probability in its topic distribution. Finally, we manually inspect the acquired topics from each model. This analysis serves as an example to show the associations between extracted topics and the discussion threads. A selection of our manual study is showcased in Table 4.

As shown in Table 4, topic A contains a list of text and log words that produce latent connections where words can be used to identify and locate related issues. For instance, log word *hang* is related to system interruption or responding errors. *Bonjour* indicates logged errors related to the *bonjour* service of *windows*, which discovers and services devices from a local network. As for text words, the co-occurrence of *dll*, *system32*, and *md5* could indicate the procedure of integrity checking of critical windows system files. *GMER* is a rootkit detector that mines potentially malicious processes that are secretly running in the background. Such a word indicates a potential solution or at least a must-try step when issues related to those log words exist. For example, thread T1¹³ reports an issue that windows failed to update.

Throughout the diagnosis, the fixing includes searching for a missing *dll* file, solving application hang, running *GMER*, etc.

From the top text words of Topic B, we infer that this topic may relate to *registry* and *dll* files under the *system32* folder. *Adwcleaner* refers to leveraging Malwarebytes AdwCleaner, a junkware and malware removal tool, for system cleaning. The corresponding log words such *targetinst*, *fault*, *stamp* are related to errors generated from different windows applications. For example, thread T2¹⁴ describes an internet connection error and a google search redirect error. Such issues are also reflected in system error logs. To resolve such issues, the expert suggests to copy a missing DLL file from a cache folder (C:/WINDOWS/system32/dllcache). Similarly, thread T3¹⁵ explains a malware infection in google chrome. The user performs diagnosis through multiple tools, with *Adwcleaner* used as the first option. We believe such knowledge assists users in quickly locating issues and finding solutions from others' expertise.

In Topic C, the top text words *network* and *connect* indicate that the topic is likely related to internet connection issues, and *MSE* (Microsoft Security Essentials) suggests that such issues may relate to virus or malware. For example, thread T5¹⁶ states that the user is unable to connect to or download updates from anti-malware sites. The system error logs indicate that Windows AntiMalware fails fail to update its signatures. Another example in thread T4¹⁷ mentions connection errors to local devices and internet with the similar signature issues in the system error log.

Our observation indicates that our approach is capable of revealing hidden associations between log and text through designated topic modeling models. Such information can be leveraged to better support and facilitate issue diagnosis procedures.

¹⁴ <https://www.bleepingcomputer.com/forums/t/422583/deleted-ipsecsysi-now-have-no-internet-connection/>

¹⁵ <https://www.bleepingcomputer.com/forums/t/525169/unknown-chrome-infection/>

¹⁶ <https://www.bleepingcomputer.com/forums/t/421442/unable-to-update-antivirus-definitions-and-windows-updates/>

¹⁷ <https://www.bleepingcomputer.com/forums/t/454549/vista-is-falling-apart/>

¹³ <https://www.bleepingcomputer.com/forums/t/452705/windows-wont-update/>

Table 4
Examples of selected threads to associate topics with representative words.

	Measure	Topic	Top text terms	Top log terms	Threads
M_+	Coherence	A	java,md5,document, gmer,dll,registri, repair,system32	update,fault,windows, hang,server,address, search,dcom,network, dll,bonjour,device	[T1] Windows Will not Update ^a
M_{2L}	Perplexity	B	registri,java,error, system32,repair,dll, adwclean,document, googl,md5	fault,windows,path, code,program,update, targetinst,schedule, stamp,bonjour	[T2] deleted ipsec.sys... I now have no internet connection ^b [T3] Unknown Chrome infection ^c
M_{bi}	Aggregated	C	essenti,connect,access, error,point,network, mse,java,machin,think	signature,type,update, engine,stage,source, current,previous, authority,user	[T4] Vista is falling apart ^d [T5] Unable to update antivirus definitions and windows updates ^e

^a<https://www.bleepingcomputer.com/forums/t/452705/windows-wont-update/>.

^b<https://www.bleepingcomputer.com/forums/t/422583/deleted-ipsecsysi-now-have-no-internet-connection/>.

^c<https://www.bleepingcomputer.com/forums/t/525169/unknown-chrome-infection/>.

^d<https://www.bleepingcomputer.com/forums/t/454549/vista-is-falling-apart/>.

^e<https://www.bleepingcomputer.com/forums/t/421442/unable-to-update-antivirus-definitions-and-windows-updates/>.

5. Follow-up study

The results that we obtained indicate that M_{2L} consistently performs best. Yet, we are left with three open issues: (i) the selection of the quality attribute of interest (i.e., internal quality, inferential power, and language alignment) has a huge impact on the final number of topics that M_{2L} produces, (ii) it is unclear which model ranks second, and possibly more importantly (iii) it is unclear how many topics would have been produced by each model if we chose the parameter set that yields the best balance between our three quality attributes of interest.

To determine the best balance between the three quality attributes, we (i) design a composite metric that takes into account coherence, perplexity, and inter-corpus and (ii) use that metric as the objective function during hyperparameter tuning. However, designing such a composite metric is not trivial. While both *coherence* and *inter-corpus ratio* range from zero to one, *perplexity* does not have boundaries. In other words, aggregating the base metrics in such a way that each base metric contributes the same to the final value is not trivial. To circumvent the problem, we artificially define the boundaries for perplexity based on empirical data. More specifically, we use the perplexity values from the previous perplexity-based hyperparameter optimization as a reference range to create a normalized perplexity value, as shown below.

$$\text{Normalized Perplexity} = 1 - \frac{p - p_{\text{refmin}}}{p_{\text{refmax}} - p_{\text{refmin}}} \quad (1)$$

p is the current perplexity value, while p_{refmax} and p_{refmin} are the maximum and minimum perplexity values of the reference range. The perplexity values outside this range will be rounded to the nearest range boundary. Next, we composite the aggregated scores using the harmonic mean of coherence, cross-lingual coverage ratio, and the normalized perplexity.

Findings. Observation (4) The M_+ model outperforms the M_{bi} model, but by a small margin. Fig. 10 depicts the aggregated scores that we obtained from the hyperparameter tuning of our three candidate models. As expected, the M_{2L} model ranks first with an aggregated score of 0.87. M_+ ranks seconds with an aggregated score of 0.77, which is 11.6% higher than that of M_{bi} . The figure also indicates that using the default parameter-set would lead generally lead to very poor models.

