

1 **LiteRev, an Automation Tool to Support**
2 **Literature Reviews: A Case Study on Acute and**
3 **Early HIV Infection in Sub-Saharan Africa**

4
5 Erol Orel ^{a,*}, Iza Ciglenecki ^{a,b}, Amaury Thiabaud ^a, Alexander Temerev ^a,
6 Alexandra Calmy ^{a,c}, Olivia Keiser ^{a,#}, Aziza Merzouki ^{a,#}

7
8 ^a Institute of Global Health, University of Geneva, Geneva, Switzerland

9 ^b Médecins Sans Frontières, Geneva, Switzerland

10 ^c HIV/AIDS Unit, Division of Infectious Diseases, Geneva University Hospital, Geneva,
11 Switzerland

12
13 * Corresponding author:

14 Erol Orel

15 Institute of Global Health, University of Geneva

16 Chemin des Mines 9, 1202 Geneva, Switzerland

17 Tel. +41 22 379 04 58

18 Erol.Orel@unige.ch

19
20 # Co-last authors

21
22 Word count: Abstract 439 words; main text 4'833 words; 1 table; 4 figures; 1 supporting material

23 Keywords: Africa, acute, early, HIV, Literature Review, Automation, Clustering, Topic, LiteRev

24

25 **Abstract**

26 **Background**

27 Literature Reviews (LRs) identify, evaluate, and synthesise relevant papers to a
28 particular research question to advance understanding and support decision making.
29 However, LRs, especially traditional systematic reviews are slow, resource intensive,
30 and are outdated quickly.

31 **Objective**

32 Using recent Natural Language Processing (NLP) and Unsupervised Machine Learning
33 (UML) methods, this paper presents a tool named LiteRev that supports researchers in
34 conducting LRs.

35 **Methods**

36 Based on the user's query, LiteRev can perform an automated search on different open-
37 access databases and retrieve relevant metadata on the resulting papers. Papers
38 (abstracts or full texts) are text processed and represented as a Term Frequency-
39 Inverse Document Frequency (TF-IDF) matrix. Using dimensionality reduction
40 (PaCMAP) and clustering (HDBSCAN) techniques, the corpus is divided into different
41 topics described by a list of keywords. The user can select one or several topics of
42 interest, enter additional keywords to refine their search, or provide key papers to the
43 research question. Based on these inputs, LiteRev performs an iterative nearest
44 neighbours search, and suggests a list of potentially interesting papers. The user can
45 tag the relevant ones and trigger a new search until no additional paper is suggested for
46 screening. To assess the performance of LiteRev, we ran it in parallel to a manual LR
47 on the burden and care for acute and early HIV infection in sub-Saharan Africa. We
48 assessed the performance of LiteRev using True and False Predictive Values, recall
49 and Work Saved over Sampling.

50 **Results**

51 We extracted, text processed and represented into a TF-IDF matrix 631 unique papers
52 from PubMed. The topic modelling module identified 5 main topics and 16 topics
53 (ranging from 13 to 98 papers) and extracted the 10 most important keywords for each.
54 Then, based on 18 key papers, we were able to identify 2 topics of interest with 7 key

55 papers in each of them. Finally, we ran the k-nearest neighbours module and LiteRev
56 suggested first a list of 110 papers for screening, among which 45 papers were
57 confirmed as relevant. From these 45 papers, LiteRev suggested 26 additional papers,
58 out of which 8 were confirmed as relevant. At the end of the iterative process (4
59 iterations), 193 papers out of 613 papers in total (31.5% of the whole corpus) were
60 suggested by LiteRev. After title/abstract screening, LiteRev identified 64 out of the 87
61 relevant papers (i.e., recall of 73.6%). After full text screening, LiteRev identified 42 out
62 of the 48 relevant papers (i.e., recall of 87.5%, and Work Saved over Sampling of
63 56.0%).

64 Conclusions

65 We presented LiteRev, an automation tool that uses NLP and UML methods to
66 streamline and accelerate LRs and to support researchers in getting quick and in-depth
67 overviews on any topic of interest.

68

69 Introduction

70 Recently, the traditional emphasis of Literature Reviews (LRs) in identifying, evaluating,
71 and synthesising all relevant papers to a particular research question has shifted
72 towards mapping research activity and consolidating existing knowledge [1]. Despite
73 this broader scope, manual LRs are still error-prone, time and resource-intensive and
74 have become ever more challenging over the years due to the increasing number of
75 papers published in academic databases. It is estimated that within two years of
76 publication, about one fourth of all LRs are outdated, as reviewers fail to incorporate
77 new papers on their topic of interest [2,3].

78 To shorten time to completion, automation tools have been developed to either fully
79 automate or semi-automate one or more specific tasks involved in conducting a LR,
80 such as screening titles and abstracts [4,5], sourcing full texts or automating data
81 extraction [6]. Also, recent advances in Natural Language Processing (NLP) and
82 Unsupervised Machine Learning (UML) have produced new techniques that can
83 accurately mimic manual LRs faster and at lower costs [7,8,9]. In Vienna, in 2015, the
84 International Collaboration for the Automation of Systematic Reviews (ICASR) was
85 initiated to establish a set of principles to enable tools to be developed and integrated
86 into toolkits [10]. Also, since 2021, an open source machine learning framework named
87 ASReview is under development for efficient and transparent systematic reviews [11].

88 In 2020, our group of researchers started developing an automation tool for LRs [12] in
89 order to obtain a comprehensive overview of the sociobehavioral factors influencing HIV
90 prevalence and incidence in Malawi. In this paper, we propose LiteRev, a new version
91 of our automation tool that overcomes some of the shortcomings of our previous tool.
92 While previously restricted to Paperity, PubMed, PubmedCentral, JSTOR, and arXiv,
93 the search now includes two primary preprint services in the field of epidemiology and
94 medical sciences, biorXiv and medRxiv, and CORE, a large collection of open-access
95 research papers. Also in our previous tool, the search was systematically performed on
96 the papers' full texts and references were included in the processed text. In LiteRev, the
97 user can choose to focus on the abstract or on the full text and include or exclude the
98 references. In addition, multiple parallel Application Programming Interfaces (APIs)
99 connections to each database have been implemented, allowing for faster retrieval of
100 papers. Since two years, NLP and UML have rapidly evolved, and LiteRev makes use
101 of the most recent text processing, embedding and clustering techniques. Finally, we
102 added a k-nearest neighbour's search module that allows the user to find papers of high
103 similarities with key papers to the research question.

104 In order to assess the performance of LiteRev, we conducted a manual LR using one
105 open-access database, PubMed, and two subscription-based databases, Embase and

106 Web of Science. All papers available by the 20th of December 2022 and related to
107 burden and care for acute and early HIV infection in sub-Saharan Africa were retrieved
108 and after removing duplicates, unique papers were screened for relevance. After
109 screening, papers from PubMed identified as relevant by the manual LR were compared
110 to the list of suggested papers by LiteRev. We discussed the performance using
111 standard classification metrics such as True and False Predictive Values, recall, and
112 Work Saved over Sampling (WSS).

113 **Methods**

114 **LiteRev**

115 **Metadata collection and text processing**

116 Based on the user's query, LiteRev can perform an automated search, using the
117 corresponding APIs, on 8 different open-access databases: PubMed, PubMedCentral,
118 CORE, JSTOR, Paperity, arXiv, biorXiv and medXriv. Available metadata, i.e., list of
119 authors and their affiliations, MesH keywords, Digital Object Identifier (DOI), title,
120 abstract, publication date, journal provider, and URL of the PDF version of the full text
121 paper, are retrieved and stored in a PostgreSQL database. If the full text is not available
122 as a metadata, it is extracted automatically from the available PDF file, then, references,
123 acknowledgements, and other unnecessary terms are removed and the remaining text
124 is checked to confirm that it still satisfies the search terms. Duplicated papers are
125 merged, collecting as much information as possible on the same paper from different
126 sources. Depending on the user's needs and requirements, LiteRev can be performed
127 on the full text or on the abstract.

