

Towards More Robust NLP System Evaluation: Handling Missing Scores in Benchmarks

Anonymous ACL submission

Abstract

The evaluation of natural language processing (NLP) systems is crucial for advancing the field, but current benchmarking approaches often assume that all systems have scores available for all tasks, which is not always practical. In reality, several factors such as the cost of running baseline, private systems, computational limitations, or incomplete data may prevent some systems from being evaluated on entire tasks. This paper formalize an existing problem in NLP research: benchmarking when some systems scores are missing on the task, and proposes a novel approach to address it. Our method utilizes a compatible partial ranking approach to impute missing data, which is then aggregated using the Borda count method. It includes two refinements designed specifically for scenarios where either task-level or instance-level scores are available. We also introduce an extended benchmark, which contains over 131 million scores, an order of magnitude larger than existing benchmarks. We validate our methods and demonstrate their effectiveness in addressing the challenge of missing system evaluation on an entire task. This work highlights the need for more comprehensive benchmarking approaches that can handle real-world scenarios where not all systems are evaluated on the entire task.

1 Introduction

Benchmarking and system evaluation are critical processes for assessing the performance of AI systems, providing a standardized means of comparing various models and techniques while keeping track of technological advancements (Ruder, 2021; Dehghani et al., 2021; Post, 2018). However, evaluating general-purpose systems, such as foundation models used for generative tasks (Lehman et al., 2023; Kocoń et al., 2023; OpenAI, 2023; Brown et al., 2020; Raffel et al., 2020), presents unique challenges. A single task, metric, or dataset may

not be sufficient to effectively gauge their capabilities (Herbrich et al., 2006; Novikova et al., 2018; Sedoc and Ungar, 2020). Therefore, it is crucial to develop tools that can benchmark these systems on a multitude of tasks (Aribandi et al., 2021), enabling a comprehensive assessment of their overall performance (Peyrard et al., 2021).

In recent years, the field of NLP has made significant strides, with frequent emergence of new models (Lehman et al., 2023; Kocoń et al., 2023; Brown et al., 2020; OpenAI, 2023; Raffel et al., 2020; Liu et al., 2019; Fan et al., 2021) and techniques (Bommasani et al., 2021; Hupkes et al., 2022). To evaluate the performance of these systems across various tasks, datasets, and metrics (Colombo et al., 2022c) have been created. However, with the increasing complexity of these benchmarks, missing scores has become a significant challenge. Missing data can arise from a variety of sources, such as benchmarks that are too large or time-consuming to run (*e.g.*, BigBench has recently introduced MiniBench for these reasons (Srivastava et al., 2022)), high costs associated with reproducing experiments (*e.g.*, see Table 3 in (Artetxe et al., 2022)), incomplete datasets (see Table 5 in (Reid and Artetxe, 2022)), data collection errors, data cleaning procedures, data privacy concerns (particularly in-house datasets (Guibon et al., 2021)), and specialized expertise required to process niche datasets (Peng et al., 2019). In recent work, two main approaches have been followed to deal with missing scores, which are discarding data (Pfeiffer et al., 2022) or ignoring certain tasks (see Table 10 in (Lin et al., 2022) and Table 5 in (Martin et al., 2020)) or evaluations. However, these approaches are unsatisfactory as they can lead to biased and unreliable evaluations.

In this work, we address benchmarking NLP systems *when one or several systems cannot be evaluated on a specific task*. We propose the development of effective methods for aggregating

metrics that can handle missing data and enable a comprehensive assessment of system performance. Our approach will ensure the reliability and validity of NLP system evaluations and contribute to the creation of benchmarks that can be used to compare and evaluate NLP systems effectively. Specifically, our contributions are listed below.

1. **Introducing a new problem with a direct impact on NLP research:** benchmarking when there are missing system evaluations for an entire task, which has practical implications (Pfeiffer et al., 2022; Lin et al., 2022; Martin et al., 2020; Guibon et al., 2021; Peng et al., 2019).

2. **A novel method for benchmarking NLP systems with missing system scores.** We present a novel method that effectively tackles the issue of missing system evaluations for entire tasks. Our work includes a novel combinatorial approach for imputing missing data in partial rankings. It allows using standard rank aggregation algorithms such as Borda and offers two refinements tailored to the availability of either task-level or instance-level scores of the systems across different tasks.

3. **An extended benchmark for a comprehensive and accurate evaluation of NLP systems:** previous works on score aggregation relied on a benchmark of 250K scores (Colombo et al., 2022a; Peyrard et al., 2021), and did not release the system’s input, output, and ground truth texts. In our work, we collected their scores and extended the benchmark by adding over 131M scores.

4. **Extensive validation of benchmarking methods:** Results show that our method effectively handles missing scores and is more robust than existing methods, affecting final conclusions.

2 Formulation & Related Work

2.1 General Considerations

Comparing systems with benchmarks. Benchmarking aims to determine the ranking of systems based on their scores to identify the best-performing systems. In this process, each system is evaluated on individual tests within a larger set and assigned a score according to a specific metric. Depending on the available information, two approaches are typically employed. When only **task-level** information is available (i.e., the system scores on each task), a **task-level aggregation** is utilized to obtain the final ranking. On the other hand, when **instance-level information** is available, i.e., the system scores on each instance of

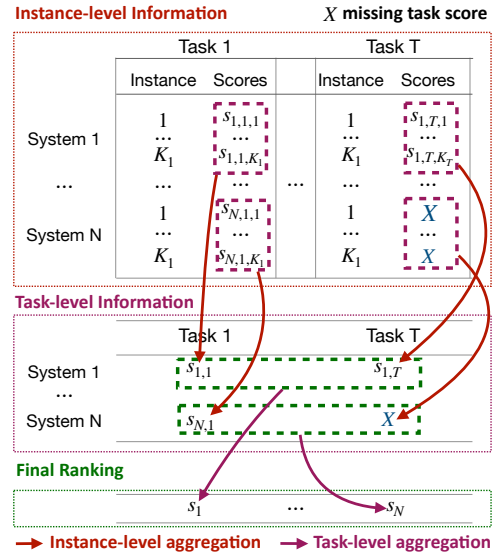


Figure 1: **Framework for benchmarking NLP systems with two information granularity:** *instance-level* (red above) and *task-level* (purple below). The final goal of benchmarking is to produce a ranking (green bottom). The instance-level aggregation allows for the derivation of task-level information, which is used to synthesize system performance via the final ranking (in green). **X** indicates the presence of missing values.

each task test set, an **instance-level aggregation** method is used to obtain the final system ranking. The mean aggregation has been adopted to aggregate at both the instance and task levels.

Benchmarking in the presence of missing data. As benchmarks and models continue to grow in size and complexity, the occurrence of missing system performance of entire tasks becomes increasingly common. This is particularly true in situations where one or more systems cannot be evaluated on a specific task due to factors such as the high cost of running the model or the extensive computational requirements of the benchmarks (Gehrmann et al., 2022a, 2021, 2022b). Fig. 1 illustrates the general framework (i.e., with instance and system level).

2.2 Problem Formulation

This discussion will use notation similar to that in the previously mentioned work (Colombo et al., 2022a). In essence, we are dealing with a scenario where N systems are being compared based on their performance on T different tasks. Each task $t \in \{1, \dots, T\}$ has a specific metric m_t associated with it and has been evaluated on k test instances with $k \in \{1, \dots, K_t\}$, where K_t is the test size of task t . The score of each system on each instance of each test set is represented by

the real number $s_{n,t,k} \in \mathbb{R}$. The final goal of benchmarking is to output a ranking of each system according to some objective criterion. We denote by \mathfrak{S}_N the symmetric group on N elements. With this objective in mind aggregating instance and task level information is equivalent to computing a permutation $\sigma \in \mathfrak{S}_N$ corresponding to the final ranking of the N systems. In this formalism, system i is the σ_i -th best system according to the considered aggregation. Equivalently, ordering $\pi = (\pi_1 \succ \pi_2 \succ \dots \succ \pi_N)$ denotes that π_i is better than system π_{i+1} for all i . Let us first define the different granularity of benchmarking depending on whether we have access to instance scores.

Aggregating with Missing Task Level Information. Given a set of scores $(s_{n,t}, 1 \leq n \leq N_t, 1 \leq t \leq T)$ where N_t is the number of systems for which we have access to the score on task t , find a proper aggregation procedure.

Thus the problem of task-level information aggregation boils down to finding f^T :

$$f^T : \underbrace{\mathcal{S}_{N_1} \times \dots \times \mathcal{S}_{N_T}}_{T \text{ times}} \longrightarrow \mathfrak{S}_N. \quad (1)$$

where $\mathcal{S}_{N_t} = (s_{n,t}, 1 \leq n \leq N_t)$ is the set of scores achieved by each system evaluated on the task t . Only N_t systems are evaluated on task t . In many cases, we not only have access to task-level performance but also individual instance-level scores. As a result, the challenge lies in effectively aggregating information at the instance level.

Aggregating Missing Instance Level Information. Given a set of scores $(s_{n,t,k}, 1 \leq n \leq N_t, 1 \leq t \leq T, 1 \leq k \leq K_t)$ where similarly as previously N_t is the number of systems for which we have access to the score on task t , find a proper aggregation procedure.

Thus the problem of instance-level information aggregation boils down to finding f^I :

$$f^I : \underbrace{\mathcal{S}_{N_1}^1 \times \dots \times \mathcal{S}_{N_1}^{K_1} \times \dots \times \mathcal{S}_{N_t}^1 \times \dots \times \mathcal{S}_{N_t}^{K_t} \times \dots \times \mathcal{S}_{N_T}^1 \times \dots \times \mathcal{S}_{N_T}^{K_T}}_{T \sum_t K_t \text{ times}} \longrightarrow \mathfrak{S}_N. \quad (2)$$

where $\mathcal{S}_{N_t}^k = (s_{n,t,k}, 1 \leq n \leq N_t)$ is the set of the score achieved by each system evaluated on the task t for the specific instance k .

Remark 1. In the context of complete ranking, which is also a classical setting for benchmarking NLP systems and has been addressed in (Colombo et al., 2022a), we have $N_t = N$ for all $t \in [1, T]$.

2.3 Handling Complete Scores in NLP System Evaluation

The literature relies on two main techniques for aggregating score information to benchmark ML systems: mean aggregation and ranking-based aggregation.

Mean aggregation (σ^μ) is the default choice for practitioners. At the task level σ^μ is defined as :

$$\sigma^\mu = \text{argsort} \left(\text{argsort} \left[\frac{1}{T} \sum_{1 \leq t \leq T} s_{n,t} \text{ for } 1 \leq n \leq N \right] \right)$$

and at the instance level :

$$\sigma^\mu = \text{argsort} \left(\text{argsort} \left[\frac{1}{T} \sum_{1 \leq t \leq T} \frac{1}{K_t} \sum_{1 \leq k \leq K_t} s_{n,t,k} \text{ for } 1 \leq n \leq N \right] \right)$$

where $\text{argsort}(\mathbf{u})$ is the permutation that sorts the items in \mathbf{u} . However, this approach has its limitations, particularly when evaluating tasks of different natures or using evaluation scores that are not on the same scale. Indeed in NLP, metrics can have different ranges (or even be unbounded) and systems are evaluated based on diverse criteria such as quality, speed, or number of parameters. In such cases, conventional rescaling or normalization techniques may not sufficiently capture the inherent difficulty of each task.

Ranking Based Aggregation To address the challenges mentioned, researchers have proposed ranking-based aggregations (Peyrard et al., 2021; Colombo et al., 2022a). These methods aggregate rankings instead of scores. In (Colombo et al., 2022a), the authors tackle the problem of generating a ranking by aggregating rankings, utilizing the Borda count method (see Ssec. E.1 for more details on Borda Count) known for its computational properties (Bartholdi et al., 1989; Dwork et al., 2001; Ali and Meilă, 2012). Extending the Borda count method is not a straightforward task either. Next, we present our aggregation procedure that can handle missing system scores on a whole task.

3 Ranking with missing system evaluation

In this section, we will outline our methodology for ranking multiple systems in multi-task benchmarks, even if some systems have not been evaluated on one or more tasks. We use the ranking and ordering notation interchangeably.

3.1 Partial Rankings

Mapping Scores to Partial Rankings To address the challenge of benchmarking with missing system evaluations, we propose a ranking-based approach that focuses on aggregating rankings rather

than directly combining scores. Suppose we have a specific task t with a task-level score denoted as S_{N_t} , or in the case of instance-level information, a task t and instance k with score $S_{N_t}^k$. In scenarios where there are missing evaluations at the task-level or instance-level, a *partial ranking* of systems is generated. A partial ordering represents an incomplete ranking that includes only a subset of items from a larger set. We denote the partial ordering of systems as $\pi^{N_t} = (\pi_1 \succ \pi_2 \succ \dots \succ \pi_{N_t})$ for the task-level scenario, and as $\pi^{N_t,k} = (\pi_1^k \succ \pi_2^k \succ \dots \succ \pi_{N_t}^k)$ for the instance-level scenario. Here, π_i represents the i -th best system according to the set S_{N_t} in the task-level scenario, while π_i^k represents the i -th best system according to π^k in the instance-level scenario.

Compatible Permutation When working with partial rankings, it is necessary to construct a complete ranking that respects the order of the evaluated systems, i.e., a linear extension of the partial ranking. This is accomplished by creating a compatible permutation (Gessel and Zhuang, 2018), which is a permutation of all systems consistent with the partial ranking. To construct a compatible permutation, we begin with the partial ranking of the evaluated systems and extend it to include the missing systems while maintaining the order of the evaluated systems. For example, let’s consider a partial ordering $\pi_1 \succ \pi_2$ based on the evaluation of only these two systems. If there is an additional system that has not been evaluated, we can construct three compatible permutations: $\pi_3 \succ \pi_1 \succ \pi_2$, $\pi_1 \succ \pi_3 \succ \pi_2$ and $\pi_1 \succ \pi_2 \succ \pi_3$. These permutations ensure that the ordering of the evaluated systems is preserved while incorporating the missing system into the complete ranking.