The optimal parameter-set for each model is shown in Table 5. The number of topics produced by the models differs substantially: 20 (0.3% of the documents), 64 (0.8% of the documents), and 780 (9.8% of the documents). Hence, although M_{2L} has

the highest aggregated score, it is also providing a substantially coarser-grained view of the themes discussed in the documents. In addition, the high coherence and the low perplexity suggest that the *document summarization* step in the M_{2L} construction pipeline (Fig. 5) is possibly too aggressive.

We now compare M_+ and M_{bi} . Analysis of the optimal parameter set for each model reveals that M_+ has a higher overall performance than M_{bi} in coherence, perplexity, and cross-lingual coverage ratio. On the other hand, previous findings indicate that these two models outperform another when a single evaluation metric (i.e., coherence, perplexity) is preferred. Hence, in a practical scenario, a choice between these two models would be driven by a preference for coherence over perplexity or vice-versa (provided that the data has similar characteristics to the one that we investigate in this study).

6. Related work

LDA models are widely studied and adopted in the field of software engineering for topic modeling in mining software repositories (Chen et al., 2016; Li et al., 2018; Agrawal et al., 2018). In this section, we discuss the research related to our study from four areas of study: interplay between natural language and code in software engineering, multilingual LDA, tuning LDA for software engineering tasks, and applications of LDA in software engineering.

Interplay between natural language and code in software engineering. Natural language information in source code such as proper comments and identifier names makes it easier for humans to understand and evolve a given piece of code. Yet, from an algorithmic/computation perspective, such information is irrelevant. That is, the semantics of the code dictate the computation. Understanding this interplay between code and natural language has led to several research studies (Casalnuovo et al., 2020; Hindle et al., 2012; Allamanis et al., 2015, 2018). For instance, Casalnuovo et al. (2020) define a theory called *Dual Channel Constraints* (DCC) to highlight the relevance and interplay between the *algorithmic channel* and the *natural language channel* during software development. The theory's name stems from the realization that the two channels interact and constrain each other. Indeed, several problems arising from the interplay between the two channels have been studied in the software engineering literature.

The dual-channel nature of source code is usually leveraged for automated code generation and suggestion. Hindle et al. (2012)

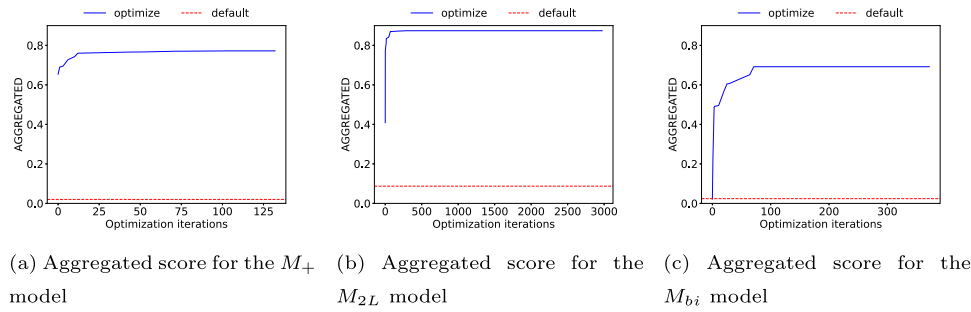


Fig. 10. Aggregated scores for our candidate LDA models.

Table 5

Optimal parameter set for each candidate model when maximizing aggregated score.

	Alpha	Beta	# Topics	Coherence	Perplexity	Crossl. ratio	Optimal aggreg. score
M_+	20.60	2.44	780	0.63	231.87	0.87	0.77
M_{2L}	24.62	32.40	64	0.70	71.26	1.00	0.87
M_{bi}	13.13	28.07	20	0.52	400.59	0.80	0.69

find that source code is analogous to natural language that is likely to be repetitive and predictable. Based on this observation, the authors propose a statistical language model that outperforms a state-of-the-art IDE in code completion. Iyer et al. (2018) propose an approach that automatically generates Java class member functions from dual information channels, which are the method documentation in natural language and the programmatic context from the rest of the code in a given class. Chen et al. (2021) introduced *Codex*, a fine-tuned GPT language model that generates standalone Python functions based on their docstrings that are documented in natural language. *Codex* now empowers GitHub Copilot¹⁸ to assist developers' daily programming tasks with automated coding suggestions. The dual-channel nature of source code is also corroborated and leveraged by Chakraborty et al. (2022) to improve the performance of NLP models that are designated for source code. Their approach learns to write coherent code that is similar to developers' code, and it achieves state-of-the-art performance in multiple code generative tasks.

Another line of research focuses on *natural language* information to enhance the quality and understanding of programs. Tan et al. (2007) mine implicit rules from code comments written in natural language to detect comments that are inconsistent with the corresponding source code that can lead to misunderstanding. Movshovitz-Attias and Cohen (2013) predict comments from Java source files using a statistical language model. They find using comment completion approaches can save up to 47% comment typing. To assist developers in writing informative and self-explanatory comments, Louis et al. (2018, 2020) build a machine learning framework that takes code snippets and natural language comments as input to detect and remove redundant comments from source code and to predict locations in source code that would benefit from adding comments. Hu et al. (2018, 2020) capture the source code semantics and structure to automatically generate code comments for Java methods. The following research by Zhou et al. (2022) introduces a novel approach to improve the effectiveness and robustness of existing deep learning models that generate code comments. Similarly, few studies focus on the challenge of automated logging text generation. He et al. (2018) investigate the characteristics of logging descriptive text and propose an information-retrieval approach to automatically generate logging description texts. Ding et al. (2022) propose a novel approach that generates texts in logging

code by translating corresponding code snippets into meaningful descriptive texts with a neural machine translation model.

The dual-channel nature of source code is also leveraged to address other DevOps tasks, including predicting function types (Malik et al., 2019; Petrescu et al., 2021; Dash et al., 2018), recovering original code from obfuscating source code (Vasilescu et al., 2017; Lachaux et al., 2021), predicting defection with comments and bug reports (Tan et al., 2007; Huo et al., 2018; Hu et al., 2020), identifying and managing self-admitted technical debt (SATD) (Yu et al., 2022), etc. In this vein, our work can be seen as a DCC study in which we aim to understand the relationship between log excerpts and natural language excerpts by means of a topic model.