128 Natural Language Processing has evolved rapidly, and, in particular, some powerful
129 tools were developed to process text data much more efficiently. We included those
130 features into LiteRev (Gensim [13] and Spacy [14]). After removing papers with empty
131 text, emails, newline characters, single quotes, internet addresses, and punctuation are
132 stripped and papers that do not fulfil the language(s) (one or multiple) chosen by the
133 user are discarded. Sentences are then split into words and lemmatised to remove as
134 many variations of the same word as possible. Words belonging to a list of stop-words
135 (i.e., words that are not informative) and words with less than three characters are also
136 removed. Next, bigrams, trigrams and four-grams (i.e., the combination of two, three
137 and four words) are created using a probabilistic measure. In practice, n-gram models
138 are highly effective in modelling language data. Finally, we remove words that are in
139 only one paper or words that occur too often (i.e., in more than 80% of the corpus) to
140 have a significant meaning.

141 Clustering and topic modelling

142 Topic modelling allows organising documents into clusters based on similarity, and
143 identifying abstract topics covered by similar papers. In LiteRev, it allows the user to
144 broaden the search strategy and to get a more comprehensive and organised overview
145 of the corpus. It can also help to quickly discard a pool of papers when searching the
146 literature for a specific topic and significantly reduce the amount of text to verify
147 manually.

148
149 After abstracts or full texts are processed, each paper's remaining words (namely bags
150 of words) are represented as a Term Frequency-Inverse Document Frequency (TF-IDF)
151 matrix which is computed using the Scikit-Learn package [15]. A TF-IDF matrix is similar
152 to a document (in row) - word (in column) co-occurrence matrix, normalised by the
153 number of papers in which the word is present. Less meaningful words, often present in
154 the corpus, get a lower score. Because of the often-high dimension of the TF-IDF matrix
155 (size of corpus x size of vocabulary), it is needed to embed the matrix using a Pairwise
156 Controlled Manifold Approximation (PaCMAP) dimensionality reduction technique [16].
157 The corpus is then divided into different clusters using the Hierarchical Density-Based
158 Spatial Clustering of Applications with Noise (HDBSCAN) algorithm [17]. While
159 HDBSCAN allows to classify some papers as noise, we have decided to consider those
160 as a cluster in itself.

161
162 PaCMAP and HDBSCAN have several important hyperparameters that need to be
163 determined. Table S1 in the Supplementary Material represents the 4 hyperparameters
164 involved and the ranges of their possible values. To find the best set of
165 hyperparameters possible, we use the Tree-structured Parzen Estimator algorithm
166 implemented by the Optuna package [18] and store the results of 500 trials in the
167 previously created PostgreSQL database. The Density-Based Clustering Validation
168 (DBCV), a weighted sum of "Validity Index" values of clusters [19], is the considered
169 performance metric to compare the different sets. Its value varies between 0 and 1
170 when used with HDBSCAN, with larger values providing better clustering solutions. This
171 metric takes the noise into account and captures the shape property of clusters via
172 densities and not distances. For coherency check, another metric is computed, the
173 Silhouette coefficient, which measures cluster cohesiveness and separation with an
174 index between -1 to 1, with larger values providing better clustering solutions [20].

175
176 If after 500 trials, the DBCV score is below 0.5, another round of 500 trials is performed,
177 and so on until a DBCV score equal or above 0.5 is reached. Once the values of the
178 hyperparameters that maximise the DBCV score are determined, obtained clusters that
179 are larger than 25% of the corpus are clustered again with the same entire procedure
180 described above (starting from the text processing). Once each cluster is smaller than

181 25% of the corpus, its 10 most important words are extracted using the YAKE package
182 [21] to ensure interpretability and define topics. This supports the user in getting a quick
183 overview of the corpus and, if desired, to select one or more topics of interest for further
184 exploration. They can then also enter additional keywords to refine his search.

185 Nearest neighbours

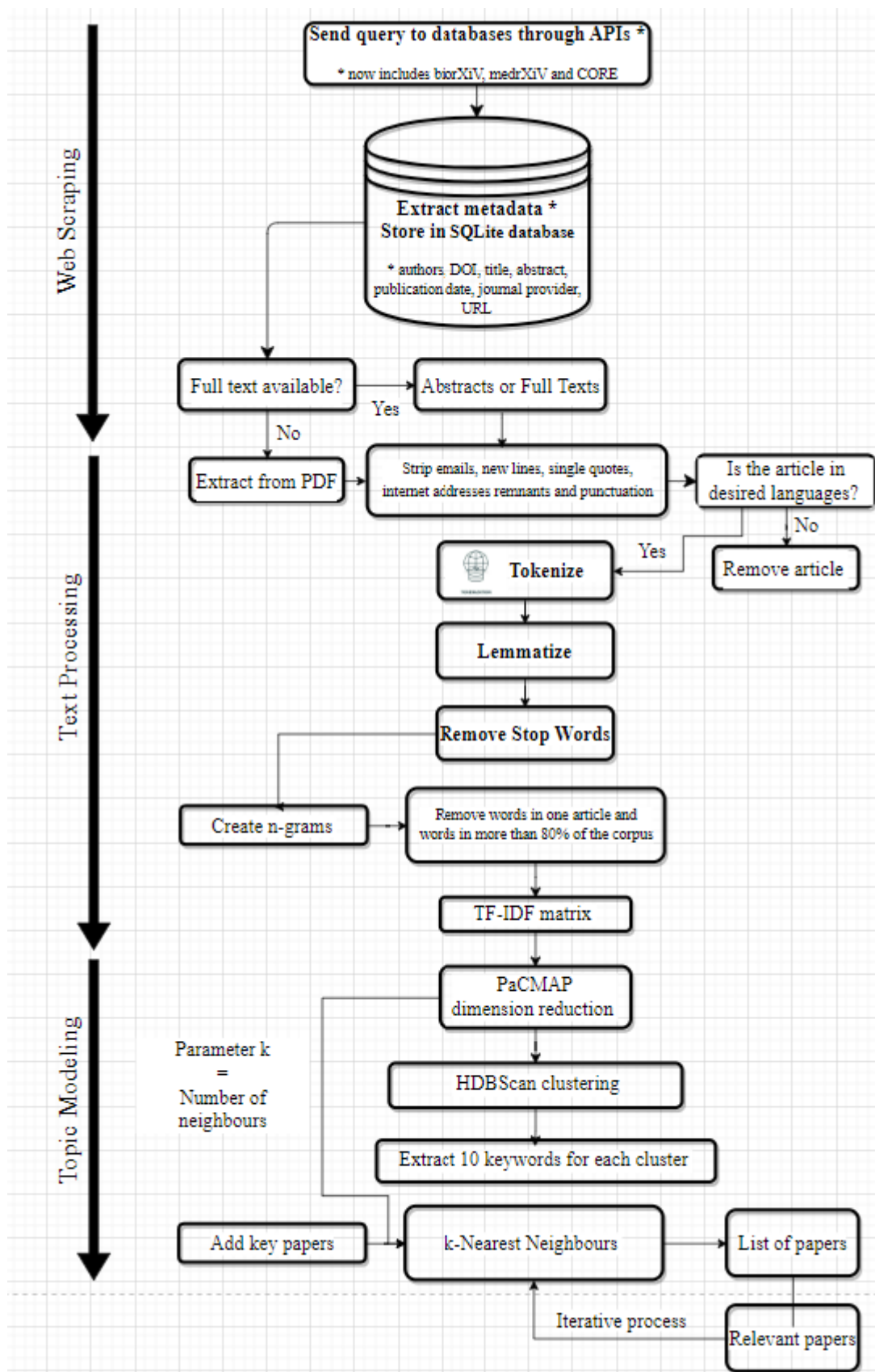
186 LiteRev allows the user to define papers in the corpus as being key to the research
187 question (or to add them). Using the k-nearest neighbour (k-NN) algorithm from the
188 Scikit-Learn package [15], a list of potentially relevant papers is provided to the user.
189 Papers deemed to be relevant are tagged by the user and considered as new key
190 papers. This process is iterated as long as relevant papers are being identified
191 (generally 3 to 4 iterations). The initial value of the hyperparameter k, which represents
192 the number of nearest neighbours to be selected, is equal to the value of the number of
193 neighbours for PaCMAP obtained at the first clustering process. The dimension space is
194 the same as the number of dimensions obtained during the embedding process by
195 PaCMAP.

