



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer Screening for

in vitro

systematic reviews: a comparison of screening methods and training of a machine learning classifier

Citation for published version:

Wilson, E, Cruz, F, Maclean, D, Ghanawi, J, Mccann, SK, Brennan, PM, Liao, J, Sena, ES & Macleod, MR 2023, 'Screening for in vitro systematic reviews: a comparison of screening methods and training of a machine learning classifier', *Clinical science*. <https://doi.org/10.1042/CS20220594>

Digital Object Identifier (DOI):

[10.1042/CS20220594](https://doi.org/10.1042/CS20220594)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Clinical science

Publisher Rights Statement:

This is the author's peer-reviewed manuscript as accepted for publication.

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 Screening for *in vitro* systematic reviews: a comparison of screening 2 methods and training of a machine learning classifier

3 Emma Wilson¹, Florenz Cruz², Duncan Maclean³, Joly Ghanawi⁴, Sarah K McCann², Paul M Brennan¹,
4 Jing Liao¹, Emily S Sena¹, Malcolm Macleod¹

5 (1) Centre for Clinical Brain Sciences, The University of Edinburgh, Edinburgh, UK

6 (2) Berlin Institute of Health at Charité-Universitätsmedizin Berlin, QUEST Center, Berlin, Germany

7 (3) University of Edinburgh Medical School, The University of Edinburgh, Edinburgh, UK

8 (4) Independent Researcher, UK

9

10 Corresponding author

11 Emma Wilson, emma.wilson@ed.ac.uk

12 ORCID iDs

- 13 • Emma Wilson, <http://orcid.org/0000-0002-8100-7508>
- 14 • Florenz Cruz
- 15 • Duncan Maclean
- 16 • Joly Ghanawi, <https://orcid.org/0000-0002-7945-2055>
- 17 • Sarah K McCann, <https://orcid.org/0000-0003-4737-2349>
- 18 • Paul Brennan, <https://orcid.org/0000-0002-7347-830X>
- 19 • Jing Liao, <https://orcid.org/0000-0002-9591-8070>
- 20 • Emily S Sena, <http://orcid.org/0000-0002-3282-8502>
- 21 • Malcolm Macleod <https://orcid.org/0000-0001-9187-9839>

22 Funding

23 SKM and FC were supported by the German Federal Ministry of Education and Research (BMBF)
24 under the Confirmatory Preclinical Studies and Systematic Reviews Initiative, grant number:
25 01KC190.

26 Ethics Statement

27 This research did not require ethical approval.

28 Conflicts of interest

29 The authors declare no conflicts of interest.

30

31 Keywords

32 Meta-research; systematic review; automation; machine learning; *in vitro* models

33

34

35

36 **CRedit**

37 MM conceptualised the study. MM and EW contributed to the methodology. EW and DM curated
38 the data used in this study. EW, FC, DM, JG, SKM, PMB, ESS, and MM contributed to the
39 investigation. EW and MM performed the formal analysis, and EW created the data visualisations. JL
40 contributed the resources and software used in this study. EW led the data-to-day project
41 administration, and MM and ESS provided supervision. EW wrote the original draft, and EW, MM,
42 SKM, PMB and JL contributed to reviewing and editing subsequent drafts.

43 **Abstract**

44 Objective: Existing strategies to identify relevant studies for systematic review may not perform
45 equally well across research domains. We compare four approaches based on either human or
46 automated screening of either title and abstract or full text; and report the training of a machine
47 learning algorithm to identify *in vitro* studies from bibliographic records. Methods: We used a
48 systematic review of oxygen-glucose deprivation (OGD) in PC-12 cells to compare approaches. For
49 human screening, two reviewers independently screened studies based on title and abstract or full
50 text, with disagreements reconciled by a third. For automated screening, we applied text mining to
51 either title and abstract or full text. We trained a machine learning algorithm with decisions from
52 2,000 randomly selected PubMed Central records enriched with a dataset of known *in vitro* studies.
53 Results: Full text approaches performed best, with human (sensitivity 0.990, specificity 1.000,
54 precision 0.994) outperforming text mining (sensitivity 0.972, specificity 0.980, precision 0.764). For
55 title and abstract, text mining (sensitivity 0.890, specificity 0.995, precision 0.922) outperformed
56 human screening (sensitivity 0.862, specificity 0.998, precision 0.975). At our target sensitivity of
57 95% the algorithm performed with specificity of 0.850 and precision of 0.700. Conclusion: In this *in*
58 *vitro* systematic review, human screening based on title and abstract erroneously excluded 14% of
59 relevant studies, perhaps because title and abstract provide an incomplete description of methods
60 used. Our algorithm might be used as a first selection phase in *in vitro* systematic reviews to limit the
61 extent of full text screening required.

62

63 Words: 248

64

65 **Clinical Perspective**

66 Systematic reviews of *in vivo* animal experimental data have made important contributions to the
67 evidence-based translation of findings from the laboratory to human clinical trials, and has informed
68 clinical trial design. Equally, *in vitro* research makes key contributions to the development of new
69 treatments and therapies. Recently, we have seen an increase in the number of systematic reviews
70 investigating *in vitro* research relevant to human health. However, the nature of the *in vitro*
71 literature may be different to *in vivo*, and it is important to determine where systematic review
72 methodologies as currently used can be simple applied or may require adaptation. Here we show
73 that title and abstract screening has low sensitivity to identify relevant *in vitro* publications, and we
74 make recommendations to optimise search and screening strategies for *in vitro* systematic reviews.

75

76 Words: 133

77	Abbreviations	
78	API	Application Programming Interface
79	ASySD	Automated Systematic Search Deduplicator
80	AUC	Area Under the Curve
81	CAMARADES	Collaborative Approach to Meta-Analysis and Review of Animal Data from
82		Experimental Studies
83	EPPI-Centre	Evidence for Policy and Practice Information and Co-ordinating Centre
84	FPR	False Positive Rate
85	LDH	Lactate Dehydrogenase
86	MTT	3-(4,5-dimethylthiazol-2-yl)-2,5-diphenyltetrazolium bromide
87	NCBI	National Center for Biotechnology Information
88	NPQIP	Nature Publication Quality Improvement Project
89	OGD	Oxygen-Glucose Deprivation
90	PC-12	Pheochromocytoma-12
91	PICO	Population, Intervention, Comparator/Control, Outcome
92	PMC	PubMed Central
93	PRISMA	Preferred Reporting Items for Systematic reviews and Meta-Analyses
94	RCT	Randomised Controlled Trial
95	RegEx	Regular Expression
96	ROC	Receiver Operating Characteristic
97	SGD	Stochastic Gradient Descent
98	SVM	Support Vector Machine
99	SYRCLE	SYstematic Review Center for Laboratory animal Experimentation
100	TiAb	Title and Abstract

101 Introduction

102 Experiments conducted *in vitro* play an invaluable role in the research pipeline. *In vitro* models,
103 including 3D organoids, have recently attracted attention as methods which might reduce and
104 eventually replace the use of animals in research [1]. However, challenges in translating findings
105 from *in vitro* research to the clinic may hinder efforts to replace animal research. Poor reporting of
106 measures to reduce the risk of bias in *in vitro* studies may contribute to this translational challenge
107 [2], and research which systematically identifies such issues [3] may lead to improvements in the
108 design, conduct and reporting of *in vitro* research, and, thereby, their adoption as alternatives to
109 animal research.

110 Systematic review is a research method used to summarise and critically appraise all available
111 published evidence related to a pre-defined research question [4]. The use of systematic review to
112 evaluate evidence from clinical trials has led to significant improvements in clinical trial design,
113 conduct and reporting [5]. The application of systematic review methodologies to *in vivo* animal
114 studies has, similarly, identified opportunities for improvement [6,7]. More recently, reviews of *in*
115 *vitro* data have suggested similar problems may be prevalent there [2,3,8].