Multilingual LDA. Multilingual topic models have been an active research area since the nineties (Yang et al., 2019; Hao et al., 2018; Balikas, 2017; Mimno et al., 2009; Zhu et al., 2013; Fung, 1998; Boyd-Graber and Blei, 2009; Jagarlamudi and Daumé, 2010; Hu et al., 2014; Yuan et al., 2018; Xie et al., 2020). Earlier studies in the area of multilingual topic models tried to align topics across languages by assuming that the documents written in different languages are parallel (e.g., direct translations of one another) or *highly comparable*. However, such an assumption seldom holds in practice. It is far more common for corpora to only have some degree of comparability (a.k.a., comparable corpora).

Recent approaches for aligning topics across comparable corpora span a variety of techniques: the polylingual topic model – *pLDA* (Mimno et al., 2009) (a seminal work that extended LDA to the multilingual context and that is implemented in Mallet), polylingual tree LDA – *ptLDA* (Hu et al., 2014) (extends the polylingual topic model with tree-based LDA – *tLDA* (Hu and Boyd-Graber, 2012)), *JointLDA* (uses a special bilingual dictionary to mine multilingual topics) (Jagarlamudi and Daumé, 2010), Multilingual Topic Anchoring – *MTAnchor* (Yuan et al., 2018) (uses the *anchor words* algorithm (Arora et al., 2012) in lieu of the more common alternatives, such as LDA (Blei et al., 2003) and LSA (Landauer et al., 1998)), and Multilingual Cultural-common Topic Analysis – *MCTA* (Shi et al., 2016) (uses a technique based on *auxiliary distributions* (Johnson and Riezler, 2000) that incorporates word distributions from the other language in an attempt to capture common semantics at the topic level).

More recently, Yang et al. (2019) introduced a new multilingual topic model specifically designed to find topics across *loosely comparable* corpora. In loosely comparable corpora, although some topics are shared, emphasis may diverge and some

¹⁸ <https://github.com/features/copilot>

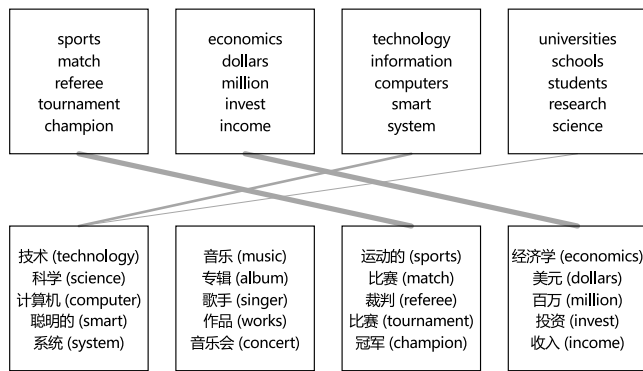


Fig. 11. Illustration of topic links (adapted from the study of Yang et al., 2019). Thicker lines denote links with higher weights. Topic pairs with many word translation pairs have high link weights. Topic pairs with partial overlap receive lower weights. A topic is unlinked if there is no corresponding topic in the other language.

topics may lack analogs in the other language. Their proposed model works by minimizing the Euclidean distances of translation pairs' transformed topic distributions. Transformations are learned from the original translation pairs' topic distributions and rely on a translation dictionary. In simpler terms, their model learns *weighted* topic links and only assigns a high link weight to topic pairs whose top words have many direct translation pairs. A topic is left unlinked if there is no matching topic in the other language (i.e., it does not enforce linking). Such a characteristic makes their model more robust for loosely comparable corpora with high topic misalignment. An example is shown in Fig. 11. The authors empirically show that their model substantially outperforms ptLDA, MTAnchor, and MCTA in both intra- and cross-lingual classification tasks for loosely comparable corpora.

In our study, natural language (which is used to describe a technical problem, investigate a technique problem, and propose candidate solutions) is *very loosely* comparable to logs (which can be seen as machine language). Differently from Yang et al. (2019), in our case, there is no pre-existing dictionary that can be leveraged to find direct correspondences between the words in natural language and those found in system logs. Hence, although Yang's model outperforms several multilingual models, we note that it is unsuitable in the context of our study. In addition, to the best of our knowledge, the implementation of Yang's model is not readily available.

Tuning LDA for Software Engineering tasks. Effectively applying LDA in the software engineering context is not trivial. Panichella et al. (2013) investigated how to effectively use LDA in three software engineering tasks, namely (i) traceability link recovery, (ii) feature location, and (iii) software artifact labeling. The authors conclude that LDA has a high sensitivity to different parameter settings and thus a careful calibration of those is essential. The authors propose to calibrate (hyperparameter tune) the LDA parameters using a genetic algorithm (GA). The decision variable in their study is the quality of an LDA-driven document clustering. This quality is measured using the *mean Silhouette coefficient* (Kogan, 2006), which is a popular clustering quality metric. The authors observe that LDA parameter calibration using GA led to accuracy improvements in experiments involving the three aforementioned software engineering tasks.

Agrawal et al. (2018) highlight that LDA is subject to *order effects* (Gennari et al., 1989) due to its stochastic nature. Having order effects implies that the produced topics might be different every time LDA is run. The authors refer to this problem as *topic instability*. Topic instability confuses users and lowers the

efficacy of text mining classifiers that rely on LDA to generate their input training data. To mitigate the problem, the authors propose an approach called LDADE that uses differential evolution (DE) (Storn and Price, 1997) to hyperparameter tune LDA with the goal of making topics look more similar across runs. The authors conduct experiments with LDADE and conclude that (i) it dramatically reduces topic instability and (ii) it leads to accuracy improvements in text mining classification tasks. However, it is not clear whether optimizing for topic stability yields better topics or even an adequate number of topics.

Treude and Wagner (2019) also recognize the need to calibrate LDA. They state that "there are only rough and sometimes conflicting guidelines available on how these parameters should be set". The authors conduct experiments with two large textual datasets: one from README files in GitHub and another from discussion threads on StackOverflow. The authors use the irace tool (López-Ibáñez et al., 2016) to calibrate LDA hyperparameters and they choose *perplexity* to be the decision variable. The authors create the corpora by segmenting the original datasets per programming language. The authors conclude that each corpus required different parameter settings in order to achieve good perplexity values (especially with regards to beta). They conclude that rules of thumb for topic modeling calibration are generally not applicable to their corpora.