196
197 The list of relevant papers from the k-NN search and/or a list of papers about one or
198 more topics can then be exported in a csv or html format and their pdf retrieved and
199 stored in a zip folder. For visualisation and further exploration, an interactive 2D
200 representation of the corpus is available in a html format. Every dot, coloured according
201 to the cluster it belongs to, represents a paper with the following available information:
202 date, title, 10 most important keywords of the cluster's topic and the cluster number.
203 When clicking on a paper (dot), a direct access to the full text is provided using the
204 URL. The diagram in Figure 1 shows the entire process flow of LiteRev.

205

206 Figure 1: Diagram of LiteRev process

207



209 Manual Literature Review

210 The manual LR aimed at summarising the current evidence on burden and care
211 provided for acute and early HIV infection (AEHI) in sub-Saharan Africa, to inform
212 policy, practice and research in future, addressing the following questions:

- 213 - What is the prevalence of AEHI in sub-Saharan Africa among people being
214 tested for HIV?
- 215 - What models of care have been used for AEHI diagnosis and care, including
216 treatment, partner notifications and behaviour change?
- 217 - What linkage to care has been reached?
- 218 - What facilitators and barriers to AEHI care were identified?

219
220 We searched all papers in PubMed, Embase and Web of Science related to burden and
221 care for acute and early HIV infection in sub-Saharan Africa that were published from
222 the inception of the databases to December 20 2022, using the query: ("early hiv" OR
223 "primary hiv" OR "acute hiv" OR "HIV Human immuno deficiency virus" OR ("Window
224 period" AND HIV)) AND ("Africa South of the Sahara" OR "Central Africa" OR "Eastern
225 Africa" OR "Southern Africa" OR "Western Africa" OR "sub-saharan africa" OR
226 "subsaharan africa" OR angola OR benini OR botswana OR "burkina faso" OR burundi
227 OR cameroon OR "cape verde" OR "central africa" OR "central african republic" OR
228 chad OR comoros OR congo OR "cote d ivoire" OR "democratic republic congo" OR
229 djibouti OR "equatorial guinea" OR eritrea OR eswatini OR ethiopia OR gabon OR
230 gambia OR ghana OR guinea OR "guinea-bissau" OR kenya OR lesotho OR liberia OR
231 madagascar OR malawi OR mali OR mayotte OR mozambique OR namibia OR niger
232 OR nigeria OR rwanada OR sahel OR "sao tome and principe" OR senegal OR "sierra
233 leone" OR somalia OR "south africa" OR "south sudan" OR sudan OR tanzania OR
234 togo OR uganda OR zambia OR zimbabwe)". This query is specific to PubMed syntax
235 and is the exact same both for the manual LR and for LiteRev. Syntax specific queries
236 for the manual LR in Embase and Web of Science are to be found in the Supplementary
237 Material. Papers retrieved from Embase and Web of Science have not been used by
238 LiteRev and will not be part of the comparison and performance assessment but its
239 results will be discussed in the Results and Discussion section.

240
241 The studies were included if they described AEHI prevalence among population tested
242 for HIV and/or describe the diagnostic strategy, model of care and/or linkage to care for
243 AEHI, including studies looking at perceptions and barriers among patients and staff.
244 Only studies conducted in sub-Saharan Africa were included. We followed the JBI
245 methodology for conducting LRs [22] and papers identified by the databases were
246 uploaded into Rayyan [23]. Duplicates were deleted and the screening process, on titles
247 and abstracts, was conducted independently by 2 reviewers (EO and IC). Selected

248 papers were further manually screened based on full text for eligibility against inclusion
249 criteria. LiteRev was run in parallel on the abstracts only but results were compared
250 both to the title/abstract screening phase and to the full text screening phase of the
251 manual LR.

252 Performance Comparison

253 In order to assess the performance of LiteRev, we compared the results from the
254 manual LR to the same review conducted using LiteRev. Relevant and not relevant
255 papers, as identified by the manual LR during the title/abstract screening phase and the
256 full text screening phase, were defined as true labels. Suggested and not suggested
257 papers by LiteRev were considered as predicted labels. Based on these figures, two
258 confusion matrices were produced. Positive and Negative Predictive Values (% of
259 relevant and not relevant papers correctly identified; PPV and NPV), recall (number of
260 relevant papers identified using LiteRev among those identified using manual review)
261 and Work Saved over Sampling (WSS) [24,25], percentage of abstracts or full texts that
262 the user did not have to read because they were not suggested for screening by
263 LiteRev, were computed and discussed.

$$\text{WSS} = \frac{(\text{True Negatives} + \text{False Negatives})}{\text{True Negatives} + \text{False Negatives} + \text{True Positives}} - (1 - \text{Recall})$$

264 where:

- 265 - True negatives is the number of non-relevant abstracts that were correctly
266 identified as non-relevant by LiteRev, i.e. that were not suggested by LiteRev for
267 screening,
- 268 - False negatives is the number of relevant abstracts incorrectly classified as non-
269 relevant by LiteRev.

270 Results

271 LiteRev

272 Text Processing and Topic Modelling

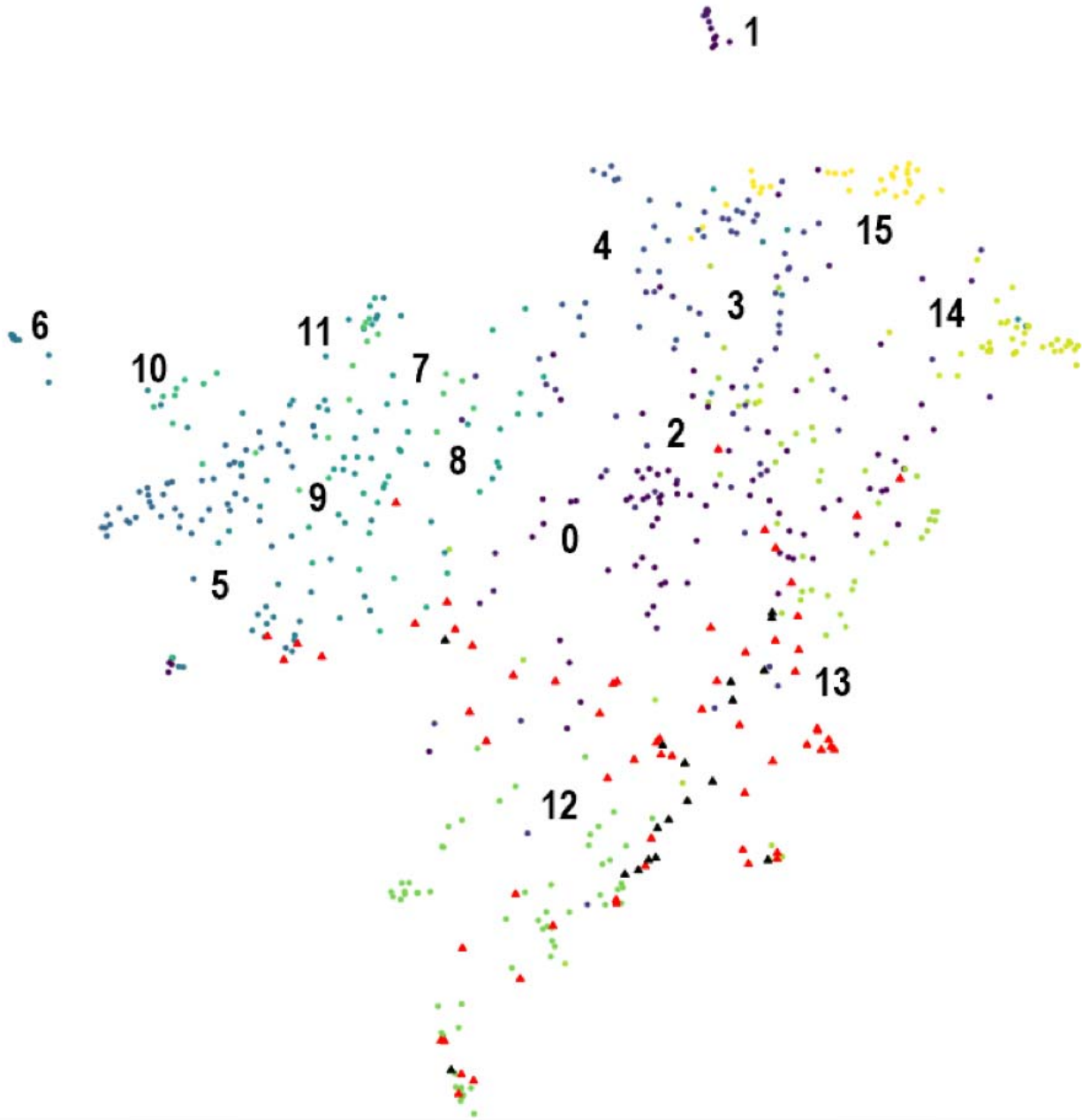
273 Based on the search strategy described in the Methods section, we obtained 653
274 papers with metadata directly from PubMed and added one key paper given by the user
275 that was not present in the list of retrieved papers. After removing duplicates (3), papers
276 with no abstract available (15), those not in english (3), and empty abstract after text

