# Makar: A Framework for Multi-source Studies based on Unstructured Data

Mathias Birrer\*, Pooja Rani\*, Sebastiano Panichella<sup>†</sup>, Oscar Nierstrasz\*

 \*Software Composition Group, University of Bern Bern, Switzerland
 Scg.unibe.ch/staff
 <sup>†</sup> Zurich University of Applied Science (ZHAW) panc@zhaw.ch

Abstract—To perform various development and maintenance tasks, developers frequently seek information on various sources such as mailing lists, Stack Overflow (SO), and Quora. Researchers analyze these sources to understand developer information needs in these tasks. However, extracting and preprocessing unstructured data from various sources, building and maintaining a reusable dataset is often a time-consuming and iterative process. Additionally, the lack of tools for automating this data analysis process complicates the task to reproduce previous results or datasets.

To address these concerns we propose *Makar*, which provides various data extraction and preprocessing methods to support researchers in conducting reproducible multi-source studies. To evaluate Makar, we conduct a case study that analyzes code comment related discussions from SO, Quora, and mailing lists. Our results show that Makar is helpful for preparing reproducible datasets from multiple sources with little effort, and for identifying the relevant data to answer specific research questions in a shorter time compared to state-of-the-art tools, which is of critical importance for studies based on unstructured data. Tool webpage: https://github.com/maethub/makar

*Index Terms*—Mining developer sources, Code comments, Stack Overflow, Mailing lists

### I. INTRODUCTION

As a software system continues to evolve, it becomes bigger and more complex, and developers need various kinds of information to perform activities such as adding features, or performing corrective maintenance [1]. Developers typically seek information on internal (available within IDE) or external sources such as Q&A forums,<sup>1</sup> Github<sup>2</sup> to satisfy their information needs as shown in Figure 1 [2].

To support developers in various activities and understand their information needs, researchers have analyzed these external sources such as Github, CVS, mailing lists, and CQA sites [3] (see Figure 1). However, extracting and preprocessing unstructured data from these sources, and maintaining the data due



Fig. 1. Developers seek various sources during software development

to lack of automated techniques pose various challenges in conducting reproducible studies [4], [5], [3]. To gain a deeper understanding of these challenges, we surveyed the literature that focuses on studying developers information needs from different external sources (see section II).

Prior works have raised and identified the crucial factors affecting the reproducibility of the mining studies such as data retrieval methodology, data processing steps, or dataset availability [6], [5], [4]. Chen *et al.* pointed out that 50% of articles do not report whether word stemming, a common text preprocessing step, is used or not [4]. Amann *et al.* pointed out that only 29% of the mining studies made their dataset available [5]. As a consequence, more tools and techniques are required to enable the preprocessing and analysis of multi-source studies to facilitate their replicability.

To address these concerns, we propose *Makar*, a tool that leverages popular data retrieval, processing, and handling techniques to support researchers in conducting reproducible studies. We establish its requirements based on the surveyed studies. To evaluate Makar, we conduct a case study that analyzes code comment related discussions from SO, Quora, and mailing lists. Thus the contribution of this paper is threefold:

- We present the challenges researchers face in mining and analyzing the unstructured data from the external sources.
- We present Makar, a tool to support researchers in conducting multi-source and reproducible empirical studies.
- We report the state-of-art tools comparison to Makar.

We gratefully acknowledge the financial support of the Swiss National Science Foundation for the project "Agile Software Assistance" (SNSF project No. 200020-181973, Feb. 1, 2019 - April 30, 2022).

<sup>&</sup>lt;sup>1</sup>www.stackoverflow.com

<sup>&</sup>lt;sup>2</sup>https://github.com/

## **II. BACKGROUND STUDY**

To identify the challenges researchers face with various sources, we surveyed the relevant papers considered in a recent systematic literature review (SLR) [7]. The SLR includes the studies in which researchers collected developer information needs by interviewing people (people-centric) or from online platforms (technology-centric) for the program comprehension tasks [7]. We included only technology-centric studies (29 studies) due to our interest in the external sources. Following the same inclusion and exclusion criteria from the SLR (e.g., studies not older than 15 years), we further included 23 additional papers that focus on studying developer information needs from other sources such as mobile app stores (e.g., user reviews) and Quora, resulting in a total of 52 papers. In particular, we included the study if it focuses in part or whole on software developer information needs related to software development and includes empirical evidence. We excluded the study if is a review, survey, or tool study, older than 15 years, not peer-reviewed, or not in English. As we aim to focus on the diversity of sources rather than on a deep overview of a particular source, we excluded the studies analyzing the same project from the same source. The list of selected papers and detailed observations are reported in the "Background Study" file in the Replication Package (RP).<sup>3</sup>

