

Anonymity at Risk? Assessing Re-Identification Capabilities of Large Language Models in Court Decisions

Alex Nyffenegger^{1,3*} Matthias Stürmer^{2,3} Joel Niklaus^{2,3,4*}

¹University of Fribourg ²Bern University of Applied Sciences

³University of Bern ⁴Stanford University

Abstract

Anonymity in court rulings is a critical aspect of privacy protection in the European Union and Switzerland but with the advent of LLMs, concerns about large-scale re-identification of anonymized persons are growing. In accordance with the Federal Supreme Court of Switzerland (FSCS), we study re-identification risks using actual legal data. Following the initial experiment, we constructed an anonymized Wikipedia dataset as a more rigorous testing ground to further investigate the findings. In addition to the datasets, we also introduce new metrics to measure performance. We systematically analyze the factors that influence successful re-identifications, identifying model size, input length, and instruction tuning among the most critical determinants. Despite high re-identification rates on Wikipedia, even the best LLMs struggled with court decisions. We demonstrate that for now, the risk of re-identifications using LLMs is minimal in the vast majority of cases. We hope that our system can help enhance the confidence in the security of anonymized decisions, thus leading the courts to publish more decisions.

1 Introduction

The swift advancements in Natural Language Processing (NLP) (Vaswani et al., 2017; Brown et al., 2020; Ouyang et al., 2022; Khurana et al., 2023) have introduced new challenges to the security of traditional legal processes (Tsarapatsanis and Aletras, 2021). As public access to data increases in tandem with digital advancements (Katz et al., 2023; EUGH, 2018; Lorenz, 2017), the potential risks associated with data disclosure have become increasingly significant. Larger and more capable Language Models (LMs), more powerful vector stores and potent embeddings together have the capacity to extract unintended information from pub-

* Equal contribution.

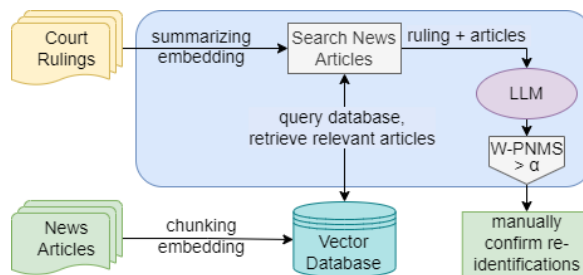


Figure 1: Re-identification framework

lic data (Borgeaud et al., 2022; Carlini et al., 2021; Roberts et al., 2020; Alkhamissi et al., 2022; Ippolito et al., 2023; Carlini et al., 2023). This poses a security risk, as identifying individuals in legal proceedings can lead to privacy breaches, leading to inequity in insurance, enabling extortion, and even risking public defamation.

Over the past decade, at least 18 requests for name changes following re-identification of convicts have been registered in Switzerland, indicating existing issues due to imprudent media coverage (Stückelberger et al., 2021). The number of cases involving unlawful disclosure of personal information is likely to rise. Therefore, the prevention of re-identification is critical not only for the protection of the accused, but also for the courts. Munz (2022) even suggests that the state could be held accountable for non-monetary damages to judged persons, underscoring the urgent need for courts to address the re-identification issue proactively. Vokinger and Mühlematter (2019) and Niklaus et al. (2023a) have shown that companies can be re-identified by simply extracting information from the court decisions with regular expressions and matching it with public databases.

We see strong parallels between re-identification and penetration testing, where cyber-security experts attempt to find and exploit vulnerabilities in a computer system (Altulaihan et al., 2023). To the best of our knowledge, we are the first to study

the re-identification task of anonymized persons from court decisions. We provide a framework for anonymization teams in courts and researchers alike to battle-test anonymizations of cases (illustrated in Figure 1).

In this work, we investigate to what extent Large Language Models (LLMs) like LLaMA-2, GPT-4 or BLOOM (Touvron et al., 2023a; OpenAI, 2023; Scao et al., 2023) can re-identify individuals in Swiss court decisions. Our main findings reveal that while top models identify persons from masked Wikipedia articles, they struggle with the harder task of court decision re-identification. Only in cases we manually re-identified in a painstaking process and thus know re-identification is possible, and using a highly curated set of manually identified relevant news articles, they are capable of identifying the anonymized defendants from cases. Finally, in detailed ablations, we identify three main factors influencing the re-identification risk: input length, model size, and instruction tuning.

With our research, we are testing whether affected parties in rulings could still be identified despite anonymization. Thus the results from our research can guide legal entities, data privacy advocates, and NLP practitioners in devising strategies to mitigate potential re-identification risks. This is relevant beyond Switzerland, as anonymization of court rulings became mandatory across the EU with the introduction of the GDPR (See Appendix F.4). The German Supreme Court even ruled that all rulings should be anonymized and published. However, in 2021 barely one percent of rulings were being published (Hamann, 2021) (See Appendix F.4). This may be partially caused by fears that publications are insufficiently anonymized and courts could be held accountable. We hope that our framework will be used to ensure privacy for anonymized documents and will therefore lead to more cases being published across Europe. In the spirit of open science, we release all datasets and code for reproducibility with permissive licenses.¹²³

¹<https://huggingface.co/datasets/rcds/wikipedia-persons-masked>

²https://huggingface.co/datasets/rcds/swiss_rulings

³<https://github.com/Skatinger/Anonymity-at-Risk-Assessing-Re-Identification-Capabilities-of-Large-Language-Models>

Main Research Questions

This study addresses three research questions:

RQ1: Performance of LLMs on re-identifications: How effectively can various LLMs re-identify masked persons within Wikipedia pages and in Swiss court rulings?

RQ2: Influential Factors: What are the key factors that influence the performance of LLMs in re-identification tasks?

RQ3: Privacy Implications: How will evolving LLM capabilities and their use in re-identifications affect the preservation of privacy in anonymized court rulings in Switzerland?

By addressing these questions, we aim to highlight LLMs' capabilities and limitations in re-identification tasks and enhance understanding of required privacy considerations in the ongoing digital transformation of legal practice.

Contributions

The contributions of this paper are threefold:

- We curate and publish a unique, large-scale Wikipedia dataset with masked entities.
- We introduce new metrics to evaluate performance of re-identifications of persons within texts. Using those metrics, we provide a thorough evaluation and benchmark of various state-of-the-art LLMs in the context of re-identifying masked entities within Wikipedia entries and Swiss court rulings. This includes an exploration of the most critical factors influencing model performance.
- We underscore and investigate the potential privacy implications of using LLMs for re-identification tasks.

2 Related Work

Chen et al. (2017) used LMs for machine reading to answer open domain questions, providing models with necessary context from Wikipedia articles for knowledge extraction.

LMs as Knowledge Bases With the advent of the transformer (Vaswani et al., 2017), more powerful models became able to store information within their parameters (Petroni et al., 2019; Alkhamissi et al., 2022) and the idea of using models directly without additional context became viable. Petroni et al. (2019) found that LMs can be used as knowledge bases, drawing information from their training set to answer open domain questions. Roberts et al. (2020) went a step further and evaluated different

sizes of T5 models (Raffel et al., 2020) showing that larger models can store more information, but unlike other Question Answering (QA) systems are not able to show where facts come from. This is especially a problem when models hallucinate an answer when they are unsure, as correctness of an answer is hard to factually check without sources (Petroni et al., 2019). With Lewis et al. (2020) finding that good results on open domain question answering heavily depends on the overlap of questions and training data, Wang et al. (2021) showed that even without overlapping data, knowledge retrieval is possible, although with much lower performance. Wang et al. (2021) discovered that knowledge exists in model parameters but is not always retrieved effectively. They introduced QA-bridge-tune, a method enabling more reliable information retrieval from model parameters.

Retrieval Augmented Generation To improve reliability of results even further Lewis et al. (2021) introduced the combination of pretrained models and a dense vector index of Wikipedia, finding that QA tasks are answered with more specific and factual knowledge than parametric models alone, while hallucinations are reduced when using Retrieval Augmented Generation (RAG) (Shuster et al., 2021). Recent research (Kassner et al., 2021) shows that multilingual models excel in knowledge retrieval tasks, particularly when questions match the language of the training data. However, inter-language retrieval underperforms, indicating lower performance for questions in a different language than the data source (Jiang et al., 2020).

Re-Identification Studies Staab et al. (2023) managed to extract personal information at scale, by using comments from Reddit users to identify clues such as age, gender or location. The exact names were not extracted. In re-identification within court rulings, Vokinger and Mühlematter (2019) linked medical keywords from public sources to those in court rulings, identifying persons through associations with drugs and medicine. This successful partial re-identification suggests language models might achieve similar results. Niklaus et al. (2023a) used regular expressions to extract project ids from court decisions which they matched with publicly available data from the simap database of public procurement tenders. Although both works manage to re-identify companies from court decisions, they are limited to very specific attack vectors. In this work, we study the risk of large scale general attacks using LLMs.

3 Collaboration with the Supreme Court

To ensure responsible research and maximize downstream usability, we collaborated closely with the Federal Supreme Court of Switzerland (FSCS). The FSCS currently uses regular expressions and a BERT-based (Devlin et al., 2018) token classifier to provide suggestions to human anonymizers for what entities should be masked. In a prior project, we improved their system’s recall on anonymization tokens from 83% to 93% by pre-training a legal specific model. In this work, we partner with their anonymization team for testing.

4 Datasets

To perform our case study, we select Switzerland for its richness in published data – both newspapers and court decisions – and its high privacy standards. We created three datasets: First, the *Court Decisions* dataset consisting of anonymized Swiss case law serves as a substantial benchmark for evaluating re-identification risks in court rulings. Only the FSCS can assess the outcomes on this dataset, as they exclusively possess knowledge of the involved individuals. The second dataset called *Legal-News Linkage* provides a small sample of manually re-identified court rulings, elucidating potential re-identification cases by LLMs. Finally, we curated a dataset consisting of *Wikipedia* biographical pages and automatically anonymized it. This extensive dataset facilitates a broad analysis of re-identification techniques using LLMs.

4.1 Court Decisions Dataset

We used the Swiss caselaw corpus by Rasiyah et al. (2023) to benchmark re-identification on court rulings. The FSCS likely rules the most publicised cases as the final body of appeal in Switzerland and offered to validate re-identifications in a limited fashion, leading us to discard cases from other courts. This decision aligned well with the fact that federal court cases occur more often in the news, elevating the likelihood of potential re-identifications. To make sure that all evaluated models have been trained on relevant data, we only used cases from the year 2019, resulting in approx. 8K rulings.

4.2 Legal-News Linkage Dataset

The Court Decisions dataset offers large scale, but no ground truth (i.e., we do not know if a re-identification is at all possible). For this reason, we created the Legal-News Linkage Dataset,

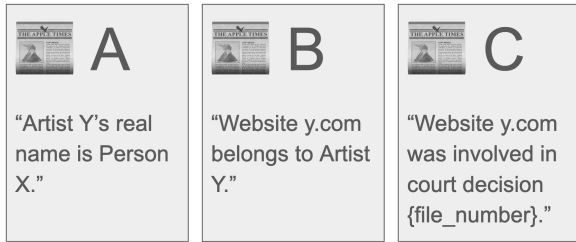


Figure 2: Simplified example of content in newspaper articles. Note that only using all three articles, the re-identification is made possible.

where we have high certainty of the anonymized person. We created this dataset by manually linking court rulings and newspaper articles using keywords like the file number of the court decision (e.g., 4A_375/2021) or the penalty (e.g., 10 years in prison). It was not possible to construct a systematic process to create this dataset at scale because of individual idiosyncrasies of each decision. The rarity of such cases in Swiss news and the intensive manual effort involved limited our dataset to these seven instances. In an iterative process we accumulated roughly 100 related newspaper articles per court decision by searching for information found in the seed newspaper articles, such as the person’s name. This accumulation was necessary because there are multiple newspaper articles for each court case mentioning different aspects of the person. One article is not enough; only in aggregation, it is possible to perform the re-identification (illustrated in Figure 2). Due to cost reasons, we were not able to use the full newspaper dataset. To represent a realistic scenario, we added 1000 unrelated newspaper articles instead. This ensures the linkage process from news articles to court rulings is successful. The curated dataset includes seven court rulings and approx. 2000 news articles. To maintain privacy, we do not publish this dataset. The news articles are proprietary and were sourced from swissdox.ch.