Why use a combinatorial approach? Imputing missing data using compatible permutations enables us to leverage the Borda aggregation, inheriting its theoretical and practical advantages. Unlike classical methods like harmonic Fourier analysis (Kondor and Barbosa, 2010; Kondor and Dempsey, 2012; Cl  men  on et al., 2011) or multi-resolution analysis (Sibony et al., 2015), our approach works, providing a distinct combinatorial solution for imputing missing data in partial rankings.

3.2 Our Ranking Procedures

Our method can be described in two steps:

Our ranking procedure in a nutshell

1. **Matrix Representation of the rankings (Sssec. 3.2.1).** To harness the full potential of the available information in partial rankings, we **efficiently** generate all compatible permutations from the given partial rankings.

2. **Final System Ranking from Matrix Representation.** To obtain the final ranking of the systems, we propose a one-level (σ^l) approach (see Sssec. 3.2.2) for both task-level and instance-level information and a two-level aggregation approach (σ^{2l}) for instance-level information (see Sssec. 3.2.3).

3.2.1 Matrix representation of the rankings

Intuition. Our algorithm first summarizes the available information in all tasks and imputes the missing information in a consistent manner. To do this, we use a matrix representation M^π for each partial ranking π . This matrix decomposes the ranking information in pairwise variables, i.e., for every pair of systems i, j there is a variable representing the probability that system i outperforms system j .

Why using matrix representation? Using pairwise information has many advantages in ranking problems with missing data since it allows decomposing the total ranking information in $N(N-1)/2$ different variables. This decomposition has been used in statistical problems on partial and complete rankings (F  rnkranz and H  llermeier, 2003; Lu and Boutilier, 2014a,b; Shah et al., 2017), for computing distances among partial rankings (Fagin et al., 2003), clustering (Ailon, 2010) and classification (H  llermeier et al., 2008) among others. However, these problems consider specific forms of missing data such as top- k rankings (Fagin et al., 2003) or bucket orderings (Achab et al., 2019). Our approach differs from the aforementioned literature in the fact that we impute the missing data in a consistent manner in order to be able to deal with arbitrary missing data.

Efficiently building M^π . Let us consider a partial ranking π and let $M^\pi \in [0, 1]^{N \times N}$ be its matrix representation. Matrix M_{ij}^π denotes the proportion of complete rankings that are compatible with π and satisfy the condition $i \succ j$, where i and j are distinct systems in the task. Formally, we can

distinguish three cases:

1. if system i is directly compared to system j in π . Set $M_{i,j}^\pi = 0$ if $i \succ j$ else $M_{i,j}^\pi = 1$.

2. if no information is provided for either system i or system j in π , meaning that both systems are unobserved in the partial ranking. In this case, $M_{i,j}^\pi = 0.5$, which is the natural choice when no information is available.

3. if we lack direct information about the comparison between system i and j in π (one system was evaluated and the other was not), we represent this situation by setting the corresponding matrix entry to the proportion of compatible permutations ranking system i higher than system j among the total number of compatible permutations (see Ap. E).

A naive algorithm for generating the matrix M^π from π would have factorial complexity and it is thus exorbitant in practice for a relatively small number of systems, say $N > 10$. **One of the contributions of our solution is to reduce the complexity to $O(n^3)$ by efficiently computing $M_{i,j}^\pi$.** The closed-form expressions for $M_{i,j}^\pi$ as well as the proof for uniformity can be found in Ap. E.

3.2.2 σ^l : A one level approach

Intuition. At this stage, we have demonstrated the construction of a matrix M^π for a given partial ranking. However, in benchmarking scenarios, systems are typically evaluated on multiple tasks (in the case of task-level evaluation) or on multiple instances and tasks (in the case of instance-level evaluation), requiring the combination of multiple partial rankings. In this section, we will describe our approach for performing the one-level aggregation to address this requirement.

Combining Multiple Partial Rankings for Benchmarking. To combine the different matrices into a single matrix M we sum over all the tasks (in the case of task-level information) or instances and tasks (in the case of instance-level information). Formally, this is achieved by performing the following operation to obtain the combined matrix $M^I = \sum_{t \in [1, T]} \sum_{k \in [1, K_t]} M^{\pi^{r_t, k}}$, where $M^{\pi^{r_t, k}}$ represents the partial ranking induced on task t and instance k . Similarly, for the task level we define $M^T = \sum_{t \in [1, T]} M^{\pi^{r_t}}$ where $M^{\pi^{r_t}}$ represents the partial ranking induced on task t .

Obtaining the final system ranking In the final step, our goal is to obtain the final system ranking

σ^l based on the matrix M^I or M^T . To achieve this, we use the Borda Count method, which involves computing the column-wise sum of the matrix and return the permutation that sorts the scores in increasing order. This step aligns with the approach proposed in (Colombo et al., 2022a). Formally:

$$\sigma^l = \text{argsort}(\text{argsort}(\sum_i M_{i,0}, \dots, \sum_i M_{i,N})). \quad (3)$$

Here, M represents the matrix M^T for task-level information, and M^I for the instance-level.

3.2.3 σ^{2l} : A two-level approach

Intuition. In the case of instance-level information, we also present a two-step procedure that draws inspiration from the widely adopted two-step mean aggregation approach.

Procedure. In the first step, we apply the task-level aggregation approach to generate individual rankings for each task t , resulting in T different permutations. In the second step, we aggregate these multiple rankings using the Borda aggregation method. Formally σ^{2l} can be computed as:

1. For each task t , compute $M^t = \sum_{k \in [1, K_t]} M^{\pi^{r_t, k}}$

2. For each task t , compute $\sigma^{2l, t} = \text{argsort}(\text{argsort}(\sum_i M_{i,0}^t, \dots, \sum_i M_{i,N}^t))$.

3. Compute the Borda count aggregation σ^{2l} of $[\sigma^{2l, 1}, \dots, \sigma^{2l, t}, \dots, \sigma^{2l, T}]$. (Colombo et al., 2022b)

3.3 Confidence Intervals for σ^l

To evaluate systems with missing data, it is crucial to measure the uncertainty of partial rankings. In the previous section, we discussed combining partial rankings into a complete ranking. In this section, we analyze the confidence of our data regarding pairwise comparisons of system performance.

Under any ranking model such as Mallows Model (Fligner and Verducci, 1986) or Plackett-Luce (Plackett, 1975), M_{ij}^π are random variables of known expected value. In the previous section, we computed the empirical value \widehat{M}_{ij}^π , approximating the true value M_{ij}^π . Here, we want to know how close these two quantities are. Formally, we are looking for a confidence interval of level δ , that is the value for c_{ij} around \widehat{M}_{ij}^π that contains M_{ij}^π with high probability, $P(|\widehat{M}_{ij}^\pi - M_{ij}^\pi| \geq c_{ij}) \leq 1 - \delta$. Noting that $0 \leq M_{ij}^\pi \leq 1$, we can use the Hoeffding inequality (Hoeffding, 1994) to compute the value of the confidence interval:

$$c_{ij} = \sqrt{\frac{-\log \delta}{2z_{ij}}}, \quad (4)$$

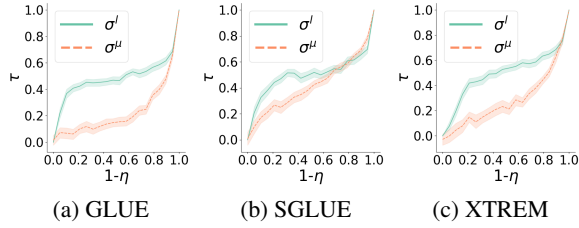


Figure 2: Task-Level Robustness Experiment. We compare the robustness of our method σ^l with the mean aggregation method σ^μ by measuring the Kendall τ correlation coefficient between their respective rankings after removing a proportion η of scores and by considering the whole scores.

where z_{ij} is the number of times the systems have been compared.

Intuition: to determine the significance of the difference in performance between system i and j , we can compare M_{ij} to 0.5. Thus, i performs better than j iff $M_{ij}^\pi > .5$. If the difference between M_{ij} and 0.5 is small, the performance difference between the two systems may not be statistically significant, indicating that we cannot determine which system performs better than the other.

The confidence interval developed above says that the true parameter M_{ij}^π is included in the interval $[\widehat{M}_{ij}^\pi - c_{ij}, \widehat{M}_{ij}^\pi + c_{ij}]$ with high probability. It follows that if 0.5 is not in this interval then we can say that one of the systems is better than the other with a high probability. Similar approaches have been proposed to find complete rankings and best-ranked systems with high probability (Busa-Fekete et al., 2014; Szörényi et al., 2015).

3.4 Baseline methods

To date, there is no established method for benchmarking NLP systems with missing data. To compare our proposed algorithm to existing methods, we consider a baseline approach that ignores missing data and relies on mean aggregation. This approach has been used in previous studies (Pfeiffer et al., 2022; Lin et al., 2022; Martin et al., 2020; Guibon et al., 2021; Peng et al., 2019), and we will refer to it as σ^μ in our experiments.

4 Synthetic Experiments

4.1 Data Generation

The analysis of a toy experiment involves synthetic scores generated from $N = 20$ systems, $T = 20$ tasks, and $K = 20$ instances. Each system’s performance is modeled by a Gumbel random variable

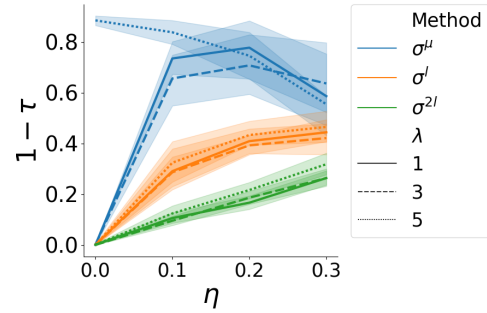


Figure 3: Synthetic experiment. Robustness for missing data η and different scaling corruptions λ .

G_n with a center at $\phi \times n$ and a scale of $\beta = 1$, where ϕ is a dispersion parameter between 0 and 1. The scores of each system, $s(n, t, k)$, are independent and identically distributed samples of G_n centered at $\phi \times n$ with a scale of $\beta = 1$. Furthermore, the scores from different systems are sampled independently. Since the difference between G_{n+1} and G_n follows a logistic distribution with a mean of ϕ and a scale of 1, the probability that system $n + 1$ performs better than system n is at least 0.5, i.e., $P(G_{n+1} - G_n > 0) \geq 0.5$. Thus, the ranking of systems for all k and t is a realization of the true ranking $[1, \dots, N]$, with a noise term controlled by the dispersion parameter ϕ . The extreme scenarios are $\phi = 0$ and $\phi = 1$, where $\phi = 0$ means that all scores $s(n, t, k)$ have the same distribution and $\phi = 1$ results in a strong consensus and a clear system ranking. Unless specifically mentioned, each experiment is repeated 100 times per data point.

4.2 Robustness To Scaling

In order to conduct a more detailed comparison of the ranking, we introduce a corruption in the scores of a specific task by rescaling them with a positive factor of λ . For this experiment, the corrupted tasks are randomly chosen. Although this corruption does not have any impact on our ranking process (since the ranking induced by a task-instance pair remains unchanged), it progressively disrupts the mean aggregation procedure as the value of λ increases (see Fig. 3 for detailed results). *This experiment further validates the use of rankings in NLP benchmarking, as these metrics involve different natures of measurements (e.g., BLEU score vs. number of parameters or speed) and can have bounded or unbounded scales.*

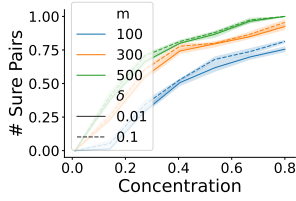


Figure 4: Confidence analysis.

4.3 Pairwise Confidence Analysis

To determine the number of system comparisons required to achieve a desired confidence level of δ , we use Eq. 4. Fig. 4 presents the results for two confidence levels (δ). The graph illustrates the number of system pairs for which 0.5 is not within the confidence interval, plotted against the number of comparisons for different values of m and ϕ . As expected, when the rankings are more concentrated (*i.e.*, when ϕ is closer to 1), fewer system comparisons are needed to achieve a high number of valid system comparisons. Real-world benchmarks usually have over 500 pairs.

5 Empirical Experiments

In this section, we benchmark our methods on real rankings. We introduce a dataset with over 100 million scores, surpassing previous datasets by several orders of magnitude (see Ssec. 5.1 and Ap. C).

5.1 A Comprehensive Collection of NLP System Scores

Our dataset builds upon the one used in (Colombo et al., 2022a) and includes two types of datasets: those with task-level information and those with instance-level information.

Datasets with Task Level Information Our datasets are based on GLUE (Wang et al., 2018), SGLUE (Wang et al., 2019), and XTREME (Hu et al., 2020), which include tasks of varying natures such as accuracy, F1-score, and mean square errors. In addition, we collected data from the GEM Benchmark (Gehrmann et al., 2021), which was an ideal use case for our methods as it encompasses missing data by design (as shown in Table 3 of (Gehrmann et al., 2021)) and includes evaluations of various natures such as lexical similarity, semantic equivalence, faithfulness evaluation, diversity, and system characterization (*i.e.*, size of the vocabulary).

Datasets with Instance Level Information We did not use the data from (Peyrard et al., 2021)

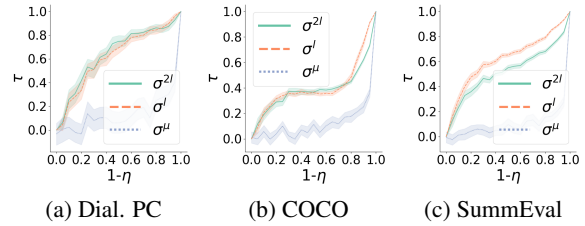


Figure 5: Instance-Level Robustness Experiment. We evaluate the robustness of our proposed aggregation methods, namely σ^{2l} , σ^l , and the mean aggregation method σ^μ , by randomly removing a proportion η of all instances on a specific task for a specific system.

for the datasets with instance-level information because they did not provide the sentence and reference test required to add more evaluation metrics or more systems.

