

Research Master's programme Methodology & Statistics for the  
Behavioural, Biomedical and Social Sciences  
Utrecht University, the Netherlands

MSc Thesis Gerbrich Ferdinand (6466982)

TITLE: Active learning for efficient systematic reviews: evaluating  
models across research areas.

May 2020

Supervisors:

Prof. Dr. Rens van de Schoot

Jonathan de Bruin, MSc.

Dr. Raoul Schram

Second grader:

Prof. dr. René Eijkemans

Preferred journal of publication: [Systematic Reviews](#)

Word count: 6254

---

## <sup>1</sup> Background

<sup>2</sup> Systematic reviews are top of the bill in research. A systematic review brings together all studies relevant  
<sup>3</sup> to answer a specific research question [1]. Systematic reviews inform practice and policy [2] and are key in  
<sup>4</sup> developing clinical guidelines [3]. However, systematic reviews are costly because they involve the manual  
<sup>5</sup> screening of thousands of titles and abstracts to identify publications relevant to answering the research  
<sup>6</sup> question.

<sup>7</sup> Conducting a systematic review typically requires over a year of work by a team of researchers [4]. Never-  
<sup>8</sup> theless, systematic reviewers are often bound to a limited budget and timeframe. Currently, the demand  
<sup>9</sup> for systematic reviews exceeds the available time and resources by far [5]. Especially when the need for  
<sup>10</sup> guidelines is urgent - such as in the context of the current COVID-19 crisis - it is almost impossible to  
<sup>11</sup> provide a review that is both timely and comprehensive.

<sup>12</sup> To ensure a timely review, reducing workload in systematic reviews is essential. With advances in Machine  
<sup>13</sup> Learning (ML), there has been wide interest in tools to reduce workload in systematic reviews [6]. Various  
<sup>14</sup> learning models have been proposed, aiming to predict whether a given publication is relevant or irrelevant to  
<sup>15</sup> the systematic review. Previous findings suggest that such models potentially reduce workload with 30-70%  
<sup>16</sup> at the cost of losing 5% of relevant publications, i.e. 95% recall [7].

<sup>17</sup> A well-established approach in increasing efficiency in title and abstract screening is screening prioritization  
<sup>18</sup> [8,9]. In screening prioritization, the learning model presents the reviewer with the publications which are  
<sup>19</sup> most likely to be relevant first, thereby expediting the process of finding all of the relevant publications. Such  
<sup>20</sup> an approach allows for substantial time-savings in the screening process as the reviewer can decide to stop  
<sup>21</sup> screening after a sufficient number of relevant publications have been retrieved [10]. Moreover, reviewing  
<sup>22</sup> relevant publications early facilitates a faster transition of those publications to the next steps in the review  
<sup>23</sup> process [8].

<sup>24</sup> Recent studies have demonstrated the effectiveness of screening prioritization by means of active learning  
<sup>25</sup> models [10–16]. With active learning, the machine learning model can iteratively improve its predictions  
<sup>26</sup> on unlabelled data by allowing the model to select the records from which it wants to learn [17]. The  
<sup>27</sup> model queries these records to a human annotator who provides them with a label, from which the model  
<sup>28</sup> then updates its predictions. The general assumption is that by letting the model select which records are  
<sup>29</sup> labelled, the model can achieve higher accuracy while requiring the human annotator to label as few records  
<sup>30</sup> as possible [18]. Active learning has proven to be an efficient strategy in large unlabelled datasets where labels  
<sup>31</sup> are expensive to obtain [18]. This makes the screening phase in systematic reviewing an ideal candidate for

32 such models because typically, labelling a large number of publications is very costly. When active learning  
33 is applied in the screening phase, the reviewer screens publications that are selected by an active learning  
34 model. Subsequently, the active learning model learns from the reviewers' decision ('relevant', 'irrelevant')  
35 and uses this knowledge to update its predictions and to select the next publication to be screened by the  
36 reviewer.

37 The application of active learning models in reducing workload of systematic reviews has been extensively  
38 studied [10–12,15,16]. Whilst previous studies have evaluated active learning models in many forms and  
39 shapes [10–12,15], all studies used the same classification technique to predict relevanc of publications,  
40 namely Support Vector Machine. Findings from outside the field of active learning show that different  
41 classification techniques can serve different needs in the retrieval of relevant publications, for example recall  
42 versus precision [19,20]. Therefore, it is essential to evaluate different classification techniques in the context  
43 of active learning models. Another component known to influence performance of the models is the way  
44 how the textual content of titles and abstracts are represented in a model, called the feature extraction  
45 strategy [21,22]. Previous studies all adopt an effective but rather simplistic 'bag of words' strategy [10–  
46 12,15]. It is of interest to evaluate models using this approach by comparing them to models adopting a more  
47 sophisticated strategy, called 'Doc2vec' [21]. Lastly, previous studies have mainly focussed on reviews from a  
48 single scientific field, like medicine [15,16] and software engineering [11]. Model replications on reviews from  
49 varying research contexts are essential to draw conclusions about the general effectiveness of active learning  
50 models [7,23]. As far as known to the authors, Miwa et al [12] were the only researchers to make a direct  
51 comparison between two systematic reviews from different research areas, namely the social and the medical  
52 sciences. They found that active learning was more difficult on data from the social sciences due to the  
53 different nature of the vocabularies used. Therefore, it is of interest to evaluate model performance across  
54 different research contexts. Taken together, evaluations of active learning models in the context of systematic  
55 reviewing are required (1) across different classification techniques, (2) feature extraction strategies, and (3)  
56 review contexts. The current study aims to address these issues by answering the following research questions:

57 **RQ1** What is the performance of several active learning models across different classification techniques?

58 **RQ2** What is the performance of several active learning models across different feature extraction strate-  
59 gies?

60 **RQ3** Does the performance of active learning models differ across systematic reviews from different re-  
61 search areas?

62 The purpose of the current paper is to increase the evidence base on active learning models for reducing  
63 workload in title and abstract screening in systematic reviews. We adopt four different classification tech-

niques (Naive Bayes, Linear Regression, Support Vector Machine, and Random Forest) and two different feature extraction strategies (TF-IDF and Doc2vec) for the purpose of maximizing the number of identified relevant publications, while minimizing the number of publications needed to screen. Model performance was assessed by conducting a simulation on six systematic review datasets. Datasets were collected from the fields of medicine [24,25], virology [26], software engineering [11], behavioural public administration [27] and psychology [28], to assess generalizability of the models across research contexts. The models, datasets and simulations are implemented in a pipeline of active learning for screening prioritization, called **ASReview** [29]. **ASReview** is an open source and generic tool such that users can adapt and add modules as they like, encouraging fellow researchers to replicate findings from previous studies. All scripts and data used are openly published to facilitate usability and acceptability of ML-assisted title and abstract screening in the field of systematic review.

The remaining part of this paper is organized as follows. The Technical details section elaborates on the characteristics of active learning models for identifying relevant publications in the context of systematic reviews. The Simulation study section describes the study that was designed to answer the research questions. The findings of the simulation study are reported in the Results section. The implications of the findings in context of previous research are discussed in the Discussion section, followed by this study's main conclusions in the Conclusion section.

## 81 **Technical details**

82 What follows is a more detailed account of the active learning models. The structure and functions of the  
83 key components of the models are introduced to clarify the choices made in the design of the current study.

## 84 **Task description**

85 The screening process of a systematic review starts with all publications obtained in the search. The task  
86 is to identify which of these publications are relevant, by screening them at the title and abstract level. In  
87 active learning for screening prioritization, the screening process proceeds as follows:

- 88 • Start with the set of all unlabelled records (titles and abstracts),  $\mathcal{U}$ .
- 89 • The reviewer provides a label for a few, for example 5-10, records  $x \in \mathcal{U}$ , creating an set of labelled  
90 records  $x \in \mathcal{L}$  such that  $x \notin \mathcal{U}$ . The label can be either Relevant  $\langle x, R \rangle$  or Irrelevant  $\langle x, I \rangle$ .

- 91        • The active learning cycle starts:
- 92            1. A classifier,  $C$ , is trained on the labelled records  $\mathcal{L}$ ,  $C = \text{train}(\mathcal{L})$
- 93            2. The classifier predicts relevancy scores for all unlabelled records  $\mathcal{U}$ ,  $C(\mathcal{U})$
- 94            3. Based on the predictions by  $C$ , the model selects the most relevant record  $x^* \in \mathcal{U}$
- 95            4. The model queries the reviewer to screen this record,  $\langle x^*, ? \rangle$
- 96            5. The reviewer screens the record and provides a label,  $\langle x^*, R \rangle$  or  $\langle x^*, I \rangle$
- 97            6. The newly labelled record is added to the training data, such that  $x \in \mathcal{L}$  and  $x \notin \mathcal{U}$
- 98            7. Back to step 1.
- 99        In this active learning cycle, the model can incrementally improve its predictions on the remaining unlabelled  
100      title and abstracts. The relevant titles and abstracts are identified as early in the process as possible. The  
101      reviewer and the model keep interacting until the reviewer decides to stop or until all records been labelled.

## 102 Class imbalance problem

103      There are two classes in the dataset: relevant and irrelevant publications. Typically, the inclusion rate is low  
104      as only a fraction of the publications belong to the relevant class (2.94%, [4]). The class imbalance causes  
105      the classifier to miss relevant publications, because there are far fewer examples of relevant than irrelevant  
106      publications to train on [7]. Moreover, classifiers can achieve high accuracy but still fail to identify any of  
107      the relevant publications [15]. This is evident in the case of a systematic review dataset where only three  
108      percent of publications are relevant. A model would achieve 97% accuracy when classifying all publications  
109      as irrelevant, even though none of the relevant papers would have been correctly identified.

110      Previous studies have addressed the class imbalance problem by rebalancing the training data in various  
111      ways [7]. To decrease the class imbalance in the training data, the models in the current study rebalance  
112      the training set by Dynamic Resampling (DR). DR undersamples the number of irrelevant publications in  
113      the training data, whereas the number of relevant publications are oversampled such that the size of the  
114      training data remains the same. The ratio between relevant and irrelevant publications in the rebalanced  
115      training data is not fixed, but dynamically updated depending on the number of publications in the available  
116      training data, the number of publications in the total data, and the ratio between relevant and irrelevant  
117      publications in the available training data.