116 Tools and guidance developed by Cochrane have contributed substantially to improving the
117 methodological quality of clinical systematic reviews [9–12]. Similar guidance has been articulated
118 for systematic reviews of animal studies including a protocol template [13], the CAMARADES
119 reporting quality checklist [14], the SYRCLE risk of bias checklist [15], and the development of a
120 Preferred Reporting Items for Systematic reviews and Meta-Analyses (PRISMA) extension for such
121 reviews is ongoing [16]. These adapted guidelines reflect important differences between clinical and
122 animal studies, including study size (many human patients per study versus few laboratory animals
123 per study) and heterogeneity between studies (lower in human than in animal studies).

124 It is possible that the methods used to plan and conduct *in vitro* systematic reviews must be further
125 adapted. One key feature is the process of searching and screening for relevant publications. In a
126 typical systematic review of animal data, search results from multiple databases are combined,
127 duplicate citations are removed, and titles and abstracts are screened for relevance. General
128 guidance is that the screeners should, if anything, be over-inclusive at this stage (i.e. perform with
129 high sensitivity, perhaps at a cost in specificity: [17,18]). This stage is followed by full text screening
130 to determine eligibility.

131 In a pilot systematic review of *in vitro* data conducted in 2019 (unpublished) we found an
132 unexpectedly low yield of included studies and hypothesised that title and abstract ([TiAb])
133 screening may not be sufficiently sensitive. Where animal and *in vitro* experiments were reported in
134 the same publication, we were concerned that a full summary of *in vitro* methods and results may
135 not always be included in the abstract. This would lead to studies being incorrectly excluded at the
136 [TiAb] screening phase. Further, as systematic searches are often conducted on [TiAb] text –
137 especially where relevant field tags such as MeSH terms may not be available – relevant *in vitro*
138 studies may not even be identified in literature searches specifically designed to identify *in vitro*-
139 related terms. These concerns are consistent with a recent finding that, in studies where multiple
140 outcomes were investigated, negative findings were less likely to be included in the abstract text and
141 therefore less likely to be included in systematic reviews [19]. In our view, for the purposes of most
142 systematic reviews, screening approaches should perform with a sensitivity of at least 95% — that is,
143 they should wrongly exclude fewer than 1 in 20 relevant studies.

144 One approach to this problem would be to conduct broader systematic searches to capture any
145 article that might contain an *in vitro* experiment and to screen studies for relevance on the basis of
146 the full text PDF article. However, this would be significantly burdensome, in a context where a
147 major limitation of current methodologies is the time and effort required to complete a systematic
148 review. This is especially true in preclinical systematic reviews, which tend to screen and include a
149 higher number of publications compared with clinical reviews.

150 Recently, automation tools have been developed to accelerate parts of the systematic review
151 process including screening [20–22], PICO extraction [23,24] and risk of bias assessment [25–27].
152 These tools allow researchers to conduct reviews more quickly and without requiring as much
153 human effort; we wondered if automation tools might address the issue of incomplete [TiAb]
154 descriptions.

155

156 **Aims**

157 Here we compare the performance of four different screening methods — (i) human screening
158 based on [TiAb] only; (ii) human screening based on full text; (iii) automated screening based on
159 [TiAb] only, and (iv) automated screening based on full text — in an exemplar systematic review of
160 ischaemic injury induced by oxygen-glucose deprivation in PC-12 cells. Then, we train a machine
161 learning algorithm, developed specifically for systematic review screening, to identify studies which
162 report the results of *in vitro* experiments.

163 **Methods**

164 **Method 1: Comparison of screening methods in an example systematic review**

165 The study protocol for the comparison of screening methods is available at <https://osf.io/cq48b/>.
166 Methods of analyses were not described in the protocol, and deviations from the protocol are
167 described in Appendix 1.

168 **Search strategy**

169 We conducted a systematic search of PubMed (accessed via NCBI) and Embase (accessed via Ovid)
170 on 16th March 2020. Full search terms are given in Appendix 2(i) and included a series of terms to
171 identify the experimental approach (e.g., “oxygen glucose deprivation”), the condition modelled
172 (e.g., “brain ischaemia”), and the experimental materials (e.g. “PC-12”). An error in implementing
173 the search terms led to our combining the first two of these phrases with “OR” rather than “AND”
174 (the errors are underlined in Appendix 2(i)) resulting in the retrieval of many more studies than had
175 the search been implemented as intended. We did not notice this error until all studies had been
176 screened, and we provide a primary analysis of the search as implemented, with a secondary
177 analysis of the search as was intended.

178 We imported each search result into EndNote X8, created a single XML file of all search results, and
179 removed duplicate citations using the Automated Systematic Search Deduplicator (ASySD) tool [28].
180 This performs automatic deduplication with limited human input and was designed specifically for
181 use in preclinical (but not necessarily *in vitro*) systematic review projects. We imported our
182 deduplicated search results to EndNote X8 and retrieved full text PDFs using EndNote’s in-built “find
183 full text” feature, then converted PDFs to plain text files using the PDF to text function from
184 XpdfReader (<https://www.xpdfreader.com/>).

185 **Eligibility criteria for analysis**

186 We included records which had both an English-language abstract and an English-language full text.
187 We excluded conference abstracts, records with no abstract, records with no English-language full
188 text, records where a full text was not retrieved by EndNote X8, and records which did not have a
189 machine-readable full text.

190 **Systematic review Inclusion and exclusion criteria**

191 The screening task was to identify controlled experiments exposing PC-12 cells to oxygen-glucose
192 deprivation (OGD) *in vitro* and reporting effects on cell death or survival (MTT assay, LDH assay, or
193 cell counting), whether investigating the effects of OGD or the impact of interventions (e.g.
194 pharmacological, genetic) intended to modulate the effects of OGD. There was no limitation by
195 publication date.

196 **Human screening**

197 For human screening, we used the Systematic Review Facility (SyRF) web application [29] to screen
198 studies against our inclusion criteria. A pool of 6 reviewers were allocated records in random order,
199 and each record was screened by at least two reviewers. Where there was disagreement, the record
200 was automatically presented to a third reviewer for arbitration. All decisions were taken blinded to
201 the decision(s) of other reviewers, and whether the task was initial screening (i.e. “reviewer 1” or
202 “2”) or reconciliation of conflicting opinions (“reviewer 3”). Reviewers first screened each study
203 based on [TiAb], and then, in the same session, were asked to screen the study again based on the

204 full text PDF. Therefore, each publication was screened twice, first on the basis of [TiAb] and then on
205 the basis of the full text.

206 **Automated screening using regular expressions**

207 For automated screening we used the R programming language and Regular Expressions (RegEx). A
208 regular expression is a sequence of characters which can be used to search for and match certain
209 patterns within text [30]. We developed a RegEx to identify relevant publications by matching terms
210 such as “oxygen-glucose deprivation”, “OGD”, “oxygen and glucose deprivation”, or “deprived of
211 oxygen and glucose”. The full RegEx is given in Appendix 3. We then used the AutoAnnotation R
212 package [31] to count the number of occurrences of regular expressions matches in the [TiAb] or full
213 text. One match meant that some form of the term oxygen-glucose deprivation was mentioned only
214 once within the text, two matches meant that some form of the term was mentioned twice, *et*
215 *cetera*.

216 **The gold standard dataset**

217 To create a dataset with the highest proportion of true decisions, we reasoned that reconciled
218 human full text screening decisions were likely to be most complete. Where there was disagreement
219 between the human full text decision and another decision, then that study was evaluated by a
220 senior experienced reviewer, and where they were not in agreement with the reconciled human full
221 text screening decision, their re-evaluated decision was used as the gold standard.

222 **Evaluation of screening performance**

223 We assessed the performance of each approach by calculating the sensitivity, specificity and
224 precision, characterising the "purity in retrieval performance" [32], (number of true positive
225 decisions divided by the number of positive decisions) using the Caret R package [33].