Therefore, we collected all the data from scratch and extended the dataset in two ways. First, we collected data from five distinct tasks - dialogue (Mehri and Eskenazi, 2020), image description (Young et al., 2014), summary evaluation (Dang et al., 2008; Owczarzak and Dang, 2011; Bhandari et al., 2020; Fabbri et al., 2021), data-to-text (Gardent et al., 2017; Zhou and Lampouras, 2020), and translation (Ranasinghe et al., 2021). For translation, we added datasets from WMT15 to WMT21 in several languages such as en, ru, ts, and others.

Secondly, we expanded the set of used metrics from 10 to 17, including Bleu (Papineni et al., 2002), MoverScore (Zhao et al., 2019) and BERTScore (Zhang et al., 2019a).

Overall, our benchmark grew from 250K scores to over 131 M score. This extensive data work is one of the core contributions of this paper, and we believe it will be valuable for future research. A more detailed description of the datasets can be found in Ap. C.

5.2 Task-Level Benchmarking in Real-World Scenarios

In this section, we explore aggregating missing data with task-level information. First, we test the robustness of our proposed method (σ^l) against the mean aggregation method (σ^μ) and then we quantify the difference between the two output rankings. σ^l is more robust than σ^μ . To compare the effectiveness of aggregation methods in handling missing values, we randomly remove a proportion η of the task-level data and measure robustness by computing the Kendall τ between the rankings obtained

with and without missing values. From Fig. 2, we see that with no missing values (*i.e.*, $\eta = 0$), the aggregation methods yield the same rankings as the full data, with $\tau = 1$. Conversely, when all values are missing (*i.e.*, $\eta = 1$), there is no correlation. Overall, we find that σ^l achieves a higher correlation, with a large improvement of more than 10 points compared to other methods. These results demonstrate that, on average, the rankings remain more stable when using our proposed method.

σ^l outputs a different ranking than σ^μ . We evaluated the correlation between rankings in the robustness experiment shown in Fig. 2. We compared the rankings from σ^l and σ^μ by computing the averaged τ across varying proportions of missing data. Results in Tab. 1 show a weak correlation between the rankings, indicating they produce different outcomes. This is further supported by Tab. 2, which shows the percentage of times the top 1 and top 3 rankings differ in the 2k generated rankings. These results demonstrate that our ranking procedure not only is more robust but also yields different conclusions when benchmarking systems with missing tasks.

	$\tau_{\sigma^l \leftrightarrow \sigma^\mu}$	Dataset	top 1	top 3
GLUE	0.17 ± 0.24	GEM	0.52	0.25
SGLUE	0.33 ± 0.27	SGLUE	0.20	0.15
XTREM	0.26 ± 0.26	GLUE	0.10	0.07
GEM	0.36 ± 0.36	XTREM	0.19	0.09

Table 1: Agreement measured by Kendall τ correlation.

Table 2: Percentage of times the top 1 and top 3 systems are the same between σ^l and σ^μ .

5.3 Instance-Level Benchmarking in Real-World Scenarios

In this section, we evaluate the robustness of σ^{2l} , σ^l , and the baseline σ^μ .

σ^{2l} and σ^l are more robust than σ^μ . Similarly to the previous robustness experiment, we randomly remove a proportion η of scores by discarding all instances of a specific task. *The goal of this missing value sampling is to simulate how missing scores may occur when certain systems are not evaluated on specific tasks.* For each method, Fig. 5 reports the τ correlation coefficient between the ranking obtained with missing values and the ranking obtained with complete scores.

Both σ^{2l} and σ^l produce highly correlated rankings, while being different from σ^μ . We conducted a replication of the agreement analysis presented in Ssec. 5.2 and present the findings in

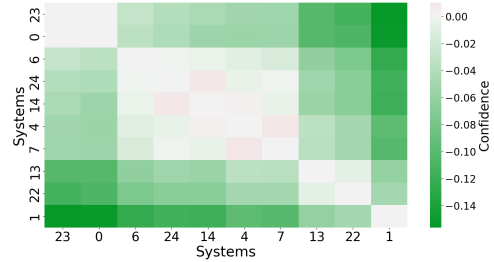


Figure 6: Confidence interval analysis on WMT en-de for a corruption level of $\eta = 0.2$ and a confidence level $\delta = 0.01$. The final ranking can be seen on the x-axis: left to right is best to worst

Tab. 3 and Tab. 4. Our results align with those of our previous experiments, demonstrating that both of our ranking-based procedures (σ^{2l} and σ^l) are more robust in the presence of missing data and yield different rankings than σ^μ .

	Corr.	σ^{2l} vs σ^l	Top 1	Top 3
$\tau_{\sigma^{2l} \leftrightarrow \sigma^l}$	0.80 ± 0.22		0.67	0.36
$\tau_{\sigma^l \leftrightarrow \sigma^\mu}$	0.20 ± 0.28	σ^l vs σ^μ	0.21	0.09
$\tau_{\sigma^\mu \leftrightarrow \sigma^{2l}}$	0.19 ± 0.28	σ^μ vs σ^{2l}	0.19	0.09

Table 3: Agreement. Table 4: Top 1 & 3 analysis.

5.4 Statistical Analysis

Confidence interval for practitioners. The confidence interval is valuable for informing additional comparisons between systems i and j . A narrow interval indicates a reliable comparison, while a wider interval suggests more uncertainty and the need for additional comparisons across tasks. For example, in Fig. 6, we report the results of applying σ^l on WMT en-de with a confidence level of $\delta = 0.1$. Green value in position $i < j$ illustrate that system $0.5 \notin [\widehat{M}_{ij}^\pi - c_{ij}, \widehat{M}_{ij}^\pi + c_{ij}]$ and $i \succ j$ with high probability. The scale of green displays the distance between 0.5 and the CI, so the greener the more $i \succ j$. The results reveal distinct blocks where top systems (*i.e.*, 9,1,16,15) significantly outperform others with high confidence. Near the diagonal, the elements indicate relatively closer performance of the systems.

6 Conclusions

Our study sheds light on the limitations of the conventional mean-aggregation when dealing with missing data. Our novel statistical perspective and aggregation procedures that are both robust and grounded in social choice theory. Overall, the one-level aggregation method (σ^l) stands out as the most robust approach.

7 Limitations

The initial limitation we pinpoint is the task’s reliance on the noise model applied to the data, which affects the outcomes. In an extreme scenario where a system lacks all measures, our method might not consistently rank it. Additional edge cases could be investigated, such as a system being poor in only one task with missing data, leading to potentially misleading ranking. To address this, we introduced the confidence interval in Section 3.3, supported by results in Section 5.4, to effectively recognize such challenging scenarios. It’s important to highlight that these edge cases can impact all ranking procedures involving missing data.

Another limitation pertains to our ranking procedure’s lack of consideration for user preferences regarding tasks. For instance, a user might emphasize certain tasks, such as A, D, and H, with task A carrying greater importance than the others. A natural approach to address this issue involves adopting a weighted variation of the Borda count or drawing inspiration from (Dwork et al., 2001). Although this avenue remains unexplored within our current work, it holds promise as a captivating direction for future investigations.

References

Mastane Achab, Anna Korba, and Stephan Cl  men  on. 2019. [Dimensionality reduction and \(bucket\) ranking: a mass transportation approach](#). volume 98, pages 64–93. PMLR.

Nir Ailon. 2010. [Aggregation of partial rankings, p-ratings and top-m lists](#). *Algorithmica (New York)*.

Khetam Al Sharou, Zhenhao Li, and Lucia Specia. 2021. Towards a better understanding of noise in natural language processing. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 53–62.

Alnur Ali and Marina Meil  . 2012. Experiments with kemeny ranking: What works when? *Mathematical Social Sciences*, 64(1):28–40.

Vamsi Aribandi, Yi Tay, Tal Schuster, Jinfeng Rao, Huaixiu Steven Zheng, Sanket Vaibhav Mehta, Honglei Zhuang, Vinh Q Tran, Dara Bahri, Jianmo Ni, et al. 2021. Ext5: Towards extreme multi-task scaling for transfer learning. *arXiv preprint arXiv:2111.10952*.

Mikel Artetxe, Itziar Aldabe, Rodrigo Agerri, Olatz Perez-de Vi  aspre, and Aitor Soroa. 2022. [Does corpus quality really matter for low-resource languages?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages

7383–7390, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2019. On the cross-lingual transferability of monolingual representations. *arXiv preprint arXiv:1910.11856*.

Mikel Artetxe and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the Association for Computational Linguistics*, 7:597–610.

Lo  c Barrault, Ondr  j Bojar, Marta R Costa-Jussa, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, et al. 2019. Findings of the 2019 conference on machine translation (wmt19). ACL.

John J Bartholdi, Craig A Tovey, and Michael A Trick. 1989. The computational difficulty of manipulating an election. *Social Choice and Welfare*, 6(3):227–241.

D Bayer and P Diaconis. 1992. [Trailing the dovetail shuffle to its lair](#). *The Annals of Applied Probability*, 2:294–313.

Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. 2009. The fifth pascal recognizing textual entailment challenge. In *TAC*.

Manik Bhandari, Pranav Gour, Atabak Ashfaq, Pengfei Liu, and Graham Neubig. 2020. Re-evaluating evaluation in text summarization. *arXiv preprint arXiv:2010.07100*.

Ondr  j Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, et al. 2017. Findings of the 2017 conference on machine translation (wmt17). Association for Computational Linguistics.

Ondr  j Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, et al. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198.

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754

755	Róbert Busa-Fekete, Eyke Hüllermeier, and Balázs Szörényi. 2014. Preference-Based Rank Elicitation using Statistical Models: The Case of Mallows . In <i>Proceedings of the 31th International Conference on Machine Learning, (ICML)</i> , pages 1071–1079.	In <i>proceedings of Sinn und Bedeutung</i> , volume 23, pages 107–124.	810
756			811
757			
758		Mostafa Dehghani, Yi Tay, Alexey A Gritsenko, Zhe Zhao, Neil Houlsby, Fernando Diaz, Donald Metzler, and Oriol Vinyals. 2021. The benchmark lottery. <i>arXiv preprint arXiv:2107.07002</i> .	812
759			813
760	Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. <i>arXiv preprint arXiv:1708.00055</i> .		814
761			815
762		William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In <i>Proceedings of the Third International Workshop on Paraphrasing (IWP2005)</i> .	816
763			817
764			818
765	Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. <i>arXiv preprint arXiv:1905.10044</i> .		819
766		Cynthia Dwork, Ravi Kumar, Moni Naor, and Dandapani Sivakumar. 2001. Rank aggregation methods for the web. In <i>Proceedings of the 10th international conference on World Wide Web</i> , pages 613–622.	820
767			821
768			822
769			823
770	Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. <i>Transactions of the Association for Computational Linguistics</i> , 8:454–470.	Alexander R Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. Summeval: Re-evaluating summarization evaluation. <i>Transactions of the Association for Computational Linguistics</i> , 9:391–409.	824
771			825
772			826
773			827
774			828
775		Ronald Fagin, Ravi Kumar, and D. Sivakumar. 2003. Comparing top k lists . <i>SIAM Journal on Discrete Mathematics</i> .	829
776	Stéphan Cléménçon, Romaric Gaudel, and Jérémie Jakubowicz. 2011. Clustering rankings in the fourier domain. In <i>Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2011, Athens, Greece, September 5-9, 2011. Proceedings, Part I 11</i> , pages 343–358. Springer.		830
777			831
778		Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021. Beyond english-centric multilingual machine translation. <i>The Journal of Machine Learning Research</i> , 22(1):4839–4886.	832
779			833
780			834
781			835
782	Pierre Colombo, Nathan Noiry, Ekhine Irurozki, and Stephan CLEMENCON. 2022a. What are the best systems? new perspectives on NLP benchmarking . In <i>Advances in Neural Information Processing Systems</i> .		836
783			837
784		Akhbardeh Farhad, Arkhangorodsky Arkady, Biesialska Magdalena, Bojar Ondřej, Chatterjee Rajen, Chaudhary Vishrav, Marta R Costa-jussa, España-Bonet Cristina, Fan Angela, Federmann Christian, et al. 2021. Findings of the 2021 conference on machine translation (wmt21). In <i>Proceedings of the Sixth Conference on Machine Translation</i> , pages 1–88. Association for Computational Linguistics.	838
785			839
786			840
787	Pierre Colombo, Nathan Noiry, Ekhine Irurozki, and Stéphan Cléménçon. 2022b. What are the best systems? new perspectives on NLP benchmarking . In <i>NeurIPS</i> .		841
788			842
789			843
790			844
791	Pierre Colombo, Maxime Peyrard, Nathan Noiry, Robert West, and Pablo Piantanida. 2022c. The glass ceiling of automatic evaluation in natural language generation. <i>arXiv preprint arXiv:2208.14585</i> .		845
792		Michael A Fligner and Joseph S Verducci. 1986. Distance based ranking models. <i>Journal of the Royal Statistical Society</i> , 48(3):359–369.	846
793			847
794			848
795	Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. Xnli: Evaluating cross-lingual sentence representations. <i>arXiv preprint arXiv:1809.05053</i> .	Johannes Fürnkranz and Eyke Hüllermeier. 2003. Pairwise preference learning and ranking. pages 145–156.	849
796			850
797			851
798		Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. The webnlg challenge: Generating text from rdf data. In <i>Proceedings of the 10th International Conference on Natural Language Generation</i> , pages 124–133.	852
799			853
800	Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In <i>Machine Learning Challenges Workshop</i> , pages 177–190. Springer.		854
801			855
802			856
803		Sebastian Gehrmann, Tosin Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi, Aremu Anuoluwapo, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna Clinciu, Dipanjan Das, Kaustubh D Dhole, et al. 2021. The gem benchmark: Natural language generation, its evaluation and metrics. <i>arXiv preprint arXiv:2102.01672</i> .	857
804	Hoa Trang Dang, Karolina Owczarzak, et al. 2008. Overview of the tac 2008 update summarization task. In <i>TAC</i> .		858
805			859
806			860
807	Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank: Investigating projection in naturally occurring discourse.		861
808			862
809			863