277 processing (2), 631 unique papers were transformed in a TF-IDF matrix comprised of
278 631 rows representing the corpus and 3'136 columns representing the unique words
279 (vocabulary), including n-grams.

280 For the first embedding and clustering process, a DBCV score of 0.533 was obtained
281 after the first 500 trials with the following best set of hyperparameters: PaCMAP: 310
282 dimensions and 18 neighbours; HDBSCAN: minimum cluster size of 30 and minimum
283 samples of 7. This resulted in 5 main clusters composed of respectively 203, 193, 169,
284 35 and 31 papers. The 3 largest main clusters contained more than 25% of the total
285 number of papers in the corpus, which triggered 3 additional text processing,
286 embedding and clustering processes. The best set of hyperparameters for these
287 additional processes can be found in Table S1 in the Supplementary Material.

288 At the end, the pool of 203 papers were splitted in 5 clusters (with respectively 98, 41,
289 25, 21 and 18 papers), the pool of 193 papers in 7 clusters (with respectively 47, 40, 37,
290 22, 20 14 and 13 papers) and the pool of 169 papers in 2 clusters (with respectively 87
291 and 82 papers). In total, the corpus of 631 papers was divided in 16 clusters ranging
292 from 13 to 98 papers. Figure 3 shows the 2D map of the corpus with the 16 clusters
293 identified. Table 1 shows the corresponding 16 topics grouped by main topics,
294 described by their 10 most important keywords and the number of papers in each.

295 Figure 2: 2D representation of the corpus with the 16 clusters. Black triangles represent
296 the 18 key papers and red triangles represent the 64 relevant papers correctly identified
297 by LiteRev
298
299



300
301
302 Table 1: The 16 topics grouped by main topics (in blue) with the 10 most important
303 keywords, the number of papers and the number of relevant papers in total (key papers)
304

Topic	Keywords	# of papers	# of relevant papers (key)
woman, patient, risk, year, treatment, associate, month, incidence, testing, care			
0	woman, risk, year, incidence, high, man, partner, transmission, sexual, testing	98	9 (1)
1	cart, month, initiation, group, treatment, viral, rna, child, infant, week	18	0 (0)
2	risk, high, health, day, score, aehi, prevalence, care, diagnosis, population	41	6 (2)
3	patient, treatment, care, late, diagnosis, associate, testing, aor, datum, initiation	21	0 (0)
4	patient, disease, adult, infect, lymphadenopathy, cell, tuberculous, lymphadenitis, associate, present	25	0 (0)
cell, viral, response, virus, subtype, individual, primary, isolate, antibody, infect			
5	antibody, response, neutralize, vaccine, isolate, neutralization, epitope, env, primary, individual	47	0 (0)
6	subtype, resistance, drug, sequence, mutation, diversity, strain, primary, patient, recombinant	40	3 (0)
7	response, specific, associate, increase, immune, ifn, early, gag, point, level	37	0 (0)
8	level, viremia, acute, associate, early, individual, infect, load, cytokine, set	20	1 (0)
9	load, early, copy, log, plasma, subtype, woman, time, african, rna	22	5 (1)

10	isolate, primary, tropic, individual, derive, clone, strain, infect, dual, sequence	13	0 (0)
11	response, immune, phi, specific, control, activation, plasma, individual, acute, cytokine	14	0 (0)
test, testing, blood, ahi, risk, acute, positive, care, sample, assay			
12	blood, assay, sample, donor, positive, risk, incidence, antibody, estimate, acute	82	21 (7)
13	ahi, care, participant, health, intervention, patient, diagnosis, early, acute, risk	87	37 (7)
infant, mother, week, transmission, child, month, age, woman, test, infect			
14	infant, mother, week, transmission, child, month, age, woman, test, infect	35	0 (0)
child, year, mortality, age, infect, treatment, patient, associate, month, clinical			
15	child, year, mortality, age, infect, treatment, patient, associate, month, clinical	31	0 (0)