4.3 Wikipedia Dataset

The Court Decisions dataset is large and realistic but offers no ground truth. The Legal-News Linkage dataset is realistic and offers ground truth but is small. With the Wikipedia dataset, offering ground truth at scale at the expense of realism, we can study the effect of various factors on model’s re-identification performance (see Section 7.2). We randomly chose 10K from 69K examples to mirror the Court Decisions dataset’s size. Construction involved three steps: 1) We filtered Wikipedia pages

marked as persons by their length (> 4K characters) as a proxy for importance/prevalence, 2) we stored paraphrased Wikipedia pages alongside original content to assess model reliance on exact training text phrasing (Carlini et al., 2021), and 3) we replaced all occurrences of the person’s name with a mask token. Further details on the construction process are in Appendix E.2.

5 Metrics

Re-identification of persons is a known problem for imaging (Karanam et al., 2018), but comparable metrics for re-identifications within texts are, to the best of our knowledge, not established. Unlike memorization verification (Carlini et al., 2023) the re-identification of persons requires the model to be able to connect knowledge over multiple datapoints (see Section 4.2). This means that information does not always exist in a single knowledge triple, but is connected over several ones or requires several ones to lead to a re-identification. To allow the quantification of produced results, we introduce the following four novel metrics to measure re-identification performance of a person in a text:

Partial Name Match Score (PNMS) evaluates predictions against a regular expression requiring any part of a persons’s name to be a match for the prediction to be considered as correct. For example, "Max Orwell" would match "George Orwell". This allows for matches with predictions that only contain a part of the name. Manual experimentation suggested that persons can be re-identified by using just a part of their name. The predicted name might be near exact, hence the allowance for partial matches. The metric accepts n predictions and deems any collection of predictions correct if at least one of the n predictions is correct.

Normalized Levenshtein Distance (NLD) is introduced to assess the precision of predictions deemed correct by PNMS. Given that there is no clear-cut distinction between correct and incorrect, using the Levenshtein distance provides a more nuanced perspective on how close the predictions are to the target. For the top five predictions, the smallest distance of all five was used. Using the best distance of n given predictions, the distance was normalized against the length of the target name to avoid distortions in results. As example, the distance between "Alice Cooper" and "Alina Cooper" would be two, and with the normalization by $len("Alina Cooper")$ applied result in 0.16.

Last Name Match Score (LNMS) works the same way as PNMS, but only the last name is considered. The last name is defined as the last whitespace-separated part of a full name string. Partial matches are accounted as correct as well meaning that the name "Mill" would also be counted as correct if the target was "Miller". This overlap might cause a very slight imprecision but does not lead to problems in evaluations as all models have the same advantage.

Weighted Partial Name Match Score (W-PNMS) blends PNMS and the LNMS using a weighted sum, emphasizing the significance of last names for re-identification. Let $\alpha = 0.35$ be the weight for PNMS. Thus, W-PNMS is calculated as $W\text{-PNMS} = \alpha \times \text{PNMS} + (1 - \alpha) \times \text{LNMS}$.

The metrics are designed to recognize both exact and partial name matches. We prompted our models to predict full names, yet texts often contain name variations such as "J. Doe" or "Mr. Doe" prompting us to accommodate partial name matches and measure the NLD. Our methodology overlooks spelling variations and multilingual representations, which, in our experience, are rare enough to safely de-prioritize.

6 Experimental Setup

We ran models using the HuggingFace Transformers library on two 80GB NVIDIA A100 GPUs, using default model configurations in 8-bit precision. For efficiency, we only used the first 1K characters of each Wikipedia page. For court rulings, we extended input length to 10K characters, maximizing model sequence lengths. Sequences exceeding maximum input length were automatically truncated. We used temperature 1 and considered the top 5 predictions. See Figure 9 for a high level overview of our code architecture.

6.1 Prompt Engineering

The effectiveness of model responses is significantly influenced by input prompt design (Liu et al., 2022; Wei et al., 2023). Various models require distinct prompting strategies to perform well. We tailored prompts for each model, but without extensive optimization, ensuring a consistent effort across models. Experimental results indicated that once a prompt communicated the re-identification task to a model, further refinement of the prompt did not substantially improve any metrics.⁴

⁴Prompt examples in Appendix F.2

6.2 Retrieval Augmented Generation

To estimate how well an LLM could use information from news articles without training one we used RAG (Lewis et al., 2021): From the 1.7K news articles gathered for the legal-news linkage dataset, we split texts into 1K-character chunks, embedded them with OpenAI's text-embedding-ada-002, and stored the embeddings in a Chroma vector database (<https://www.trychroma.com/>). To re-identify a ruling, we fed it to GPT-3.5-turbo-16k, prompting it to summarize the decision, emphasizing facts in news articles and retaining key details, including masked entities.

We then embedded this shorter version the same way as the articles and matched against the stored article chunks using the similarity search provided by Chroma. The top five retrieved documents together with the shortened version of the ruling were given to GPT-4 with the prompt to use the information given in the documents to re-identify the person referred to as <mask>. This method skips the large training effort required to store knowledge in LLMs while still demonstrating the capability of LLMs to comprehend multi-hop information from news articles and apply it to re-identification.

6.3 Evaluated Models

For the rulings dataset, we utilized models that were specifically trained on news articles and court rulings, alongside the two multilingual models, GPT-4 and mT0. The selection of these models, as detailed in Table 3, was informed by their pre-training on relevant news content. For the Wikipedia dataset, we used various models with different pre-training datasets and architectures. By using a large and diverse selection of models, prominent factors for good performance can be found more easily and results are more reliable. A full list is available in Table 3. All models except the commercial models ChatGPT and GPT-4 are publicly available on the HuggingFace Hub.

6.4 Baselines

We propose two baselines for easier interpretation:

Random Name Guessing Baseline predicts for every example five first and last names paired up to full names at random. This gives a good impression on predictive performance when models understand the task or at least guess while not actually knowing the entities name. Names were chosen from a GPT-3.5-generated list of 50 names.

Majority Name Guessing Baseline predicts the top five common first and last names for the English language, with the names being paired up to full names in their order of commonness. First names were sourced from the US Social Security Administration⁵ and last names from Wiktionary⁶.

7 Results

7.1 Performance on Court Rulings

Re-identifications on Rulings Test Set We show results in Figure 3. Among all evaluated models, only `legal_xlm_roberta` (561M) and `legal_swiss_roberta` (561M)⁷ re-identified a single person from 7673 rulings. As discussed later in Section 7.2, this aligns with expectations since evaluated models, excluding GPT-4 and mT0, do not meet key factors for effective re-identification: input length, model size, and instruction tuning. Despite their smaller size and lack of instruction tuning, these models made some reasonable guesses. Conversely, larger multilingual models like GPT-4 and mT0 failed to give credible guesses. We tested GPT-4 on the top 50 most reasonably predicted examples from other models. Potentially reflecting OpenAI’s commitment to privacy alignment, GPT-4 consistently indicated that the person was not present in the text, refraining from leaking training data or making speculative guesses. mT0, trained on mC4 likely containing Swiss news articles, underperformed despite strong performance on the Wikipedia dataset, treating the text as cloze test instead of attempting to guess names. While mT0’s predictions lacked meaningful output, the success of smaller models to predict some believable speculations suggests they might not have been relying solely on chance but made informed guesses. Most predictions corresponded to words already present in the ruling or were not a name. Excluding the few viable predictions (titled *good*), the others consisted of empty predictions or single letters.

Re-identification with Retrieval Applying the same models on the legal-news linkage dataset, the results were not better even though for this small dataset we had the confirmation that all rulings were re-identifiable with the information in the training data. None of the models were able to

⁵<https://www.ssa.gov/oact/babynames/decades/century.html>

⁶[https://en.wiktionary.org/wiki/Appendix:English_surnames_\(England_and_Wales\)](https://en.wiktionary.org/wiki/Appendix:English_surnames_(England_and_Wales))

⁷Model details in Appendix 3

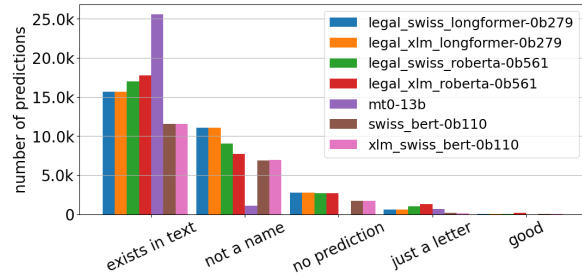


Figure 3: Prediction categories on rulings dataset. "good" are the only possibly correct predictions.

predict any person correctly. However, using the RAG approach worked much better. When passing the relevant news articles and the corresponding court ruling to the context, GPT-3.5-turbo-16k was able to identify 4 out of 7 entities, with the full name for one example. GPT-4 performed even better, correctly identifying 5 out of 7, with the full name for one example. Interestingly, the two cases which were easiest for us humans to identify were not identified by either model. This result not only suggests that re-identification by training on enough news articles could be possible, but that models powerful enough to understand the task and the given information are capable of using not only their training data information, but simultaneously ingest relevant additional information. It could even be possible to re-identify decisions without any pre-training by ingesting the full news dataset and embed information on a large scale, leading to large scale re-identifications in the worst case.

7.2 Factors for Re-identification on Wikipedia

Performance in re-identification tasks varied significantly across models (see Table 4 for the full results). Some larger models such as Flan_T5 or mT0 reach scores above 0.3 or for GPT-4 even above 0.6 for W-PNMS with very low NLD while models like Pythia or Cerebras-GPT failed completely, below the guessing baseline even. Table 1 lists the top performers on the Wikipedia dataset.

Original vs paraphrased In Table 2 we compare the effect of paraphrases on re-identification performance. We find models to perform slightly better on the original text, both when we constrain the input by the number of characters and by a number of sentences (to ensure that the same amount of information is given). Note that the average paraphrased sentence is significantly shorter than the average original sentence (95 vs 141 characters, see Appendix F.1). We see two possible reasons:

Model	Size [B]	PNMS \uparrow	NLD \downarrow	W-PNMS \uparrow
GPT-4	1800	0.71	0.17	0.65
GPT-3.5	175	0.52	0.23	0.46
mT0	13	0.37	0.42	0.31
Flan_T5	11	0.37	0.45	0.30
incite	3	0.37	0.53	0.30
Flan_T5	3	0.35	0.48	0.29
BLOOMZ	7.1	0.34	0.45	0.29
T0	11	0.34	0.45	0.28

Table 1: Models w/ W-PNMS ≥ 0.28 on Wikipedia dataset

Data Config	PNMS \uparrow	NLD \downarrow	LNMS \uparrow	W-PNMS \uparrow
input constrained to 1000 characters				
original	0.35 ± 0.04	0.52 ± 0.05	0.25 ± 0.03	0.29 ± 0.03
paraphrased	0.33 ± 0.03	0.48 ± 0.03	0.24 ± 0.02	0.27 ± 0.02
input constrained to eight sentences				
original	0.33 ± 0.05	0.57 ± 0.11	0.22 ± 0.04	0.26 ± 0.05
paraphrased	0.28 ± 0.03	0.51 ± 0.04	0.19 ± 0.03	0.22 ± 0.03

Table 2: Mean and std over top performers (incite_instruct, Flan_T5, T0, BLOOMZ, mT0)

1) information is lost in paraphrasing due to shorter outputs, and 2) it is harder for the models to retrieve the information because of changed surface form compared to the training data. To simulate a more realistic scenario closer to re-identifying court decisions, we use the paraphrased texts henceforth.

Model Size Comparing differently sized versions of a model as shown in Figure 4, we observed a clear performance boost as model size increases, consistent with prior research suggesting better knowledge retrieval with larger models (Roberts et al., 2020). Performance typically improves significantly when transitioning from smaller to medium-sized models, though the gains diminish for larger models. While not all models performed the same for the larger model sizes, the general performance progression indicates that performance gains stagnate when models are scaled beyond their sweet spot. On average this turning point appears to be at around 3B parameters but varies for different models with some models still reaching better performances for much larger sizes. Models with low performance show only a minor improvement with increased size. The small increase might be due to the model understanding the task better but still not being able to retrieve the requested name, but by chance giving more diverse answers and coincidentally matching some predictions.