864	Sebastian Gehrmann, Abhik Bhattacharjee, Abinaya Mahendiran, Alex Wang, Alexandros Papangelis, Aman Madaan, Angelina McMillan-Major, Anna Shvets, Ashish Upadhyay, Bingsheng Yao, et al. 2022a. Gemv2: Multilingual nlg benchmarking in a single line of code. <i>arXiv preprint arXiv:2206.11249</i> .	Donald E Knuth. 1970. Permutations, matrices and generalized young tableaux. <i>Pacific Journal of Mathematics</i> , 34:709–727.	921
865			922
866			923
867			
868		Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniec, Marcin Gruz, Arkadiusz Janz, Kamil Kanclerz, Anna Kocoń, Bartłomiej Koptyra, Wiktor Mieleszczenko-Kowszewicz, Piotr Miłkowski, Marcin Oleksy, Maciej Piasecki, Łukasz Radliński, Konrad Wojtasik, Stanisław Woźniak, and Przemysław Kazienko. 2023. Chatgpt: Jack of all trades, master of none. <i>Preprint</i> , arXiv:2302.10724.	924
869			925
870	Sebastian Gehrmann, Elizabeth Clark, and Thibault Selam. 2022b. Repairing the cracked foundation: A survey of obstacles in evaluation practices for generated text. <i>arXiv preprint arXiv:2202.06935</i> .		926
871			927
872			928
873			929
874	Ira M Gessel and Yan Zhuang. 2018. Shuffle-compatible permutation statistics. <i>Advances in Mathematics</i> , 332:85–141.		930
875			931
876			932
877	Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and William B Dolan. 2007. The third pascal recognizing textual entailment challenge. In <i>Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing</i> , pages 1–9.	R Kondor and M Barbosa. 2010. Ranking with kernels in fourier space.	933
878			934
879			
880		Risi Kondor and Walter Dempsey. 2012. Multiresolution analysis on the symmetric group. <i>Advances in Neural Information Processing Systems</i> , 25.	935
881			936
882	Gaël Guibon, Matthieu Labeau, H�el�ene Flamein, Luce Lefeuvre, and Chlo�e Clavel. 2021. Few-shot emotion recognition in conversation with sequential prototypical networks. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6858–6870, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		937
883			
884			938
885			939
886			940
887			941
888			942
889			
890	Ralf Herbrich, Tom Minka, and Thore Graepel. 2006. Trueskill™: A bayesian skill rating system. In <i>Advances in Neural Information Processing Systems</i> , volume 19. MIT Press.	Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In <i>Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning</i> .	943
891			944
892			945
893			946
894	Wassily Hoeffding. 1994. Probability inequalities for sums of bounded random variables. <i>The collected works of Wassily Hoeffding</i> , pages 409–426.	Patrick Lewis, Barlas O�guz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2019. Mlqa: Evaluating cross-lingual extractive question answering. <i>arXiv preprint arXiv:1910.07475</i> .	947
895			948
896			949
897			950
898	Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In <i>International Conference on Machine Learning</i> , pages 4411–4421. PMLR.	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	951
899			952
900			953
901			
902		Chin-Yew Lin, Guihong Cao, Jianfeng Gao, and Jian-Yun Nie. 2006. An information-theoretic approach to automatic evaluation of summaries. In <i>Proceedings of the Human Language Technology Conference of the NAACL, Main Conference</i> , pages 463–470.	954
903	Diewke Hupkes, Mario Giulianelli, Verna Dankers, Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Christos Christodoulopoulos, Karim Lasri, Naomi Saphra, Arabella Sinclair, et al. 2022. State-of-the-art generalisation research in nlp: a taxonomy and review. <i>arXiv preprint arXiv:2210.03050</i> .		955
904			956
905			957
906			958
907		Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. <i>arXiv preprint arXiv:2109.07958</i> .	959
908			960
909	Eyke H�ullermeier, Johannes F�urnkranz, Weiwei Cheng, and Klaus Brinker. 2008. Label ranking by learning pairwise preferences. <i>Artificial Intelligence</i> , 172:1897–1916.		961
910			
911			962
912			963
913	Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 252–262.	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. 2022. Few-shot learning with multilingual generative language models. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9019–9052.	964
914			965
915			966
916			967
917			968
918			969
919			970
920			971
		Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	972
			973

974	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Noth-	1027
975	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	man, Kevin Knight, and Heng Ji. 2017. Cross-lingual	1028
976	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	name tagging and linking for 282 languages. In <i>Pro-</i>	1029
977	Roberta: A robustly optimized bert pretraining ap-	<i>ceedings of the 55th Annual Meeting of the Associa-</i>	1030
978	proach. <i>arXiv preprint arXiv:1907.11692</i> .	<i>tion for Computational Linguistics (Volume 1: Long</i>	1031
		<i>Papers)</i> , pages 1946–1958.	1032
979	Barrault Loïc, Biesialska Magdalena, Bojar Ondřej, Fe-	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	1033
980	dermann Christian, Graham Yvette, Grundkiewicz	Jing Zhu. 2002. Bleu: a method for automatic evalu-	1034
981	Roman, Haddow Barry, Huck Matthias, Joanis Eric,	ation of machine translation. In <i>Proceedings of the</i>	1035
982	Kocmi Tom, et al. 2020. Findings of the 2020 con-	<i>40th annual meeting of the Association for Computa-</i>	1036
983	ference on machine translation (wmt20). In <i>Proceed-</i>	<i>tional Linguistics</i> , pages 311–318.	1037
984	<i>ings of the Fifth Conference on Machine Translation</i> ,		
985	pages 1–55. Association for Computational Linguis-	Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. Trans-	1038
986	tics,.	fer learning in biomedical natural language process-	1039
		ing: An evaluation of BERT and ELMo on ten bench-	1040
987	Tyler Lu and Craig Boutilier. 2014a. Effective sam-	marking datasets . In <i>Proceedings of the 18th BioNLP</i>	1041
988	pling and learning for mallows models with pairwise-	<i>Workshop and Shared Task</i> , pages 58–65, Florence,	1042
989	preference data. <i>Journal of Machine Learning Re-</i>	Italy. Association for Computational Linguistics.	1043
990	<i>search</i> .		
991	Tyler Lu and Craig Boutilier. 2014b. Effective sam-	Maxime Peyrard, Teresa Botschen, and Iryna Gurevych.	1044
992	pling and learning for mallows models with pairwise-	2017. Learning to score system summaries for better	1045
993	preference data. <i>Journal of Machine Learning Re-</i>	content selection evaluation. In <i>Proceedings of the</i>	1046
994	<i>search</i> .	<i>Workshop on New Frontiers in Summarization</i> , pages	1047
		74–84.	1048
995	Louis Martin, Benjamin Muller, Pedro Javier Or-	Maxime Peyrard, Wei Zhao, Steffen Eger, and Robert	1049
996	tiz Suárez, Yoann Dupont, Laurent Romary, Éric	West. 2021. Better than average: Paired evaluation	1050
997	de la Clergerie, Djamel Seddah, and Benoît Sagot.	of nlp systems. <i>arXiv preprint arXiv:2110.10746</i> .	1051
998	2020. CamemBERT: a tasty French language model .		
999	In <i>Proceedings of the 58th Annual Meeting of the As-</i>	Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James	1052
1000	<i>sociation for Computational Linguistics</i> , pages 7203–	Cross, Sebastian Riedel, and Mikel Artetxe. 2022.	1053
1001	7219, Online. Association for Computational Lin-	Lifting the curse of multilinguality by pre-training	1054
1002	guistics.	modular transformers . In <i>Proceedings of the 2022</i>	1055
		<i>Conference of the North American Chapter of the</i>	1056
1003	Shikib Mehri and Maxine Eskenazi. 2020. Ustr: An	<i>Association for Computational Linguistics: Human</i>	1057
1004	unsupervised and reference free evaluation metric for	<i>Language Technologies</i> , pages 3479–3495, Seattle,	1058
1005	dialog generation. <i>arXiv preprint arXiv:2005.00456</i> .	United States. Association for Computational Lin-	1059
		guistics.	1060
1006	Jun-Ping Ng and Viktoria Abrecht. 2015. Better sum-	Mohammad Taher Pilehvar and Jose Camacho-Collados.	1061
1007	marization evaluation with word embeddings for rouge.	2018. Wic: the word-in-context dataset for evaluat-	1062
1008	<i>arXiv preprint arXiv:1508.06034</i> .	ing context-sensitive meaning representations. <i>arXiv</i>	1063
		<i>preprint arXiv:1808.09121</i> .	1064
1009	Joakim Nivre, Mitchell Abrams, Zeljko Agic, Lars	Robin L Plackett. 1975. The analysis of permutations.	1065
1010	Ahrenberg, Lene Antonsen, et al. 2018. Univers-	<i>Journal of the Royal Statistical Society: Series C</i>	1066
1011	al dependencies 2.2 (2018). <i>URL http://hdl. han-</i>	<i>(Applied Statistics)</i> , 24(2):193–202.	1067
1012	<i>dle.net/11234/1-1983xxx</i> . LINDAT/CLARIN digital		
1013	library at the Institute of Formal and Applied Linguis-	Maja Popović. 2017. chrF++: words helping character	1068
1014	tics, Charles University, Prague, <i>http://hdl. handle-</i>	n-grams. In <i>Proceedings of the second conference on</i>	1069
1015	<i>net/11234/1-1983xxx</i> .	<i>machine translation</i> , pages 612–618.	1070
1016	Jekaterina Novikova, Ondřej Dušek, and Verena	Matt Post. 2018. A call for clarity in reporting BLEU	1071
1017	Rieser. 2018. Rankme: Reliable human ratings	scores . In <i>Proceedings of the Third Conference on</i>	1072
1018	for natural language generation. <i>arXiv preprint</i>	<i>Machine Translation: Research Papers</i> , pages 186–	1073
1019	<i>arXiv:1803.05928</i> .	191, Belgium, Brussels. Association for Computa-	1074
1020	OpenAI. 2023. Gpt-4 technical report . <i>Preprint</i> ,	tional Linguistics.	1075
1021	<i>arXiv:2303.08774</i> .		
1022	Karolina Owczarzak and Hoa Trang Dang. 2011.	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	1076
1023	Overview of the tac 2011 summarization track:	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	1077
1024	Guided task and aesop task. In <i>Proceedings of the</i>	Wei Li, and Peter J Liu. 2020. Exploring the limits	1078
1025	<i>Text Analysis Conference (TAC 2011)</i> , Gaithersburg,	of transfer learning with a unified text-to-text trans-	1079
1026	<i>Maryland, USA, November</i> .	former. <i>The Journal of Machine Learning Research</i> ,	1080
		21(1):5485–5551.	1081

1082	Afshin Rahimi, Yuan Li, and Trevor Cohn. 2019. Mas-	semantic compositionality over a sentiment treebank.	1137
1083	sively multilingual transfer for ner. <i>arXiv preprint</i>	In <i>Proceedings of the 2013 conference on empirical</i>	1138
1084	<i>arXiv:1902.00193</i> .	<i>methods in natural language processing</i> , pages	1139
		1631–1642.	1140
1085	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,	1141
1086	Percy Liang. 2016. Squad: 100,000+ questions	Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,	1142
1087	for machine comprehension of text. <i>arXiv preprint</i>	Adam R Brown, Adam Santoro, Aditya Gupta,	1143
1088	<i>arXiv:1606.05250</i> .	Adrià Garriga-Alonso, et al. 2022. Beyond the	1144
		imitation game: Quantifying and extrapolating the	1145
1089	Tharindu Ranasinghe, Constantin Orasan, and Ruslan	capabilities of language models. <i>arXiv preprint</i>	1146
1090	Mitkov. 2021. An exploratory analysis of multilin-	<i>arXiv:2206.04615</i> .	1147
1091	gual word-level quality estimation with cross-lingual		
1092	transformers. <i>arXiv preprint arXiv:2106.00143</i> .		
		Lars St, Svante Wold, et al. 1989. Analysis of variance	1148
1093	Machel Reid and Mikel Artetxe. 2022. PARADISE:	(anova). <i>Chemometrics and intelligent laboratory</i>	1149
1094	Exploiting parallel data for multilingual sequence-	<i>systems</i> , 6(4):259–272.	1150
1095	to-sequence pretraining. In <i>Proceedings of the 2022</i>		
1096	<i>Conference of the North American Chapter of the</i>	R P Stanley. 1986. <i>Enumerative Combinatorics</i> .	1151
1097	<i>Association for Computational Linguistics: Human</i>	Wadsworth Publishing Company.	1152
1098	<i>Language Technologies</i> , pages 800–810, Seattle,		
1099	United States. Association for Computational Lin-	Miloš Stanojević, Amir Kamran, Philipp Koehn, and	1153
1100	guistics.	Ondřej Bojar. 2015. Results of the wmt15 metrics	1154
		shared task. In <i>Proceedings of the Tenth Workshop</i>	1155
1101	Ond rej Bojar, Christian Federmann, Mark Fishel,	<i>on Statistical Machine Translation</i> , pages 256–273.	1156
1102	Yvette Graham, Barry Haddow, Matthias Huck,		
1103	Philipp Koehn, and Christof Monz. 2018. Findings of	Balázs Szörényi, Róbert Busa-Fekete, Adil Paul, and	1157
1104	the 2018 conference on machine translation (wmt18).	Eyke Hüllermeier. 2015. Online rank elicitation for	1158
1105	In <i>Proceedings of the Third Conference on Machine</i>	plackett-luce: A dueling bandits approach. <i>Advances</i>	1159
1106	<i>Translation</i> , volume 2, pages 272–307.	<i>in Neural Information Processing Systems</i> , 28.	1160
1107	Mario Rodríguez-Cantelar, Chen Zhang, Chengguang	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-	1161
1108	Tang, Ke Shi, Sarik Ghazarian, João Sedoc, Luis Fer-	preet Singh, Julian Michael, Felix Hill, Omer Levy,	1162
1109	nando D’Haro, and Alexander Rudnicky. 2023.	and Samuel R Bowman. 2019. Superglue: A stickier	1163
1110	Overview of robust and multilingual automatic eval-	benchmark for general-purpose language understand-	1164
1111	uation metrics for open-domain dialogue systems at	ing systems. <i>arXiv preprint arXiv:1905.00537</i> .	1165
1112	dstc 11 track 4. <i>arXiv preprint arXiv:2306.12794</i> .		
		Alex Wang, Amanpreet Singh, Julian Michael, Felix	1166
1113	Melissa Roemmele, Cosmin Adrian Bejan, and An-	Hill, Omer Levy, and Samuel R Bowman. 2018.	1167
1114	drew S Gordon. 2011. Choice of plausible alter-	Glue: A multi-task benchmark and analysis platform	1168
1115	natives: An evaluation of commonsense causal rea-	for natural language understanding. <i>arXiv preprint</i>	1169
1116	soning. In <i>2011 AAAI Spring Symposium Series</i> .	<i>arXiv:1804.07461</i> .	1170
1117	Sebastian Ruder. 2021. Challenges and Opportuni-	Alex Warstadt, Amanpreet Singh, and Samuel R Bow-	1171
1118	ties in NLP Benchmarking. http://ruder.io/	man. 2019. Neural network acceptability judgments.	1172
1119	nlp-benchmarking .	<i>Transactions of the Association for Computational</i>	1173
		<i>Linguistics</i> , 7:625–641.	1174
1120	João Sedoc and Lyle Ungar. 2020. Item response theory	Herbert S. Wilf. 1999. <i>East Side, West Side . . . - an</i>	1175
1121	for efficient human evaluation of chatbots. In <i>Pro-</i>	<i>introduction to combinatorial families-with Maple</i>	1176
1122	<i>ceedings of the First Workshop on Evaluation and</i>	<i>programming</i> . arxiv.	1177
1123	<i>Comparison of NLP Systems</i> , pages 21–33, Online.		
1124	Association for Computational Linguistics.	Adina Williams, Nikita Nangia, and Samuel R Bow-	1178
		man. 2017a. A broad-coverage challenge corpus	1179
1125	Nihar B. Shah, Sivaraman Balakrishnan, Adityanand	for sentence understanding through inference. <i>arXiv</i>	1180
1126	Guntuboyina, and Martin J. Wainwright. 2017.	<i>preprint arXiv:1704.05426</i> .	1181
1127	Stochastically transitive models for pairwise com-		
1128	parisons: Statistical and computational issues. <i>IEEE</i>	Adina Williams, Nikita Nangia, and Samuel R Bow-	1182
1129	<i>Transactions on Information Theory</i> .	man. 2017b. A broad-coverage challenge corpus	1183
		for sentence understanding through inference. <i>arXiv</i>	1184
1130	Eric Sibony, Stéphan Cléménçon, and Jérémie Jakubow-	<i>preprint arXiv:1704.05426</i> .	1185
1131	icz. 2015. Mra-based statistical learning from in-		
1132	complete rankings. In <i>International Conference on</i>	Yinfei Yang, Yuan Zhang, Chris Tar, and Jason	1186
1133	<i>Machine Learning</i> , pages 1432–1441. PMLR.	Baldrige. 2019. Paws-x: A cross-lingual adversarial	1187
		dataset for paraphrase identification. <i>arXiv preprint</i>	1188
1134	Richard Socher, Alex Perelygin, Jean Wu, Jason	<i>arXiv:1908.11828</i> .	1189
1135	Chuang, Christopher D Manning, Andrew Y Ng, and		
1136	Christopher Potts. 2013. Recursive deep models for		

