Maria Maistro^a, Timo Breuer^b, Philipp Schaer^b and Nicola Ferro^c

^aUniversity of Copenhagen, Denmark ^bTH Köln - University of Applied Sciences, Germany ^cUniversity of Padua, Italy

ARTICLE INFO

Keywords: Reproducibility Information Retrieval Evaluation

ABSTRACT

Science is facing a so-called reproducibility crisis, where researchers struggle to repeat experiments and to get the same or comparable results. This represents a fundamental problem in any scientific discipline because reproducibility lies at the very basis of the scientific method. A central methodological question is how to measure reproducibility and interpret different measures. In Information Retrieval (IR), current practices to measure reproducibility rely mainly on comparing averaged scores. If the reproduced score is close enough to the original one, the reproducibility experiment is deemed successful, although the identical scores can still rely on entirely different result lists. Therefore, this paper focuses on measures to quantify reproducibility in IR and their behavior. We present a critical analysis of IR reproducibility measures by synthetically generating runs in a controlled experimental setting, which allows us to control the amount of reproducibility error. These synthetic runs are generated by a deterioration algorithm based on swaps and replacements of documents in ranked lists. We investigate the behavior of different reproducibility measures with these synthetic runs in three different scenarios. Moreover, we propose a normalized version of Root Mean Square Error (RMSE) to quantify reproducibility better. Experimental results show that a single score is not enough to decide whether an experiment is successfully reproduced because such a score depends on the type of effectiveness measure and the performance of the original run. This study highlights how challenging it can be to reproduce experimental results and quantify the amount of reproducibility.

1. Introduction

Researchers in all areas of science are confronted with the so-called replication or reproducibility crisis that became apparent in more and more disciplines during the last years (Open Science Collaboration, 2015). While initially being discussed in psychological science (Pashler and Wagenmakers, 2012), it quickly became clear that researchers continuously fail to reproduce previous experiments and findings, regardless of being their own or foreign work. As reported by Baker (2016), approximately 70% of researchers in physics and engineering are not able to reproduce someone else's experiments, and roughly 50% fail to reproduce even their own experiments. Computational and data-intensive sciences (Freire, Fuhr and Rauber, 2016; National Academies of Sciences, Engineering, and Medicine, 2019) are no exception, as shown in Artificial Intelligence (AI) and Machine Learning (ML) (Gibney, 2020) and most recently in IR (Breuer, Ferro, Fuhr, Maistro, Sakai, Schaer and Soboroff, 2020).

In IR, the well-known Cranfield paradigm (Cleverdon, 1962, 1967) has helped researchers in the last decades to live in a state of felt comfort and security as test collections allowed a relatively high level of standardization and best practice. In this context, reproducibility was taken for granted. However, reproducibility was far too long a "close enough" attitude. Researchers put reasonable effort into understanding how an approach was implemented and how an experiment was conducted. Then, after several iterations, when they obtain performance scores that somehow resemble the original ones, they decide that an experimental result is reproduced.

Obviously this approach is not enough to guarantee the reproducibility of IR evaluation settings as many different aspects of the systems and experimental settings are neglected, even more today in the current deep learning and neural era (Lin, 2018; Yang, Lu, Yang and Lin, 2019; Marchesin, Purpura and Silvello, 2020). To discuss the different facets

Smm@di.ku.dk (M. Maistro); timo.breuer@th-koeln.de (T. Breuer); philipp.schaer@th-koeln.de (P. Schaer); ferro@dei.unipd.it (N. Ferro)

ORCID(s): 0000-0002-7001-4817 (M. Maistro); 0000-0002-1765-2449 (T. Breuer); 0000-0002-8817-4632 (P. Schaer); 0000-0001-9219-6239 (N. Ferro)



Figure 1: Toy example with rankings that result in the same P@10 but different reproducibility scores. The numbers in the ranking list for each query represent the document identifiers. Ranking items with a darker color represent relevant documents.

of reproducibility, to create awareness within the community, and to formulate concrete guidelines and best practices, a number of events took place in recent years. In 2011 the DESIRE workshop focused on data infrastructures for managing experimental IR data (Agosti, Ferro and Thanos, 2012). A year later a SIGIR workshop on open source IR systems (Trotman, Clarke, Ounis, Culpepper, Cartright and Geva, 2012) discussed how to support reproducibility from the developers' perspective. Workshops on evaluation as a service (Hopfgartner, Hanbury, Müller, Kando, Mercer, Kalpathy-Cramer, Potthast, Gollub, Krithara, Lin, Balog and Eggel, 2015) and reproducible baselines (Arguello, Crane, Diaz, Lin and Trotman, 2015) followed a little later. In SIGIR 2022, a tutorial provided an introduction to reproducible experiments in IR (Lucic, Bleeker, de Rijke, Sinha, Jullien and Stojnic, 2022). The European Conference on Information Retrieval (ECIR) in 2015 and the International ACM Conference on Research and Development in Information Retrieval (SIGIR) in 2022 picked up these ideas and introduced a special Reproducibility Track, where the community is encouraged to submit papers that repeat, reproduce, generalize, and analyze prior work. Here a successful reproduction is not a hard requirement, but the authors are invited to provide an in-depth evaluation of the reproduction process. Next to ECIR and SIGIR, reproducibility is also considered as part of the reviewing process in other major venues.

Different models and definitions were introduced to form a general understanding and to find a common terminology while discussing reproducibility issues in IR. As a general framework the Dagstuhl Seminar on "Reproducibility of Data-Oriented Experiments in e-Science" introduced the Platform, Research goal, Implementation, Method, Actor, and Data (PRIMAD) model in 2016 to describe the different aspects of reproducibility. The ACM "Artifact Review and Badging" guidelines¹ introduce a consistent terminology, which is heterogeneous across disciplines. For the case of reproducibility the ACM guidelines describe "...[reproducibility] means that an independent group can obtain the same result using the author's own artifacts", but does not clarify how to compare these results and what measurement should be used. In recent years the IR community came to the conclusion that simply comparing averages of performance

¹https://www.acm.org/publications/policies/artifact-review-and-badging-current

scores does not tell us anything about the actual degree of reproducibility of the systems at hand. If we simply compare the averages across topics with the "close enough" attitude, we disregard per-topics scores. This means that two systems with the same average score, might behave very differently on different topics. Moreover, per-topic retrieval performances are typically compared by looking at relevance assessments for a given position in the result list, without looking at the actual documents and their position itself, so two completely distinct result lists can potentially produce the same performance score (see Figure 1).

Therefore Breuer et al. (2020) introduced different measures which allow for comparing experimental results at different levels from most specific to most general: the ranked lists of retrieved documents; the actual scores of effectiveness measures; the observed effects; and significant differences. These measures illustrate that nearly identical performance scores do not necessarily mean that the results were "bitwise" reproduced (National Academies of Sciences, Engineering, and Medicine, 2019), which in IR would mean reproducing the same ranked lists of results. With the help of the repro_eval toolkit (Breuer, Ferro, Maistro and Schaer, 2021), the proposed measures can be computed for any given IR experiment. While this is a huge step towards providing a robust environment and setting for reproducible IR experiments, applying and interpreting these measures remains an open question. Indeed, the fundamental methodological question is how to measure reproducibility and provide guidelines on reading and interpreting the different measures.

In this paper, we investigate these issues by considering a set of controlled experiments based on artificially deteriorated runs to find relations between the measure score and the amount of detriment in each experiment. Our aim is to discover reproducibility patterns and achieve a better understanding of how to interpret measures which quantify reproducibility. While in our previous work we investigated both reproducibility (different team, same dataset) and replicability (different team, different dataset) setting, this study's approach – in terms of the ACM guidelines – is only applicable for reproducibility experiments. By focusing on reproducibility only we gain a fixed evaluation environment that allow this research to investigate factors that affect reproducibility of IR experiments, how to measure, and to provide insights on reproducibility measures and their behavior.

The main contributions of this paper are: (1) a critical analysis of IR measures to quantify reproducibility: to the best of our knowledge this is the first study that investigates IR reproducibility with simulated runs; (2) a deterioration algorithm to simulate reproducibility runs: this algorithm is based on 2 operations, i.e., replacing and swapping documents in a ranked list, and allows to control the amount of reproducibility error and to generate different deterioration archetypes; (3) a normalization procedure for RMSE: this mitigates the variability in RMSE score due to the run type and always returns scores in [0, 1], which are easier to interpret.

We conduct experiments with 3 different scenarios and more than 3 millions simulated reproducibility runs (source code publicly available²). Experiments show that quantifying and interpreting reproducibility results is not trivial because reproducibility scores depend on different factors, as the type of reproducibility run under investigation and the measure used to quantify reproducibility, together with its underlying effectiveness measure. This highlights once more that simple comparisons of average scores, the "close enough" attitude, is not enough to determine whether an experiment is successfully reproduced or not.

The remainder of this paper is organized as follows: Section 2 presents related work, Section 3 revises IR reproducibility measures and describes our approach to deteriorate runs; Section 4 reports on the experimental results, and Section 5 presents our conclusions and directions for future work.

2. Related Work

As recently analyzed, the overall scientific progress driven by computational experiments is limited and reproducibility issues together with weak baselines are often seen as a reason for stagnation (Dacrema, Cremonesi and Jannach, 2019; Dacrema, Boglio, Cremonesi and Jannach, 2021; Lv, Ding, Liu, Chen, Feng, He, Zhou, Jiang, Dong and Tang, 2021). Yang et al. (2019) re-evaluated this circumstance and confirmed that it is still true for IR research almost ten years after it had already been pointed out by Armstrong, Moffat, Webber and Zobel (2009). With special regards to the computational sciences, Ivie and Thain (2018) argue that non-reproducible results "stem from a need to simultaneously satisfy the needs of both the computer and a human". According to them, compromises between the concrete operations of a computer and the human reasoning about the experiment on a more abstract level have to be made. These compromises such as software, workflows, and statistics can introduce barriers to reproducibility.

²https://github.com/irgroup/ipm-reproducibility

While there is scientific discourse on how reproducibility and replicability are defined and how the terminology should be applied to the scientific experiments (De Roure, 2014; Plesser, 2017), we follow the definitions by the updated ACM Policy on Artifact and Review Badging. The definition of reproducibility aligns with the International Vocabulary for Metrology (VIM). More specifically, *reproducibility* describes the validation of results by a different team using the same experimental setup. Analogously, we consider an experiment to be reproduced if the results are validated with the same test collection.

Recently, the topic of reproducible research came into focus of the IR community (Ferro, 2017). In the following, we review initiatives and attempts towards reproducible IR experiments. Overall, attempts were made on a conceptual level in the form of workshops, conference tracks, technical platforms and tools, as well as by the explicit review of experimental resources and artifacts.

Conceptually, the *Platform, Research goal, Implementation, Method, Actor, and Data (PRIMAD)* model (Ferro, Fuhr, Järvelin, Kando, Lippold and Zobel, 2016) focuses on defining which experimental components are kept the same and which one are changed in a reproducibility study and, especially, on what is gained and understood in reproducing or failing to reproduce a study.

Since 2015, ECIR hosts a dedicated reproducibility track for analyzing the extent to which previous studies are valid and reproducible. SIGIR introduced its reproducibility track in 2022. Likewise, there have been workshops like RIGOR (Arguello et al., 2015) that analyze reproducibility, inexplicability, and generalizability with a special focus on open-source software, or OSIRRC (Clancy, Ferro, Hauff, Sakai and Wu, 2019) that evaluates reproducible retrieval systems packaged with a technical framework based on Docker, or CENTRE (Soboroff, Ferro, Maistro and Sakai, 2019; Sakai, Ferro, Soboroff, Zeng, Xiao and Maistro, 2019; Ferro, Fuhr, Maistro, Sakai and Soboroff, 2019a) that validates previous experiments at CLEF, NTCIR, and TREC, and a tutorial on best practices to ease reproducibility of IR research (Lucic et al., 2022).

Recently, several open-source retrieval toolkits were introduced that facilitate the implementation of commonly used (keyword-based) retrieval methods. Anserini builds up on the Lucene library and offers a variety of regression guides for experiments with frequently used test collections (Yang, Fang and Lin, 2018). Likewise, the Terrier toolkit implements common retrieval methods. Both toolkits have Python bindings (Lin, Ma, Lin, Yang, Pradeep and Nogueira, 2021; Macdonald and Tonellotto, 2020) that not only make it possible to implement the experiment by a scripting language but also offer interfaces to state-of-the-art machine/deep learning libraries. Even more, the Python toolkit ir_datasets (MacAvaney, Yates, Feldman, Downey, Cohan and Goharian, 2021) offers an interface to a catalog that comprises ad-hoc test collections and related resources like topic, relevance judgments, and benchmarks.

Ferro and Kelly (2018) surveyed the SIGIR community about what reproducibility is in system-oriented and useroriented IR, and more recently, the SIGIR community enforced artifact evaluations as part of the ACM SIGIR Artifact Review and Badging³. The goal is to account for experimental setups, i.e., software and data, that are made available to the community and allow reproductions and replications of the results presented in the original publication.

Besides curated test collections, the IR community promotes reproducibility by archiving experimental data from evaluation campaigns at TREC (Voorhees, Rajput and Soboroff, 2016) or CLEF (Agosti, Di Nunzio, Ferro and Silvello, 2019). Thus, the results of previous experiments, i.e., system runs, can be used as points of reference to which we compare our reimplementations. When it comes to software archival, Potthast, Gollub, Wiegmann and Stein (2019) implemented the TIRA Integrated Research Architecture (TIRA), which has been supporting several shared task events since 2012 (Gollub, Stein and Burrows, 2012a; Gollub, Stein, Burrows and Hoppe, 2012b). TIRA handles software submissions with a cluster to store and host virtual machines. This allows not only storing software permanently but also executing it at any later time. Based on public code repositories and a Docker container environment, the STELLA framework (Breuer and Schaer, 2021; Schaer, Breuer, Castro, Wolff, Schaible and Tavakolpoursaleh, 2021) offered a distributed approach to deploy reproducible software artefacts for IR experiments.

Nonetheless, there are few methods to measure the extent of reproducibility. As an answer, reproducibility measures for the system-oriented IR experiments were introduced (Breuer et al., 2020). The framework comprises a set of different measures that quantify reproducibility (and replicability) with different levels of specificity ranging from finegrained comparisons of document rankings to more general comparisons of topic score distributions. The introduced reproducibility measures were validated with the help of reimplemented results in reference to the original systems runs. Since no dedicated reproducibility dataset was available, the retrieval method was altered in a principled way to simulate a researcher's attempt to reimplement a system with different parameters and configurations. While it is

³http://sigir.org/general-information/acm-sigir-artifact-badging/



Figure 2: The overall approach to investigate the reproducibility measures is based on the comparison of a *reproduced ranking* to a *reference ranking*. Swapping and replacing replacing documents give us control of the deterioration in the reproduced ranking.

possible to control the overall retrieval performance, it is not always intuitive how specific parameters might influence the final document ranking. Since the reimplementations are analyzed by their final outputs, i.e., the system runs, different reimplementations could likewise be simulated by modifying the final ranking only to have even more control of how the ranking is modified. Similar simulated outputs were recently exploited by Parapar and Radlinski (2021) as systematic perturbations of an ideal ranking to evaluate novel metrics for Recommender Systems research.

Score distribution has been studied since the early days of IR (Arampatzis and Robertson, 2011; Swets, 1963; Robertson, 2001) and applied to various tasks such as, for example, rank fusion (Manmatha, Rath and Feng, 2001), thresholds in ranking and filtering (Arampatzis and Kamps, 2009; Arampatzis and van Hameren, 2001), studying the interaction between topics and systems (Robertson and Kanoulas, 2012), and query performance prediction (Cummins, 2014). Recently, Parapar, Losada, Presedo Quindimil and Barreiro (2020) relied on score distribution to simulate perturbed runs. This approach assumes that relevant and non relevant documents follow different score distributions. The score of an IR system is then simulated with a mixture of 2 distributions, in (Parapar et al., 2020) these are 2 lognormal distributions that are fitted with observed scores from IR systems. We did not follow this approach to generate our simulated runs, but we propose a different algorithm (see Section 3.2), which allows to generate deteriorated runs in a principled way, by swapping or replacing relevant documents in a ranked result list. Theoretical results by Ferrante, Ferro and Maistro (2015) ensure that, for example, a positive swap does not decrease the measure score. In this way, we can generate decreasing perturbed runs and keep full control of what modifications in the ranked list produced that decrease. This latter aspect is important to understand and explain the behaviour of perturbed runs and to relate patterns to specific modifications in the ranked result list; it is also matters when it comes to quantify the amount of reproducibility when using measures that consider the ranked result list, as those by Breuer et al. (2020). On the other side, approaches based on score distributions directly generate modified scores, without giving the possibility to link them back to modifications in the ranked result list. Moreover, they have a stochastic component, therefore it is not possible to determine to what extent the variations in measure scores are due to the actual perturbations and to some random effects.

3. Approach

Breuer et al. (2020) presented a number of measures for reproducibility and replicability. These measures were validated with a constellation of runs, generated by systematically changing some set of parameters, e.g., the vocabulary size, the tolerance stopping criterion, the regularization strength, etc. While this work provides an overall analysis of reproducibility and replicability evaluation measures, it still remains unclear how to interpret or quantify the scores of such measures. For example, consider a reproducibility experiment where the RMSE score between the original

and the reproduced run is equal to 0.05. How different are the original and reproduced runs? Can we consider the experiment as successfully reproduced?

In the following sections, we investigate these questions and dive deeper down in the problem of reproducibility experiments. Figure 2 illustrates our overall approach. By starting with a *reference ranking*, we deteriorate the ranking by *swaps and replacements*. Swapping and replacing documents give us better control of the deterioration in the *reproduced ranking*. Both the reference and reproduced rankings are required for the *reproducibility analysis*, which we use in this study to investigate the reproducibility measures. First, we present reproducibility measures from Breuer et al. (2020) and we propose a new approach to normalize RMSE for reproducibility (Section 3.1). Then, we present our approach to generate deteriorated runs (Section 3.2). Starting from any run, we apply two basic operations, swap and replacements, to gradually modify both the order of ranked documents and the set of retrieved documents. Artificially deteriorating the runs allows us to define the amount of deterioration in each run, thus controlling how different is the deteriorated run with respect to the original run. Finally, we describe our algorithm to generate such deteriorated runs.

3.1. Reproducibility Measures

In the following we describe IR reproducibility measures in details. In Breuer et al. (2020), reproducibility measures were categorized in 4 main groups: measures based on the *ordering of documents, effectiveness measures, statistical tests*, and *overall effect over a baseline*. We use the same categorization and focus on the first 3 groups. Measures that consider the overall effect over a baseline, i.e., Effect Ratio (ER) and Delta Relative Improvement (DeltaRI), are excluded because they require both a baseline and advanced run to be reproduced, thus a different deterioration algorithm. Their investigation is left for future work. To the best of our knowledge, there are no other reproducibility measures that could have been included in this work.

Let *r* be the original run and *r'* be the reproduced run. The set of topics is $\mathcal{T} = \{1, \dots, T\}$, where $t \in \mathcal{T}$ denotes a specific topic. The original ranked list of documents in response to topic *t* is $r_t = (d_1, \dots, d_N)$, where d_k is the document at rank position *k*. Similarly, r'_t is the reproduced ranked list of documents for topic *t* and d'_k denotes the document at rank position *k*. Any IR evaluation measure such as Average Precision (AP) or normalized Discounted Cumulated Gain (nDCG) (Järvelin and Kekäläinen, 2002) is denoted by *M*.