118 **Classification**

119 To make predictions on the unlabelled publications, a classifier is trained on features from the training data.  
120 A technique widely used in classification tasks is the Support Vector Machine (SVM). SVMs separate the  
121 data into classes by finding a multidimensional hyperplane [30,31]. SVMs have been proven to be effective  
122 in active learning models for screening prioritization [11,12]. Moreover, SVMs are the currently the only  
123 classifier implemented in ready-to-use software tools implementing active learning for screening prioritization  
124 (Abstrackr [32], Colandr [33], FASTREAD [11], Rayyan [34], and RobotAnalyst [35]).  
125 Whilst the performance of several classification techniques has been investigated in the ML-aided title-  
126 and-abstract screening field in general [19,20], the relatively new subfield of active learning for screening  
127 prioritization has not yet studied the performance of classifiers other than SVMs [10–15]. The current study  
128 aims to address this gap by exploring performance of three classifiers besides SVM:

- 129 • L2-regularized Logistic Regression (LR) models the probabilities describing the possible outcomes by a  
130 logistic function. The L2 penalty is imposed on the coefficients to reduce the number of features upon  
131 which the given solution is dependent [36].  
132 • Naive Bayes (NB) is a supervised learning algorithm often used in text classification. Based on Bayes'  
133 Theorem, with the ‘naive’ assumption that all features are independent given the class value [37].  
134 • Random Forests (RF) is a supervised learning algorithm where a large number of decision trees are fit  
135 on bootstrapped samples of the original data. All trees cast a vote on the class, which are aggregated  
136 into a class prediction for each instance [38].

137 These three classification techniques were selected because they are widely adopted methods in text classi-  
138 fication [39]. Moreover, these techniques can be run on a personal computer as they require a relatively low  
139 amount of processing power.

140 **Feature extraction**

141 To predict publication class, the classifier uses information from the publications in the dataset. Examples  
142 of such information are titles and abstracts. However, a model cannot make predictions from the titles and  
143 abstracts as they are; their textual content needs to be represented numerically. The textual information  
144 needs to be mapped to feature vectors. This process of numerically representing textual content is referred  
145 to as ‘feature extraction’.

146 A classical example of feature extraction is a ‘bag of words’ (bow) representation. For each text in the data  
147 set, the term frequency - the number of occurrences of each word - is stored. This leads to  $n$  features, where  
148  $n$  is the number of distinct words in the texts [36]. A serious weakness of this method is that it highly  
149 values often occurring but otherwise meaningless words such as “the”. A more sophisticated bow approach  
150 is Term-Frequency Inverse Document Frequency (TF-IDF), which circumvents this problem by adjusting  
151 the term frequency in a text with the inverse document frequency, the frequency of a given word in the  
152 entire data set [40]. A downside of TF-IDF and other bow methods is that they do not take into account  
153 the ordering of words, thereby ignoring semantics. An example of an approach that aims to overcome this  
154 weakness is Doc2vec (D2V). Doc2vec extracts features of the texts by a neural network, capable of grasping  
155 semantics by learning to predict the words in the texts [21].

## 156 **Query strategy**

157 The active learning model can adopt different strategies in selecting the next publication to be screened  
158 by the reviewer. A strategy mentioned before is selecting the publication with the highest probability of  
159 being relevant. In the active learning literature this is referred to as certainty-based active learning [17].  
160 Another well-known strategy is uncertainty-based active learning, where the instances that are presented  
161 next are those instances on which the model’s classifications are the least certain, i.e. close to 0.5 probability  
162 [17]. Traditionally, this strategy trains the most accurate model because the model can learn the most from  
163 instances it is uncertain about. However, a study comparing performance of both strategies in detecting  
164 relevant publications found that the accuracy gain of uncertainty-based screening was not significant [12].

165 Certainty-based active learning is the preferred strategy for the task at hand. Firstly, this strategy is far  
166 better suited to the goal of prioritizing relevant publications compared to uncertainty-based active learning,  
167 in which the publications are prioritized that the model is most uncertain about. Secondly, certainty-based  
168 active learning is far better equipped at dealing with imbalanced data in active learning, as it aims to present  
169 only records that belong to the relevant class [41].

## 170 **Simulation study**

171 The section below describes the simulation study that was carried out to answer the research questions.

<sub>172</sub> **Set-up**

<sub>173</sub> To address **RQ1**, four models combining every classifier with TF-IDF feature extraction were investigated:

- <sub>174</sub> 1. SVM + TF-IDF
- <sub>175</sub> 2. NB + TF-IDF
- <sub>176</sub> 3. RF + TF-IDF
- <sub>177</sub> 4. LR + TF-IDF

<sub>178</sub> To address **RQ2**, the classifiers were combined with Doc2vec feature extraction, leading to the following  
<sub>179</sub> three models:<sup>1</sup>

- <sub>180</sub> 5. SVM + Doc2vec
- <sub>181</sub> 6. RF + Doc2vec
- <sub>182</sub> 7. LR + Doc2vec

<sub>183</sub> The combination NB + D2V could not be tested because the Multinomial Naive Bayes classifier<sup>2</sup> can only  
<sub>184</sub> handle a feature matrix with positive values, whereas the Doc2vec feature extraction approach<sup>3</sup> produces  
<sub>185</sub> a feature matrix that can also contain negative values. Performance of the seven models was evaluated by  
<sub>186</sub> simulating every model on six systematic review datasets, addressing **RQ3**. Hence, 42 simulations were  
<sub>187</sub> carried out, representing all model-dataset combinations. To account for variance, every simulation was  
<sub>188</sub> repeated for 15 trials. For every simulation, hyperparameters were optimized through a random search to  
<sub>189</sub> arrive at maximum model performance. Simulations were run using ASReview's simulation mode [29]. There  
<sub>190</sub> was no need for a human reviewer as the model could query the labels in the data instead.

<sub>191</sub> Every simulation started with an initial training set of one relevant and one irrelevant publication to represent  
<sub>192</sub> a 'worst case scenario' where the reviewer has minimal prior knowledge on the publications in the data. To  
<sub>193</sub> account for bias due to the content of the initial publications, the initial training set was randomly sampled  
<sub>194</sub> from the dataset for every of the 15 trials. Although varying over trials, the 15 initial training sets were kept  
<sub>195</sub> constant over datasets to allow for a direct comparison of models within datasets. A seed value was set to  
<sub>196</sub> ensure reproducibility. The classifier was retrained every time after a publication had been labelled. The  
<sub>197</sub> simulation ended after all publications in the dataset had been labelled.

---

<sup>1</sup>

<sup>2</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html#sklearn.naive\\_bayes.MultinomialNB](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html#sklearn.naive_bayes.MultinomialNB)

<sup>3</sup><https://radimrehurek.com/gensim/models/doc2vec.html>

198 This study has been approved by the Ethics Committee of the Faculty of Social and Behavioural Sciences  
199 of Utrecht University, filed as an amendment under study 20-104. Simulations were run in ASReview,  
200 version 0.9.3 [29]. Analyses were carried out using R, version 3.6.1 [42]. Scripts and data are stored in the  
201 GitHub repository for this thesis<sup>4</sup>. The output resulting from the simulation was stored on the Open Science  
202 Framework page of this thesis,<sup>5</sup> as their size exceeded the storage limit of GitHub by far. Due to their  
203 large number, the simulations were carried out on Cartesius, the Dutch national supercomputer. Access was  
204 granted by SURF via a grant (ID EINF-156).

## 205 Datasets

206 The models were simulated on a convenience sample of six systematic review datasets. The data selection  
207 process was driven by two factors. Firstly, datasets were selected based on their background, given the need  
208 for datasets from diverse research areas. Secondly, datasets were selected by their availability, given the  
209 limited timespan of the current project. Thirdly, all original data files should be openly published with a  
210 CC-BY license, and are available through the ASReview GitHub page.

211 Datasets were collected from the fields of medicine, virology, software engineering, behavioural public ad-  
212 ministration, and psychology to assess generalizability of the models across research contexts. The Wilson  
213 dataset [43] is on a review on effectiveness and safety of treatments of Wilson Disease, a rare genetic disorder  
214 of copper metabolism [25]. The Ace dataset contains publications on the efficacy of Angiotensin-converting  
215 enzyme (ACE) inhibitors, a drug treatment for heart disease [24]. The Virus dataset is from a systematic  
216 review on studies that performed viral Metagenomic Next-Generation Sequencing (mNGS) in farm animals  
217 [26]. From the field of *software engineering*, the Software dataset contains publications from a review on  
218 fault prediction in source code [44]. The Nudging dataset [45] belongs to a systematic review on nudg-  
219 ing healthcare professionals [27], stemming from the area of *behavioural public administration*. The PTSD  
220 dataset contains publications from the field of *psychology*. The corresponding systematic review is on studies  
221 applying latent trajectory analyses on posttraumatic stress after exposure to traumatic events [28]. Of these  
222 six datasets, Ace, and Software have been used for model simulations in previous studies on ML-aided title  
223 and abstract screening, respectively [24] and [11].

224 Data were preprocessed from their original source into a test dataset, containing title and abstract of the  
225 publications obtained in the initial search. Candidate studies with missing abstracts and duplicate instances  
226 were removed from the data. Preprocessing scripts and resulting datasets can be found on the GitHub

---

<sup>4</sup><https://github.com/GerbrichFerdinands/asreview-thesis>

<sup>5</sup><https://osf.io/7mr2g/>

repository for this thesis. Test datasets were labelled to indicate which candidate studies were included in the systematic review, thereby indicating relevant publications. All test datasets consisted of thousands of candidate studies, of which only a fraction was deemed relevant to the systematic review. For the Virus and the Nudging dataset, the inclusion rate was about 5 percent. For the remaining six datasets, inclusion rates were centered around 1-2 percent. (Table 1).

Table 1: Statistics on the test datasets obtained from six original systematic reviews.

Dataset	Candidate publications	Relevant publications	Inclusion rate (%)
Nudging	1,847	100	5.4
PTSD	5,031	38	0.8
Software	8,896	104	1.2
Ace	2,235	41	1.8
Virus	2,304	114	5.0
Wilson	2,333	23	1.0

## Evaluating performance

Model performance was assessed by three different measures, Work Saved over Sampling (WSS), Relevant References Found (RRF), and Average Time to Discovery (ATD).

