# Validating Search Query Simulations:
# A Taxonomy of Measures

Andreas Konstantin Kruff[1][0009−0002−8350−154X], Nolwenn
Bernard[1][0009−0007−0565−3210], and Philipp Schaer[1][0000−0002−8817−4632]

TH Köln - University of Applied Sciences, Germany
`firstname.lastname@th-koeln.de`

**Abstract.** Assessing the validity of user simulators when used for the
evaluation of information retrieval systems remains an open question,
constraining their effective use and the reliability of simulation-based re-
sults. To address this issue, we conduct a comprehensive literature review
with a particular focus on methods for the validation of simulated user
queries with regard to real queries. Based on the review, we develop a
taxonomy that structures the current landscape of available measures.
We empirically corroborate the taxonomy by analyzing the relationships
between the different measures applied to four different datasets rep-
resenting diverse search scenarios. Finally, we provide concrete recom-
mendations on which measures or combinations of measures should be
considered when validating user simulation in different contexts. Fur-
thermore, we release a dedicated library with the most commonly used
measures to facilitate future research.

**Keywords:** User simulation · Query simulation · Validation · Informa-
tion retrieval

## 1 Introduction

Validation of user simulators remains one of the main blocking points for their
wide adoption in the field of (interactive) information retrieval [10]. Indeed,
validation contributes to building trust in the simulation and its results. To the
best of our knowledge, there are no standardized methodology and measures to
validate user simulators. Validation is a faceted problem, as previous work show
different interpretations of validation [7, 32], which may explain this absence of a
unified framework. For example, Breuer et al. [9] validate simulation with regards
to retrieval performance, characteristics of simulated search sessions, and term-
based similarity, while Huurnink et al. [22] validate simulation by comparing
system rankings based on performance metrics.

Simulation in information retrieval can focus on different parts of the search
process, for example, query variants generation [1, 2, 9] and clicks [13]. In this
work, we specifically focus on query simulation and its validation. Query simu-
lation aims to produce queries mimicking real user queries that can be used to

---

[*] Equal contribution.

generate synthetic data for training and evaluating information retrieval systems in a scalable and reproducible manner. Despite the importance of validation and the long existence of the field, to the best of our knowledge, a comprehensive overview of validation facets and measures for the query simulation task is still missing. Hence, the choice of validation facets and measures often remains ad-hoc and based on intuition given the specific context of the simulation.

In this work, we provide a comprehensive overview of validation facets and associated measures in the form of a taxonomy based on a literature review. This taxonomy provides a starting point for practitioners to select the facets and measures that are relevant for their simulation approach and the available data. The taxonomy is built using a bottom-up approach, starting from the measures used in previous work, we group them into common facets. The taxonomy is at most four levels deep, providing fine-grained facets divided between two main meta-facets: (1) facets related to the inability to distinguish between simulated and real data, and (2) facets related to performance prediction capabilities of a simulator. To corroborate the identified facets, we investigate the relationships between the different measures in multiple use cases. For this analysis, we look at correlations and mutual information between the measures. It aims to identify potential complementarities and redundancies between the measures, intuitively guiding the selection of measures for future work. Across the datasets, the correlations were predominantly linear and monotonic. Only one dataset showed notable non-linear relationships, based on the applied thresholds. The exploratory factor analysis generally indicated two to three underlying factors, with the measures loading consistently on these factors.

In summary, the contributions of this work are threefold. First, we propose a taxonomy of validation facets and associated measures for query simulation. Second, we provide an analysis of the relationships between the measures, which inherently reflects the relationships between facets. Third, we release a library to automatically compute measures from the taxonomy, facilitating future research in the field.

## 2   Literature Review

This work aims to provide a general overview of previous work on validation approaches applied to query simulation in the context of (interactive) information retrieval. Therefore, we ask the following research questions:

**RQ1:** What are the different facets considered when validating query simulation in (interactive) information retrieval?

**RQ2:** What are the measures associated with these facets?

To answer these questions, we review the literature, as described in Section 2.1, on query simulation with a focus on validation approaches. The findings of this review (Section 2.2) serve as a foundation for our proposed taxonomy of validation facets and measures, which is presented in Section 3.