Ordering of Documents Kendall's τ Union (KTU) (Breuer et al., 2020; Ferro, Fuhr, Maistro, Sakai and Soboroff, 2019b; Ferro, Maistro, Sakai and Soboroff, 2018) and Rank-Biased Overlap (RBO) (Webber, Moffat and Zobel, 2010) compare the actual rankings of documents in the original and reproduced runs. Since Kendall's τ is not defined for rankings that contain different sets of documents, we compute KTU instead. For each topic, KTU computes $r_t \cup r'_t$, the union of the original and reproduced rankings, by removing duplicate documents. Then, from the original and reproduced rankings, we create a list of rank positions, l_t and l'_t , where the item at position k in l_t is the rank position of d_k in $r_t \cup r'_t$. Finally, the item at position k in l'_t is the rank position of d'_k in $r_t \cup r'_t$. Finally, we compute Kendall's τ (Kendall, 1945) between the lists of rank positions as usual:

$$\begin{aligned} \text{KTU}_{t}(r_{t}, r_{t}') &= \tau(l_{t}, l_{t}') = \frac{P - Q}{\sqrt{(P + Q + U)(P + Q + V)}} \\ \text{KTU}(r, r') &= \frac{1}{T} \sum_{t=1}^{T} \tau(l_{t}, l_{t}') \end{aligned} \tag{1}$$

where P is the total number of concordant pairs (document pairs that are ranked in the same order in both rankings), Q is the total number of discordant pairs (document pairs that are ranked in the opposite order in the two rankings), U and V are the numbers of ties, in l_t and l'_t respectively. Finally, KTU(r, r') is the average computed across topics. Figure 3 illustrates how to compute KTU for two rankings.

Differently from Kendall's τ , RBO can deal with rankings that retrieve a different set of documents. RBO assumes that r_t and r'_t are infinite rankings and is computed as follows:

$$RBO_{t}(r_{t}, r_{t}') = (1 - \phi) \sum_{k=1}^{\infty} \phi^{k-1} \cdot A_{k}$$

$$RBO(r, r') = \frac{1}{T} \sum_{t=1}^{T} RBO_{t}(r_{t}, r_{t}')$$
(2)

Maistro et al.: Preprint submitted to Elsevier



Figure 3: Toy example of computing KTU: r_t and r'_t are the original and reproduced rankings, $r_t \cup r'_t$ is their union, and l_t and l'_t are the lists of rank positions of documents from r_t and r'_t in the union list.

where $\phi \in [0, 1]$ is a parameter to account for top-heaviness: the lower ϕ the higher the weight at the top of the ranking. A_k is the overlap between r_t and r'_t , it is defined as the cardinality of the intersection of r_t and r'_t up to rank position k, divided by k. Intuitively, RBO computes the proportion of overlap between the original and reproduced rankings at each cut-off and then applies a discount based on the cut-off threshold: the higher the cut-off, the higher the overlap, so the more severe the discount.

Effectiveness Measures Root Mean Square Error (RMSE) (Breuer et al., 2020; Ferro et al., 2019b, 2018) computes the quadratic error per topic. Given an IR evaluation measure *M*, it is defined as follows:

RMSE
$$(r, r', M) = \sqrt{\frac{1}{T} \sum_{t=1}^{T} (M(r_t) - M(r'_t))^2}$$
 (3)

Differently from KTU and RBO, RMSE does not consider the order of documents, but just their labels: if two documents with the same label are swapped, RMSE is not affected, since the underlying IR evaluation measure M is not affected. Therefore, we can consider RMSE less sensitive or less stricter than ranking based measures.

One limitation of RMSE, as defined in Equation (3), is the interpretability of its score. RMSE ranges in [0, *B*], where the upper bound $B = \max_{r'} \{RMSE(r, r', M)\}$ depends on the original run *r* and the IR measure *M*. For example, assume that a reproduced run *r'* achieves an RMSE score of 0.15 with respect to AP. We can not infer whether this is successfully reproduced, unless we know the score of the worst case, i.e., RMSE upper bound. If the upper bound was 0.9, then the reproduced run would be a successful experiment, since its score is reasonably far from the worst case. Conversely, if the upper bound was 0.2, then the reproduced run could possibly be improved.

We propose to normalize RMSE by its maximum score, to ease the interpretability of this measure. The maximum value of RMSE for a given run is computed as follows:

max-RMSE(r, M) =
$$\sqrt{\frac{1}{T} \sum_{t=1}^{T} (\max\{M(r_t), 1 - M(r_t)\})^2}$$
 (4)

and normalized Root Mean Square Error (nRMSE) is the ratio between RMSE score and its upper bound:

$$nRMSE(r, r', M) = \frac{RMSE(r, r', M)}{max-RMSE(r, M)}$$
(5)

If we consider the example above, in the first case (upper bound 0.9) nRMSE is approximated by \sim 0.17, while in the second case (upper bound 0.2) it is 0.75. This immediately informs us on the quality of the reproduced run: in

the first case the reproduced run is close to the original run, in the second case more work is needed to get a better reproducibility run.

Equation (4) holds for any IR evaluation measure M ranging in [0, 1] as explained next. For a topic t, the maximum difference is achieved by considering the worst case, i.e., the reproduced ranking scoring 0 or 1:

- if $M(r_t) \ge 0.5$, then the maximum difference is achieved by the reproduced ranking scoring 0: $|M(r_t) M(r'_t)| = M(r_t)$;
- if $M(r_t) < 0.5$, then the maximum difference is achieved by the reproduced ranking scoring 1: $|M(r_t) M(r'_t)| = 1 M(r_t)$.

If the measure M does not range in [0, 1], e.g., Discounted Cumulated Gain (DCG), the above can be adjusted by replacing 0 and 1 with the measure lower and upper bounds, and the decision threshold 0.5 with the middle point of the measure range.

Statistical Tests Breuer et al. (2020) propose to use statistical tests as a mean to compare the original and reproduced runs. A paired two-tailed *t*-test is performed between the scores of the original and reproduced runs for each topic. Then, the *p*-value is used as an indicator of the quality of the reproducibility experiment: the smaller the *p*-value, the stronger the evidence that the reproduced run is different from the original run, thus the reproduced run does not succeed in its purpose.

Note that the *p*-value is the least sensitive measure among those listed above. Not only it is not affected by swapping documents with the same label, as RMSE based measure, but it is also not sensitive to (random) fluctuations on the distribution of scores across topics.

3.2. Run Deterioration

We investigate the behavior of the reproducibility measures in Section 3.1 in a synthetic experimental setting. This allows us to control the amount of "reproducibility error", i.e., quantify how much different is the original run from the reproduced run.

We started from the typical scenario of an IR system implemented as a multi-stage retrieval pipeline: a first retrieval stage, where candidate relevant documents are selected in the complete collection, is followed by a re-ranking stage, where the goal is to push the most relevant documents at the beginning of the ranking. For instance, a cost-efficient key-word based method retrieves the top-k results that are re-ranked by more costly relevance feedback or deep learning approaches.

We identify these two stages in the retrieval pipeline with two types of operations⁴:

- *Stage 1 Replacement*: a document in the ranking is replaced by another document outside the ranking, i.e., a not retrieved document;
- Stage 2 Swap: two documents in the ranking are swapped.

When reproducing the first retrieval stage, a researcher determines the set of retrieved documents. In terms of comparison between the original and reproduced runs, this translates in the replacement of some documents of the original run with others that were not retrieved, or were retrieved below the top k. Analogously, when reproducing the second re-ranking stage, a researcher determines the order of retrieved documents. In terms of comparison between the original and reproduced runs, this translates in swapping pairs of documents in the original run.

Depending on their impact on the ranking, we categorize replacements and swaps as *positive* or *negative*. A positive replacement increases the number of relevant documents retrieved, i.e., a non-relevant document is replaced with a relevant document that was not retrieved. Conversely, negative replacements replace a relevant document with a non-relevant one. A positive swap moves a relevant document towards the top of the ranking, i.e., a relevant document at the bottom of the ranking is swapped with a non-relevant document at the top of the ranking. Conversely, negative swaps move relevant documents from the top to the bottom of the ranking.

Based on the different combinations of swaps and replacements, we define four different archetypes, which basically correspond to the four quadrants of the Cartesian plan in Figure 4. These archetypes relate to the corresponding implementations in a multi-stage retrieval pipeline as explained in the following:

⁴The same operations were identified by Ferrante et al. (2015) in a different context, i.e., they defined a partial order among rankings via swaps and replacements. In this paper we focus on the differences between rankings rather than their mutual order.



Figure 4: Representation of the four archetypes in terms of replacements and swaps.

- *Archetype I*: Replacing non-relevant with relevant documents. Swapping non-relevant with relevant documents. *Example:* First stage ranker improves, as well as the second stage re-ranker.
- Archetype II: Replacing non-relevant with relevant documents. Swapping relevant with non-relevant documents. *Example:* First stage ranker improves, but the second stage re-ranker deteriorates.
- *Archetype III*: Replacing relevant with non-relevant documents. Swapping relevant with non-relevant documents. *Example:* First stage ranker deteriorates, as well as the second stage re-ranker.
- *Archetype IV*: Replacing relevant with non-relevant documents. Swapping non-relevant with relevant documents. *Example:* First stage ranker deteriorates, but the second stage re-ranker improves.

In Breuer et al. (2020), runs were deteriorated in a different way by modifying parameters such as the vocabulary size, the tolerance of the stopping criterion, adding/removing pre-processing steps, etc. This mimics the behavior of a researcher, who tests different values of some parameters in order to find the best combination to reproduce the original work. Even if more realistic, this approach has some limitations. First, it is not straightforward to determine how varying some parameters will affect the original run. For example, assuming that the original run was generated by keeping stopwords, it is hard to estimate a priori the impact of stopwords removal: will the reproduced run score better or worse than the original run? Moreover, it is even harder to predict the impact at the document level. How many replacements and/or swaps will occur when removing stopwords? Second, for some parameters it might be intuitive to estimate their impact in terms of evaluation scores, but this tends to represent only the case of a reproduced run which is worse than the original run (archetype III). For example, reducing the vocabulary size will almost certainly decrease the quality of the reproduced runs. Finally, in both the above situations, we can not quantify ahead the amount of error introduced in the reproduced run.

Algorithm to Deteriorate Runs Algorithm 1 provides the pseudo-code in the case of archetype I, i.e., positive swaps and positive replacements, for a single ranked list of documents, i.e., one topic. The same algorithm can be applied to all the other archetypes with some simple adjustments by performing the appropriate checks (in Figure 5) on the number of non-relevant and relevant documents and changing the direction of swaps and replacements. The provided source code includes the implementation of all archetypes and a variable number of topics.

The algorithm takes as input: *r* the original ranking which will be deteriorated; *s* and *p*, the maximum number of positive swaps and replacements to perform; [a, b] and [c, d], the source and destination intervals, where a < b, c < d and b < c, i.e., [a, b] is closer to the top of the ranking than [c, d] and $[a, b] \cap [c, d] = \emptyset$. The source and destination interval are the regions of the ranking, i.e., rank positions, where we want to perform swaps and replacements. For

Algorithm 1 Deterioration with positive swaps and replacements

Input: Original ranking r, Number of swaps s, Number of Replacements p, Source interval [a, b], Destination interval [c,d]**Output:** Deteriorated ranking r'1: NR src \leftarrow number of non-relevant documents in [a, b] 2: R dst \leftarrow number of relevant documents in [c, d]3: R nrt ← number of relevant and not retrieved documents 4: $r' \leftarrow r$ ▷ initialize the deteriorated ranking 5: if NR src > 0 then 6: $\hat{s} \leftarrow \min(s, \mathbb{R}_{dst})$ ▷ actual number of swaps $\hat{p} \leftarrow \min(p, \mathbb{R}_nrt)$ ▷ actual number of replacements 7: ⊳ total number of operations 8: tot op $\leftarrow \hat{s} + \hat{p}$ if tot_op > NR_src then 9: $c \leftarrow \left[\frac{s \cdot \text{NR}_\text{src}}{s+p}\right]$ 10: $\hat{s} \leftarrow \min(c, \mathbb{R}_{dst})$ 11: $\hat{p} \leftarrow \min(\text{NR src} - \hat{s}, \text{R nrt})$ 12: $tot_op \leftarrow \hat{s} + \hat{p}$ 13: end if $14 \cdot$ $rk_NR_src \leftarrow rank positions of non-relevant documents in [a, b]$ 15: rk_NR_src.random_shuffle() 16: 17: $rk_R_dst \leftarrow rank positions of relevant documents in [c, d]$ rk_R_dst.random_shuffle() 18: 19: for $k = 0, ..., \hat{s}$ do $i \leftarrow \mathrm{rk} \ \mathrm{NR} \ \mathrm{src}[k]$ 20: $j \leftarrow rk_R_dst[k]$ 21: r'[i] = r[j] \triangleright Document at rank *j* becomes the document at rank *i* 22: r'[i] = r[i] \triangleright Document at rank *i* becomes the document at rank *j* 23: end for 24. for $k = \hat{p}, \ldots, \text{tot}$ op do 25: list R nrt ←list of relevant and not retrieved documents 26: $i \leftarrow \mathrm{rk} \ \mathrm{NR} \ \mathrm{src}[k]$ 27: $r'[i] \leftarrow \text{list}_R_\text{nrt.pop}()$ \triangleright Replacement at rank *i* 28: end for 29. 30: end if 31: return Deteriorated ranking r'

positive swaps, the algorithm swaps non-relevant documents from the source interval with relevant documents in the destination interval. For positive replacements, the algorithm replaces non-relevant documents from the source interval with relevant not retrieved documents. Note that swaps and replacements will not interfere with each other, i.e., the rank positions where swaps and replacements are performed do not overlap. This means that a document that was swapped can not be replaced, and a document that was replaced can not be swapped.

The algorithm starts by checking that the number of non-relevant documents in the source interval is greater than zero (line 5). If all documents in the source interval are relevant, we can not perform any positive swap or replacement, therefore the ranking can not be deteriorated. In this case, we simply return the original ranking r.

If there are non-relevant documents in the source interval, the algorithm determines the actual number of swaps and replacements, i.e., the number of swaps and replacements that can be performed. The algorithm always performs the highest number of actual swaps and replacements, which is lower or equal to the inputted number of swaps and replacements. This is done in 3 steps, detailed in the following.

First, the algorithm checks that the number of swaps to perform does not exceed the number of relevant documents in the destination interval. Indeed, relevant documents in the destination interval are swapped with non-relevant documents in the source interval. Therefore, the number of actual swaps is set to the minimum between the inputted



Figure 5: Upper bounds on the number of operations for each archetype.

number of swaps and the number of relevant documents in the destination interval (line 6). If there are no relevant documents in the destination interval, no swaps are performed.

Second, the algorithm checks that the number of replacements does not exceed the number of relevant documents that are not retrieved. Therefore, the actual number of replacements is set to the minimum between the inputted number of replacements and the number of relevant documents not retrieved (line 7). If all the relevant documents are retrieved, no replacements are performed.

Third, the total number of operations performed by the algorithm is set to the sum of the number of actual swaps and replacements (line 8). The total number of operations has to be lower or equal to the number of non-relevant documents in the source interval. If this requirement is not met (line 9), then the actual number of swaps and replacement is decided so that the original proportion of inputted swaps and replacements is preserved. This is determined by solving the following system of equations with variables \hat{s} and \hat{p} :

$$\begin{cases} \frac{\hat{s}}{\hat{p}} = \frac{s}{p} \\ \hat{s} + \hat{p} = \text{NR_src} \end{cases}$$
(6)

where NR_src is the total number of non-relevant documents in the source interval. The integer solution is given in Equation (7):

$$\hat{s} = \left\lfloor \frac{s \cdot \text{NR_src}}{s+p} \right] \qquad \hat{p} = \text{NR_src} - \hat{s}$$
(7)

where $\lfloor \cdot \rfloor$ denotes the closest integer (round function).

Finally, the actual number of swaps and replacements needs to be compared with the available slots. Therefore, the actual number of swaps is set to the minimum between the value in Equation (7) and the number of relevant documents in the destination interval (line 11). The actual number of replacements, is set to the minimum between the value in Equation (7) and the number of relevant documents that are not retrieved (line 12).

Figure 5 reports the upper bounds for the number of operations with respect to each combination of swap and replacements: $rel[\cdot, \cdot]$ and $nrel[\cdot, \cdot]$ are the number of relevant and non-relevant documents in the specified interval, R_rt and R_nrt are the number of relevant documents retrieved and not retrieved, and NR_rt is the number of non-relevant documents retrieved. Note that there are two different cases, defined by the location where operations are performed: (1) swap and replacement operate on the same set of relevant or non-relevant documents (archetypes I

 Table 1

 Experimental scenarios defined by the number and types of topics, number and types of runs, and recall score.

Scenario 1	Scenario 2	Scenario 3
1 simulated topic 3 simulated runs Recall = 1	1 simulated topic 3 simulated runs Recall = 0.5	50 real topics 2 real runs -

and III); (2) swap and replacement operate on disjoint sets of documents, relevant versus non-relevant documents (archetypes II and IV). The case determines the upper bound on the total number of operations: in case (1) the sum of swaps and replacements should be always lower than the available slots, i.e., the number of relevant documents in [a, b] for archetype I, and the number of non-relevant documents in [a, b] for archetype IV; in case (2), if we ignore all the other constraints, we can perform a swap and a replacement for each document in [a, b], so the greatest upper bound is the number of rank positions in the interval, i.e., b - a + 1.

The rest of the algorithm performs swaps (lines 19 - 24) and replacements (lines 25 - 29). Given the list of rank positions of non-relevant documents in the source interval (line 15) and the list of rank positions of relevant documents in the destinations interval (line 17), \hat{s} rank positions are selected randomly by shuffling those lists (lines 16 and 18). Then, non-relevant documents are swapped from the source interval with relevant documents in the destination interval (lines 29 - 24).

Similarly, \hat{p} rank positions are selected randomly in the list of rank positions of non-relevant documents in the source interval. This set does not contain rank positions where a swap was performed, so a swapped document can not be replaced. Then, given the list of relevant documents that are not retrieved, non-relevant documents are replaced by relevant ones (line 28).

In the experiments (see Section 4), Algorithm 1 is applied iteratively with the inputs p and s in specified ranges. This generates deteriorated rankings with an increasing number of replacements and swaps. In each iteration, the documents that are swapped and replaced are chosen randomly and independently of the previous iterations. Let $r'_{p,s}$ be a ranking deteriorated with p replacements and s swaps. Then $r'_{p+1,s+1}$ is not equal to $r'_{p,s}$ with an additional replacement and swap, but it is a new ranking where the rank positions of swaps and replacements are randomly chosen and might not overlap with the rank positions chosen in the previous iteration.

4. Experiments

Next, we present the experimental analysis of reproducibility measures when runs are deteriorated with swaps and replacements. First, we describe the experimental set-up and implementation details (Section 4.1), then we analyze the properties of different experimental scenarios (Section 4.2), finally we report and discuss the experimental results (Section 4.3). All the source code to run the experiments is publicly available 5 .

4.1. Experimental Set-up

Experimental Scenarios We define 3 different experimental scenarios (Table 1), where we vary the original run, i.e., the input run of the deterioration algorithm. In the first 2 scenarios, we use simulated runs with one topic, i.e., a single ranking of documents. This allows us to better understand the behavior of reproducibility measures in a simple controlled setting. We generate 3 different types of simulated rankings:

- Perfect ranking: all relevant documents are at the top of the ranking;
- Realistic Ranking: relevant documents are placed with decreasing probability while going down in the ranking (exponential distribution);
- Reversed ranking: all relevant documents are placed at the end of the ranking.

We differentiate scenarios 1 and 2 by the recall of the rankings. In scenario 1, the simulated rankings retrieve all the relevant documents, i.e., recall equal to 1. In Scenario 2, the simulated rankings do not retrieve all the relevant documents, i.e., we assume recall equal to 0.5.