WSS indicates the reduction in publications needed to be screened, at a given level of recall [24]. Typically measured at a recall level of 0.95 [24], WSS@95 yields an estimate of the amount of work that can be saved at the cost of failing to identify 5% of relevant publications. In the current study, WSS is computed at 0.95 recall. RRF statistics are computed at 10%, representing the proportion of relevant publications that are found after screening 10% of all publications.

Both RRF and WSS are sensitive to random effects as these statistics are strongly dependent on the position of the cutoff value. Moreover, WSS makes assumptions about acceptable recall levels whereas this level might depend on the research question at hand [7]. A statistic that is not dependent on some arbitrary cutoff value is the ATD, which indicates the average proportion of publications needed to screen to find a relevant publication.

Furthermore, model performance was visualized by plotting recall curves. Plotting recall as a function of the proportion of screened publications offers insight in model performance throughout the entire screening process [11,13]. The curves give information in two directions. On the one hand they display the number of publications that need to be screened to achieve a certain level of recall (1-WSS), but on the other hand they

250 present how many relevant publications are identified after screening a certain proportion of all publications  
251 (RRF). Moreover, the recall curves relate to the ATD in such a way that the area above the curve is equal  
252 to the ATD.

253 For every simulation, the RRF@10, WSS@95, and ATD are reported as means over 15 trials. To indicate the  
254 spread of performance within simulations, the means are accompanied by an estimated<sup>6</sup> standard deviation  
255  $\hat{s}$ . To compare overall performance across datasets, median performance is reported for every dataset,  
256 accompanied by the Median Absolute Deviation (MAD), indicating variability between models within a  
257 certain dataset. Recall curves are plot for every simulation, representing the average recall over 15 trials  $\pm$   
258 the standard error of the mean.

## 259 Results

260 This section proceeds as follows. Firstly, the results of the Nudging dataset are discussed in detail to provide  
261 a basis for answering the research questions. Secondly, the results are presented for each research question  
262 over all datasets.

### 263 Evaluation on the Nudging dataset

264 Figure 1a shows the recall curves for all simulations on the Nudging dataset. As described in the previous  
265 section, these curves plot recall as a function of the proportion of publications screened. The curves represent  
266 the average recall over 15 trials  $\pm$  the standard error of the mean in the direction of the y-axis. The x-axis is  
267 cut off at 40% since all for simulations, the models reached 95% recall after screening 40% of the publications.  
268 The dashed horizontal lines indicate the RRF@10 values, the dashed vertical lines the WSS@95 values. The  
269 recall curves relate to the ATD such that ATD is equal to the area above the curve. The dashed grey diagonal  
270 line corresponds to the expected recall curve when publications are screened in a random order. Desirable  
271 model performance was defined as maximizing recall while minimizing the number of publications needed to  
272 screen.

273 The recall curves were used to examine model performance throughout the entire screening process and  
274 to make a visual comparison between models within datasets. For example in Figure 1a, after screening  
275 about 30% of the publications all models had already found 95% of the relevant publications. Moreover,

<sup>6</sup>The metrics for all individual 15 trials deviate slightly from the overall mean over 15 trials because of pre-averaging in the ASReview source code. As the analyses across all trials did not produce information on the 15 separate runs, the standard deviation of the mean,  $\hat{s}$ , was estimated by computing the standard deviation within the individual 15 trials.

<sup>276</sup> after screening 5% the green curve - representing the RF + TF-IDF model - splits away from the others  
<sup>277</sup> and remains to be the lowest of all curves until about 30% of publications have been screened. Hence, from  
<sup>278</sup> screening 5 to 30 percent of publications, the RF + TF-IDF model was the slowest in finding the relevant  
<sup>279</sup> publications. The ordering of the remaining recall curves changes throughout the screening process, but  
<sup>280</sup> maintain to show relatively similar performance at face value.

<sup>281</sup> Figure 1b shows a subset of the recall curves in Figure 1a, namely the curves for the first four models only  
<sup>282</sup> to allow for a visual comparison across classification techniques adopting the TF-IDF feature extraction  
<sup>283</sup> strategy. Figure 1c shows recall curves for the remaining three models to compare the models using Doc2vec  
<sup>284</sup> feature extraction. Figures 1d-f plot recall curves for models adopting the TF-IDF feature extraction strategy  
<sup>285</sup> to recall curves for their Doc2vec-using counterparts to allow for a comparison between models adopting TF-  
<sup>286</sup> IDF and models adopting Doc2vec feature extraction.

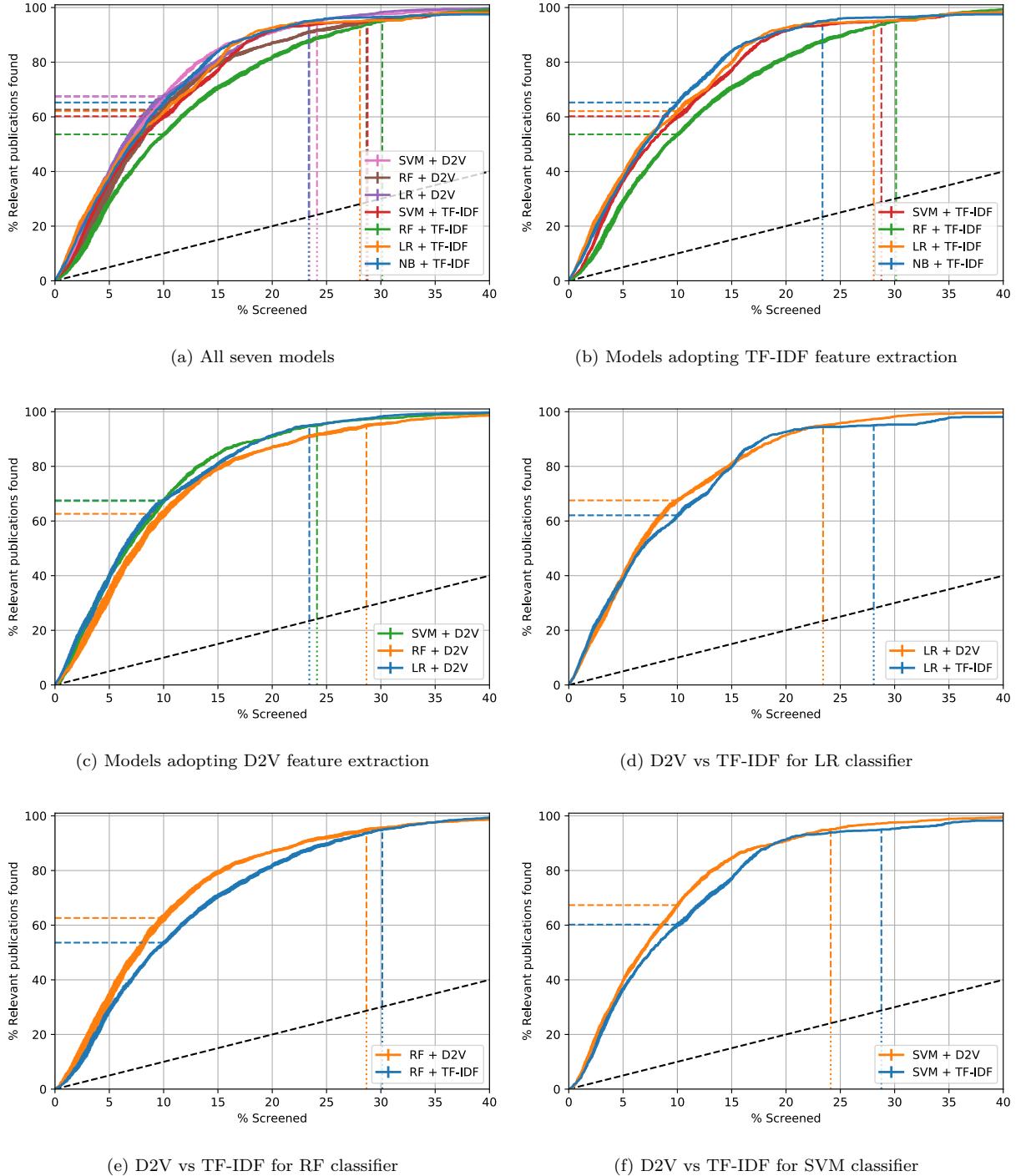


Figure 1: Recall curves for the Nudging dataset.

287 It can be seen from the data in the first column of Table 2 that in terms of ATD, the best performing models  
 288 on the Nudging dataset were SVM + D2V and LR + D2V, both with an ATD of 8.9%. This indicates that

289 the average proportion of publications needed to screen to find a relevant publication was 8.9% for both  
 290 models. In the SVM + D2V model, the standard deviation was 0.33 , whereas for the LR + D2V model  
 291  $\hat{s} = 0.47$ . This indicates that for the SVM + D2V model, the ATD values of individual trials were closer  
 292 to the overall mean compared to the LR + D2V model, meaning that the SVM + D2V model performed  
 293 more stable across different initial training datasets. Median ATD for this dataset was 9.6% with an MAD  
 294 of 1.06, indicating that for half of the models, the ATD was within 1.06 distance from the median ATD.

295 Table 2: ATD values ( $\bar{x}(\hat{s})$ ) for all model-dataset combinations, and median (MAD) for all datasets.

	Nudging	PTSD	Software	Ace	Virus	Wilson
SVM + TF-IDF	10.2 (0.19)	2.1 (0.13)	1.9 (0.04)	7.3 (1.18)	8.5 (0.17)	4.2 (0.33)
NB + TF-IDF	9.4 (0.29)	1.8 (0.11)	1.5 (0.03)	5.0 (0.53)	8.2 (0.22)	4.1 (0.37)
RF + TF-IDF	11.8 (0.44)	3.4 (0.27)	2.0 (0.09)	7.0 (0.76)	10.6 (0.42)	5.9 (1.20)
LR + TF-IDF	9.6 (0.19)	1.7 (0.10)	1.4 (0.02)	6.1 (1.20)	8.4 (0.24)	4.5 (0.34)
SVM + D2V	8.9 (0.33)	2.1 (0.15)	1.4 (0.05)	6.2 (0.34)	8.5 (0.21)	4.7 (0.31)
RF + D2V	10.4 (0.88)	3.1 (0.34)	1.6 (0.09)	7.3 (1.29)	9.3 (0.43)	7.5 (1.56)
LR + D2V	8.9 (0.47)	1.9 (0.17)	1.4 (0.04)	5.6 (0.18)	8.4 (0.41)	4.9 (0.32)
median (MAD)	9.6 (1.06)	2.1 (0.49)	1.5 (0.12)	6.2 (1.14)	8.5 (0.18)	4.7 (0.66)

296 As Table 3 shows, the highest WSS@95 value on the Nudging dataset was achieved by the NB + TF-IDF  
 297 model with a mean of 71.7, meaning that this model reduced the number of publications needed to screen  
 298 with 71.7% at the cost of losing 5% of relevant publications. The estimated standard deviation of 1.37  
 299 indicates that in terms of WSS@95, this model performed the most stable across trials. The model with the  
 300 lowest WSS@95 value was RF + TF-IDF ( $\bar{x} = 64.9$ ,  $\hat{s} = 2.50$ ). Median WSS@95 of these models was 66.9%,  
 301 with a MAD of 3.05%, indicating that WSS@95 values of models varied the most within this dataset.