In our paper, we acknowledge the high sensitivity of LDA to parameter choices. Differently from prior work, we not only optimize multiple variables (coherence, perplexity, and cross-lingual coverage ratio) but also an aggregated one. Our rationale is that all these variables are simultaneously important. In terms of the hyperparameter tuning tool/algorithm, we observe that there is no consensus in the literature at this point: researchers seem to choose the tool/algorithm based on the problem that they are trying to solve and the decision variable. To circumvent this problem, we use OpenTuner. OpenTuner is an ensemble technique that chooses the best search technique dynamically based on how they perform. More generally, evaluation of LDA models is actually research in itself (Wallach et al., 2009; Chang et al., 2009; Stevens et al., 2012; Layman et al., 2016).

Applications of LDA in Software Engineering. LDA has been applied to support a plethora of software engineering tasks. In 2016, Chen et al. (2016) published a survey that discussed enumerated these tasks: traceability link recovery, concept location, bug triaging, search engine usage analysis, requirements traceability, auditor support for exploring the implementation of requirements, analyzing bug report quality, estimating the number of topics in source code, clone detection, finding related code on the web, web service discovery, source code summarization and labeling, and refactoring. Since then, LDA has continued to be used to support software engineering tasks. More recent uses of LDA include: understanding the effects of testing on code quality (Chen et al., 2017), understanding system logs (Li et al., 2018), and discovering the themes discussed by programmers in the context of blockchain technology (Wan et al., 2021).

7. Threats to validity

Construct Validity. In our data collection, we extract discussion threads from the Bleeping Computer forum (Section 2). It is possible that a discussion thread will be preemptively closed with a reference to a prior thread that discusses the same issue (a *reference thread*). Topic models such as LDA operate directly on word distributions over documents and do not account for document links (cross-references). Future work should investigate whether reference threads should be weighted higher during topic discovery.

In RQ1 (Section 4.1), we use the C_V measure to quantify topic coherence. However, there exist other popular measures for quantifying topic coherence, including: UCI coherence (Newman et al., 2010) (which relies on the notion of *pointwise mutual information* – PMI (Church and Hanks, 1990)), UMass coherence (Mimno et al., 2011), and normalized PMI coherence (NPMI) (Bouma, 2009; Aletras and Stevenson, 2013). It is likely that using any of these measures would yield different coherence results compared to those that we obtained. Nevertheless, the C_V measure was found to be the top-performing measure in a systematic study of coherence measures using different datasets (Röder et al., 2015).

In this study, we adopt community-recognized thresholds for the interpretation of the C_V measure. Such thresholds derive from the knowledge and experience of practitioners in reasoning about C_V values. Hence, we believe that these thresholds are capable of providing insights into the evaluation of topic coherence. To the best of our knowledge, there is no gold standard in the literature for interpreting absolute C_V values. We thus call for future research to propose systematic ways to evaluate the C_V values of a topic model.

In our follow-up study (Section 5), we define an aggregated score that takes into account our three LDA qualities of interest (internal quality, inferential power, and language alignment). We assign the same weight to the three qualities of interest. Assigning different weights to different qualities would likely yield different results. Future research may also explore more sophisticated (including more computationally expensive) approaches for finding good trade-offs between the qualities of interest in an LDA model, such as multi-objective optimization (Deb, 2001). Moreover, we normalize the perplexity measure in order to make it bounded and compatible with the other measures. Using other normalization techniques (e.g., $\tanh()$ with adjustments) would possibly influence the final results as well. Ultimately, the take-home message of Section 5 is that multiple quality metrics should be simultaneously considered when training an LDA model.

Internal Validity. In Section 3, we describe our three candidate LDA models. Our two-layer LDA (M_{2L}) summarizes documents in order to counterbalance the prevalence of English words compared to logs. In practice, such a summarization can be seen as an algorithm that *transforms* the original set of documents into another set of documents that have a different word distribution. That is, M_{2L} alters the data. Hence, even though we observe that LDA performs the best compared to the other candidate models, our results should be carefully scrutinized. An alternative design for M_{2L} would consist of (i) chaining the layer 1 LDA models as priors (input) to the third LDA and (ii) performing the Bayesian Inference through all three models. This is theoretically feasible, since the Dirichlet Process distribution satisfies prior conjugacy (i.e., the posterior distribution is also a Dirichlet process). Nevertheless, such an approach does not tackle the text/log imbalance problem (c.f., Section 3.2.2). We thus invite future work to adapt M_{2L} by replacing the summarization step with some other suitable preprocessing technique (e.g., term reweighing schemes Wilson and Chew, 2010).

Agrawal et al. (2018) discusses the problem of topic instability, which arises from the stochastic nature of LDA. In preliminary runs of our experiments, we configured OpenTuner to (i) build the same candidate model 10 times during every iteration and (ii) use the *median value* of the decision variable (e.g., C_V coherence) to inform its tuning process. We did not observe a significant change in results compared to running the tuning process in its standard form (single model construction per iteration). In other words, we did not observe significant topic instability in our candidate models. Nevertheless, future work should take topic instability

into account when assessing our candidate models with other datasets.

The results from our paper heavily depend on the OpenTuner tool. OpenTuner is a robust and actively maintained autotuning tool with an active community around it. Using other autotuning tools (e.g., iRace (Treude and Wagner, 2019; López-Ibáñez et al., 2016)) or specific tuning algorithms (e.g., genetic algorithms (Panichella et al., 2013; Holland, 1975)) may produce a different result.

External Validity. Although our results derive from a specific dataset from Bleeping Computer, our proposed approach is generalizable to evaluations on other datasets. Our study should be seen as a first attempt at linking English text to system logs using LDA-based models. We strongly encourage future work to assess our models on other datasets and compare the results to those that we obtained in this paper. In particular, future work should also assess our candidate models from a more qualitative perspective (e.g., determine whether discovered topics align with the intuition of domain experts).

8. Conclusion

In this paper, we investigate the problem of linking natural language (English) to machine language (system logs) in the (technical) support threads. Uncovering such links would enable the construction of knowledge base that connects technical problem discussions to log excerpts that are often associated (requested/provided) alongside those discussions.

As an initial effort towards addressing the aforementioned problem, we evaluate three LDA candidate models. Two of them are designed by us, which we refer to *enriched LDA* (M_+) and *two-layer LDA* (M_{2L}). The third one is an *off-the-shelf bilingual LDA* (M_{bi}) implementation made available in the Mallet package.

Our evaluation of the LDA candidate models is multi-faceted. More specifically, as opposed to recent LDA research within software engineering (Treude and Wagner, 2019; Agrawal et al., 2018; Panichella et al., 2013), we evaluate models alongside three axes instead of a single one, namely: internal quality, inferential power, and language alignment. These axes are operationalized by coherence, perplexity, and cross-lingual coverage ratio respectively. We evaluate our models using discussion threads from a Bleeping Computer forum.

