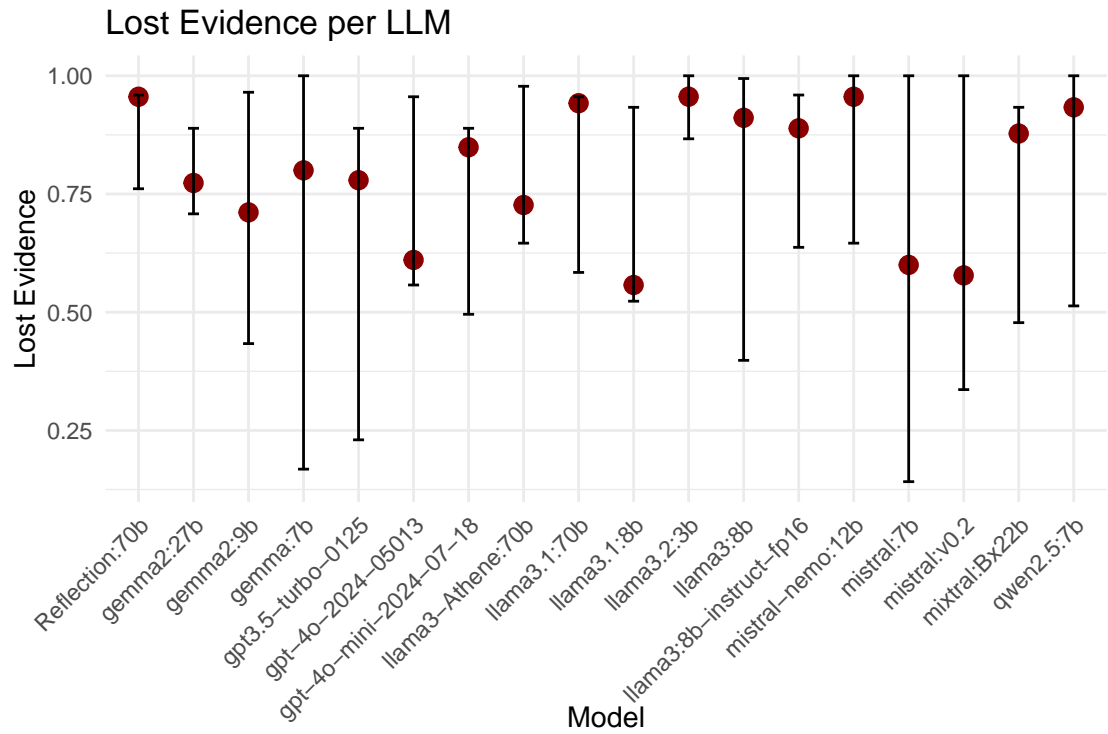# Graphical Abstract

**LLM4SCREENLIT: Recommendations on Assessing the Performance of Large Language Models for Screening Literature in Systematic Reviews**

Lech Madeyski, Barbara Kitchenham, Martin Shepperd



Lost Evidence per LLM

# Highlights

**LLM4SCREENLIT: Recommendations on Assessing the Performance of Large Language Models for Screening Literature in Systematic Reviews**

Lech Madeyski, Barbara Kitchenham, Martin Shepperd

- Use distilled good practices & recommendations for evaluating LLMs in SR screening

- Report confusion matrices enabling (meta-)analyses & alternative metric computation

- Prioritize lost evidence/recall in SR screening evaluations and Weighted MCC (WMCC)

- Use cost-benefit analysis where lost evidence is a critical issue

# LLM4SCREENLIT: Recommendations on Assessing the Performance of Large Language Models for Screening Literature in Systematic Reviews

Lech Madeyski[a,*], Barbara Kitchenham[b], Martin Shepperd[c]

[a]*Wrocław University of Science and Technology, Wyb. Wyspianskiego 27, Wrocław, 50-370, Poland*
[b]*Keele University, UK*
[c]*Brunel University of London, UK*

## Abstract

*Context*: Large language models (LLMs) are released faster than users' ability to evaluate them rigorously. When LLMs underpin research, such as identifying relevant literature for systematic reviews (SRs), robust empirical assessment is essential.
*Objective*: We identify and discuss key challenges in assessing LLM performance for selecting relevant literature, identify good (evaluation) practices, and propose recommendations.
*Method*: Using a recent large-scale study as an example, we identify problems with the use of traditional metrics for assessing the performance of Gen-AI tools for identifying relevant literature in SRs. We analyzed 27 additional papers investigating this issue, extracted the performance metrics, and found both good practices and widespread problems, especially with the use and reporting of performance metrics for SR screening.
*Results*: Major weaknesses included: i) a failure to use metrics that are robust to imbalanced data and do not directly indicate whether results are better than chance, e.g., the use of Accuracy, ii) a failure to consider the impact of lost evidence when making claims concerning workload savings, and iii) pervasive failure to report the full confusion matrix (or performance metrics from which it can be reconstructed) which is essential for future meta-analyses. On the positive side, we extract good (evaluation) practices on which our recommendations for researchers and practitioners, as well as policymakers, are built.
*Conclusions*: SR screening evaluations should prioritize lost evidence/recall alongside chance-anchored and cost-sensitive Weighted MCC (WMCC) metric, report complete confusion matrices, treat unclassifiable outputs as referred-back positives for assessment, adopt leakage-aware designs with non-LLM baselines and open artifacts, and ground conclusions in cost–benefit analysis where FNs carry higher penalties than FPs.

*Keywords:* large language models, LLM, classification metrics, class imbalance, systematic reviews, lost evidence, cost-sensitive

## 1. Introduction

Large language models (LLMs) are increasingly employed to automate the challenging task of paper screening in systematic reviews (SRs), promising substantial reductions in human workload and faster evidence synthesis in software engineering [1, 2, 3] and beyond [4, 5, 6, 7, 8, 9, 10, 11]. To quantify their effectiveness, standard confusion-matrix metrics, e.g., accuracy, precision, recall, specificity, and F1-score, are adopted by comparing model decisions (include/exclude) against human reference labels. Although these metrics offer a familiar evaluation framework, their uncritical application can yield misleading conclusions. In particular, we argue it is essential to consider the features of the problem domain and which metrics best address them. This contrasts with an evaluation using all the 'usual' metrics in the hope that embedded within the results will be useful insights.

---

*Corresponding author

We argue that there are four key features relating to screening papers for SRs. First, the data will tend to be extremely imbalanced so that the negative class (i.e, papers that are irrelevant to the SR) will considerably outnumber the positive class (i.e., relevant papers). Second, the costs of misclassifications are not equal. A relevant paper wrongly excluded, i.e., a false negative (FN), will likely have a far greater impact upon the quality of the SR than an irrelevant paper that passes the screening, i.e., a false positive (FP), and then wastes human effort subsequently rectifying the situation. Third, resources are limited, so we are concerned about the overall costs and benefits of deploying different screening tools. Fourth, we would like to be reassured that sophisticated, yet essentially black-box methods such as LLM screening tools are actually doing better than guessing.

Delgado-Chaves et al. [4] (hereafter referred to as DC+) recently evaluated 18 LLMs for screening studies in three SRs, comparing selections with human reviewers and using confusion matrix metrics. They also emphasize comparing LLMs with more traditional machine learning methods, specifically the random forest method. While DC+ is a key step in assessing LLMs for SR screening, methodological issues (using simple counts of correctly classified papers—referred to as Accuracy in classification studies—as performance metrics, reporting biased metrics, and excluding unclassifiable papers from performance assessment) obscure its outcomes.

We also report the results of reviewing 27 other papers that studied the use of LLMs to screen literature for SLRs. The papers were obtained from the primary studies included in two systematic reviews ([12] and [13]), together with papers we found from informal searches at the beginning of our investigation. We investigated these papers for three purposes: i) to confirm that the problems we observed in the DC+ paper are not unique to that paper, ii) to assess whether any other problems exist, and iii) to investigate whether there are additional good practices to recommend.

## 2. Issue with Confusion Matrix Metrics

This section explains why correctness inadequately measures LLM performance in SRs, highlighting Recall and Lost Evidence as critical metrics for literature screening. Throughout this paper, we refer to the counts from confusion matrices as True Positives (TPs), True Negatives (TNs), False Positives (FPs), and False Negatives (FNs). The formulas used to construct the confusion matrix metrics discussed in this section can be found in the Appendix.

### 2.1. The Fallacy of Correctness

The DC+ study abstract indicates their results favour using LLMs, summarizing their findings as follows:

> "on average, the 18 LLMs classified 4,294 (min 4,130; max 4,329), 1,539 (min 1,449; max 1,574), and 27 (min 22; max 37) of the titles and abstracts correctly as either included or excluded for the three SRs, respectively."