**Table 1.** Overview of the results from the literature review with digital libraries.

| Source | # Papers | # Relevant | Papers incl. validation |
|---|---|---|---|
| ACM Digital Library | 66 | 12 | [1, 2, 4, 12, 14, 15, 21, 24, 25] |
| IEEE Xplore | 26 | 5 | [11, 27, 31] |
| SpringerLink | 82 | 8 | [8, 9, 16, 20, 22, 30, 34] |
| arXiv | 26 | 2 | [17] |
| **Total** | 200 | 27 | 20 |

### 2.1  Methodology

For this literature review, we select four source databases: ACM Digital Library,[1] IEEE Xplore,[2] Springer Nature Link,[3] and arXiv.[4] Our choice is mainly motivated by the research areas indexed by these sources, i.e., the first two sources index papers from computer science, while the last two sources also index papers in other areas [19]. To have a broad coverage of the literature, we design a search query that captures papers related to query simulation in general. Therefore, we use the following query over all the different fields of the sources: `"query simulation" OR "simulated quer*"`. Note that the wildcard '*' operator is not supported by arXiv, so we adapt the query to: `"query simulation" OR "simulated queries" OR "simulated query"`.

### 2.2  Results

The search query was executed on the four sources on August 1, 2025, and returned a total of 200 papers (including duplicates). To keep the papers relevant to our research questions, we applied a two-step filtering process performed by two of the authors independently; uncertainties were resolved through discussion.[5] First, we filtered the papers based on their title, abstract, and introduction to ensure they were related to query simulation in the context of (interactive) information retrieval. Second, after reading the full text of the 27 remaining papers, we filtered out those that did not include any form of validation of the simulation. At the end of this process, 20 papers are kept for further analysis, as summarized in Table 1. Additionally, we also added 4 papers [18, 28, 33, 35] that were not retrieved by the search query but were known to us and matched the selection criteria. The final set of selected papers is analyzed in the next section to identify the different facets and measures used to validate query simulation.

---

[1] https://dl.acm.org/
[2] https://ieeexplore.ieee.org/
[3] https://link.springer.com/
[4] https://arxiv.org/
[5] The list of annotated papers is made public at: http://bit.ly/4sxDYWT.

## 3   Taxonomy of Validation Facets and Measures

The 24 papers selected are analyzed with regard to the validation facets and associated measures they consider. Following a bottom-up approach, we iteratively group the specific facets into increasingly abstract categories, resulting in higher-level nodes that capture broader validation dimensions. This process leads to the taxonomy of validation facets and measures presented in Figure 1. We observe two main meta-facet in the taxonomy, one focusing on the *indistinguishability* of the simulated data from real user data and the other focusing on the *performance approximation* abilities of the simulation. We argue that the former tends to promote a user-centric perspective for the validation, while the latter is more system-centric. Indeed, we find that *indistinguishability* can be assessed automatically with regards to specific characteristics such as the similarity between the simulated and real queries, e.g., [2, 9, 25, 35] or the structure of the queries, e.g., [31]. Furthermore, it can also be assessed by human evaluation, for example, by asking human judges to distinguish between search sessions with real and simulated queries [18] or to evaluate the semantic similarity of the simulated queries to human queries [1]. On the other hand, *performance approximation* is more concerned with the analysis of the outcome produced by the simulated queries, i.e., are the results obtained similar to the ones obtained with real queries? A common approach to answer this question consists of computing performance metrics, such as traditional information retrieval metrics, for the simulated queries and the real queries with one or more systems; then, comparing the proximity of the results or, in the case of multiple systems, the ranking of the systems based on the simulated queries to the ranking based on real queries corresponding to a tester-based approach [26]. For example, Zerhoudi and Granitzer [34] use mean square logarithmic error to assess the proximity between isoquants of simulated and real queries, while, [21, 22] look at the correlation between mean reciprocal rank-based rankings of the systems based on real and simulated queries. Another approach is to analyze the overlap between the search engine results pages (SERPs) produced for the simulated and real queries. For example, Traub et al. [31] analyse the overlap between the retrieved documents for the simulated and real queries. We acknowledge that additional measures have been used in the reviewed literature, but are not represented in the taxonomy. This choice is motivated by the fact that these measures are either highly domain-specific or have not been used in multiple papers, indicating their limited adoption to multiple application scenarios. It includes the average NewsGuard score [14] and the Earth Mover Distance between Tip-of-the-Tongue linguistic codes [21].