Our results indicate that M_{2L} performs best overall. However, it provides a substantially coarser-grained view of the themes discussed in the documents, since it produces very few topics (20, 0.3% of the documents). M_+ performs slightly better than M_{bi} . Comparing the two reveals that M_+ produces topics with higher coherence, while M_{bi} produces topics with lower perplexity.

Given the scale of the study and the heterogeneity of the discussion threads, it is difficult for us to qualitatively determine whether the topics produced by the models are adequate to be used in practice. We invite future studies to assess the qualities of topics produced by LDA models in a practical scenario. Researchers may also refer our evaluation approach to domain-specific datasets with the availability of domain experts, to determine the feasibility of employing our proposed models in practice.

CRedit authorship contribution statement

Kundi Yao: Conceptualization, Methodology, Software, Validation, Writing – original draft, Writing – review & editing. **Gustavo A. Oliva:** Conceptualization, Methodology, Writing – original draft, Writing – review & editing. **Ahmed E. Hassan:** Supervision, Project administration. **Muhammad Asaduzzaman:** Data curation, Investigation. **Andrew J. Malton:** Supervision. **Andrew Walenstein:** Supervision.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

I have attached an external link to our replication package. However, we are not permitted to share the data used in this study due to copyright policy.

References

- Agrawal, A., Fu, W., Menzies, T., 2018. What is wrong with topic modeling? And how to fix it using search-based software engineering. *Inf. Softw. Technol.* 98, 74–88. <http://dx.doi.org/10.1016/j.infsof.2018.02.005>, URL <https://www.sciencedirect.com/science/article/pii/S0950584917300861>.
- Aletras, N., Stevenson, M., 2013. Evaluating topic coherence using distributional semantics. In: *Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013) – Long Papers*. Association for Computational Linguistics, Potsdam, Germany, pp. 13–22. URL <https://aclanthology.org/W13-0102>.
- Allamanis, M., Barr, E.T., Devanbu, P.T., Sutton, C., 2018. A survey of machine learning for big code and naturalness. *ACM Comput. Surv.* 51 (4), 81:1–81:37. <http://dx.doi.org/10.1145/3212695>.
- Allamanis, M., Tarlow, D., Gordon, A.D., Wei, Y., 2015. Bimodal modelling of source code and natural language. In: *Bach, F.R., Blei, D.M. (Eds.), Proceedings of the 32nd International Conference on Machine Learning. ICML 2015, Lille, France, 6–11 July 2015*. In: *JMLR Workshop and Conference Proceedings*, vol. 37, JMLR.org, pp. 2123–2132. URL <http://proceedings.mlr.press/v37/allamanis15.html>.
- Ansel, J., Kamil, S., Veeramachaneni, K., Ragan-Kelley, J., Bosboom, J., O'Reilly, U.-M., Amarasinghe, S., 2014. OpenTuner: An extensible framework for program autotuning. In: *Proceedings of the 23rd International Conference on Parallel Architectures and Compilation. PACT '14*, Association for Computing Machinery, New York, NY, USA, pp. 303–316. <http://dx.doi.org/10.1145/2628071.2628092>.
- Arora, S., Ge, R., Moitra, A., 2012. Learning topic models – going beyond SVD. In: *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*. IEEE Computer Society, Los Alamitos, CA, USA, pp. 1–10. <http://dx.doi.org/10.1109/FOCS.2012.49>, URL <https://doi.ieeecomputersociety.org/10.1109/FOCS.2012.49>.
- Balikas, G., 2017. Mining and Learning from Multilingual Text Collections Using Topic Models and Word Embeddings (theses). Grenoble 1 UGA - Université Grenoble Alpes, URL <https://hal.archives-ouvertes.fr/tel-01706347>.
- Barua, A., Thomas, S.W., Hassan, A.E., 2014. What are developers talking about? An analysis of topics and trends in stack overflow. *Empir. Softw. Eng.* 19 (3), 619–654. <http://dx.doi.org/10.1007/s10664-012-9231-y>.
- Bettenburg, N., Adams, B., Hassan, A.E., Smidt, M., 2011. A lightweight approach to uncover technical artifacts in unstructured data. In: *2011 IEEE 19th International Conference on Program Comprehension*, pp. 185–188. <http://dx.doi.org/10.1109/ICPC.2011.36>.
- Blei, D.M., Ng, A.Y., Jordan, M.I., 2003. Latent Dirichlet allocation. *J. Mach. Learn. Res.* 3, 993–1022.
- Bouma, G., 2009. Normalized (pointwise) mutual information in collocation extraction. In: *From Form to Meaning: Processing Texts Automatically, Proceedings of the Biennial GSCS Conference 2009, Vol. Normalized*. Tübingen, pp. 31–40.
- Boyd-Graber, J., Blei, D.M., 2009. Multilingual topic models for unaligned text. In: *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence. UAI '09*, AUAI Press, Arlington, Virginia, USA, pp. 75–82.
- Casalnuovo, C., Barr, E.T., Dash, S.K., Devanbu, P., Morgan, E., 2020. A theory of dual channel constraints. In: *Rothermel, G., Bae, D. (Eds.), ICSE-NIER 2020: 42nd International Conference on Software Engineering, New Ideas and Emerging Results*. Seoul, South Korea, 27 June – 19 July, 2020, ACM, pp. 25–28. <http://dx.doi.org/10.1145/3377816.3381720>.