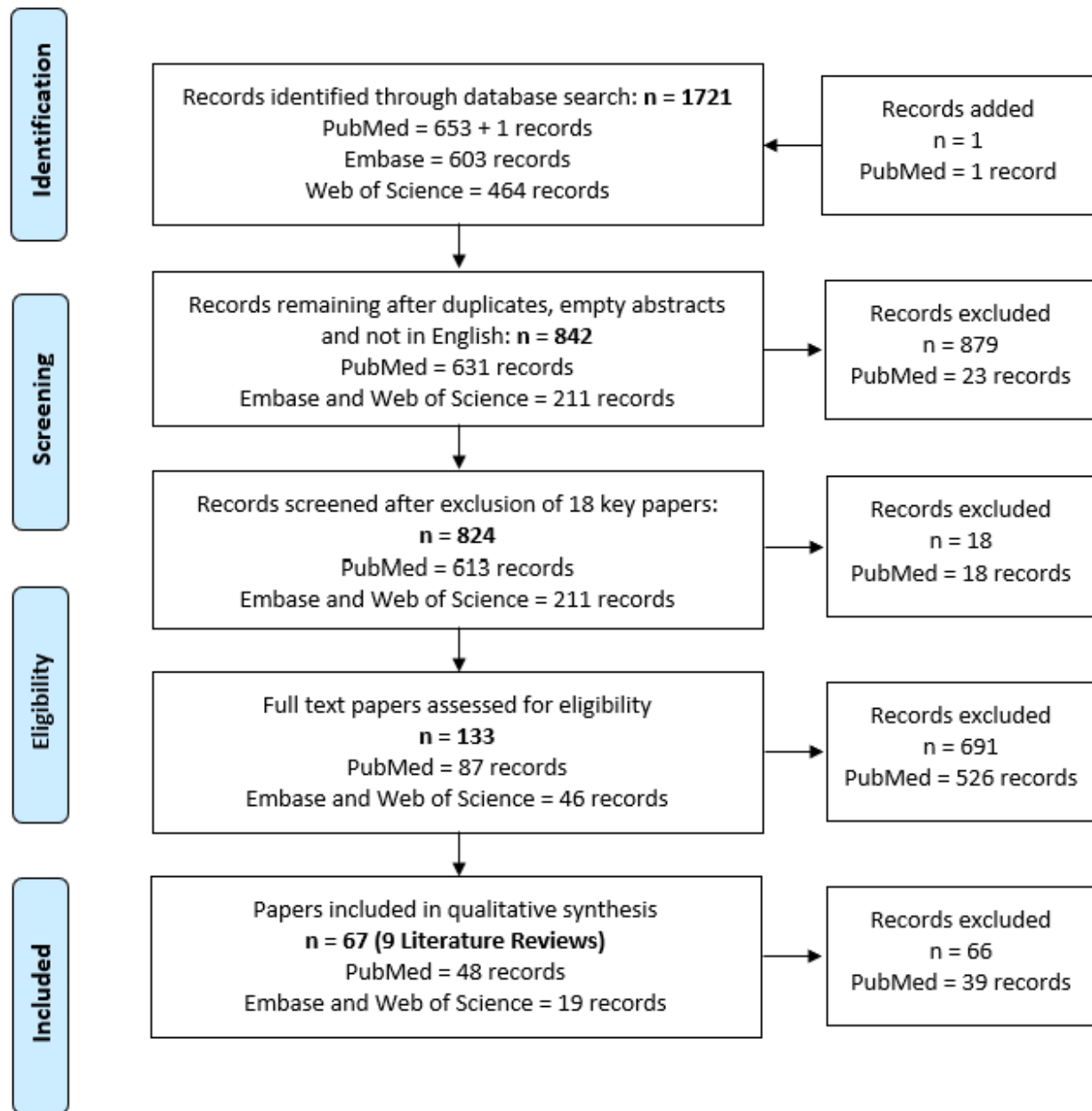
305 Manual Literature Review

306 Using the search query described in the Methods section, 1'721 records were retrieved,
307 among which 653 were from PubMed and 1'067 records from 2 subscription-based
308 databases, namely Embase and Web of Science. 879 records were excluded after
309 removing duplicates, empty abstracts and papers that were not written in English. This
310 resulted in 631 unique papers in PubMed and 211 unique papers in Embase and Web
311 of Science. We also removed the 18 key papers out of the PubMed corpus before the
312 screening phases. In total, 613 papers in PubMed were screened at the title and
313 abstract level and 87 of them were relevant to the research question. After the full text
314 screening phase on these 87 relevant papers, we found 48 papers to be relevant with
315 the manual LR.

316 Out of the 211 unique papers from Embase and Web of Science, 46 papers were found
317 relevant to the research question after the title/abstract screening phase (i.e., 34.6% of
318 all relevant papers), and 19 after the full text screening phase (i.e., 28.4% of all relevant
319 papers) (Figure 3). From these 19 relevant papers, 3 were conference abstracts and 1
320 paper was kept only based on its title and abstract as the full text couldn't be found.
321 These 221 papers were not part of PubMed, and hence, not available to LiteRev.

322 Figure 3: PRISMA diagram of the manual LR related to burden and care for acute and
323 early HIV infection in sub-Saharan Africa

324



325

326 Nearest Neighbours Search and Performance Comparison

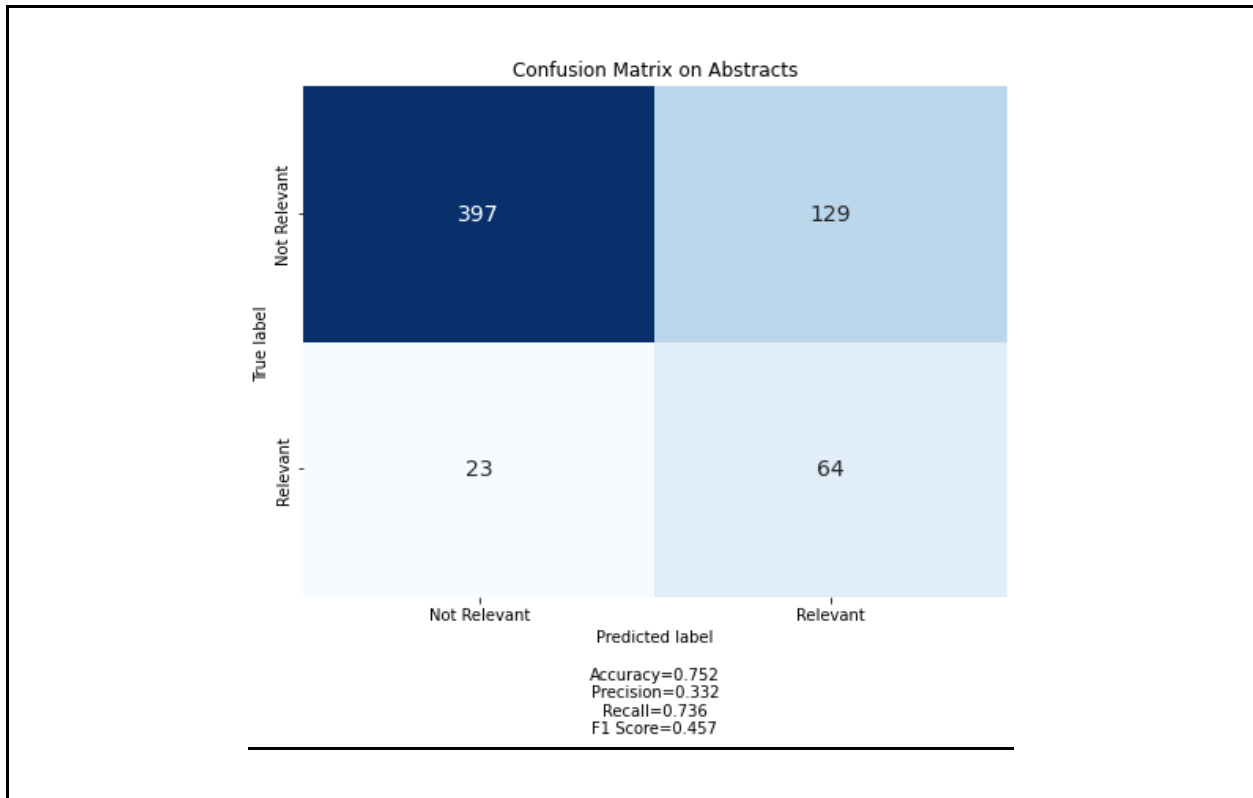
327 We were provided by the user (IC) a list of 18 key papers on the topic. With these 18
 328 papers, we performed a k-nearest neighbours search on the corpus, embedded into 310
 329 dimensions, with k=18, the number of the nearest neighbours for PaCMAP that
 330 maximised the DBCV score of the first clustering process. The first k-NN search
 331 suggested 110 papers, including 45 of the relevant papers identified by the manual LR
 332 title/abstract screening (precision of 41%). Based on these 45 relevant papers, the
 333 second k-NN iteration suggested 26 additional papers out of which 8 were confirmed as

