
Datasheet for Dataset: Papyrus

Frédéric Piedboeuf
RALI, Diro
Université de Montréal
frederic.piedboeuf@umontreal.ca

Philippe Langlais
RALI, Diro
Université de Montréal
felipe@iro.umontreal.ca

1 Motivation

For what purpose was the dataset created? The dataset was created in order to address the problem of French, as well as multilingual, keyphrase generation.

Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)? The papyrus repository belongs to Université de Montréal (UdeM) which gave the RALI research lab of UdeM the right to prepare and distribute the present dataset for research purposes. The dataset has been prepared as part of the first author doctoral research.

Who funded the creation of the dataset? The dataset was created independently as part of the first author doctoral research, in cooperation with the Papyrus team from Université de Montréal.

Any other comments? This dataset was created and curated in partnership with the Papyrus team of Université de Montréal. We thank them for letting us use and publish this data, which will without a doubt be a great use to the research community.

2 Composition

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Each instance represent an abstract, a title, and keyphrases associated with a document that was uploaded on Papyrus. Most of these documents are theses, but there are additional documents as well, such as institutional public reports, presentations, or documents uploaded by members of the faculty.

How many instances are there in total (of each type, if appropriate)? There are 26508 documents in the *non-curated* version of the dataset, among them 16249 valid ones. By *curated* we mean any preprocessing/filtering that was done after the collection step. Each document is identified by a unique index which corresponds to the index in Papyrus. Invalid documents include indexes that point to error pages, as well as documents that have no keyphrases, titles, or abstracts.

Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? The dataset itself contains all documents from Papyrus, but we also release four curated subsets of it that do not contain all entries. Papyrus-f contains only French keyphrases and abstracts, and Papyrus-e only English ones. Papyrus-m entries are one abstract as well as associated keyphrases that have the same language, and Papyrus-a contains all data, with one entry being all abstracts and keyphrases of a document.

34 **What data does each instance consist of?** Each instance consist of a title, one or more
35 abstracts in several languages, as well as several keyphrases as written by the author of the
36 document.

37 **Is there a label or target associated with each instance?** Yes, the targets/labels are
38 the keyphrases.

39 **Is any information missing from individual instances?** Some of the documents do
40 not have keyphrases in all languages, which will lead to invalid entries when separated for
41 the various tasks.

42 **Are relationships between individual instances made explicit (e.g., users' movie
43 ratings, social network links)?** In some subtasks, multiples translation of the same
44 abstracts will find itself as multiple inputs. Other than that, some documents have been
45 written by the same authors, for examples students who write more than one thesis. Also,
46 professors typically supervise several students and therefore there is a possible relation
47 between theses (and therefore abstracts) written by the students supervised by a same
48 professor.

49 **Are there recommended data splits (e.g., training, development/validation,
50 testing)?** Yes, we provide the data already split.

51 **Are there any errors, sources of noise, or redundancies in the dataset?** Keyphrases
52 are automatically assigned to a language using a simple heuristic: For keyphrases present in
53 the abstract, we assign them to all the abstracts that contain them. For absent keyphrases,
54 we use `fastText` and assign it to the most likely language. If `fastText` cannot identify the
55 languages¹, such as for the keyphrase "間 ma", "1000", or "Leptoquark", we assign it to all
56 abstracts for that entry. This causes some source of error, but we estimated an accuracy of
57 98.9% of accuracy on the keyphrases using 100 random entries from Papyrus-m. There is
58 also possible sources of redundancies in the data. The same abstract and keyphrases were
59 sometimes assigned to two or more documents. We remove the additional copies in the
60 curated data.

61 **Is the dataset self-contained, or does it link to or otherwise rely on external
62 resources (e.g., websites, tweets, other datasets)?** The dataset is self-contained.

63 **Does the dataset contain data that might be considered confidential (e.g., data
64 that is protected by legal privilege or by doctor– patient confidentiality, data
65 that includes the content of individuals' non-public communications)?** No. All
66 documents are published openly.

67 **Does the dataset contain data that, if viewed directly, might be of- fensive,
68 insulting, threatening, or might otherwise cause anxiety?** Potentially : The dataset
69 is composed of theses covering a wild range of subjects, including but not limited to rape,
70 police brutality, or the n word.

71 **Does the dataset identify any subpopulations (e.g., by age, gen- der)?** Some
72 instances of the dataset are theses discussing specific populations, such as women, queer
73 people, or some other subpopulations.

74 **Is it possible to identify individuals (i.e., one or more natural per- sons), either
75 directly or indirectly (i.e., in combination with other data) from the dataset?**
76 Yes, one simply has to research the title or the abstract on a search engine to find the
77 original thesis and the author(s) who wrote it.

78 **Does the dataset contain data that might be considered sensitive in any way (e.g.,
79 data that reveals race or ethnic origins, sexual ori- entations, religious beliefs,
80 political opinions or union member- ships, or locations; financial or health data;**

¹If the top-15 languages from `fastText` are not concordant with the languages of the abstracts.

81 biometric or genetic data; forms of government identification, such as social
82 security numbers; criminal history)? No.
83 Any other comments? No.

84 3 Collection Process

85 **How was the data associated with each instance acquired?** Each document came in
86 Papyrus with some meta-data that we used to extract the abstracts, titles, and keyphrases.
87 Separation into languages for the various tasks was however conducted automatically with
88 a process that we estimate accurate at 98.9%.