Input length Figure 5 reveals that performance improves with increasing input size, though the

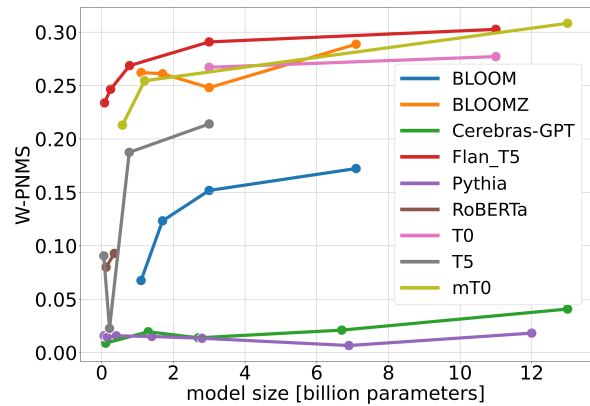


Figure 4: Re-identification score by parameter count

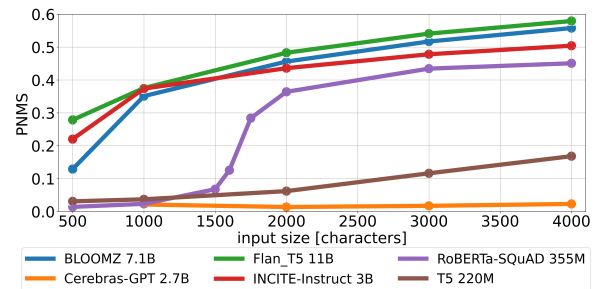


Figure 5: Re-Identification score across input lengths

degree of improvement varies among models. For most models, performance increased strongly until 2K characters (approx. 500 tokens) and then flattened. The model roberta_squad which is only 355M parameters but fine-tuned on a QA dataset was able to gain a strong increase in performance nearly matching the top performers.

Instruction tuning As shown in Figure 6, instruction tuned models perform much better at re-identification. Even though both versions of each model were pretrained on the same datasets and contain the same knowledge, the instruction tuned models were far more likely to understand the task and retrieve the correct name, which is consistent with previous research (Longpre et al., 2023; Ouyang et al., 2022; Muennighoff et al., 2023).

Decoding strategies We see in Figure 7 that

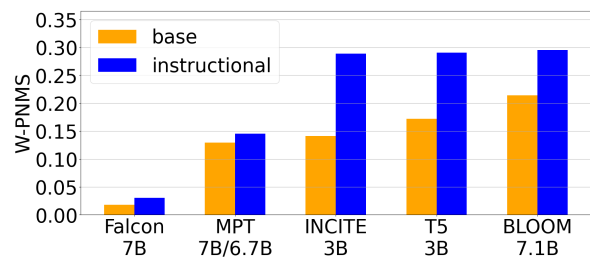


Figure 6: Base vs. instruction tuned performance

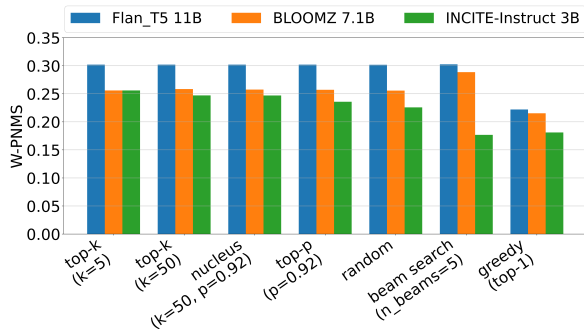


Figure 7: Decoding strategies of top performing models

overall the variation in performance across decoding strategies is small. Greedy decoding performed much worse, likely because it naturally only considers the top-1 prediction. Performance varies most for beam search: Incite_instruct performed worst, while BLOOMZ achieved its best results. Looking at the precision of decisions, the NLD is better for predictions produced with beam search, meaning beam search can deliver more precise re-identifications, while top-k might find generally more likely names, but not necessarily the exact full name. With two out of three evaluated models performing best with beam search and NLD being best with this sampling strategy we used beam search for all other experiments.

Re-Identification methods In Figure 8 we compare fill mask, QA and text generation models across model sizes. We excluded text generation models below the random name guessing baseline because they failed to follow the instructions (i.e., Pythia, Cerebras-GPT, Falcon, Falcon-Instruct, GPT-J). We find models performing the fill mask and QA tasks to underperform the text generation models across the board, and even at the same model size. While performance increases for models performing fill mask, the opposite happens for models doing QA when scaling up model size. Given that most large-scale models are text generation models, they tend to outperform fill mask and QA counterparts. The improved performance of these models can be attributed to their ability to retain more information, a characteristic inherent to larger models (Roberts et al., 2020).

8 Conclusions and Future Work

8.1 Answering the Main Research Questions

RQ1: Performance of LLMs on re-identifications: How effectively can various LLMs re-identify masked persons within Wikipedia pages

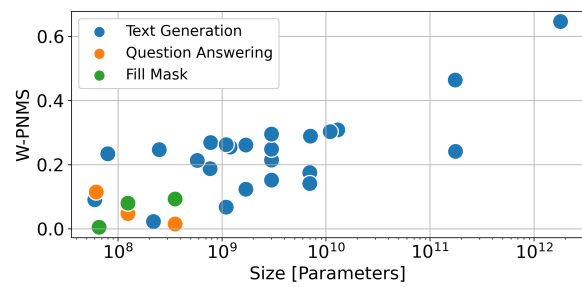


Figure 8: Relation of re-identification score to model size across model types

and in Swiss court rulings?

We find that vanilla LLMs can not re-identify individuals in Swiss court rulings. Additionally, relatively small models trained on Swiss news articles and court rulings respectively can barely guess credible names. Finally, by augmenting strong LLMs with retrieval on a manually curated dataset, a small subset of individuals can be re-identified.

RQ2: Influential factors: What are the key factors that influence the performance of LLMs in re-identification tasks?

We identified three influential factors affecting the performance of LLMs in re-identification tasks: model size, input length, and instruction tuning.

RQ3: Privacy Implications: How will evolving LLM capabilities and their use in re-identifications affect the preservation of privacy in anonymized court rulings in Switzerland?

We demonstrate that, for now, significant privacy breaches using LLMs on a large scale are unattainable without considerable resources. Yet, the Wikipedia benchmark revealed that larger models, when exposed to adequate pre-training information, can proficiently identify anonymized persons. As LLMs get more powerful and integrated with tools like retrieval (Lewis et al., 2021), coding and arbitrary API access (Schick et al., 2023), we fear heightened re-identification risks. Therefore, we urge courts to perform checks like outlined in our study on a regular basis before publication to safeguard privacy. To set an example, we are in close contact with the FSCS to transfer insights into their anonymization practice. Risks of the courts not having sufficient access to trained personnel with the necessary skills for such testing remain.

8.2 Conclusions

Similar to penetration testing in cyber-security, we battle-tested the anonymization of Swiss court

cases using LLMs. Currently, the risk of vanilla LLMs re-identifying individuals in Swiss court rulings is limited. However, if a malicious actor were to invest significant resources by pre-training on relevant data and augmenting the LLM with retrieval, we fear increased re-identification risk. We identified three major factors influencing re-identification performance: the model’s size, input length, and instruction tuning. As technology progresses, the implications for privacy become more pronounced. It is imperative to tread cautiously to ensure sanctity of privacy in court cases remains uncompromised.

8.3 Future Work

Liu et al. (2023) showed that models extract information better if it is located at the start or end of large contexts. For the large models which can ingest full court rulings, this could mean that ordering parts of the rulings by their relevancy for re-identifications could improve chances for successful re-identifications. Further research is required to analyze which parts of rulings are the most relevant for re-identification. Specific pre-training of large models on relevant data and sophisticated prompting techniques such as chain of thought (Wei et al., 2023) may increase re-identification risk. In this work, we only considered information in textual form, either embedded in the weights by pretraining or put into the context with retrieval. Future work may also investigate the use of more structured information, such as structured databases or knowledge graphs. We believe the Swiss court system serves as an ideal candidate for studying re-identification due to the high privacy standards and data richness both in newspapers and published court decisions. In future work, we would like to extend our analysis to other countries with similar concerns, such as many from the EU.

Ethics and Broader Impact

Abundant publication of court rulings is crucial for judicial accountability and thus for a functioning democratic state. Additionally, it greatly facilitates legal research by removing barriers to case documents access. However, courts hesitate to publish rulings, fearing repercussions due to possible privacy breaches. Solid automated anonymization is key for courts publishing decisions more plentiful, faster, and regularly. Strong re-identification methods can be a valuable tool to stress-test anonymization systems in the absence of formal guarantees

of security. However, re-identification techniques, akin to penetration testing in security, are dual-use technologies by nature and thus pose a certain risk if misused. Fortunately, our findings indicate that without a significant investment of resources and expertise, large scale re-identification using LLMs is currently not feasible.

Limitations

The metrics employed to gauge the re-identification risk present inherent ambiguities. By comparing exact name matches and assessing the general similarity to the target name, we can infer the likelihood of manual re-identification. Yet, for lesser-known individuals or those with widespread names (such as the common Swiss first-name Simon or last-name Schmid), a generic first name paired with a surname might be insufficient for precise identification. Thus, manual scrutiny remains necessary to distill the correct person from the model’s suggested candidates. Essentially, while models scoring highly on our metrics can suggest potential identities, they might not always identify a person with certainty, especially when common names or lesser-known individuals are involved. In this work, we always checked possible re-identifications with high scores manually and therefore recommend this to future researchers and practitioners.

Additional to our ablations on input length, instruction tuning, decoding strategies, re-identification methods, paraphrasing, and model size, we would like to investigate the effect of tokenization on re-identification risk. The hidden challenge here is that constructing a controlled experiment to isolate the effect of tokenization requires access to models pretrained with identical architectures but varying vocabularies/tokenizers, which, to our knowledge, are not available (neither in LLAMA, BLOOMZ, Flan-T5, etc.). This, together with the enormous costs of pretraining such models, limited the feasibility of such an investigation in this work.

Acknowledgements

We thank Daniel Brunner from the Swiss federal Supreme court for evaluating our predictions on court rulings. A kind thanks goes to Dominique Schläfli, the law student who helped us navigate the complicated texts of court rulings to curate a dataset of re-identified court rulings allowing us to benchmark different re-identification strategies.

We thank the anonymous reviewers for their detailed feedback.

References

- Together AI. 2023. [Releasing 3B and 7B RedPajama-INCITE family of models including base, instruction-tuned & chat models.](#)
- Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. 2022. [A Review on Language Models as Knowledge Bases.](#) *arXiv:2204.06031 [cs]*. ArXiv: 2204.06031.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [Falcon-40B: an open large language model with state-of-the-art performance.](#)
- Esra Abdullatif Altulaihian, Abrar Alismail, and Mounir Frikha. 2023. [A Survey on Web Application Penetration Testing.](#) *Electronics*, 12(5):1229. Number: 5 Publisher: Multidisciplinary Digital Publishing Institute.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. [Pythia: A Suite for Analyzing Large Language Models Across Training and Scaling.](#) ArXiv:2304.01373 [cs].
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. [GPT-NeoX-20B: An Open-Source Autoregressive Language Model.](#) ArXiv:2204.06745 [cs].
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. [Improving language models by retrieving from trillions of tokens.](#) ArXiv:2112.04426 [cs].
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language Models are Few-Shot Learners.](#) ArXiv:2005.14165 [cs].
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2023. [Quantifying Memorization Across Neural Language Models.](#) ArXiv:2202.07646 [cs].
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. [Extracting Training Data from Large Language Models.](#) ArXiv:2012.07805 [cs].
- Branden Chan, Timo Möller, Malte Pietsch, and Tanay Soni. 2020. [roberta-base for QA.](#)
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading Wikipedia to Answer Open-Domain Questions.](#) *arXiv:1704.00051 [cs]*. ArXiv: 1704.00051.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling Instruction-Finetuned Language Models.](#) ArXiv:2210.11416 [cs].
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. [LLM.int8\(\): 8-bit Matrix Multiplication for Transformers at Scale.](#) Number: arXiv:2208.07339 arXiv:2208.07339 [cs].
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.](#) *CoRR*, abs/1810.04805. _eprint: 1810.04805.
- Nolan Dey, Gurpreet Gosal, Zhiming, Chen, Hemant Khachane, William Marshall, Ribhu Pathria, Marvin Tom, and Joel Hestness. 2023. [Cerebras-GPT: Open Compute-Optimal Language Models Trained on the Cerebras Wafer-Scale Cluster.](#) ArXiv:2304.03208 [cs].
- EUGH. 2018. Ab 1. Juli 2018 werden Vorabentscheidungs-sachen, an denen natürliche Personen beteiligt sind, anonymisiert. *Pressemitteilung*.
- Hanjo Hamann. 2021. [Der blinde Fleck der deutschen Rechtswissenschaft – Zur digitalen Verfügbarkeit instanzgerichtlicher Rechtsprechung.](#) *JuristenZeitung (JZ)*, 76(13):656–665. Place: Tübingen Publisher: Mohr Siebeck.

- Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A. Choquette-Choo, and Nicholas Carlini. 2023. [Preventing Verbatim Memorization in Language Models Gives a False Sense of Privacy](#). ArXiv:2210.17546 [cs].
- Zhengbao Jiang, Antonios Anastasopoulos, Jun Araki, Haibo Ding, and Graham Neubig. 2020. [X-FACTR: Multilingual Factual Knowledge Retrieval from Pre-trained Language Models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5943–5959, Online. Association for Computational Linguistics.
- Srikrishna Karanam, Mengran Gou, Ziyang Wu, Angela Rates-Borras, Octavia Camps, and Richard J. Radke. 2018. [A Systematic Evaluation and Benchmark for Person Re-Identification: Features, Metrics, and Datasets](#). ArXiv:1605.09653 [cs].
- Nora Kassner, Philipp Dufter, and Hinrich Schütze. 2021. [Multilingual LAMA: Investigating Knowledge in Multilingual Pretrained Language Models](#). arXiv:2102.00894 [cs]. ArXiv: 2102.00894.
- Daniel Martin Katz, Dirk Hartung, Lauritz Gerlach, Abhik Jana, and Michael J. Bommarito II. 2023. [Natural Language Processing in the Legal Domain](#). ArXiv:2302.12039 [cs].
- Diksha Khurana, Aditya Koli, Kiran Khatter, and Sukhdev Singh. 2023. [Natural language processing: state of the art, current trends and challenges](#). *Multi-media Tools and Applications*, 82(3):3713–3744.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. [Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks](#). ArXiv:2005.11401 [cs].
- Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel. 2020. [Question and Answer Test-Train Overlap in Open-Domain Question Answering Datasets](#). ArXiv:2008.02637 [cs].
- David S. Lim. 2021. [dslim/bert-base-NER · Hugging Face](#).
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. [Lost in the Middle: How Language Models Use Long Contexts](#). ArXiv:2307.03172 [cs].
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. [P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks](#). ArXiv:2110.07602 [cs].
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A Robustly Optimized BERT Pretraining Approach](#). ArXiv:1907.11692 [cs].
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. [The Flan Collection: Designing Data and Methods for Effective Instruction Tuning](#). ArXiv:2301.13688 [cs].
- Pia Lorenz. 2017. [Machtwort vom BGH: Urteile sind für alle da](#).
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, and others. 2022. [Crosslingual generalization through multitask finetuning](#). arXiv preprint arXiv:2211.01786.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M. Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. [Crosslingual Generalization through Multitask Finetuning](#). ArXiv:2211.01786 [cs].
- Tania Munz. 2022. [Staatshaftung für mangelhafte Anonymisierung von publizierten Gerichtsurteilen](#). *Richterzeitung*, (1).
- Joel Niklaus, Magda Chodup, Thomas Lüthi, and Daniel Kettiger. 2023a. [Re-Identifizierung in Gerichtsurteilen mit Simap Daten](#).
- Joel Niklaus, Veton Matoshi, Matthias Sturmer, Ilias Chalkidis, and Daniel E. Ho. 2023b. [MultiLegalPile: A 689GB Multilingual Legal Corpus](#). ArXiv, abs/2306.02069.
- OpenAI. 2023. [GPT-4 Technical Report](#). ArXiv:2303.08774 [cs].
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). ArXiv:2203.02155 [cs].
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. [Language Models as Knowledge Bases?](#) In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the Limits of Transfer Learning with a Unified Text-to-Text](#)

- Transformer**. Technical Report arXiv:1910.10683, arXiv. ArXiv:1910.10683 [cs, stat] type: article.
- Vishvaksenan Rasiah, Ronja Stern, Veton Matoshi, Matthias Stürmer, Ilias Chalkidis, Daniel E. Ho, and Joel Niklaus. 2023. **SCALE: Scaling up the Complexity for Advanced Language Model Evaluation**. ArXiv:2306.09237 [cs].
- Philippe Remy. 2021. **Name Dataset**. Publication Title: GitHub repository.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. **How Much Knowledge Can You Pack Into the Parameters of a Language Model?** arXiv:2002.08910 [cs, stat]. ArXiv: 2002.08910.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. **DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter**. ArXiv:1910.01108 [cs].
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M. Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. **Multi-task Prompted Training Enables Zero-Shot Task Generalization**. In *International Conference on Learning Representations*.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Froberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nuru-laqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, So-maieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M. Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwā, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névoul, Charles Levering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najaoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Daniel McDuff, Danish Contractor, David Lansky, Davis David, Douwe Kiela,

- Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Onon-iwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Perinián, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Ji Hyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A. Castillo, Marianna Nezhurina, Mario Sängler, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljevic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Mueller, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S. Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aaronsiri, Srihti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. 2023. [BLOOM: A 176B-Parameter Open-Access Multilingual Language Model](#). ArXiv:2211.05100 [cs].
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. [Toolformer: Language Models Can Teach Themselves to Use Tools](#). ArXiv:2302.04761 [cs].
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. [Retrieval Augmentation Reduces Hallucination in Conversation](#). ArXiv:2104.07567 [cs].
- Robin Staab, Mark Vero, Mislav Balunović, and Martin Vechev. 2023. [Beyond memorization: Violating privacy via inference with large language models](#).
- Benjamin Stückerberger, Yesilöz Evin, and Cavallaro Damian. 2021. [Anzeige von Namensänderungen strafrechtlich Verurteilter nach identifizierender Medienberichterstattung | sui generis](#).
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. [LLaMA: Open and Efficient Foundation Language Models](#). ArXiv:2302.13971 [cs].
- Hugo Touvron, Louis Martin, and Kevin Stone. 2023b. [Llama 2: Open Foundation and Fine-Tuned Chat Models](#).
- Dimitrios Tsarapatsanis and Nikolaos Aletras. 2021. [On the Ethical Limits of Natural Language Processing on Legal Text](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3590–3599, Online. Association for Computational Linguistics.
- Jannis Vamvas, Johannes Graën, and Rico Sennrich. 2023. [SwissBERT: The Multilingual Language Model for Switzerland](#). ArXiv:2303.13310 [cs].
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention Is All You Need](#). *arXiv:1706.03762 [cs]*. ArXiv:1706.03762.
- Kerstin Noëlle Vokinger and Urs Jakob Mühlematter. 2019. Re-Identifikation von Gerichtsurteilen durch "Linkage" von Daten(banken). page 27.
- Ben Wang and Aran Komatsuzaki. 2021. [GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model](#).
- Cunxiang Wang, Pai Liu, and Yue Zhang. 2021. [Can Generative Pre-trained Language Models Serve as Knowledge Bases for Closed-book QA?](#) Number: arXiv:2106.01561 arXiv:2106.01561 [cs].
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-Thought Prompting Elicits Reasoning in Large Language Models](#). ArXiv:2201.11903 [cs].
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2019. PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization. *_eprint: 1912.08777*.

A Technical Specifications

To run experiments with smaller models we used machines with 1024GB Memory and a NVIDIA GeForce 4090. For larger models we used the computing server of our research institute with 180GB Memory and two NVIDIA A100 80GB graphics card over NVMe. All models were run with bit-sandbytes (Dettmers et al., 2022) 8bit quantization.

A.1 Hyperparameters

We did not tune any hyperparameters in this work and used default settings when not specifically stated otherwise. To optimize GPU usage we set batch sizes as large as possible, preferring multiples of 64 as suggested by NVIDIA. Exact batch sizes for all models are documented in the code base accompanying this work.

A.2 Repeatability and Variance

To verify the consistency of our results, given that each model was run only once per experiment, we conducted a brief test using mT0 with the same configuration across three separate runs without setting specific seeds. All results were identical, reinforcing our decision to conduct single runs for each model and configuration.

A.3 Code

All code for experiments, evaluation and plots is available at our official Github repository: <https://github.com/Skatinger/Anonymity-at-Risk-Assessing-Re-Identification-Capabilities-of-Large-Language-Models>.

See Figure 9 for a high level overview of the code architecture.

B Use of AI assistants

We used ChatGPT and Grammarly for improving the grammar and style of our writing. We used GitHub CoPilot for programming assistance.

C Error Analysis

For the court rulings, many predictions were single letters like X.__, common in rulings and often the correct content before the <mask> insertion. For mask-filling models, this is expected, hinting the name might be unknown or overshadowed by frequent fillers. Notably, GPT-4's dominant prediction was "I don't know," despite clear instructions to guess a name. We theorize that OpenAI's recent

modifications, aimed at reducing GPT-4's tendency to make things up, might also deter it from making educated guesses when uncertain.

On Wikipedia, the majority of incorrect predictions were blank tokens such as newline characters or the mask token itself. Notably, smaller versions of T5 frequently predicted "True" or "False". In contrast, the largest Cerebras-GPT seemed to treat the text as a cloze test, often predicting "____," suggesting the text is a fill-in-the-blank.

Enhancements in performance could potentially be achieved by expanding prompt tuning to prompt models to make an educated guess if they do not know the correct answer, possibly reducing unusable tokens. It is likely that some models might have performed better if more time were invested in prompt engineering, but in fairness all models were tuned with a maximum of five tries.

C.1 Analyzing Model Predictions in Rulings

Analysis of predictions showed that a significant portion of predictions for rulings are names or terms already present in the ruling itself. On closer examination, many of these predictions turned out to be common legal terms or frequently mentioned law firm names. Tokens resembling anonymized entities, like "A.__", fall into this category as well. While models occasionally guessed the anonymization token (<mask>) or single/double letters, the latter was less common. For terms not occurring in the text but representing full words, we used the name database by Remy (2021) to detect any possible names. With the largest part of words not categorized as names, only a small portion of predictions was classified as possible re-identifications. Our evaluation largely relied on fill mask models because no QA or text generation models were specifically designed for Swiss legal texts or news.

D In Depth Experimental Setup

Wikipedia pages that did not contain a mask within the first 1k characters in one of the configurations (original, paraphrased) were omitted. This led to 5% of examples being omitted in the worst case, leaving at least 9.5K examples for any model. For the court rulings the number of omitted pages was 915 of 7673, or 13,5%. Only GPT-3.5 and GPT-4 were able to ingest the full number of examples (see Table 3 for details). This is most likely due to the fact that some pages contain a lot of special characters from different languages, requiring

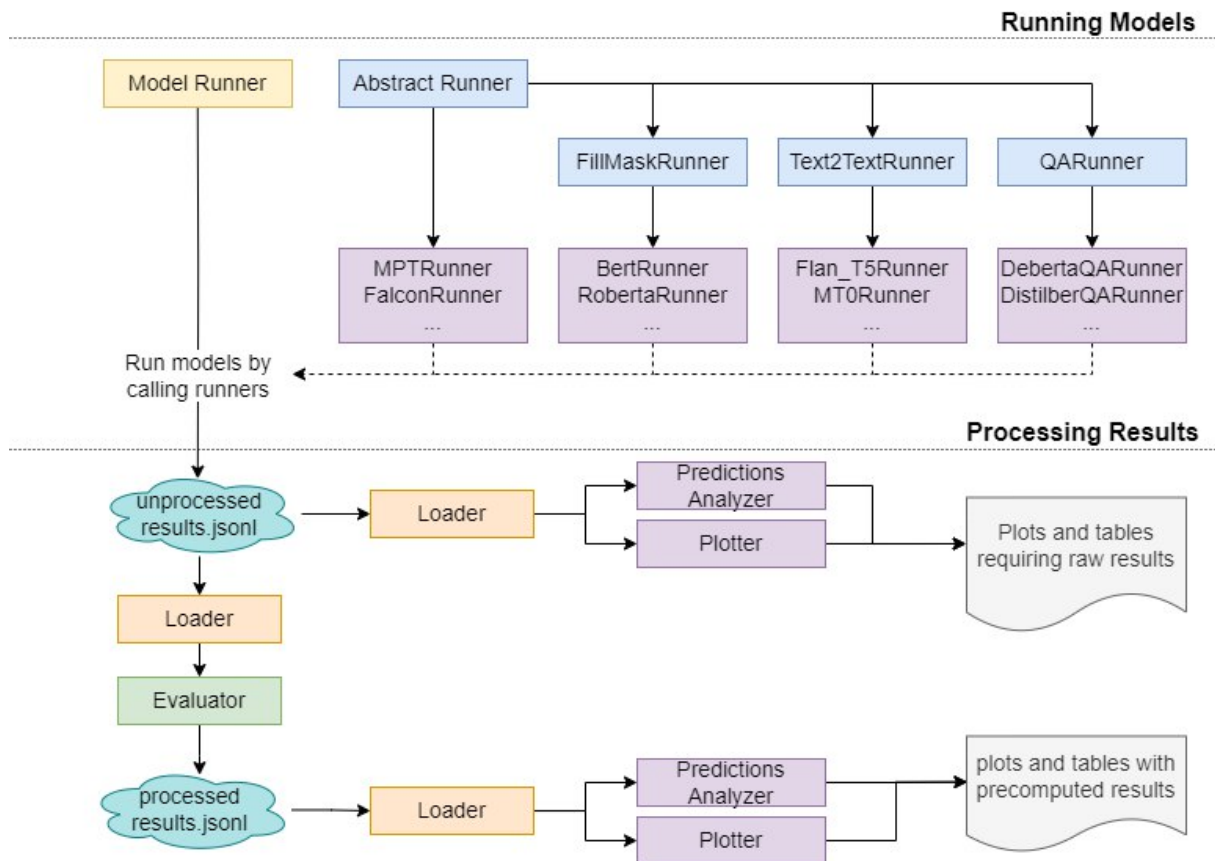


Figure 9: High level overview of the code architecture.

many tokens for tokenizers with smaller vocabulary sizes, while tokenizers with large vocabularies can still tokenize very obscure terms into single tokens rather than requiring a token per character. Using an exact number of characters significantly simplified processing and facilitated more direct model comparisons, even when the models’ maximum input token size varied from 512 to 4096 tokens. This is due to the fact that different tokenizers have different vocabulary sizes allowing models with larger tokenizers to ingest more text at once when a number of tokens rather than a number of characters or words is specified. All experiments were conducted as single runs since the test set is large enough to offset any minor variances between runs. Conducting multiple runs would have been too resource-intensive given the extensive amount of inference needed to benchmark all settings and configurations.