This statement is misleading because the reviews SR-I and SR-II are highly imbalanced (many irrelevant vs. few relevant studies), meaning rejecting all studies would still achieve high scores for correctness and percentage correctness. For example, gemma:7b in Table 1 scores 96.17% for Accuracy but found *none* of the relevant studies (TP=0). Similarly, mistral-nemo:12b scored 96.11% for Accuracy found none of the relevant studies, and identified 3 FPs (i.e., irrelevant studies) as relevant. In contrast, two models that at least identified some of the relevant papers (TP>0) had lower Accuracy values. This means that if we optimise on the Accuracy metric, we would select worse, or in some cases, completely ineffectual LLMs that failed to detect any relevant primary studies.

In addition, Specificity, which is the proportion of all negatives correctly identified, is also misleading for imbalanced data dominated by negatives. As can be seen in Table 1, the two LLMs that did not classify any of the positives correctly had perfect Specificity values.

| Metrics: | Models: gemma:7b | llama3-Athene:70b | llama3.1:8b | mistral-nemo:12b |
|---|---|---|---|---|
| True Negatives (TNs) | 4324 | 4242 | 4048 | 4326 |
| False Negatives (FNs) | 172 | 125 | 90 | 172 |
| True Positives (TPs) | 0 | 47 | 82 | 0 |
| False Positives (FPs) | 0 | 82 | 281 | 3 |
| Total Articles (N*) | 4496 | 4496 | 4501 | 4501 |
| Evidence Lost | 100% | 73% | 52%* | 100% |
| Accuracy | 96.17%* | 95.40% | 91.80% | 96.11% |
| MCC | NaN | 0.29* | 0.29* | -0.005 |
| Weighted MCC (WMCC)** | NaN | 0.40 | 0.48* | -0.014 |
| Precision | NaN | 0.36* | 0.23 | 0.00 |
| Recall | 0.00 | 0.27 | 0.48* | 0.00 |
| Specificity | 1.00* | 0.98 | 0.94 | 1.00* |
| F1 | NaN | 0.31* | 0.31* | NaN |
| Cost | 1720 | 1332 | 1181* | 1723 |

Table 1: Performance metrics for four of the LLMs used in SR-I, revealing problems with using Accuracy and other non-chance adjusted metrics, and ignoring relative costs. NB The asterisks '*' denote the 'best' LLM by metric, i.e., row-wise. The double asterisks '**' denotes that we used weight $w = 10$ to calculate WMCC.

## 2.2. The Critical Importance of Lost Evidence

In SRs, falsely rejecting relevant studies (FNs) loses evidence, potentially irretrievably. Recall, also referred to as Sensitivity (the proportion of relevant studies identified), and Lost Evidence (1-Recall) are therefore fundamental performance metrics. DC+ Figure 1 reveals all reviews—even the balanced SR-III—had problematic Lost Evidence scores. Figure 1 shows that Lost Evidence was serious across all models. Across all SRs, Lost Evidence ranged from 14% (best-case in SR-II) to 100% (worst case in SR-I), with 46 out of 54 LLM classifications missing more than 50% of positive papers. No model performed well on all datasets. The only consistency was llama3.2:3b delivered extremely poor predictions on all three data sets, likely because it had the fewest parameters.

Unlike FNs, FPs (irrelevant studies incorrectly included) only waste effort in subsequent screening. The dominant risk to SR validity comes from missing papers that should be included. However, allowing an extremely large number of FPs, in order to minimize FNs, would mean that there was little value in using the LLM classification. This can be modelled as the FN/FP cost ratio, requiring a subjective assessment of the cost of missing evidence versus the cost saving involved in not processing irrelevant studies. The cost ratio will vary by domain (healthcare SRs may value lost evidence more highly than software engineering SRs) and review type (systematic scoping reviews tolerate missing studies better than formal SRs).

Table 1 uses a plausible 10:1 cost ratio, which indicates that Llama3.1:8b provides better classifications than Llama3-Athene:70b, because although it delivers substantially more FPs, it also delivers more TPs and fewer FNs. In addition, Llama3-Athene:70b fails to classify 5 items.

## 2.3. Reporting Biased and Unsuitable Performance Metrics

DC+ reported analysis using six performance metrics (Precision, Recall, Specificity, F1, MCC, and PABAK) calculated from the confusion matrices, stating that "multiple evaluation metrics offers a comprehensive perspective on their performance and robustness". However, Precision, Recall, Specificity, and F1 are all biased for imbalanced data. While Recall directly relates to Lost Evidence assessment, the value of the other metrics can all be significantly impacted by imbalanced data [14]. PABAK, Prevalence Adjusted Bias Adjusted Kappa, [15, 16] is equivalent to the centred version of the Accuracy metric, so it is not unbiased in any meaningful way.

Of the metrics DC+ deploy, only the Matthews Correlation Coefficient (MCC) [17] is unbiased because it considers all four elements of the confusion matrix without any bias towards TPs or TNs. MCC is an application of the Pearson correlation coefficient, ranging from -1 to 1, with values near zero indicating
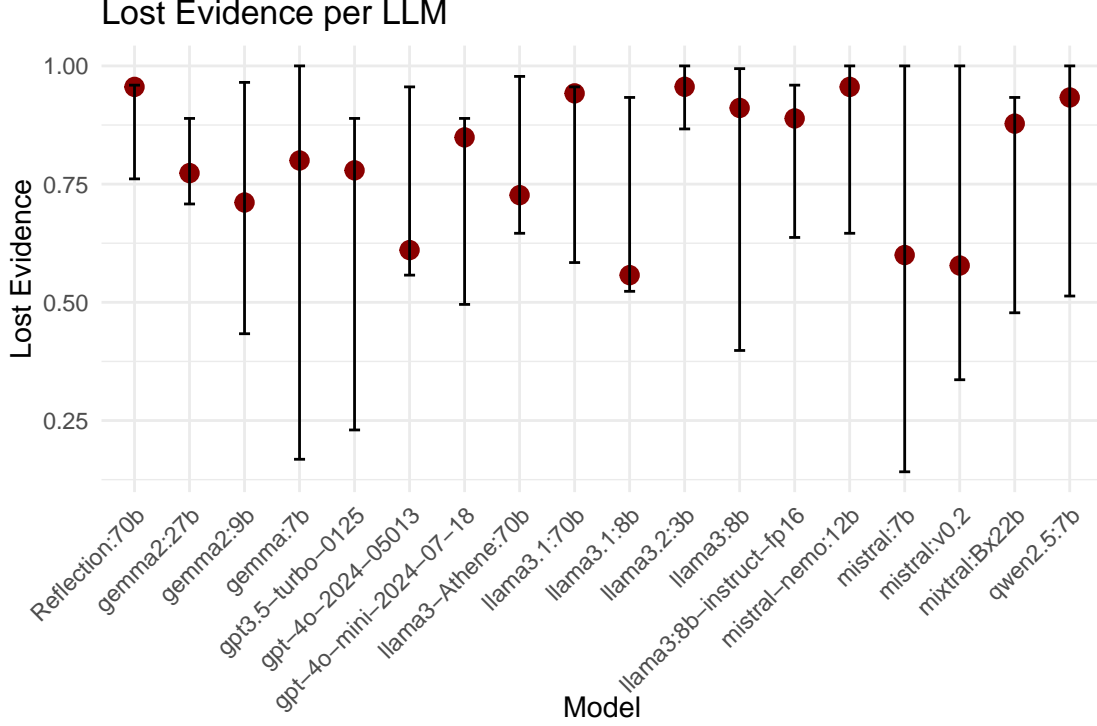
Figure 1: Lost Evidence per Model for three SRs (the median of Lost Evidence is presented as a point) and the min/max show the extremes)

chance-level performance. Like any correlation coefficient, MCC remains reasonably robust to imbalanced datasets [17, 18]. The problem is that MCC does not address the differential costs of FPs and FNs. However, it is possible to construct an appropriate weighted measure of MCC (WMC), which we discuss in Section 3.3.

In addition, reporting multiple performance metrics is inherently misleading because all these metrics are derived from the same four confusion matrix elements, and the elements of the confusion matrix are not themselves independent. This is clear because we only require limited information about the overall classification process in order to generate all four elements. For example, if we know the number of negative papers (N), the number of positive papers (P) as defined by the baseline (gold standard) classification process, together with the number of papers classified as negative by the LLM (n) and the number of those n papers that were TNs, then we can to construct the remaining three elements of the confusion matrix because:

$$FP = N - TN \tag{1}$$

and

$$FN = n - TN \tag{2}$$

and

$$TP = P - FN \tag{3}$$

Thus, a large number of related performance metrics are unhelpful because they are functionally correlated in ways that are more often more difficult to understand than the basic confusion matrix elements.

