

DenseReviewer: A Screening Prioritisation Tool for Systematic Review based on Dense Retrieval

Xinyu Mao¹, Teerapong Leelanupab¹,
Harrison Scells², and Guido Zuccon¹

¹ The University of Queensland

² University of Kassel and hessian.AI

Abstract. Screening is a time-consuming and labour-intensive yet required task for medical systematic reviews, as tens of thousands of studies often need to be screened. Prioritising relevant studies to be screened allows downstream systematic review creation tasks to start earlier and save time. In previous work, we developed a dense retrieval method to prioritise relevant studies with reviewer feedback during the title and abstract screening stage. Our method outperforms previous active learning methods in both effectiveness and efficiency. In this demo, we extend this prior work by creating (1) a web-based screening tool that enables end-users to screen studies exploiting state-of-the-art methods and (2) a Python library that integrates models and feedback mechanisms and allows researchers to develop and demonstrate new active learning methods. We describe the tool’s design and showcase how it can aid screening. The tool is available at <https://densereviewer.ielab.io>. The source code is also open sourced at <https://github.com/ielab/densereviewer>.

Keywords: Systematic Reviews · Dense Retrieval · Relevance Feedback.

1 Introduction and Related Work

Medical systematic reviews (SRs) synthesise evidence from the literature, requiring high recall to avoid missing relevant studies. The screening process is critical to ensure high recall and is a two-stage process: Firstly, the title and abstract of studies are assessed by medical researchers or librarians for relevance, followed by the full text. The former title and abstract (T&A) screening generally involves tens of thousands of studies [1], leading to a high workload and cost. Several tools and products have been developed to reduce this workload, including ASReview [17],¹ Covidence,² DistillerSR,³ and RobotAnalyst [11].⁴ These tools classify studies using classical machine learning. Each study is suggested for inclusion or exclusion or labelled with a confidence score by the model.

¹ <https://asreview.nl/>

² <https://www.covidence.org/>

³ <https://www.distillersr.com/>

⁴ <https://nactem.ac.uk/robotanalyst/>

Table 1: Comparison of key features between DenseReviewer and popular SR tools. ‘Full Screen’ shows the title and abstract of one study at a time, while ‘Ranking List’ presents studies with their titles and abstracts for screening in an order, typically by relevance.

Screening Tool	Core Technology	Interface	Code
ASReview	Machine Learning and Active Learning (model training)	Full Screen	Open Source
Covidence		Ranking List	Proprietary
DistillerSR		Ranking List	Proprietary
RobotAnalyst		Ranking List	Proprietary
DenseReviewer	Dense Retrieval and Relevance Feedback (query vector updating)	Ranking List and Full Screen	Open Source

Prior work has proposed to use active learning (AL) [16] to strategically select studies for manual judgement in order to iteratively train models to more effectively prioritise relevant studies. The use of AL in systematic review automation tools is so far limited [17]. Furthermore, recent studies [19,9] showed that neural models such as BERT have the potential to prioritise studies much more effectively than previous approaches using AL, especially when pre-trained on domain-specific data (e.g., bio-medicine) [9]. However, one downside to these highly effective models, and in fact all AL methods, is that they still require bootstrapping in the form of pre-selected relevant studies. Their computational cost also makes them considerably slower than traditional classification methods.

In this paper, we demonstrate DenseReviewer, a screening tool leveraging dense retrievers and tailored queries (i.e., PICO: patient/population, intervention, comparison, and outcome [13]) for T&A screening. Key features of DenseReviewer are summarised in table 1, and compared with popular SR tools mentioned above. DenseReviewer iteratively updates a PICO query efficiently via the Rocchio’s algorithm for dense retrieval [8] based on the screener’s feedback (i.e., the judgement of each screened study). Our previous work [10] showed that this dense retrieval-based approach is more effective for screening than logistic regression-based and BERT-based active learning workflows, while maintaining efficiency comparable to traditional machine learning-based methods.

2 Overview of DenseReviewer

DenseReviewer offers two modes: ranking mode and focus mode. Figure 1 provides an overview of ranking mode, allowing users to browse and make assessments on studies from the perspective of a paginated, ranked list. Figure 2 provides an overview of focus mode, allowing users to review each study individually and efficiently. In this mode, keyboard controls can be used to quickly move between studies and make assessments.