1190	Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. <i>Transactions of the Association for Computational Linguistics</i> , 2:67–78.
1191	
1192	
1193	
1194	
1195	Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. 2018. Record: Bridging the gap between human and machine commonsense reading comprehension. <i>arXiv preprint arXiv:1810.12885</i> .
1196	
1197	
1198	
1199	
1200	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019a. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> .
1201	
1202	
1203	
1204	Yuan Zhang, Jason Baldridge, and Luheng He. 2019b. Paws: Paraphrase adversaries from word scrambling. <i>arXiv preprint arXiv:1904.01130</i> .
1205	
1206	
1207	Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M Meyer, and Steffen Eger. 2019. Moverscore: Text generation evaluating with contextualized embeddings and earth mover distance. <i>arXiv preprint arXiv:1909.02622</i> .
1208	
1209	
1210	
1211	
1212	Giulio Zhou and Gerasimos Lampouras. 2020. Webnlg challenge 2020: Language agnostic delexicalisation for multilingual rdf-to-text generation. In <i>Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)</i> , pages 186–191.
1213	
1214	
1215	
1216	
1217	
1218	Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp. 2017. Overview of the second bucc shared task: Spotting parallel sentences in comparable corpora. In <i>Proceedings of the 10th Workshop on Building and Using Comparable Corpora</i> , pages 60–67.
1219	
1220	
1221	
1222	
1223	Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp. 2018. Overview of the third bucc shared task: Spotting parallel sentences in comparable corpora. In <i>Proceedings of 11th Workshop on Building and Using Comparable Corpora</i> , pages 39–42.
1224	
1225	
1226	
1227	

Appendices

1228

A	Extended Related Work and Other Baselines methods	15	1229
A.1	More on missing data in NLP . . .	15	1230
A.2	Why not directly imputing data in the score?	15	1231
A.3	Ablation experiments	15	1232
A.4	Other ideas of future work	15	1233
A.5	More on the technical contribution of the algorithm.	16	1234
			1235
			1236
			1237
B	Ethical Statement & Limitation of our work	17	1238
			1239
C	Dataset Description	17	1240
C.1	Task Level Information	17	1241
C.2	Instance Level Information	17	1242
C.3	Data Statistics	18	1243
D	Additional Real-Data Experiments	18	1244
D.1	Example of Ranking with missing data on XTREM	18	1245
D.2	Additional Robustness Experiment on task level datasets	18	1246
D.3	Additional Robustness Experiment on instance level datasets	18	1247
D.4	Additional Confidence Analysis on Task Level	19	1248
			1249
			1250
			1251
			1252
E	On the Rankings	21	1253
E.1	Borda Count on permutations (in vector notation)	21	1254
E.2	Borda Count on permutations in pairwise matrix notation	21	1255
E.3	Generating all compatible rankings	22	1256
E.4	Proof of uniformity	22	1257
			1258
			1259

1260
1261

A Extended Related Work and Other Baselines methods

1262

A.1 More on missing data in NLP.

1263
1264
1265
1266
1267
1268
1269
1270

Another approach to handling missing data in benchmarks would be to create new datasets, however, it can be a sluggish, costly, and expertise-demanding process (see footnote 5 in (Lin et al., 2021)). Moreover, there are situations where collection becomes infeasible, such as when working with private datasets, calling for the need to develop tools that can rank systems with missing scores.

1271
1272

A.2 Why not directly imputing data in the score?

1273
1274
1275
1276
1277
1278
1279

Directly imputing values as scores is not the current practice in NLP. In fact, this approach would be inadequate due to potential variations in metric scale and difficulty, leading to a failure in accurately capturing task difficulty (as mentioned above and in (Dolan and Brockett, 2005)). To illustrate this to we present some experiments.

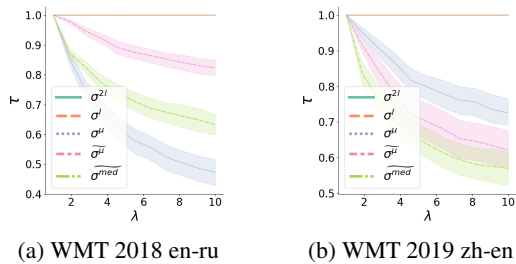


Figure 7: **Imputation methods are not robust to scaling.** To further compare the ranking, we corrupt the scores of a given task by re-scaling them by a factor λ . Whereas it does not affect our ranking procedure (every ranking induced by a task-instance pair remains the same), it increasingly perturbs the mean aggregation and other imputation procedures as λ increases. $\tilde{\sigma}^{med}$ corresponds to the median imputation and $\tilde{\sigma}^{\mu}$ corresponds to the mean imputation.

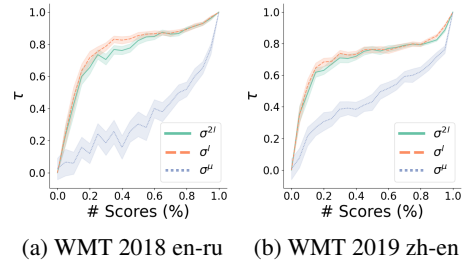


Figure 8: **Additional Instance-Level Robustness Experiment (see Figure 17).** We evaluate the robustness of our proposed aggregation methods, namely σ^{2l} , σ^l , and the mean aggregation method σ^{μ} , by randomly removing a proportion η of all instances on a specific task for a specific system.

A.3 Ablation experiments

1280

An interesting experiment is checking how the performance of the different methods vary when changing the number of systems, tasks or utterances. In the following experiments, we take the WebNLG2020 English dataset. For each ablation experiment—whether it’s concerning systems, tasks, or metrics—and for each η value within the set 0.7, 0.8, 0.9, we randomly select the systems, tasks, or utterances to retain. Then, in the same way as the robustness experiment, we compute correlations and iterate this process 20 times. Results are presented in Figures 9, 12 and 11. As we can see, the correlations are consistent when changing the number of these parameters.

1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294

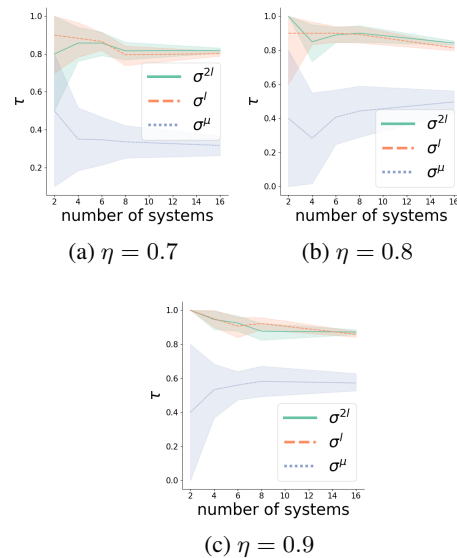


Figure 9: Ablation experiment on the number of systems

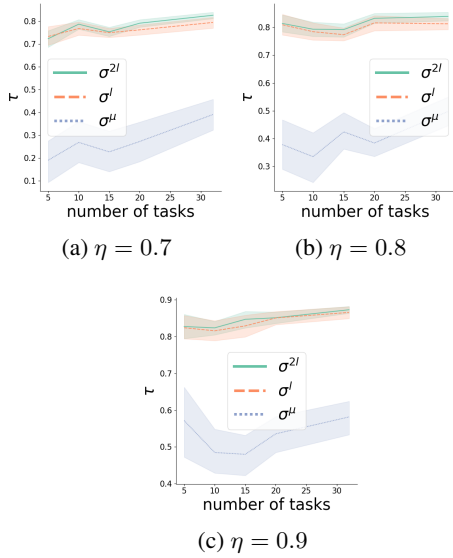


Figure 10: Ablation experiment on the number of tasks

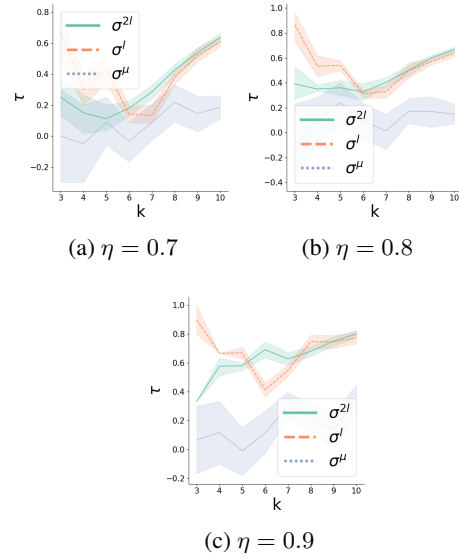


Figure 12: Ablation experiment on the top k systems

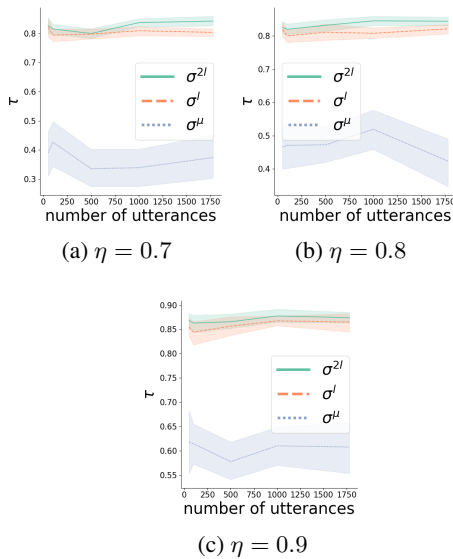


Figure 11: Ablation experiment on the number of utterances

A.4 Other ideas of future work

In the future, we would like to explore several refinements of the method:

- **Impact of the inter-task correlation.** The task correlation can impact the choice of the best system. In the future, we would like to study the impact of the choice of the ranking procedure in depth.
- **Impact of misleading evaluation.** Evaluation in NLP can be noisy due to the variety in language and lack of metric robustness (Al Sharou et al., 2021; Rodríguez-Cantelar et al., 2023). Future work will include the consideration of this factor when choosing the aggregation method.
- **Comparison with the ANOVA method (St et al., 1989).** Although this is slightly outside the scope of the paper, we would like to compare our confidence interval with the one obtained with the ANOVA method.

A.5 More on the technical contribution of the algorithm.