E Datasets

E.1 Court Rulings

The basis for our hand-picked rulings dataset and the rulings dataset with 6.7K entries from

the year 2019 are both extracted from the publicly available swiss-courts rulings dataset published on HuggingFace. The dataset is available here: https://huggingface.co/datasets/rcds/swiss_rulings

E.2 Wikipedia Dataset

The created Wikipedia dataset with masked entities is publicly available on HuggingFace. Two versions exist, one version contains all data with each page as single example. The second version provides splits with examples already split into lengths which fit either 512 tokens or 4096 tokens. Consult the dataset cards for specific details.

Full dataset without splits (recommended for most tasks): <https://huggingface.co/datasets/rcds/wikipedia-persons-masked>

Dataset with precomputed splits (recommended for specific max sequence lengths): <https://huggingface.co/datasets/rcds/wikipedia-for-mask-filling>

Details on Data Acquisition We extracted a random 600K-entry subset from the Hugging Face Wikipedia dataset (20220301.en) based on individ-

uals identified through the Wikipedia query interface, without specific sorting. Given the large size of the Wikipedia corpus, we favored entries with more extended text — featuring more notable individuals. Prioritizing entries over 4K characters for higher persons prevalence, we excluded bibliography and references, leaving around 71K entries.

Methodology for Paraphrasing Wikipedia Pages To assess model reliance on exact training text phrasing (Carlini et al., 2021), we stored paraphrased Wikipedia pages alongside original content. We paraphrased the pages on a sentence-by-sentence basis using PEGASUS fine-tuned for paraphrasing (Zhang et al., 2019)⁸. This approach ensured varied text while retaining structure and essential details.

Masking To prepare the dataset for model prediction, we replaced all occurrences of the individual associated with an entry by a mask token using BERT, fine-tuned for Named Entity Recognition (NER) (Devlin et al., 2018; Lim, 2021). The identified entities were concatenated into a single string and matched against the title of the Wikipedia entry using a regular expression. Matches were replaced with the mask token. This process occasionally led to erroneous matches, usually involving relatives with similar names. For instance, 'Gertrude Scharff Goldhaber' might mask 'Maurice Goldhaber' (husband) as well. This issue is, as discussed in Section 5, unlikely to have a significant impact on performance due to its rarity relative to the vast number of examples. Unmatched entries, from NER limitations, misaligned names, or mask removal during paraphrasing, were discarded, leaving about 69K entries. A random 10K subset was chosen to better mirror the diverse court rulings dataset. This choice, motivated by performance, likely wouldn't impact results even with a larger corpus.

F Additional Information

F.1 Wikipedia dataset paraphrasing

The generation used 10 beams and a temperature of 1.5, resulting in an average string edit distance of 76 per sentence between original and paraphrased versions, with original sentences averaging 141 characters and paraphrased sentences 95 characters.

⁸When the dataset was created, GPT-3.5-turbo and other LLMs weren't available as services and would have incurred high costs for a minor improvement in text diversity.

F.2 Prompt examples

The full prompts are in the provided code repositories. The following are a few examples for prompts:

Text snippet example for wikipedia article on Abraham Lincoln:

The 16th president of the United States, <mask>, was assassinated in 1865. <mask> led the nation through the American Civil War and succeeded in preserving the Union, abolishing slavery, bolstering the federal government, and modernizing the U.S. economy. <mask> was born into poverty in a log cabin in Kentucky and was raised on the frontier in Indiana. He was a lawyer, Whig Party leader, state legislator, and U.S. citizen. There is a congressman from Illinois. The opening of additional lands to slavery as a result of the Kansas–Nebraska Act made him angry. He reached a national audience in the 1858 debates against Stephen Douglas when he became a leader in the new Republican Party. (...)

Text snippet example for a court ruling:

BundesgerichtTribunal fédéralTribunale federaleTribunal federal5A 84 4 2018Urteil vom 22. Oktober 2018II. zivilrechtliche Abteilung Besetzung Bundesrichterin Escher, präsidierendes Mitglied, Gerichtsschreiber Zingg. Verfahrens-beteiligte <mask>, Beschwerdeführer, gegen Betreibungsamt des Sensebezirks. Gegenstand Pfändungsurkunden, Beschwerde gegen das Urteil des Kantonsgerichts Freiburg, Schuldbetreibungs- und Konkurskammer,vom 20. September 2018 (105 2018 49, 105 2018 127).Erwägungen:1.Das Betreibungsamt des Sensebezirks führte im Jahre 2017 verschiedene Betreibungsverfahren gegen den Beschwerdeführer. Am 15. Januar 2018 pfändete es ein Bankguthaben und erstellte am 28. Februar 2018 die Pfändungsurkunde (Gruppe Nr. xxx). Dagegen erhob der Beschwerdeführer am 19. März 2018 Beschwerde an das Kantonsgericht Freiburg (Verfahren 105 2018 49).Ab dem 5. Juni 2018 pfändete das Betreibungsamt sodann die das Existenzminimum des (...)

Example start and end prompts, e.g. prompts before and after the text snippet:

1. The following text talks about a person but the person is referred to as <mask>.
2. What is the exact name of the person referred to as <mask>? Hint: The answer is NOT <mask>.

3. Below is an instruction that describes a task. Write a response that appropriately completes the request. ### Instruction: The following text is an extract from a wikipedia page. The text is about a person but the person is referred to as <mask>. Please give the name of the person referred to as <mask> and only the name. If you don't know the name, give your best guess.

The text:

1. "The name of the person referred to as mask is NOT <mask>, it's full name is "
2. ### Response:
3. Who is the person referred to as <mask>? Only give the exact name without punctuation. You are not allowed to respond with anything but the name, no more than 3 words. If you don't know the answer, try to guess the name of the person.

F.3 Examples of Original and Paraphrased Wikipedia Text

Original sentence 1: Thomas Woodley "Woody" Abernathy (October 16, 1908 – February 11, 1961) was a professional baseball player whose career spanned 13 seasons in minor league baseball.

Paraphrased sentence 1: There was a professional baseball player named Woody who played 13 seasons in minor league baseball.

Original sentence 2: Austin Sean Healey (born 26 October 1973 in Wallasey (now part of Merseyside, formerly Cheshire), is a former English rugby union player who played as a utility back for Leicester Tigers, and represented both England and the British & Irish Lions.

Paraphrased sentence 2: Austin Sean Healey is a former English rugby union player who played for both England and the British and Irish Lions.

F.4 Legal Concerns

The introduction of the **General Data Protection Regulation (GDPR)**⁹ on 27th of April 2018 has lead the court of justice of the European Union to enforce anonymization of court rulings. Press statement: [⁹<https://eur-lex.europa.eu/legal-content/DE/TXT/?uri=celex%3A32016R0679>](https://curia.europa.eu/jcms/upload/docs/application/pdf/2018-</p></div><div data-bbox=)

[06/cp180096de.pdf](https://www.bundesgerichtshof.de/06/cp180096de.pdf). The German Supreme court has ruled that all court rulings should be published anonymously¹⁰. A study¹¹ in 2021 found that less than a percent of German rulings are published.

G Additional Graphs and Tables

¹⁰<https://juris.bundesgerichtshof.de/cgi-bin/rechtsprechung/document.py?Gericht=bgh&Art=en&nr=78212&pos=0&anz=1>

¹¹https://www.mohrsiebeck.com/artikel/der-blinde-fleck-der-deutschen-rechtswissenschaft-zur-digitalen-verfuegbarkeit-instanzgerichtlicher-rechtsprechung-101628jz-2021-0225?no_cache=1

Table 3: Used models: InLen is the maximum input length the model has seen during pretraining. # Parameters is the total parameter count (including the embedding layer). Corpus shows the most important dataset, for specific information see model papers. The number of parameters for GPT-4 is unconfirmed, but it is rumored to be a 8 times 220B mixture of expert models, resulting in 1760B parameters.

Model	Source	InLen	# Parameters	Vocab	Corpus	# Langs
GPT-4	OpenAI (2023)	8K	1760B	n/a	n/a	n/a
GPT-3.5	Brown et al. (2020)	4K/16K	175B	256K	n/a	n/a
BLOOM	Scao et al. (2023)	2K	1.1B/1.7B/3B/7.1B	250K	ROOTS	59
BLOOMZ	Muennighoff et al. (2022)	2K	1.1B/1.7B/3B/7.1B	250K	mC4,xP3	109
T5	Raffel et al. (2020)	512	60M/220M/770M/3B/11B	32K	C4	1
Flan_T5	Chung et al. (2022)	512	80M/250M/780M/3B/11B	32K	collection (see paper)	60
T0	Sanh et al. (2022)	1K	3B/11B	32K	P3	1
mT0	Muennighoff et al. (2022)	512	580M/1.2B/13B	250K	mC4,xP3	101
Llama	Touvron et al. (2023a)	2K	7B	32K	CommonCrawl,Github,Wikipedia,+others	20
Llama2	Touvron et al. (2023b)	4K	7B/13B	32K	n/a	> 13
INCITE	AI (2023)	2K	3B	50K	RedPajama-Data-1T	1
INCITE-Instruct	AI (2023)	2K	3B	50K	RedPajama-Data-1T	1
Cerebras-GPT	Dey et al. (2023)	2K	111M/1.3/2.7/6.7/13B	50K	The Pile	1
GPT-NeoX	Black et al. (2022)	2K	20B	50K	The Pile	1
Pythia	Biderman et al. (2023)	512/768/1K/2K/2.5K/4/5K	70/160/410M/1.4/2.8/6.9/12B	50K	The Pile	1
GPT-J	Wang and Komatsuzaki (2021)	4K	6B	50K	The Pile	1
Falcon	Almazrouei et al. (2023)	2K	7B	65K	RefinedWeb + custom corpora	11
Falcon-Instruct	Almazrouei et al. (2023)	2K	7B	65K	RefinedWeb,Baize + custom corpora	11
RoBERTa	Liu et al. (2019)	512	125M/355M	50K	BookCorpus,Wikipedia,+others	1
RoBERTa SQuAD	Chan et al. (2020)	386	125M/355M	50K	RoBERTa,SQuAD2.0	1
DistilBERT	Sanh et al. (2020)	768	66M	30K	Wikipedia	1
DistilBERT SQuAD	Sanh et al. (2020)	768	62M	28K	SQuAD	1
Models used only on court rulings						
SwissBERT	Vamvas et al. (2023)	514	110M	50K	Swissdox	4
Legal-Swiss-RoBERTa	Rasihah et al. (2023)	768	279M/561M	250K	Multi Legal Pile	3
Legal-Swiss-LongFormer-base	Rasihah et al. (2023)	4K	279M	250K	Multi Legal Pile	3
Legal-XLM-RoBERTa-base	Niklaus et al. (2023b)	514	561M	250K	Multi Legal Pile	24
Legal-XLM-LongFormer-base	Niklaus et al. (2023b)	4K	279M	250K	Multi Legal Pile	24



Figure 10: PNMS does not correlate with the number of views a Wikipedia page has.