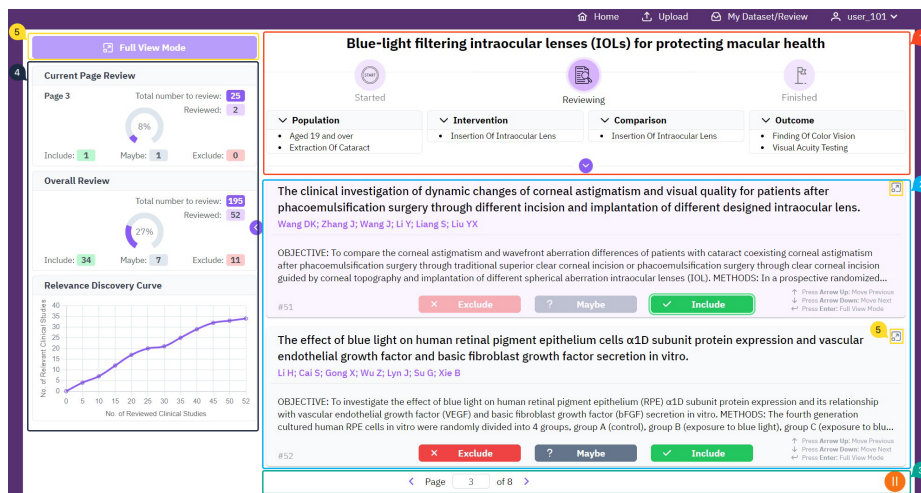


Fig. 1: Ranking mode. ① contains the PICO query. ② lists the studies; users can use keyboard or mouse controls to expand a study to read and judge. Assessed studies are highlighted in purple. ③ contains page controls, with a pause button to save the review’s progress and toggles to stop upon reaching the last page. ④ shows two pie charts that display the ratio of reviewed to unreviewed studies and the distribution of judgements and a line chart that displays the relevance discovery curve, indicating the saturation of relevant studies throughout the screening progresses. ⑤ allows users to enter focus mode (see Figure 2).

Users can upload their corpus, retrieved from PubMed in `nbib` format, and submit a structured PICO query [13]. Uploaded studies must include the pmid (PubMed identifier), title, abstract and a list of authors. A dense retriever ranks uploaded studies with respect to the PICO query. After assessing each page of studies, the remaining unjudged studies are re-ranked based on Rocchio’s algorithm with positive (include) and negative (exclude) feedback [10]. DenseReviewer is open source, and all its components, including the front-end, API back-end, and database, are packaged within Docker containers. The containers can be downloaded from a single git repository and self-hosted with ease. Aside from the GUI-based system for screening, we provide a Python library⁵ to enable experimentation using existing datasets [18,14,5,6,7], and training or loading dense retrievers with customised backbone models.

3 Architecture and Library

DenseReviewer comprises six Docker containers that can all be deployed on a single cloud instance: (i) a web-based front end, (ii) a REST API back end, (iii) a database for storing information such as user activity, uploaded corpora, and ranking studies, (iv) a message broker for managing asynchronous task queues,

⁵ <https://github.com/ielab/dense-screening-feedback>



Fig. 2: Focus mode. ❶ displays the study’s title, list of authors, and PubMed ID. ❷ contains the full abstract. ❸ includes the three assessment options, a ranking index of the current study on the page, and the page number, positioned centrally, to the left, and right, respectively. Buttons ❹ and ❺ allow users to navigate between studies.

<pre>python tevatron_pipe.py \ --collection_split clef19_dta_train \ --model_path biolinkbert \ --q_max_len 128 \ --p_max_len 256 \ --train_n 11 \ --train_epoch 60</pre>	<pre>python dense_query_tar.py \ --collection_split clef19_dta_test \ --model biolinkbert_128_256_11 \ --n_iteration 20 \ --top_k 20 \ --output_path ourput_dir \ --alpha 1.0 \ --beta 0.8 \ --gamma 0.2</pre>
(a) Training	(b) Screening

Fig. 3: Command line usage of DenseReviewer for training and screening.

(v) a service dedicated to tasks such as parsing, encoding, indexing, and initial ranking when new datasets are uploaded to review, and (vi) a service responsible for handling re-ranking. We deployed DenseReviewer on an AWS EC2 accelerated computing instance (g4dn.xlarge: 1 T4 Tensor GPU, 16 GB GPU memory, 4 vCPUs, and 16 GB instance memory). The architecture can scale to more powerful instances, such as the G5 or P4 series with NVIDIA A10G and A100 GPUs, without requiring major modifications.

Figure 3 shows the command line usage for (i) training dense retrievers based on Tevatron (Figure 3a) and (ii) running experiments with the trained dense retrievers and feedback methods (Figure 3b).

4 Conclusion and Future Work

In this demonstration paper, we introduced DenseReviewer, a framework for title and abstract screening for medical systematic reviews with dense retrieval and active learning. We presented a web application, which showcased two interfaces to screen studies. We also presented the underlying Python library which showcased dense retrieval active learning experimentation.

Going forward, we have several new functionalities planned. The first planned functionality is to extend the PICO query to allow users to express exclusion criteria. The second planned functionality is the automatic extraction and highlighting of potential words or sentences associated with (non-)relevant PICO elements. Several studies have shown the potential of large language models (LLMs) in supporting this task [4,12]. The third planned functionality is to automate relevance judgements during screening. We will investigate LLM-assisted screening to further reduce screeners' workload. Several studies have shown the applicability of LLMs at this task [3,2,15].