⁵https://github.com/irgroup/ipm-reproducibility

In scenario 3, we used 2 real runs, BM25 (Robertson, Walker, Jones, Hancock-Beaulieu and Gatford, 1994; Robertson and Zaragoza, 2009) and RM3 (Lavrenko and Croft, 2001), with 50 topics each. This allows us to generalize the behavior from simulated simple settings (1 and 2) to real runs.

Experimental Settings In our analysis, we consider all the reproducibility measures in Section 3.1: Kendall's τ , RBO, RMSE, nRMSE and *p*-values, instantiated as in Breuer et al. (2020). Specifically, RBO is computed with $\phi = 0.8$, which roughly corresponds to greater weight on the top 5 rank positions (the smaller ϕ , the more top-heavy the measure) (Webber et al., 2010). RMSE and nRMSE are instantiated with AP, nDCG, and Precision@10 (P@10). Recall that in scenarios 1 and 2 there is a single topic, so we can not compute the *p*-value.

The simulated runs in scenarios 1 and 2 are the same, they retrieve 1000 documents, out of which 100 are relevant. In scenario 2, 100 extra relevant documents are added outside the ranking (relevant and not retrieved documents), so the recall is equal to 0.5. 32In scenario 3, the real runs are derived from the Washington Post Corpus (v2) and the topics of TREC Common Core 2018^6 .

We chose this collection because: (1) deep pools were used to collect relevance labels, so there is a good amount of relevant documents that can be swapped or replaced; (2) the same collection was used in (Breuer et al., 2020), so the results from our experiments with deteriorated runs can be compared with the results with real reproducibility runs in (Breuer et al., 2020); (3) the collection is recent and better reflects the performance of state-of-the-art IR systems. On this collection, we validate two runs derived with either the baseline method BM25 (Robertson et al., 1994) or the advanced method BM25 combined with RM3 (Jaleel, Allan, Croft, Diaz, Larkey, Li, Smucker and Wade, 2004). For both indexing and retrieval, we make use of the Anserini toolkit (Yang et al., 2018). The document collection is indexed with Porter stemming and Indri's stopword list (Strohman, Metzler, Turtle and Croft, 2005). During retrieval, we use the topic's title and description as the query string and instantiate both BM25 ($k_1 = 0.9$, b = 0.4) and RM3 (10 feedback terms, 10 feedback documents, an original query weight of 0.5) with Anserini's default settings.

The deterioration algorithm is implemented in Python 3.9 and run with source interval [1, 500] (first half) and destination interval [501, 1000] (second half). The number of swaps and replacements is varied from 0 to 250 with step 1. Indeed, having 500 available rank positions, we can perform at most 250 swaps and 250 replacements. In scenario 3, the same input parameters are applied to all topics.

The deterioration algorithm is applied iteratively for each pair (p, s) of replacements and swaps. In scenarios 1 and 2 this corresponds to 62 500 randomly deteriorated rankings for each run and to 3 125 000 randomly deteriorated rankings for each run in scenario 3.

All the reproducibility measures are computed with the Python toolkit repro_eval (Breuer et al., 2021). Measure scores are reported with heatmaps, as illustrated in the visual representation of the 4 archetypes in Figure 4. The x-axis represents replacements and the y-axis swaps. Negative integers denote negative operations, i.e., negative replacements or swaps. All heatmaps are generated with seaborn⁷ in Python and are down-sampled with step 5. The complete figures with step 1 are not included in this paper due to their large size (~ 5 MB each figure), but are stored in a public Zenodo archive⁸. For the sake of better readability and understandability, we have included the corresponding Tables with numerical results in the Appendix A. Note that for RMSE, nRMSE and the p-values, the Tables only contain AP scores. The remaining results of nDCG and P@10 can be found in the public code repository.

4.2. Analysis of Experimental Scenarios

Next we analyse the properties of the experimental scenarios. This provides the context to better understand the experimental results in Section 4.3.

4.2.1. Scenario 1: Simulated Runs with Recall 1

As summarized in Table 1, scenario 1 considers 3 simulated runs with one topic. All the 3 runs retrieve all the 100 relevant documents, therefore they have recall = 1.

The Recall Base (RB), i.e., the total number of relevant documents, sets upper bounds on the number of replacements that can be performed, as detailed in Section 3.2. Specifically, no positive replacements can be performed, since there are no relevant documents that are not retrieved. For negative replacements, we can at most replace RB relevant documents with non-relevant ones, thus a maximum of RB negative replacements.

⁶https://trec-core.github.io/2018/

⁷https://seaborn.pydata.org/

⁸https://zenodo.org/record/5902542

Similarly, the type of run, sets upper bounds on the number of swaps that can be performed. For the perfect ranking, no positive swaps can be performed, since all the relevant documents are already placed at the beginning of the ranking. This, combined with the absence of positive replacements, means that for the perfect ranking there is no possible deterioration within archetype I (first quadrant). Therefore, in all figures of the perfect ranking in scenario 1, the first quadrant achieves the best score for all measures.

Conversely, for the reversed ranking, we can not perform negative swaps, since all the relevant documents are already placed at the end of the ranking. Moreover, we can not perform negative replacements, since in the source interval (first half of the ranking) there are no relevant documents to replace. Therefore, for the reversed ranking we can not perform any type of replacement and negative swaps, thus the run can not be deteriorated with respect to archetypes III and IV. This means that the third and fourth quadrants achieve the best score for all measures for the reversed ranking in scenario 1.

Finally, the total number of operations, i.e., sum of swaps and replacements, can not exceed the available slots, as detailed in Figure 5. For example, for the perfect ranking, we can perform at most RB, i.e., 100, negative replacements and negative swaps. This explains why in all figures there is a central region where we can see variations in the measure score, this is the region where the sum of operations is lower than the available slots. Moving further from the center, the number of operations reaches the upper bound, so there is no more room to deteriorate the run and the measure reaches a plateau.

4.2.2. Scenario 2: Simulated Runs with Recall Less than 1

Scenario 2 exploits exactly the same runs as scenario 1, but 100 extra relevant and not retrieved documents are added, so recall drops from 1 to 0.5. This allows us to consider extra operations, as positive replacements, that can not be performed in scenario 1. Therefore, the upper bound on the number of replacements corresponds to RB/2 = 100 both for positive and negative replacements.

As discussed for scenario 1, the type of run determines further constraints on the number of swaps and replacements. For the perfect ranking, positive swaps can not be performed, since all the documents are already at the top of the ranking. However, positive replacements can happen, since there are 100 relevant not retrieved documents. Therefore, differently from scenario 1, in scenario 2, the measure scores for the perfect ranking change in all quadrants.

For the reversed ranking, negative swaps can not occur, since all the relevant documents are already placed at the end of the ranking. Negative replacements can not occur as well, since there are no relevant documents in the source interval, i.e., the first half of the ranking. However, differently from scenario 1, we can perform positive replacements. This means that there is no deterioration for archetype III (negative replacements and swaps), while for all other archetypes we can perform some operations, so the measure scores can change.

As in scenario 1, the total number of operations (replacements and swaps) is limited by the number of available slots, as illustrated in Figure 5. Therefore all figures corresponding to scenario 2 exhibit a central region, where the measure scores decrease, and an outer region, where the measure scores reach a plateau.

4.2.3. Scenario 3: Real Runs

In scenario 3 we use real runs with 50 topics computed with BM25 and RM3. In principle, all the operations can be performed and the total number of operations is constrained by the type of run, i.e., number of relevant documents and their location, as illustrated in Figure 5. Therefore, all figures corresponding to real runs exhibit a trend similar to scenarios 1 and 2, i.e., a central region with varying scores and an outer region with a plateau score.

On average across topics, RM3 run performs slightly better than BM25 run with respect to all IR evaluation measures under consideration (AP, nDCG and P@10). This means that, on average, RM3 retrieves more relevant documents closer to the top of the ranking than BM25. As a consequence, the total number of available positive operations is lower for RM3 than BM25. For example, since RM3 has a higher recall, there are fewer relevant not retrieved documents that can be used for positive replacements. Conversely, the number of negative operations is higher for RM3 than BM25. For example, since RM3 has a higher recall, there are more relevant retrieved documents that can be used for positive replacements.

4.3. Experimental Results

Next, we present the experimental results for the three different scenarios: (1) simulated rankings that retrieve all relevant documents; (2) simulated rankings that retrieve half of the relevant documents; (3) real runs computed with BM25 and RM3. The results are organized by measure, starting with ranking-based measures with KTU and RBO



Figure 6: Scenario 1 - Kendall's τ Union: Simulated runs with a single topic. Each run retrieves all relevant documents (Recall = 1). The runs are deteriorated with replacements in [1,500] and swaps from [1,500] to [501,1000].



Figure 7: Scenario 2 - Kendall's τ **Union:** Simulated runs with a single topic. Each run retrieves only half of the relevant documents (Recall = 0.5). The runs are deteriorated with replacements in [1, 500] and swaps from [0, 500] to [501, 1000].

(Section 4.3.1), effectiveness-based measures with RMSE and nRMSE (Section 4.3.2), and statistical test measures with *p*-values (Section 4.3.3).

4.3.1. Ranking-based Measures: KTU and RBO

Figure 6 shows KTU values for the deteriorated runs in scenario 1. As expected, we can see that KTU reaches the best score equal to 1 for archetype I with the perfect ranking (Figure 6a) and for archetypes III and IV with the reversed ranking (Figure 6c). KTU for the realistic ranking (Figure 6b) looks very similar to the perfect ranking. This happens because most of the relevant documents are placed in [1, 500], as in the perfect ranking, and Kendall's τ is not top-heavy, so it does not account for the rank position where a swap or replacement occurred.

Figure 7 reports KTU values for scenario 2, where positive replacements can occur. The effect of positive replacements is clearly seen on the right part of each subfigure, indeed Figure 7 can be obtained from Figure 6 by overlaying the effect of positive replacements in quadrants I and IV. Moreover, the minimum KTU score drops from scenario 1 to 2, since positive replacements represent additional operations that we can perform to deteriorate the runs.

The perfect and realistic rankings exhibit a similar behavior, as observed in scenario 1. This happens for the same reason, i.e., the perfect ranking retrieves most of the relevant documents in the source interval [1, 500] and KTU does not account for the rank position.

The bright region where the measure score is close to the best score is reduced from the whole archetype I in Figure 6 to an horizontal thick line corresponding to the zero replacements axis and positive swaps in Figure 7. This happens because: (1) we can perform positive replacements, which affects archetype I; (2) the number of positive swaps is upper bounded by the number of relevant documents in the interval [501, 1000], which is 0 or close to 0 for the perfect and realistic ranking. Therefore, for the perfect and realistic ranking we can mostly observe the effect of



Figure 8: Scenario 3 - Kendall's τ Union: BM25 and RM3 runs on TREC Common Core 2018. The runs are deteriorated with replacements in [1, 500] and swaps from [0, 500] to [501, 1000].

replacements in archetypes I and II. For the reversed ranking in Figure 7c, the bright region is reduced to archetype III, since no negative replacements and swaps can occur.

For the perfect and realistic rankings the worst region in terms of KTU is archetype IV, which corresponds to the region where the highest number of operations can be performed (up to 100 positive replacements and 100 negative swaps). This is followed by archetype III, where the 100 available slots are split between negative replacements and swaps, archetype II, where no positive swaps can be performed, and archetype I, where there are only positive replacements. The reversed ranking behaves in a somehow specular way, with archetype I, the one with the greatest number of operations, being the worst, and archetype III, no available operations, being the best.

Figure 8 shows KTU values for the real runs of scenario 3. In this case we report KTU averaged across all 50 topics. At a first glance, the behavior of the real runs is very similar to the perfect and realistic rankings in Figure 7. The main difference is that the vertical bright line is thinner for real runs and can be seen also for the horizontal axis corresponding to zero swaps and positive replacements. This might happen because of the upper bounds in terms of positive operations, i.e., there are not so many relevant documents that are not retrieved that we can use for positive replacements, or relevant documents in [501, 1000] that we can move towards the beginning of the ranking.

As for the scenario 2, the worst region is archetype IV, where we can perform the highest number of operations because positive replacements and negative swaps are somehow complementary (see Figure 5). The best region is archetype I, where the total number of operations is constrained by the number of relevant documents not retrieved and the number of relevant documents retrieved at the end of the ranking. When we compare BM25 and RM3, Figures 8a and 8b, we can see that archetype I is slightly brighter for RM3 than BM25. This happens because, on average, RM3 retrieves more relevant documents and closer to the top of the ranking, so there are less available slots for positive operations.

In general, all scores in KTU figures (Figures 6, 7, and 8) are higher then those reported in Breuer et al. (2020), where KTU reaches values close to 0. This is due to two main reasons: (1) the upper bound on the number of operations; (2) the type of operations. First, with real reproducibility runs there is no upper bound to the number of negative replacements: for example we can always replace a document with a non-relevant one. Second, real runs include other types of deterioration operations, as replacements and swaps in the same relevance class. For example, swapping two non-relevant or relevant documents affects KTU, but it is not considered in our analysis. As already mentioned, this synthetic experimental set-up allows us to control for variables such as the number of replacements and swaps, while keeping the level of complexity low, which allow us to interpret the results.

As seen from the darker regions in the heatmaps of the simulated rankings and real runs, archetype IV has the lowest KTU scores. It means that a less effective first-stage ranker combined with a more effective re-ranker leads to the least reproduced version of the original ranking in terms of KTU. If the first stage ranking deteriorates, the exact document order cannot be reproduced even if the re-ranking is better than in the original experiment.



Figure 9: Scenario 1 - RBO: Simulated runs with a single topic. Each run retrieves all relevant documents (Recall = 1). The runs are deteriorated with replacements in [1, 500] and swaps from [1, 500] to [501, 1000].



Figure 10: Scenario 2 - **RBO:** Simulated runs with a single topic. Each run retrieves only half of the relevant documents (Recall = 0.5). The runs are deteriorated with replacements in [1, 500] and swaps from [0, 500] to [501, 1000].

Figure 9 reports RBO in scenario 1. Even if both KTU and RBO considers the rank position of documents, we can clearly see the difference in their behavior, due to RBO being top-heavy. Indeed, the frontier region is more clear for RBO than KTU and RBO shows a more diverse behavior across runs.

The perfect and the realistic rankings (Figures 9a and 9b) exhibit again a similar behavior. RBO is equal or close to the perfect score for archetype I and there is a similar frontier region, where RBO transitions from the best score to lower scores. The size of the frontier region corresponds to the maximum number of operations that can be performed: a straight line corresponding to -100 in quadrant II and IV and the bisector with sum -100 in quadrant III. This happens because we can perform up to 100 negative replacements and no positive swap in quadrant II; up to 100 negative swaps and replacements in quadrant III; and up to 100 negative swaps and no positive replacement in quadrant IV.

The main difference between the perfect and realistic rankings is the color of the outer region, i.e., the value of RBO when the total number of deterioration operations is performed. For the perfect ranking this score is close to 0, the worst score, for the realistic ranking this score is close to 0.56. This behavior is due to RBO being top-heavy: for the perfect ranking, all the documents at rank positions [1, 100] are replaced by non-relevant documents or moved at the end of the ranking, both these operations have a great impact on RBO. For the realistic ranking the relevant documents are spread in the ranking, so their swaps or replacement has less impact on RBO score.

The reversed ranking (Figure 9c) exhibits a different behavior. There is no change in RBO scores for quadrant III and IV, since we can not perform any negative swap nor positive or negative replacement. Therefore, the upper half of Figure 9c shows only the effect of positive swaps, from the end of the ranking to [1, 500]. RBO does not reach a plateau for the reversed ranking because its value heavily depends on the rank positions where relevant documents are swapped, which are randomly chosen with up to 100 choices over 500 rank positions in each iteration.



Figure 11: Scenario 3 - RBO: BM25 and RM3 runs on TREC Common Core 2018. The runs are deteriorated with replacements in [1,500] and swaps from [0,500] to [501,1000].

Figure 10 shows RBO values in scenario 2, i.e., the same runs of scenario 1, but with 100 extra relevant and not retrieved documents. We can clearly see the effect of positive replacements for the realistic and reversed rankings (Figures 10b and 10c): as for KTU, Figure 10 can be obtained from Figure 9 by overlaying the layer generated by positive replacements on quadrants I and IV. The effect of positive replacements can not be seen for the perfect ranking (Figure 10a) due to RBO top heaviness. Indeed the top of the ranking is already filled with relevant documents and positive replacements affect rank positions chosen at random in [101, 500], which affect RBO score only to a small extent.

Conversely, for the realistic ranking, even a few replacements in the interval [1, 100] can considerably decrease RBO score. Moreover, for the realistic ranking, archetype I does not exhibit a uniform color, i.e., a plateau score, because the drops in RBO heavily depends on the rank position where a positive replacement occurs. The same reasoning applies also to the reversed ranking. For the realistic ranking we can still see the vertical bright region corresponding to few replacements, which is also visible for KTU (see Figure 7a).

When we compare archetypes, the difference is not marked as clearly as for KTU. Archetypes IV and I are still the worst regions for the realistic and reversed rankings respectively. These correspond to the regions where the maximum number of operation can be performed, accordingly to the constraints in Figure 5. For the perfect ranking there is no clear distinction among archetypes where RBO reaches very low values close to 0. This happens because negative operations affects the very top of the raking in [1, 100], thus they decrease RBO score to a great extent.

Figure 11 shows RBO averaged across 50 topics for the real runs in scenario 3. The behavior of both runs is somehow in between the perfect and the realistic ranking: there is a bright area corresponding to archetype I and archetype IV is slightly darker than all other regions. The impact of positive replacements can still be seen in the quadrants I and IV of both figures, but it is not as strong as for the realistic ranking in Figure 10b. The reason is that RBO is averaged across topics, and for several of them there are only a few relevant documents not retrieved that can be added.

In general, RBO scores in all figures (Figures 9, 10, and 11) is within the ranges reported by (Breuer et al., 2020). The only exception is the perfect ranking, which exhibits lower values. Indeed, the type of deterioration that can be performed for the perfect ranking is quite extreme and not likely in a real case, e.g., moving/replacing all documents at the top of the ranking.

RBO behaves differently from KTU because of its top-heaviness. For RBO, the deterioration operations that occur at the beginning of the run have greater impact on the scores. Therefore, rather the the quantity, as for KTU, for RBO the location of those operations is the most important factor.

Compared to KTU, the RBO can be parameterized as a more top-heavy rank correlation measure. Depending on the parametrization, RBO reaches a plateau for archetypes II, III, and IV, as can be seen for the simulated rankings and real runs. It means that a more effective first-stage ranker cannot compensate for a less effective re-ranker (and vice versa) regarding the reproduction quality of a top-weighted ranking once they deviate too much from the original ranking methods.



Figure 12: Scenario 1 - **RMSE:** Simulated runs with a single topic. Each run retrieves all relevant documents (Recall = 1). The runs are deteriorated with replacements in [1, 500] and swaps from [1, 500] to [501, 1000].

4.3.2. Effectiveness-based Measures: RMSE and nRMSE

Figure 12 shows RMSE value in scenario 1 for all effectiveness measures: AP (first row), nDCG (second row), and P@10, (third row). Recall that, for the perfect ranking (first column), positive replacements and swaps can not be performed, so archetype I corresponds to RMSE = 0. Conversely, for the reversed ranking (third column), negative swaps and replacements can not occur, so archetype III and IV have RMSE = 0. Finally, since recall = 1 in scenario 1, there are no positive replacements.

Scenarios 1 and 2 consider simple simulated runs with a single topic. These represent special cases for RMSE, where RMSE score is the absolute difference between the effectiveness score of the original and reproduced runs.