Challenges in various sources. Table I reports our main findings, where the column Source represents the source, and the columns Data Extraction, Data Relevancy, and Data Preprocessing reports the major challenges associated with handling data from each source, and the column Makar reports the sources Makar supports in Data Extraction, Data *Preprocessing* currently  $(\checkmark)$  or those planned for the future (FW). As the challenges of selecting relevant data given in Data Relevancy column depends on a research context, Makar supports exploring and filtering data to select relevant data. The results show that while a few sources offer convenient ways for extracting data (e.g., SO), there are other sources (e.g., Quora) that are more prohibitive and complex to acquire the required data. Similarly, extracting and processing data from mailing lists require manual efforts. Therefore, extracting the data manually is still widely adopted in practice. However, the use of manual extraction methods can lead to inconsistent collection and processing of data across sources, which impacts the reproducibility of the studies.

Requirements. Based on the gathered challenges in the survey, we identified relevant functional and non-functional requirements for Makar. The tool intends to cover the common use cases found in the survey while being extensible to support additional or more specific scenarios encountered in the case study. It can also be used by developers to manage their information in development as depicted in Figure 1.

We identified five main functional requirements: data import, data management, data processing, data querying, and data export. Data import focuses on the ways to import the data into the tool, Data management on building and maintaining the data, Data processing focuses on the need to preprocess the data (HTML removal, stop word removal), Data Querying on searching the data, and Data export focuses on exporting the data from the tool in order to support further analysis. We also identified non-functional requirements for Makar. It should be easily extensible in areas where the projects have different technical requirements, such as import adapters, preprocessing steps, or export adapters. The tool should be able to handle large amounts of data (scale of 100k records) and still have acceptable usage performance (e.g., for search queries).

## **III. MAKAR ARCHITECTURE**

Makar has been developed as a web application so that it can be hosted on accessible and possibly powerful servers. Thus, it allows multiple users to work with the same dataset. It is a Ruby on Rails (RoR)<sup>4</sup> web application with a Postgresql<sup>5</sup> database in the back end. To have minimal technical requirements to run the tool, to provide maximal compatibility and ease of installation on different platforms and operating system, Makar runs in a Docker container. <sup>6</sup> We provide instructions to run the tool on the tool repository<sup>7</sup> and its demonstration on Youtube.8 We show its architecture and features in Figure 2 and the next paragraphs.





• Data import: the user can import data from diverse sources such as CSV and JSON directly. The tool also supports direct import adapters for the following sources: Apache Mailing List Archive,<sup>9</sup> Github Pull Requests (via Github Archive),<sup>10</sup> Github Issues Via the Github API,<sup>11</sup> and Stack Overflow Search Excerpts.<sup>12</sup> The import

4https://rubyonrails.org/

5https://www.postgresql.org/

<sup>6</sup>https://www.docker.com/

<sup>7</sup>https://github.com/maethub/makar

<sup>8</sup>https://youtu.be/Yqj1b4Bv-58

10http://www.gharchive.org/

<sup>11</sup>https://developer.github.com/v3/search/#search-issues 12https://api.stackexchange.com/docs/excerpt-search

578

<sup>&</sup>lt;sup>3</sup>https://doi.org/10.5281/zenodo.4434822

<sup>9</sup>https://mail-archives.apache.org/mod\_mbox/

Source	Data Extraction	Data Relevancy	Data Preprocessing	Makar
SO	Public API & Data dumps	Selecting relevant and pertinent ques- tions [2], [8]	Removing noisy data such as HTML tags, code snippets [2], [3]	~
Quora	No official API available to access data and no publicly available dataset [9]	Finding relevant topics and ques- tions [9]	Preprocessing data for the study	~
Mailing lists	MBOX file, if available otherwise the unstructured text in the mails requires manual data extraction	Contains unstructured text ( <i>i.e.</i> , no tags or assigned topics)	Consists of heterogeneous types of in- formation (stack traces, simple text, footers)	~
Bug Reports	Data extraction is performed manu- ally [10]	Information overload and it requires human interpretation to select relevant data [10]	Contains stack traces, simple text, and code snippets	~
User Reviews	No public API to access and extract user review data [11], [12]	Require human interpretation to select relevant data [12], suffers from a sam- pling bias [13]	lack of source code artifacts makes the preprocessing straightforward	FW
Combining Sources	Extracting data from multiple sources [14]	Require human interpretation to map data across sources [14]	Consistent preprocessing of data [3]	~

TABLE I Challenges from each data source

adapters can be extended easily using ImportAdapter component for other sources shown in Table I.