### 2.4. Dropping Unclassified Papers

DC+ chose to exclude papers that could not be classified from confusion matrices. In real SRs, difficult-to-classify papers typically undergo further screening [19, 20]. To align with standard SR practices, LLMs

should identify unclassifiable papers as *referred-back* to a human reviewer for further assessment ([21]). For the purposes of assessing LLM performance, referred-back papers are a mixture of FPs and TPs (as any referred-back paper will be included in the next screening round, i.e., it will be treated as a positive) and should be classified appropriately in confusion matrices, rather than reducing the total number of classified papers to ignore referred-back papers.

### 2.5. Good Practices

Although we have criticised its choice of performance metrics, the DC+ paper adopts two extremely useful good (evaluation) practices:

**(P1) Reporting full confusion matrices**: The paper provides the full confusion matrices for each LLM and SR in publicly accessible supplementary materials. Access to the full confusion matrices means that other researchers and meta-analysts can easily construct any performance metric of interest in their own context, in particular, the unbiased MCC metric or the weighted MCC metric.

**(P2) Comparing with non-LLM baselines**: The paper also assesses the performance results obtained when using a more traditional classification technique (i.e., the random forest method). This is important because researchers wanting to improve the efficiency of their SR process are not concerned only about the relative effectiveness of different LLMs, but also need to know whether LLMs outperform other techniques that do not share the problems associated with LLM use, such as the risk of hallucinations.

## 3. Current LLM Literature Screening Evaluation Practices

To assess whether the performance metric problems in the DC+ paper were representative of current research practice, we also reviewed another 27 papers addressing literature screening, see Table 2.

Table 2: Summary of the Screening Papers

| ID | Performance Metrics | CM | Origin | Type | Other Baselines |
|---|---|---|---|---|---|
| Akinseloyin-2024 [22] | L-Rel; MAP; Rec; WS | No | Sandner | J | No |
| Attri-2024 [23] | Acc; %Pos; Rec; Spec | No | Kim | A | No |
| Cai-2023 [24] | F1; Prec; Rec; WS | No | Both | J | No |
| Cao-2024 [25] | Acc; Rec; Spec | No | Both; A | GL | No |
| Castillo-2023 [26] | Acc; F1; NegPred; Nulls; Prec; Rec; Spec | Yes | A | C | No |
| Datta-2024 [27] | Acc; F1; Prec; Rec | No | Kim | A | No |
| DC+[4] | Corr; F1; MCC; PABAK; Prec; Rec; Spec | Yes | A | A | Yes |
| Dennstadt-2024 [5] | Acc; F1; Prec; Sens; Spec | Yes | Kim | J | Yes |
| Du-2024 [28] | Acc; F1; Prec; Rec | No | Kim | J | Yes |
| Felizado-2024 [3] | Acc | Yes | A | J | No |
| Gargari-2024 [29] | Acc; F1; Rec; Spec | No | Sandner | L | No |
| Guo-2024 [30] | Acc; F1; Kappa; PABAK; RecExc; RecInc | Yes | Both | J | No |
| Huotala-2024 [1] | %Exc; %Inc; Prec; Rec | No | Kim; A | J | No |
| Issaiy-2024 [31] | BalAcc; Jaccard; Kappa; NegPred; Pos & Neg Likelihood; Prec; PropMissed; Rec; Spec; WS | No | Both | J | No |
| Kaur-2024 [32] | Acc; Rec; Spec | No | Kim | A | No |
| Khraisha-2024 [6] | Acc; Kappa; PABAK; Rec; Spec; Weighted Kappa | No | Kim | J | No |
| Li-2024 [33] | Acc; Rec; Spec | No | Sandner | J | No |
| Lin-2023 [34] | Acc; F1; Prec; Rec; ROC; Spec | No | Kim | J | Yes |
| Rai-2024 [35] | Acc; Prec | No | Kim | A | Yes |
| Robinson-2024 [36] | Acc; Prec; Rec | No | A | GL | Yes |
| Royer-2023 [37] | Rec; Spec | No | Kim | A | Yes |
| Spillias-2024 [38] | Kappa | Yes | Sandner | J | No |

| ID | Performance Metrics | CM | Origin | Type | Other Baselines |
|---|---|---|---|---|---|
| Syriani-2023 [39] | BalAcc; F2; Fleiss' Kappa; MCC; NegPred; Prec; Rec; Spec | No | A | GL | Yes |
| Syriani-2024 [40] | BalAcc; MCC; NegPred; Prec; Rec; Spec | No | A | J | Yes |
| Thode-2025 [2] | Prec; Rec | No | A | J | No |
| Tran-2023 [41] | Rec; Spec; WS | No | Sandner | GL | No |
| Wang-2024 [10] | BalAcc; F3; Prec; Rec; Success Rate; WS | No | Both | C | No |
| Wilkins-2023 [42] | Acc; Kappa; Weighted Kappa; Weighted Rec; Weighted Spec | No | A | GL | No |

Table 2 reports details about the papers included in this review:

1. The papers were assembled from three different sources, which are shown in the columns labelled *Origin*. Papers were obtained from two systematic reviews of papers reporting empirical studies of LLM support for literature selection: Kim et al. [12] included 14 papers reporting 33 separate studies in their meta-analysis of the performance metrics F1, precision and Recall/Sensitivity, and Sander et al. [13] included 11 papers reporting 13 separate studies in their systematic review of Recall/Sensitivity and workload reduction. Five papers were included in both studies (labelled as Both in the Origin column). In addition, we identified 10 other papers (including DC+) from our own informal searches. These are labelled A (for Authors) in the Origin column. Only two of the papers were also identified by Kim et al. [12] or Sander et al. [13].

2. The column labelled *Performance Metrics*, identifies the metrics reported in each paper, excluding simple confusion metric counts. The performance metrics referred to by shortened labels are Acc (Accuracy or Correctness), BalAcc (Balanced Accuracy), Rec (Recall or Sensitivity), Spec (specificity), Prec (Precision), NegPred (Negative Prediction), RecInc (Sensitivity Included), RecExc (Sensitivity Excluded, Pos (Positive), Neg (Negative), Null (Number of null or otherwise invalid outcomes), PropMissed (Proportion Missed). The term *WS* was used to refer to some form of work saved metric irrespective of the specific term used by the authors. In addition, Akinseloyin et al [22] investigated prioritising papers in terms of relevance, rather than classifying them. They used metrics appropriate to that task, but did not fully define them. Their Recall statistics and Work saved statistics were calculated relative to the percentage of prioritised papers evaluated.

3. The six papers that reported confusion metric counts or percentages are identified in the column labelled *CM*.

4. The column labelled *Type* identifies whether the paper was published in a journal (J), conference proceedings (C), was only available as an Abstract (A), was a letter to the editor (L), or was grey literature (GL).

5. The column labelled *Other Baselines* identifies the papers that compared LLM performance with other types of machine learning algorithms such as logistic regression.

Based on the Performance Metrics column,Figure 2 reports a summary of the performance metric usage.