   We also observe that statistical comparison is often used in complement to the other measures. For example, correlation coefficients are often used when comparing rankings of systems based on real and simulated queries, e.g., [8, 9, 21, 30]. Additionally, statistical significance tests can be employed to assess if simulated and real data follow the same distribution (indicating indistinguishability). Statistical comparison can strengthen the validation by showing that the results are not due to chance.

| Indistinguishability | | | Performance Approximation | | |
|---|---|---|---|---|---|
| **Characteristics** | | | **System-result similarity** | | |
| | **Query Similarity** | | | **Traditional IR** | |
| | | **Semantic** | | | Recall@K |
| | | | BERT Score | | Precision@K |
| | | | Cosine Similarity | | F-Measure |
| | | | Wordnet Similarity | | MAP@K |
| | | **Syntactic** | | | RBP |
| | | | Jaccard Similarity | | nDCG@K |
| | **Behavioral similarity** | | | | MRR@K |
| | | Flesch-Kincaid grade level | | | ARP |
| | | Query Intent Distribution | | **Session-based Measures** | |
| | | Lexical Diversity (TTR) | | | sDCG |
| | **Data** | | | | sRBP |
| | | Query length | | | Effort/Effect |
| | | # of terms | | **SERP Overlap** | |
| | | # of named entities | | | SERP comparison |
| | | Distribution of POS Tags | | | Pool properties |
| **Human Assessment** | | | | | SERP Jaccard Similarity |
| | Turing Test | | | | RBO |
| | Topicality | | | | |
| | Semantic similarity | | | | |

**Fig. 1.** Taxonomy of validation facets (in blue) and measures (in white) for query user simulation.

We note that not all the facets and measures proposed in the taxonomy are necessarily applicable to all query simulation approaches. This is particularly the case for the *performance approximation* facet, for which some of the measures require relevance judgments, e.g., normalized discounted cumulative gain or mean average precision. However, we argue that the taxonomy can serve as a starting point to decide which facets and measures can be used to validate a specific query simulation approach. Furthermore, it can easily be extended with additional facets and measures, including some from other fields, such as natural language generation.

## 4   Experimental setup

In order to support the proposed taxonomy of validation facets and measures, we conduct a comprehensive analysis of the relationships between different measures. It allows us to investigate to what extent the various measures complement

each other in their informational value or overlap due to redundant information. The intuition is that measures belonging to the same facet should share similar information, while measures from different facets should capture complementary aspects of the simulation.

This section first describes the method applied for this analysis (Section 4.1) and then presents the datasets used for the experiments (Section 4.2).

### 4.1   Methodology

We propose a methodology in five steps to analyze the relationships between different measures in our taxonomy and test our initial intuition regarding the facets. The first step is optional, depending on the data and resources available. We describe each step as follows:

0. (optional) **Augment data** with search engine results pages (SERPs) given that a retrieval system with the document collection from the original dataset indexed is available. This step allows the computation of performance prediction measures, especially those related to the analysis of overlap between the results obtained with the simulated and real queries (e.g., SERP Jaccard, RBO).
1. **Compute measures** from the taxonomy that are applicable given the available data. Measures could be computed for each pair of simulated and real queries (one-to-one) or for a set of simulated queries corresponding to a real query (one-to-many). In this work, we consider only one-to-one comparisons to ensure consistency across the datasets.
2. **Conduct exploratory factor analysis (EFA)** on the computed measures to identify underlying latent factors. We inspect the factor loadings to see if measures that load highly on the same factor correspond to the same facet in the taxonomy.
3. **Analyze correlation matrices** of the computed measures to identify linear and monotonic relationships. We compute Pearson $\rho$ and Kendall's $\tau$ correlation coefficients between all pairs of measures. The former captures linear relationships, while the latter captures monotonic relationships. We look at the strength and direction of the correlations to identify clusters of measures that share similar information and compare them to the facets in the taxonomy.
4. **Analyze mutual information** between all pairs of measures to identify potential nonlinear dependencies. We compute normalized mutual information (NMI) to quantify the amount of shared information between two measures, regardless of the nature of their relationship. We look for pairs of measures with high NMI but low Pearson and Kendall correlations, as they may capture complementary aspects of the simulation not reflected in linear or monotonic relationships.