302 Table 3: WSS@95 values ( $\bar{x}(\hat{s})$ ) for all model-dataset combinations, and median (MAD) for all datasets.

	Nudging	PTSD	Software	Ace	Virus	Wilson
SVM + TF-IDF	66.2 (2.90)	91.0 (0.41)	92.0 (0.10)	75.8 (1.95)	69.7 (0.81)	79.9 (2.09)
NB + TF-IDF	71.7 (1.37)	91.7 (0.27)	92.3 (0.08)	82.9 (0.99)	71.2 (0.62)	83.4 (0.89)
RF + TF-IDF	64.9 (2.50)	84.5 (3.38)	90.5 (0.34)	71.3 (4.03)	63.9 (3.54)	81.6 (3.35)
LR + TF-IDF	66.9 (4.01)	91.7 (0.18)	92.0 (0.10)	81.1 (1.31)	70.3 (0.65)	80.5 (0.65)
SVM + D2V	70.9 (1.68)	90.6 (0.73)	92.0 (0.21)	78.3 (1.92)	70.7 (1.76)	82.7 (1.44)
RF + D2V	66.3 (3.25)	88.2 (3.23)	91.0 (0.55)	68.6 (7.11)	67.2 (3.44)	77.9 (3.43)
LR + D2V	71.6 (1.66)	90.1 (0.63)	91.7 (0.13)	77.4 (1.03)	70.4 (1.34)	84.0 (0.77)
median (MAD)	66.9 (3.05)	90.6 (1.53)	92.0 (0.47)	77.4 (5.51)	70.3 (0.90)	81.6 (2.48)

303 As can be seen from the data in Table 4, LR + D2V was the best performing model in terms of RRF@10,  
 304 with a mean of 67.5 indicating that after screening 10% of publications, on average 67.5% of all relevant

305 publications had been identified, with a standard deviation of 2.59. The worst performing model was RF +  
 306 TF-IDF ( $\bar{x} = 53.6$ ,  $\hat{s} = 2.71$ ). Median performance was 62.6, with an MAD of 3.89 indicating again that  
 307 RRF@10 values were most dispersed for models within this dataset.

308 Table 4: RRF@10 values ( $\bar{x}$ , ( $\hat{s}$ )) for all model-dataset combinations, and median (MAD) for all datasets.

	Nudging	PTSD	Software	Ace	Virus	Wilson
SVM + TF-IDF	60.2 (3.12)	98.6 (1.40)	99.0 (0.00)	86.2 (5.25)	73.4 (1.62)	90.6 (1.17)
NB + TF-IDF	65.3 (2.61)	99.6 (0.95)	98.2 (0.34)	90.5 (1.40)	73.9 (1.70)	87.3 (2.55)
RF + TF-IDF	53.6 (2.71)	94.8 (1.60)	99.0 (0.00)	82.3 (2.75)	62.1 (3.19)	86.7 (5.82)
LR + TF-IDF	62.1 (2.59)	99.8 (0.70)	99.0 (0.00)	88.5 (5.16)	73.7 (1.48)	89.1 (2.30)
SVM + D2V	67.3 (3.00)	97.8 (1.12)	99.3 (0.44)	84.2 (2.78)	73.6 (2.54)	91.5 (4.16)
RF + D2V	62.6 (5.47)	97.1 (1.90)	99.2 (0.34)	80.8 (5.72)	67.3 (3.19)	75.5 (14.35)
LR + D2V	67.5 (2.59)	98.6 (1.40)	99.0 (0.00)	81.7 (1.81)	70.6 (2.21)	90.6 (5.00)
median (MAD)	62.6 (3.89)	98.6 (1.60)	99.0 (0.00)	84.2 (3.71)	73.4 (0.70)	89.1 (2.70)

## 309 Overall evaluation

310 Recall curves for the simulations on the five remaining datasets are presented in Figure 2. For the sake of  
 311 conciseness, recall curves are only plotted once per dataset, like in Figure 1a. Please refer to Additional file  
 312 1 for figures presenting subsets of recall curves for the remaining datasets, like in Figure 1b-f.

313 First of all, for all datasets, the models were able to detect the relevant publications much faster compared  
 314 to when screening publications at random order as the recall curves exceed the expected recall at screening  
 315 at random order by far. Even the worst results outperform this reference condition. Across simulations, the  
 316 ATD was at maximum 11.8% (in the Nudging dataset), the WSS@95 at least 63.9% (in the Virus dataset),  
 317 and the lowest RRF@10 was 53.6% (in the Nudging dataset). Interestingly, all these values were achieved  
 318 by the RF + TF-IDF model.

319 Similar to the simulations on the Nudging dataset (Figure 1b), the ordering of recall curves changes through-  
 320 out the screening process, indicating that model performance is dependent on the number of publications  
 321 that have been screened. Moreover, the ordering of models in the Nudging dataset (Figure 1b) does not  
 322 replicate on the remaining five datasets (Figure 2).

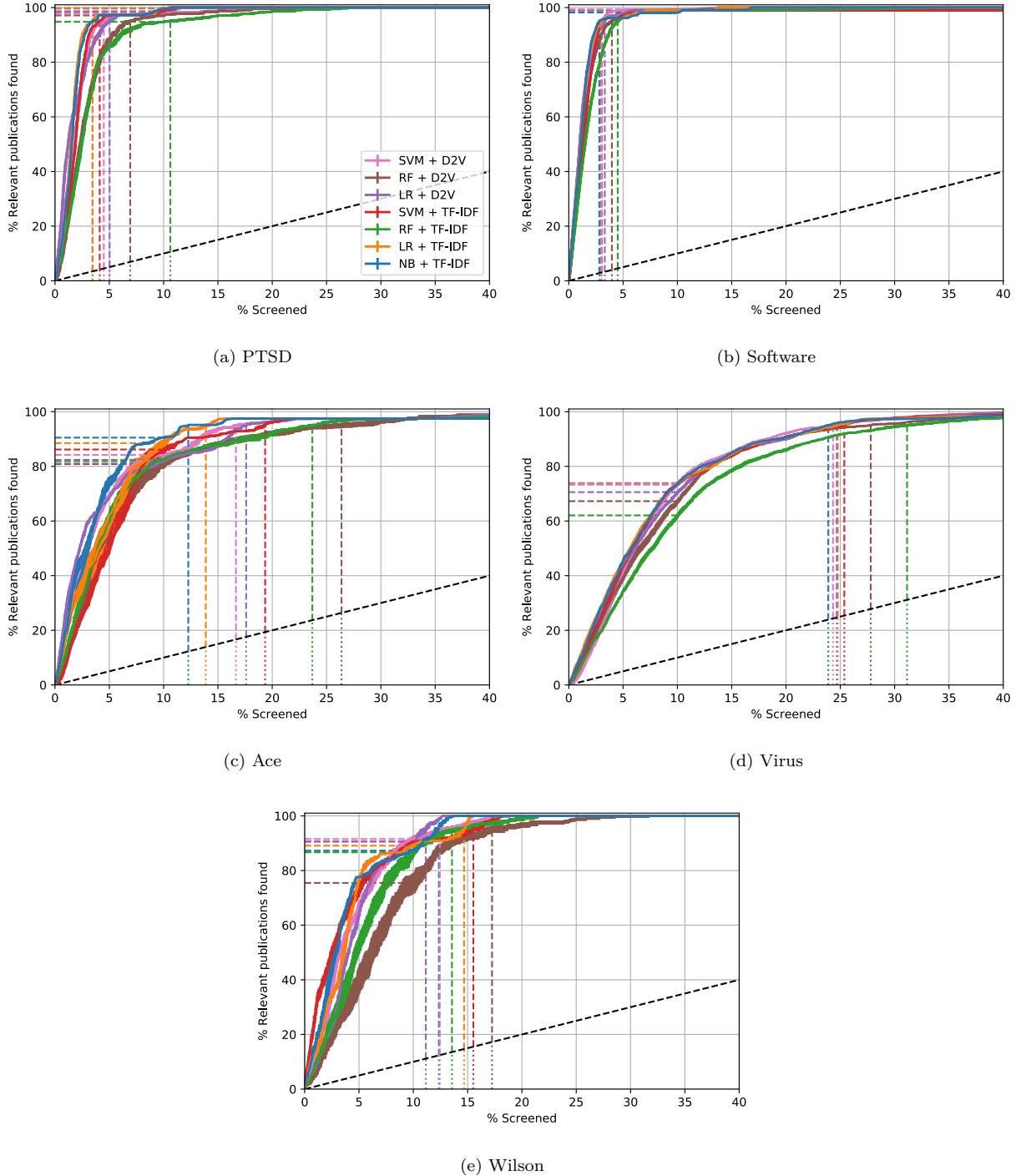


Figure 2: Recall curves for all seven models on (a) the PTSD, (b) Software, (c) Ace, (d) Virus, and (e) Wilson dataset.

323 **RQ1 - Comparison across classification techniques**

324 The first research question was aimed at evaluating the first four models adopting either the NB, SVM, LR  
325 or RF classification technique, combined with TF-IDF feature extraction. When comparing ATD-values of  
326 the models (Table 2), the NB + TF-IDF model ranked first in the Ace, Nudging, PTSD, Virus and Wilson  
327 dataset, and second in the PTSD and the Software dataset, in which the LR + TF-IDF model achieved the  
328 lowest ATD value. The RF + TF-IDF ranked last in all of the datasets except in the Ace dataset, where  
329 the SVM + TF-IDF model achieved the highest ATD-value.