In scenario 1, RMSE behaves similarly to KTU and RBO (see Figures 6 and 9): the perfect and realistic rankings exhibit a similar behavior across measures, due to the position of the relevant documents (mostly) at the top of the rankings. Moreover, we can observe the same shape of the frontier region, with the boundary around 100 negative





Figure 13: Scenario 2 - **RMSE:** Simulated runs with a single topic. Each run retrieves only half of the relevant documents (Recall = 0.5). The runs are deteriorated with replacements in [1,500] and swaps from [0,500] to [501,1000].

replacements and swaps, and the diagonal line y = -x - 100 for quadrant III. Recall that this happens because there are 100 relevant and retrieved documents, so we can perform up to 100 negative replacements or swaps.

In terms of RMSE scores, the perfect ranking assumes a slightly darker tone than the realistic ranking. This is observed also for RBO (Figure 9) and is due to: (1) the top-heaviness of IR effectiveness measures, and (2) the rank position of relevant documents. Since all relevant documents are placed at the top of the ranking, the perfect ranking can be deteriorated to a greater extent and top-heavy measures are more affected.

The main difference among IR effectiveness measures is the variation of scores across archetypes: nDCG seems more sensitive to the type of deterioration, while AP and P@10 do not exhibit any difference for archetypes II, III, and IV. This is due to the different discount functions: P@10 focuses only on the top 10 positions, thus it does not detect what happens after the cut-off threshold; AP has a steeper discount function than nDCG, thus it is less sensible to replacements and swaps at lower rank positions.

Figure 13 shows RMSE scores for scenario 2, i.e., the same rankings as in scenario 1, but retrieving half of the relevant documents (recall = 0.5). In this scenario we can perform positive replacements, since there are 100 relevant and not retrieved documents. Therefore, the only archetype where RMSE score does not change is archetype III for the reversed ranking, where negative replacements and swaps can not occur.

Similarly to KTU and RBO, we can clearly see the effect of positive replacements. For the realistic and reversed rankings, we can still identify this as a sort of overlaying layer in the positive replacement region, i.e., archetypes I and IV. This is more evident for the reversed ranking than the realistic ranking. The only exception, where the effect of positive replacements is not visible at all, is P@10 for the perfect ranking. In this case there is no effect of positive replacements because all the positive replacements happen below the cut-off threshold, indeed the perfect ranking is filled with relevant documents up to position 100 and all positive replacements occur below.

Again AP and nDCG behave in a similar way and in general the perfect ranking is slightly darker than the realistic ranking due to the more extreme deterioration that can be applied to the perfect ranking. Moreover, as in scenario 1, nDCG is more sensitive to different types of operations, assigning different colors to different archetypes. The worst region is located in archetype II, corresponding to negative replacements and positive swaps. Symmetrically, the best region is located within archetype IV, with positive replacements and negative swaps. In this latter region, 100 relevant documents are swapped with documents in the second half of the ranking, but their effect is somehow compensated by the 100 non-relevant documents replaced with relevant ones in the first half of the ranking.

Wile KTU and RBO exhibit a bright vertical region corresponding to the positive *y*-axis (few replacements and positive swaps), the same region has a different shape for RMSE computed with AP and nDCG. Indeed, the vertical bright line corresponding to the positive *y*-axis continues in quadrant IV, where it starts on the diagonal and then divides itself in two branches, one vertical and one horizontal. This corresponds to the boundary where negative swaps are compensated by positive replacements. The horizontal branch is higher for AP than nDCG since AP discount function is harsher than nDCG one, therefore a higher number of replacements is needed to compensate relevant documents that are moved at the end of the ranking.

Figure 14 shows RMSE scores for real runs in scenario 3. Both runs consider 50 topics, thus RMSE is not simply the absolute difference in effectiveness scores, as in scenario 1 and 2. Since scores are aggregated across topics, colors in the heatmaps are less variable and the horizontal and vertical bright lines are barely visible.

In scenario 3, we can observe similar trends as in the other two scenarios. First, AP and nDCG behave in a similar way. They both exhibit a thin bright line corresponding to the positive *y*-axis. The region in archetype IV where positive replacements compensate negative swaps is not visible for real runs, except for a thin line just below the positive *x*-axis. Second, nDCG is more sensitive to different archetypes while AP treats them equally. The worst region is still archetype II, where negative replacements and positive swaps occur. Since there are very few relevant documents at the end of the run, they can not compensate the impact of negative replacements in the first half of the run. Third, RMSE is darker for RM3 than BM25, as it happens for the perfect and realistic rankings. Since RM3 performs slightly better than BM25, deterioration can be more extreme and lead to higher differences in RMSE scores.

When comparing different IR effectiveness measures, Breuer et al. (2020) report that P@10 suffers of higher variance than AP and nDCG. This is clearly shown in Figures 12 and 13 in the frontier region, where the transition to the plateau score is smoother for AP and nDCG, and in archetype I and II for the reversed ranking, where RMSE does not reach a plateau value. In scenario 3, this effect is less strong but still visible for archetype I.

Compared to the rank correlation measures, the RMSE does not depend on the actual documents but on the relevance labels in the reproduced ranking. RMSE can be instantiated with different retrieval measures, and as such, it shows a different sensitivity behaviour across the archetypes depending on the retrieval measure. When instantiated with P@10, the heatmaps are similar to those of RBO, focusing on the reproduction quality of the top-ranked results. Comparing AP and nDCG reveals that the discount function has an impact on the sensitivity of RMSE, while nDCG is more sensitive (cf. archetype II) due to the less harsh discount function. As can be seen from the simulated rankings for a single topic, a more effective re-ranker can indeed compensate for a less effective reproduced first-stage ranker to a certain extent, whereas the required replacements depend on the discount function of the retrieval measure.

Figure 15 reports nRMSE scores for AP, nDCG, and P@10 in scenario 1. Recall that nRMSE is a version of RMSE normalized by the maximum RMSE score that can be achieved by the reproduced run, i.e., the worst possible reproduced run. Therefore, nRMSE general trend, i.e., appearance of the heatmaps, should not be different from



Figure 14: Scenario 3 - RMSE: BM25 and RM3 runs on TREC Common Core 2018. The runs are deteriorated with replacements in [1,500] and swaps from [0,500] to [501,1000].

standard RMSE, even if there is a difference in the actual score. Indeed, the heatmaps in Figure 15 exhibit the same trend of those in Figure 12, i.e., the same behavior across archetypes and the same shape of the frontier region. As for standard RMSE, nDCG is the only measure able to detect the difference among archetypes (see Figures 15d and 15e).

The effect of normalization is clearly visible when comparing Figures 15 and 12. Indeed, all heatmaps, except for the perfect ranking, are darker with nRMSE, meaning that the relative error (nRMSE) is higher than the absolute error (RMSE). The difference between the perfect and realistic rankings (first and second column) is almost imperceptible. The two rankings are somehow comparable and the same deterioration is applied, which leads to similar nRMSE scores. nRMSE and RMSE for the perfect ranking are the same (Figure 15 and Figure 16). This happens because the normalization factor for the perfect ranking is equal to 1 (see Equation (4)).

Figure 16 reports nRMSE in scenario 2, i.e., the same runs as in scenario 1 but with the addition of 100 relevant and not retrieved documents. As expected, the general appearance of nRMSE is very similar to RMSE (compare Figures 16 and 13): same shape of the frontier region and same behavior across measures and runs. nDCG is once more the best



Figure 15: Scenario 1 - nRMSE: Simulated runs with a single topic. Each run retrieves all relevant documents (Recall = 1). The runs are deteriorated with replacements in [1, 500] and swaps from [1, 500] to [501, 1000].

measure to distinguish among different archetypes. The bright region where negative swaps are counterbalanced by positive replacements is still visible for nDCG and AP with respect to the perfect and realistic ranking. Again, the line where positive replacements compensate negative swaps is closer to the positive *x*-axis for AP than nDCG. This is due to AP requiring more positive replacements to compensate negative swaps because of its steeper discount function.

The impact of positive replacements can be identified as an additional layer in quadrants I and IV, more visible for P@10 (third row) and the reversed ranking with all measures (third column). This happens also for RMSE, RBO, and KTU. Recall that no positive replacements can occur at rank positions [1, 10], so the effect of positive replacements is not visible for P@10 and the perfect ranking (Figure 16g).

The effect of normalization is still visible as a general darkening of colors in the heatmaps. In scenario 2, this is true also for the perfect ranking, where the normalization factor is < 1 and not 1 as in scenario 1, because there are relevant documents that are not retrieved. The only exception is P@10 with respect to the perfect ranking, because the measure considers only rank positions up to 10. The effect of normalization does not completely remove difference



Figure 16: Scenario 2 - nRMSE: Simulated runs with a single topic. Each run retrieves only half of the relevant documents (Recall = 0.5). The runs are deteriorated with replacements in [1, 500] and swaps from [0, 500] to [501, 1000].

among the perfect and realistic ranking, as in scenario 1, except for P@10. Again, because the perfect ranking places all the retrieved relevant documents at rank positions [1, 100], we can deteriorate this ranking in a more extreme way. This is further combined with the top-heaviness of IR measures, indeed AP is more top-heavy than nDCG, i.e., the difference between the perfect and realistic ranking is more noticeable for AP than nDCG.

Figure 17 shows nRMSE for scenario 3 with real runs, both with 50 topics. Again the general behavior of nRMSE is similar to RMSE (compare Figures 17 and 14) and the same considerations are valid: (1) there is a thin bright line corresponding to the positive *x*-axis and the line where positive replacements counterbalance negative swaps is barely visible; (2) nDCG is better in distinguishing among archetypes than AP; (3) nRMSE is slightly darker for RM3 than BM25.

The effect of normalization can be clearly detected as darker colors across all measures and runs. Again, the relative error shown by nRMSE is higher than absolute differences shown by RMSE.



Figure 17: Scenario 3 - nRMSE: BM25 and RM3 runs on TREC Common Core 2018. The runs are deteriorated with replacements in [1, 500] and swaps from [0, 500] to [501, 1000].

In comparison to RMSE, the nRMSE is normalized by the worst possible ranking. Overall, similar trends are observable for the archetypes, but the darker heatmaps indicate larger relative errors in comparison to the absolute differences.

4.3.3. Statistical tests: p-values

Next we present the same heatmaps for *p*-values. These are computed only for scenario 3 because this is the only scenario where there are 50 topics. Results are reported in Figure 18.

At a first glance we can see that all figures are mainly black meaning that *p*-value is very sensitive to deterioration in the runs. Indeed, the darker the color, the smaller the *p*-value, thus the higher the evidence that the original and reproduced runs are significantly different.



Figure 18: Scenario 3 - p-value: BM25 and RM3 runs on TREC Common Core 2018. The runs are deteriorated with replacements in [1,500] and swaps from [0,500] to [501,1000].

Furthermore, all the figures can be interpreted as a sort of "photographic negative" of RMSE and nRMSE in Figures 14 and 17. The bright line corresponding to the positive *y*-axis, could be seen also for RMSE and nRMSE, even if not so strongly. The horizontal line close to the positive *x*-axis represents the border where negative swaps are compensated by positive replacements. This line is closer to the *x*-axis for AP than nDCG. The same behavior is observed in scenario 2 for RMSE and nRMSE and is due to the different discount functions of these measures.

Among different IR effectiveness measures, nDCG (second row) seems the most sensitive, indeed all archetypes are mostly dark. AP (first row) is less sensitive with respect to archetype I. This happens because for real runs we can perform more negative operations, since the number of relevant documents in [0, 500] is in general higher than the number of relevant documents in [501, 1000] or the number of relevant not retrieved. P@10 (third row) is the measure with highest variance, especially with respect to archetype I. Indeed, adding or moving relevant documents at the top of the ranking, either by positive swaps or replacements, has a stronger effect than removing them.

4.4. Final Remarks

As reproducibility measures, KTU and RBO compare the ranking of documents for the original and reproduced runs. These are the most stringent reproducibility measures, because they require that the reproduced run ranks the documents in the same exact order as the original run. Both measures account more for the amount of deterioration operations and RBO also for their location, rather than the operation type, i.e., replacement or swap. Up to 50 replacements or swaps can cause KTU drop from 1 to 0.8, which is still a reasonably high correlation in IR (Voorhees, 1998, 2001). However, we can observe much lower correlations in studies with real runs (Breuer et al., 2020). This suggests that the reproduced runs, in a non-controlled experimental setting, can be very different from the original runs, with many more than 50 swaps or replacements.

The main difference between KTU and RBO is due to RBO top-heaviness. This means that even less than 50 replacements or swaps can cause RBO drop to the lowest score, i.e., 0. This happens when swaps and replacements occur at the top of the ranking, because RBO is less sensitive to changes towards the end of the ranking. Therefore, an RBO score of 0.8, might denote a great overlap at the top of the ranking, but further analyses are needed when the whole ranking has to be reproduced.

RMSE, nRMSE and statistcal tests depend on the underlying effectiveness measure, e.g., AP, nDCG, etc. Therefore, the reproducibility scores returned by these measures should be considered in combination with: (1) the type of effectiveness measures and its property; (2) the type of original run, the more extreme the run (very good or very bad performing run), the greater the impact in terms of reproducibility. For example, with RMSE, 50 replacements or swaps can increase RMSE up to: 0.5 with AP, 0.6 with nDCG, and 1 with P@10. If we consider only AP, 50 replacements or swaps can increase RMSE in a range from 0.5 (perfect ranking) to 0.2 (realistic ranking).

While nRMSE can mitigate the impact of the run type, not much can be done in terms of difference between measures, as they have different properties and they are intended to evaluate runs with different point of views. Topheavy rank based measures (AP and nDCG) exhibit a region where negative swaps are counterbalanced by positive replacements, thus relatively low RMSE and nRMSE scores (lower than 0.1) or high *p*-values (greater than 0.8) might not be enough to conclude that a run was successfully reproduced. On the other side, RMSE, nRMSE and statistical tests with a set based measure as P@10, exhibit high variance, so a few swaps or replacements in the top 10 rank positions can cause very high scores. Finally, Table 2 summarizes our experimental results and provides an overview of the reproducibility measures and the corresponding use cases.

5. Conclusions, Limitations and Future Work

This paper presents an analysis of IR reproducibility measures. We consider rank based measures, i.e., KTU and RBO, which compare the original and reproduced lists of documents. As effectiveness based measures, we consider RMSE, which compares the difference in terms of IR effectiveness measures, e.g., AP, between the original and reproduced runs. We propose nRMSE, a normalized version of RMSE, to account for the relative difference instead of the absolute one. Finally, we consider statistical tests by looking at the *p*-value.

We propose a deterioration algorithm to deteriorate an input run and generate a reproduced run, where we can control the amount of deterioration. The deterioration algorithm exploits 2 operations: replacement of a retrieved document with a non retrieved one; swaps of two documents in the run.

We compute all reproducibility measures with the output of the deterioration algorithm. We propose 3 experimental scenarios with different types of runs: (1) synthetic runs with a single topic and recall = 1; (2) synthetic runs with a single topic and recall = 0.5; (3) real runs with 50 topics.

Once more, our experiments confirm that reproducibility is not a trivial problem. Not only it is hard to reproduce the results of a published paper, but it is even harder to decide when such results are reproduced to an acceptable extent. For example, given a reproducibility experiment with RMSE score equal to 0.05, we can not conclude by solely looking at this score that reproducibility is achieved. This happens because, even if our deterioration algorithm allows to control the amount of noise, we can not completely disentangle the effects due to the type of run, amount of relevant documents, and the IR effectiveness measure. Reproducibility measures such as RMSE or statistical tests depend on the underlying effectiveness measure and on the performance of the original run. Our normalized version of RMSE can mitigate the variations due to the original run but the dependency on the IR effectiveness measure can not be removed. Other measures, such as KTU and RBO do not depend on any effectiveness measure but are much stricter, as they require the same exact order of documents for each topic. Moreover, RBO top heaviness can be misleading, as a few changes in the very top of the ranking can cause a severe drop in its score.

Table 2

Summary and overview of the reproduciblity measures. For each measure, the level of specifity and a corresponding use case are given. The experimental findings are summarized once again and examples of possible evaluation criteria are given.

Measure	Level of specificity	Use case
KTU	Document order	KTU should be used for the most specific level of measuring reproducibility. Our experimental results suggest that if the first stage ranking deteriorates, a re-ranker cannot compensate for the deterioration and recover the exact document ranking. Use this measure for the highest degree of rigor.
RBO	Document order	Similar to KTU, the RBO is determined by the document ranking. In comparison, it implies a user model and can be parameterized to put higher weights on top-ranked documents. As a result, it is less strict than KTU, i.e., if parameterized accordingly, it will result in higher scores, when there are overlaps between the top-ranked documents in the original and the reproduced run. Similar to KTU, our experiments show that a deteriorated reproduced first stage ranker cannot be compensated by an improved reproduced re-ranker in terms of RBO. Use this measure when it is important to account for the correlation between the document rankings that likely only have a few overlaps.
RMSE / nRMSE	Effectiveness	RMSE is determined by the distance between the topic score distributions. It is a document-agnostic reproducibility measure, i.e., a good or perfect reproduction does not depend on particular documents but on the corre- sponding relevance labels instead. It means that a perfect reproduction (RMSE=0) could be achieved with different documents that comply with the relevance labels in the original rankings. Our experimental results suggest that an improved reproduced re-ranker can indeed compensate a deteriorated first stage ranker as different document can result in the same performance scores at the topic-level. nRMSE is normalized by the worst possible ranking, which results in larger relative errors. Use these measures, if it is less important to reproduce the exact document
p-values	Statistical comparison	ordering but it is more important to reproduce the effectiveness without neglecting the score distribution over the topics. At the most general level, the topic score distributions can be compared by paired t-tests. The intuition follows the idea of low p-values indicating a higher probability of failing the reproduction. In turn, higher p-values indicate more similar topic score distributions. Our experimental results suggest that the p-values are very sensitive to deteriorations in the runs. Use this approach, if it is important to account for the reproduction of topic score distributions beyond RMSE / nRMSE.

Moreover, our reproducibility study is limited by the choice of the experimental collection and the number of topics for real run. It would be beneficial to validate our results with a larger number topics, however publicly available IR collections usually include 50 or fewer assessed documents. One exception is the TREC Million Query Track, which includes 1, 755 judged topics in 2007 Allan, Carterette, Dachev, Aslam, Pavlu and Kanoulas (2008), 782 in 2008 Allan, Aslam, Pavlu, Kanoulas and Carterette (2009) and 684 in 2009 Carterette, Pavlu, Fang and Kanoulas (2010). However, these collections are not good choices for our experimental set-up because: (1) shallow pools were used so there is only a small number of relevant documents per topic that can be swapped or replaced; (2) statistical approaches were used to select the documents to judge, therefore only special variations of AP should be used to account for the estimated relevance.

The present work is also limited by the number and type of operations considered by our deterioration algorithm. As future work, we plan to extend the deterioration algorithm in different directions: first, we can consider different types of operations to better mimic real reproducibility runs, for example, swaps and replacements with multi-graded relevance; second, we can consider different experimental set-ups, for example, focus the deterioration on smaller intervals at the top of the run instead of considering the first and second half. Moreover, we can extend the present work and consider other measures as Effect Ratio (ER) and other tasks as replicability or generalizability.

Finally, our study does not consider the user perspective and what is the impact of reproducibility errors on the user experience. We plan to conduct a user study to understand the impact of reproducibility on the final user, for example, what is the impact of replacing or swapping a document in terms of user behavior, e.g., clicks.

As for any evaluation task, the final remark is to use different evaluation measures of different categories, depending on the final goal and/or the application domain. Comparing scores averaged across topics is not enough to conclude that an experiment is successfully reproduced. As our experiments show, even very similar per-topic scores can correspond to different rankings. Therefore, together with the source code, we should consider the release of runs computed on publicly available datasets, following the idea of open runs (Voorhees et al., 2016) and, if possible, annotated in a standardized way (Breuer, Keller and Schaer, 2022).