Based on the results from steps 2–4, we consider two criteria to support our proposed taxonomy. First, the clusters of measures identified by EFA should

**Table 2.** Datasets overview.

| Dataset | Queries | Simulations | Document collection | Qrels |
|---------|---------|-------------|---------------------|-------|
| Sim4IA 2025 [23] | 35 | $31 + 19$ | CORE$^{\dagger}$ | No |
| UQV100 [2, 5] | 100 | 3 | ClueWeb12-B$^{\ddagger}$ | Yes |
| UQV subset [9] | 50 | 21 | Common Core 2017 [3] | Yes |
| DL seed queries [1] | 126 | 19 | MS MARCO Passage v2 [6] | Yes |

$^{\dagger}$ https://core.ac.uk/

$^{\ddagger}$ https://www.lemurproject.org/clueweb12.php/

mainly align with the facets in the taxonomy. Second, measures within the same facet should demonstrate high correlation and mutual information, while measures from different facets should show lower correlation and mutual information, indicating that they capture complementary aspects of the simulation.

### 4.2 Datasets

We applied our proposed methodology on four public datasets (see Table 2):

– **Sim4IA 2025**: The Sim4IA dataset originates from the Sim4IA 2025 Micro-Shared Task Workshop [29]. It is based on original search sessions derived from CORE log files and provides one query per session as the gold standard for evaluation. Participants are instructed to generate a ranked list of 10 candidate queries for prediction (i.e., Task A). Consequently, the analysis can be carried out either one-to-one based on the ranking or one-to-many, averaged across all candidate queries per session.

– **UQV100**: The original UQV100 collection comprises 100 topics, each accompanied by a backstory and associated queries. Here, GPT-3.5 is used to generate query variants with temperatures of 0.0, 0.5, and 1. Since the number of variants produced per topic varied, the first generated variant is consistently selected to avoid skewing the evaluation measures [2].

– **UQV subset**: This dataset contains simulated query variants using the UQV collection from TREC Common Core 2017. It involves 21 simulators, each of which generated ten query variants, that are grouped into two methods, namely the TREC Topic Searcher (TTS) and the Known Item Searcher (KIS). In addition, eight modification strategies are applied that vary in the degree of term variation and selection [9]. Furthermore, eight human annotators produced up to ten variants per original query. This setup enables multiple one-to-one and one-to-many comparisons between the human annotators and the simulators. Neither the human user queries nor the simulated queries are explicitly ranked; however, we implicitly assume that the first generated variant corresponds to the most natural reformulation.

– **DL seed queries**: This dataset is based on seed queries from the Deep Learning Tracks 2021 with 53 queries and 2022 with 73 queries. Query variants were generated with GPT-4 using a temperature of 1. The variants

are created with different personas, user groups, and textual transformation strategies, as well as in a neutral setting without any persona. For each configuration, three variants are generated, but they are also not explicitly ranked [1].

The Sim4IA dataset is unique as there are no comparable datasets available that focus on next query prediction based on real user sessions using simulation. For further evaluation, we therefore rely on datasets derived from the UQV collections to apply and assess various query variant strategies. Although the underlying tasks differ, as contextual information for the predictions is missing, we argue that this should affect only the absolute quality of the results, but not the correlations between the measures. Due to the given topic descriptions or backstories, the simulators have different yet related contexts for the query simulation. While the underlying test collection for the Sim4IA dataset unfortunately does not include relevance judgments, which prevents the calculation of traditional IR measures, the three other collections do provide such judgments.

## 5   Analysis

In this section, we summarize the findings from the exploratory factor analysis (EFA), the correlation matrices analysis, and mutual information (i.e., steps 2–4 of the method). These different analyses reveal consistent patterns across the different datasets.

The EFA consistently identifies three to four main dimensions underlying query evaluation. Classical information retrieval (IR) performance metrics (i.e., nDCG, Precision, Recall, MAP, and MRR) are loaded together, confirming that they capture highly overlapping aspects of retrieval effectiveness and are largely redundant. Query similarity measures (i.e., BERT Score, Jaccard, Cosine, and WordNet similarities) form a separate cluster, indicating a complementary perspective on conceptual overlap between the simulated and the real queries. SERP overlap measures (i.e., SERP Jaccard similarity and RBO) tend to either form a distinct factor or cluster together with query similarity measures. This reflects their role in capturing similarity in result rankings, which is related to but not fully determined by query-level similarities. Behavioral and descriptive query features (i.e, query length, number of unique terms, Fleisch-Kincaid grade level, lexical diversity, number of named entities) generally constitute a distinct, loosely defined factor.