330 Additionally, in terms of WSS@95 (Table 3) the ranking of models was strikingly similar across all datasets.  
331 In the Ace, Nudging, Software, and Virus dataset, the highest WSS@95 value was always achieved by the  
332 NB + TF-IDF model, followed by LR + TF-IDF, SVM + TF-IDF, and RF + TF-IDF. In the PTSD dataset  
333 this ranking applied as well, except that the LR + TF-IDF and NB + TF-IDF showed equal WSS@95 values.  
334 The ordering of the models for the Wilson dataset was NB + TF-IDF, RF + TF-IDF, LR + TF-IDF and  
335 SVM + TF-IDF.

336 Moreover, in terms of RRF@10 (Table 4) the NB + TF-IDF model achieved the highest RRF@10 value in  
337 the Ace, Nudging, and Wilson dataset. LR + TF-IDF ranked first in the PTSD dataset, SVM + TF-IDF  
338 was the best performing model within the Wilson dataset. The RF + TF-IDF model was again the worst  
339 performing model within all datasets, with one exception for the Software dataset. In this dataset, NB +  
340 TF-IDF ranked fourth, the remaining three models achieved an equal RRF@10 score.

341 Taken together, these results show that while all four models perform quite well, the NB + TF-IDF shows  
342 high performance on all measures across all datasets, whereas the RF + TF-IDF model never performed best  
343 on any of the measures across all datasets.

344 **RQ2 - Comparison across feature extraction techniques**

345 The following section is concerned with the question of how models using different feature extraction strate-  
346 gies relate to each other. The recall curves for the Nudging data (Figure 1d-f) show a clear trend of the models  
347 adopting Doc2vec feature extraction outperforming their TF-IDF counterparts. This trend also shows from  
348 the WSS@95 and RRF@10 values indicated by the vertical and horizontal lines in the figure. Likewise, the  
349 ATD values (Table 2) indicate that for the models adopting a particular classification technique, the model  
350 adopting Doc2vec feature extraction always achieved a smaller ATD-value than the model adopting TF-IDF  
351 feature extraction.

352 In contrast, this pattern of models adopting Doc2vec outperforming their TF-IDF counterparts in the Nudging  
353 dataset does not replicate across other datasets. Whether evaluated in terms of recall curves, WSS@95,  
354 RRF@10 or ATD, the findings were mixed. Neither one of the feature extraction strategies showed superior  
355 performance within certain datasets nor within certain classification techniques.

356 **RQ3 - Comparison across research contexts**

357 First of all, models showed much higher performance for some datasets than for others. While performance  
358 on the PTSD (Figure 2a) and the Software dataset (Figure 2b) was quite high, performance was much lower  
359 across models for the Nudging (Figure 1a) and Virus (Figure 2d) datasets. There does not seem to be a  
360 clear distinction between the datasets from the biomedical sciences (Ace, Virus, and Wilson) and datasets  
361 from other fields (Nudging, PTSD, Software). The PTSD, Software and Nudging dataset also demonstrated  
362 high performance in terms of the median ATD, WSS@95 and RRF@10 values for these models (Table 2, 3,  
363 and 4).

364 Secondly, variability of model performance differed across datasets. For the PTSD (Figure 2a), Software  
365 (Figure 2b), and the Virus (Figure 2d) datasets, recall curves form a tight group meaning that within  
366 these datasets, the models perform relatively similar. For the Nudging (Figure 1a), Ace (Figure 2c), and  
367 Wilson (Figure 2e) dataset, the recall curves are much further apart, indicating that model performance is  
368 much more dependent on the classification technique and feature extraction strategy. The MAD values  
369 of the ATD, WSS@95 and RRF@10 confirm that within the PTSD, Software and Virus datasets, model  
370 performance is less spread out than within the Nudging, Ace and Wilson dataset.

371 Moreover, the curves for the Ace (Figure 2c) and Wilson (Figure 2e) datasets show a larger standard error of  
372 the mean compared to other the other datasets. For these datasets, model performance seemed to be more  
373 dependent on the initial training data set compared to other datasets.

374 **Discussion**

375 The current study set out to evaluate performance of active learning models for the purpose of identifying  
376 relevant publications in systematic review datasets. It has been one of the first attempts to examine different  
377 classification strategies and feature extraction strategies in active learning models for systematic reviews.  
378 Moreover, this study has provided a deeper insight into the performance of active learning models across  
379 research contexts.

380 Overall, the findings confirm the great potential of active learning models in reducing workload for systematic  
381 reviewers. All models were able to detect 95% of the relevant publications after screening less than 40%  
382 of the total number of publications, indicating that active learning models can save more than half of the  
383 workload in the screening process. In a previous study, the Ace dataset was used to simulate a model that  
384 did not use active learning, finding a WSS@95 value of 56.61% [24], whereas the models in the current study  
385 achieved far superior WSS@95 values varying from 68.6% to 82.9% in this dataset. Active learning models  
386 clearly outperformed a model which did not use active learning. In addition, the Software dataset was used  
387 to simulate an active learning model [11] and reached WSS@95 of 91%, strikingly similar the WSS@95 values  
388 found in the current study which ranged from 90.5% to 92.3%.

### 389 **Classification techniques**

390 The first research question in this study sought to evaluate models adopting different classification techniques.  
391 The most obvious finding to emerge from these evaluations was that the NB + TF-IDF model consistently  
392 performed as one of the best models. The results suggest that whilst the widely used SVM-classifier performed  
393 fairly well, LR and NB classification strategies are interesting if not superior alternatives to the standard in  
394 this field.

### 395 **Feature extraction strategy**

396 The overall results on models adopting Doc2vec versus TF-IDF feature extraction strategy remain inconclu-  
397 sive. According to these findings, adopting Doc2vec instead of the well-established TF-IDF feature extraction  
398 strategy does not lead to better performing models. Given these results, although preliminary, preference  
399 goes out to teh TF-IDF feature extraction technique as this relatively simplistic technique will lead to more  
400 interpretable model.

### 401 **Research contexts**

402 Simulating models on a heterogenous collection of systematic review datasets revealed that model perfor-  
403 mance is very data-dependent. Within some datasets, models achieved much higher overall performance  
404 than within other datasets. Moreover, for some datasets, differences between models were much larger than  
405 for other datasets. It has been suggested that active learning is more difficult for datasets from the social  
406 sciences compared to data from the medical sciences [12]. This does not appear to be the case as performance  
407 within the biomedical datasets (Wilson, Virus, Ace) was not in any way superior to performance within the

408 datasets from the social sciences (PTSD and Nudging). An issue that emerges from these findings is that  
409 difficulty of active learning was not confined to any particular research area. A possible explanation for this is  
410 that difficulty of active learning could be attributed to factors more directly related to the systematic review  
411 at hand, such as the inclusion rate and the complexity of inclusion criteria used to identify relevant publica-  
412 tions [16,46]. Although the current study did not investigate the inclusion criteria of systematic reviews, the  
413 datasets on which the active learning models performed worst, Nudging and Virus, were interestingly also  
414 the datasets with the highest inclusion rates, 5.4% and 5.0%, respectively.

## 415 **Limitations and future research**

416 When applied in systematic review practice, the success of active learning models stands or falls down with the  
417 generalizability of model performance across unseen datasets. It is important to bear in mind that model  
418 hyperparameters were optimized for each model-dataset combination. Thus, the observed results reflect  
419 maximum model performance for the datasets at hand. The question remains whether model performance  
420 generalizes to datasets for which the hyperparameters were not optimized. Further research should be  
421 undertaken to determine the sensitivity of model performance to the hyperparameter values.

422 Screening publications in systematic reviews is typically a two-step process. First, titles and abstracts are  
423 screened to identify potentially relevant publications, called abstract inclusions. Second, the fulltexts of  
424 these publications are read to identify the relevant publications. This implies that the relevant publications  
425 are selected based on information that the models do not have. To truly assess the added value of active  
426 learning models in title-and-abstract screening, models should be evaluated on their capability of detecting  
427 the abstract inclusions instead of relevant publications only. However, this data is typically not available.  
428 Hence, greater efforts are needed to provide information on the abstract inclusions in openly published  
429 systematic review datasets.

430 An unanticipated finding was that the runtime of simulations varied widely across models, indicating that  
431 some models need more time to retrain after a publication has been labelled than other models. This finding  
432 has important implications for the practical application of such models, as an efficient model should be able  
433 to keep up with the decision-making speed of the reviewer. Further studies taking into account retraining  
434 time will need to be undertaken.

<sup>435</sup> **Conclusions**

<sup>436</sup> Overall, the findings of this study confirm that active learning models show great potential of finding relevant  
<sup>437</sup> publications in a systematic review dataset, while minimizing the number of publications needed to screen.  
<sup>438</sup> The results shed new light on the performance of different classification techniques, indicating that the Naive  
<sup>439</sup> Bayes classification technique is superior to the widely used Support Vector Machine. As model performance  
<sup>440</sup> differs vastly across datasets, this study raises the question what causes models to yield more workload  
<sup>441</sup> savings for some systematic review datasets than for others. In order to facilitate the applicability of active  
<sup>442</sup> learning models in systematic review practice, it is essential to identify how dataset characteristics relate to  
<sup>443</sup> model performance.

<sup>444</sup> **Declarations**

<sup>445</sup> **Ethics approval and consent to participate**

<sup>446</sup> Was reported in the main text.

<sup>447</sup> **Consent for publication**

<sup>448</sup> Not applicable.

<sup>449</sup> **Availability of data and materials**

<sup>450</sup> As reported in the main text, all data and materials are available through the GitHub repository for this  
<sup>451</sup> thesis, <https://github.com/GerbrichFerdinands/asreview-thesis>. This repository contains all systematic re-  
<sup>452</sup> view datasets used during this study and their preprocessing scripts, scripts and data on the hyperparamter  
<sup>453</sup> optimization, scripts on the simulations, scripts for analyzing the results of the simulations, and the source  
<sup>454</sup> files for this manuscript. All output files of the simulation study are stored on the Open Science Framework  
<sup>455</sup> page of this thesis, <https://osf.io/7mr2g/>.

<sup>456</sup> **Competing interests**

<sup>457</sup> The author declares that they has no competing interests.