Acknowledgments

Xinyu Mao is supported by a UQ Earmarked PhD Scholarship and this research is funded by the Australian Research Council Discovery Projects programme ARC DP DP210104043. We extend our gratitude to the engineering team of AI DETA Technologies Co.⁶ (previously named Thaibiogenix Co.), including Kanlayakorn Yeenang, Weeravat Buachoom, Tasanai Srisawat, Warangkhana Sukpartcharoen, and Thuchpun Apivitcholachat for their consultation and support in developing the web application for deploying the `DenseReviewer` core.

References

1. Borah, R., Brown, A.W., Capers, P.L., Kaiser, K.A.: Analysis of the time and workers needed to conduct systematic reviews of medical interventions using data from the prospero registry. *BMJ open* **7**(2) (2017)
2. Bron, M.P., Greijn, B., Coimbra, B.M., van de Schoot, R., Bagheri, A.: Combining large language model classifications and active learning for improved technology-assisted review (2024)
3. Cao, C., Sang, J., Arora, R., Kloosterman, R., Cecere, M., Gorla, J., Saleh, R., Chen, D., Drennan, I., Teja, B., et al.: Prompting is all you need: Llms for systematic review screening. *medRxiv* pp. 2024–06 (2024)
4. Ghosh, M., Mukherjee, S., Ganguly, A., Basuchowdhuri, P., Naskar, S.K., Ganguly, D.: Alpapico: Extraction of pico frames from clinical trial documents using llms. *Methods* **226**, 78–88 (2024)
5. Kanoulas, E., Li, D., Azzopardi, L., Spijker, R.: Clef 2017 technologically assisted reviews in empirical medicine overview. In: *CEUR Workshop Proceedings*. vol. 1866, pp. 1–29 (2017)

⁶ <https://aideta.com/>

6. Kanoulas, E., Li, D., Azzopardi, L., Spijker, R.: Clef 2018 technologically assisted reviews in empirical medicine overview. In: CEUR Workshop Proceedings. vol. 2125 (2018)
7. Kanoulas, E., Li, D., Azzopardi, L., Spijker, R.: Clef 2019 technology assisted reviews in empirical medicine overview. In: CEUR Workshop Proceedings. vol. 2380 (2019)
8. Li, H., Mourad, A., Zhuang, S., Koopman, B., Zuccon, G.: Pseudo relevance feedback with deep language models and dense retrievers: Successes and pitfalls. *ACM Transactions on Information Systems* **41**(3), 1–40 (2023)
9. Mao, X., Koopman, B., Zuccon, G.: A reproducibility study of goldilocks: Just-right tuning of bert for tar. In: European Conference on Information Retrieval. pp. 132–146. Springer (2024)
10. Mao, X., Zhuang, S., Koopman, B., Zuccon, G.: Dense retrieval with continuous explicit feedback for systematic review screening prioritisation. In: Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 2357–2362 (2024)
11. Przybyła, P., Brockmeier, A.J., Kontonatsios, G., Le Pogam, M.A., McNaught, J., von Elm, E., Nolan, K., Ananiadou, S.: Prioritising references for systematic reviews with robotanalyst: a user study. *Research synthesis methods* **9**(3), 470–488 (2018)
12. Reason, T., Langham, J., Gimblett, A.: Automated mass extraction of over 680,000 picos from clinical study abstracts using generative ai: A proof-of-concept study. *Pharmaceutical Medicine* pp. 1–8 (2024)
13. Scells, H., Zuccon, G., Koopman, B., Deacon, A., Azzopardi, L., Geva, S.: Integrating the framing of clinical questions via pico into the retrieval of medical literature for systematic reviews. In: CIKM'17 (2017)
14. Scells, H., Zuccon, G., Koopman, B., Deacon, A., Azzopardi, L., Geva, S.: A test collection for evaluating retrieval of studies for inclusion in systematic reviews. In: Proceedings of the 40th international ACM SIGIR conference on research and development in information retrieval. pp. 1237–1240 (2017)
15. Scherbakov, D., Hubig, N., Jansari, V., Bakumenko, A., Lenert, L.A.: The emergence of large language models (llm) as a tool in literature reviews: an llm automated systematic review. *arXiv preprint arXiv:2409.04600* (2024)
16. Settles, B.: Active learning literature survey (2009)
17. Van De Schoot, R., De Bruin, J., Schram, R., Zahedi, P., De Boer, J., Weijdema, F., Kramer, B., Huijts, M., Hoogerwerf, M., Ferdinands, G., et al.: An open source machine learning framework for efficient and transparent systematic reviews. *Nature machine intelligence* **3**(2), 125–133 (2021)
18. Wang, S., Scells, H., Clark, J., Koopman, B., Zuccon, G.: From little things big things grow: A collection with seed studies for medical systematic review literature search. In: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 3176–3186 (2022)
19. Yang, E., MacAvaney, S., Lewis, D.D., Frieder, O.: Goldilocks: Just-right tuning of bert for technology-assisted review. In: European Conference on Information Retrieval. pp. 502–517. Springer (2022)