The analysis of the correlation matrices across the datasets aligns with these findings. Pearson correlation matrices indicate strong linear relationships among classical IR metrics ($\bar{\rho} = 0.77$) and basic query characteristics ($\bar{\rho} = 0.91$). The two SERP Overlap measures are also highly correlated across datasets ($\bar{\rho} = 0.85$). Query similarity measures exhibit a moderate internal correlation ($\bar{\rho} = 0.47$) and a similarly moderate correlation with SERP Overlap measures ($\bar{\rho} = 0.41$). The average correlation coefficients were obtained by first averaging over all measure pairs and then across the four datasets. Figure 2 shows a representative heatmap of the observed patterns, with the BERT score as a notable
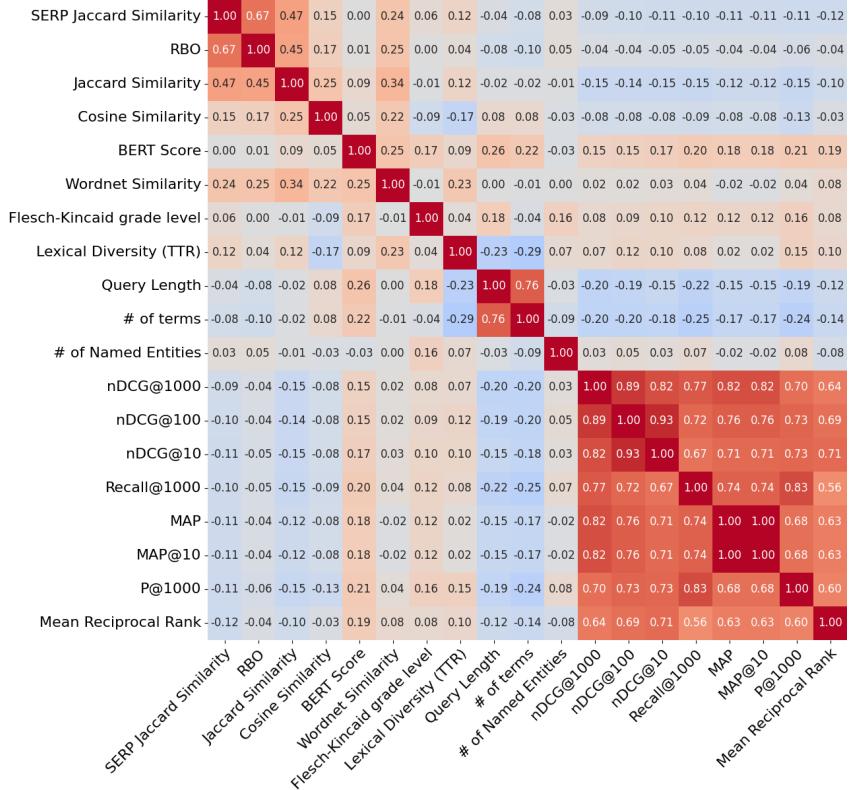
**Fig. 2.** Pearson correlation matrix for the DL 2021 seed queries.

exception. Kendall's $\tau$ correlation matrices largely mirror these patterns, though associations are slightly weaker. Mutual information analysis indicates no relevant nonlinear dependencies for most datasets. In the Sim4IA dataset, limited nonlinear associations are observed, mainly involving semantic similarity, lexical diversity, and Fleisch-Kincaid grade level.

Overall, the analysis indicates that while classical IR measures are highly redundant, query similarity measures and SERP overlap measures provide additional, complementary dimensions for evaluating simulated query validity. Nonlinear relationships appear to be rare and dataset-specific, reinforcing the dominance of linear and monotonic associations.

Unsurprisingly, the analysis of the measures across the four datasets reveals a consistently high correlation among traditional information retrieval (IR) measures. These measures are also the most frequently used in the studies included in our literature review, with many studies employing multiple traditional IR measures within the same study. This indicates that researchers place considerable value on system-based query performance similarity. However, our analysis

suggests that using multiple traditional IR measures in this way does not provide substantial additional insights for evaluation purposes.

In contrast, SERP overlap measures, which are also categorized as "performance approximation" measures and reflect system-side evaluation, offer complementary information. Their predictive power is not redundant with that of traditional IR measures, and they tend to correlate more with query similarity measures. Although both clusters evaluate queries from a system perspective, we suggest using them together, as they capture different aspects of system performance. Furthermore, the partial correlation observed between SERP overlap and query similarity measures provides additional information about string-based query similarity.