- Chakraborty, S., Ahmed, T., Ding, Y., Devanbu, P.T., Ray, B., 2022. NatGen: generative pre-training by “naturalizing” source code. In: *Roychoudhury, A., Cadar, C., Kim, M. (Eds.), Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering. ESEC/FSE 2022, Singapore, Singapore, November 14–18, 2022*, ACM, pp. 18–30. <http://dx.doi.org/10.1145/3540250.3549162>.
- Chang, J., Boyd-Graber, J., Gerrish, S., Wang, C., Blei, D.M., 2009. Reading tea leaves: How humans interpret topic models. In: *Proceedings of the 22nd International Conference on Neural Information Processing Systems. NIPS '09*, Curran Associates Inc., Red Hook, NY, USA, pp. 288–296.
- Chen, T.-H., Thomas, S.W., Hassan, A.E., 2016. A survey on the use of topic models when mining software repositories. *Empir. Softw. Eng.* 21 (5), 1843–1919. <http://dx.doi.org/10.1007/s10664-015-9402-8>.
- Chen, T.-H.P., Thomas, S.W., Hemmati, H., Nagappan, M., Hassan, A.E., 2017. An empirical study on the effect of testing on code quality using topic models: A case study on software development systems. *IEEE Trans. Reliab.* 66 (3), 806–824. <http://dx.doi.org/10.1109/TR.2017.2699938>.
- Chen, T.-H., Thomas, S.W., Nagappan, M., Hassan, A.E., 2012. Explaining software defects using topic models. In: *Proceedings of the 9th IEEE Working Conference on Mining Software Repositories. MSR '12*, IEEE Press, pp. 189–198.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H.P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F.P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W.H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A.N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W., 2021. Evaluating large language models trained on code. *CoRR* abs/2107.03374. [arXiv:2107.03374](https://arxiv.org/abs/2107.03374).
- Church, K.W., Hanks, P., 1990. Word association norms, mutual information, and lexicography. *Comput. Linguist.* 16 (1), 22–29. URL <https://aclanthology.org/J90-1003>.
- Dash, S.K., Allamanis, M., Barr, E.T., 2018. RefiNym: using names to refine types. In: *Leavens, G.T., Garcia, A., Pasareanu, C.S. (Eds.), Proceedings of the 2018 ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. ESEC/SIGSOFT FSE 2018, Lake Buena Vista, FL, USA, November 04–09, 2018*, ACM, pp. 107–117. <http://dx.doi.org/10.1145/3236024.3236042>.
- Deb, K., 2001. *Multi-Objective Optimization Using Evolutionary Algorithms*. In: *Wiley-Interscience Series in Systems and Optimization*, Wiley.
- Ding, Z., Li, H., Shang, W., 2022. LoGenText: Automatically generating logging texts using neural machine translation. In: *IEEE International Conference on Software Analysis, Evolution and Reengineering. SANER 2022, Honolulu, HI, USA, March 15–18, 2022*, IEEE, pp. 349–360. <http://dx.doi.org/10.1109/SANER53432.2022.00051>.
- Fung, P., 1998. A statistical view on bilingual lexicon extraction: From parallel corpora to non-parallel corpora. In: *Proceedings of the Third Conference of the Association for Machine Translation in the Americas on Machine Translation and the Information Soup. AMTA '98*, Springer-Verlag, Berlin, Heidelberg, pp. 1–17.
- Gennari, J.H., Langley, P., Fisher, D., 1989. Models of incremental concept formation. *Artificial Intelligence* 40 (1), 11–61. [http://dx.doi.org/10.1016/0004-3702\(89\)90046-5](http://dx.doi.org/10.1016/0004-3702(89)90046-5), URL <https://www.sciencedirect.com/science/article/pii/0004370289900465>.
- Hao, S., Boyd-Graber, J., Paul, M.J., 2018. Lessons from the Bible on modern topics: Low-resource multilingual topic model evaluation. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, pp. 1090–1100. <http://dx.doi.org/10.18653/v1/N18-1099>, URL <https://aclanthology.org/N18-1099>.
- He, P., Chen, Z., He, S., Lyu, M.R., 2018. Characterizing the natural language descriptions in software logging statements. In: *Huchard, M., Kästner, C., Fraser, G. (Eds.), Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering. ASE 2018, Montpellier, France, September 3–7, 2018*, ACM, pp. 178–189. <http://dx.doi.org/10.1145/3238147.3238193>.
- Hindle, A., Barr, E.T., Su, Z., Gabel, M., Devanbu, P.T., 2012. On the naturalness of software. In: *Glinz, M., Murphy, G.C., Pezzè, M. (Eds.), 34th International Conference on Software Engineering. ICSE 2012, June 2–9, 2012, Zurich, Switzerland*, IEEE Computer Society, pp. 837–847. <http://dx.doi.org/10.1109/ICSE.2012.6227135>.
- Holland, J.H., 1975. *Adaptation in Natural and Artificial Systems*, second ed. University of Michigan Press, Ann Arbor, MI, 1992.
- Hu, Y., Boyd-Graber, J., 2012. Efficient tree-based topic modeling. In: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers - Volume 2. ACL '12*, Association for Computational Linguistics, USA, pp. 275–279.
- Hu, X., Li, G., Xia, X., Lo, D., Jin, Z., 2018. Deep code comment generation. In: *Khomh, F., Roy, C.K., Siegmund, J. (Eds.), Proceedings of the 26th Conference on Program Comprehension. ICPC 2018, Gothenburg, Sweden, May 27–28, 2018*, ACM, pp. 200–210. <http://dx.doi.org/10.1145/3196321.3196334>.
- Hu, X., Li, G., Xia, X., Lo, D., Jin, Z., 2020. Deep code comment generation with hybrid lexical and syntactical information. *Empir. Softw. Eng.* 25 (3), 2179–2217. <http://dx.doi.org/10.1007/s10664-019-09730-9>.