458 **Funding**

459 Computing hours on the Cartesius supercomputer were funded by SURFsara. SURFsara had no role whatsoever  
460 ever in the design of the current study, nor in the data collection, analysis and interpretation, nor in writing  
461 the manuscript.

462 **Acknowledgements**

463 I am grateful for all researchers who have made great efforts to openly publish the data on their systematic  
464 reviews, special thanks go out to Rosanna Nagtegaal. I would also like to thank Caroline van Baal for  
465 supporting me in writing this thesis, and prof. dr. René Eijkemans, for being the second grader of this  
466 thesis. Finally, I would like to express my appreciation to my supervisors prof. dr. Rens van de Schoot,  
467 Jonathan de Bruin, and dr. Raoul Schram. Your door was always open and your enthusiasm was contagious.

## References

- [1] PRISMA-P Group, Moher D, Shamseer L, Clarke M, Ghersi D, Liberati A, et al. Preferred reporting items for systematic review and meta-analysis protocols (PRISMA-P) 2015 statement. *Syst Rev* 2015;4:1. <https://doi.org/10.1186/2046-4053-4-1>.
- [2] Gough D, Elbourne D. Systematic Research Synthesis to Inform Policy, Practice and Democratic Debate. *Soc Policy Soc* 2002;1:225–36. <https://doi.org/10/bdmp7h>.
- [3] Chalmers I. The lethal consequences of failing to make full use of all relevant evidence about the effects of medical treatments: The importance of systematic reviews. In: Treating individuals—from randomised trials to personalised medicine., Lancet; 2007, pp. 37–58.
- [4] Borah R, Brown AW, Capers PL, Kaiser KA. Analysis of the time and workers needed to conduct systematic reviews of medical interventions using data from the PROSPERO registry. *BMJ Open* 2017;7:e012545. <https://doi.org/10/f9tf57>.
- [5] Lau J. Editorial: Systematic review automation thematic series. *Syst Rev* 2019;8:70. <https://doi.org/10/ggsmwf>.
- [6] Harrison H, Griffin SJ, Kuhn I, Usher-Smith JA. Software tools to support title and abstract screening for systematic reviews in healthcare: An evaluation. *BMC Med Res Methodol* 2020;20:7. <https://doi.org/10.1186/s12874-020-0897-3>.

- [7] O'Mara-Eves A, Thomas J, McNaught J, Miwa M, Ananiadou S. Using text mining for study identification in systematic reviews: A systematic review of current approaches. *Syst Rev* 2015;4:5. <https://doi.org/10.1186/2046-4053-4-5>.
- [8] Cohen AM, Ambert K, McDonagh M. Cross-Topic Learning for Work Prioritization in Systematic Review Creation and Update. *J Am Med Inform Assoc* 2009;16:690–704. <https://doi.org/10/c3shq2>.
- [9] Shemilt I, Simon A, Hollands GJ, Marteau TM, Ogilvie D, O'Mara-Eves A, et al. Pinpointing needles in giant haystacks: Use of text mining to reduce impractical screening workload in extremely large scoping reviews. *Res Synth Methods* 2014;5:31–49. <https://doi.org/10.1002/jrsm.1093>.
- [10] Yu Z, Menzies T. FAST2: An intelligent assistant for finding relevant papers. *Expert Syst Appl* 2019;120:57–71. <https://doi.org/10.1016/j.eswa.2018.11.021>.
- [11] Yu Z, Kraft NA, Menzies T. Finding better active learners for faster literature reviews. *Empir Softw Eng* 2018;23:3161–86. <https://doi.org/10.1007/s10664-017-9587-0>.
- [12] Miwa M, Thomas J, O'Mara-Eves A, Ananiadou S. Reducing systematic review workload through certainty-based screening. *J Biomed Inform* 2014;51:242–53. <https://doi.org/10.1016/j.jbi.2014.06.005>.
- [13] Cormack GV, Grossman MR. Evaluation of machine-learning protocols for technology-assisted review in electronic discovery. In: Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval, Gold Coast, Queensland, Australia: Association for Computing Machinery; 2014, pp. 153–62. <https://doi.org/10.1145/2600428.2609601>.
- [14] Cormack GV, Grossman MR. Autonomy and Reliability of Continuous Active Learning for Technology-Assisted Review 2015.
- [15] Wallace BC, Trikalinos TA, Lau J, Brodley C, Schmid CH. Semi-automated screening of biomedical citations for systematic reviews. *BMC Bioinform* 2010;11:55. <https://doi.org/10.1186/1471-2105-11-55>.
- [16] Gates A, Johnson C, Hartling L. Technology-assisted title and abstract screening for systematic reviews: A retrospective evaluation of the Abstrackr machine learning tool. *Syst Rev* 2018;7:45. <https://doi.org/10/ggpx4>.
- [17] Settles B. Active Learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning* 2012;6:1–114. <https://doi.org/10.2200/S00429ED1V01Y201207AIM018>.
- [18] Settles B. Active Learning Literature Survey. University of Wisconsin-Madison Department of Computer Sciences; 2009.

- [19] Kilicoglu H, Demner-Fushman D, Rindflesch TC, Wilczynski NL, Haynes RB. Towards Automatic Recognition of Scientifically Rigorous Clinical Research Evidence. *J Am Med Inform Assn* 2009;16:25–31. <https://doi.org/10/bjkh9>.
- [20] Aphinyanaphongs Y. Text Categorization Models for High-Quality Article Retrieval in Internal Medicine. *J Am Med Inform Assoc* 2004;12:207–16. <https://doi.org/10/cpv52>.
- [21] Le QV, Mikolov T. Distributed Representations of Sentences and Documents 2014.
- [22] Zhang W, Yoshida T, Tang X. A comparative study of TF\*IDF, LSI and multi-words for text classification. *Expert Syst Appl* 2011;38:2758–65. <https://doi.org/10/dp7268>.
- [23] Marshall IJ, Johnson BT, Wang Z, Rajasekaran S, Wallace BC. Semi-Automated evidence synthesis in health psychology: Current methods and future prospects. *Health Psychol Rev* 2020;14:145–58. <https://doi.org/10/ggjv98>.
- [24] Cohen AM, Hersh WR, Peterson K, Yen P-Y. Reducing Workload in Systematic Review Preparation Using Automated Citation Classification. *J Am Med Inform Assoc* 2006;13:206–19. <https://doi.org/10.1197/jamia.M1929>.
- [25] Appenzeller-Herzog C, Mathes T, Heeres MLS, Weiss KH, Houwen RHJ, Ewald H. Comparative effectiveness of common therapies for Wilson disease: A systematic review and meta-analysis of controlled studies. *Liver Int* 2019;39:2136–52. <https://doi.org/10.1111/liv.14179>.
- [26] Kwok KTT, Nieuwenhuijse DF, Phan MVT, Koopmans MPG. Virus Metagenomics in Farm Animals: A Systematic Review. *Viruses* 2020;12:107. <https://doi.org/10.3390/v12010107>.
- [27] Nagtegaal R, Tummers L, Noordegraaf M, Bekkers V. Nudging healthcare professionals towards evidence-based medicine: A systematic scoping review. *J Behav Public Adm* 2019;2. <https://doi.org/doi.org/10.30636/jbpa.22.71>.
- [28] van de Schoot R, Sijbrandij M, Winter SD, Depaoli S, Vermunt JK. The GRoLTS-Checklist: Guidelines for reporting on latent trajectory studies. *Struct Equ Model Multidiscip J* 2017;24:451–67. <https://doi.org/10/gdpcw9>.
- [29] van de Schoot R, de Bruin J, Schram R, Zahedi P, Kramer B, Ferdinand G, et al. ASReview: Active learning for systematic reviews 2020. <https://doi.org/10/ggssnj>.
- [30] Tong S, Koller D. Support vector machine active learning with applications to text classification. *J Mach Learn Res* 2001;2:45–66.

- [31] Kremer J, Steenstrup Pedersen K, Igel C. Active learning with support vector machines. *WIREs Data Min Knowl Discov* 2014;4:313–26. <https://doi.org/10/f6fss7>.
- [32] Wallace BC, Small K, Brodley CE, Lau J, Trikalinos TA. Deploying an interactive machine learning system in an evidence-based practice center: Abstrackr. In: Proceedings of the 2nd ACM SIGHIT International Health Informatics Symposium, Miami, Florida, USA: Association for Computing Machinery; 2012, pp. 819–24. <https://doi.org/10.1145/2110363.2110464>.
- [33] Cheng SH, Augustin C, Bethel A, Gill D, Anzaroot S, Brun J, et al. Using machine learning to advance synthesis and use of conservation and environmental evidence. *Conserv Biol* 2018;32:762–4. <https://doi.org/10.1111/cobi.13117>.
- [34] Ouzzani M, Hammady H, Fedorowicz Z, Elmagarmid A. Rayyan—a web and mobile app for systematic reviews. *Syst Rev* 2016;5:210. <https://doi.org/10.1186/s13643-016-0384-4>.
- [35] Przybyła P, Brockmeier AJ, Kontonatsios G, Pogam M-AL, McNaught J, Erik von Elm, et al. Prioritising references for systematic reviews with RobotAnalyst: A user study. *Res Synth Methods* 2018;9:470–88. <https://doi.org/10.1002/jrsm.1311>.
- [36] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine learning in Python. *J Mach Learn Res* 2011;12:2825–30.
- [37] Zhang H. The Optimality of Naive Bayes. In: vol. 2, 2004.
- [38] Breiman L. Random Forests. *Machine Learning* 2001;45:5–32. <https://doi.org/10.1023/A:1010933404324>.
- [39] Aggarwal CC, Zhai C. A Survey of Text Classification Algorithms. In: Aggarwal CC, Zhai C, editors. *Mining Text Data*, Boston, MA: Springer US; 2012, pp. 163–222. [https://doi.org/10.1007/978-1-4614-3223-4\\_6](https://doi.org/10.1007/978-1-4614-3223-4_6).
- [40] Ramos J, others. Using tf-idf to determine word relevance in document queries. In: Proceedings of the first instructional conference on machine learning, vol. 242, Piscataway, NJ; 2003, pp. 133–42.
- [41] Fu JH, Lee SL. Certainty-Enhanced Active Learning for Improving Imbalanced Data Classification. In: 2011 IEEE 11th International Conference on Data Mining Workshops, Vancouver, BC, Canada: IEEE; 2011, pp. 405–12. <https://doi.org/10.1109/ICDMW.2011.43>.
- [42] R Core Team. R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing; 2019.