Similarly, query similarity and SERP overlap measures show moderate correlation depending on the dataset. Nevertheless, these measures complement each other by capturing distinct dimensions of system ranking and query similarity. Within query similarity measures themselves, moderate intercorrelation exists, yet each measure carries unique information. Depending on the evaluation task, one might prioritize lexical similarity, semantic similarity, or a combination of both to achieve a more comprehensive assessment of query similarity.

For behavioral and descriptive data statistic measures, no consistent correlations could be observed across the datasets. Even within the clusters, correlations are generally negligible, except for query length and query number of terms, which show high correlation, as might be expected. This indicates that these measures capture distinct aspects of the query and user behavior that are not reflected in other clusters.

We acknowledge that restricting the analysis to a one-to-one comparison based on the top-1 simulated query per original query could potentially bias the observed correlations. To assess the robustness of our findings, we applied a bootstrapping procedure over 1,000 iterations, randomly selecting simulated queries referring to the same topic within the same simulator and across different simulators for each iteration. The results show that both Pearson and Kendall correlations exhibit only minor deviations, with maximum absolute differences of approximately 0.15 and the majority of deviations being substantially lower. These findings indicate that the choice of the specific simulated query within a topic has a negligible impact on the overall correlation analysis. Notably, measures based on query statistics often show relatively higher deviations in their correlations with other measures. This can be explained by the fact that semantic similarity measures (e.g., cosine similarity) are not necessarily related to query length, so depending on which simulated query is selected, correlations between length-dependent and semantic measures can fluctuate, even if the queries are semantically similar.

## 6   Conclusion

In this paper, we investigate the methods and measures used to validate user query simulation approaches. To provide an overview of the current landscape

of validation facets and measures, we propose a taxonomy based on a literature review and corroborate it with an empirical analysis of the relationships between the measures. All measures computed in this analysis are bundled into a complementary software library, allowing for reproducibility and further experimentation on the matter of query simulation validation.

While the taxonomy offers a practical conceptual framework, our analysis shows that the actual relationships are more nuanced and context-dependent. In particular, traditional IR performance metrics are found to be highly redundant, capturing largely overlapping aspects of retrieval effectiveness. In contrast, semantic similarity measures provide complementary information about conceptual overlap, and SERP-based overlap measures often capture unique patterns not reflected by either traditional or semantic metrics.

We acknowledge that some measures included in the current taxonomy could not be evaluated in this study. Indeed, each measure has specific data requirements, and not all datasets provide the necessary data. For example, some measures require relevance judgments that are not available in all datasets, while others need human assessments that are costly and difficult to reproduce.

Although we were unable to compare all available measures included in the taxonomy, we argue that the validation of user query simulation should consider different measures, preferably from both user- and system-centric facets, such as query similarity and traditional IR measures, to obtain a comprehensive assessment. A more fine-grained analysis or complementary user studies focusing on human assessment methods, such as Turing tests, are left as future work. Moreover, investigating if our findings generalize to other tasks, adjacent to query simulation, such as (next) utterance prediction (e.g., Sim4IA 2025 Task B), and how these or alternative measures could be employed to assess similarity across entire interaction sessions is another interesting direction for future work. Finally, the taxonomy could be extended with additional facets and measures, such as query performance prediction measures that propose an alternative way to predict query performance without requiring relevance judgments.

**Disclosure of Interests.** The authors have no competing interests to declare that are relevant to the content of this article.

## Appendix

### Software Library

We implement a Python library to compute the measures listed in the taxonomy and support one-to-one and one-to-many comparisons.[6] It provides a unified

---

[6] https://github.com/irgroup/query_sim_validation

framework for validating query simulation approaches, and was used to compute the comparisons in Section 5. The library includes implementation the following measures:

– Basic query statistics like length, number of terms, or named entities
– Flesch-Kincaid grade scores
– Type-Token Ratio (TTR)
– Jaccard similarity
– Cosine similarity for different embedding models
– BERT score
– WordNet-based similarity
– Various retrieval performance metrics
– SERP overlap based on Jaccard index and RBO

The library compares both original and simulated sessions, encoded in JSON. We introduce a session data model that includes a session ID, an ID, interactions, and, optionally, a rank for a simulated session. Note that original and simulated sessions must contain matching session IDs for comparison. An interaction comprises a query, a search engine result page, and, if available, clicked document IDs.