Our technical contribution boils down to extending the Borda aggregation to the case of missing data (aka incomplete rankings). In the ranking literature, two types of approaches can be identified to deal with partial rankings: Relying on top-k rankings. In this case, all the systems are evaluated but only those that are ranked in the first k positions are provided. There are many methods to aggregate

Another similar experiment we investigate here is when only keeping the top k systems (determined for each aggregation method using the complete dataset). Here are the results :

and all in this setting, for example, (Ailon, 2010). This is different from our scenario where some systems cannot be evaluated on particular tasks. Relying on incomplete rankings. In this case, only k systems are evaluated on a specific task. This fits our scenario. Rank aggregation/statistical analysis in the case of $k=2$ is called pairwise ranking and is well handled by the literature (Knuth, 1970; Lu and Boutilier, 2014a; Plackett, 1975; Popović, 2017; Zhang et al., 2018). These approaches are limited and only use pairwise comparisons which can lead to paradoxes when ranking more systems. In this paper, we introduce an aggregation procedure for arbitrary values of k . Our main technical contribution is to extend the Borda aggregation to incomplete rankings. To the best of our knowledge, this is the only paper dealing with aggregation -not specifically Borda- of incomplete rankings.

B Ethical Statement & Limitation of our work

It is important to consider the potential ethical implications and limitations of our work. One ethical concern is the potential bias in the reranking process, as the selection of the "best" hypothesis may favor certain perspectives or reinforce existing biases present in the training data. Care should be taken to ensure fairness and mitigate any potential bias before applying our methods.

C Dataset Description

C.1 Task Level Information

We provide additional details on the data collection for Task Level Information.

We gathered data from four benchmark studies, namely GLUE (General Language Understanding Evaluation) (Wang et al., 2018), SGLUE (SuperGLUE) (Wang et al., 2019)¹, XTREME (Hu et al., 2020) and GEM. In the GLUE dataset, there were a total of 105 systems evaluated across nine different tasks: CoLA, SST-2, MRPC, STS-B, QQP, MNLI, QNLI, RTE, and WNLI (Warstadt et al., 2019; Socher et al., 2013; Dolan and Brockett, 2005; Cer et al., 2017; Rajpurkar et al., 2016; Williams et al., 2017a; Dagan et al., 2005; Giampiccolo et al., 2007; Bentivogli et al., 2009; Levesque et al., 2012). The SGLUE dataset consisted of 24 systems evaluated on 10 different tasks: BoolQ, CB, COPA, MultiRC, ReCoRD, RTE, WiC, WSC, AX-b, and AX-g

¹Results can be accessed at <https://super.gluebenchmark.com/>

(Clark et al., 2019; De Marneffe et al., 2019; Roemle et al., 2011; Khashabi et al., 2018; Zhang et al., 2018; Levesque et al., 2012; Pilehvar and Camacho-Collados, 2018). The XTREME benchmark comprised 15 systems and included tasks such as sentence classification (XNLI and PAXS-X), structured prediction (Universal Dependencies v2.5 and Wikiann), sentence retrieval (BUCC and Tatoeba), and question answering (XQuAD, MLQA, TyDiQA-GoldP) (Conneau et al., 2018; Williams et al., 2017b; Yang et al., 2019; Zhang et al., 2019b; Nivre et al., 2018; Rahimi et al., 2019; Pan et al., 2017; Zweigenbaum et al., 2018, 2017; Artetxe and Schwenk, 2019; Artetxe et al., 2019; Rajpurkar et al., 2016; Lewis et al., 2019; Clark et al., 2020).

Each benchmark employed a variety of metrics with different scales, including accuracy, f1, and correlation. Additionally, the GEM benchmark involved 22 systems evaluated using diverse metrics such as prediction length, vocabulary size, entropy, Rouge, NIST, Bleu', Meteor', Bleurt, Nubia, and Bertscore.

C.2 Instance Level Information

In this particular setting, our primary focus is on evaluating the performance of natural language generation (NLG) systems, as these scores are among the easiest to collect. We concentrate on five different tasks: summary evaluation, image description, dialogue, and translation. For *summary evaluation*, we utilize the TAC08 (Dang et al., 2008), TAC10, TAC11 (Owczarzak and Dang, 2011), RSUM (Bhandari et al., 2020), and SEVAL (Fabbri et al., 2021) datasets. Regarding *sentence-based image description*, we rely on the FLICKR dataset (Young et al., 2014). For *dialogue*, we make use of the PersonaChat (PC) and TopicalChat (TC) datasets (Mehri and Eskenazi, 2020). For the translation part, we added datasets from WMT15 (Stanojević et al., 2015), WMT16 (Bojar et al., 2016), WMT17 (Bojar et al., 2017), WMT18 (rej Bojar et al., 2018), WMT19 (Barrault et al., 2019), WMT20 (Loïc et al., 2020), and WMT21 (Farhad et al., 2021) in several languages such as en, ru, ts, and others. For all datasets except MLQE, we consider automatic metrics based on S3 (both variant pyr/resp) (Peyrard et al., 2017), ROUGE (Lin, 2004) (including five of its variants (Ng and Abrecht, 2015)), JS [1-2] (Lin et al., 2006), Chrfpp (Popović, 2017), BLEU, BERTScore (Zhang et al., 2019a), and MoverScore (Zhao et al., 2019). For

the MLQE dataset, we solely consider several versions of BERTScore, MoverScore, and ContrastScore. Additionally, we incorporate human evaluation, which is specific to each dataset.

C.3 Data Statistics

To give to the reader a better sense of the richness of our benchmark, we report in Fig. 13 the statistics on our dataset. We demonstrate a diverse distribution of system counts across various datasets, ranging from a minimum of 2 systems to a maximum of 60 systems. Regarding the total number of sentences (instances) and the average number per system, as depicted in Fig. 14 and Fig. 15, the smaller datasets consist of several hundred sentences in total, while the larger datasets encompass up to several hundred thousand sentences in total.

D Additional Real-Data Experiments

In this dedicated section, we aim to provide curious readers with a deeper understanding of the capabilities of our methods by presenting additional figures and experimental results. Through these supplementary materials, we intend to shed more light on the effectiveness and potential of our approaches, enabling readers to gain valuable insights into our methods.

D.1 Example of Ranking with missing data on XTREM

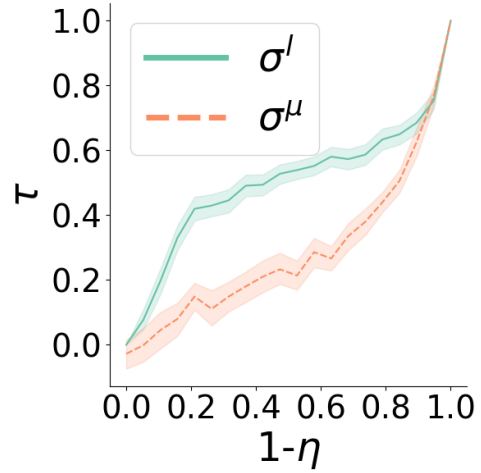
In this section, we aim to illustrate the distinction between different rankings obtained using σ^l and σ^μ on XTREM dataset for a specific noise realization. Using Tab. 5, we obtain the following rankings:

- σ^l gives the following ranking : $M0 > M3 > M2 > M1 > M7 > M5 > M4 > M8 > M9$
- σ^μ gives the following ranking : $M7 > M4 > M0 > M6 > M9 > M2 = M3 > M1 > M8 > M5$.

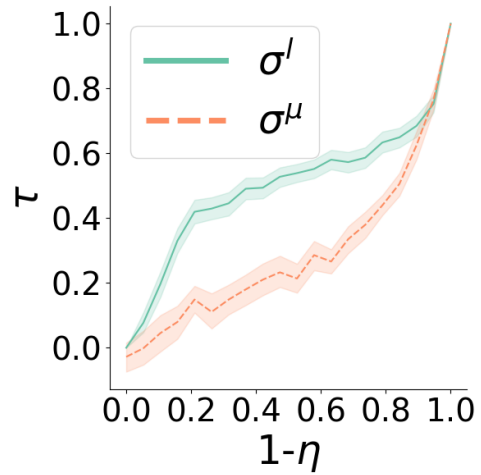
We can see that the two methods disagree on the best systems in this case. However, as can be seen in our experiments, the ranking-based method is more robust.

D.2 Additional Robustness Experiment on task level datasets

In this section, we report additional experiments on the task level robustness.



(a) XTREM



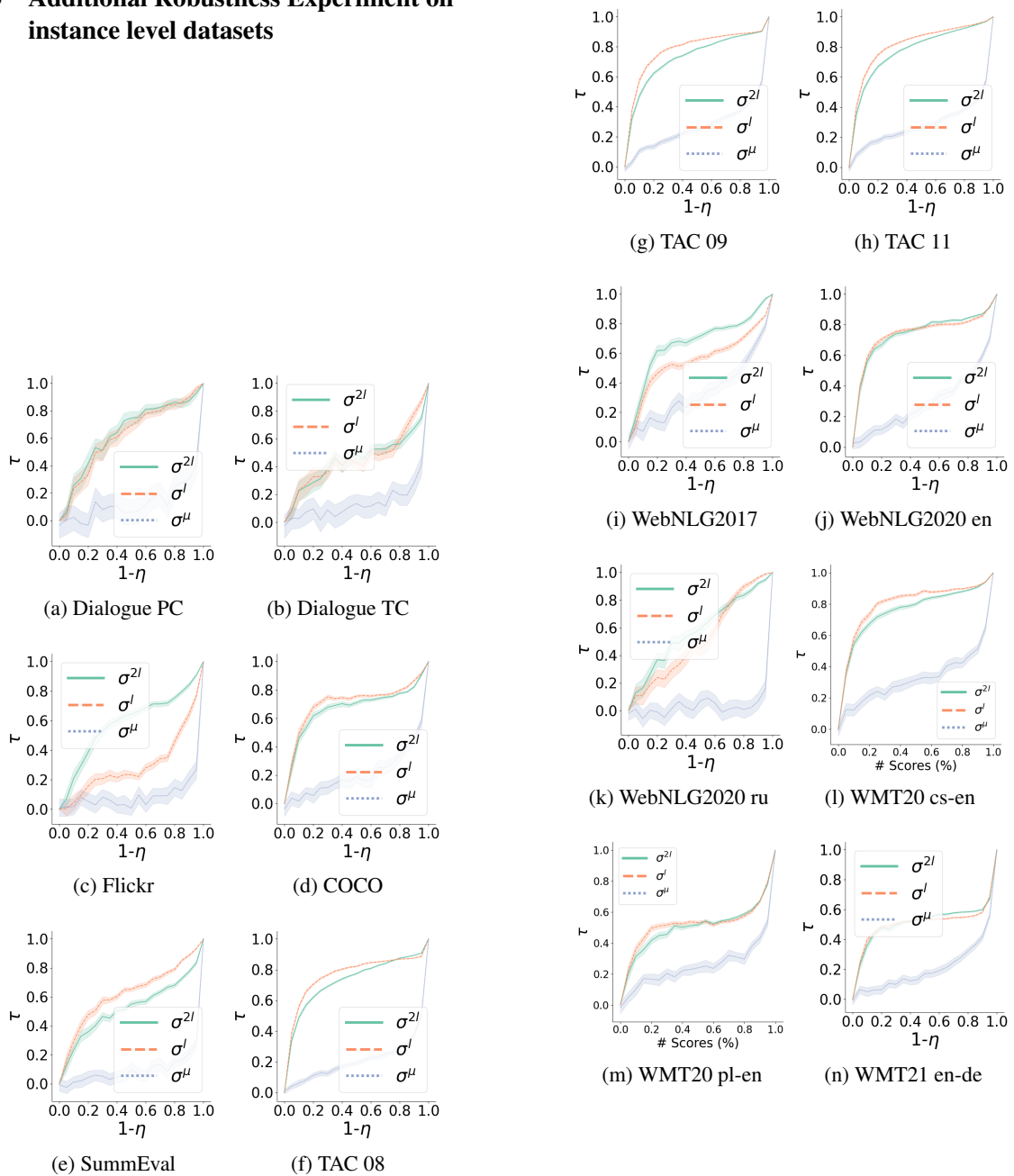
(b) GEM

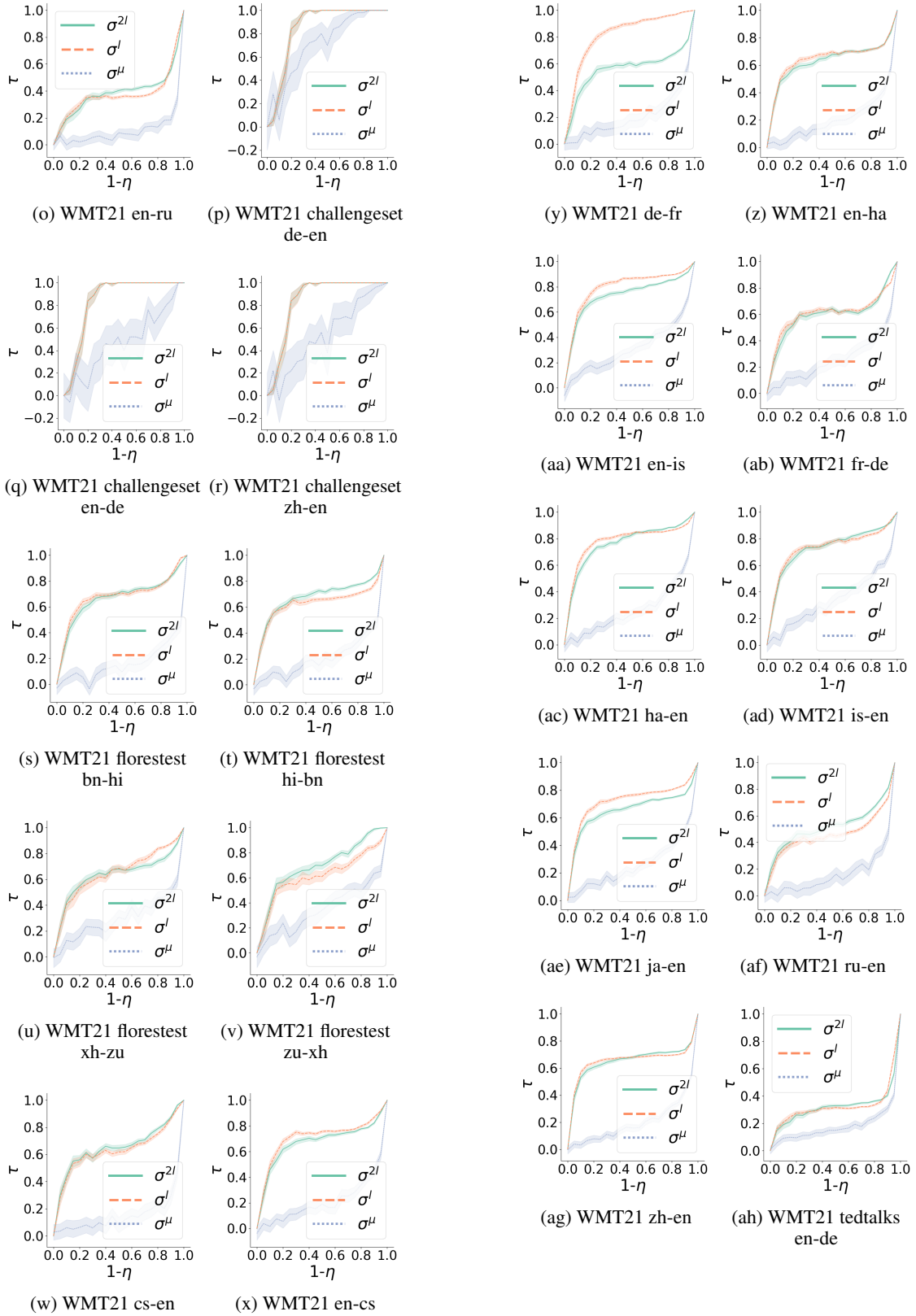
Figure 16: Robustness experiments on XTREM and GEM