Implementation and Data

The implementation of all our experiments is publicly available at https://github.com/irgroup/ipm-repro ducibility. The GitHub repository includes instructions on how to run the code. The experiments in this paper are conducted only on publicly available data.

Acknowledgements

This paper was partially supported by the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 893667 and the German Research Foundation (DFG) under project no. 407518790.

CRediT authorship contribution statement

Maria Maistro: Conceptualization, Methodology, Software, Writing, Project administration. Timo Breuer: Conceptualization, Methodology, Software, Writing. Philipp Schaer: Conceptualization, Methodology, Writing, Supervision. Nicola Ferro: Conceptualization, Methodology, Writing, Supervision.

A. Appendix

Table 3Scenario 1 - Kendall's τ Union

						Pe	rfect rank	ing				
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.8111	0.8111	0.8111	0.8111	0.8999	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
200		0.8111	0.8111	0.8111	0.8111	0.8960	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
150		0.8111	0.8111	0.8111	0.8111	0.8970	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
100		0.8111	0.8111	0.8111	0.8111	0.9080	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
50		0.8111	0.8111	0.8111	0.8111	0.9086	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
0		0.8111	0.8111	0.8111	0.8111	0.8945	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-50		0.8073	0.8128	0.7775	0.7987	0.7767	0.8788	0.8687	0.8979	0.8567	0.8919	0.8883
-100		0.7900	0.7842	0.7017	0.7904	0.7870	0.7000	0.7484	0.7424	0.7601	0.7503	0.7420
-150		0.7605	0.7993	0.7941	0.7649	0.7023	0.7125	0.7058	0.7401	0.7520	0.7365	0.7555
-250		0.7954	0.7620	0.7841	0.7740	0.7432	0.7559	0.7300	0.7520	0.7021	0.7303	0.7598
		0.1.501	0.1001	0.1012	0.1110	D			0.1010	0.1000	0.1.105	0.1000
						Kea	austic rank	king				
\sim	Replace-						_					
<u> </u>	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.8111	0.8111	0.8111	0.8111	0.8919	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
200		0.8111	0.8111	0.8111	0.8111	0.9063	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
150		0.8111	0.8111	0.8111	0.8111	0.9010	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
100		0.8111	0.8111	0.8111	0.8111	0.8980	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
50		0.0111	0.0111	0.0111	0.0111	0.9012	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
50		0.0111	0.0111	0.0111	0.0111	0.0905	0.8674	1.0000	1.0000	1.0000	1.0000	1.0000
-100		0.7930	0.0010	0.7094	0.7934	0.7900	0.0074	0.0900	0.0079	0.8050	0.0520	0.0700
-150		0.7645	0.7665	0.7887	0.7906	0.7722	0.7608	0.7590	0.7508	0.7154	0.7459	0.7613
-200		0.7831	0.7691	0.7653	0.7775	0.7669	0.7518	0.7528	0.7787	0.7189	0.7635	0.7370
-250		0.7806	0.7886	0.7418	0.7756	0.7717	0.7413	0.7417	0.7313	0.7297	0.7394	0.7337
						Rev	ersed ran	king				
	Roplace							5				
	ments	-250	-200	-150	_100	-50	Ο	50	100	150	200	250
Swaps	ments	250	200	150	100	50	0	50	100	150	200	230
250		0 77/12	0 7300	0 77/18	0 7315	0 7550	0 7/79	0 7344	0 7/3/	0 7357	0 7680	0 7834
200		0.7742	0.7509	0.7740	0.7313	0.7520	0.7470	0.7544	0.7434	0.7337	0.7602	0.7054
150		0.7520	0.7710	0.7232	0 7365	0.7150	0.7407	0.7535	0 7393	0.7650	0.7632	0.7460
100		0.7470	0.7612	0.7613	0.7794	0.7306	0.7654	0.7411	0.7658	0.7626	0.7188	0.7636
50		0.8880	0.8860	0.8943	0.8582	0.8777	0.8807	0.8622	0.8954	0.8702	0.8724	0.8616
0		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-50		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-100		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-150		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-200		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-250		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000

Table 4Scenario 2 - Kendall's τ Union

						Pe	rfect rank	ing				
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.8111	0.8111	0.8111	0.8111	0.8980	1.0000	0.9314	0.8732	0.8722	0.8666	0.8609
200		0.8111	0.8111	0.8111	0.8111	0.9064	1.0000	0.9440	0.8681	0.8636	0.8711	0.8625
150		0.8111	0.8111	0.8111	0.8111	0.8977	1.0000	0.9277	0.8656	0.8798	0.8773	0.8680
100		0.8111	0.8111	0.8111	0.8111	0.9005	1.0000	0.9279	0.8726	0.8720	0.8688	0.8801
50		0.8111	0.8111	0.8111	0.8111	0.9069	1.0000	0.9370	0.8626	0.8799	0.8636	0.8815
0		0.8111	0.8111	0.8111	0.8111	0.9015	1.0000	0.9304	0.8786	0.8676	0.8753	0.8630
-50		0.8139	0.8001	0.7983	0.7949	0.7714	0.8623	0.8044	0.7571	0.7414	0.7641	0.7504
-100		0.8015	0.7782	0.7750	0.7980	0.7494	0.7608	0.6991	0.6419	0.0315	0.6403	0.6350
-150		0.7984	0.7726	0.7898	0.7038	0.7744	0.7590	0.0755	0.0321	0.0501	0.0020	0.0549
-200		0.8000	0.7640	0.7705	0.7000	0.7419	0.7519	0.7112	0.0341	0.0437	0.0100	0.0303
-230		0.1900	0.7091	0.7015	0.7050	0.7001	0.1550	0.0190	0.0495	0.0393	0.0409	0.0307
						Rea	alistic ran	king				
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.8111	0.8111	0.8111	0.8111	0.8971	1.0000	0.9392	0.8754	0.8754	0.8890	0.8669
200		0.8111	0.8111	0.8111	0.8111	0.9040	1.0000	0.9443	0.8848	0.8813	0.8611	0.8716
150		0.8111	0.8111	0.8111	0.8111	0.8865	1.0000	0.9409	0.8736	0.8695	0.8782	0.8817
100		0.8111	0.8111	0.8111	0.8111	0.8960	1.0000	0.9281	0.8789	0.8709	0.8707	0.8614
50		0.8111	0.8111	0.8111	0.8111	0.9049	1.0000	0.9392	0.8751	0.8741	0.8720	0.8911
0		0.8111	0.8111	0.8111	0.8111	0.9014	1.0000	0.9367	0.8900	0.8749	0.8704	0.8702
-50		0.7904	0.8066	0.8051	0.7825	0.7717	0.8628	0.8107	0.7404	0.7389	0.7578	0.7449
-100		0.7733	0.7880	0.7560	0.7889	0.7812	0.7471	0.7050	0.0205	0.0453	0.6526	0.0022
-100		0.7070	0.7043	0.7072	0.7507	0.7471	0.7203	0.0901	0.0134	0.0304	0.0520	0.0409
-250		0.754	0.7755	0.7564	0 7763	0.7668	0.7209	0.6530	0.6034	0.0439	0.0333	0.6312
			0.1100	0.1001	0.1100				0.0001	0.0202	0.0010	0.0011
						Kev	ersed ran	king				
	Replace-						-					
~	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.7541	0.7523	0.7639	0.7216	0.7618	0.7456	0.7137	0.6144	0.6638	0.6627	0.6028
200		0.7324	0.7825	0.7722	0.7823	0.7293	0.7719	0.7218	0.6472	0.6217	0.6439	0.6756
150		0.7235	0.7580	0.7264	0.7581	0.7361	0.7146	0.6724	0.6630	0.6202	0.6542	0.6520
100		0.7769	0.7483	0.7540	0.7514	0.7632	0.7547	0.7430	0.6458	0.6263	0.6268	0.6445
50		0.8974	0.8///	0.8///	0.8922	0.8530	0.8766	0.8012	0.7428	0.7727	0.7327	0.7562
U 50		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9410	0.8700	0.070	0.8728	U.8/14
-50		1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	0.9370	0.0704	0.0079	0.0790	0.0095
-150		1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	0.9302	0.8823	0.8781	0.8850	0.8835
-200		1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	0.9460	0.8689	0.8786	0.8799	0.8631
-250		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9451	0.8702	0.8773	0.8630	0.8689
		1.0000	1.0000	1.0000	2.0000	1.0000	2.0000	5.5.51	5.0.02	5.01.0	5.0000	5.0005

Table 5Scenario 3 - Kendall's τ Union

							BM25					
	Replace- ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.9018	0.9023	0.9031	0.9049	0.9159	0.9781	0.9508	0.9439	0.9398	0.9378	0.9396
200		0.9026	0.9025	0.9023	0.9046	0.9145	0.9780	0.9535	0.9441	0.9393	0.9407	0.9402
150		0.9029	0.9027	0.9027	0.9037	0.9144	0.9792	0.9516	0.9421	0.9394	0.9382	0.9407
100		0.9025	0.9024	0.9029	0.9042	0.9147	0.9789	0.9519	0.9421	0.9400	0.9402	0.9400
50		0.9042	0.9044	0.9044	0.9065	0.9162	0.9806	0.9546	0.9448	0.9404	0.9419	0.9403
0		0.9227	0.9227	0.9227	0.9250	0.9361	1.0000	0.9735	0.9644	0.9588	0.9613	0.9604
-50		0.9175	0.9175	0.9146	0.9130	0.9123	0.9132	0.8856	0.8731	0.8753	0.8742	0.8736
-100		0.9142	0.9133	0.9117	0.9098	0.9024	0.8955	0.8746	0.8657	0.8599	0.8568	0.8588
-150		0.9110	0.9073	0.9083	0.9059	0.9042	0.8943	0.8700	0.8623	0.8590	0.8569	0.8584
-200		0.9100	0.9075	0.9069	0.9022	0.9000	0.8931	0.8698	0.8621	0.8553	0.8529	0.8558
-250		0.9068	0.9064	0.9042	0.9014	0.8997	0.8893	0.8696	0.8614	0.8549	0.8560	0.8564
							RM3					
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.8976	0.8972	0.8993	0.9022	0.9170	0.9805	0.9599	0.9514	0.9499	0.9477	0.9459
200		0.8987	0.8990	0.8985	0.9034	0.9182	0.9799	0.9585	0.9529	0.9483	0.9467	0.9474
150		0.8974	0.8985	0.8998	0.9030	0.9180	0.9801	0.9567	0.9527	0.9484	0.9484	0.9462
100		0.8978	0.8979	0.8980	0.9022	0.9175	0.9808	0.9596	0.9527	0.9484	0.9479	0.9463
50		0.8980	0.8999	0.9012	0.9029	0.9179	0.9814	0.9601	0.9530	0.9494	0.9489	0.9478
0		0.9158	0.9167	0.9182	0.9215	0.9366	1.0000	0.9779	0.9723	0.9687	0.9669	0.9668
-50		0.9106	0.9089	0.9083	0.9058	0.9052	0.9130	0.8925	0.8869	0.8853	0.8824	0.8798
-100		0.9036	0.9049	0.9021	0.8974	0.8948	0.8942	0.8695	0.8632	0.8618	0.8570	0.8566
-150		0.9040	0.9008	0.8966	0.8938	0.8895	0.8845	0.8663	0.8602	0.8529	0.8537	0.8527
-200		0.8993	0.8990	0.8946	0.8944	0.8866	0.8849	0.8627	0.8565	0.8546	0.8501	0.8497
-250		0.8962	0.8958	0.8930	0.8920	0.8881	0.8821	0.8583	0.8586	0.8492	0.8529	0.8475

Table 6 Scenario 1 - RBO

						Pe	rfect rank	ing				
	Replace-											
<u> </u>	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.0000	0.0000	0.0000	0.0000	0.7455	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
200		0.0000	0.0000	0.0000	0.0000	0.3113	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
150		0.0000	0.0000	0.0000	0.0000	0.3262	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
100		0.0000	0.0000	0.0000	0.0000	0.6531	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
50		0.0000	0.0000	0.0000	0.0000	0.3277	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
0		0.0000	0.0000	0.0000	0.0000	0.8534	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-50		0.0000	0.0000	0.0000	0.0000	0.0000	0.5425	0.3000	0.1005	0.2700	0.0002	0.0393
-100		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		0.0000	0.0000	0.0000	0.0000	D			0.0000	0.0000	0.0000	0.0000
						Kea	instic rank	king				
	Replace-						_			. – -		
~	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.5636	0.5636	0.5636	0.5636	0.7925	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
200		0.5636	0.5636	0.5636	0.5636	0.6542	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
150		0.5636	0.5636	0.5636	0.5636	0.6853	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
100		0.5636	0.5636	0.5636	0.5636	0.9435	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
50		0.5636	0.5636	0.5636	0.5636	0.6645	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
0		0.5636	0.5636	0.5636	0.5636	0.8030	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-50		0.5030	0.5030	0.5030	0.5030	0.5030	0.0498	0.8701	0.8738	0.7951	0.7178	0.7020
-100		0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050
-130		0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050	0.5050
-250		0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636
						De		lina				
						Kev	ersed ran	king				
	Replace-						_			. – -		
~	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.8188	0.8773	0.3387	0.9764	0.4923	0.8383	0.7831	0.7705	0.9413	0.9572	0.8120
200		0.9997	0.7226	0.7720	0.8590	0.5825	0.9749	0.9990	0.8643	0.4903	0.9378	0.9456
150		0.7834	0.9771	0.7591	0.9128	0.8126	0.9801	0.9354	0.7729	0.9420	0.3013	0.6279
100		0.7404	0.3980	0.8966	0.8811	0.7819	0.7901	0.4352	0.9234	0.5900	0.8156	0.9976
50		0.9949	0.8734	0.9999	0.9203	0.9974	0.9507	0.4129	0.9997	0.9804	0.9893	0.9773
0		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-5U 100		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-100		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-100		1 0000	1.0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000
-200 -250		1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000	1 0000
		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000