We provide a script to compute some of the measures implemented as an example of how to operate the library. Upon execution, the results are provided in JSONL format, where each line corresponds to one simulator and contains, for each measure, a list of all calculated values in ranking order. These results can then be further processed for descriptive statistics or statistical testing.

## Bibliography

[1] Alaofi, M., Ferro, N., Thomas, P., Scholer, F., Sanderson, M.: Demographically-inspired query variants using an LLM. In: Proceedings of the 2025 International ACM SIGIR Conference on Innovative Concepts and Theories in Information Retrieval (ICTIR), pp. 390–400, ICTIR '25 (2025)

[2] Alaofi, M., Gallagher, L., Sanderson, M., Scholer, F., Thomas, P.: Can generative LLMs create query variants for test collections? An exploratory study. In: Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 1869–1873, SIGIR '23 (2023)

[3] Allan, J., Harman, D., Kanoulas, E., Li, D., Van Gysel, C., Voorhees, E.M.: TREC 2017 common core track overview. In: TREC, TREC '17 (2017)

[4] Azzopardi, L., de Rijke, M., Balog, K.: Building simulated queries for known-item topics: an analysis using six european languages. In: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 455–462, SIGIR '07 (2007)

[5] Bailey, P., Moffat, A., Scholer, F., Thomas, P.: UQV100: A test collection with query variability. In: Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 725–728, SIGIR '16 (2016)

[6] Bajaj, P., Campos, D., Craswell, N., Deng, L., Gao, J., Liu, X., Majumder, R., McNamara, A., Mitra, B., Nguyen, T., Rosenberg, M., Song, X., Stoica, A., Tiwary, S., Wang, T.: MS MARCO: A human generated machine reading comprehension dataset. cs.CL/1611.09268 (2018)

[7] Balog, K., Zhai, C.: User simulation for evaluating information access systems. Foundations and Trends in Information Retrieval **18**(1-2), 1–261 (2024), ISSN 1554-0669

[8] Berendsen, R., Tsagkias, M., de Rijke, M., Meij, E.: Generating pseudo test collections for learning to rank scientific articles. In: Information Access Evaluation. Multilinguality, Multimodality, and Visual Analytics, pp. 42–53, CLEF '12 (2012)

[9] Breuer, T., Fuhr, N., Schaer, P.: Validating simulations of user query variants. In: Advances in Information Retrieval, pp. 80–94, ECIR '22 (2022)

[10] Breuer, T., Kreutz, C.K., Fuhr, N., Balog, K., Schaer, P., Bernard, N., Frommholz, I., Gohsen, M., Ji, K., Jones, G.J.F., Keller, J., Liu, J., Mladenov, M., Pasi, G., Trippas, J., Wang, X., Zerhoudi, S., Zhai, C.: Report on the 1st workshop on simulations for information access (Sim4IA 2024) at SIGIR 2024. SIGIR Forum **58**(2), 1–14 (2025)

[11] Cai, J., Shao, X., Ma, W.: Ontology driven semantic search over structure p2p network. In: 2009 Ninth International Conference on Hybrid Intelligent Systems, pp. 29–34, HIS '09 (2009)

[12] Carterette, B., Bah, A., Zengin, M.: Dynamic test collections for retrieval evaluation. In: Proceedings of the 2015 International Conference on The Theory of Information Retrieval, pp. 91–100, ICTIR '15 (2015)

[13] Chuklin, A., Markov, I., de Rijke, M.: Click Models for Web Search. Synthesis Lectures on Information Concepts, Retrieval, and Services, Springer Cham (2015)

[14] Elsweiler, D., Ateia, S., Bink, M., Donabauer, G., Fernández Pichel, M., Frummet, A., Kruschwitz, U., Losada, D.E., Ludwig, B., Meyer, S., Pascual Presa, N.: Query smarter, trust better? exploring search behaviours for verifying news accuracy. In: Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 515–526, SIGIR '25 (2025)

[15] Elsweiler, D., Losada, D.E., Toucedo, J.C., Fernandez, R.T.: Seeding simulated queries with user-study data for personal search evaluation. In: Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 25–34, SIGIR '11 (2011)

[16] Engelmann, B., Breuer, T., Friese, J.I., Schaer, P., Fuhr, N.: Context-driven interactive query simulations based on generative large language models. In: Advances in Information Retrieval, pp. 173–188, ECIR '24 (2024)