226 **Assessing best performance**

227 Perfect performance is achieved when sensitivity and specificity are both 100%. 100% sensitivity is
228 achieved when all relevant publications are included during screening, and 100% specificity is
229 achieved when all non-relevant publications are excluded during screening. We calculated the
230 Euclidian distance (d) between the performance achieved and this optimum performance as

231

$$d = \sqrt{(1 - \text{Sensitivity})^2 + (1 - \text{Specificity})^2}$$

232

233 and the screening method with the smallest value of d was considered optimal. For the automation
234 approaches, we used the same approach to calculate the point on the Receiver Operating
235 Characteristic (ROC) curve closest to peak performance. As a second measure of performance, we
236 used the area under the ROC curve.

237 **Method 2: Developing a trained machine learning classifier for *in vitro* systematic review**
238 **screening**

239 The study protocol for the development of a machine learning classifier is available at
240 <https://osf.io/bjisp2/>. Deviations from the protocol methods are described in Appendix 1.

241 **Definition of *in vitro* research**

242 For the purposes of developing this screening tool, we define *in vitro* research as involving the
243 manipulation of biomolecules (including enzymes, genes, genomes), cells, tissues, or organs in a
244 controlled, artificial environment such as a petri dish, well or test tube.

245 Our definition includes samples which may also be described as *ex vivo* (tissues originating from
246 experimental animals) if the experimental intervention under investigation was applied to the
247 specimen after derivation rather than being applied *in vivo* pre mortem or before tissue collection.

248 **Generation of a screened dataset**

249 Using the PMC API we downloaded 2,000 randomly sampled records from PubMed Central (PMC) on
250 the 19th of December 2019 [34]. We used no search terms, filters or restrictions to generate this
251 sample.

252 We uploaded all 2,000 PMC records to the SyRF web application for full text screening based on our
253 definition of *in vitro* research, given above. Each study was screened by two independent reviewers
254 and disagreements were reconciled by a third independent reviewer.

255 We then supplemented our 2,000 screened records with 453 known *in vitro* studies previously
256 screened as part of the Nature Publication Quality Improvement Project (NPQIP) study [2]. The
257 merged dataset included a unique identifier for each study, the [TiAb] text, and a binary flag
258 indicating the include or exclude screening decision.

259 **Training the machine learning algorithm**

260 We used the binary screening decisions (“include” or “exclude”) from our merged dataset to train a
261 machine learning algorithm hosted by our collaborators at The Evidence for Policy and Practice
262 Information and Co-ordinating Centre (EPPI-Centre), University College London. The algorithm uses a
263 tri-gram ‘bag-of-words’ model for feature selection and implements a linear support vector machine
264 (SVM) with stochastic gradient descent (SGD), as described in Approach 1 used by Bannach-Brown et
265 al. [21]. The algorithm associates the training set screening decisions with features it identifies in the
266 relevant [TiAb] text, and uses these features to predict the inclusion or exclusion status for new
267 unseen studies.

268 The dataset was randomly split into a training set (80%), validation set (20%) to ensure the algorithm
269 performed optimally.

270 **Error correction and retraining classifier**

271 After algorithm training, we performed a round of error correction as described by Bannach-Brown
272 et al. [21]. We identified the 100 studies with the largest discrepancy between human screening and
273 algorithm score, and had humans rescreen these studies to identify if there had been a human error
274 during screening. We then retrained the machine learning algorithm using the set of 2453 screened
275 records thus corrected.

276 Results

277 Performance of different screening methods for case study: *in vitro* OGD systematic 278 review

279 Search results

280 Figure 1 shows the PRISMA flow diagram. Our systematic search as implemented retrieved a total of
281 9,952 records (4,219 from NCBI PubMed and 5,733 from Ovid Embase). Following deduplication, we
282 identified 6,380 unique records.

283 We were able to retrieve full text PDFs for 5,362 (84%) of the unique records identified from our
284 search. From this, we included a total of 5,172 records in our analysis. We excluded 119 records
285 which were conference abstracts, 42 records where the PDF was not machine-readable, and 29
286 records which had no abstract.

287 Performance of different screening methods

288 Human reviewers identified 282 of 5172 records for inclusion based on [TiAb], and 318 of 5172 when
289 screening against full text. The number of RegEx matches was between 0 and 15 for [TiAb], and
290 between 0 and 281 for full text (Figure 2). We then calculated the sensitivity and specificity at each
291 RegEx threshold (that is, including studies based on N RegEx matches, with $N = 1-281$) and set
292 thresholds for inclusion of 1 match for [TiAb] screening and 2 matches for full text screening (Figure
293 3). Finally, we re-examined those records where there was a discrepancy between human full text
294 screening and one of the other screening approaches. This focussed review identified 3 records
295 which had been omitted by human full text screening but identified by the full text RegEx, and 2
296 records included in error by human full text screening. This gave 319 included studies (6.2% of 5172).

297 Compared to this gold standard, human [TiAb] screening correctly identified 275 of 319 studies, and
298 wrongly included an additional 7 studies ($d = 0.138$). Human full text screening correctly identified
299 316 of 319 studies, wrongly including 2 studies ($d = 0.009$). RegEx of [TiAb] correctly identified 284 of
300 319 studies and wrongly included 24 studies ($d = 0.110$), and RegEx of full text (with an optimal
301 threshold of 2 matches) correctly identified 310 of 319 studies, wrongly including 96 ($d = 0.034$)
302 (Table 1). Area under the curve (AUC) was 0.944 for RegEx applied to [TiAb] and 0.986 for RegEx of
303 full text.

304 1060 citations were excluded from the full text analysis because we were unable to retrieve (1018)
305 or to process (42) the full text. Within this additional corpus, human [TiAb] and RegEx [TiAb]
306 screening respectively identified 57 and 66 additional studies which appeared relevant. Without
307 access to the full text, we cannot determine how many of these might be false positives; and given
308 the sensitivity of these approaches in the main cohort of studies it is likely that further relevant
309 studies will have been excluded.

Screening Method	Number of True Positives	Number of True Negatives	Number of False Positives	Number of False Negatives	Sensitivity	Specificity	Precision	Euclidian distance (<i>d</i>)	AUC
Human Title/Abstract	275	4846	7	44	0.862	0.998	0.975	0.138	0.930
Human Full text	316	4851	2	3	0.990	1.000	0.994	0.009	0.995
RegEx Title/Abstract	284	4829	24	35	0.890	0.995	0.922	0.110	0.944
RegEx Full text	310	4757	96	9	0.972	0.980	0.764	0.034	0.986

310

311 **Table 1: Performance of different screening methods.** A total of 5,172 records were screened using each method. For sensitivity, specificity, and precision,
 312 the optimal performance value is 1. For RegEx title/abstract, the optimal threshold shown is 1 match. For RegEx full text, the optimal threshold shown is 2
 313 matches. A lower Euclidian distance (*d*) indicates better performance.

314 **Analysis of the ‘intended’ search strategy**

315 The error in implementing our search strategy had a profoundly beneficial effect on our ability to
 316 detect relevant articles. On 5th May 2022 we searched NCBI PubMed and Ovid Embase using our
 317 intended search strategy (Appendix 2(ii)), limited by date of record creation to 16th March 2020 (to
 318 align with the initial search), and retrieved 910 unique records (438 from NCBI PubMed and 700
 319 from Ovid Embase, compared with 4,219 and 5,733 respectively in the ‘incorrectly’ implemented
 320 search). Remarkably, only 133 (or 42%) out of the 319 studies we identified using our ‘incorrect’
 321 search were identified by our planned search strategy. If we had used this approach, and if
 322 subsequent human [TiAb] had been conducted, the performance of human [TiAb] screening would
 323 have been overinflated, giving an apparent sensitivity of 0.925 and specificity of 0.999.

324

325 **Training a machine learning classifier for *in vitro* systematic review screening**

326 **Dataset of screened studies**

327 Of the 2,000 articles randomly selected from PMC, after full text screening we judged 296 to
 328 describe *in vitro* research. Combining these with 453 *in vitro* studies from NPQIP, gave a complete
 329 dataset with 749 included studies and 1704 excluded studies (total $N = 2,453$). We randomly divided
 330 these into training ($N = 1962$) and validation ($N = 491$) sets.