Model	Classification	Structured Prediction	Question Answering	Sentence Retrieval
M0	90.3	X	76.3	93.7
M1	90.1	X	75.0	X
M2	89.3	75.5	75.2	92.4
M3	89.0	76.7	73.4	93.3
M4	88.3	X	X	X
M5	X	X	X	X
M6	87.9	75.6	X	91.9
M7	X	X	X	92.6
M8	X	75.4	X	X
M9	88.2	74.6	X	89.0

Table 5: XTREM dataset with 10 systems and 18 missing values ($\eta = 0.45$)

D.3 Additional Robustness Experiment on instance level datasets





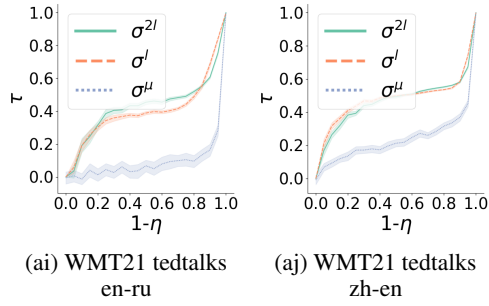


Figure 17: Instance-Level Robustness Experiment. We evaluate the robustness of our proposed aggregation methods, namely σ^{2l} , σ^l , and the mean aggregation method σ^μ , by randomly removing a proportion η of all instances on a specific task for a specific system. Each experiment is repeated 100 times for each proportion.

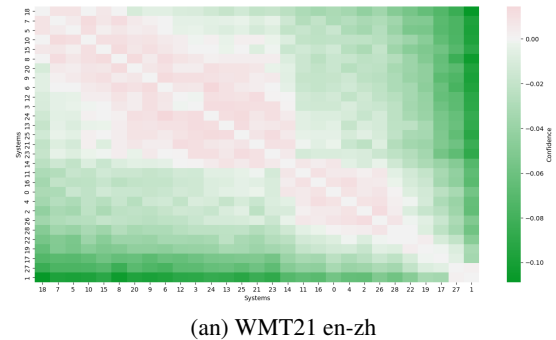
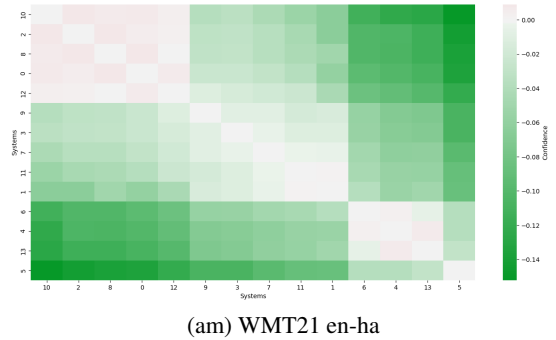
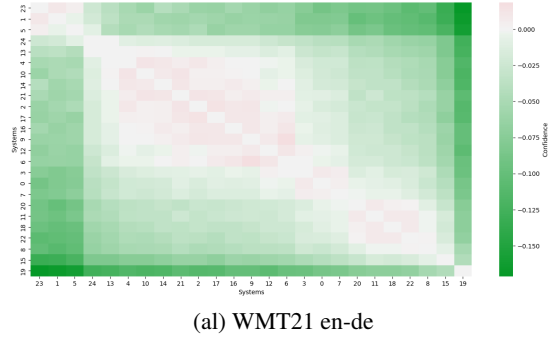
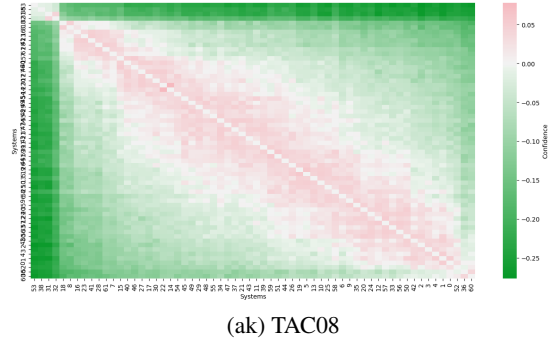


Figure 18: Confidence intervals for various instance level datasets with $\eta = 0.2$ and $\delta = 0.01$

D.4 Additional Confidence Analysis on Task Level

In this section, we present additional experiments conducted on four instance-level datasets. We computed confidence intervals for the instance-level, similar to the approach used in Section ???. Consistent with the main findings in the paper, our observations reveal that closer performance among systems is indicated near the diagonal and we can clearly observe group of systems. This analysis of confidence intervals provides valuable insights into the relative performance of different systems.

E On the Rankings

This section gathers technical considerations on the ranking methods used in our algorithm.

E.1 Borda Count on permutations (in vector notation)

Remark 2. The Borda count is a ranking system that aggregates a set of permutations $\sigma^1, \dots, \sigma^L \in$

1494 \mathfrak{S}_N by summing the ranks of each system and then
 1495 ranking the obtained sums. The procedure is as
 1496 follows:

1497 1. Compute $\text{sum}_n := \sum_{l=1}^L \sigma_n^l$ for every $1 \leq n \leq$
 1498 N ,

1499 2. Output $\sigma := \text{Borda}(\sigma^1, \dots, \sigma^L) \in$
 1500 \mathfrak{S}_N that ranks the sums, sum_n
 1501 $(\text{argsort}(\text{argsort}(\text{sum}_1, \dots, \text{sum}_T)))$.

1502 E.2 Borda Count on permutations in pairwise 1503 matrix notation

1504 In Sssec. 3.2.1 we argue that a ranking $\sigma \in \mathfrak{S}_N$
 1505 can also be written as a pairwise matrix and in
 1506 Sssec. 3.2.2 and Sssec. 3.2.3 we further elaborate
 1507 on how to write ranking data-set D in pairwise
 1508 matrix form $M^D \in [0, 1]^{N \times N}$. Under this notation,
 1509 the final aggregated ranking σ for the Borda count
 1510 algorithm can be shown to be equivalent to the
 1511 permutation that sorts the sum of the columns in
 1512 M^D ,

$$1513 \quad \sigma = \text{argsort} \left(\text{argsort} \left[\sum_i M_{i,0}^D, \dots, \sum_i M_{i,N}^D \right] \right). \quad (5)$$

1514 E.3 Generating all compatible rankings

1515 In this section, we detail the computation of the
 1516 $M_{i,j}^\pi$ when item i is not evaluated and item j is
 1517 evaluated. Let us fix some notation first. For the
 1518 following, k is the number of observed systems in
 1519 π , item i is not evaluated, item j is evaluated and
 1520 r is the (partial) rank of item j . Under this setting,
 1521 we set $M_{i,j}^\pi = p(n, k, r)$, i.e., the proportion of
 1522 compatible rankings that rank i before j when π
 1523 has k items. The closed-form expressions for these
 1524 quantities are given in Eq. 6. Here we note that
 1525 $t(n, k)$ is the total number of rankings of n items
 1526 compatible with π , S_b^a is the number of shuffles
 1527 of two lists of lengths a and b and V_b^a denotes the
 1528 variations of a out of b items, i.e., the number of
 1529 possible arrangements of selections of a objects
 1530 out of b , where the order of the selected objects
 1531 matters.

$$\begin{aligned}
 p(n, k, r) &= \sum_{i=0}^{n-k-1} V_{n-k-1}^i * (i+1) \\
 &\quad * S_{i+1}^r (n-k-i-1)! \\
 &\quad * S_{k-r-1}^{n-k-i-1} / t(n, k) \quad (6) \\
 t(n, k) &= (n-k)! * S_{n-k}^k \\
 S_b^a &= (a+b)! / (a! + b!) \\
 V_b^a &= a! / (b-a)!
 \end{aligned}$$

1532 **Remark 3.** A naive algorithm for generating the
 1533 matrix M^π from $\sigma \in S_{N-r_{tk}}$ would have factorial
 1534 complexity and it is thus exorbitant in practice for
 1535 a relatively small number of systems, say $N > 10$.
 1536 However, our solution has a complexity of $O(n^3)$
 1537 and can be precomputed once at the beginning of
 1538 the benchmarking process to efficiently generate
 1539 the pairwise matrix M^π from partial ranking π .
 1540

1541 E.4 Proof of uniformity

1542 In this section, we give the intuition and the proof
 1543 for Eq. 6. This section follows a classic strategy on
 1544 Enumerative Combinatorics (Stanley, 1986; Wilf,
 1545 1999): if we can define an algorithm to gener-
 1546 ate compatible permutations uniformly at random
 1547 (such as that in Algorithm 2), we can easily adapt
 1548 it to count those permutations to yield an efficient
 1549 counting expression, as we do in Eq. 6.

1550 We start by introducing 2 basic operations of
 1551 *permute* and *shuffle*, along with the number of
 1552 possible outcomes of these operations.

1553 **Permute a list - permute(l)** Given a list of n ob-
 1554 jects, generate a permutation of these items. There
 1555 are $n!$ possible ways of permuting n items. An
 1556 efficient way for generating random permutations
 1557 is the Fisher-Yates-Knuth algorithm (Knuth, 1970).

1558 **Shuffle two lists - shuffle(A, B)** Given two
 1559 disjoint lists of distinct elements A, B of lengths
 1560 a, b respectively, generate a permutation σ of the
 1561 two lists of length $a + b$ in such a way that the
 1562 relative order of the items in the lists A and B is
 1563 respected in σ . This name and idea is based on
 1564 the popular way of shuffling two decks of cards
 1565 (Bayer and Diaconis, 1992). Its easy to see that
 1566 Algorithm ?? generates every possible shuffling
 1567 with equal probability. The total number of shuffles
 1568 of lists A, B is given in Eq. 6 as S_b^a .

1569 **Counting complete, compatible rankings** At
 1570 this point, we are ready to detail the expression of
 1571 $p(n, k, r)$ in Eq. 6, both the intuition and the proof
 1572 of uniformity. For this, we propose in Algorithm 2

Algorithm 1: Generate a random shuffle of lists A and B

```

1 for  $i \in [a + b]$  do
2    $rand \leftarrow$  random number in  $[0, 1]$ ;
3   if  $rand > 0.5 \vee B$  is empty  $\wedge A$  is non
   empty then
4      $\sigma(i) = pop(A)$ ;
5   else
6      $\sigma(i) = pop(B)$ ;
7   end
8 end

```

1573 to sample complete, compatible rankings and then
1574 adapt this sampling algorithm to a counting algo-
1575 rithm in Theorem 1.

1576 **Notation** We start by fixing the notation. Let
1577 β be a partial ranking of length k which includes
1578 item j in rank r , $\beta_1 \succ \dots \succ \beta_r = j \succ \dots \succ \beta_k$.
1579 Let η be a disjoint set of $n - k$ items that have not
1580 been ranked and which includes the unobserved
1581 item i . The goal is to generate (i) a compatible
1582 ranking with β (a ranking σ of all the items in such
1583 a way that the relative ordering of the items of β is
1584 maintained) and (ii) which ranks item i before item
1585 j . We denote the " s -head" of a list to the items in
1586 the first s positions in that list.

1587 **Intuition** We are now ready to explain the intu-
1588 ition. Each of the possible compatible permutations
1589 that rank i before j is generated in the following
1590 way:

1591 Algorithm 2 generates permutations that rank
1592 item j at position s , item i before j and we iterate
1593 for all possible values of s . First, in line 2
1594 we select $s - 1$ items randomly from η , where the
1595 order of the items matter (i.e., a variation). Then,
1596 we insert item i in a random position of this list,
1597 denoted η_{head} in line 3. In line 4 we shuffle these
1598 two lists, i.e., η_{head} and the r -head of β , β_{head} ,
1599 i.e., the sublist with the items that are ranked before
1600 j . The result of the shuffling process is the $s + r$ -
1601 head of the output permutation σ . We permute
1602 the rest of the unobserved items denoting these list
1603 η_{tail} , in line 6. Finally, we shuffle this list η_{tail}
1604 and the $k - r$ -tail of η in line 7. The result of
1605 this shuffle is the tail of σ . Finally, in line 8 we
1606 return the concatenation of $\sigma_{head}, j, \sigma_{tail}$, which
1607 is clearly a compatible permutation with β as the
1608 relative order of the items in β is maintained in the
1609 output.