89 **What mechanisms or procedures were used to collect the data (e.g., hardware
90 apparatuses or sensors, manual human curation, software programs, software
91 APIs)?** A scrapping process was used to collect the data off the internet and then curated
92 with heuristics as described before.

93 **If the dataset is a sample from a larger set, what was the sampling strategy
94 (e.g., deterministic, probabilistic with specific sampling probabilities)?** N/A

95 **Who was involved in the data collection process (e.g., students, crowdworkers,
96 contractors) and how were they compensated (e.g., how much were
97 crowdworkers paid)?** Only the authors of the papers were involved in the data collection
98 process.

99 **Over what timeframe was the data collected?** The latest version of the dataset
100 was collected on April 7, 2022. The whole scrapping and curating of the dataset took
101 approximately two weeks.

102 **Were any ethical review processes conducted (e.g., by an institu- tional review
103 board)?** No, but we asked and obtained permission from the people in charge of Papyrus
104 for the collection and distribution of the dataset.

105 **Did you collect the data from the individuals in question directly, or obtain it
106 via third parties or other sources (e.g., websites)?** We collected it from the Papyrus
107 website.

108 **Were the individuals in question notified about the data collec- tion?** No, but the
109 people in charge of Papyrus were notified and gave consent.

110 **Did the individuals in question consent to the collection and use of their data?**
111 N/A.

112 **If consent was obtained, were the consenting individuals pro- vided with a
113 mechanism to revoke their consent in the future or for certain uses?** N/A.

114 **Has an analysis of the potential impact of the dataset and its use on data
115 subjects (e.g., a data protection impact analysis) been con- ducted?** N/A.

116 Any other comments? No.

117 4 Preprocessing/cleaning/labeling

118 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization
119 or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction,
120 removal of instances, processing of miss- ing values)?** Language assignment of
121 both abstract and keyphrases was done using the method described earlier. Then when
122 separating into the different tasks, some instances were dropped due to lack of labels. All
123 original data is however provided with the dataset.

124 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data**
125 **(e.g., to support unanticipated future uses)?** Yes, the raw data is available in a
126 jsonlines format along the curated dataset.

127 **Is the software that was used to preprocess/clean/label the data available?** Yes,
128 all preprocessing is available in the github².

129 **Any other comments?** No.

130 5 Uses

131 **Has the dataset been used for any tasks already?** The dataset has been used for
132 keyphrases generation in a paper submitted to Neurips dataset 2022.

133 **Is there a repository that links to any or all papers or systems that use the**
134 **dataset?** Not at the moment.

135 **What (other) tasks could the dataset be used for?** The dataset could be used in a
136 reverse direction, that is to say for abstract or title generation from the keyphrases, as well
137 as any downstream tasks that use keyphrases. One could also not use the keyphrases and
138 do for example a task of title generation from the abstract or the reverse.

139 **Is there anything about the composition of the dataset or the way it was collected**
140 **and preprocessed/cleaned/labeled that might impact future uses?** No.

141 **Are there tasks for which the dataset should not be used?** No.

142 **Any other comments?** No.

143 6 Distribution

144 **Will the dataset be distributed to third parties outside of the entity (e.g.,**
145 **company, institution, organization) on behalf of which the dataset was created?**
146 The dataset will be available to all for research purposes under the creative common license.

147 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?**
148 The dataset will be hosted on the RALI website³, and downloading instructions will be
149 available on the github linked above.

150 **When will the dataset be distributed?** The dataset is available for download since
151 may 2022.

152 **Will the dataset be distributed under a copyright or other intellectual property**
153 **(IP) license, and/or under applicable terms of use (ToU)?** The dataset is distributed
154 under the creative commons public license.

155 **Have any third parties imposed IP-based or other restrictions on the data**
156 **associated with the instances?** No.

157 **Do any export controls or other regulatory restrictions apply to the dataset or**
158 **to individual instances?** No.

159 **Any other comments?** No.

160 7 Maintenance

161 **Who will be supporting/hosting/maintaining the dataset?** The dataset will be
162 hosted, supported, and maintained by the RALI team and the first author.

²<https://github.com/smolPixel/French-keyphrase-generation>

³<http://rali.iro.umontreal.ca/rali/en>

163 **How can the owner/curator/manager of the dataset be contacted (e.g., email**
164 **address)?** The author can be contacted at frederic.piedboeuf@umontreal.ca

165 **Will the dataset be updated (e.g., to correct labeling errors, add new instances,**
166 **delete instances)?** If enough new instances are available that it becomes pertinent to
167 create a new dataset and repeat the study, a follow-up will be made on this subject, and a
168 new paper describing the differences between the old and new datasets submitted. If errors
169 are detected, the dataset will be corrected and changes described on both the RALI website
170 and the Github.

171 **If the dataset relates to people, are there applicable limits on the retention of**
172 **the data associated with the instances (e.g., were the individuals in question told**
173 **that their data would be retained for a fixed period of time and then deleted)?**
174 No.

175 **Will older versions of the dataset continue to be supported/hosted/maintained?**
176 Yes. In both cases or errors or follow up studies, the changed will be described on Github
177 and the RALI website.

178 **If others want to extend/augment/build on/contribute to the dataset, is there**
179 **a mechanism for them to do so?** We actively encourage others to extend the dataset
180 by getting more multilingual data. Please contact the first author if you have an extension
181 to the dataset.

182 **Any other comments?** No.