331 **Machine learning performance**

332 We trained the machine learning algorithm on title and abstract [TiAb] in the training set, and then
 333 applied the algorithm to the validation set, attributing each citation with a decimal score between 0
 334 to 1, where higher scores reflect a stronger machine prediction of inclusion. We then established a
 335 threshold such that 95% of relevant studies in the validation set would be retrieved (i.e. sensitivity =
 336 0.950 or higher). A machine score threshold for inclusion of 0.25 (Figure 4) gave specificity of 0.824
 337 at sensitivity of 0.951, and precision of 0.692 (Table 2). We then checked human decisions for the
 338 100 citations with greatest mismatch between human decisions and machine predictions. 35
 339 citations had the human decision reversed, with 31 citations included by human decision now
 340 excluded, and 4 citations excluded by human decision now included. Retraining on this revised
 341 corpus gave specificity of 0.850 (increase of 0.026) at sensitivity of 0.954, and precision of 0.700
 342 (increase of 0.008) (Table 2), with a machine score threshold for inclusion of 0.29 (Figure 4).

	Sensitivity	Specificity	Precision
Initial	0.951	0.824	0.692
Corrected	0.954	0.850	0.700

343

344 **Table 2: Performance of the trained machine learning algorithm before and after error correction.**

345 Discussion

346 Screening in *in vitro* systematic reviews

347 In typical biomedical systematic reviews, a systematic search of [TiAb] text retrieves potentially
348 relevant articles, which are then screened by two independent reviewers, and any disagreements
349 reconciled by a third reviewer. The broader the terms of the systematic search the higher will be the
350 sensitivity, but because of the inevitably high total number of citations returned, this will come at
351 the cost of an increased burden of human screening. Here we show that in a systematic review of
352 the effects of oxygen-glucose deprivation in PC-12 cells, human screening of [TiAb] was the least
353 sensitive (0.862) of four approaches tested and would have wrongly excluded around one in every 7
354 relevant manuscripts. Human full text screening performs with a sensitivity of 0.990, wrongly
355 excluding only 1 in 1000 manuscripts. However, because of the time involved this is not a feasible
356 approach for most systematic reviews.

357 While we did not formally compare the time taken by human reviewers and the RegEx algorithms,
358 there is a substantial reduction in time taken, even accounting for the requirement to develop the
359 regular expressions and convert PDF to text. Dual human screening of 5000 [TiAb], even at 30
360 seconds per record, would take over 80 hours, and full text screening around 800 hours; compared
361 with less than one day for the RegEx approach.

362 The RegEx approach achieved higher sensitivity than human screening when applied to [TiAb] text.
363 For full text, sensitivity was slightly lower (0.972) than human screening (0.990). For both RegEx
364 approaches, specificity was lower than human screening ([TiAb]: 0.995 versus 0.998: full text 0.980
365 versus 1.000). For contrast with human [TiAb] screening, RegEx full text screening identifies an extra
366 35 studies (13%) which should be included, at a cost of increasing the number included in error from
367 7 to 96. This could therefore serve as a useful first step before human full text screening, which
368 could take place at the data extraction stage. However, the usefulness of RegEx full text screening
369 will be heavily dependent on the quality of that RegEx, and we strongly advise researchers carefully
370 to consider synonyms, alternate spellings and different combinations of target words or phrases.

371 The benefits of this approach were highlighted, inadvertently, by our mis-formed search strategy.
372 Our intended search would only have returned 42% of the relevant articles identified in the search
373 as implemented, for a maximum sensitivity of 0.42 if subsequent citation screening performed
374 perfectly. While the work of human full text screening these 910 citations would be less than that
375 required for the 5,172 citations included by our broader search, combining that broader search with
376 RegEx applied to full text would achieve sensitivity of 0.972 while requiring human full text review of
377 406 of 5,172 citations.

378

379 Automation in *in vitro* systematic reviews

380 In the first stage we applied automated full text screening to the results of a search strategy which
381 largely interrogates title and abstract. It is therefore likely that additional relevant publications will
382 have been omitted from those search returns, for the same reason as they were not detected by our
383 [TiAb] RegEx. This is confirmed by the very poor performance of what we had considered to be a
384 well-constructed search strategy.

385 While conceptually attractive, applying the full text RegEx approach to all of NCBI PubMed is
386 currently impractical, requiring access to the full text of over 30 million scientific publications. We

387 therefore explored an intermediate approach, where we trained a machine learning algorithm to
388 detect reports of *in vitro* research, that these might then be interrogated by the full text RegEx. In a
389 random sample of PubMed Central records, 14.8% included reports of *in vitro* research (based on
390 human full text screening), and the *in vitro* algorithm, applied to Title and Abstract only, performs
391 with sensitivity of 0.954. However, across a corpus of 30m publications, the specificity of 0.85
392 suggests that of 8.1m publications labelled as reporting *in vitro* research, 3.8m would have been
393 wrongly included, and 200,000 would have been excluded in error.

394 The performance of the full text RegEx in unselected reports of *in vitro* research is not known, but
395 we estimate a prevalence for inclusion of around 0.01% (~400 from ~ 4 million). Estimating
396 sensitivity and specificity in this context would require full text screening of several hundreds of
397 thousands of articles and is not currently practicable. However, performance of this approach
398 against the “gold standard” performance identified here may be feasible. We think that some
399 combination of broad but “conventional” search strategies, combined with algorithmic identification
400 of the *in vitro* literature and RegEx interrogation of selected full text articles, will prove an effective
401 approach.

402

403 **Limitations**

404 Due to lack of full text availability, it was not possible for us to generate a gold standard dataset of
405 all the studies which should be included in the complete corpus of 6,232 studies (5,172 included in
406 the main analysis + 1,060 additional studies). Examining [TiAb] of these additional studies identified
407 an additional 66 potentially relevant studies, but we were not able to confirm this because we were
408 unable to access the full texts. Given a sensitivity for the [TiAb] RegEx of 0.890 as an estimate
409 suggests an additional 10 studies not included by the TiAb RegEx. Taken together, we estimate the
410 total number of relevant studies in the corpus of 6,232 to be 76 more than we have identified,
411 suggesting that there are around 395 relevant studies in that corpus.

412 We can therefore provide rough estimates of the overall sensitivity of various approaches; [TiAb]
413 approaches can be applied to all 6232 and we predict would have identified 332 of the estimated
414 total of 395 studies (sensitivity = 0.841). RegEx [TiAb] would identify 350 (sensitivity = 0.886).
415 Because full text approaches can only be used where we have access to full text, the sensitivity falls
416 from 316 of 319 to 316 of 395 (human, sensitivity = 0.800) and from 310 of 319 to 310 of 395 (RegEx,
417 sensitivity = 0.785) respectively. Our preferred approach is therefore to use full text RegEx where full
418 text is available, supplemented by [TiAb] RegEx when only abstracts are available. In the example
419 provided, this approach identifies 376 studies (310 from RegEx of full text and 66 from RegEx of
420 [TiAb] when only [TiAb] available). With an estimated 395 relevant studies this represents a
421 sensitivity of 0.952.

422 One limitation of the RegEx based approach is that – unlike human screening – it cannot be used
423 where files are not machine readable or where no abstract is provided.

424 A limitation of the machine learning algorithm for detecting *in vitro* research is that it was trained on
425 only English-language [TiAb]s, and so performance in texts in other languages is not known.
426 Excluding non-English language texts may introduce bias and reduce the generalisability of
427 systematic reviews; although in clinical systematic reviews this has been found to have little or no
428 impact on the conclusions of the review [35], we do not yet know the extent or the impact of this
429 potential bias in reviews of *in vitro* experimental data. The algorithm may also perform poorly in

430 contexts where cell preparations are used as therapies in human studies, for instance CAR-T cells in
431 cancer or stem cell transplantation in neurodegenerative diseases.