1610 It is easy to see that Algorithm 2 generates the

Algorithm 2: Generate a random ranking among those compatible with β

```

1 for  $s \in [n]$  do
2    $\eta_{head} \leftarrow$   $s - 1$  items from  $\eta$  where the
   order matters ;
3    $\eta_{head} \leftarrow$  insert  $i$  in  $\eta_{head}$  ;
4    $\sigma_{head} \leftarrow$  shuffle( $\eta_{head}, \beta_{head}$ ) ;
5    $\eta_{tail} \leftarrow \eta \setminus \eta_{head}$  ;
6    $\eta_{tail} \leftarrow$  permute( $\eta_{tail}$ ) ;
7    $\sigma_{tail} \leftarrow$  shuffle( $\eta_{tail}, \beta_{tail}$ ) ;
8   return ( $\sigma_{head} \succ j \succ \sigma_{tail}$ ) ;
9 end

```

1611 target permutations uniformly at random. Follow-
1612 ing a classic strategy on Enumerative Combina-
1613 torics (Stanley, 1986; Wilf, 1999) we use this algo-
1614 rithm as a proof for $p(n, k, r)$.

1615 **Theorem 1.** The number of complete permutations
1616 of n items compatible with partial ranking β that
1617 rank the unobserved item i before the observed
1618 item j is given by the following expression, 1618

$$1619 \quad p(n, k, r) = \sum_{i=0}^{n-k-1} V_{n-k-1}^i * (i + 1) \quad (7)$$

$$1620 \quad * S_{i+1}^r (n - k - i - 1)! \quad (8)$$

$$1621 \quad * S_{k-r-1}^{n-k-i-1} / t(n, k). \quad (9)$$

1622 *Proof.* It is easy to see that in Algorithm 2 there is
1623 a bijection between the permutations in the target
1624 (that is, the permutations compatible with β for
1625 which $i \succ j$) and each outcome of Algorithm 2.
1626 Clearly, for uniform at random outcomes of the
1627 *shuffle* and *permute* operations, the outcome of
1628 Algorithm 2 will be random as well. Therefore,
1629 the number of possible outcomes of the algorithm
1630 equals the number of permutations in the target. 1631

1632 It follows that each term in $p(n, k, r)$ Each term
1633 in the previous expression comes from a different
1634 line in 2: 1634

- Line 2: The number of variations of i items
1635 out of $n - k - 1$ is V_{n-k-1}^i . 1636

- Line 3: There are $s + 1$ ways of inserting item
1637 i , thus the term $(r + 1)$. 1638

- Line 4: There are S_{s+1}^r ways of shuffling
1639 η_{head} and β_{head} . 1640

1641 • Line 6: There are $(n - k - s - 1)!$ possible
1642 permutations of the items in η_{tail} .

1643 • Line 7: There are $S_{k-r-1}^{n-k-s-1}$ ways of shuffling
1644 the two tails.

1645 • Line 8: Finally, since we compute the propor-
1646 tion by dividing among the total number of
1647 compatible permutations.

1648 By repeating this process for all $s < n - k - 1$
1649 the proof is completed.

1650 \square

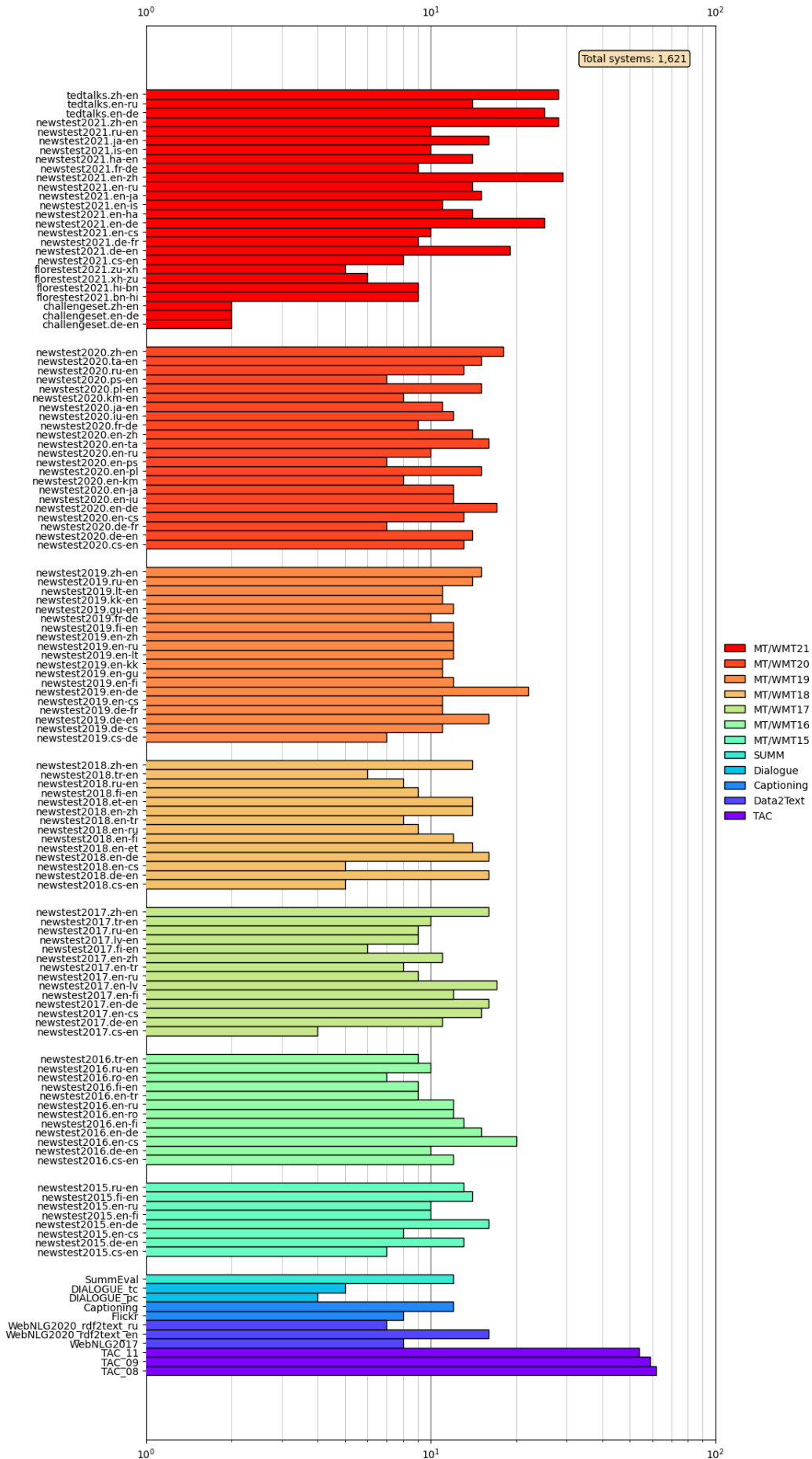


Figure 13: Number of systems in each dataset (log scale)

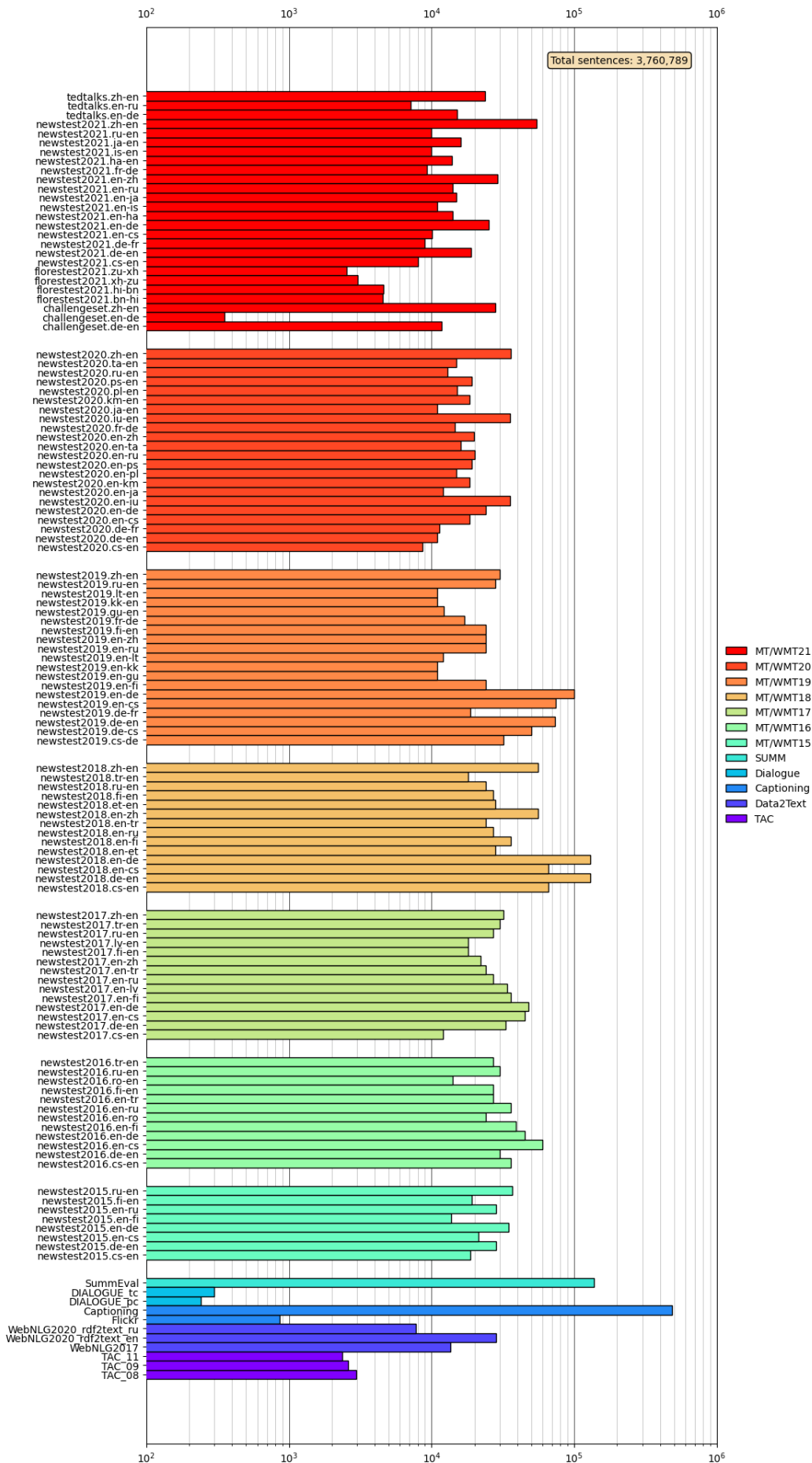


Figure 14: Number of sentences in each dataset (log scale)

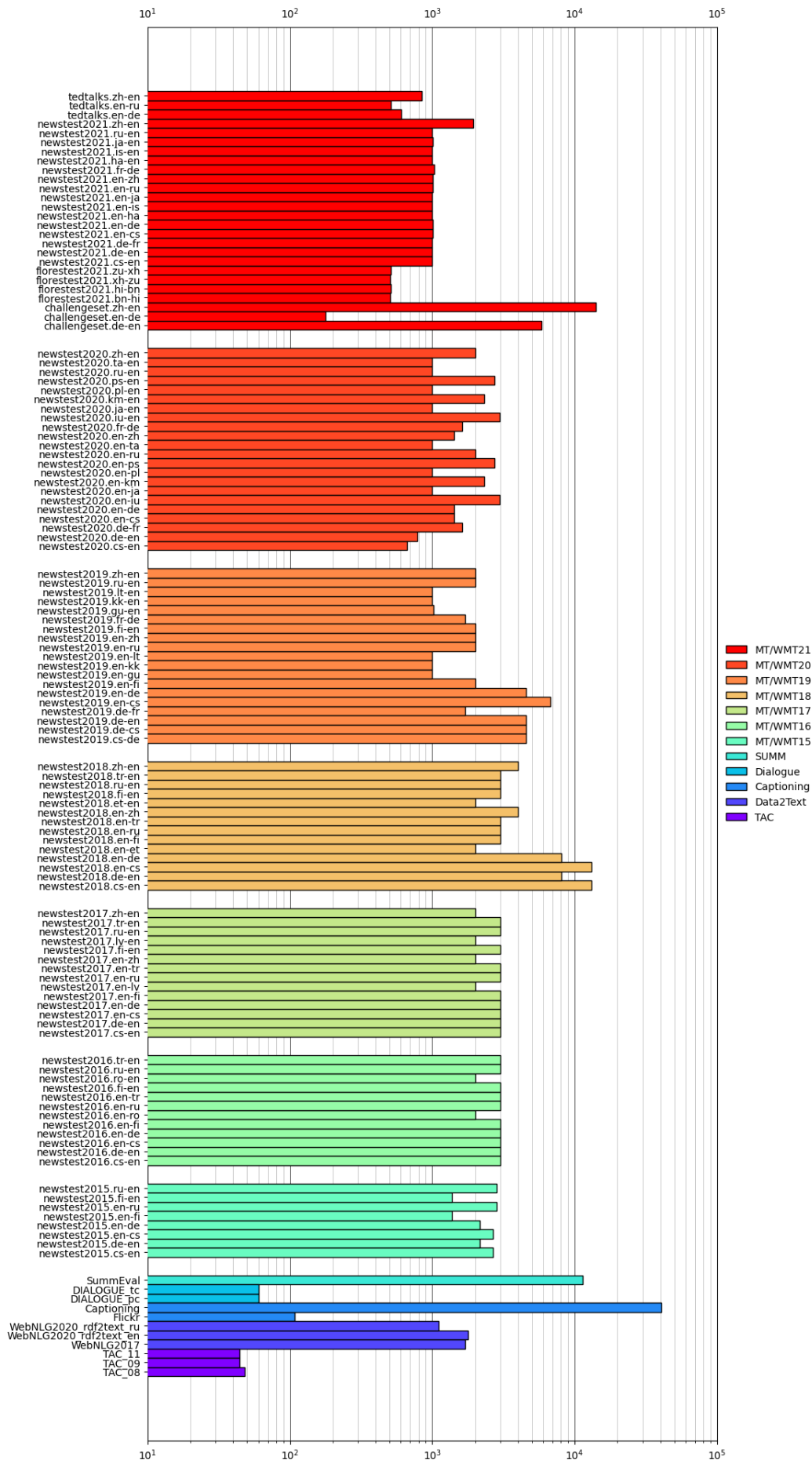


Figure 15: Average number of sentences per system in each dataset (log scale)