[17] Erbacher, P., Denoyer, L., Soulier, L.: Interactive query clarification and refinement via user simulation. In: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 2420–2425, SIGIR '22 (2022)

[18] Günther, S., Hagen, M.: Assessing query suggestions for search session simulation. In: Causality in Search and Recommendation (CSR) and Simulation

of Information Retrieval Evaluation (Sim4IR) workshops at SIGIR 2021, CSR-Sim4IR '21 (2021)

[19] Gusenbauer, M., Haddaway, N.R.: Which academic search systems are suitable for systematic reviews or meta-analyses? evaluating retrieval qualities of google scholar, pubmed, and 26 other resources. Research Synthesis Methods **11**(2), 181–217 (2020)

[20] He, B., Ounis, I.: Term frequency normalisation tuning for bm25 and dfr models. In: Advances in Information Retrieval, pp. 200–214, ECIR '05 (2005)

[21] He, Y., Kim, T.E., Diaz, F., Arguello, J., Mitra, B.: Tip of the Tongue query elicitation for simulated evaluation. In: Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 3398–3407, SIGIR '25 (2025)

[22] Huurnink, B., Hofmann, K., de Rijke, M., Bron, M.: Validating query simulators: An experiment using commercial searches and purchases. In: Multilingual and Multimodal Information Access Evaluation, pp. 40–51, CLEF '10 (2010)

[23] Kruff, A.K., Kreutz, C.K., Breuer, T., Schaer, P., Balog, K.: Sim4ia-bench: A user simulation benchmark suite for next query and utterance prediction. In: Advances in Information Retrieval, ECIR '26 (2026)

[24] Labhishetty, S., Zhai, C.: An exploration of tester-based evaluation of user simulators for comparing interactive retrieval systems. In: Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 1598–1602, SIGIR '21 (2021)

[25] Labhishetty, S., Zhai, C.: PRE: A precision-recall-effort optimization framework for query simulation. In: Proceedings of the 2022 ACM SIGIR International Conference on Theory of Information Retrieval, pp. 51–60, ICTIR '22 (2022)

[26] Labhishetty, S., Zhai, C.: RATE: A reliability-aware tester-based evaluation framework of user simulators. In: Advances in Information Retrieval, pp. 336–350, ECIR '22 (2022)

[27] Morrison, D., Marchand-Maillet, S., Bruno, É.: Query log simulation for long-term learning in image retrieval. In: 2011 9th International Workshop on Content-Based Multimedia Indexing, pp. 55–60, CBMI '11 (2011)

[28] Rahmani, H.A., Ramineni, V., Yilmaz, E., Craswell, N., Mitra, B.: Towards understanding bias in synthetic data for evaluation. In: Proceedings of the 34th ACM International Conference on Information and Knowledge Management, pp. 5166–5170, CIKM '25 (2025)

[29] Schaer, P., Kreutz, C.K., Balog, K., Breuer, T., Kruff, A.K.: Second SIGIR workshop on simulations for information access (Sim4IA 2025). In: Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 4172–4175, SIGIR '25 (2025)

[30] Sinha, A., Mall, P.R., Roy, D.: Exploring the nexus between retrievability and query generation strategies. In: Advances in Information Retrieval, pp. 177–193, ECIR '24 (2024)

[31] Traub, M.C., Samar, T., van Ossenbruggen, J., He, J., de Vries, A., Hardman, L.: Querylog-based assessment of retrievability bias in a large newspaper corpus. In: Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries, pp. 7–16, JCDL '16 (2016)

[32] Zeigler, B.P., Muzy, A., Kofman, E.: Theory of Modeling and Simulation: Discrete Event & Iterative System Computational Foundations, Third Edition. Elsevier (2019)

[33] Zendel, O., Al Lawati, S.F.D., Rashidi, L., Scholer, F., Sanderson, M.: A comparative analysis of linguistic and retrieval diversity in llm-generated search queries. In: Proceedings of the 34th ACM International Conference on Information and Knowledge Management, pp. 4014–4023, CIKM '25 (2025)

[34] Zerhoudi, S., Granitzer, M.: Simulating user querying behavior using embedding space alignment. In: Linking Theory and Practice of Digital Libraries, pp. 386–394, TPDL '22 (2022)

[35] Zhang, E., Wang, X., Gong, P., Yang, Z., Mao, J.: Exploring human-like thinking in search simulations with large language models. In: Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 2669–2673, SIGIR '25 (2025)