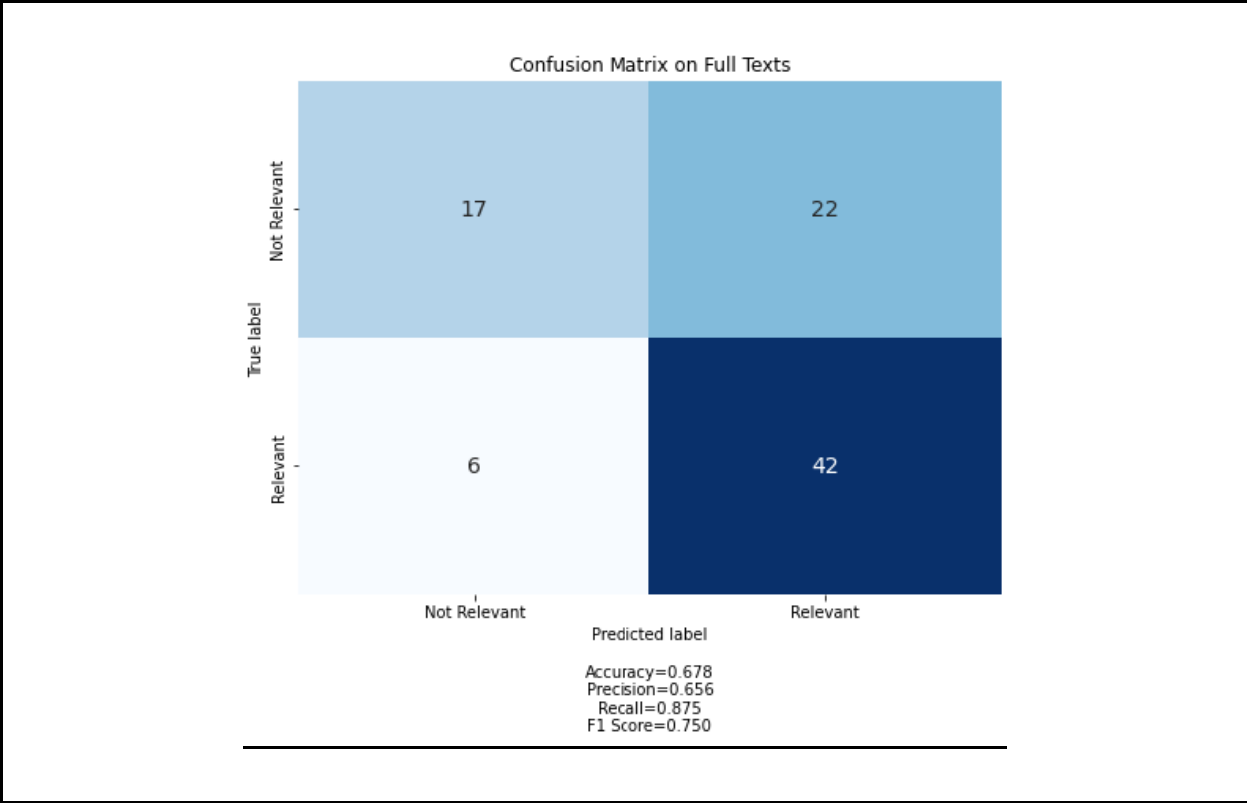
334 relevant (precision of 31%). The third iteration found 9 more relevant papers out of 38
335 papers suggested (precision of 24%). The fourth and last iteration suggested 19 papers
336 out of which 1 was relevant (precision of 5%).

337 In total, 193 papers out of the 613 papers were suggested by LiteRev. Suggested
338 papers included 64 of the 87 papers identified as relevant during the title/abstract
339 screening of the manual LR. Figure 3 maps the key papers (black triangles) and the
340 relevant papers (red triangles) identified at the title/abstract screening level of the
341 manual LR, and that were correctly classified as relevant by LiteRev. Table 1 indicates
342 the number of key papers and the number of relevant papers in each topic.

343 Figure 4 (top panel) summarises the above results and represents the confusion matrix
344 between LiteRev (Predicted labels) and the manual LR (True labels) after the
345 title/abstract screening phase. Based on these numbers, the PPV was 33.2%, the NPV
346 was 94.5% and the recall was 73.6%, which led to a WSS of 42.1%.

347 Figure 4: Confusion Matrices based on the results of (top panel) the title/abstract
348 screening, and (bottom panel) full-text screening performed during the manual LR.





349 The 64 relevant papers found by LiteRev belonged essentially to 2 topics (30 relevant
 350 papers in one and 14 relevant papers in the other). The topic that contained 30 relevant
 351 papers had 87 papers in total, and covered early diagnosis, care seeking and
 352 interventions during Acute HIV Infection stage (keywords: « ahi, care, participant,
 353 health, intervention, patient, diagnosis, early, acute, risk »). The topic that contained 14
 354 relevant papers had 82 papers, and covered the detection of AHI by antibody assays
 355 and incidence estimate (keywords: « blood, assay, sample, donor, positive, risk,
 356 incidence, antibody, estimate, acute »). Screening 53 additional papers (those not
 357 suggested by the nearest neighbours search) from these 2 topics would allow the user
 358 to identify 3 additional relevant papers.

359 After the full text screening phase of the manual LR, 48 out the 87 relevant papers from
 360 the title and abstract screening phase were deemed relevant to the research question.
 361 The list of (64) relevant papers suggested by LiteRev (based on abstracts only),
 362 included 42 out of the 48 papers confirmed as relevant after the full text screening
 363 phase of the manual LR. Figure 4 (bottom panel) summarises the above results and
 364 represents the confusion matrix between LiteRev (Predicted labels) and the manual LR
 365 (True labels) after the full text screening phase. Based on these numbers, the PPV was
 366 65.6%, the NPV was 26.1% and the recall was 87.5%, which led to an additional WSS
 367 of 13.9% for an overall WSS of 56.0% compared to the manual LR.

368 Processing time

369 The processing time represents the overall computation time taken by LiteRev to
370 complete the entire process of metadata retrieval, processing, clustering, and neighbour
371 search. It doesn't include the time that the user took to check the relevance of the
372 suggested papers. The percentage of time saved by the user, is expressed by the work
373 saved over sampling (WSS) metric.

374 It took 5 minutes for LiteRev to retrieve the metadata of the 653 papers and to text
375 process the remaining 631 abstracts and transform it into a TF-IDF matrix. It took an
376 additional two days for the main clustering and the 3 additional clustering processes.
377 Each trial of the optimization process with a specific set of hyperparameters required on
378 average 1 minute of computation. With 3'000 trials in total (500 for the main clustering
379 process, 1'000 for the first two additional clustering processes and 500 for the last one)
380 run sequentially, this led to an additional 50 hours, i.e., roughly 2 days, to complete the
381 entire optimisation process. This computation time can be substantially reduced by
382 running the trials in parallel. Finally, the nearest neighbours are obtained almost
383 instantaneously.

384 Discussion

385 Principal Results

386 We presented LiteRev, an automation tool that uses NLP and UML methods to support
387 researchers in different steps of a manual LR. The identification of papers to be
388 included in a LR is a critical and time-intensive process, with the majority of time spent
389 in screening thousands of papers for relevance. By combining text processing, literature
390 mapping, topic modelling and similarity-based search, LiteRev provides a fast and
391 efficient way to remove duplicates, select papers from specific languages, visualise the
392 corpus on a 2D map, identify the different topics covered when addressing the research
393 question and suggest a list of potentially relevant papers to the user based on their input
394 (e.g., prior knowledge of key papers).

395
396 Preliminary usage of LiteRev showed that it significantly reduced the researcher's
397 workload and overall time required to perform a LR. Compared to a manual LR, LiteRev
398 correctly identified 87.5% of relevant papers (recall), by screening only 31.5% of the
399 whole corpus, which corresponds to a total Work Saved over Sampling of 56.0% (WSS)
400 at the end of the full text screening phase. In addition, the actual time spent on running
401 LiteRev and retrieving the results was relatively short, and the user was free in the

402 meantime to focus on other work. The text processing and the nearest neighbours
403 search took no more than 5 minutes of computation for 631 papers.

404

405 With its topic modelling capability, LiteRev aims at summarising current evidence on a
406 specific research question, to inform policy, practice and research. For our use case,
407 LiteRev identified 5 main topics and 16 different topics related to acute and early HIV
408 infection in Sub-Saharan Africa, allowing the researcher to have an overview on the
409 different perspectives related to this research question. Finding 61 out of the 105
410 relevant papers after the title and abstract screening phase (including the key papers) in
411 only 2 topics validates the quality of the clustering.