- Data management: Makar provides schemas, collections, filters and records to manage datasets as shown in Figure 2. Schemas define the structure of a dataset and its records, and records are rows of the dataset (similar to schemas and records in databases). Collections are arbitrary selections of records, which can be used to manage various subsets of the data. A record can belong to multiple collections. Filters are the search queries that help one to filter data from existing collections or schemas, and can be saved to provide efficient querying and rebuilding of the dataset. For example, a study analyzing SO questions imports the SO dataset into Makar. The study design requires only questions having the word "javadoc" in the question title. To fulfill this requirement, the user can create a filter (e.g., "All Questions with Javadoc in Title" filter) by searching the question titles for "javadoc" as shown in Figure 3. The user can create a collection that uses this filter and use the collection as their dataset for further analysis. In the case, the user add more data from SO to update her dataset (collection), Makar facilitates syncing the collection using the Autofilter option (reapplying the same filter) as shown in Figure 4.
- Data processing. The user can preprocess the data in Makar through *transformation steps*. A *transformation step* is a single operation that is applied to all records in a collection. Currently, the tool supports operations such as *text cleaning*, *natural language processing*, *data restructuring*, and *arithmetic and counting*.
  - In *text cleaning*, the user can strip all HTML tags, or selected HTML tags, or replace records with custom values *e.g.*, remove HTML tags from questions in SO.
  - In natural language processing, the user can apply



Fig. 3. Search Interface of Makar

Export			
Attributes			
Id × Title × Body × Tags × has_code ×			
Export format			
CSV			
Export			
	Laport M ∺ Title ∺ Body ∺ bgs ∺ bas.code ∺ Coper format* Cov		

Fig. 4. Dataset Preparation Interface of Makar

word stemming,<sup>13</sup> remove all stop words,<sup>14</sup> or remove all punctuation.

In *data restructuring*, the user can merge records having same value, create new records, remove duplicates, split text on defined substring, add a static value. In addition, the user can create a new dataset with

<sup>13</sup> https://snowballstem.org/algorithms/porter/stemmer.html

<sup>14</sup> http://snowball.tartarus.org/algorithms/english/stop.txt

a randomized sample, which is widely performed in manual analysis studies.

- In *arithmetic and counting*, the user can also perform simple arithmetic steps *e.g.*, counting frequent occurrences of a particular value or a word.
- Data export. The user can select which attributes are to be selected for the export, and then export the data in the required format as shown in Figure 4. Currently, the tool supports *CSV*, *JSON*, and .txt (file) formats. Makar also supports adding more complicated export formats via ExportAdapter. To perform LDA (Latent Dirichlet Allocation) analysis using Mallet, we added the Mallet adapter (custom export adapter).<sup>15</sup>

## IV. MULTI-SOURCE ANALYSIS USING MAKAR

Code comments play a crucial role in program comprehension and maintenance [15]. However, their semi-structured nature and the availability of multiple commenting conventions confront developers with numerous ways to write them. Consequently, developers often post questions to learn about different conventions on various sources such as Q&A websites [2]. To identify such concerns, we conducted an empirical study on SO, Quora, and mailing lists using Makar.

Methodology. We manually identified ten relevant tags from SO by searching comment and convention keywords on its tag page.<sup>16</sup> The selected tags are: comments, commenting, code-comments, block-comments, autocommenting, commentconventions, convention, conventions, coding-style. Based on a heuristics-based approach proposed by Ying et al., we added five more relevant tags: documentation, todo, codedocumentation, naming, readability [8]. We used the relevant tags from SO as keywords to find relevant topics on Quora. As a result, we obtained five topics from Quora: Code Comments, Source Code, Coding Style, Programming Languages, Comment (computer programming)]. We mined mailing lists of five Apache projects that we selected based on the top programming languages, Line Of Code, and number of commits from the Apache statistics report.<sup>17</sup> Thus, we considered Lucene (Java), Ambari (JS), OpenOffice (C++), Cloudstack (Python), and Subversion (C). From these projects, we mined @dev, @users and @docs mailing lists. The resulting data from each source is shown in Table II.

To obtain the high-level overview for SO questions, we used the popular topic analysis method, LDA [4]. To obtain a more detailed view of each source, we extracted a statistically significant sample set of discussions from each source (reaching 95% confidence level and an error margin of 5%) to analyze manually. Makar supported us in preparing the dataset suitable for the LDA analysis and manual analysis.