### 3.1. Review Results

Comparing the 27 papers we included in this review with DC+ reveals significant methodological gaps in evaluating Gen-AI tools for systematic review screening, with an underutilization of appropriate performance metrics (see Figure 2) and limited approaches to cost-benefit considerations:

1. **Limited MCC adoption**: Only 3 papers (11% of the sample) employed MCC, including DC+ and the two papers by Syriani et al. [39, 40], despite MCC's advantages for handling imbalanced datasets common in systematic review screening. However, Syriani et al. rescaled MCC from its standard
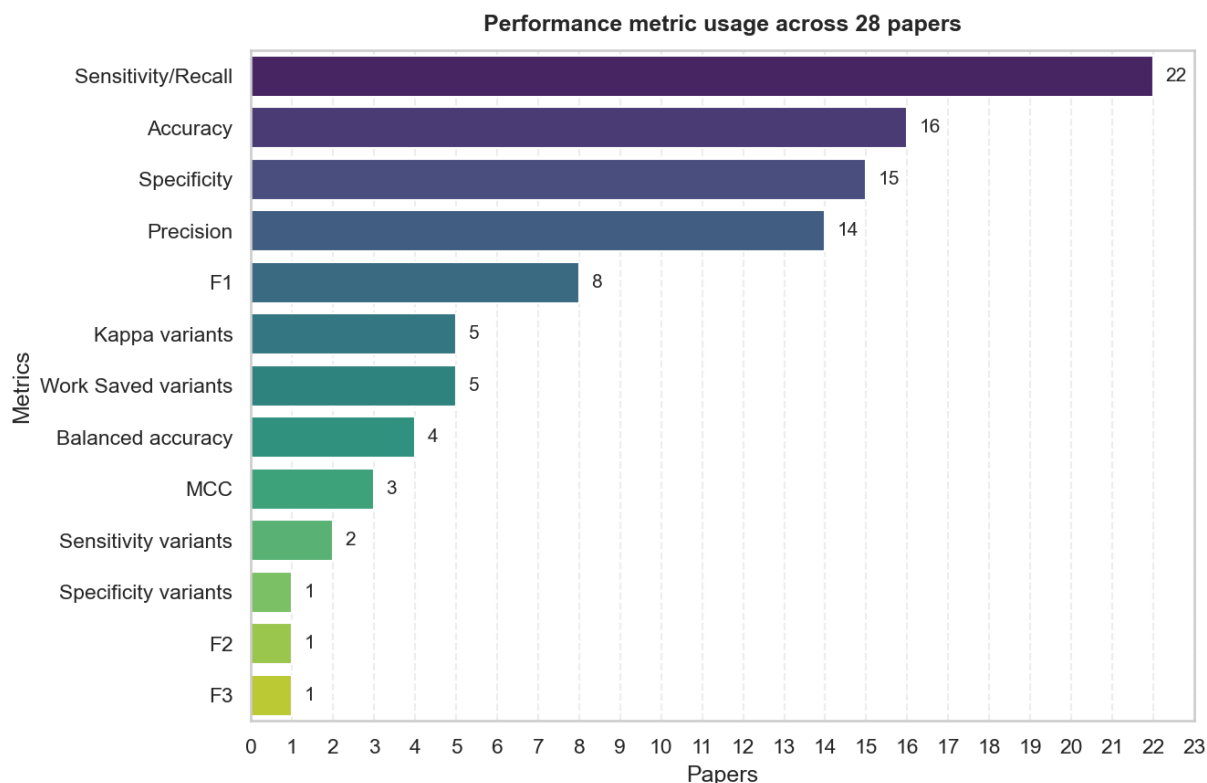
Figure 2: Distribution of evaluation metrics used across 27 papers analyzing Gen-AI tools for systematic review screening

$[-1, 1]$ range to $[0, 1]$, which may complicate cross-study comparisons. Failure to use MCC is a missed opportunity for robust performance evaluation, particularly given the inherent class imbalance in screening tasks.

2. **Insufficient confusion matrix reporting**: Only 6 papers (22%) reported complete confusion matrices, though 4 additional papers provided sufficient information (in terms of total positives, total negatives, sensitivity, and specificity) to enable reconstruction. This deficiency hampers meta-analysis and prevents readers from computing alternative metrics or conducting independent cost-benefit assessments tailored to their specific screening contexts.

3. **Sensitivity/Recall and Accuracy as dominant metrics**:

   - The prevalence of sensitivity/recall (22 papers, 78%) reflects the field's appropriate concern with minimizing FNs (failing to identify relevant studies would undermine SR validity). However, this focus often occurred without complementary metrics to assess the trade-offs with overall efficiency.

   - The widespread use of accuracy (16 papers, 59%) alongside high sensitivity reporting suggests many researchers may not fully appreciate accuracy's limitations for imbalanced datasets.

   - The very modest adoption of balanced accuracy (only 4 papers, 15%) and extremely low MCC usage indicates insufficient awareness of metrics specifically designed for class imbalance scenarios or the need for chance-anchoring.

   - Five papers used Kappa variants (some employing multiple versions), which is appropriate when there is no well-defined baseline. However, it is inappropriate to use standard performance metrics and kappa variants on the same confusion matrix unless you are comparing two potentially fallible

Adoption of good practices by papers

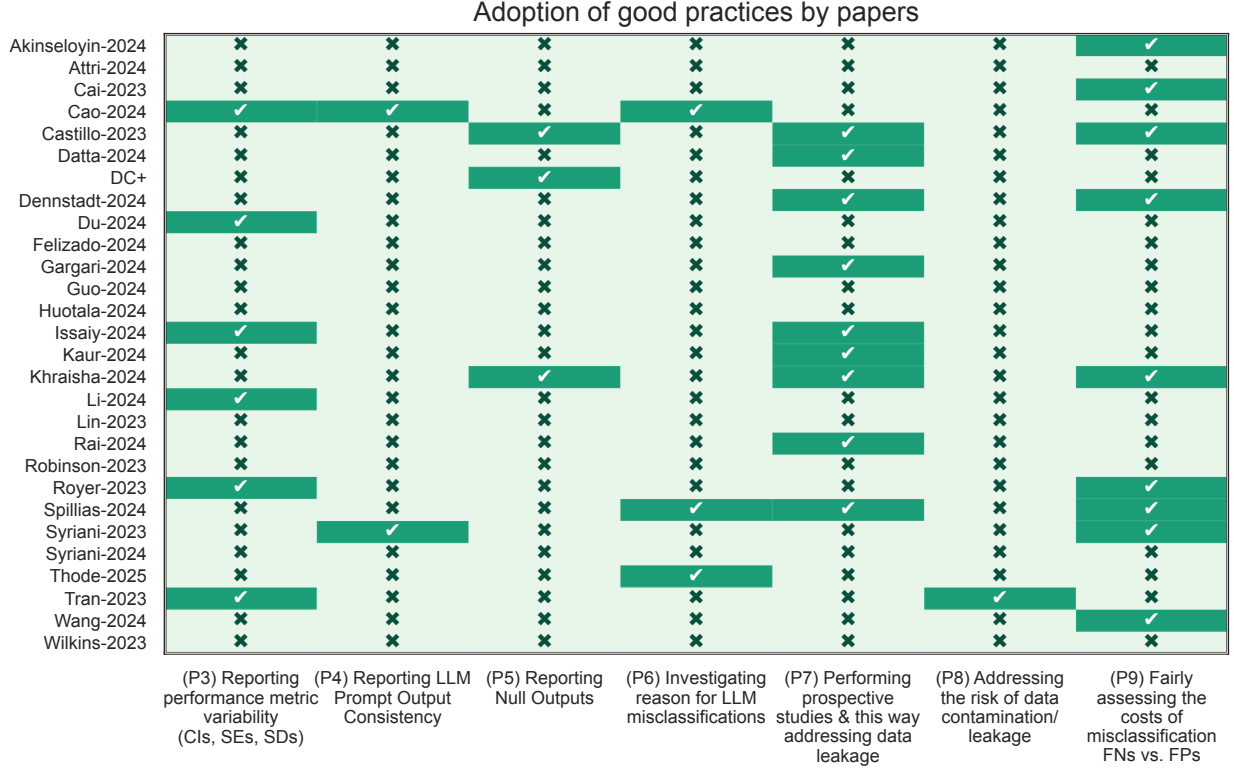| Paper | (P3) Reporting performance metric variability (CIs, SEs, SDs) | (P4) Reporting LLM Prompt Output Consistency | (P5) Reporting Null Outputs | (P6) Investigating reason for LLM misclassifications | (P7) Performing prospective studies & this way addressing data leakage | (P8) Addressing the risk of data contamination/ leakage | (P9) Fairly assessing the costs of misclassification FNs vs. FPs |
|---|---|---|---|---|---|---|---|
| Akinseloyin-2024 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Attri-2024 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Cai-2023 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Cao-2024 | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Castillo-2023 | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ |
| Datta-2024 | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| DC+ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Dennstadt-2024 | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ |
| Du-2024 | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Felizado-2024 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Gargari-2024 | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Guo-2024 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Huotala-2024 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Issaiy-2024 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Kaur-2024 | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Khraisha-2024 | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ |
| Li-2024 | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Lin-2023 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Rai-2024 | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Robinson-2023 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Royer-2023 | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Spillias-2024 | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Syriani-2023 | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Syriani-2024 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Thode-2025 | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Tran-2023 | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ |
| Wang-2024 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Wilkins-2023 | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |

Figure 3: Adoption of good practices

classifications (e.g., classifications made by two LLMs, or a classification made by a single human researcher and an LLM). In addition, the PABAK variant does not seem to be a useful metric in any circumstances.

4. **Failure to address costs as well as benefits:** None of the five papers that explicitly claimed to measure workload savings, nor the Sandner SR [13] that explicitly investigated workload savings, considered the cost of False Negatives when specifying their WS metric.