432

433 **Conclusion**

434 Firstly, we show that in an *in vitro* systematic review, human screening based on title and abstract
435 erroneously excluded 14% of relevant studies. This may be due to an incomplete description in the
436 abstract of all experiments described in a publication, and this may be more likely in the pre-clinical
437 literature, where several experiments are often presented in a single publication. We then describe
438 a machine learning algorithm which detects publications reporting *in vitro* research with high
439 sensitivity. We propose this tool may be used as a first selection phase in *in vitro* systematic reviews
440 to limit the extent of full text screening which our first finding suggests is necessary.

441 **Acknowledgements**

442 We are grateful to Professor James Thomas (EPPI-Centre) for providing an API to the machine
443 learning algorithm used in this study, and to Dr Alexandra Bannach-Brown (Berlin Institute of Health,
444 QUEST Center, Charité Universitätsmedizin Berlin, Germany) for providing feedback on the initial
445 draft of this manuscript.

446

447 **Data Availability Statement**

448 Data and code used in the analysis are available on GitHub ([https://github.com/emma-wilson/in-](https://github.com/emma-wilson/in-vitro-screening)
449 [vitro-screening](https://github.com/emma-wilson/in-vitro-screening)) and are shared under a Creative Commons Attribution 4.0 International License. We
450 do not have permission to share the API key required to run the machine learning, however, further
451 information about access is available at: Thomas J, Brunton J, Graziosi S (2010) EPPI-Reviewer 4.0:
452 software for research synthesis. EPPI-Centre Software. London: Social Science Research Unit,
453 Institute of Education.

References

- 455 [1] van Berlo D, Nguyen VVT, Gkouzioti V, Leineweber K, Verhaar MC, van Balkom BWM. Stem
456 cells, organoids, and organ-on-a-chip models for personalized in vitro drug testing. *Curr Opin*
457 *Toxicol* 2021;28:7–14. <https://doi.org/10.1016/j.cotox.2021.08.006>.
- 458 [2] The NPQIP Collaborative group. Did a change in Nature journals' editorial policy for life
459 sciences research improve reporting? *BMJ Open Sci* 2019;3:e000035.
460 <https://doi.org/10.1136/bmjos-2017-000035>.
- 461 [3] Sander T, Ghanawi J, Wilson E, Muhammad S, Macleod M, Kahlert UD. Meta-analysis on
462 reporting practices as a source of heterogeneity in in vitro cancer research. *BMJ Open Sci*
463 2022;6:e100272. <https://doi.org/10.1136/bmjos-2021-100272>.
- 464 [4] Egger M, editor. *Systematic reviews in health care: meta-analysis in context*. 2. ed., [Nachdr.].
465 London: BMJ Books; 2009.
- 466 [5] Plint AC, Moher D, Morrison A, Schulz K, Altman DG, Hill C, et al. Does the CONSORT checklist
467 improve the quality of reports of randomised controlled trials? A systematic review. *Med J*
468 *Aust* 2006;185:263–7. <https://doi.org/10.5694/j.1326-5377.2006.tb00557.x>.
- 469 [6] Crossley NA, Sena E, Goehler J, Horn J, van der Worp B, Bath PMW, et al. Empirical Evidence of
470 Bias in the Design of Experimental Stroke Studies: A Metaepidemiologic Approach. *Stroke*
471 2008;39:929–34. <https://doi.org/10.1161/STROKEAHA.107.498725>.
- 472 [7] Hirst JA, Howick J, Aronson JK, Roberts N, Perera R, Koshiaris C, et al. The Need for
473 Randomization in Animal Trials: An Overview of Systematic Reviews. *PLoS ONE* 2014;9:e98856.
474 <https://doi.org/10.1371/journal.pone.0098856>.
- 475 [8] Emmerich CH, Gamboa LM, Hofmann MCJ, Bonin-Andresen M, Arbach O, Schendel P, et al.
476 Improving target assessment in biomedical research: the GOT-IT recommendations. *Nat Rev*
477 *Drug Discov* 2021;20:64–81. <https://doi.org/10.1038/s41573-020-0087-3>.
- 478 [9] Jadad AR, Cook DJ, Jones A, Klassen TP, Tugwell P, Moher M, et al. Methodology and Reports of
479 Systematic Reviews and Meta-analyses: A Comparison of Cochrane Reviews With Articles
480 Published in Paper-Based Journals. *JAMA* 1998;280:278.
481 <https://doi.org/10.1001/jama.280.3.278>.
- 482 [10] Shea B, Moher D, Graham I, Pham B, Tugwell P. A Comparison of the Quality of Cochrane
483 Reviews and Systematic Reviews Published in Paper-Based Journals. *Eval Health Prof*
484 2002;25:116–29. <https://doi.org/10.1177/0163278702025001008>.
- 485 [11] Fleming PS, Seehra J, Polychronopoulou A, Fedorowicz Z, Pandis N. Cochrane and non-
486 Cochrane systematic reviews in leading orthodontic journals: a quality paradigm? *Eur J Orthod*
487 2013;35:244–8. <https://doi.org/10.1093/ejo/cjs016>.
- 488 [12] Dosenovic S, Jelacic Kadic A, Vucic K, Markovina N, Pieper D, Puljak L. Comparison of
489 methodological quality rating of systematic reviews on neuropathic pain using AMSTAR and
490 R-AMSTAR. *BMC Med Res Methodol* 2018;18:37. <https://doi.org/10.1186/s12874-018-0493-y>.
- 491 [13] de Vries RBM, Hooijmans CR, Langendam MW, van Luijk J, Leenaars M, Ritskes-Hoitinga M, et
492 al. A protocol format for the preparation, registration and publication of systematic reviews of
493 animal intervention studies: Protocol format for animal systematic reviews. *Evid-Based Preclin*
494 *Med* 2015;2:e00007. <https://doi.org/10.1002/ebm2.7>.
- 495 [14] Macleod MR, O'Collins T, Howells DW, Donnan GA. Pooling of Animal Experimental Data
496 Reveals Influence of Study Design and Publication Bias. *Stroke* 2004;35:1203–8.
497 <https://doi.org/10.1161/01.STR.0000125719.25853.20>.
- 498 [15] Hooijmans CR, Rovers MM, de Vries RB, Leenaars M, Ritskes-Hoitinga M, Langendam MW.
499 SYRCLE's risk of bias tool for animal studies. *BMC Med Res Methodol* 2014;14:43.
500 <https://doi.org/10.1186/1471-2288-14-43>.
- 501 [16] Hunniford VT, Montroy J, Fergusson DA, Avey MT, Wever KE, McCann SK, et al. Epidemiology
502 and reporting characteristics of preclinical systematic reviews. *PLOS Biol* 2021;19:e3001177.
503 <https://doi.org/10.1371/journal.pbio.3001177>.