Table 7 Scenario 2 - RBO

Replacements -250 -200 -150 -100 -50 0 50 100 150 200 250 550 0.0000 0.0000 0.0000 0.0000 0.0000 1.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000							Pe	rfect rank	ing				
ments -250 -200 -150 -100 -50 0 50 100 150 200 200 250 0.0000 0.0000 0.0000 0.0000 0.2975 1.0000 0.0000 0.0000 0.0000 0.0000 0.0000		Replace-											
Swaps 250 0.0000 0.0000 0.0000 0.0000 0.2975 1.0000 1.0		ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	Swaps												
200 0.0000 0.0000 0.0000 0.0000 1.0000 <td>250</td> <td></td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.2975</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td>	250		0.0000	0.0000	0.0000	0.0000	0.2975	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
150 0.0000 0.0000 0.0000 0.6763 1.0000 <td>200</td> <td></td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.2957</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td>	200		0.0000	0.0000	0.0000	0.0000	0.2957	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	150		0.0000	0.0000	0.0000	0.0000	0.6763	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
50 0.0000 0.0000 0.0000 0.6000 1.0000 0.0000	100		0.0000	0.0000	0.0000	0.0000	0.1627	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
0 0.0000 0.0000 0.0000 0.0000 1.0000 1.0000 1.0000 1.0000 -50 0.0000	50		0.0000	0.0000	0.0000	0.0000	0.6891	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-50 0.0000 <td>0</td> <td></td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.5992</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td> <td>1.0000</td>	0		0.0000	0.0000	0.0000	0.0000	0.5992	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
-100 0.0000 <td>-50</td> <td></td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> <td>0.4035</td> <td>0.3270</td> <td>0.9357</td> <td>0.2121</td> <td>0.2886</td> <td>0.1245</td>	-50		0.0000	0.0000	0.0000	0.0000	0.0000	0.4035	0.3270	0.9357	0.2121	0.2886	0.1245
150 0.0000 <td>-100</td> <td></td> <td>0.0000</td>	-100		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-200 0.0000 <td>-150</td> <td></td> <td>0.0000</td>	-150		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-250 0.0000 <td>-200</td> <td></td> <td>0.0000</td>	-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Realistic ranking Replace- ments -250 -200 -150 -100 -50 0 50 100 150 200 250 250 0.5636 <	-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Replacements -250 -200 -150 -100 -50 0 50 100 150 200 250 250 0.5636							Rea	alistic rank	king				
ments -250 -200 -150 -100 -50 0 50 100 150 200 250 250 0.5636		Replace-							-				
Swaps Survey Survey <thsurvey< th=""> <thsurvey< td="" th<=""><td></td><td> ments</td><td>-250</td><td>-200</td><td>-150</td><td>-100</td><td>-50</td><td>0</td><td>50</td><td>100</td><td>150</td><td>200</td><td>250</td></thsurvey<></thsurvey<>		ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	Swaps												
200 0.5636 <td>250</td> <td></td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.9444</td> <td>1.0000</td> <td>0.9999</td> <td>0.8776</td> <td>0.5976</td> <td>1.0000</td> <td>0.5969</td>	250		0.5636	0.5636	0.5636	0.5636	0.9444	1.0000	0.9999	0.8776	0.5976	1.0000	0.5969
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	200		0.5636	0.5636	0.5636	0.5636	0.8307	1.0000	0.5972	0.9993	0.9995	0.9995	0.8770
100 0.5636 <td>150</td> <td></td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.8918</td> <td>1.0000</td> <td>0.9999</td> <td>0.5972</td> <td>0.8776</td> <td>0.9995</td> <td>0.9999</td>	150		0.5636	0.5636	0.5636	0.5636	0.8918	1.0000	0.9999	0.5972	0.8776	0.9995	0.9999
50 0.5636 0.5636 0.5636 0.7027 1.0000 0.9995 0.9623 0.5976 0.9989 1.0000 0 0.5636 <	100		0.5636	0.5636	0.5636	0.5636	0.9394	1.0000	0.5972	0.5970	0.9617	1.0000	0.5970
0 0.5636	50		0.5636	0.5636	0.5636	0.5636	0.7027	1.0000	0.9995	0.9623	0.5976	0.9989	1.0000
-50 0.5636 <td>0</td> <td></td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.6544</td> <td>1.0000</td> <td>1.0000</td> <td>0.9618</td> <td>0.9995</td> <td>1.0000</td> <td>0.8770</td>	0		0.5636	0.5636	0.5636	0.5636	0.6544	1.0000	1.0000	0.9618	0.9995	1.0000	0.8770
-100 0.5636 <td>-50</td> <td></td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.5636</td> <td>0.8806</td> <td>0.8363</td> <td>0.6129</td> <td>0.7744</td> <td>0.6844</td> <td>0.3930</td>	-50		0.5636	0.5636	0.5636	0.5636	0.5636	0.8806	0.8363	0.6129	0.7744	0.6844	0.3930
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	-100		0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5630	0.4408	0.1235	0.1235	0.4413
-200 0.5636 0.5637 0.200	-150		0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5259	0.5258	0.4413	0.4412
-250 0.5636 0.5636 0.5636 0.5636 0.5636 0.5636 0.5635 0.1607 0.5253 0.1608 Reversed ranking Replacements -250 -200 -150 -100 -50 0 50 100 150 200 250 Swaps 250 0.8776 0.6695 0.7280 0.9956 0.8717 0.9730 0.6580 0.4023 0.8390 0.8775 0.2124 200 0.8158 0.9665 0.3960 0.9577 0.5963 0.9392 0.9366 0.1607 0.1553 0.3721 0.4337 150 0.9837 0.7668 0.9194 0.7405 0.8697 0.3151 0.7288 0.3326 0.5141 0.2673 0.3118 100 0.9565 0.7962 0.9399 0.8605 0.4223 0.5476 0.9291 0.9208 0.9862 0.4508 0.8407 50 0.8161 0.9719 0.9694 0.8661 <t< td=""><td>-200</td><td></td><td>0.5636</td><td>0.5636</td><td>0.5636</td><td>0.5636</td><td>0.5636</td><td>0.5636</td><td>0.5636</td><td>0.5636</td><td>0.5253</td><td>0.4413</td><td>0.0388</td></t<>	-200		0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5253	0.4413	0.0388
Reversed ranking Replace- ments -250 -200 -100 -50 0 50 100 150 200 250 -100 100 100 -250 -250 -250 -250 -250 -250 -250 -250 -250 -100 -250 -250 -250 -260 -260 -260 -250	-250		0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5636	0.5635	0.1607	0.5253	0.1608
Replace- ments -250 -200 -150 -100 -50 0 50 100 150 200 250 Swaps 250 0.8776 0.6695 0.7280 0.9956 0.8717 0.9730 0.6580 0.4023 0.8390 0.8775 0.2124 200 0.8158 0.9665 0.3960 0.9657 0.5963 0.9392 0.9366 0.1607 0.1553 0.3721 0.4337 150 0.9837 0.7668 0.9194 0.7405 0.8697 0.3151 0.7288 0.3326 0.5141 0.2673 0.3118 100 0.9565 0.7962 0.9399 0.8605 0.4223 0.5476 0.9291 0.9208 0.9862 0.4508 0.8407 50 0.8161 0.9719 0.9694 0.8661 0.5923 0.5714 0.8412 0.6009 0.8974 0.8458 0.8022 0 1.0000 1.0000 1.0000 1.0000 1.0000 0.9999 0.8735 0.9184							Rev	ersed ran	king				
ments-250-200-150-100-50050100150200250Swaps2502500.87760.66950.72800.99560.87170.97300.65800.40230.83900.87750.21242000.81580.96650.39600.96570.59630.93920.93660.16070.15530.37210.43371500.98370.76680.91940.74050.86970.31510.72880.33260.51410.26730.31181000.95650.79620.93990.86050.42230.54760.92910.92080.98620.45080.8407500.81610.97190.96940.86610.59230.57140.84120.60090.89740.84580.802201.00001.00001.00001.00001.00000.99990.87350.91840.65830.9466-1001.00001.00001.00001.00001.00000.97320.83670.77240.66900.9655-1501.00001.00001.00001.00001.00000.94530.91370.95530.96500.8105-2001.00001.00001.00001.00001.00000.94530.91370.95530.91850.9213-2001.00001.00001.00001.00001.00000.94530.91370.95530.91850.9215-2501.00001.00001.0000 <td></td> <td>Replace-</td> <td></td>		Replace-											
Swaps 250 0.8776 0.6695 0.7280 0.9956 0.8717 0.9730 0.6580 0.4023 0.8390 0.8775 0.2124 200 0.8158 0.9665 0.3960 0.9657 0.5963 0.9392 0.9366 0.1607 0.1553 0.3721 0.4337 150 0.9837 0.7668 0.9194 0.7405 0.8697 0.3151 0.7288 0.3326 0.5141 0.2673 0.3118 100 0.9565 0.7962 0.9399 0.8605 0.4223 0.5476 0.9291 0.9208 0.9862 0.4508 0.8407 50 0.8161 0.9719 0.9694 0.8661 0.5923 0.5714 0.8412 0.6009 0.8974 0.8458 0.8022 0 1.0000 1.0000 1.0000 1.0000 0.8774 0.7944 0.8728 0.9198 0.5297 -50 1.0000 1.0000 1.0000 1.0000 1.0000 0.9999 0.8735 0.9184		ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
2500.87760.66950.72800.99560.87170.97300.65800.40230.83900.87750.21242000.81580.96650.39600.96570.59630.93920.93660.16070.15530.37210.43371500.98370.76680.91940.74050.86970.31510.72880.33260.51410.26730.31181000.95650.79620.93990.86050.42230.54760.92910.92080.98620.45080.8407500.81610.97190.96940.86610.59230.57140.84120.60090.89740.84580.802201.00001.00001.00001.00001.00000.87740.79440.87280.91980.5297-501.00001.00001.00001.00001.00000.97320.83670.77240.66900.9655-1001.00001.00001.00001.00001.00000.79430.78360.49610.56840.5923-2001.00001.00001.00001.00001.00000.94530.91370.95530.91850.9215-2501.00001.00001.00001.00001.00000.99230.94560.37020.91850.9215	Swaps												
2000.81580.96650.39600.96570.59630.93920.93660.16070.15530.37210.43371500.98370.76680.91940.74050.86970.31510.72880.33260.51410.26730.31181000.95650.79620.93990.86050.42230.54760.92910.92080.98620.45080.8407500.81610.97190.96940.86610.59230.57140.84120.60090.89740.84580.802201.00001.00001.00001.00001.00000.87740.79440.87280.91980.5297-501.00001.00001.00001.00001.00000.99990.87350.91840.65830.9466-1001.00001.00001.00001.00001.00000.97320.83670.77240.66900.9655-1501.00001.00001.00001.00001.00000.79430.78360.49610.56840.5923-2001.00001.00001.00001.00001.00000.94530.91370.95530.91850.9185-2501.00001.00001.00001.00001.00000.99230.94560.37020.91850.9215	250		0.8776	0.6695	0.7280	0.9956	0.8717	0.9730	0.6580	0.4023	0.8390	0.8775	0.2124
150 0.9837 0.7668 0.9194 0.7405 0.8697 0.3151 0.7288 0.3326 0.5141 0.2673 0.3118 100 0.9565 0.7962 0.9399 0.8605 0.4223 0.5476 0.9291 0.9208 0.9862 0.4508 0.8407 50 0.8161 0.9719 0.9694 0.8661 0.5923 0.5714 0.8412 0.6009 0.8974 0.8458 0.8022 0 1.0000 1.0000 1.0000 1.0000 1.0000 0.8774 0.7944 0.8728 0.9198 0.5297 -50 1.0000 1.0000 1.0000 1.0000 1.0000 0.9732 0.8367 0.7724 0.6690 0.9655 -100 1.0000 1.0000 1.0000 1.0000 1.0000 0.9732 0.8367 0.7724 0.6690 0.9655 -150 1.0000 1.0000 1.0000 1.0000 1.0000 0.9453 0.9137 0.9553 0.9650 0.8105	200		0.8158	0.9665	0.3960	0.9657	0.5963	0.9392	0.9366	0.1607	0.1553	0.3721	0.4337
100 0.9565 0.7962 0.9399 0.8605 0.4223 0.5476 0.9291 0.9208 0.9862 0.4508 0.8407 50 0.8161 0.9719 0.9694 0.8661 0.5923 0.5714 0.8412 0.6009 0.8974 0.8458 0.8022 0 1.0000 1.0000 1.0000 1.0000 1.0000 0.8774 0.7944 0.8728 0.9198 0.5297 -50 1.0000 1.0000 1.0000 1.0000 1.0000 0.9732 0.8367 0.7724 0.6690 0.9655 -100 1.0000 1.0000 1.0000 1.0000 1.0000 0.9732 0.8367 0.7724 0.6690 0.9655 -150 1.0000 1.0000 1.0000 1.0000 1.0000 0.7943 0.7836 0.4961 0.5684 0.5923 -200 1.0000 1.0000 1.0000 1.0000 1.0000 0.9453 0.9137 0.9553 0.9650 0.8105 -250	150		0.9837	0.7668	0.9194	0.7405	0.8697	0.3151	0.7288	0.3326	0.5141	0.2673	0.3118
50 0.8161 0.9719 0.9694 0.8661 0.5923 0.5714 0.8412 0.6009 0.8974 0.8458 0.8022 0 1.0000 1.0000 1.0000 1.0000 1.0000 0.8774 0.7944 0.8728 0.9198 0.5297 -50 1.0000 1.0000 1.0000 1.0000 1.0000 0.9999 0.8735 0.9184 0.6583 0.9466 -100 1.0000 1.0000 1.0000 1.0000 1.0000 0.9732 0.8367 0.7724 0.6690 0.9655 -150 1.0000 1.0000 1.0000 1.0000 1.0000 0.9733 0.7836 0.4961 0.5684 0.5923 -200 1.0000 1.0000 1.0000 1.0000 1.0000 0.9453 0.9137 0.9553 0.9650 0.8105 -250 1.0000 1.0000 1.0000 1.0000 0.9923 0.9456 0.3702 0.9185 0.9215	100		0.9565	0.7962	0.9399	0.8605	0.4223	0.5476	0.9291	0.9208	0.9862	0.4508	0.8407
0 1.0000 1.0000 1.0000 1.0000 1.0000 0.8774 0.7944 0.8728 0.9198 0.5297 -50 1.0000 1.0000 1.0000 1.0000 1.0000 0.9999 0.8735 0.9184 0.6583 0.9466 -100 1.0000 1.0000 1.0000 1.0000 1.0000 0.9732 0.8367 0.7724 0.6690 0.9655 -150 1.0000 1.0000 1.0000 1.0000 1.0000 0.7943 0.7836 0.4961 0.5684 0.5923 -200 1.0000 1.0000 1.0000 1.0000 1.0000 0.9453 0.9137 0.9553 0.9650 0.8105 -250 1.0000 1.0000 1.0000 1.0000 1.0000 0.9923 0.9456 0.3702 0.9185 0.9215	50		0.8161	0.9719	0.9694	0.8661	0.5923	0.5714	0.8412	0.6009	0.8974	0.8458	0.8022
-501.00001.00001.00001.00001.00000.99990.87350.91840.65830.9466-1001.00001.00001.00001.00001.00000.97320.83670.77240.66900.9655-1501.00001.00001.00001.00001.00000.79430.78360.49610.56840.5923-2001.00001.00001.00001.00001.00000.94530.91370.95530.96500.8105-2501.00001.00001.00001.00001.00000.99230.94560.37020.91850.9215	0		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.8774	0.7944	0.8728	0.9198	0.5297
-100 1.0000 1.0000 1.0000 1.0000 1.0000 0.9732 0.8367 0.7724 0.6690 0.9655 -150 1.0000 1.0000 1.0000 1.0000 1.0000 0.7943 0.7836 0.4961 0.5684 0.5923 -200 1.0000 1.0000 1.0000 1.0000 1.0000 0.9453 0.9137 0.9553 0.9650 0.8105 -250 1.0000 1.0000 1.0000 1.0000 0.9923 0.9456 0.3702 0.9185 0.9215	-50		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9999	0.8735	0.9184	0.6583	0.9466
-1501.00001.00001.00001.00001.00000.79430.78360.49610.56840.5923-2001.00001.00001.00001.00001.00001.00000.94530.91370.95530.96500.8105-2501.00001.00001.00001.00001.00001.00000.99230.94560.37020.91850.9215	-100		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9732	0.8367	0.7724	0.6690	0.9655
-2001.00001.00001.00001.00001.00000.94530.91370.95530.96500.8105-2501.00001.00001.00001.00001.00000.99230.94560.37020.91850.9215	-150		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.7943	0.7836	0.4961	0.5684	0.5923
-250 1.0000 1.0000 1.0000 1.0000 1.0000 0.9923 0.9456 0.3702 0.9185 0.9215	-200		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9453	0.9137	0.9553	0.9650	0.8105
	-250		1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9923	0.9456	0.3702	0.9185	0.9215

Table 8 Scenario 3 - RBO

							BM25					
	Replace- ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.4545	0.4661	0.4574	0.4709	0.5074	0.9932	0.9743	0.9534	0.9485	0.9491	0.9495
200		0.4677	0.4675	0.4533	0.4613	0.5241	0.9775	0.9718	0.9266	0.9425	0.9205	0.9515
150		0.4606	0.4708	0.4590	0.4716	0.5123	0.9897	0.9626	0.9645	0.9659	0.9501	0.9585
100		0.4596	0.4650	0.4504	0.4849	0.5085	0.9886	0.9817	0.9396	0.9502	0.9493	0.9619
50		0.4590	0.4641	0.4590	0.4714	0.5024	0.9885	0.9740	0.9587	0.9607	0.9417	0.9699
0		0.4713	0.4713	0.4713	0.4741	0.5206	1.0000	0.9829	0.9506	0.9448	0.9715	0.9638
-50		0.4713	0.4713	0.4713	0.4713	0.4817	0.5237	0.5010	0.4685	0.5018	0.4883	0.4926
-100		0.4713	0.4713	0.4713	0.4713	0.4713	0.4845	0.4677	0.4644	0.4341	0.4295	0.4527
-150		0.4713	0.4713	0.4713	0.4713	0.4713	0.4713	0.4406	0.4347	0.4306	0.4109	0.4404
-200		0.4713	0.4713	0.4713	0.4713	0.4713	0.4713	0.4213	0.4398	0.4256	0.4404	0.4201
-250		0.4713	0.4713	0.4713	0.4713	0.4713	0.4713	0.4421	0.4249	0.4489	0.4059	0.4110
							RM3					
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.5635	0.5699	0.5578	0.5959	0.6418	0.9928	0.9706	0.9627	0.9481	0.9332	0.9562
200		0.5598	0.5605	0.5642	0.5779	0.6430	0.9884	0.9649	0.9636	0.9426	0.9555	0.9703
150		0.5600	0.5593	0.5714	0.5682	0.6277	0.9940	0.9456	0.9720	0.9515	0.9447	0.9421
100		0.5617	0.5638	0.5571	0.5843	0.6443	0.9980	0.9775	0.9692	0.9536	0.9550	0.9492
50		0.5621	0.5676	0.5621	0.5579	0.6362	0.9830	0.9685	0.9812	0.9540	0.9604	0.9574
0		0.5652	0.5671	0.5798	0.5797	0.6643	1.0000	0.9892	0.9730	0.9692	0.9715	0.9489
-50		0.5652	0.5652	0.5697	0.5673	0.5821	0.6581	0.6426	0.6524	0.5867	0.6088	0.6074
-100		0.5652	0.5652	0.5652	0.5679	0.5738	0.5885	0.5751	0.5643	0.5467	0.5271	0.5232
								0 5460	0 5410	0 5460	0 = 100	0 5010
-150		0.5652	0.5652	0.5652	0.5652	0.5655	0.5700	0.5462	0.5418	0.5469	0.5499	0.5218
-150 -200		0.5652 0.5652	0.5652 0.5652	0.5652 0.5652	0.5652 0.5652	0.5655 0.5652	0.5700 0.5704	0.5462 0.5389	0.5418	0.5469 0.5303	0.5499 0.5294	0.5218 0.5087

Table 9Scenario 1 - RMSE - AP

						Pe	rfect rank	ing				
	Replace- ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		1.0000	1.0000	1.0000	1.0000	0.7254	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
200		1.0000	1.0000	1.0000	1.0000	0.7503	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
150		1.0000	1.0000	1.0000	1.0000	0.7160	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
100		1.0000	1.0000	1.0000	1.0000	0.6961	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
50		1.0000	1.0000	1.0000	1.0000	0.7262	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0		1.0000	1.0000	1.0000	1.0000	0.7179	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-50		0.9981	0.9973	0.9960	0.9931	0.9845	0.7123	0.6874	0.7077	0.7154	0.0881	0.0745
-100		0.9945	0.9933	0.9902	0.9846	0.9708	0.9389	0.9350	0.9373	0.9383	0.9387	0.9373
-150		0.9909	0.9887	0.9840	0.9778	0.9058	0.9381	0.9378	0.9377	0.9300	0.9382	0.9381
-200		0.9003	0.9043	0.9001	0.9714	0.9002	0.9309	0.9305	0.9307	0.9365	0.9394	0.9361
-230		0.9040	0.9001	0.9701	0.9092	0.9301	0.9570	0.9590	0.9519	0.9501	0.9302	0.9500
						Rea	alistic rank	king				
	Replace-											
~	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.7081	0.7081	0.7081	0.7081	0.5415	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
200		0.7081	0.7081	0.7081	0.7081	0.5436	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
150		0.7081	0.7081	0.7081	0.7081	0.5263	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
100		0.7081	0.7081	0.7081	0.7081	0.5053	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
50		0.7081	0.7081	0.7081	0.7081	0.5310	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0		0.7081	0.7081	0.7081	0.7081	0.5355	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-50		0.7003	0.7054	0.7041	0.7010	0.0910	0.4828	0.4842	0.4502	0.4925	0.4008	0.4947
-100		0.7029	0.7012	0.0962	0.0929	0.0010	0.0447	0.0439	0.0449	0.0404	0.0451	0.0450
-200		0.0991	0.0904	0.6917	0.0007	0.6680	0.0402	0.0441	0.0474	0.0400	0.0459	0.0434
-250		0.6929	0.6893	0.6833	0.6765	0.6651	0.6435	0.6447	0.6467	0.6464	0.6456	0.6455
						Rev	ersed ran	king				
						T(C)		Killg				
	Replace-	050	200	150	100	50	0	50	100	150	000	050
Curana	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.1508	0.1528	0.1656	0.1377	0.1981	0.1799	0.1989	0.1733	0.1402	0.1310	0.1512
200		0.1329	0.1600	0.1898	0.1439	0.1658	0.1603	0.1587	0.1701	0.1851	0.1647	0.1493
150		0.1487	0.1277	0.1392	0.1435	0.1577	0.1420	0.1490	0.1712	0.1680	0.1564	0.2059
100		0.1496	0.1/41	0.1398	0.1439	0.1770	0.1408	0.1593	0.1383	0.1438	0.1763	0.1731
5U 0		0.0380	0.04/3	0.0315	0.0347	0.0303	0.0400	0.0741	0.0301	0.0335	0.0397	0.0440
50		0.0000		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-50				0.0000	0.0000	0.0000	0.0000		0.0000			0.0000
-150		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