412 Limitations

413 LiteRev is currently limited to open-access databases that provide free APIs to abstract
414 or full text papers. Databases often used for LRs, such as Embase or Web of Science
415 do not provide APIs access, require a subscription for accessing full text papers or do
416 not allow for text mining and machine learning analysis. Hence, 19 relevant papers
417 identified in Embase or Web of Science were not available to LiteRev. Also, when
418 performed on full texts, LiteRev currently works on digitally-generated PDFs, but not on
419 image-only (scanned) PDFs.

420

421 Another limitation concerns the possibility of sharing the list of potentially relevant
422 papers with other users/reviewers. LiteRev does not offer this functionality yet, hence
423 double screening of papers and comparison of results is not possible at the moment. To
424 overcome this limitation, the user has the option to export its list of papers into a csv
425 format uploadable on Rayyan or other similar softwares for systematic reviews.

426

427 As of today, LiteRev is still intended to complement rather than replace full systematic
428 reviews. Finally, by January 2023, no public web-based User Interface (UI) is available
429 yet.

430 Future work

431 O'Connor et al. [26] found that overall, many of the automation tools were not
432 compatible with current practice, because they were not easily integrated into current
433 workflows and not particularly easy to use for nontechnical persons. Also, there was not
434 enough evidence of accuracy to earn the trust of reviewers. LiteRev is developed in an
435 iterative and interactive way by continuously integrating feedback from users and its
436 modules can easily be updated or replaced depending on the needs of the users and
437 the technical evolutions. We are further developing LiteRev by proposing a web
438 application with a user-friendly interface and by adding more functionality in order to

439 better automate the different stages of a LR. We are also planning to implement a living
440 review [27] by retrieving new papers on each research questions in our database (e.g.,
441 “HIV” AND “Africa”) on a regular basis (e.g., every month) and each new paper will be
442 text processed and assigned to the topic it belongs to using a predictive algorithm.
443 Although we compared the performance of LiteRev with one manual LR in this paper,
444 we plan to perform additional similar comparisons and performance evaluations in the
445 future, using other published LRs covering different topics.

446 **Conclusions**

447 We presented LiteRev, an automation tool that uses NLP and UML techniques to
448 support, facilitate and accelerate the conduction of Literature Reviews providing aid and
449 automation to different steps involved in this process. Its different modules (retrieval of
450 papers’ metadata from open-access databases using a search query, processing of
451 texts, embedding and clustering, and finding of nearest neighbours) can easily be
452 updated or replaced depending on the needs of the users and the technical evolutions.
453 As more papers are published every year, LiteRev not only has the potential to simplify
454 and accelerate LRs, but it also has the capability of helping the researcher get a quick
455 and in-depth overview on any topic of interest.

456 **Financial disclosure**

457 We acknowledge the support of the Swiss National Science Foundation (SNF
458 professorship grants n°196270 and n°202660 to Prof. O. Keiser), which funded this
459 study. The funder had no role in study design, data collection and analysis, decision to
460 publish, or manuscript preparation.

461 **Authors' Contributions**

462 EO and ATh wrote the code in Python. EO and AM obtained and analysed the results it
463 produced. EO wrote the first draft of the paper. IC conducted the manual LR of the use
464 case and IC and EO identified the relevant papers. EO, IC, Ath and AM helped write the
465 paper and EO, AM, OK, AC and IC reviewed the paper.

466 **Conflicts of Interest**

467 None declared.

468 Acknowledgments

469 Ms. Mafalda Vieira Burri, the librarian from the library of the University of Geneva who
470 helped define the search queries, Mr. Alexander Temerev, doctoral student at the
471 Institute of Global Health, who helped write part of the Python code on parallelising the
472 APIs.

473 References

- 474 1. Sutton, A., Clowes, M., Preston, L. and Booth, A. (2019), Meeting the review
475 family: exploring review types and associated information retrieval requirements.
476 Health Info Libr J, 36: 202-222. <https://doi.org/10.1111/hir.12276>
- 477 2. Van Dinter R, Tekinerdogan B, Catal C. Automation of systematic literature
478 reviews: A systematic literature review. Information and Software Technology.
479 Volume 136. 2021. 106589. ISSN 0950-5849.
480 <https://doi.org/10.1016/j.infsof.2021.106589>
- 481 3. K.G. Shojania, M. Sampson, M.T. Ansari, J. Ji, S. Doucette, D. Moher. How
482 quickly do systematic reviews go out of date? A survival analysis. Ann. Intern.
483 Med., 147 (2007), pp. 224-233
- 484 4. Harrison H, Griffin SJ, Kuhn I, Usher-Smith JA. Software tools to support title and
485 abstract screening for systematic reviews in healthcare: an evaluation. BMC
486 Medical Research Methodology. 2020 Jan 13;20(1):7.
- 487 5. Olofsson H, Brolund A, Hellberg C, Silverstein R, Stenstrom K, Osterberg M,
488 Dagerhamn J. Can abstract screening workload be reduced using text mining?
489 User experiences of the tool Rayyan. Research Synthesis Methods.
- 490 6. S.R. Jonnalagadda, P. Goyal, M.D. Huffman. Automating data extraction in
491 systematic reviews: a systematic review. Syst. Rev., 4 (2015), p. 78
- 492 7. Marshall C, Brereton P. Systematic review toolbox: a catalogue of tools to
493 support systematic reviews. In: Proceedings of the 19th International Conference
494 on Evaluation and Assessment in Software Engineering: ACM; 2015. p. 23
- 495 8. Clark J, Glasziou P, Del Mar C, Bannach-Brown A, Stehlik P, Scott AM. A full
496 systematic review was completed in 2 weeks using automation tools: a case
497 study. Journal of Clinical Epidemiology. Volume 121. 2020. Pages 81-90. ISSN
498 0895-4356. <https://doi.org/10.1016/j.jclinepi.2020.01.008>
- 499 9. Clark J, McFarlane C, Cleo G, Ishikawa Ramos C, Marshall S. The Impact of
500 Systematic Review Automation Tools on Methodological Quality and Time Taken
501 to Complete Systematic Review Tasks: Case Study. JMIR Med Educ.
502 2021;7(2):e24418. Published 2021 May 31. doi:10.2196/24418
- 503 10. Beller, E., Clark, J., Tsafnat, G. et al. Making progress with the automation of
504 systematic reviews: principles of the International Collaboration for the