- Hu, Y., Zhai, K., Eidelman, V., Boyd-Graber, J., 2014. Polylingual tree-based topic models for translation domain adaptation. In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Baltimore, Maryland, pp. 1166–1176. <http://dx.doi.org/10.3115/v1/P14-1110>, URL <https://aclanthology.org/P14-1110>.
- Huo, X., Yang, Y., Li, M., Zhan, D., 2018. Learning semantic features for software defect prediction by code comments embedding. In: IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17–20, 2018, IEEE Computer Society, pp. 1049–1054. <http://dx.doi.org/10.1109/ICDM.2018.00133>.
- Iyer, S., Konstantas, I., Cheung, A., Zettlemoyer, L., 2018. Mapping language to code in programmatic context. In: Riloff, E., Chiang, D., Hockenmaier, J., Tsujii, J. (Eds.), Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. Brussels, Belgium, October 31 – November 4, 2018, Association for Computational Linguistics, pp. 1643–1652. <http://dx.doi.org/10.18653/v1/d18-1192>.
- Jagarlamudi, J., Daumé, H., 2010. Extracting multilingual topics from unaligned comparable corpora. In: Gurrin, C., He, Y., Kazai, G., Kruschwitz, U., Little, S., Roelleke, T., Rüger, S., van Rijsbergen, K. (Eds.), Advances in Information Retrieval. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 444–456.
- Johnson, M., Riezler, S., 2000. Exploiting auxiliary distributions in stochastic unification-based grammars. In: Proceedings of the 1st North American Chapter of the Association for Computational Linguistics Conference. NAACL 2000, Association for Computational Linguistics, USA, pp. 154–161.
- Kogan, J., 2006. Introduction to Clustering Large and High-Dimensional Data. Cambridge University Press.
- Lachaux, M., Rozière, B., Szafraniec, M., Lample, G., 2021. DOBF: A deobfuscation pre-training objective for programming languages. In: Ranzato, M., Beygelzimer, A., Dauphin, Y.N., Liang, P., Vaughan, J.W. (Eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021. NeurIPS 2021, December 6–14, 2021, Virtual, pp. 14967–14979, URL <https://proceedings.neurips.cc/paper/2021/hash/7d6548bdc0082aacc950ed35e91fccb-Abstract.html>.
- Landauer, T., Foltz, P., Laham, D., 1998. An introduction to latent semantic analysis. *Discourse Process.* 25, 259–284.
- Layman, L., Nikora, A.P., Meek, J., Menzies, T., 2016. Topic modeling of NASA space system problem reports: Research in practice. In: Proceedings of the 13th International Conference on Mining Software Repositories. MSR '16, Association for Computing Machinery, New York, NY, USA, pp. 303–314. <http://dx.doi.org/10.1145/2901739.2901760>.
- Li, H., Chen, T.-H.P., Shang, W., Hassan, A.E., 2018. Studying software logging using topic models. *Empir. Softw. Eng.* 23 (5), 2655–2694. <http://dx.doi.org/10.1007/s10664-018-9595-8>.
- López-Ibáñez, M., Dubois-Lacoste, J., Pérez Cáceres, L., Birattari, M., Stütze, T., 2016. The irace package: Iterated racing for automatic algorithm configuration. *Oper. Res. Perspect.* 3, 43–58. <http://dx.doi.org/10.1016/j.orp.2016.09.002>, URL <https://www.sciencedirect.com/science/article/pii/S2214716015300270>.
- Louis, A., Dash, S.K., Barr, E.T., Ernst, M.D., Sutton, C., 2020. Where should I comment my code?: a dataset and model for predicting locations that need comments. In: Rothermel, G., Bae, D. (Eds.), ICSE-NIER 2020: 42nd International Conference on Software Engineering, New Ideas and Emerging Results. Seoul, South Korea, 27 June – 19 July, 2020, ACM, pp. 21–24. <http://dx.doi.org/10.1145/3377816.3381736>.
- Louis, A., Dash, S.K., Barr, E.T., Sutton, C., 2018. Deep learning to detect redundant method comments. *CoRR abs/1806.04616* [arXiv:1806.04616](https://arxiv.org/abs/1806.04616).
- Malik, R.S., Patra, J., Pradel, M., 2019. NL2type: inferring JavaScript function types from natural language information. In: Atlee, J.M., Bultan, T., Whittle, J. (Eds.), Proceedings of the 41st International Conference on Software Engineering. ICSE 2019, Montreal, QC, Canada, May 25–31, 2019, IEEE / ACM, pp. 304–315. <http://dx.doi.org/10.1109/ICSE.2019.00045>.
- Mimno, D., Wallach, H.M., Naradowsky, J., Smith, D.A., McCallum, A., 2009. Polylingual topic models. In: Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, Singapore, pp. 880–889, URL <https://aclanthology.org/D09-1092>.
- Mimno, D., Wallach, H.M., Talley, E., Leenders, M., McCallum, A., 2011. Optimizing semantic coherence in topic models. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing. EMNLP '11, Association for Computational Linguistics, USA, pp. 262–272.
- Movshovitz-Attias, D., Cohen, W.W., 2013. Natural language models for predicting programming comments. In: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, ACL 2013, 4–9 August 2013, Sofia, Bulgaria, Volume 2: Short Papers. The Association for Computer Linguistics, pp. 35–40, URL <https://aclanthology.org/P13-2007/>.
- Newman, D., Lau, J.H., Grieser, K., Baldwin, T., 2010. Automatic evaluation of topic coherence. In: Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics. HLT '10, Association for Computational Linguistics, USA, pp. 100–108.
- Panichella, A., Dit, B., Oliveto, R., Di Penta, M., Poshvanyk, D., De Lucia, A., 2013. How to effectively use topic models for software engineering tasks? An approach based on genetic algorithms. In: Proceedings of the 2013 International Conference on Software Engineering. ICSE '13, IEEE Press, pp. 522–531.
- Petrescu, C.C., Smith, S., Giavrimis, R., Dash, S.K., 2021. Do names echo semantics? A large-scale study of identifiers used in c++'s named casts. *CoRR abs/2111.01577* [arXiv:2111.01577](https://arxiv.org/abs/2111.01577).
- Röder, M., Both, A., Hinneburg, A., 2015. Exploring the space of topic coherence measures. In: Proceedings of the Eighth ACM International Conference on Web Search and Data Mining. WSDM '15, Association for Computing Machinery, New York, NY, USA, pp. 399–408. <http://dx.doi.org/10.1145/2684822.2685324>.
- Shi, B., Lam, W., Bing, L., Xu, Y., 2016. Detecting common discussion topics across culture from news reader comments. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Berlin, Germany, pp. 676–685. <http://dx.doi.org/10.18653/v1/P16-1064>, URL <https://aclanthology.org/P16-1064>.
- Stevens, K., Kegelmeyer, P., Andrzejewski, D., Buttler, D., 2012. Exploring topic coherence over many models and many topics. In: Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning. EMNLP-CoNLL '12, Association for Computational Linguistics, USA, pp. 952–961.
- Storn, R., Price, K., 1997. Differential evolution – A simple and efficient heuristic for global optimization over continuous spaces. *J. Global Optim.* 11 (4), 341–359. <http://dx.doi.org/10.1023/A:1008202821328>, URL <https://doi.org/10.1023/A:1008202821328>.
- Tan, L., Yuan, D., Krishna, G., Zhou, Y., 2007. /!comment: Bugs or bad comments?*. *SIGOPS Oper. Syst. Rev.* 41 (6), 145–158. <http://dx.doi.org/10.1145/1323293.1294276>.
- Treude, C., Wagner, M., 2019. Predicting good configurations for GitHub and stack overflow topic models. In: Proceedings of the 16th International Conference on Mining Software Repositories. MSR '19, IEEE Press, pp. 84–95. <http://dx.doi.org/10.1109/MSR.2019.00022>.
- Vasilescu, B., Casalnuovo, C., Devanbu, P.T., 2017. Recovering clear, natural identifiers from obfuscated JS names. In: Bodden, E., Schäfer, W., van Deursen, A., Zisman, A. (Eds.), Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering. ESEC/FSE 2017, Paderborn, Germany, September 4–8, 2017, ACM, pp. 683–693. <http://dx.doi.org/10.1145/3106237.3106289>.
- Wallach, H.M., Murray, I., Salakhutdinov, R., Mimno, D., 2009. Evaluation methods for topic models. In: Proceedings of the 26th Annual International Conference on Machine Learning. ICML '09, Association for Computing Machinery, New York, NY, USA, pp. 1105–1112. <http://dx.doi.org/10.1145/1553374.1553515>.
- Wan, Z., Xia, X., Hassan, A.E., 2021. What do programmers discuss about blockchain? A case study on the use of balanced LDA and the reference architecture of a domain to capture online discussions about blockchain platforms across stack exchange communities. *IEEE Trans. Softw. Eng.* 47 (07), 1331–1349. <http://dx.doi.org/10.1109/TSE.2019.2921343>.
- Wang, C., Blei, D.M., 2011. Collaborative topic modeling for recommending scientific articles. In: Apté, C., Ghosh, J., Smyth, P. (Eds.), Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. San Diego, CA, USA, August 21–24, 2011, ACM, pp. 448–456. <http://dx.doi.org/10.1145/2020408.2020480>.
- Wilson, A.T., Chew, P.A., 2010. Term weighting schemes for latent Dirichlet allocation. In: Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings. June 2–4, 2010, Los Angeles, California, USA, The Association for Computational Linguistics, pp. 465–473, URL <https://aclanthology.org/N10-1070/>.
- Xie, Q., Zhang, X., Ding, Y., Song, M., 2020. Monolingual and multilingual topic analysis using LDA and BERT embeddings. *J. Informetr.* 14 (3), 101055. <http://dx.doi.org/10.1016/j.joi.2020.101055>, URL <https://www.sciencedirect.com/science/article/pii/S1571157719305127>.
- Yang, W., Boyd-Graber, J., Resnik, P., 2019. A multilingual topic model for learning weighted topic links across corpora with low comparability. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing. EMNLP-IJCNLP, Association for Computational Linguistics, Hong Kong, China, pp. 1243–1248. <http://dx.doi.org/10.18653/v1/D19-1120>.
- Yao, K., Li, H., Shang, W., Hassan, A.E., 2020. A study of the performance of general compressors on log files. *Empir. Softw. Eng.* 25 (5), 3043–3085. <http://dx.doi.org/10.1007/s10664-020-09822-x>.