Table 10Scenario 2 - RMSE - AP

						Pe	rfect rank	ing				
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.5000	0.5000	0.5000	0.5000	0.3709	0.0000	0.1352	0.2692	0.2808	0.2722	0.2653
200		0.5000	0.5000	0.5000	0.5000	0.3670	0.0000	0.1050	0.2602	0.2614	0.2804	0.2819
150		0.5000	0.5000	0.5000	0.5000	0.3717	0.0000	0.1114	0.2650	0.2731	0.2739	0.2747
100		0.5000	0.5000	0.5000	0.5000	0.3819	0.0000	0.1186	0.2721	0.2845	0.2861	0.2834
50		0.5000	0.5000	0.5000	0.5000	0.3707	0.0000	0.1122	0.2826	0.2905	0.2764	0.2730
0		0.5000	0.5000	0.5000	0.5000	0.3743	0.0000	0.1192	0.2623	0.2970	0.2879	0.2771
-50		0.4990	0.4986	0.4980	0.4965	0.4918	0.3507	0.2654	0.1085	0.1284	0.1476	0.1411
-100		0.4970	0.4965	0.4948	0.4921	0.4860	0.4691	0.4146	0.3201	0.3268	0.3232	0.3154
-150		0.4955	0.4943	0.4918	0.4881	0.4826	0.4687	0.4158	0.3226	0.3249	0.3279	0.3209
-200		0.4940	0.4921	0.4899	0.4858	0.4799	0.4691	0.4159	0.3283	0.3238	0.3223	0.3252
-250		0.4923	0.4902	0.4881	0.4843	0.4792	0.4080	0.4119	0.3279	0.3232	0.3273	0.3200
						Rea	alistic rank	king				
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.3541	0.3541	0.3541	0.3541	0.2633	0.0000	0.1287	0.3110	0.2990	0.2919	0.3258
200		0.3541	0.3541	0.3541	0.3541	0.2549	0.0000	0.1495	0.2904	0.2961	0.2972	0.3254
150		0.3541	0.3541	0.3541	0.3541	0.2597	0.0000	0.1287	0.3087	0.3091	0.3136	0.3050
100		0.3541	0.3541	0.3541	0.3541	0.2610	0.0000	0.1592	0.3460	0.3302	0.2807	0.3276
50		0.3541	0.3541	0.3541	0.3541	0.2587	0.0000	0.1356	0.3067	0.3047	0.2917	0.2769
0		0.3541	0.3541	0.3541	0.3541	0.2704	0.0000	0.1247	0.2997	0.3000	0.3102	0.3158
-50		0.3531	0.3527	0.3519	0.3506	0.3460	0.2455	0.1412	0.0410	0.0250	0.0294	0.0325
-100		0.3514	0.3503	0.3489	0.3464	0.3404	0.3235	0.2683	0.1694	0.1671	0.1653	0.1782
-150		0.3494	0.3482	0.3463	0.3425	0.3367	0.3233	0.2688	0.1758	0.1662	0.1664	0.1737
-200		0.3477	0.3460	0.3438	0.3402	0.3339	0.3224	0.2691	0.1822	0.1798	0.1694	0.1680
-250		0.3464	0.3442	0.3420	0.3379	0.3327	0.3231	0.2669	0.1787	0.1748	0.1716	0.1659
						Rev	ersed ran	king				
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.0658	0.0880	0.0762	0.0600	0.0718	0.0657	0.2210	0.4197	0.3361	0.3639	0.4014
200		0.0851	0.0626	0.1092	0.0740	0.0675	0.0760	0.1882	0.4067	0.3932	0.4081	0.3685
150		0.0779	0.0836	0.0734	0.0713	0.0660	0.0816	0.2166	0.4048	0.3835	0.3575	0.4052
100		0.0925	0.0821	0.0888	0.0762	0.0966	0.0850	0.1804	0.3472	0.3809	0.4038	0.4272
50		0.0369	0.0186	0.0232	0.0264	0.0307	0.0251	0.1081	0.2836	0.2470	0.2418	0.2314
0		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0561	0.1423	0.1536	0.1603	0.1739
-50		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0543	0.1720	0.1443	0.1664	0.1565
-100		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0517	0.1471	0.1662	0.1776	0.1550
-150		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0568	0.1541	0.1578	0.1749	0.1720
-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0566	0.1557	0.1477	0.1650	0.1596
-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0506	0.1461	0.1739	0.1454	0.1641

Table 11 Scenario 3 - RMSE - AP

							BM25					
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.3126	0.3126	0.3124	0.3110	0.3013	0.0262	0.1012	0.1449	0.1666	0.1697	0.1699
200		0.3127	0.3125	0.3124	0.3108	0.3001	0.0279	0.1034	0.1458	0.1684	0.1768	0.1678
150		0.3125	0.3126	0.3125	0.3112	0.3003	0.0254	0.1107	0.1412	0.1609	0.1677	0.1717
100		0.3126	0.3126	0.3126	0.3109	0.2996	0.0252	0.1047	0.1470	0.1659	0.1685	0.1662
50		0.3125	0.3125	0.3122	0.3112	0.3015	0.0227	0.1017	0.1400	0.1618	0.1616	0.1636
0		0.3135	0.3135	0.3135	0.3126	0.3029	0.0000	0.0803	0.1138	0.1346	0.1340	0.1368
-50		0.3124	0.3121	0.3115	0.3104	0.3062	0.2837	0.2745	0.2785	0.2809	0.2836	0.2839
-100		0.3110	0.3104	0.3092	0.3075	0.3036	0.2928	0.2823	0.2839	0.2884	0.2901	0.2896
-150		0.3097	0.3088	0.3074	0.3053	0.3013	0.2933	0.2832	0.2849	0.2899	0.2896	0.2894
-200		0.3086	0.3074	0.3059	0.3036	0.2998	0.2933	0.2827	0.2845	0.2893	0.2893	0.2902
-250		0.3074	0.3062	0.3047	0.3023	0.2990	0.2933	0.2831	0.2851	0.2881	0.2906	0.2897
							RM3					
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.5379	0.5379	0.5372	0.5318	0.4984	0.0144	0.1726	0.2133	0.2345	0.2390	0.2393
200		0.5381	0.5374	0.5358	0.5318	0.4980	0.0165	0.1703	0.2157	0.2367	0.2416	0.2349
150		0.5369	0.5375	0.5368	0.5322	0.4992	0.0161	0.1791	0.2179	0.2304	0.2365	0.2389
100		0.5371	0.5371	0.5374	0.5326	0.4939	0.0142	0.1717	0.2150	0.2341	0.2368	0.2322
50		0.5368	0.5368	0.5355	0.5317	0.5003	0.0129	0.1682	0.2101	0.2333	0.2349	0.2356
0		0.5438	0.5438	0.5438	0.5382	0.5048	0.0000	0.1621	0.2054	0.2257	0.2247	0.2276
-50		0.5039	0.4929	0.4805	0.4601	0.4135	0.2836	0.2618	0.2833	0.2903	0.2955	0.2973
-100		0.4715	0.4614	0.4416	0.4180	0.3778	0.2950	0.2670	0.2821	0.2999	0.3023	0.3002
-150		0.4496	0.4364	0.4170	0.3931	0.3563	0.2963	0.2673	0.2853	0.2999	0.2995	0.2982
-200		0.4327	0.4166	0.4018	0.3776	0.3440	0.2965	0.2687	0.2845	0.2994	0.2989	0.3025
-250		0.4174	0.4034	0.3887	0.3667	0.3366	0.2963	0.2691	0.2873	0.2965	0.3035	0.3002

Table 12 Scenario 1 - nRMSE - AP

						Pe	rfect rank	ing				
	Replace- ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		1.0000	1.0000	1.0000	1.0000	0.7254	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
200		1.0000	1.0000	1.0000	1.0000	0.7503	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
150		1.0000	1.0000	1.0000	1.0000	0.7160	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
100		1.0000	1.0000	1.0000	1.0000	0.6961	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
50		1.0000	1.0000	1.0000	1.0000	0.7262	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0		1.0000	1.0000	1.0000	1.0000	0.7179	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-50		0.9981	0.9973	0.9960	0.9931	0.9845	0.7123	0.6874	0.7077	0.7154	0.0881	0.0745
-100		0.9945	0.9933	0.9902	0.9846	0.9708	0.9389	0.9350	0.9373	0.9383	0.9387	0.9373
-150		0.9909	0.9887	0.9840	0.9778	0.9058	0.9381	0.9378	0.9377	0.9300	0.9382	0.9381
-200		0.9003	0.9043	0.9001	0.9714	0.9002	0.9309	0.9305	0.9307	0.9365	0.9394	0.9361
-230		0.9040	0.9001	0.9701	0.9092	0.9301	0.9510	0.9590	0.9519	0.9501	0.9302	0.9500
						Rea	alistic rank	king				
	Replace-											
~	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		1.0000	1.0000	1.0000	1.0000	0.7647	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
200		1.0000	1.0000	1.0000	1.0000	0.7677	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
150		1.0000	1.0000	1.0000	1.0000	0.7433	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
100		1.0000	1.0000	1.0000	1.0000	0.7136	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
50		1.0000	1.0000	1.0000	1.0000	0.7499	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0		1.0000	1.0000	1.0000	1.0000	0.7503	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-50		0.9975	0.9902	0.9943	0.9900	0.9707	0.0817	0.0838	0.0358	0.0955	0.0593	0.0980
-100		0.9927	0.9902	0.9001	0.9765	0.9029	0.9105	0.9122	0.9107	0.9120	0.9111	0.9117
-200		0.9075	0.9033	0.9700	0.9005	0.9301	0.9120	0.9095	0.9145	0.9125	0.9122	0.9113
-250		0.9785	0.9735	0.9650	0.9554	0.9393	0.9087	0.9105	0.9133	0.9128	0.9118	0.9110
						Rev	ersed ran	king				
						T(C)		Killg				
	Replace-	050	200	150	100	50	0	50	100	150	000	050
Swane	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.1591	0.1612	0.1748	0.1453	0.2090	0.1898	0.2098	0.1828	0.1479	0.1383	0.1595
200		0.1402	0.1688	0.2003	0.1519	0.1749	0.1691	0.1674	0.1794	0.1953	0.1738	0.1576
150		0.1569	0.1347	0.1468	0.1514	0.1663	0.1498	0.1572	0.1806	0.1772	0.1650	0.2172
100		0.1579	0.183/	0.14/5	0.1519	0.100/	0.1480	0.1001	0.1459	0.0323	0.1800	0.1827
0		0.0401	0.0499	0.0333	0.0300	0.0000	0.0405	0.0702	0.0402	0.0333	0.0419	0.0404
-50				0.0000	0.0000	0.0000	0.0000		0.0000			0.0000
- <u>10</u> 0		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-150		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

Table 13Scenario 2 - nRMSE - AP

		Perfect ranking										
	Replace- ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		1.0000	1.0000	1.0000	1.0000	0.7419	0.0000	0.2704	0.5384	0.5615	0.5444	0.5307
200		1.0000	1.0000	1.0000	1.0000	0.7340	0.0000	0.2101	0.5203	0.5228	0.5607	0.5637
150		1.0000	1.0000	1.0000	1.0000	0.7434	0.0000	0.2228	0.5301	0.5463	0.5477	0.5493
100		1.0000	1.0000	1.0000	1.0000	0.7639	0.0000	0.2372	0.5443	0.5690	0.5722	0.5667
50		1.0000	1.0000	1.0000	1.0000	0.7413	0.0000	0.2244	0.5652	0.5810	0.5527	0.5460
0		1.0000	1.0000	1.0000	1.0000	0.7486	0.0000	0.2384	0.5245	0.5939	0.5759	0.5541
-50		0.9980	0.9973	0.9961	0.9931	0.9836	0.7015	0.5308	0.2170	0.2568	0.2953	0.2821
-100		0.9940	0.9931	0.9895	0.9841	0.9719	0.9382	0.8293	0.6402	0.6537	0.6464	0.6308
-150		0.9909	0.9885	0.9837	0.9702	0.9051	0.9373	0.0310	0.0451	0.0498	0.0558	0.0418
-200		0.9001	0.9042	0.9790	0.9710	0.9599	0.9301	0.0310	0.0507	0.0470	0.0447	0.0504
-230		0.9040	0.9005	0.9702	0.9000	0.9303	0.9512		0.0339	0.0403	0.0340	0.0332
		Realistic ranking										
	Replace-											
~	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.5481	0.5481	0.5481	0.5481	0.4076	0.0000	0.1992	0.4814	0.4629	0.4518	0.5044
200		0.5481	0.5481	0.5481	0.5481	0.3947	0.0000	0.2315	0.4495	0.4583	0.4601	0.5038
150		0.5481	0.5481	0.5481	0.5481	0.4020	0.0000	0.1992	0.4780	0.4785	0.4855	0.4722
100		0.5481	0.5481	0.5481	0.5481	0.4040	0.0000	0.2465	0.5356	0.5112	0.4345	0.5072
50		0.5481	0.5481	0.5481	0.5481	0.4005	0.0000	0.2099	0.4748	0.4717	0.4516	0.4287
0		0.5481	0.5481	0.5481	0.5481	0.4185	0.0000	0.1930	0.4640	0.4645	0.4802	0.4889
-50		0.5400	0.5400	0.5448	0.5428	0.5357	0.3801	0.2180	0.0035	0.0387	0.0454	0.0503
-100		0.5440	0.5424	0.5402	0.5302	0.5270	0.5008	0.4154	0.2022	0.2507	0.2556	0.2759
-130		0.5400	0.5391	0.5301	0.5302	0.5215	0.3000	0.4102	0.2721	0.2373	0.2570	0.2090
-250		0.5363	0.5329	0.5294	0.5232	0.5150	0.5003	0.4132	0.2766	0.2707	0.2657	0.2568
						Pau	orcod ran	ling				
						Nev	erseu ran	king				
\sim	Replace-	050	000	150	100	F.0	~	F.0	100	150	000	050
Curana	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.0675	0.0903	0.0782	0.0616	0.0737	0.0675	0.2270	0.4310	0.3451	0.3737	0.4121
200		0.0874	0.0643	0.1121	0.0760	0.0693	0.0781	0.1932	0.4176	0.4037	0.4191	0.3784
150		0.0800	0.0859	0.0754	0.0732	0.0678	0.0838	0.2224	0.4156	0.3938	0.3671	0.4160
100		0.0950	0.0844	0.0912	0.0783	0.0992	0.0873	0.1853	0.3565	0.3911	0.4147	0.4387
5U 0		0.0379	0.0191	0.0238	0.0271	0.0315	0.0258	0.1109	0.2912	0.253/	0.2483	0.23/0
U 50		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0570	0.1401	0.1578	0.1040	0.1780
-50				0.0000	0.0000	0.0000	0.0000	0.0557	0.1700	0.1402	0.1700	0.1007
-150		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0551	0.1520	0.1620	0.1024	0.1391
-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0503	0.1502	0.1020	0.1790	0.1630
-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0520	0 1501	0 1786	0 1493	0 1685
200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0020	0.1001	0.1100	0.1455	0.1000

Table 14Scenario 3 - nRMSE - AP

							BM25					
	Replace- ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.3882	0.3882	0.3879	0.3862	0.3741	0.0325	0.1257	0.1799	0.2069	0.2107	0.2110
200		0.3883	0.3881	0.3879	0.3860	0.3727	0.0347	0.1284	0.1810	0.2091	0.2196	0.2084
150		0.3881	0.3882	0.3881	0.3864	0.3729	0.0316	0.1374	0.1753	0.1999	0.2083	0.2133
100		0.3882	0.3882	0.3882	0.3861	0.3720	0.0313	0.1300	0.1826	0.2060	0.2093	0.2063
50		0.3881	0.3881	0.3878	0.3865	0.3744	0.0282	0.1263	0.1738	0.2009	0.2007	0.2032
0		0.3893	0.3893	0.3893	0.3882	0.3761	0.0000	0.0997	0.1414	0.1672	0.1664	0.1699
-50		0.3880	0.3876	0.3868	0.3854	0.3803	0.3524	0.3409	0.3458	0.3488	0.3522	0.3525
-100		0.3863	0.3854	0.3840	0.3818	0.3770	0.3636	0.3506	0.3526	0.3582	0.3603	0.3596
-150		0.3846	0.3835	0.3817	0.3792	0.3742	0.3642	0.3516	0.3538	0.3600	0.3596	0.3594
-200		0.3833	0.3818	0.3799	0.3770	0.3723	0.3642	0.3510	0.3533	0.3593	0.3593	0.3604
-250		0.3817	0.3803	0.3784	0.3754	0.3713	0.3642	0.3516	0.3540	0.3578	0.3609	0.3598
							RM3					
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.4582	0.4567	0.4530	0.4427	0.4135	0.0343	0.1167	0.1374	0.1630	0.1708	0.1776
200		0.4581	0.4566	0.4523	0.4438	0.4135	0.0363	0.1165	0.1419	0.1601	0.1694	0.1765
150		0.4581	0.4565	0.4526	0.4435	0.4114	0.0306	0.1227	0.1352	0.1653	0.1710	0.1825
100		0.4583	0.4563	0.4523	0.4428	0.4153	0.0339	0.1169	0.1344	0.1625	0.1644	0.1776
50		0.4584	0.4564	0.4529	0.4446	0.4101	0.0365	0.1127	0.1275	0.1597	0.1623	0.1760
0		0.4594	0.4580	0.4536	0.4460	0.4158	0.0000	0.0815	0.1050	0.1277	0.1392	0.1548
-50		0.4575	0.4570	0.4543	0.4497	0.4367	0.3900	0.3786	0.3842	0.3900	0.3899	0.3975
-100		0.4552	0.4540	0.4520	0.4472	0.4376	0.4147	0.4035	0.4066	0.4115	0.4147	0.4186
-150		0.4529	0.4512	0.4488	0.4451	0.4368	0.4210	0.4097	0.4128	0.4181	0.4211	0.4240
-200		0.4507	0.4489	0.4461	0.4420	0.4355	0.4230	0.4127	0.4153	0.4203	0.4225	0.4277
-250		0.4488	0.4468	0.4438	0.4397	0.4338	0.4238	0.4142	0.4177	0.4214	0.4248	0.4286

Table 15

Scenario 3 - p-values - AP

							BM25					
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.0000	0.0000	0.0000	0.0000	0.0000	0.7578	0.0308	0.0013	0.0021	0.0083	0.0030
200		0.0000	0.0000	0.0000	0.0000	0.0000	0.6860	0.0572	0.0088	0.0096	0.0031	0.0064
150		0.0000	0.0000	0.0000	0.0000	0.0000	0.7157	0.0526	0.0088	0.0124	0.0105	0.0051
100		0.0000	0.0000	0.0000	0.0000	0.0000	0.6861	0.0443	0.0069	0.0046	0.0132	0.0025
50		0.0000	0.0000	0.0000	0.0000	0.0000	0.7343	0.0501	0.0171	0.0082	0.0036	0.0060
0		0.0000	0.0000	0.0000	0.0000	0.0000	1.0000	0.1148	0.0095	0.0411	0.0389	0.0389
-50		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-100		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-150		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
							RM3					
	Replace-											
	ments	-250	-200	-150	-100	-50	0	50	100	150	200	250
Swaps												
250		0.0000	0.0000	0.0000	0.0000	0.0000	0.7396	0.1037	0.1164	0.0996	0.0635	0.0940
200		0.0000	0.0000	0.0000	0.0000	0.0000	0.7496	0.1946	0.1264	0.0733	0.0737	0.0815
150		0.0000	0.0000	0.0000	0.0000	0.0000	0.7698	0.1896	0.1583	0.0787	0.0519	0.0430
100		0.0000	0.0000	0.0000	0.0000	0.0000	0.7702	0.2478	0.1263	0.0945	0.0462	0.0498
50		0.0000	0.0000	0.0000	0.0000	0.0000	0.7325	0.1986	0.1450	0.0850	0.0943	0.0916
0		0.0000	0.0000	0.0000	0.0000	0.0000	1.0000	0.3217	0.2514	0.1759	0.1777	0.1234
-50		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-100		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-150		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-200		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
-250		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