- 505 Automation of Systematic Reviews (ICASR). *Syst Rev* 7, 77 (2018).
506 <https://doi.org/10.1186/s13643-018-0740-7>
- 507 11. van de Schoot R, de Bruin J, Schram R, Zahedi P, de Boer J, Weijdemans F, et al.
508 An open source machine learning framework for efficient and transparent
509 systematic reviews. *Nat Mach Intell*. 2021 Feb;3(2):125–33.
- 510 12. Thiabaud A, Triulzi I, Orel E, Tal K, Keiser O. Social, Behavioral, and Cultural
511 factors of HIV in Malawi: Semi-Automated Systematic Review. *J Med Internet*
512 *Res* 2020;22(8):e18747. DOI: 10.2196/18747
- 513 13. Rehurek R, Sojka P. Gensim–python framework for vector space modelling. NLP
514 Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic.
515 2011;3(2)
- 516 14. Honnibal M, Montani I. spaCy 2: Natural language understanding with Bloom
517 embeddings, convolutional neural networks and incremental parsing. 2017
- 518 15. Pedregosa F, Varoquaux, Ga"el, Gramfort A, Michel V, Thirion B, Grisel O, et al.
519 Scikit-learn: Machine learning in Python. *Journal of machine learning research*.
520 2011;12(Oct):2825–30
- 521 16. Wang Y., Huang H., Rudin C., and Shaposhnik Y., (2021), Understanding How
522 Dimension Reduction Tools Work: An Empirical Approach to Deciphering t-SNE,
523 UMAP, TriMap, and PaCMAP for Data Visualization. *Journal of Machine*
524 *Learning Research*. Volume 22. 201:1-73. [http://jmlr.org/papers/v22/20-](http://jmlr.org/papers/v22/20-1061.html)
525 [1061.html](http://jmlr.org/papers/v22/20-1061.html)
- 526 17. McInnes et al, (2017), hdbSCAN: Hierarchical density based clustering, *Journal of*
527 *Open Source Software*, 2(11), 205, doi:10.21105/joss.00205
- 528 18. Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A Next-
529 generation Hyperparameter Optimization Framework. *Proceedings of the 25th*
530 *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*
- 531 19. Moulavi, Davoud, et al. Density-based clustering validation. *Proceedings of the*
532 *2014 SIAM International Conference on Data Mining*. Society for Industrial and
533 *Applied Mathematics*, 2014
- 534 20. Rousseeuw PJ. Silhouettes: a Graphical Aid to the Interpretation and Validation
535 of Cluster Analysis. *Computational and Applied Mathematics*. 1987. 20: 53–65.
536 doi:10.1016/0377-0427(87)90125-7
- 537 21. Campos, R., Mangaravite, V., Pasquali, A., Jatowt, A., Jorge, A., Nunes, C. and
538 Jatowt, A. (2020). YAKE! Keyword Extraction from Single Documents using
539 Multiple Local Features. In *Information Sciences Journal*. Elsevier, Vol 509, pp
540 257-289.
- 541 22. Peters M, Godfrey C, McInerney P, Munn Z, Trico A, Khalil H. Chapter 11:
542 Scoping Reviews. In: Aromataris E, Munn Z, editors. *JBIM Manual for Evidence*
543 *Synthesis [Internet]*. JBI; 2020 [cited 2022 Jan 22].

- 544 23. Ouzzani M, Hammady H, Fedorowicz Z, Elmagarmid A. Rayyan—a web and
545 mobile app for systematic reviews. *Syst Rev* [Internet]. 2016 Dec [cited 2022 Jan
546 23];5(1):210.
- 547 24. Matwin S, Kouznetsov A, Inkpen D, Frunza O, O'Brien P. A new algorithm for
548 reducing the workload of experts in performing systematic reviews. *J Am Med
549 Inform Assoc.* 2010;17(4):446–53.
- 550 25. Cohen AM, Hersh WR, Peterson K, Yen PY. Reducing Workload in Systematic
551 Review Preparation Using Automated Citation Classification. *J Am Med Inform
552 Assoc.* 2006;13(2):206–19.
- 553 26. O'Connor AM, Tsafnat G, Thomas J, Glasziou P, Gilbert SB, Hutton B. A
554 question of trust: can we build an evidence base to gain trust in systematic
555 review automation technologies? *Syst Rev* [Internet]. 2019;8(1):143. Available
556 from: <http://dx.doi.org/10.1186/s13643-019-1062-0>
- 557 27. Brooker J, Synnot A, McDonald S, et al.: Guidance for the production and
558 publication of Cochrane living systematic reviews: Cochrane reviews in living
559 mode. Cochrane Collaboration. 2019

560 **Abbreviations**

- 561
- 562 - AEHI: Acute and early hiv infection
- 563
- 564 - API: Application programming interface
- 565
- 566 - DBCV: Density-based clustering validation
- 567
- 568 - HDBSCAN: Hierarchical density-based spatial clustering of applications with
- 569 noise
- 570
- 571 - HIV: Human immunodeficiency virus
- 572
- 573 - k-NN: k-nearest neighbours
- 574
- 575 - LR: Literature review
- 576
- 577 - NLP: Natural language processing
- 578
- 579 - PacMAP: Pairwise controlled manifold approximation
- 580
- 581 - TF-IDF: Term frequency - Inverse document frequency
- 582
- 583 - UML: Unsupervised machine learning
- 584
- 585 - WSS: Work saved over sampling
- 586

587 **Supplementary Material**

588

589 Query for Embase:

590

591 ('acute hiv':ab,ti OR 'early hiv':ab,ti OR 'primary hiv':ab,ti OR ('window period':ab,ti AND
 592 'human immunodeficiency virus':ab,ti)) AND ('africa south of the sahara'/exp OR 'sub-
 593 saharan africa':ab,ti OR 'subsaharan africa':ab,ti OR 'africa south of the sahara':ab,ti OR
 594 angola:ab,ti OR benin:ab,ti OR botswana:ab,ti OR 'burkina faso':ab,ti OR burundi:ab,ti
 595 OR cameroon:ab,ti OR 'cape verde':ab,ti OR 'central africa':ab,ti OR 'central african
 596 republic':ab,ti OR chad:ab,ti OR comoros:ab,ti OR congo:ab,ti OR 'cote divoire':ab,ti OR
 597 'democratic republic congo':ab,ti OR djibouti:ab,ti OR 'equatorial guinea':ab,ti OR
 598 eritrea:ab,ti OR eswatini:ab,ti OR ethiopia:ab,ti OR gabon:ab,ti OR gambia:ab,ti OR
 599 ghana:ab,ti OR guinea:ab,ti OR 'guinea-bissau':ab,ti OR kenya:ab,ti OR lesotho:ab,ti
 600 OR liberia:ab,ti OR madagascar:ab,ti OR malawi:ab,ti OR mali:ab,ti OR mayotte:ab,ti
 601 OR mozambique:ab,ti OR namibia:ab,ti OR niger:ab,ti OR nigeria:ab,ti OR rwanda:ab,ti
 602 OR sahel:ab,ti OR 'sao tome and principe':ab,ti OR senegal:ab,ti OR 'sierra leone':ab,ti
 603 OR somalia:ab,ti OR 'south africa':ab,ti OR 'south sudan':ab,ti OR sudan:ab,ti OR
 604 tanzania:ab,ti OR togo:ab,ti OR uganda:ab,ti OR zambia:ab,ti OR zimbabwe:ab,ti)

605

606 Query for Web of Science:

607

608 TS=((("acute hiv" OR "early hiv" OR "primary hiv" OR ("window period" AND "human
 609 immunodeficiency virus")) AND ("africa south of the sahara"/exp OR "sub-saharan
 610 africa" OR "subsaharan africa" OR "africa south of the sahara" OR angola OR benin OR
 611 botswana OR"burkina faso" OR burundi OR cameroon OR "cape verde" OR "central
 612 africa" OR "central african republic" OR chad OR comoros OR congo OR "cote divoire"
 613 OR "democratic republic congo" OR djibouti OR "equatorial guinea" OR eritrea OR
 614 eswatini OR ethiopia OR gabon OR gambia OR ghana OR guinea OR "guinea-bissau"
 615 OR kenya OR lesotho OR liberia OR madagascar OR malawi OR mali OR mayotte OR
 616 mozambique OR namibia OR niger OR nigeria OR rwanda OR sahel OR "sao tome and
 617 principe" OR senegal OR "sierra leone" OR somalia OR "south africa" OR "south
 618 sudan" OR sudan OR tanzania OR togo OR uganda OR zambia OR zimbabwe))

619 Table S1: Hyperparameters definition, range and values for each clustering

Parameters	Description	Range	First clustering
PaCMAP n dimensions	Determine the dimensionality of the reduced dimension space that the data will be embedded into	2 - 400	310
PaCMAP n	Controls how PaCMAP balances local	2 - 400	18

neighbours	versus global structure in the data		
hdbscan min cluster size	Represents the minimum number of papers that takes a cluster	2 - 400	30
hdbscan min samples	Represents how conservative the clustering is. A large value increases the amount of points labelled as noise; therefore, clusters will be more separated from each other	2 - 400	7

620

Parameters	Second clustering	Third clustering	Fourth clustering
PaCMAP n dimensions	87	2	62
PaCMAP n neighbours	8	14	12
hdbscan min cluster size	18	12	9
hdbscan min samples	2	8	5

621