- 504 [17] CAMARADES Berlin. Preclinical Systematic Reviews & Meta-Analysis Wiki 2021.
505 <https://www.camarades.de/>. (accessed March 21, 2022).
- 506 [18] Higgins J, Thomas J, Chandler J, Cumpston M, Li T, Page M, et al. Cochrane Handbook for
507 Systematic Reviews of Interventions version 6.3 (updated February 2022). Cochrane; 2022.
- 508 [19] Duyx B, Swaen GMH, Urlings MJE, Bouter LM, Zeegers MP. The strong focus on positive results
509 in abstracts may cause bias in systematic reviews: a case study on abstract reporting bias. *Syst*
510 *Rev* 2019;8:174. <https://doi.org/10.1186/s13643-019-1082-9>.
- 511 [20] Marshall IJ, Noel-Storr A, Kuiper J, Thomas J, Wallace BC. Machine learning for identifying
512 Randomized Controlled Trials: An evaluation and practitioner’s guide. *Res Synth Methods*
513 2018;9:602–14. <https://doi.org/10.1002/jrsm.1287>.
- 514 [21] Bannach-Brown A, Przybyła P, Thomas J, Rice ASC, Ananiadou S, Liao J, et al. Machine learning
515 algorithms for systematic review: reducing workload in a preclinical review of animal studies
516 and reducing human screening error. *Syst Rev* 2019;8:23. <https://doi.org/10.1186/s13643-019-0942-7>.
- 517 [22] Marshall IJ, Wallace BC. Toward systematic review automation: a practical guide to using
518 machine learning tools in research synthesis. *Syst Rev* 2019;8:163, s13643-019-1074–9.
519 <https://doi.org/10.1186/s13643-019-1074-9>.
- 520 [23] Wallace BC, Kuiper J, Sharma A, Zhu MB, Marshall IJ. Extracting PICO Sentences from Clinical
521 Trial Reports using Supervised Distant Supervision. *J Mach Learn Res JMLR* 2016;17:132.
- 522 [24] Wang Q, Liao J, Lapata M, Macleod M. PICO Entity Extraction For Preclinical Animal Literature.
523 In Review; 2021. <https://doi.org/10.21203/rs.3.rs-1008099/v1>.
- 524 [25] Marshall IJ, Kuiper J, Wallace BC. RobotReviewer: evaluation of a system for automatically
525 assessing bias in clinical trials. *J Am Med Inform Assoc* 2016;23:193–201.
526 <https://doi.org/10.1093/jamia/ocv044>.
- 527 [26] Bahor Z, Liao J, Macleod MR, Bannach-Brown A, McCann SK, Wever KE, et al. Risk of bias
528 reporting in the recent animal focal cerebral ischaemia literature. *Clin Sci* 2017;131:2525–32.
529 <https://doi.org/10.1042/CS20160722>.
- 530 [27] Wang Q, Liao J, Lapata M, Macleod M. Risk of bias assessment in preclinical literature using
531 natural language processing. *Res Synth Methods* 2022;13:368–80.
532 <https://doi.org/10.1002/jrsm.1533>.
- 533 [28] Hair K, Bahor Z, Macleod M, Liao J, Sena ES. The Automated Systematic Search Deduplicator
534 (ASySD): a rapid, open-source, interoperable tool to remove duplicate citations in biomedical
535 systematic reviews. *Bioinformatics*; 2021. <https://doi.org/10.1101/2021.05.04.442412>.
- 536 [29] Bahor Z, Liao J, Currie G, Ayder C, Macleod M, McCann SK, et al. Development and uptake of an
537 online systematic review platform: the early years of the CAMARADES Systematic Review
538 Facility (SyRF). *BMJ Open Sci* 2021;5:e100103. <https://doi.org/10.1136/bmjos-2020-100103>.
- 539 [30] Bui DDA, Zeng-Treitler Q. Learning regular expressions for clinical text classification. *J Am Med*
540 *Inform Assoc* 2014;21:850–7. <https://doi.org/10.1136/amiajnl-2013-002411>.
- 541 [31] Liao J. Shihikoo/Autoannotation Release 2018. <https://doi.org/10.5281/ZENODO.1188823>.
- 542 [32] Buckland M, Gey F. The relationship between Recall and Precision. *J Am Soc Inf Sci* 1994;45:12–
543 9. [https://doi.org/10.1002/\(SICI\)1097-4571\(199401\)45:1<12::AID-ASIS2>3.0.CO;2-L](https://doi.org/10.1002/(SICI)1097-4571(199401)45:1<12::AID-ASIS2>3.0.CO;2-L).
- 544 [33] Kuhn M. Caret R Package 2019.
- 545 [34] Bethesda (MD). Entrez Programming Utilities Help. National Center for Biotechnology
546 Information (US); 2010.
- 547 [35] Dobrescu A, Nussbaumer-Streit B, Klerings I, Wagner G, Persad E, Sommer I, et al. Restricting
548 evidence syntheses of interventions to English-language publications is a viable methodological
549 shortcut for most medical topics: a systematic review. *J Clin Epidemiol* 2021;137:209–17.
550 <https://doi.org/10.1016/j.jclinepi.2021.04.012>.
- 551
- 552

553 **Appendix 1: Deviations from protocol**

554 **Method 1: Comparison of screening methods in an example systematic review**

- 555 • **Full text PDF retrieval** – due to time constraints, we did not conduct hand searching for PDFs
556 not retrieved by EndNote X8.
- 557 • **Inclusion and exclusion criteria for screening** – due to the RegEx being written in English, we
558 could only include records with an English-language full text in our analysis. However, our
559 search did not retrieve any records which had to be excluded solely due to this reason.
- 560 • **Regular expressions** – our protocol included both a RegEx for OGD and PC-12 cells.
561 However, we only used the OGD RegEx in our analysis.

562 **Method 2: Developing a trained machine learning classifier for *in vitro* systematic review** 563 **screening**

- 564 • **Risk of bias assessment** – we initially planned to additionally develop a tool to identify risk
565 of bias reporting (randomisation, blinding, and sample size calculation) but did not due to
566 time constraints and a lack of studies reporting randomisation, blinding, and sample size
567 calculation. Since publishing our protocol, such a tool has been developed for *in vivo*
568 research (Wang et al., 2021b). However it has not yet been validated on *in vitro* research.
- 569 • **Supplemented data from NPQIP** – we originally stated we would supplement our machine
570 learning training set with 640 records from NPQIP. This number was written in error, as only
571 453 records fit our definition of *in vitro* research.

572 **Appendix 2: Systematic Search Terms**

573 **(i) The incorrect strategy implemented in error: fragments containing errors are** 574 **underlined**

575 **NCBI PubMed**

576 (“oxygen-glucose deprivation/reoxygenation” OR “oxygen-glucose deprivation” OR “OGD” OR
577 “OGD/R” OR “oxygen and glucose-deprived model” OR “glutamate” OR “N-methyl-D-aspartate” OR
578 “NMDA” OR “H2O2” OR “hydrogen peroxide” OR “sodium nitroprusside” OR “SNP” OR “brain
579 ischemia” OR “brain ischaemia” OR “brain ischemic” OR “brain infarctions” OR “brain infarction” OR
580 “cerebral infarction” OR “cerebral infarctions” OR stroke OR “ischemic stroke” OR
581 “neuroprotection”) AND “PC12” OR “PC-12” OR “PC 12”

582 **Ovid Embase**

583 oxygen-glucose deprivation reoxygenation or oxygen-glucose deprivation or OGD or OGDR or oxygen
584 and glucose-deprived model or glutamate or N-methyl-D-aspartate or NMDA or H2O2 or hydrogen
585 peroxide or sodium nitroprusside or SNP or brain ischemia or brain ischaemia or brain ischemic or
586 brain infarctions or brain infarction or cerebral infarction or cerebral infarctions or stroke or ischemic
587 stroke or neuroprotection AND PC12 or PC-12 or PC 12

588

589 **(ii) The “correct” strategy, only deployed in our analysis of the intended search strategy.**

590 **NCBI PubMed**

591 (“oxygen-glucose deprivation/reoxygenation” OR “oxygen-glucose deprivation” OR “OGD” OR
592 “OGD/R” OR “oxygen and glucose-deprived model” OR “glutamate” OR “N-methyl-D-aspartate” OR
593 “NMDA” OR “H2O2” OR “hydrogen peroxide” OR “sodium nitroprusside” OR “SNP”)

594 **AND**

595 (“brain ischemia” OR “brain ischaemia” OR “brain ischemic” OR “brain infarctions” OR “brain
596 infarction” OR “cerebral infarction” OR “cerebral infarctions” OR stroke OR “ischemic stroke” OR
597 “neuroprotection”)

598 **AND**

599 (“PC12” OR “PC-12” OR “PC 12”)

600 **Ovid Embase**

601 (oxygen-glucose deprivation reoxygenation or oxygen-glucose deprivation or OGD or OGDR or
602 oxygen and glucose-deprived model or glutamate or N-methyl-D-aspartate or NMDA or H2O2 or
603 hydrogen peroxide or sodium nitroprusside or SNP) **and** (brain ischemia or brain ischaemia or brain
604 ischemic or brain infarctions or brain infarction or cerebral infarction or cerebral infarctions or stroke
605 or ischemic stroke or neuroprotection) **and** (PC12 or PC-12 or PC 12)