References

- Agosti, M., Di Nunzio, G.M., Ferro, N., Silvello, G., 2019. An Innovative Approach to Data Management and Curation of Experimental Data Generated through IR Test Collections, in: Ferro and Peters (2019). pp. 105–122. pp. 105–122.
- Agosti, M., Ferro, N., Thanos, C., 2012. DESIRE 2011 Workshop on Data infrastructurEs for Supporting Information Retrieval Evaluation. SIGIR Forum 46, 51–55.
- Allan, J., Aslam, J.A., Pavlu, V., Kanoulas, E., Carterette, B., 2009. Million Query Track 2008 Overview, in: Voorhees, E.M., Buckland, L.P. (Eds.), The Seventeenth Text REtrieval Conference Proceedings (TREC 2008), National Institute of Standards and Technology (NIST), Special Publication 500-277, Washington, USA.
- Allan, J., Carterette, B., Dachev, B., Aslam, J.A., Pavlu, V., Kanoulas, E., 2008. Million Query Track 2007 Overview, in: Voorhees, E.M., Buckland, L.P. (Eds.), The Sixteenth Text REtrieval Conference Proceedings (TREC 2007), National Institute of Standards and Technology (NIST), Special Publication 500-274, Washington, USA.
- Arampatzis, A., van Hameren, A., 2001. The Score Distributional Threshold Optimization for Adaptive Binary Classification Tasks, in: Kraft, Croft, Harper and Zobel (2001). pp. 285–293. pp. 285–293.
- Arampatzis, A., Kamps, J., 2009. A Signal-to-Noise Approach to Score Normalization, in: Cheung, Song, Chu, Hu and Lin (2009). pp. 797–806. pp. 797–806.
- Arampatzis, A., Robertson, S.E., 2011. Modeling score distributions in information retrieval. Information Retrieval 14, 26-46.
- Arguello, J., Crane, M., Diaz, F., Lin, J., Trotman, A., 2015. Report on the SIGIR 2015 Workshop on Reproducibility, Inexplicability, and Generalizability of Results (RIGOR). SIGIR Forum 49, 107–116.
- Armstrong, T.G., Moffat, A., Webber, W., Zobel, J., 2009. Improvements That Don't Add Up: Ad-Hoc Retrieval Results Since 1998, in: Cheung et al. (2009). pp. 601–610. pp. 601–610.
- Baker, M., 2016. 1,500 Scientists Lift the Lid on Reproducibility. Nature 533, 452-454.
- Breuer, T., Ferro, N., Fuhr, N., Maistro, M., Sakai, T., Schaer, P., Soboroff, I., 2020. How to Measure the Reproducibility of System-oriented IR Experiments, in: Chang, Y., Cheng, X., Huang, J., Lu, Y., Kamps, J., Murdock, V., Wen, J.R., Diriye, A., Guo, J., Kurland, O. (Eds.), Proc. 43rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020), ACM Press, New York, USA. pp. 349–358.
- Breuer, T., Ferro, N., Maistro, M., Schaer, P., 2021. repro_eval: A Python Interface to Reproducibility Measures of System-Oriented IR Experiments, in: Hiemstra, D., Moens, M., Mothe, J., Perego, R., Potthast, M., Sebastiani, F. (Eds.), Advances in Information Retrieval - 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28 - April 1, 2021, Proceedings, Part II, Springer. pp. 481–486. URL: https://doi.org/10.1007/978-3-030-72240-1_51, doi:10.1007/978-3-030-72240-1_51.
- Breuer, T., Keller, J., Schaer, P., 2022. ir_metadata: An extensible metadata schema for IR experiments, in: Amigó, E., Castells, P., Gonzalo, J., Carterette, B., Culpepper, J.S., Kazai, G. (Eds.), SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022, ACM. pp. 3078–3089. URL: https://doi.org/10.1145/3477495.3531738, doi:10.1145/3477495.3531738.
- Breuer, T., Schaer, P., 2021. A living lab architecture for reproducible shared task experimentation, in: Wolff, C., Schmidt, T. (Eds.), Information between Data and Knowledge: Information Science and its Neighbors from Data Science to Digital Humanities - Proceedings of the 16th International Symposium of Information Science, ISI 2021, Regensburg, Germany, March 8-10, 2021, Werner Hülsbusch. pp. 348–362. URL: https://doi.org/10.5283/epub.44953, doi:10.5283/epub.44953.
- Carterette, B., Pavlu, V., Fang, H., Kanoulas, E., 2010. Million Query Track 2009 Overview, in: Voorhees, E.M., Buckland, L.P. (Eds.), The Eighteenth Text REtrieval Conference Proceedings (TREC 2009), National Institute of Standards and Technology (NIST), Special Publication 500-278, Washington, USA.
- Cheung, D.W.L., Song, I.Y., Chu, W.W., Hu, X., Lin, J.J. (Eds.), 2009. Proc. 18th International Conference on Information and Knowledge Management (CIKM 2009), ACM Press, New York, USA.
- Clancy, R., Ferro, N., Hauff, C., Sakai, T., Wu, Z.Z., 2019. The SIGIR 2019 Open-Source IR Replicability Challenge (OSIRRC 2019), in: Piwowarski, B., Chevalier, M., Gaussier, E., Maarek, Y., Nie, J.Y., Scholer, F. (Eds.), Proc. 42nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2019), ACM Press, New York, USA. pp. 1432–1434.
- Cleverdon, C.W., 1962. Report on the Testing and Analysis of an Investigation into the Comparative Efficiency of Indexing Systems. Aslib Cranfield Research Project, College of Aeronautics, Cranfield, UK.
- Cleverdon, C.W., 1967. The Cranfield Tests on Index Languages Devices. Aslib Proceedings 19, 173–194.
- Cummins, R., 2014. Document Score Distribution Models for Query Performance Inference and Prediction. ACM Transactions on Information System (TOIS) 32, 2:1–2:28.
- Dacrema, M.F., Boglio, S., Cremonesi, P., Jannach, D., 2021. A Troubling Analysis of Reproducibility and Progress in Recommender Systems Research. ACM Transactions on Information Systems 39, 20:1–20:49.
- Dacrema, M.F., Cremonesi, P., Jannach, D., 2019. Are we really making much progress? A worrying analysis of recent neural recommendation approaches, in: Bogers, T., Said, A., Brusilovsky, P., Tikk, D. (Eds.), Proceedings of the 13th ACM Conference on Recommender Systems, RecSys 2019, Copenhagen, Denmark, September 16-20, 2019, ACM. pp. 101–109. URL: https://doi.org/10.1145/3298689.3347058, doi:10.1145/3298689.3347058.
- De Roure, D., 2014. The future of scholarly communications. Insights 27, 233-238.
- Ferrante, M., Ferro, N., Maistro, M., 2015. Towards a Formal Framework for Utility-oriented Measurements of Retrieval Effectiveness, in: Allan, J., Croft, W.B., de Vries, A.P., Zhai, C., Fuhr, N., Zhang, Y. (Eds.), Proc. 1st ACM SIGIR International Conference on the Theory of Information Retrieval (ICTIR 2015), ACM Press, New York, USA. pp. 21–30.
- Ferro, N., 2017. Reproducibility Challenges in Information Retrieval Evaluation. ACM Journal of Data and Information Quality (JDIQ) 8, 8:1-8:4.
- Ferro, N., Fuhr, N., Järvelin, K., Kando, N., Lippold, M., Zobel, J., 2016. Increasing Reproducibility in IR: Findings from the Dagstuhl Seminar on "Reproducibility of Data-Oriented Experiments in e-Science". SIGIR Forum 50, 68–82.

- Ferro, N., Fuhr, N., Maistro, M., Sakai, T., Soboroff, I., 2019a. CENTRE@CLEF2019: Overview of the Replicability and Reproducibility Tasks, in: Cappellato, L., Ferro, N., Losada, D.E., Müller, H. (Eds.), CLEF 2019 Working Notes, CEUR Workshop Proceedings (CEUR-WS.org), ISSN 1613-0073, http://ceur-ws.org/Vol-2380/.
- Ferro, N., Fuhr, N., Maistro, M., Sakai, T., Soboroff, I., 2019b. Overview of CENTRE@CLEF 2019: Sequel in the Systematic Reproducibility Realm, in: Crestani, F., Braschler, M., Savoy, J., Rauber, A., Müller, H., Losada, D.E., Heinatz Bürki, G., Cappellato, L., Ferro, N. (Eds.), Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Tenth International Conference of the CLEF Association (CLEF 2019), Lecture Notes in Computer Science (LNCS) 11696, Springer, Heidelberg, Germany. pp. 287–300.
- Ferro, N., Kelly, D., 2018. SIGIR Initiative to Implement ACM Artifact Review and Badging. SIGIR Forum 52, 4-10.
- Ferro, N., Maistro, M., Sakai, T., Soboroff, I., 2018. Overview of CENTRE@CLEF 2018: a First Tale in the Systematic Reproducibility Realm, in: Bellot, P., Trabelsi, C., Mothe, J., Murtagh, F., Nie, J.Y., Soulier, L., SanJuan, E., Cappellato, L., Ferro, N. (Eds.), Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Nineth International Conference of the CLEF Association (CLEF 2018), Lecture Notes in Computer Science (LNCS) 11018, Springer, Heidelberg, Germany. pp. 239–246.
- Ferro, N., Peters, C. (Eds.), 2019. Information Retrieval Evaluation in a Changing World Lessons Learned from 20 Years of CLEF. volume 41 of *The Information Retrieval Series*, Springer International Publishing, Germany.
- Freire, J., Fuhr, N., Rauber, A. (Eds.), 2016. Report from Dagstuhl Seminar 16041: Reproducibility of Data-Oriented Experiments in e-Science. Dagstuhl Reports, Volume 6, Number 1, Schloss Dagstuhl-Leibniz-Zentrum für Informatik, Germany.
- Gibney, E., 2020. This AI Researcher is Trying to Ward off a Reproducibility Crisis. Nature 577, 14.
- Gollub, T., Stein, B., Burrows, S., 2012a. Ousting Ivory Tower Research: Towards a Web Framework for Providing Experiments as a Service, in: Hersh, Callan, Maarek and Sanderson (2012), pp. 1125–1126. pp. 1125–1126.
- Gollub, T., Stein, B., Burrows, S., Hoppe, D., 2012b. TIRA: Configuring, Executing, and Disseminating Information Retrieval Experiments, in: Hameurlain, A., Tjoa, A.M., Wagner, R.R. (Eds.), Proc. 23rd International Workshop on Database and Expert Systems Applications (DEXA 2012), IEEE Computer Society. pp. 151–155.
- Hersh, W., Callan, J., Maarek, Y., Sanderson, M. (Eds.), 2012. Proc. 35th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2012), ACM Press, New York, USA.
- Hopfgartner, F., Hanbury, A., Müller, H., Kando, N., Mercer, S., Kalpathy-Cramer, J., Potthast, M., Gollub, T., Krithara, A., Lin, J., Balog, K., Eggel, I., 2015. Report on the Evaluation-as-a-Service (EaaS) Expert Workshop. SIGIR Forum 49, 57–65.
- Ivie, P., Thain, D., 2018. Reproducibility in scientific computing. ACM Comput. Surv. 51, 63:1–63:36. URL: https://doi.org/10.1145/31 86266, doi:10.1145/3186266.
- Jaleel, N.A., Allan, J., Croft, W.B., Diaz, F., Larkey, L.S., Li, X., Smucker, M.D., Wade, C., 2004. Umass at TREC 2004: Novelty and HARD, in: Voorhees, E.M., Buckland, L.P. (Eds.), Proceedings of the Thirteenth Text REtrieval Conference, TREC 2004, Gaithersburg, Maryland, USA, November 16-19, 2004, National Institute of Standards and Technology (NIST). URL: http://trec.nist.gov/pubs/trec13/papers/u mass.novelty.hard.pdf.
- Järvelin, K., Kekäläinen, J., 2002. Cumulated Gain-Based Evaluation of IR Techniques. ACM Transactions on Information Systems (TOIS) 20, 422–446.
- Kendall, M.G., 1945. The Treatment of Ties in Ranking Problems. Biometrika 33, 239-251.
- Kraft, D.H., Croft, W.B., Harper, D.J., Zobel, J. (Eds.), 2001. Proc. 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2001), ACM Press, New York, USA.
- Lavrenko, V., Croft, W.B., 2001. Relevance-Based Language Models, in: Kraft et al. (2001). pp. 120-127. pp. 120-127.
- Lin, J., 2018. The Neural Hype and Comparisons against Weak Baselines. SIGIR Forum 52, 40-51.
- Lin, J., Ma, X., Lin, S., Yang, J., Pradeep, R., Nogueira, R., 2021. Pyserini: A python toolkit for reproducible information retrieval research with sparse and dense representations, in: Diaz, F., Shah, C., Suel, T., Castells, P., Jones, R., Sakai, T. (Eds.), SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021, ACM. pp. 2356–2362. URL: https://doi.org/10.1145/3404835.3463238, doi:10.1145/3404835.3463238.
- Lucic, A., Bleeker, M.J.R., de Rijke, M., Sinha, K., Jullien, S., Stojnic, R., 2022. Towards Reproducible Machine Learning Research in Information Retrieval, in: Amigó, E., Castells, P., Gonzalo, J., Carterette, B., Culpepper, J.S., Kazai, G. (Eds.), Proc. 45th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2022), ACM Press, New York, USA. pp. 3459–3461.
- Lv, Q., Ding, M., Liu, Q., Chen, Y., Feng, W., He, S., Zhou, C., Jiang, J., Dong, Y., Tang, J., 2021. Are we really making much progress?: Revisiting, benchmarking and refining heterogeneous graph neural networks, in: Zhu, F., Ooi, B.C., Miao, C. (Eds.), KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14-18, 2021, ACM. pp. 1150–1160. URL: https://doi.org/10.1145/3447548.3467350, doi:10.1145/3447548.3467350.
- MacAvaney, S., Yates, A., Feldman, S., Downey, D., Cohan, A., Goharian, N., 2021. Simplified Data Wrangling with ir_datasets, in: Diaz, F., Shah, C., Suel, T., Castells, P., Jones, R., Sakai, T. (Eds.), SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021, ACM. pp. 2429–2436. URL: https://doi.org/10.1145/3404835.34 63254, doi:10.1145/3404835.3463254.
- Macdonald, C., Tonellotto, N., 2020. Declarative experimentation in information retrieval using pyterrier, in: Balog, K., Setty, V., Lioma, C., Liu, Y., Zhang, M., Berberich, K. (Eds.), ICTIR '20: The 2020 ACM SIGIR International Conference on the Theory of Information Retrieval, Virtual Event, Norway, September 14-17, 2020, ACM. pp. 161–168. URL: https://dl.acm.org/doi/10.1145/3409256.3409829.
- Manmatha, R., Rath, T., Feng, F., 2001. Modelling Score Distributions for Combining the Outputs of Search Engines, in: Kraft et al. (2001). pp. 267–275. pp. 267–275.
- Marchesin, S., Purpura, A., Silvello, G., 2020. Focal Elements of Neural Information Retrieval Models. An Outlook through a Reproducibility Study. Information Processing & Management 57, 102109.
- National Academies of Sciences, Engineering, and Medicine, 2019. Reproducibility and Replicability in Science. The National Academies Press, Washington, USA.

Open Science Collaboration, 2015. Estimating the Reproducibility of Psychological Science. Science 349, 943–952.

- Parapar, J., Losada, D.E., Presedo Quindimil, M.A., Barreiro, A., 2020. Using Score Distributions to Compare Statistical Significance Tests for Information Retrieval Evaluation. Journal of the Association for Information Science and Technology (JASIST) 71, 98–113. URL: https://doi.org/10.1002/asi.24203, doi:10.1002/asi.24203.
- Parapar, J., Radlinski, F., 2021. Towards Unified Metrics for Accuracy and Diversity for Recommender Systems, in: Pampín, H.J.C., Larson, M.A., Willemsen, M.C., Konstan, J.A., McAuley, J.J., Garcia-Gathright, J., Huurnink, B., Oldridge, E. (Eds.), Proc. 15th ACM Conference on Recommender Systems (RecSys 2021), ACM. pp. 75–84.
- Pashler, H., Wagenmakers, E., 2012. Editors' Introduction to the Special Section on Replicability in Psychological Science: A Crisis of Confidence? Perspectives on Psychological Science 7, 528–530. URL: http://journals.sagepub.com/doi/10.1177/1745691612465253, doi:10.1177/1745691612465253.
- Plesser, H.E., 2017. Reproducibility vs. replicability: A brief history of a confused terminology. Frontiers Neuroinformatics 11, 76. URL: https://doi.org/10.3389/fninf.2017.00076, doi:10.3389/fninf.2017.00076.
- Potthast, M., Gollub, T., Wiegmann, M., Stein, B., 2019. TIRA Integrated Research Architecture, in: Ferro and Peters (2019).
- Robertson, S.E., 2001. On Score Distributions and Relevance, in: Kraft et al. (2001). pp. 40-51. pp. 40-51.
- Robertson, S.E., Kanoulas, E., 2012. On Per-topic Variance in IR Evaluation, in: Hersh et al. (2012). pp. 891–900. pp. 891–900.
- Robertson, S.E., Walker, S., Jones, S., Hancock-Beaulieu, M.M., Gatford, M., 1994. Okapi at TREC–3, in: Harman, D.K. (Ed.), The Third Text REtrieval Conference (TREC-3), National Institute of Standards and Technology (NIST), Special Publication 500-225, Washington, USA. pp. 109–126.
- Robertson, S.E., Zaragoza, U., 2009. The Probabilistic Relevance Framework: BM25 and Beyond. Foundations and Trends in Information Retrieval (FnTIR) 3, 333–389.
- Sakai, T., Ferro, N., Soboroff, I., Zeng, Z., Xiao, P., Maistro, M., 2019. Overview of the NTCIR-14 CENTRE Task, in: Ishita, E., Kando, N., Kato, M.P., Liu, Y. (Eds.), Proc. 14th NTCIR Conference on Evaluation of Information Access Technologies, National Institute of Informatics, Tokyo, Japan. pp. 494–509.
- Schaer, P., Breuer, T., Castro, L.J., Wolff, B., Schaible, J., Tavakolpoursaleh, N., 2021. Overview of lilas 2021 living labs for academic search (extended overview), in: Faggioli, G., Ferro, N., Joly, A., Maistro, M., Piroi, F. (Eds.), Proceedings of the Working Notes of CLEF 2021 -Conference and Labs of the Evaluation Forum, Bucharest, Romania, September 21st - to - 24th, 2021, CEUR-WS.org. pp. 1668–1699. URL: http://ceur-ws.org/Vol-2936/paper-143.pdf.
- Soboroff, I., Ferro, N., Maistro, M., Sakai, T., 2019. Overview of the TREC 2018 CENTRE Track, in: Voorhees, E.M., Ellis, A. (Eds.), The Twenty-Seventh Text REtrieval Conference Proceedings (TREC 2018), National Institute of Standards and Technology (NIST), Special Publication 500-331, Washington, USA.
- Strohman, T., Metzler, D., Turtle, H., Croft, W.B., 2005. Indri: A language model-based search engine for complex queries, in: Proceedings of the international conference on intelligent analysis, Citeseer. pp. 2–6.
- Swets, J.A., 1963. Information Retrieval Systems. Science 141, 245-250.
- Trotman, A., Clarke, C.L.A., Ounis, I., Culpepper, J.S., Cartright, M.A., Geva, S., 2012. Open Source Information Retrieval: a Report on the SIGIR 2012 Workshop. ACM SIGIR Forum 46, 95–101.
- Voorhees, E.M., 1998. Variations in relevance judgments and the measurement of retrieval effectiveness, in: Croft, W.B., Moffat, A., van Rijsbergen, C.J., Wilkinson, R., Zobel, J. (Eds.), Proc. 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 1998), ACM Press, New York, USA. pp. 315–323.
- Voorhees, E.M., 2001. Evaluation by Highly Relevant Documents, in: Kraft et al. (2001). pp. 74-82. pp. 74-82.
- Voorhees, E.M., Rajput, S., Soboroff, I., 2016. Promoting Repeatability Through Open Runs, in: Yilmaz, E., Clarke, C.L.A. (Eds.), Proc. 7th International Workshop on Evaluating Information Access (EVIA 2016), National Institute of Informatics, Tokyo, Japan. pp. 17–20.
- Webber, W., Moffat, A., Zobel, J., 2010. A Similarity Measure for Indefinite Rankings. ACM Transactions on Information Systems (TOIS) 4, 20:1–20:38.
- Yang, P., Fang, H., Lin, J., 2018. Anserini: Reproducible ranking baselines using lucene. ACM J. Data Inf. Qual. 10, 16:1–16:20. URL: https://doi.org/10.1145/3239571, doi:10.1145/3239571.
- Yang, W., Lu, K., Yang, P., Lin, J., 2019. Critically examining the "neural hype": Weak baselines and the additivity of effectiveness gains from neural ranking models, in: Piwowarski, B., Chevalier, M., Gaussier, É., Maarek, Y., Nie, J., Scholer, F. (Eds.), Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019, ACM. pp. 1129–1132. URL: https://doi.org/10.1145/3331184.3331340, doi:10.1145/3331184.3331340.