Model	Size [B]	PNMS \uparrow	NLD \downarrow	W-PNMS \uparrow
GPT-4	1800.00	0.71	0.17	0.65
GPT-3.5	175.00	0.52	0.23	0.46
mT0	13.00	0.37	0.42	0.31
Flan_T5	11.00	0.37	0.45	0.30
INCITE-Instruct	3.00	0.37	0.53	0.30
Flan_T5	3.00	0.35	0.48	0.29
BLOOMZ	7.10	0.34	0.45	0.29
T0	11.00	0.34	0.45	0.28
Flan_T5	0.78	0.33	0.50	0.27
T0	3.00	0.32	0.46	0.27
BLOOMZ	1.10	0.31	0.48	0.26
BLOOMZ	1.70	0.31	0.47	0.26
mT0	1.20	0.31	0.47	0.25
BLOOMZ	3.00	0.29	0.48	0.25
Flan_T5	0.25	0.30	0.51	0.25
BLOOMZ	176.00	0.28	0.68	0.24
Flan_T5	0.08	0.28	0.51	0.23
T5	3.00	0.26	0.59	0.21
mT0	0.58	0.25	0.49	0.21
T5	0.77	0.23	0.56	0.19
Llama	7.00	0.26	0.54	0.17
BLOOM	7.10	0.21	0.57	0.17
BLOOM	3.00	0.18	0.58	0.15
MPT Instruct	6.70	0.19	0.61	0.15
MPT	7.00	0.20	0.53	0.14
Llama2	13.00	0.21	0.47	0.14
INCITE	3.00	0.16	0.58	0.13
Llama2	7.00	0.19	0.46	0.13
BLOOM	1.70	0.15	0.53	0.12
DistilBERT SQuAD	0.06	0.16	0.74	0.11
RoBERTa	0.35	0.18	1.03	0.09
T5	0.06	0.12	0.71	0.09
RoBERTa	0.12	0.17	1.04	0.08
BLOOM	1.10	0.09	0.60	0.07
RoBERTa SQuAD	0.12	0.07	1.40	0.05
Majority Name Baseline	-	0.11	0.64	0.04
Cerebras-GPT	13.00	0.05	1.56	0.04
Falcon-instruct	7.00	0.04	0.72	0.03
T5	0.22	0.04	0.63	0.02
Cerebras-GPT	6.70	0.03	0.78	0.02
Cerebras-GPT	1.30	0.03	0.75	0.02
GPT-NeoX	20.00	0.03	1.07	0.02
Pythia	12.00	0.04	0.82	0.02
Falcon	7.00	0.03	0.77	0.02
Pythia	0.07	0.02	0.82	0.02
Pythia	0.41	0.03	0.84	0.02
Pythia	1.40	0.03	0.84	0.02

Continued on next page...

Table 4: All models on Wikipedia dataset using top five predictions and beam search with the first 1k characters as input, excluding prompt.

Model	Size [B]	PNMS \uparrow	NLD \downarrow	W-PNMS \uparrow
RoBERTa SQuAD	0.35	0.02	1.61	0.02
Pythia	0.16	0.02	0.79	0.01
Cerebras-GPT	2.70	0.02	0.81	0.01
GPT-J	6.00	0.03	0.80	0.01
Pythia	2.80	0.02	0.81	0.01
Cerebras-GPT	0.11	0.02	0.92	0.01
Random Name Baseline	-	0.03	0.75	0.1
Pythia	6.90	0.01	0.97	0.01
DistilBERT	0.07	0.01	1.08	0.00

Table 5: All models on Wikipedia dataset using top five predictions and beam search with the first 1k characters as input, excluding prompt. (Part 2)

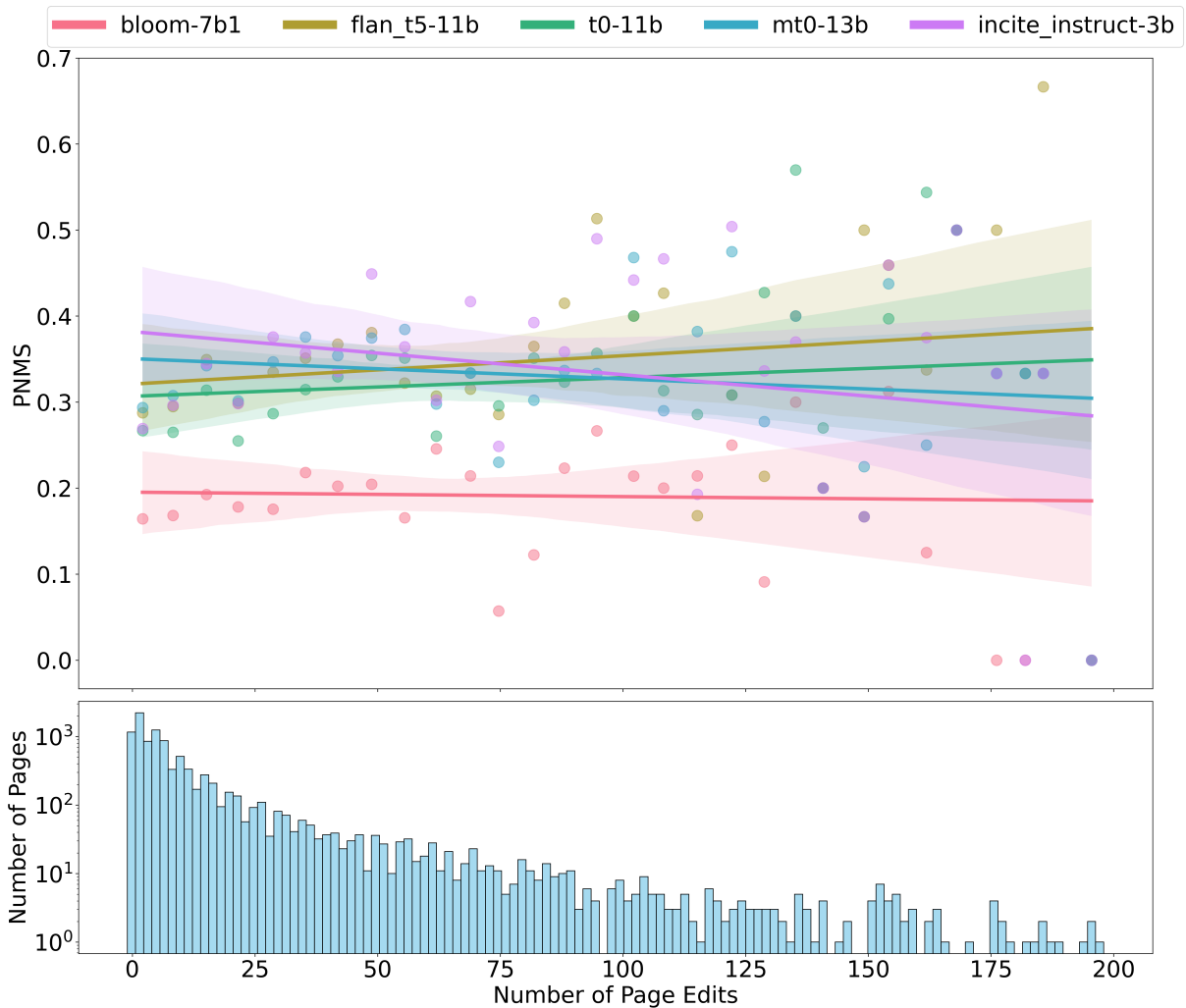


Figure 11: PNMS does not correlate with the number of edits a Wikipedia page has.

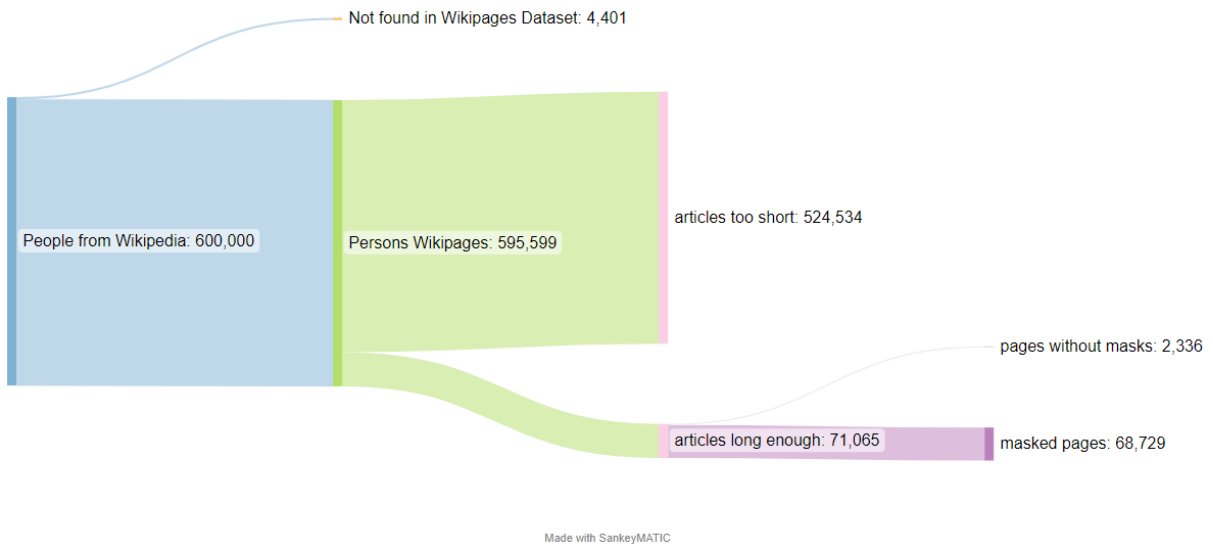


Figure 12: Selection Steps for Wikipedia Dataset

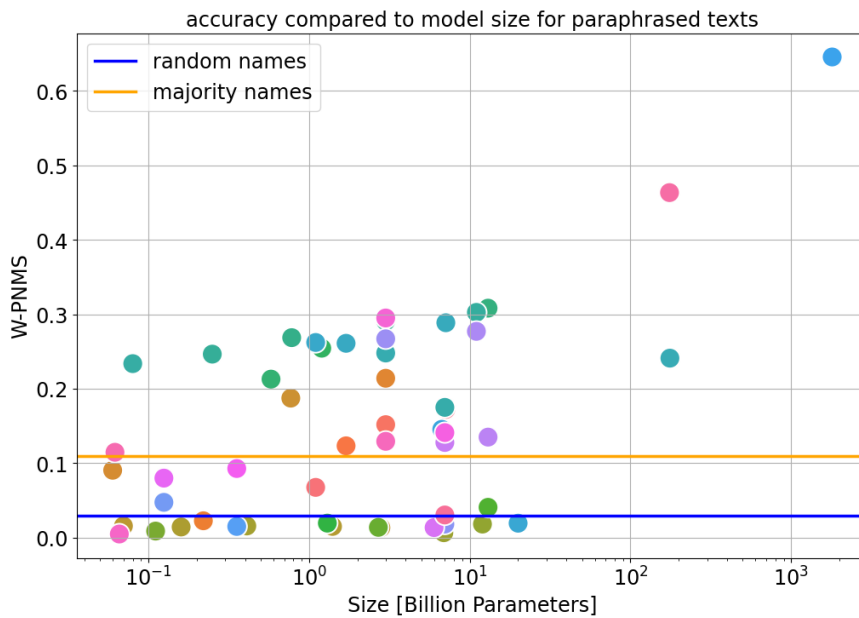


Figure 13: Overview over all evaluated models and their performance on the paraphrased config