606

607

608

609 **Appendix 3: Regular Expression for Oxygen-Glucose Deprivation**

610 \bOGD\b|(?i)(oxygen|glucose)(\s|-| and)(glucose|oxygen) depriv(ation|ed)|deprived of (oxygen
611 and glucose|glucose and oxygen)

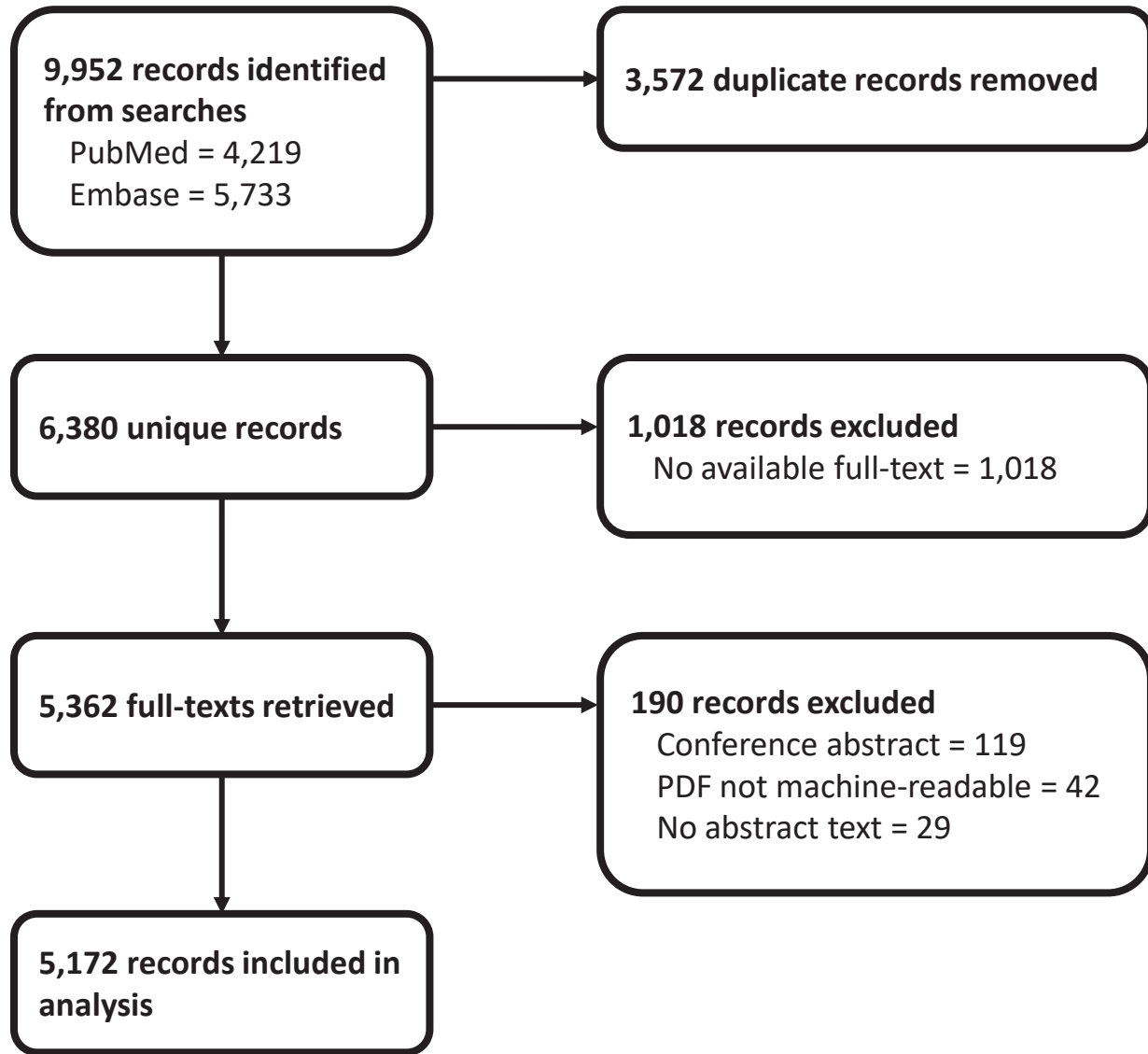


Figure 1: Flowchart of records retrieved from systematic searches, full texts retrieved, and records included in screening comparison analysis.

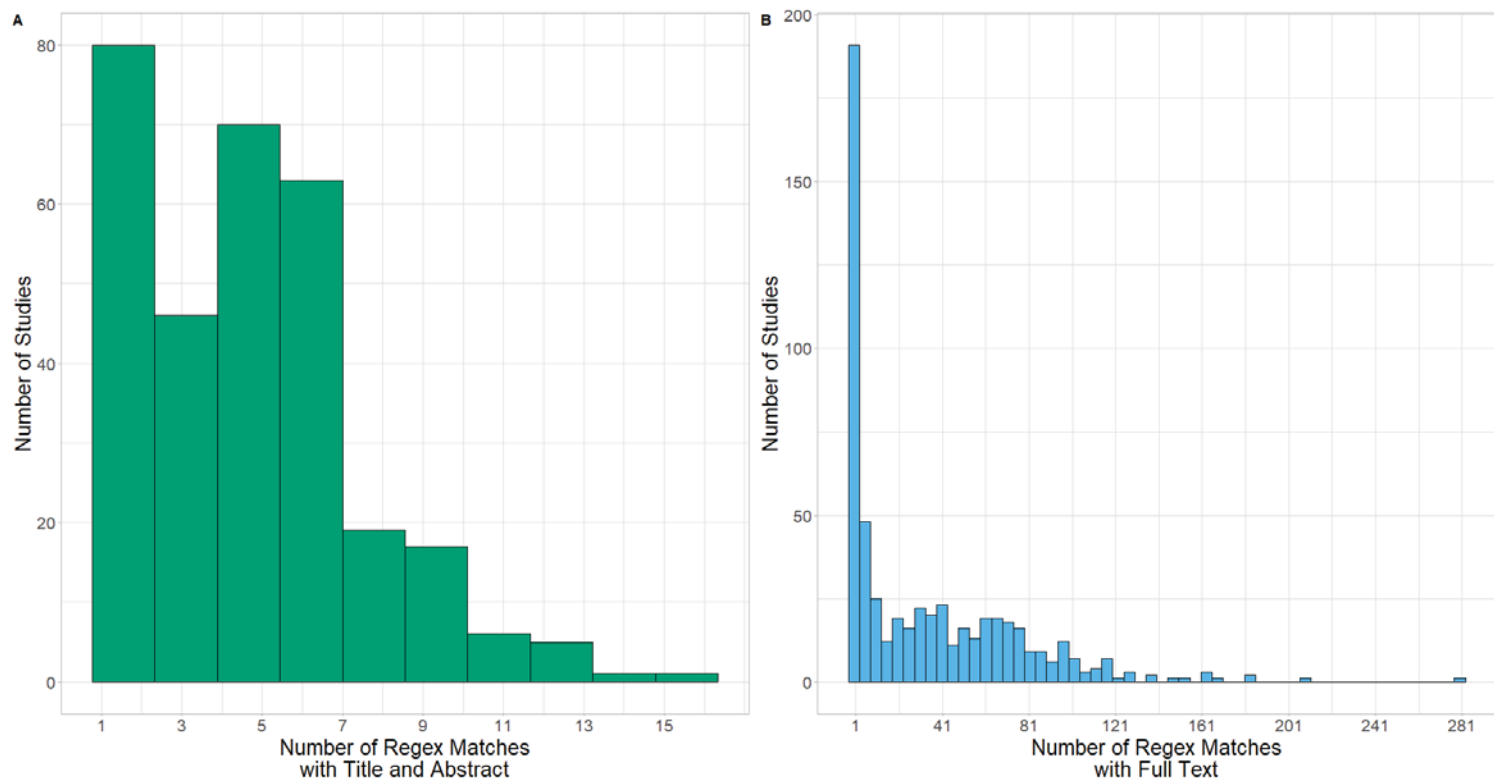


Figure 2: Histograms showing the number of studies against the number of regex matches with (A) title and abstract and (B) full text.