- [43] Appenzeller-Herzog C. Data from Comparative effectiveness of common therapies for Wilson disease: A systematic review and meta-analysis of controlled studies 2020.
- [44] Hall T, Beecham S, Bowes D, Gray D, Counsell S. A Systematic Literature Review on Fault Prediction Performance in Software Engineering. *IEEE Trans Softw Eng* 2012;38:1276–304. <https://doi.org/10.1109/TSE.2011.103>.
- [45] Nagtegaal R, Tummers L, Noordegraaf M, Bekkers V. Nudging healthcare professionals towards evidence-based medicine: A systematic scoping review 2019.
- [46] Rathbone J, Hoffmann T, Glasziou P. Faster title and abstract screening? Evaluating Abstrackr, a semi-automated online screening program for systematic reviewers. *Systematic Reviews* 2015;4:80. <https://doi.org/10/f7ms4w>.

## **Additional file 1**

Recall curves plot separately for the PTSD, Software, Ace, Virus and Wilson datasets.

## PTSD

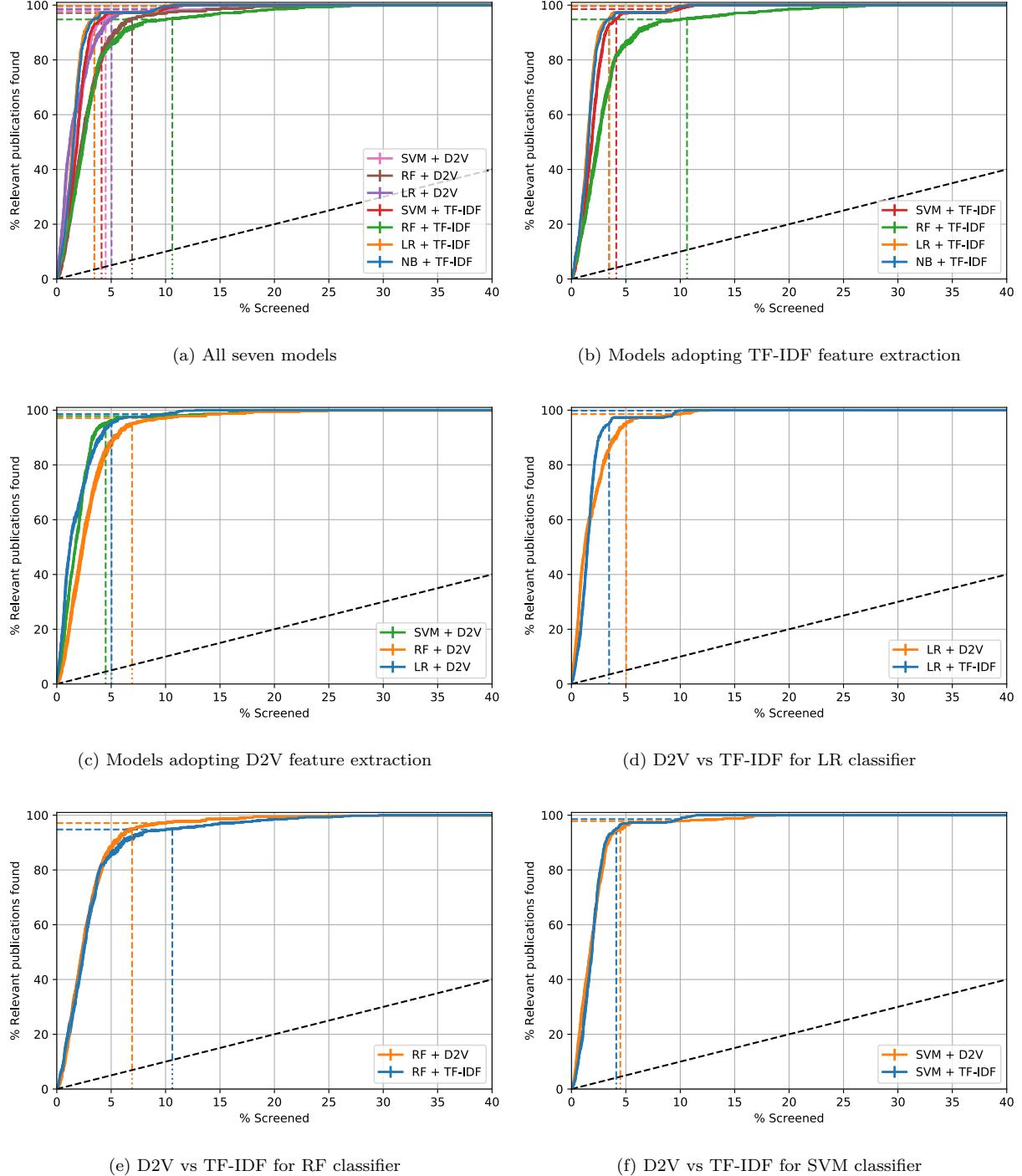


Figure 3: Recall curves for the PTSD dataset.

## Software

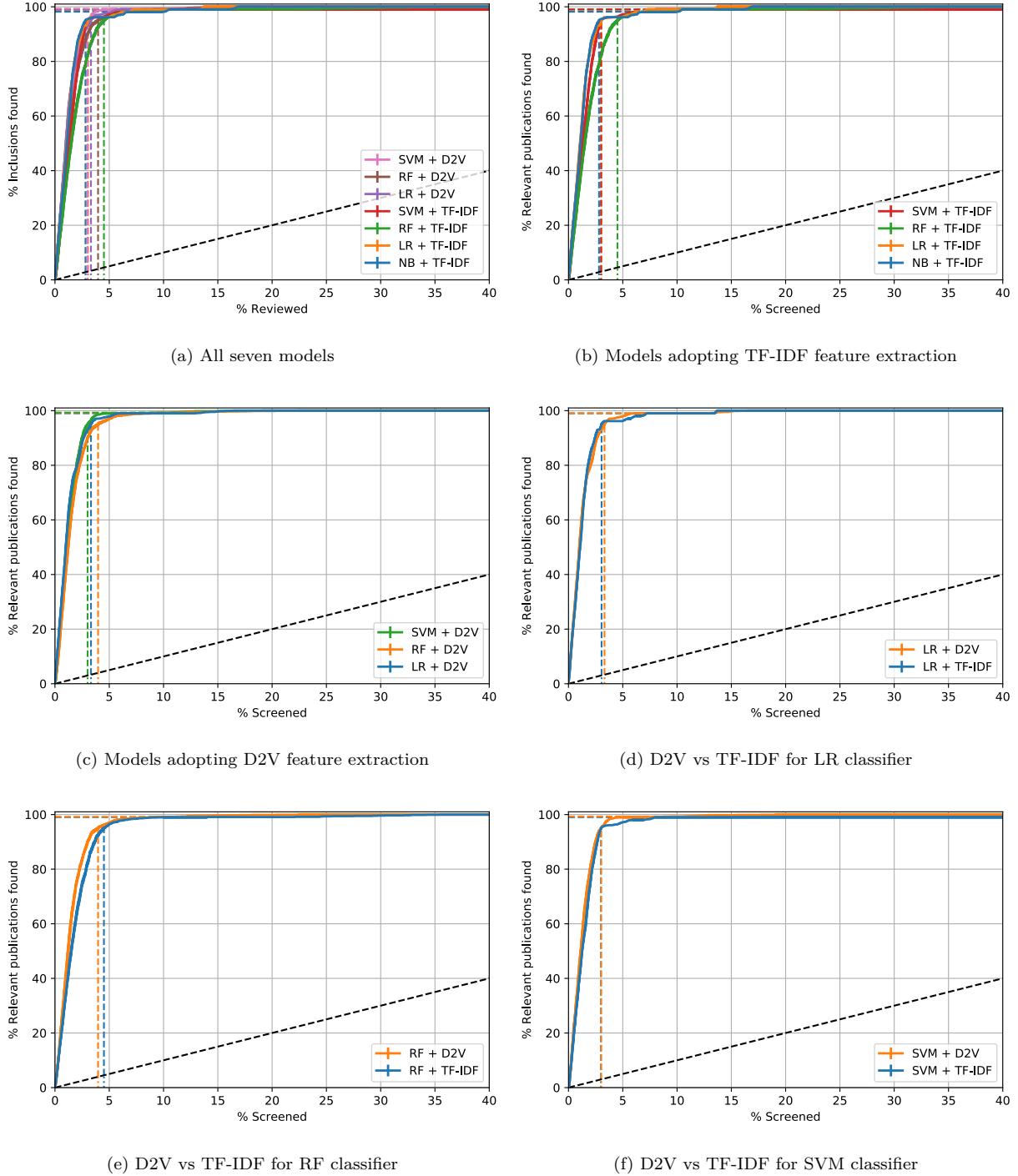


Figure 4: Recall curves for the Software dataset.