- Yu, Z., Fahid, F.M., Tu, H., Menzies, T., 2022. Identifying self-admitted technical debts with jitterbug: A two-step approach. *IEEE Trans. Softw. Eng.* 48 (5), 1676–1691. <http://dx.doi.org/10.1109/TSE.2020.3031401>.
- Yuan, M., Van Durme, B., Ying, J.L., 2018. Multilingual anchoring: Interactive topic modeling and alignment across languages. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems*, Vol. 31. Curran Associates, Inc., URL <https://proceedings.neurips.cc/paper/2018/file/28b9f8aa9f07db88404721af4a5b6c11-Paper.pdf>.
- Zhou, Y., Zhang, X., Shen, J., Han, T., Chen, T., Gall, H.C., 2022. Adversarial robustness of deep code comment generation. *ACM Trans. Softw. Eng. Methodol.* 31 (4), 60:1–60:30. <http://dx.doi.org/10.1145/3501256>.
- Zhu, Z., Li, M., Chen, L., Yang, Z., 2013. Building comparable corpora based on bilingual LDA model. In: *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Sofia, Bulgaria, pp. 278–282, URL <https://aclanthology.org/P13-2050>.



Kundi Yao is a Ph.D. student in the Software Analysis and Intelligence Lab (SAIL) at Queen's University, Canada. He obtained his M.sc. from Concordia University (Canada), and B.Eng. from Wuhan University of Technology (China). His research interests lie within software log compression, software log analytics, software performance engineering, and mining software repositories. Contact him at kundi@cs.queensu.ca.



Gustavo A. Oliva is an Adjunct Professor at the School of Computing of Queen's University in Canada. Gustavo leads the blockchain research team at the Software Analysis and Intelligence Lab (SAIL). In addition to blockchain, he also conducts empirical studies in the fields of software ecosystems (e.g., npm), software evolution, and explainable AI.

Gustavo received his M.Sc. and Ph.D. degrees from the University of São Paulo (USP) in Brazil under the supervision of professor Dr. Marco Aurélio Gerosa. More information at: <https://www.gaoliva.com>.



Ahmed E. Hassan is an IEEE Fellow, an ACM SIGSOFT Influential Educator, an NSERC Steacie Fellow, the Canada Research Chair (CRC) in Software Analytics, and the NSERC/BlackBerry Software Engineering Chair at the School of Computing at Queen's University, Canada. His research interests include mining software repositories, empirical software engineering, load testing, and log mining. He received a Ph.D. in Computer Science from the University of Waterloo. He spearheaded the creation of the Mining Software Repositories (MSR) conference and its research community. He also serves/d on the editorial boards of *IEEE Transactions on Software Engineering*, *Springer Journal of Empirical Software Engineering*, and *PeerJ Computer Science*. Contact ahmed@cs.queensu.ca. More information at: <http://sail.cs.queensu.ca/>.



Muhammad Asaduzzaman is an assistant professor in the Department of Computer Science at Lakehead University, Canada. His research interests include software maintenance and evolution, mining software repositories, program comprehension, reverse engineering, and recommendation systems. His work has been published at premier software engineering venues such as the TSE, EMSE, ASE, ICSME, SANER, and MSR. Before moving to Canada, he studied at Khulna University in Bangladesh, where he received his B.Sc. He obtained his M.Sc. and Ph.D. from the Department of Computer Science at University of Saskatchewan, Canada. More about Muhammad can be read on his website: <https://muhammad-asaduzzaman.com>.



Andrew Walenstein is the Senior Director of BlackBerry Labs. Previously he was an Assistant Professor of Computer Science at the University of Louisiana at Lafayette. His research interests are in cybersecurity, software engineering tools, and humancomputer cooperation.