TABLE II DATA EXTRACTED FROM VARIOUS SOURCES

Source	Fields extracted	Candidate posts	Manually analyzed
SO	id, title,x body, tags, creation date, view count	11931	373
Quora	url, title, body, topics, answers	689	689
Mailing lists	all	140 667	385

- **Data import**: We imported the SO data using the CSV import adapter, Quora data with the JSON adapter, and mailing lists with the *Apache Mailing List Archive* adapter. The CSV files of the dataset are provided in the RP.<sup>18</sup>
- Data processing: The data from SO contains HTML, code snippets, links and natural language text. To get meaningful results from LDA analysis, the data need to be cleaned, with text cleaning and language cleaning steps. All preprocessing steps such as removing code, HTML, punctuations, and stop words,<sup>19</sup> and stemming words<sup>20</sup> are performed by Makar using its built-in transformations as shown in Figure 5. In the figure, the Transformation as described in section III, shows various built-in transformations of Makar and Attributes shows the list of selected fields (e.g., Title, Body) from the sources. Each transformation is designed to produce a new attribute (a column) in the data records, allowing us to retrace the changes applied to the data. As it is generally uncertain in the beginning of a study which combination of preprocessing steps would lead to the best results, the flexible approach of Makar supported us in trying several scenarios efficiently.

	Code	HTML	Punctuation	Stop Word	Word Stemming	
Transformation	extract_code strip_html		string_replace	remove_stopwords	word_stemming	
Attributes	Attributes - Question   Body - Question   Body		<ul> <li>Question   Body</li> <li>Question   Title</li> </ul>	- Question   Body - Question   Title	<ul> <li>Question   Body</li> <li>Question   Title</li> </ul>	

Fig. 5. Preprocessing steps with the transformation in the tool

• **Data export**: The dataset from the case study has been exported as CSV and provided in the RP.<sup>21</sup>

## V. TOOL COMPARISON

We compare similar state-of-art tools reported in Table III based on the functionality defined in section III: extracting(*Data import*), preprocessing(*Data processing*), and managing data (*Data management*) from multiple sources in a

18 https://doi.org/10.5281/zenodo.4434822

- 19http://snowball.tartarus.org/algorithms/english/stop.txt
- <sup>20</sup>https://snowballstem.org/algorithms/porter/stemmer.html

<sup>21</sup>"Data" folder in https://doi.org/10.5281/zenodo.4434822

<sup>15</sup>http://mallet.cs.umass.edu/

<sup>&</sup>lt;sup>16</sup>https://stackoverflow.com/tagsverifiedon20Nov2020

<sup>17</sup> https://projects.apache.org/statistics.html

TABLE III

Tool	Costs	Extract	Process	Manage
Octoparse	Commercial	<b>v</b>	×	×
Knime	Free	Extension	~	~
Rapidminer	Commercial	×	Limited	Limited
ELKI	Free	×	~	×
Keel	Free	×	×	×
WEKA	Free	×	×	×
TrifactaWrangler	Commercial	×	~	~
Boa	Free	Limited	~	~
OpenRefine	Free	~	~	~
Makar	Free	~	~	~

reproducible way. In Table III, the column Extract focuses on mining data from various sources such as mailing lists, or Q&A forums, the column Process focuses on various preprocessing operations on the data, such as removing noisy HTML tags, and stop words, and the column Manage focuses on importing, exploring, and filtering the multi-source data into the tool. The links to access the tools are provided in the RP and the tool page due to space constraints.<sup>22</sup> Our direct comparison shows that the majority of previous tools (except OpenRefine, TrifactaWrangler, Octoparse) provide pipelines facilitating the process of building machine-learning based data analysis and visualizing their results. However, they lack the ability to manually explore, extract, and map the data from various sources as well as to investigate small samples, or perform ad-hoc searches on intermediate data. Researchers interested in using new sources of data or combining multiple sources using various heuristics to map the sources [14] with state-of-art-tools are limited in their decision making. Human interpretation, ad-hoc testing of simple hypotheses, rescaling a dataset, or the assessment of data quality is often required to determine the plausible approach or methodology for using the new data sources and combining the multiple sources. Compared to data cleaning tools (OpenRefine, TrifactaWrangler), Makar focuses on tailoring specific use cases for software engineering researchers, allowing them to perform a wide range of feasibility analyses or quality assessment steps on the data in a reproducible way.