### 3.2. Additional Good Practices

In addition to two good practices (P1 and P2) already found in DC+ [4] and reported in Section 2.5, we also observed several other good (evaluation) practices which are shown in Figure 3:

**(P3) Reporting performance metric variability**: Six papers reported performance metric confidence intervals, standard errors, or standard deviations. However, there was little consistency in the methods used to calculate them. Two papers used properties of the binomial distribution, one used bootstrapping, and the others did not specify any particular method. Whenever using samples to make decisions about future LLM use, it is important to specify the variability of performance metrics in order to assess both whether performance is better than chance and whether there are significant differences between the performance of different LLMs.

**(P4) Reporting LLM prompt output consistency & (P5) Reporting null outcomes**: Two papers considered the internal consistency of LLMs, and three other papers reported the rate of null outcomes. This is necessary because LLMs can sometimes fail to deliver any output, can provide nonsensical outcomes, or can provide different outcomes for the same prompt. The rate of such problems should be reported, and the researchers should know how to calculate performance metrics fairly when such problems occur.

**(P6) Investigating the reason for LLM misclassifications**: Three papers assessed the reason for LLM classifications and found systematic problems with their prompts. Thus, when validating LLM performance with a view to using an LLM for an SR, analysing classification errors found during LLM validation can provide a means to refine and improve prompts.

**(P7) Performing prospective studies & this way addressing data leakage**: Nine studies were prospective studies (as opposed to retrospective studies, which re-analyse previously published studies). Generally, prospective studies are preferable because there is less chance that researchers have oriented the study goals to suit the available data. In addition, prospective studies do not suffer from the risk of data leakage, also known as data contamination [21].

**(P8) Addressing the risk of data contamination/leakage**: In the context of LLM testing, there is always a risk that publicly available data used to train an LLM could be part of the data used to test LLM performance [21]. This is referred to as data leakage or data contamination. Only Tran et al.[41] reported that their retrospective study used data published after the LLM they studied was released. Dennstadt et al. [5], who used one prospective study and 10 benchmark studies, mentioned the issue, but did not suggest any solution. However, in general, there was little appreciation of the issue. None of the other seven papers that used data from public benchmarks mentioned the issue, and two papers that used retrospective studies that were not obtained from public benchmarks suggested that their data sets could be used as benchmarks.

**(P9) Fairly assessing the costs of misclassification FNs vs. FPs**: Nine papers (30%) explicitly considered differential costs of FNs vs FPs, but with substantial variations in their approaches:

- Khraisha-2024 [6] assigned FNs a weight 30 times greater than FPs, representing the most aggressive cost difference.

- Wang-2024 [10] mandated a minimum 95% recall threshold, prioritizing sensitivity over other metrics.

- Syriani et al. [39] employed F2 scores (weighting recall twice as heavily as precision), although they did not use F2 as a performance metric in their subsequent paper [40]. However, in both papers, their prompts requested Gen-AI systems to be lenient towards inclusions.

- The other six papers mentioned either the critical importance of Recall or the danger of missing evidence, but did not suggest any specific evaluation practices to address the issue.

### 3.3. Addressing Performance Metrics Limitations

In this paper, we have identified the problems that arise when confusion matrices are strongly imbalanced, performance metrics do not consider all elements of the confusion matrix, and do not have a meaningful zero that corresponds to a classifier that is not performing better than chance. We have also noted that the Matthews Correlation Coefficient addresses these issues. In addition, it also allows formal statistical tests to indicate whether or not a given MCC value is better than a random classifier.

However, we have also criticised performance metrics that do not address the cost asymmetry of FNs vs FPs and, MCC, as specified inEquation (15), clearly does not address this issue. Thus, to address cost asymmetry, we propose using a **Weighted Matthews Correlation Coefficient (WMCC)** that builds upon ordinary MCC.

The idea behind WMCC is that it preserves MCC's chance-anchored[1], imbalance-robust correlation meaning and directly addresses the cost asymmetry, although at the cost of losing the opportunity to perform the customary MCC statistical tests of significance. The general WMCC formula is presented in Equation (4):

$$WMCC = \frac{(TP_w * TN_w - FP_w * FN_w)}{\sqrt{(TP_w + FP_w) * (TP_w + FN_w) * (TN_w + FP_w) * (TN_w + FN_w)}} \tag{4}$$

---

[1]Performance metrics are considered chance-anchored if a defined and stable point corresponds to random performance, so for a correlation metric this is zero and for AUC this is 0.5.

Constructing a class-weighted version of MCC, i.e., WMCC, we assign weight $w$ (e.g., $w = 10$) to each positive example, i.e., TP and FN, and weight 1 to each negative example, i.e., TN and FP (when positives are $w$-times more consequential than negatives), compute the weighted confusion matrix counts, and plug them into the standard MCC formula.

Hence, weighted counts are as in Equation (5):

$$TP_w = w \cdot TP, \ \ FN_w = w \cdot FN, \ \ TN_w = 1 \cdot TN, \ \ FP_w = 1 \cdot FP \tag{5}$$

and WMCC can be simplified to the form presented in Equation (6):

$$WMCC = \frac{(w * TP * TN - FP * w * FN)}{\sqrt{(w * TP + FP) * (w * TP + w * FN) * (TN + FP) * (TN + w * FN)}} \tag{6}$$

A simple, working example (for two LLMs, `llama3-Athene:70b` and `llama3.1:8b`, from Table 3) of how to calculate WMCC under class imbalance, with asymmetric costs reflected by a weight of $w = 10$, is presented below.

Raw counts for `llama3-Athene:70b` LLM from Table 1: $TP = 47$, $FN = 125$, $TN = 4242$, $FP = 82$. Assuming $w = 10$, we may calculate WMCC for this LLM:

$$WMCC^{w=10}_{llama3-Athene:70b} = \frac{(10 * 47 * 4242 - 82 * 10 * 125)}{\sqrt{(10 * 47 + 82) * (10 * 47 + 10 * 125) * (4242 + 82) * (4242 + 10 * 125)}}$$
$$= \frac{1891240}{4748341} = 0.398 \tag{7}$$

Raw counts for `llama3.1:8b` LLM from Table 1: $TP = 82$, $FN = 90$, $TN = 4048$, $FP = 281$. Assuming $w = 10$, we may calculate WMCC for this LLM:

$$WMCC^{w=10}_{llama3.1:8b} = \frac{(10 * 82 * 4048 - 281 * 10 * 90)}{\sqrt{(10 * 82 + 281) * (10 * 82 + 10 * 90) * (4048 + 281) * (4048 + 10 * 90)}}$$
$$= \frac{3066460}{6368931} = 0.481 \tag{8}$$

This example shows the impact of weighting MCC. WMCC allows us to distinguish between `llama3-Athene:70b` and `lama3.1:8b`, which both had the same MCC value (see Table 1), indicating that `lama3.1:8b` is a better classifier because $WMCC^{w=10}_{llama3.1:8b} > WMCC^{w=10}_{llama3-Athene:70b}$ due to `lama3.1:8b` having fewer FNs than `llama3-Athene:70b`, although it has substantially more FPs.

### 3.4. Reporting Variation and Confidence Intervals

The set of papers we reviewed included six papers that reported measures of variance or confidence intervals, see Figure 3. However, when we have true and predicted classification for all abstracts related to a specific SR, we have the complete *population* not a sample, so the concept of sampling error is meaningless. Thus, if we report confidence intervals, we need to be very clear to which population and experimental hypotheses they apply.