## Ace

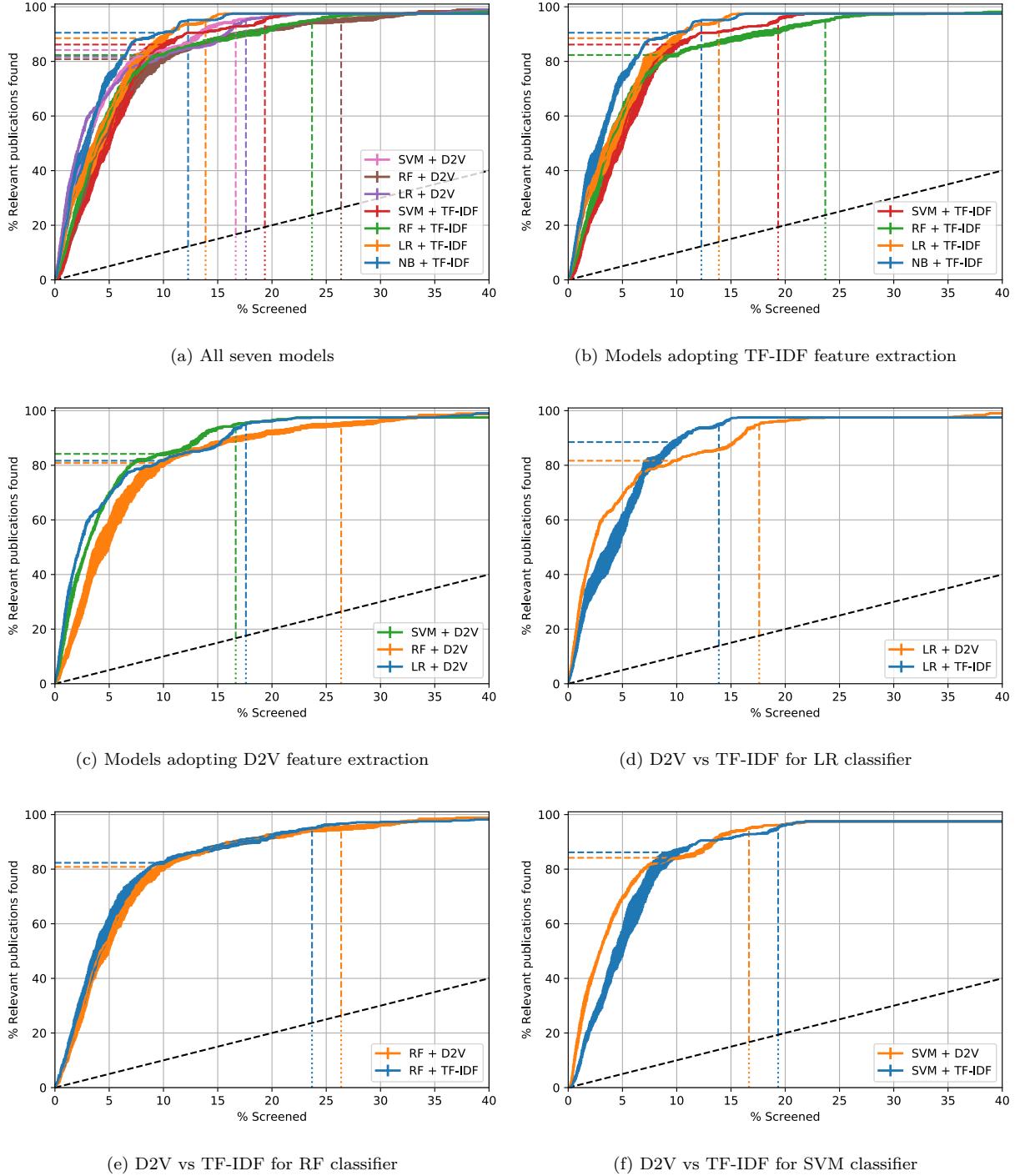


Figure 5: Recall curves for the Ace dataset.

## Virus

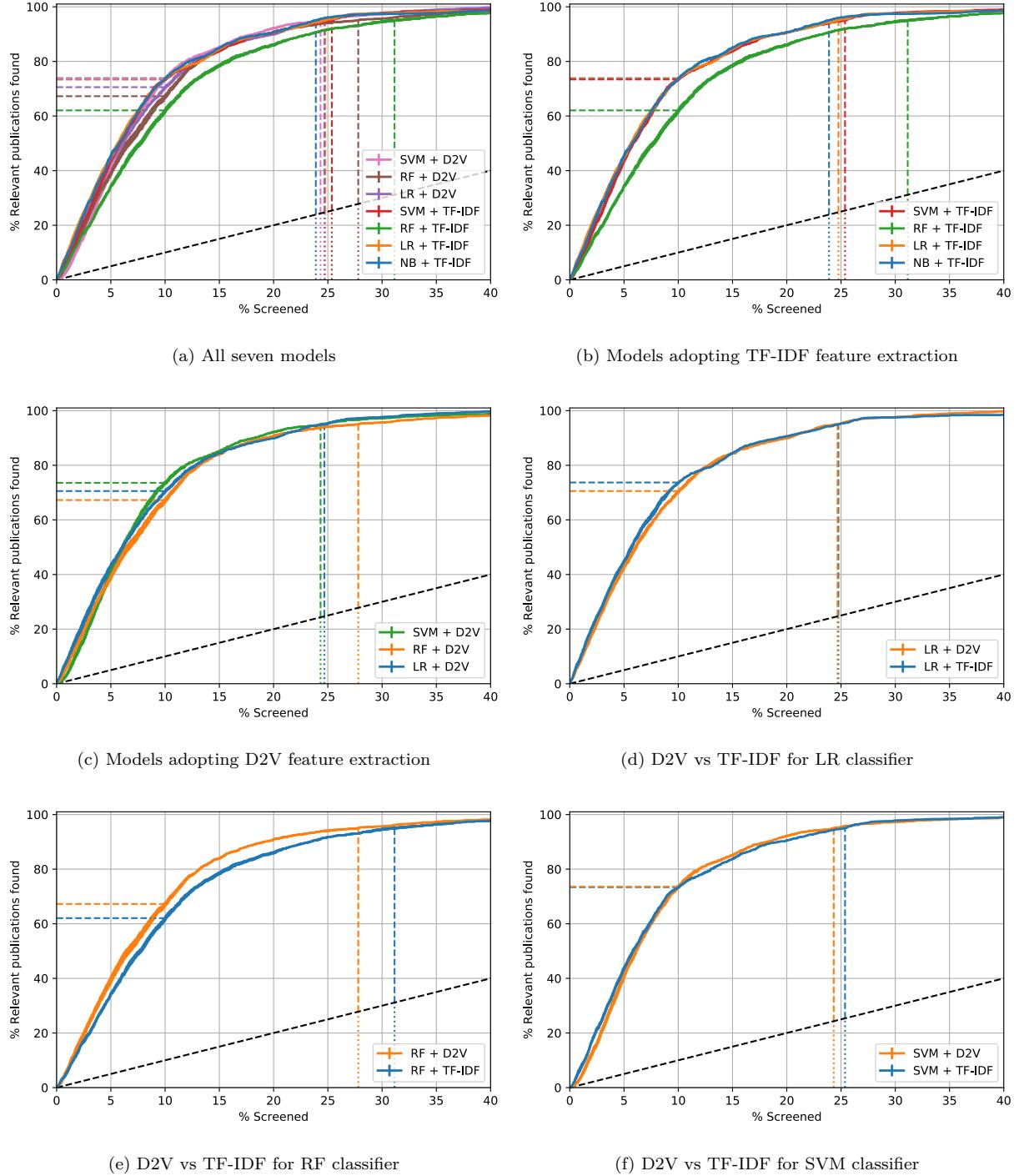


Figure 6: Recall curves for the Virus dataset.

## Wilson

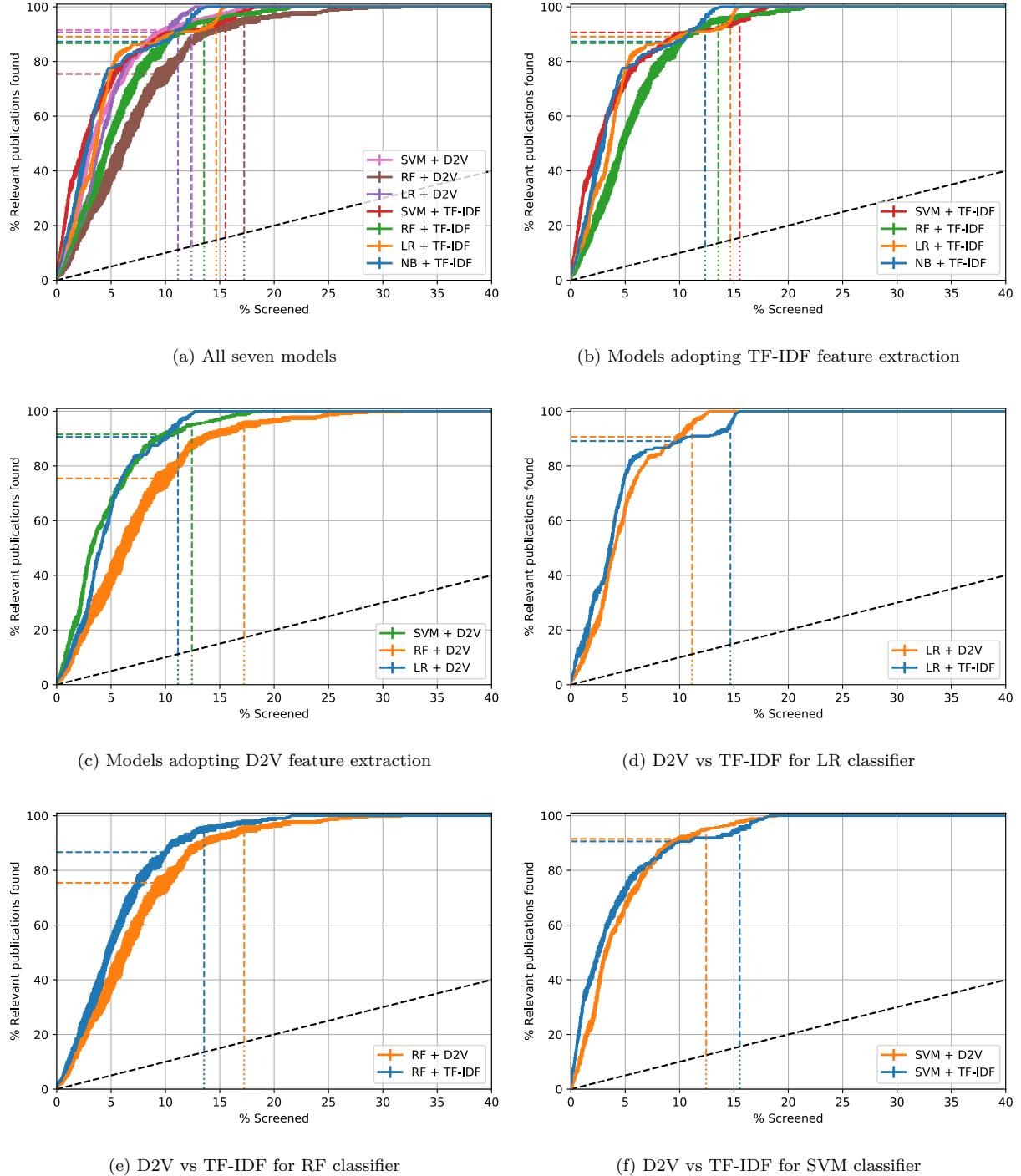


Figure 7: Recall curves for the Wilson dataset.