### CONCLUSION

In this paper, we presented Makar, a tool supporting and enabling multi-source empirical studies. The performance of the tool has been assessed through an empirical study involving 11 931 questions from Stack Overflow, 140 667 mails from mailing lists, and 700 from Quora. Makar helped us to process the multi-source data in a uniform way and to investigate various combinations of features for both LDA analysis and manual analysis. Moreover, Makar provides an extensible framework to support custom requirements, so further textual analysis techniques can be integrated to perform advanced text operations.

### REFERENCES

- M. Lehman, D. Perry, J. Ramil, W. Turski, and P. Wernick, "Metrics and laws of software evolution-the nineties view," in *Proceedings IEEE International Software Metrics Symposium (METRICS'97)*. Los Alamitos CA: IEEE Computer Society Press, 1997, pp. 20–32.
- [2] H. Gujral, A. Sharma, S. Lal, A. Kaur, A. Kumar, and A. Sureka, "Empirical analysis of the logging questions on the Stack Overflow website," in 2018 Conference On Software Engineering & Data Sciences (CoSEDS)(in-press), 2018.
- [3] G. Bavota, "Mining unstructured data in software repositories: Current and future trends," in 2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER), vol. 5. IEEE, 2016, pp. 1–12.
- [4] T.-H. Chen, S. W. Thomas, and A. E. Hassan, "A survey on the use of topic models when mining software repositories," *Empirical Softw. Engg.*, vol. 21, no. 5, pp. 1843–1919, oct 2016. [Online]. Available: https://doi.org/10.1007/s10664-015-9402-8
- [5] S. Amann, S. Beyer, K. Kevic, and H. Gall, Software Mining Studies: Goals, Approaches, Artifacts, and Replicability. Springer International Publishing, 2015, pp. 121–158. [Online]. Available: https://doi.org/10.1007/978-3-319-28406-4\_5
- [6] J. M. González-Barahona and G. Robles, "On the reproducibility of empirical software engineering studies based on data retrieved from development repositories," *Empirical Software Engineering*, vol. 17, no. 1, pp. 75–89, 2012. [Online]. Available: http://dx.doi.org/10.1007/ s10664-011-9181-9
- [7] J. Richner, "Software developers' information needs," University of Bern, Bachelor's thesis, Feb. 2019. [Online]. Available: http: //scg.unibe.ch/archive/projects/Rich19a.pdf
- [8] X.-L. Yang, D. Lo, X. Xia, Z.-Y. Wan, and J.-L. Sun, "What security questions do developers ask? a large-scale study of Stack Overflow posts," *Journal of Computer Science and Technology*, vol. 31, no. 5, pp. 910–924, 2016. [Online]. Available: https: //doi.org/10.1007/s11390-016-1672-0
- [9] S. Patil and K. Lee, "Detecting experts on quora: by their activity, quality of answers, linguistic characteristics and temporal behaviors," *Social network analysis and mining*, vol. 6, no. 1, p. 5, 2016.
- [10] S. Breu, R. Premraj, J. Sillito, and T. Zimmermann, "Information needs in bug reports: improving cooperation between developers and users," in *Proceedings of the 2010 ACM conference on Computer supported cooperative work.* ACM, 2010, pp. 301–310.
- [11] B. Fu, J. Lin, L. Li, C. Faloutsos, J. Hong, and N. Sadeh, "Why people hate your app: Making sense of user feedback in a mobile app store," in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 1276–1284.
- [12] D. Pagano and W. Maalej, "User feedback in the appstore: An empirical study," in 2013 21st IEEE international requirements engineering conference (RE). IEEE, 2013, pp. 125–134.
- [13] W. Martin, M. Harman, Y. Jia, F. Sarro, and Y. Zhang, "The app sampling problem for app store mining," in 2015 IEEE/ACM 12th Working Conference on Mining Software Repositories. IEEE, 2015, pp. 123–133.
- [14] A. Zagalsky, C. G. Teshima, D. M. German, M.-A. Storey, and G. Poo-Caamaño, "How the R community creates and curates knowledge: A comparative study of Stack Overflow and mailing lists," in *Proceedings of the 13th International Conference on Mining Software Repositories*, ser. MSR '16. New York, NY, USA: ACM, 2016, pp. 441–451. [Online]. Available: http://doi.acm.org/10.1145/2901739.2901772
- [15] S. C. B. de Souza, N. Anquetil, and K. M. de Oliveira, "A study of the documentation essential to software maintenance," in *Proceedings of* the 23rd annual international conference on Design of communication: documenting & designing for pervasive information, ser. SIGDOC '05. New York, NY, USA: ACM, 2005, pp. 68–75.

<sup>&</sup>lt;sup>22</sup>https://github.com/maethub/makar/Similar-Tools.md