Machine learning outperforms manual screening in literature reviews: a case study with rare disease natural history data

Maria Rapoport<sup>1</sup>, Antoinette Randet<sup>1</sup>, Chris Cooper<sup>2</sup> WICKENSTONES RS IN BUILDING VALL <sup>1</sup>Wickenstones Ltd, Oxford, UK; <sup>2</sup>Independent Researcher, London, UK Title and abstract screening for systematic reviews can be time consuming. A variety of machine learning applications are now available to assist researchers. Figure 1. Screening flow with ti/ab screening tools that use active learning Machine learning (ML) tools for the title and abstract (ti/ab) screening phase of systematic reviews Researcher decides to stop screening when they think that all relevant (SRs) include, among others, generative artificial intelligence and active learning-based tools. Researcher ► Active learning tools, such as ASReview,<sup>1</sup> allow reviewers to (Figure 1): records have been found. Stopping labels some Retain control over screening decisions records from the rules help to make that decision Only screen a subset of the ti/ab records by pre-defining a stopping rule that helps determine when set as relevant Researcher or irrelevant. most relevant records likely have been found decides if a However, concerns remain about: record is Set of ML learns relevant. Reproducibility of the screening using ML tools unscreened using the ti/ab records labelled ▶ Recall (i.e. proportion of relevant records identified during screening), especially in complex and Labelled records. varied record sets (i.e. sets containing various study types and populations) records Choice of optimal stopping rule, as multiple options, ranging from simple (e.g. 10% of the set of consecutive excludes) to complex (e.g. Screen-Apply-Find-Evaluate [SAFE]<sup>2</sup> procedure) have been ML orders all unlabelled records from most to least likely relevant and shows the researcher the next most likely relevant record from the set proposed We tested an ML tool and 5 stopping rules for ti/ab screening for an SR of the natural history of a rare genetic disorder to 1111 YYXX assess the reliability and performance on a variable dataset and with a complex research question. SET UP PILOT SCREEN TI/AB SCREEN **Reviewer 2: ASReview (continued)** Tested five stopping rules: **Reviewer 1: Excel** We conducted an SR to Both reviewers assessed 100 describe the natural history of random reports to align their Screened the set in random order and stopped when all records SAFE procedure, which includes multiple criteria and a second a rare genetic disorder. interpretation of the inclusion were reviewed screen of unlabeled records with a more complex model criteria. We searched multiple **Reviewer 2: ASReview** ▶ 50 consecutive irrelevant records The number of expected databases Screened the set in ASReview<sup>1</sup> and stopped when stopping rules ▶ 2.5% and 10% of the set of unscreened records in includes was calculated We deduplicated search were reached consecutively irrelevant records based the inclusion rate in results using Endnote X8. ASReview was set up for screening using the recommended the pilot screening after ▶ 95% of estimated relevant records found (based on inclusion model settings (TF-IDF, naïve Bayes, maximum, dynamic Two reviewers screened the conflicts were resolved rate from pilot screen)4 sample. Reviewer 1 used sampling)2 Microsoft Excel; Reviewer 2 Decisions from the pilot screen were used to train the algorithm used ASReview. initially With ASReview, 51% of relevant records were found after reviewing only 12% of the set. In Excel, 52% needed to be screened to find the same number of relevant records. Figure 2. Recall during ti/ab screening in Excel and ASReview 120 SAFE TI/AB SCREEN SEARCH RESULTS 50 consecutive procedure 10% consecutive Of 1,347 records, Reviewer 1 marked 109 and Reviewer 2 ► 2,521 reports retrieved second irrelevant/SAFE 100 marked 90 records as relevant (Figure 2). irrelevant 2.5% procedure initial screen ▶ 1,074 reports deduplicated

► 1.447 reports included in the title and abstract screen

#### PILOT SCREEN

- ► Substantial 93% (κ=0.68) agreement betweer reviewers
- 12% inclusion rate estimated, resulting in 161 expected relevant reports remaining in the unscreened set
- The agreement between reviewers was lower than in the pilot screen (91%; κ=0.37 [fair agreement]).
- After conflict resolution, 97 reports in addition to the pilot screen were included (overall inclusion rate: 7.5%).
  - Screening in ASReview found relevant reports notably faster than screening in Excel
  - In ASReview, most (51%) relevant reports were found after screening only 12% of the sample
  - ▶ In Excel, 52% of the sample needed to be screened to achieve the same recall
- The stopping rule of finding 95% of estimated relevant records was not reached



RWD63

-ASReview screen with recommended settings -ASReview SAFE procedure second screen Screening in random order in Excel

Note: Stopping rules apply to ASReview screen only

ML-supported screening can substantially decrease reviewer workload while maintaining high recall, even when applied to complex research topics.

### RECOMMENDATIONS FOR TI/AB SCREENING USING ACTIVE LEARNING ML

### Tailor the screening experience to your preferences

Screening in a browser-based application (e.g. ASReview) allows the use of browser add-ons (e.g. multi highlight add on in Google Chrome) that can increase screener comfort and efficiency

V&

- Ensure that the screening criteria are clear and reliably applicable before screening
- ► Because every reviewer decision affects which reports are shown to the reviewer next, a thorough pilot screening is imperative for reliable screening results.

### Choose a stopping rule that fits your project

Stopping rules balance reviewer workload and recall and range from simple to complex. Different rules will be needed to suit differing projects and objectives.

### Follow the evolving best practice

ML-supported SR methods have become more widely accepted in the last year, with, for example, the National Institute for Care Excellence (NICE) issuing guidance for company submissions.

## CONCLUSION

- ► Use of ASReview can greatly limit screening time during the ti/ab screening stage, as 51% of relevant records were found after screening 12% of the sample (vs screening 52% of the sample to find the same number of relevant records in Excel).
- ► Despite the variability of the set, ASReview reliable presented relevant records early in the screen.

i Use of AI in Evidence Generation: NICE position statement, 2024. Accessed n 29.10.2024 via https://www.nice.org.uk/about/what-we-do/our-research-ork/use-of-ai-in-evidence-generation--nice-position-statement

ASReview (or other ML applications) may supplement or replace manual ti/ab screening in future.

### REFERENCES

# an de Schoot *et al., 2021. 10.1038/s*42256-020-00287-7. erdinands *et al., 2023. 10.1186/s*13643-023-02257-7. an Haastrecht *et al., 2021. 10.3389/fma.2021.685591.* soetje, van de Schoot, 2024. 10.1186/s13643-024-02502-7

**ACKNOWLEDGEMENTS** 

The SR project was supported by Mereo BioPharma, UK

Presented at: ISPOR EU; Barcelona, Spain; 17-20 November 2024