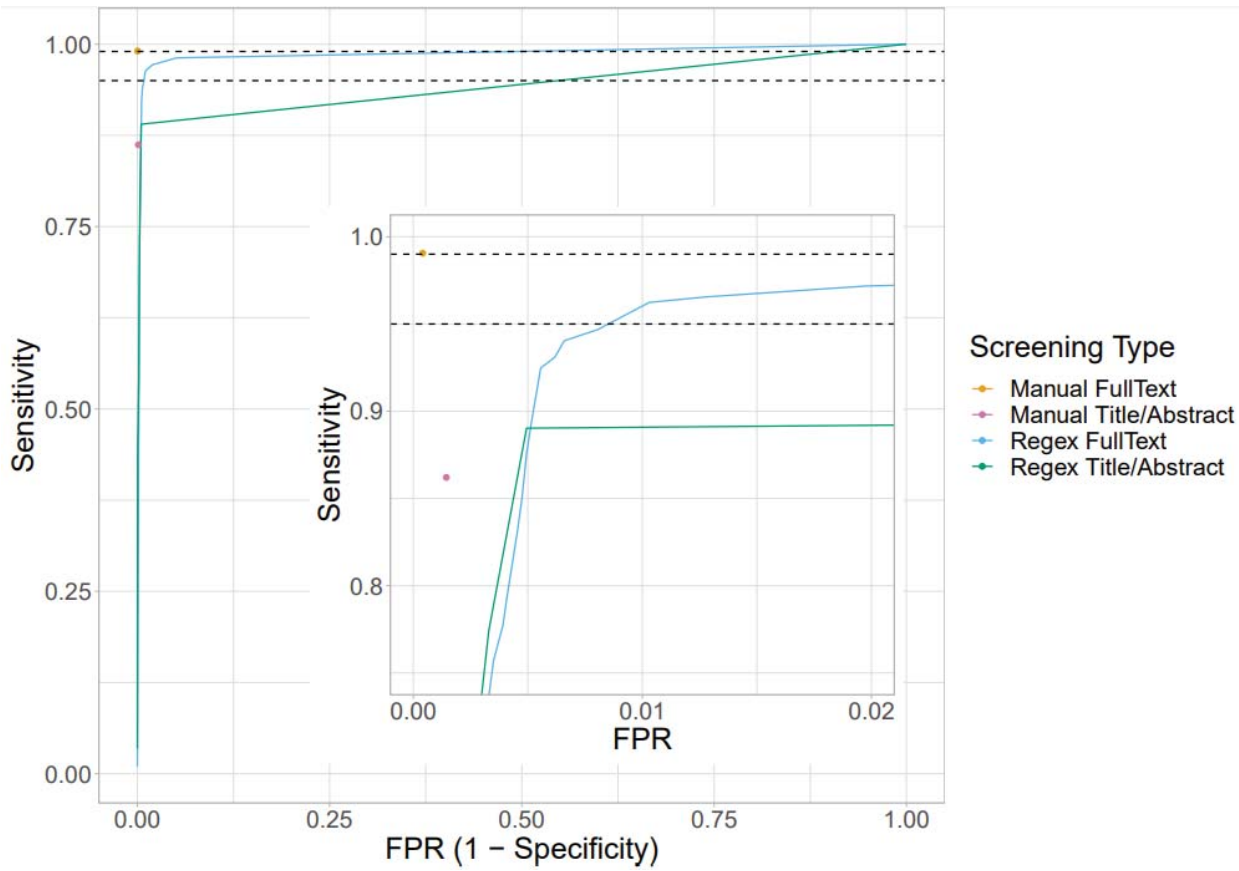


Figure 3: Receiver Operating Characteristic (ROC) curve showing the performance of all screening types at all thresholds. Horizontal dashed lines show 99% (0.99) and 95% (0.95) sensitivity. FPR = false positive rate. Inset shows the top left of the graph in more detail.

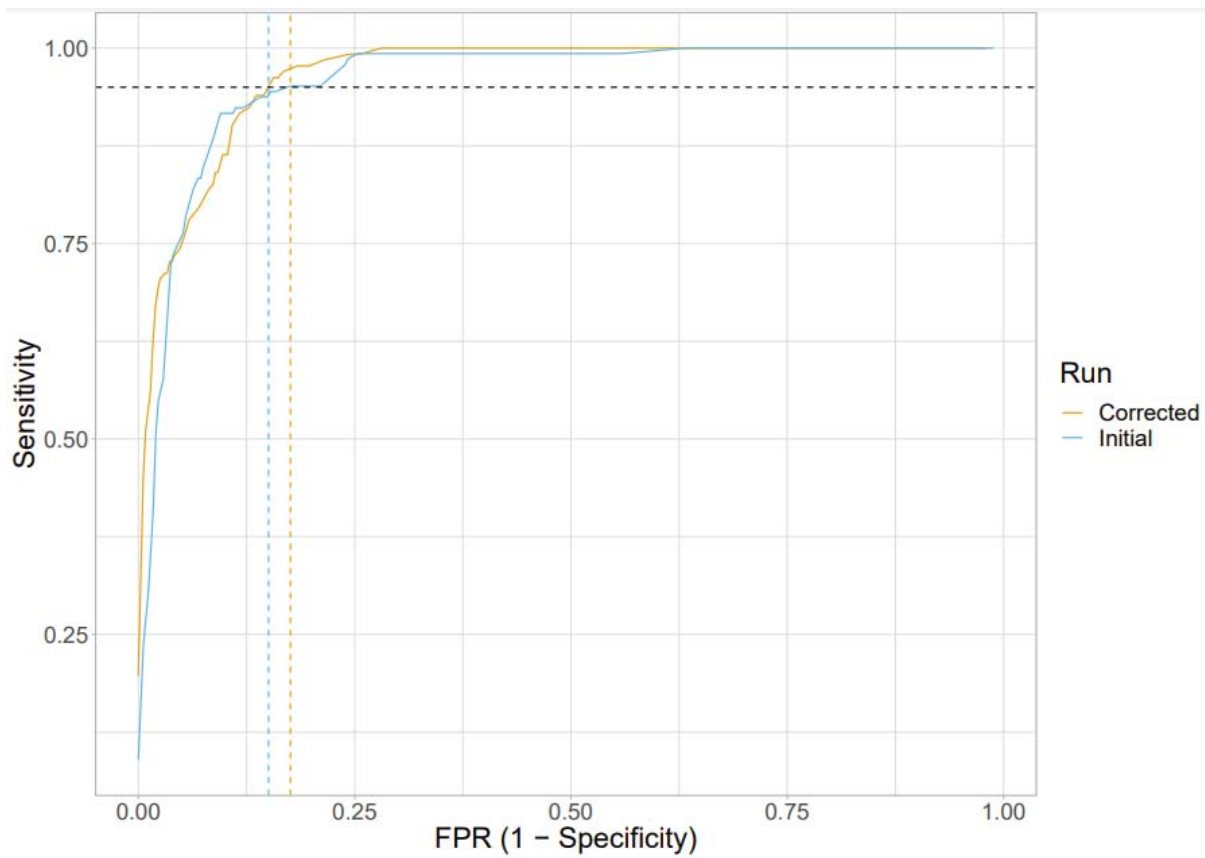


Figure 4: Receiver Operating Characteristic (ROC) curve showing both the initial and corrected performance of the machine learning algorithm at all thresholds. The vertical dashed lines show the optimal threshold (0.25 for the initial performance and 0.29 for the corrected performance). FPR = false positive rate.