One situation where CIs are important is when researchers intend to use results of a validation exercise based on a random subset of the available abstracts to decide whether an LLM can be used to analyse the remaining abstracts. In this case, we need to calculate the mean and variance of appropriate performance metrics from the validation sample results. For this purpose, we would suggest using multiple resampling of the validation sample *without replacement*, see the example visualisation in Figure 4 produced by our R simulation script, in which we reused the known performance metrics of the two models (Model A:
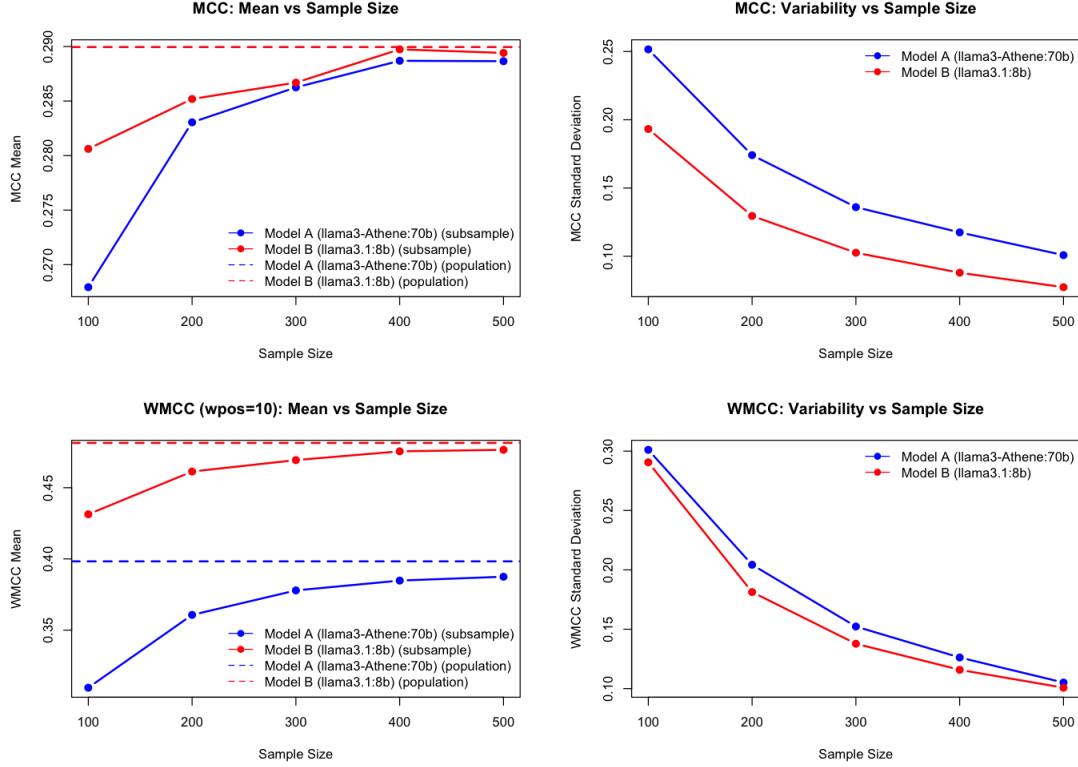
Figure 4: Subsampling stability with 100 to 500 observations

llama3-Athene:70b, and Model B: llama3.1.:8b) reported in Table 1[2]. Such visualisations may help to decide whether LLMs can be used to analyse the remaining abstracts, and which models to choose.

It is also worth mentioning, that we do not recommend CIs based on the binomial distribution because this assumes that all the data items are independently and identically distributed (iid), which is unlikely for a set of abstracts, and, for the same reason, we do we recommend bootstrapping (i.e., sampling with replacement).

Other situations where CIs are useful include:

1. Assessing whether one LLM generally performs better than another. Here we need to assess the variation between appropriate performance metrics for the LLMs across multiple SRs. This is usually the goal of meta-analysis. It requires analysts to ensure that all metrics are derived from different SRs. A similar approach can be used to assess whether different prompt strategies generally perform better than others.

2. Assessing the extent of LLM inconsistency. For such an analysis, we need individual prompts for the same abstract to be repeated.

However, in either case, it is important to design the evaluation such that repeated analyses of the same prompts or the same abstracts in order to assess performance metric variability do not lead to any data leakage.

---

[2]We developed an R script that generated 10000 subsample distributions of different sizes (100, 200, 300, 400, 500) for MCC and WMCC that are in line with the performance metrics reported in Table 1 and visualized the results.

*3.5. Review Limitations*

A limitation of this review is that it is based on a convenience sample of primary studies. We sourced papers from our own knowledge, but in order to address the criticism that we might have explicitly selected papers that support our own opinions, we also included all the primary studies included in two systematic reviews. Furthermore, the problems we found in the DC+ paper were also found in papers from all three sources. Likewise, good practices varied across papers from different sources. It is perhaps surprising how little overlap there was between the three data sources. This would suggest that there are potentially many more papers on the topic than the 28 used in this review.

Another limitation is that the data extraction was performed by one researcher (Kitchenham). However, most of the extracted data were objective and easy to find (or confirm were missing). The most subjective element was deciding whether or not the paper recognized the difference between FPs and FNs in the context of literature screening. Some papers explicitly reported that FNs were more important than FPs, while others reported that Recall was the most important performance metric without explaining why.

## 4. Discussion and Conclusions

In this paper, we have demonstrated the critical challenges associated with evaluating LLMs for literature screening in SRs. We have used the recent evaluation study DC+ [4] as a means of illustrating some of the challenges of choosing appropriate classification metrics. DC+ authors report and discuss thoroughly misleading metrics such as Accuracy, which suggest poorly performing tools are to be preferred. They also omit any sense of cost or the idea that FPs are less problematic than FNs. Notably, across three SR datasets and 18 LLMs in DC+, screening performance was frequently undermined by substantial Lost Evidence (i.e., high FN rates), see Figure 1, demonstrating that correctness/accuracy and related metrics can be seriously misleading under class imbalance and asymmetric misclassification costs. Our review of practices in 27 other papers confirms that these problems are not unique to DC+.

We believe performance metrics problems arise because researchers do not always select metrics that are mapped to the actual needs of the problem domain. A better approach is to report the individual confusion matrix counts together with a small number of performance metrics focussed on the specific research questions being investigated. This means that evaluations of literature screening must prioritize chance-anchored metrics such as MCC and, additionally, explicitly reflect FN vs. FP asymmetry via cost-sensitive analysis. For this purpose, we propose Weighted MCC (WMCC) as a principled extension that retains MCC's correlation meaning and robustness to class imbalance, while addressing the challenge of asymmetric misclassification costs by encoding domain-specific cost ratios.

*4.1. Implications*

There are many implications stemming from this paper for researchers and practitioners, as well as for those responsible for setting journal and conference policies. Researchers should prioritize prospective, leakage-aware benchmarks and standardize WMCC reporting to enable credible conclusions and robust meta-analytic synthesis across SRs. Practitioners should adopt a reporting kit of their LLM validation exercises that includes complete confusion matrices, Lost Evidence, and WMCC with confidence intervals. They could also set in advance a maximum acceptable Lost Evidence (minimum acceptable recall) threshold that reflects the risk tolerance and objectives of the specific review (they would likely be less stringent for software engineering SRs than for high-stakes clinical SRs), so that Lost Evidence from FNs is bounded by design. In practice, this means that if preliminary validation of LLMs suggests that LLM performance cannot meet or exceed acceptable levels for a specific SR task, and there are no clear indications of how prompts could be usefully refined, the task must be performed by human researchers

Journals, conferences, and SR guidelines should require confusion matrices, uncertainty estimates, baseline comparators, explicit leakage/contamination statements, and open artifacts (prompts, seeds, and any materials/artifacts), while discouraging accuracy-focused reporting.

*4.2. Recommendations*

Although, in Sections 2.5 and 3.2, we have discussed a variety of good practices for evaluating LLM screening performance, our main goal was to provide actionable recommendations for (i) Researchers and practitioners, as well as (ii) Policymakers (e.g., journals, conferences, guideline authors) on how to deal with the observed challenges associated with evaluating the performance of LLMs for screening literature for SRs. In addition to focusing our recommendations on specific target audiences, we decided to organize recommendations into decision-centric themes to improve comprehension. As a result, we have the following recommendations organized by target audience and themes:

---

**Target Audience: Researchers and Practitioners**

**Theme: Metrics and cost-sensitive evaluation**

**(R1) Standardize reporting on Lost Evidence (Recall), MCC, and Weighted MCC (WMCC) with explicit justification of FN:FP cost ratios, and report only relevant metrics while avoiding Accuracy/PABAK as primary metrics** (origin: (P9) in Section 2.5 & Section 2.3).

**(R2) Base comparative conclusions on cost-sensitive analyses that reflect asymmetric misclassification costs, using WMCC to combine chance-correction with cost asymmetry and avoiding over-optimizing Recall alone** (origin: (P9) in Section 2.5).

**(R3\*) Predefine acceptable Lost Evidence (minimum Recall) thresholds as guardrails for the review's risk tolerance and objectives, aligned to review type (e.g., SR, Mapping/Scoping Study, Rapid Reviews) and domain (e.g., healthcare, software engineering)** (origin: Section 4.1).

---
*An asterisk '\*' means that the recommendation is optional.

**Theme: Reporting and transparency**

**(R4) Publish complete confusion matrices for every model, dataset, and prompt to enable recomputation of necessary metrics like Lost Evidence, MCC, WMCC, and alternative (e.g., cost-benefit) analyses and future meta-analyses**[a] (origin: (P1) in Section 2.5).