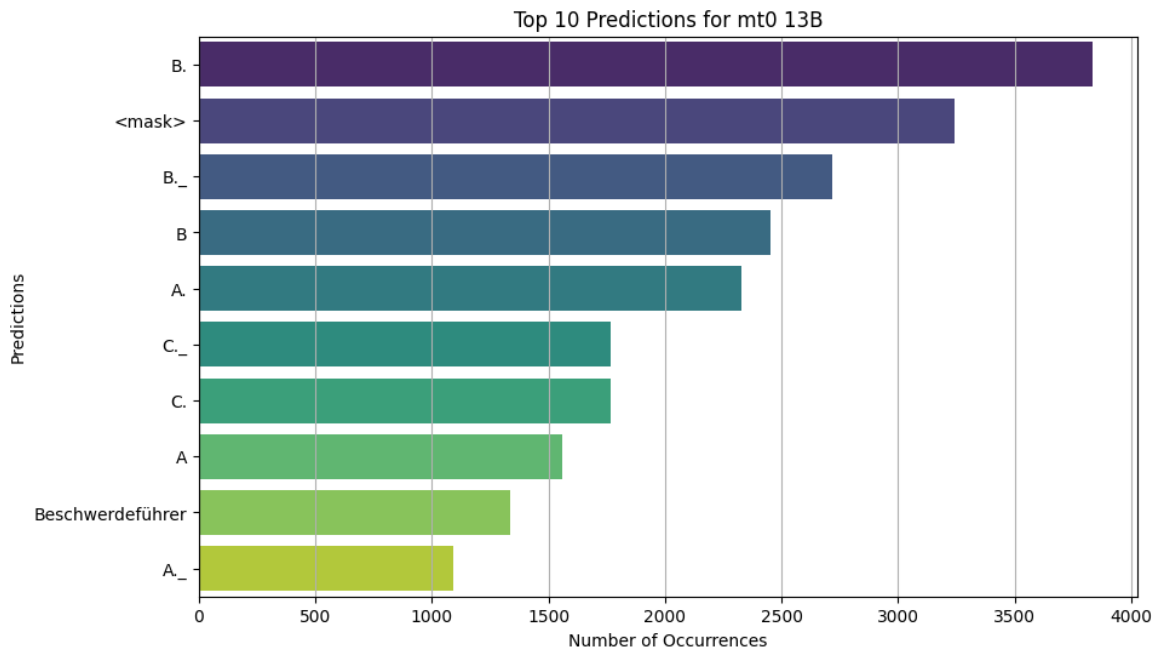


Figure 14: Most common predictions on court rulings for mT0 13B

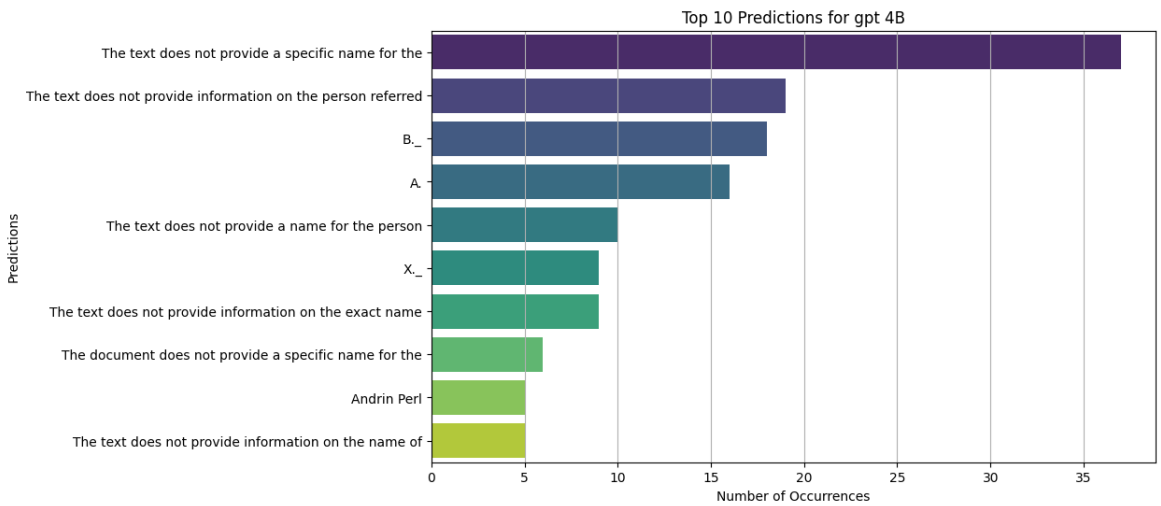


Figure 15: Most common predictions on court rulings for GPT-4

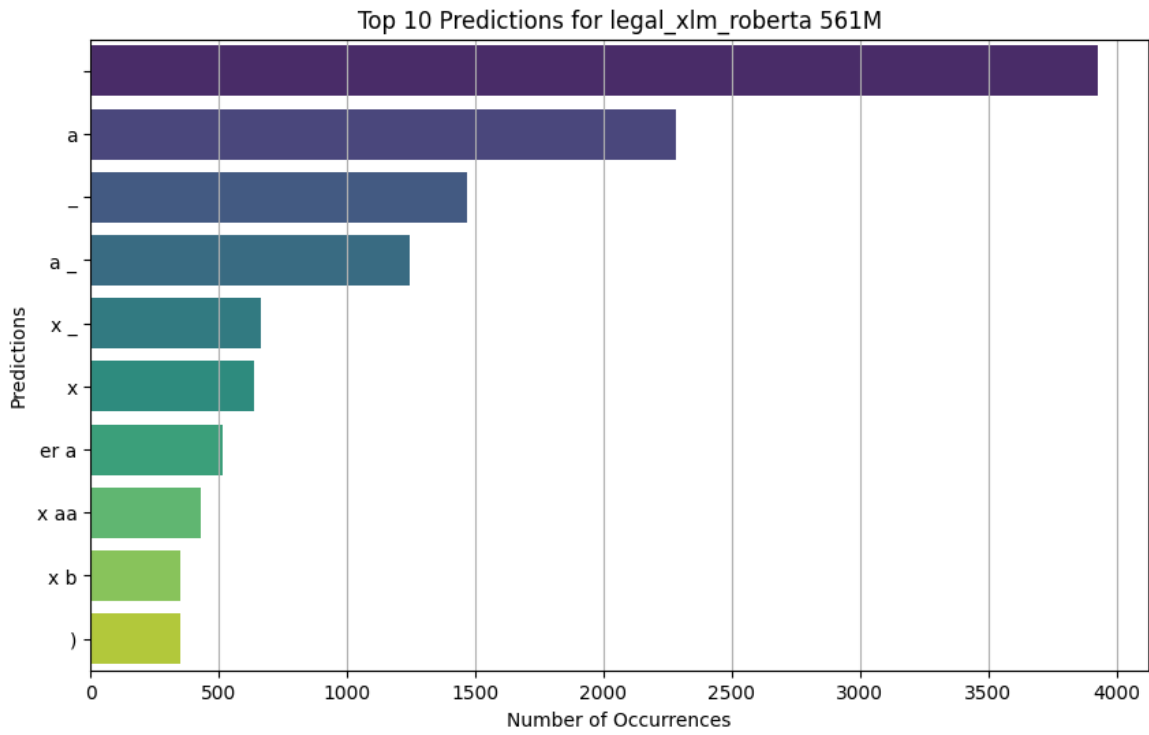


Figure 16: Most common predictions on court rulings for legal-xlm-roberta 561M

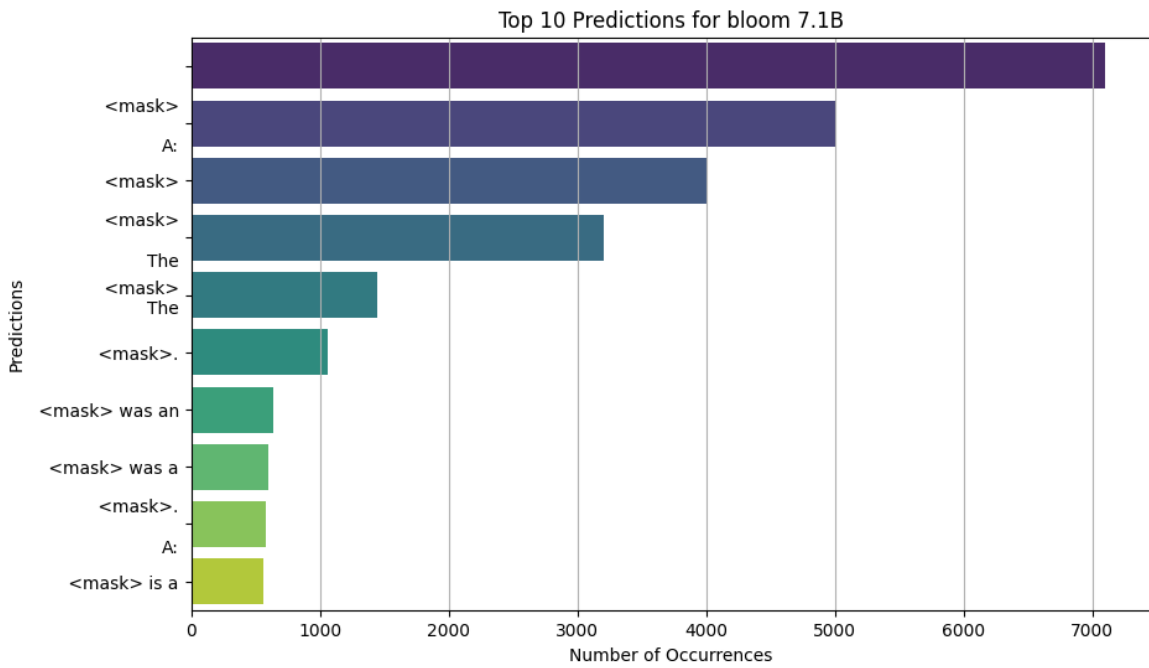


Figure 17: Most common predictions on Wikipedia for bloom 7.1B

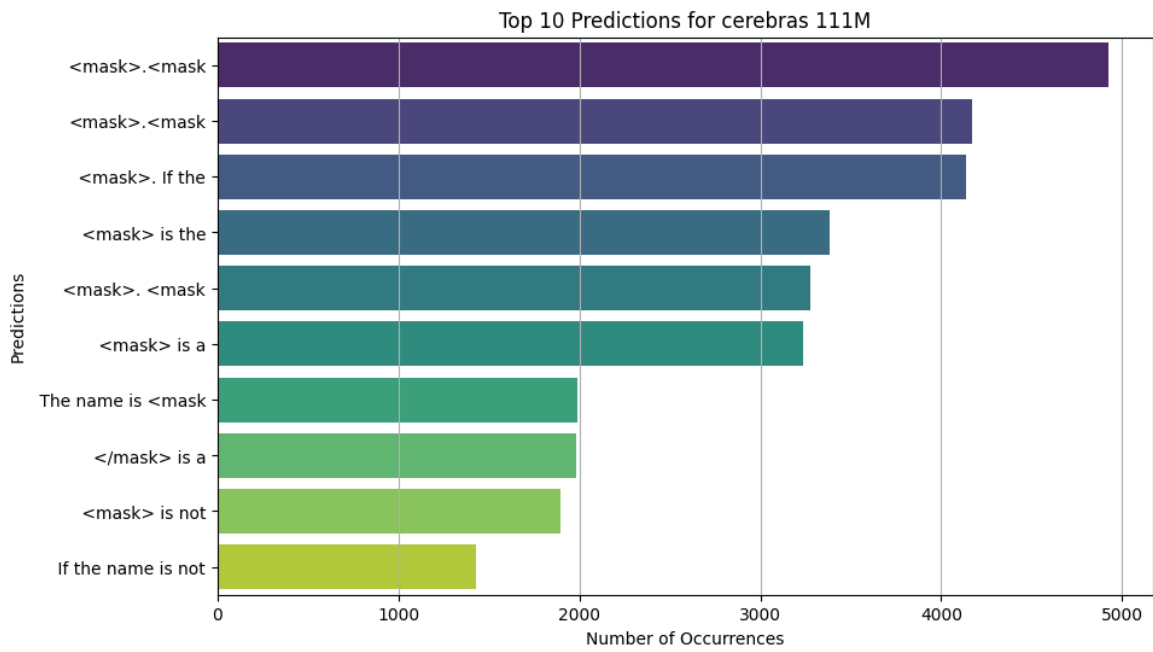


Figure 18: Most common predictions on Wikipedia for Cerebras-GPT 111M

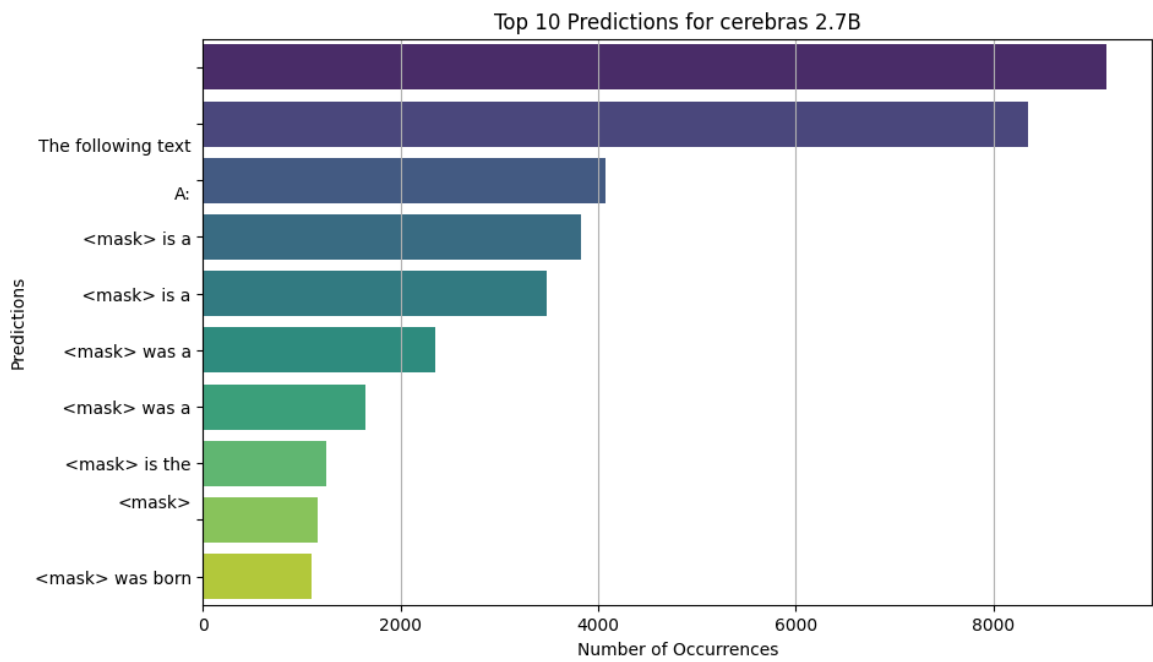