**(R5) For validation tests, report uncertainty for each performance metric via confidence intervals and document the estimation method used. When testing LLMs on a random sample of abstracts to decide whether to deploy them on the remaining abstracts, use resampling without replacement to estimate confidence intervals for likely performance on the full dataset** (origin: (P3) in Section 3.2; resampling method proposed by the authors in Section 3.4 and illustrated in Figure 4).

**(R6) Quantify LLM output consistency and null or invalid outputs, specify the evaluations rule that treats unclassifiable or referred-back items for fair metric computation** (origin: (P4) and (P5) in Section 3.2).

**(R7) Release open artifacts, including prompts, seeds, code, and curated data, as well as include non-LLM baselines for fair comparison and to support open science** (origin: (P2) in Section 2.5).

---
[a]DC+ revealed a consistent problem with lost evidence across three reviews and 18 LLMs. This important result is only clear because they reported all their confusion matrices.

### Theme: Study design and validity

**(R8) Use prospective or temporally safeguarded evaluations and explicitly state contamination/leakage risks and mitigations to prevent training-test overlap** (origin: (P7) & (P8) in Section 3.2).

**(R9) Include non-LLM baselines and assess LLMs (as black-box tools) against conventional methods to ensure practical value and fair benchmarking** (origin: (P2) in Section 2.5).

### Theme: Decision thresholds and operations

**(R9) When observed Lost Evidence exceeds the pre-specified threshold, escalate to human review or adjust prompts/models to maintain SR validity** (origin: Section 4.1).

## Target Audience: Policymakers (Journals, Conferences, Guideline authors)

### Theme: Metrics and cost-sensitive evaluation

**(R1$_{PM}$) Require reporting of Lost Evidence (Recall), MCC, and WMCC with declared FN:FP cost ratios, and discourage accuracy-centric or PABAK-focused reporting as primary evidence** (origin: Section 4.1).

**(R2$_{PM}$) Mandate cost-sensitive evaluation narratives that explain trade-offs between efficiency and Lost Evidence, referencing WMCC or equivalent methods** (origin: Section 4.1).

### Theme: Reporting and transparency

**(R3$_{PM}$) Require complete confusion matrices for all reported metrics to enable recomputation and meta-analytic synthesis** (origin: Section 4.1).

**(R4$_{PM}$) Require disclosure of LLM output consistency and null-output rates, with explicit rules for handling unclassifiable or referred-back items in evaluation** (origin: Sections 2.4 and 4.1).

**(R5$_{PM}$) Require open artifacts (prompts, seeds, code, data, and materials) and baseline comparators to support independent verification and fair comparison** (origin: Section 4.1 and (P2) in Section 2.5).

### Theme: Study design and validity

**(R4$_{PM}$) Require explicit leakage/contamination statements and temporal or provenance safeguards in retrospective or benchmark-based studies** (origin: Section 4.1, (P7) and (P8) in Section 3.2).

**(R6$_{PM}$) Require inclusion of non-LLM baselines for claims about efficiency or effectiveness, and disallow claims based on accuracy-only evidence** (origin: Section 4.1).

> **Theme: Decision thresholds and governance**
>
> **(R7$_{PM}$) Encourage pre-registration of acceptable Lost Evidence (minimum Recall) thresholds and escalation rules as part of protocol submissions, aligned to domain risk** (origin: Section 4.1).

### 4.3. Future research

This paper consolidates nine good practices (see Sections 2.5 and 3.2) observed across the reviewed literature, demonstrates the pitfalls of accuracy-centric reporting under class imbalance, and proposes WMCC to integrate chance-correction with cost asymmetry, thereby turning fragmented results into operational guidance for SR screening, offering a coherent evaluation framework that supports more credible decisions. In spite of this, further research should investigate the effectiveness of combining multiple good practices identified by us (P1-P9) with ones identified by other researchers to keep the evaluation framework as robust as possible.

### CRediT statement

**Lech Madeyski**: Conceptualisation, Data curation, Methodology, Software, Formal analysis, Investigation, Writing – original draft, Writing – review & editing, Visualisation.
**Barbara Kitchenham**: Conceptualisation, Methodology, Validation, Investigation, Writing – review & editing.
**Martin Shepperd**: Conceptualisation, Methodology, Software, Formal analysis, Writing – review & editing.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Data availability

Data will be available on request.

### Appendix

This section reports the formulas used to calculate the metrics discussed in Section 2. The formulas are based on counts obtained from a confusion matrix as shown in Table 3.

|  | Gold Standard True | Gold Standard False | Total |
|---|---|---|---|
| Predicted True | TP | FP | TP+FP |
| Predicted False | FN | TN | FN+TN |
| Total | TP+FN | FP+TN | N |

Table 3: A Confusion Matrix based on the Classifications Assumed to be True and the Classifications produced by the Prediction Model

Accuracy measures the proportion of all items that are correctly classified:

$$Accuracy = \frac{TP + TN}{TN + TP + FN + FP} \tag{9}$$

Recall, which is also referred to as Sensitivity, measures the proportion of all positives correctly classified

$$Recall = \frac{TP}{TP + FN} \tag{10}$$

Precision measures proportion of all items that were classified as positive that were correctly classified.

$$Precision = \frac{TP}{TP + FP} \tag{11}$$

Specificity, which is also referred to as the True Negative rate, measures the proportion of all negatives that were correctly classified:

$$Specificity = \frac{TN}{TN + FP} \tag{12}$$

F1 is a confusion matrix metric designed to assess retrieval from search engine queries, where the number of true negatives (TNs) cannot be counted:

$$F1 = \frac{2 \times TP}{2 \times TP + FP + FN} \tag{13}$$

PABAK, Prevalence Adjusted Bias Adjusted Kappa [15, 16], is defined as:

$$PABAK = 2 \times p_o - 1 \tag{14}$$

where $p_o$ is the observed agreement i.e., the proportion of identical classifications, also known as Accuracy. So if Accuracy=1, PABAK=1, if Accuracy=0, PABAK=-1 and if Accuracy=0.5 PABAK=0. This means that PABAK is simply a centred version of Accuracy, and is just as unreliable as Accuracy for imbalanced datasets.

The Matthews Correlation Coefficient (MMC) is a form of correlation coefficient calculated as:

$$\text{MCC} = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP) \times (TP + FN) \times (TN + FP) \times (TN + FN)}} \tag{15}$$

**References**

[1] A. Huotala, M. Kuutila, P. Ralph, M. Mäntylä, The Promise and Challenges of Using LLMs to Accelerate the Screening Process of Systematic Reviews, in: Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering (EASE'24), ACM, New York, NY, USA, 2024, pp. 262–271.

[2] L. Thode, U. Iftikhar, D. Mendez, Exploring the use of LLMs for the selection phase in systematic literature studies, Information and Software Technology 184 (2025) 107757.

[3] K. R. Felizardo, M. S. Lima, et al., ChatGPT application in Systematic Literature Reviews in Software Engineering: an evaluation of its accuracy to support the selection activity, in: Proceedings of the 18th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM'24), ACM, New York, NY, USA, 2024, pp. 25–36.

[4] F. M. Delgado-Chaves, M. J. Jennings, et al., Transforming literature screening: The emerging role of large language models in systematic reviews, Proceedings of the National Academy of Sciences 122 (2025) e2411962122.

[5] F. Dennstädt, J. Zink, P. M. Putora, J. Hastings, N. Cihoric, Title and abstract screening for literature reviews using large language models: an exploratory study in the biomedical domain., Syst Rev 13 (2024) 158.

[6] Q. Khraisha, S. Put, J. Kappenberg, A. Warraitch, K. Hadfield, Can large language models replace humans in systematic reviews? Evaluating GPT-4's efficacy in screening and extracting data from peer-reviewed and grey literature in multiple languages., Res Synth Methods 15 (2024) 616–626.

[7] F. Trad, R. Yammine, J. Charafeddine, M. Chakhtoura, M. Rahme, G. El-Hajj Fuleihan, A. Chehab, Streamlining systematic reviews with large language models using prompt engineering and retrieval augmented generation, BMC Medical Research Methodology 25 (2025) 130.

[8] R. Sanghera, A. J. Thirunavukarasu, M. El Khoury, J. O'Logbon, Y. Chen, A. Watt, M. Mahmood, H. Butt, G. Nishimura, A. A. S. Soltan, High-performance automated abstract screening with large language model ensembles, Journal of the American Medical Informatics Association 32 (2025) 893–904.

[9] D. Scherbakov, N. Hubig, V. Jansari, A. Bakumenko, L. A. Lenert, The emergence of large language models as tools in literature reviews: a large language model-assisted systematic review, Journal of the American Medical Informatics Association 32 (2025) 1071–1086.

[10] S. Wang, H. Scells, S. Zhuang, M. Potthast, B. Koopman, G. Zuccon, Zero-shot generative large language models for systematic review screening automation, in: N. Goharian, N. Tonellotto, Y. He, A. Lipani, G. McDonald, C. Macdonald, I. Ounis (Eds.), Advances in Information Retrieval, Springer Nature Switzerland, Cham, 2024, pp. 403–420.

[11] T. Oami, Y. Okada, T.-a. Nakada, Performance of a Large Language Model in Screening Citations, JAMA Network Open 7 (2024) e2420496–e2420496.

[12] J. K. Kim, M. Rickard, P. Dangle, N. Batra, M. Chua, A. Khondker, K. Szymanski, R. Misseri, A. Lorenzo, Evaluating large language models for title/abstract screening: A systematic review and meta-analysis & development of new tool, Journal of Medical Artificial Intelligence (2025).

[13] E. Sandner, L. Fontana, K. Kothari, A. Henriques, I. Jakovljevic, A. Simniceanu, A. Wagner, C. Gütl, Evaluating large language models for literature screening: A systematic review of sensitivity and workload reduction, in: Proceedings of the 14th International Conference on Data Science, Technology and Applications - Volume 1: DATA, INSTICC, SciTePress, 2025, pp. 508–517. doi:10.5220/0013562900003967.

[14] D. M. W. Powers, Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation, International Journal of Machine Learning Technology 2 (2011). URL: https://api.semanticscholar.org/CorpusID:3770261.

[15] T. Byrt, J. Bishop, J. B. Carlin, Bias, prevalence and kappa, Journal of Clinical Epidemiology 46 (1993) 423–429.

[16] G. Chen, P. Faris, B. Hemmelgarn, R. L. Walker, H. Quan, Measuring agreement of administrative data with chart data using prevalence unadjusted and adjusted kappa, BMC Medical Research Methodology 9 (2009) 1–8.

[17] B. Matthews, Comparison of the predicted and observed secondary structure of t4 phage lysozyme, Biochimica et Biophysica Acta (BBA)-Protein Structure 405 (1975) 442–451.

[18] A. Luque, A. Carrasco, A. Martín, A. de Las Heras, The impact of class imbalance in classification performance metrics based on the binary confusion matrix, Pattern Recognition 91 (2019) 216–231.

[19] B. Kitchenham, L. Madeyski, D. Budgen, SEGRESS: Software Engineering Guidelines for REporting Secondary Studies, IEEE Transactions on Software Engineering 49 (2023) 1273–1298. doi:TSE.2022.3174092.

[20] B. Kitchenham, D. Budgen, P. Brereton, Evidence-Based Software Engineering and Systematic Reviews, CRC Press, 2016.

[21] T. Woelfle, J. Hirt, et al., Benchmarking human–AI collaboration for common evidence appraisal tools, Journal of Clinical Epidemiology 175 (2024) 111533.

[22] O. Akinseloyin, X. Jiang, V. Palade, A question-answering framework for automated abstract screening using large language models, Journal of the American Medical Informatics Association 31 (2024) 1939–1952.

[23] S. Attri, R. Kaur, B. Singh, P. Rai, Msr57 transforming systematic literature reviews: unleashing the potential of gpt-4: a cutting-edge large language model, to elevate research synthesis, Value in Health 27 (2024) S270.

[24] X. Cai, Y. Geng, Y. Du, B. Westerman, D. Wang, C. Ma, J. J. G. Vallejo, Utilizing chatgpt to select literature for meta-analysis shows workload reduction while maintaining a similar recall level as manual curation, medRxiv (2023) 2023–09.

[25] C. Cao, J. Sang, R. Arora, R. Kloosterman, M. Cecere, J. Gorla, R. Saleh, D. Chen, I. Drennan, B. Teja, et al., Prompting is all you need: LLMs for systematic review screening, medRxiv (2024) 2024–06.

[26] P. Castillo-Segura, C. Alario-Hoyos, C. D. Kloos, C. F. Panadero, Leveraging the potential of generative ai to accelerate systematic literature reviews: an example in the area of educational technology, in: 2023 World Engineering Education Forum-Global Engineering Deans Council (WEEF-GEDC), IEEE, 2023, pp. 1–8.

[27] S. Datta, K. Lee, H. Paek, M. Mojarad, V. Prabhu, J. Zhang, E. Foley, J. Glasgow, C. Liston, Y. Zheng, et al., Msr103 optimizing systematic literature reviews in endometrial cancer: Leveraging ai for real-time article screening and data extraction in clinical trials, Value in Health 27 (2024) S279.

[28] J. Du, E. Soysal, D. Wang, L. He, B. Lin, J. Wang, F. J. Manion, Y. Li, E. Wu, L. Yao, Machine learning models for abstract screening task-a systematic literature review application for health economics and outcome research, BMC Medical Research Methodology 24 (2024) 108.

[29] O. K. Gargari, M. H. Mahmoudi, M. Hajisafarali, R. Samiee, Enhancing title and abstract screening for systematic reviews with gpt-3.5 turbo, BMJ Evidence-based Medicine 29 (2024) 69–70.

[30] E. Guo, M. Gupta, J. Deng, Y.-J. Park, M. Paget, C. Naugler, Automated paper screening for clinical reviews using large language models: data analysis study, Journal of Medical Internet Research 26 (2024) e48996.

[31] M. Issaiy, H. Ghanaati, S. Kolahi, M. Shakiba, A. H. Jalali, D. Zarei, S. Kazemian, M. A. Avanaki, K. Firouznia, Methodological insights into chatgpt's screening performance in systematic reviews, BMC medical research methodology 24 (2024) 78.

[32] R. Kaur, P. Rai, S. Attri, G. Kaur, B. Singh, Msr15 revolutionizing systematic literature reviews: harnessing the power of large language model (gpt-4) for enhanced research synthesis, Value in Health 27 (2024) S262.

[33] M. Li, J. Sun, X. Tan, Evaluating the effectiveness of large language models in abstract screening: a comparative analysis, Systematic Reviews 13 (2024) 219.

[34] Y. Lin, J. Li, H. Xiao, L. Zheng, Y. Xiao, H. Song, J. Fan, D. Xiao, D. Ai, T. Fu, et al., Automatic literature screening using the pajo deep-learning model for clinical practice guidelines, BMC Medical Informatics and Decision Making 23 (2023) 247.

[35] P. Rai, R. Kaur, S. Pandey, S. Attri, G. Kaur, B. Singh, Msr59 advancing systematic literature reviews: The integration of ai-powered nlp models in data collection processes, Value in Health 27 (2024) S270.

[36] A. Robinson, W. Thorne, B. P. Wu, A. Pandor, M. Essat, M. Stevenson, X. Song, Bio-sieve: exploring instruction tuning large language models for systematic review automation, arXiv preprint arXiv:2308.06610 (2023).

[37] J. Royer, E. Wu, R. Ayyagari, S. Parravano, U. Pathare, M. Kisielinska, Msr131 prospects for automation of systemic literature reviews (slrs) with artificial intelligence and natural language processing, Value in Health 26 (2023) S418.

[38] S. Spillias, P. Tuohy, M. Andreotta, R. Annand-Jones, F. Boschetti, C. Cvitanovic, J. Duggan, E. A. Fulton, D. B. Karcher, C. Paris, et al., Human-ai collaboration to identify literature for evidence synthesis, Cell Reports Sustainability 1 (2024).

[39] E. Syriani, I. David, G. Kumar, Assessing the ability of ChatGPT to screen articles for systematic reviews, arXiv preprint arXiv:2307.06464 (2023).

[40] E. Syriani, I. David, G. Kumar, Screening articles for systematic reviews with ChatGPT, Journal of Computer Languages 80 (2024) 101287.

[41] V.-T. Tran, G. Gartlehner, S. Yaacoub, I. Boutron, L. Schwingshackl, J. Stadelmaier, I. Sommer, F. Aboulayeh, S. Afach, J. Meerpohl, et al., Sensitivity, specificity and avoidable workload of using a large language models for title and abstract screening in systematic reviews and meta-analyses, medRxiv (2023) 2023–12.

[42] D. Wilkins, Automated title and abstract screening for scoping reviews using the gpt-4 large language model, arXiv preprint arXiv:2311.07918 (2023).