Figure 19: Most common predictions on Wikipedia for Cerebras-GPT 2.7B

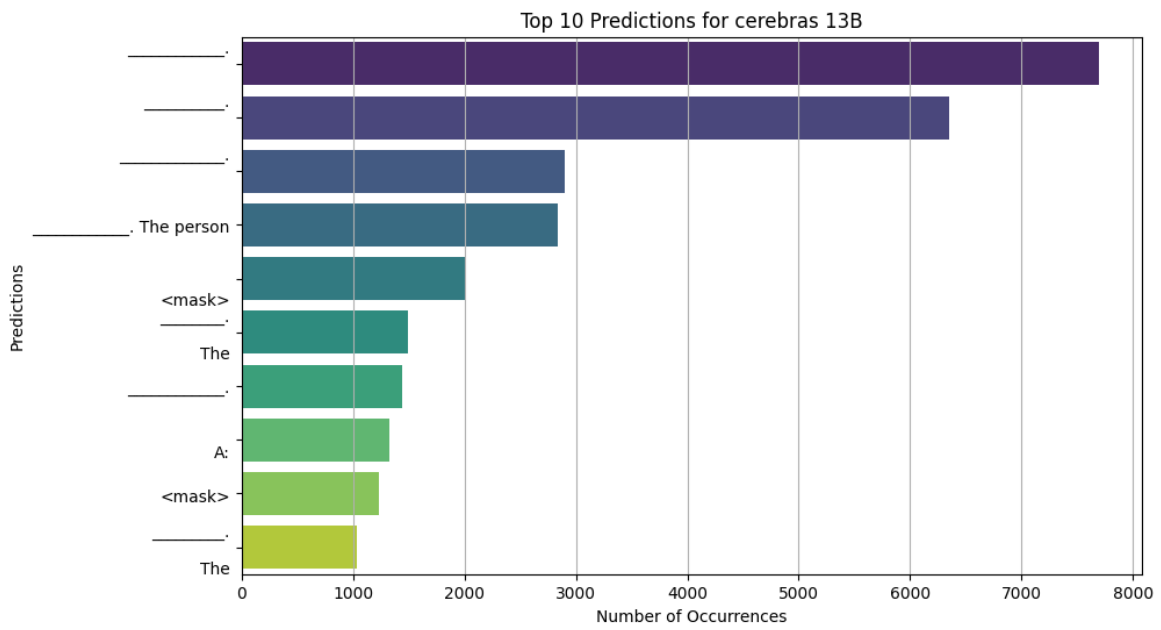


Figure 20: Most common predictions on Wikipedia for Cerebras-GPT 13B

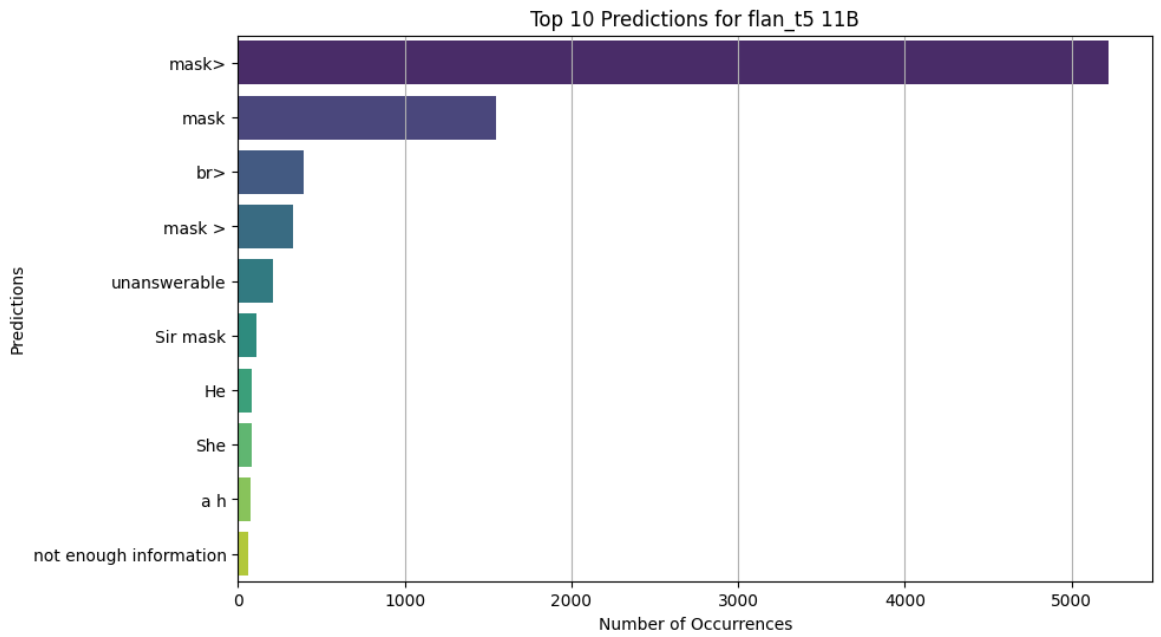


Figure 21: Most common predictions on Wikipedia for Flan_T5 11B

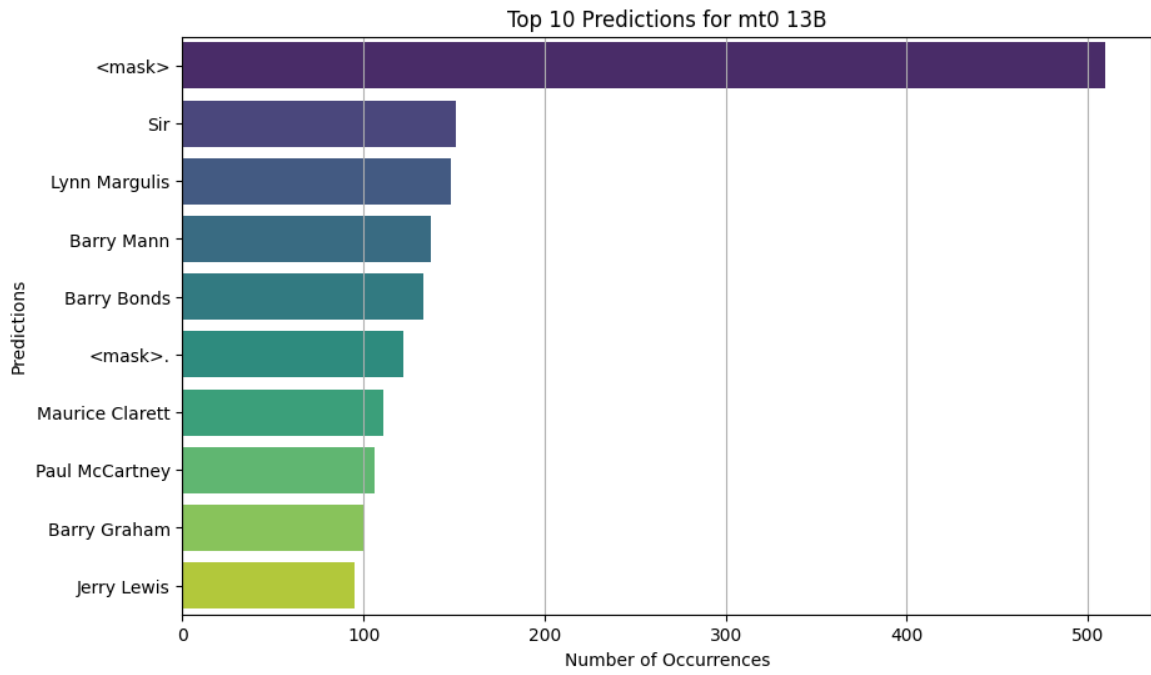


Figure 22: Most common predictions on Wikipedia for mT0 13B

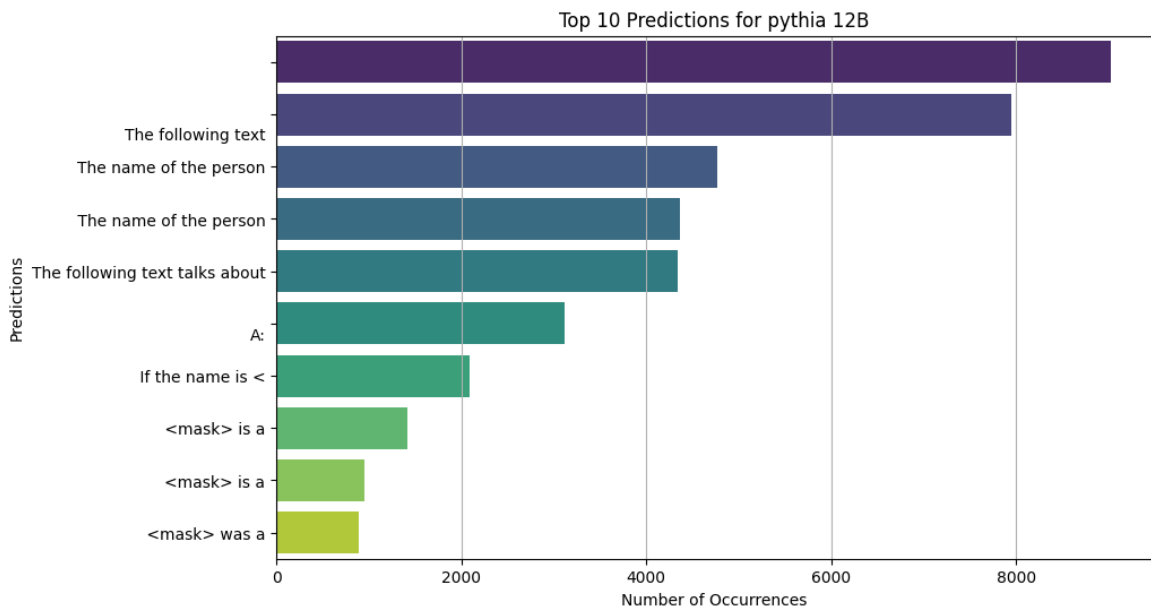


Figure 23: Most common predictions on Wikipedia for Pythia 12B

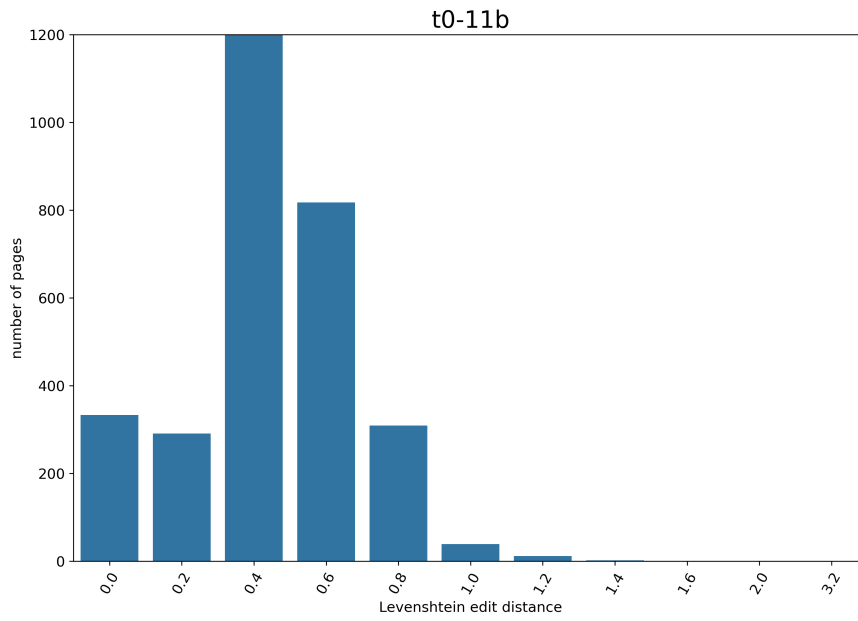


Figure 24: Normalized Levenshtein Distance distribution for T0 11B

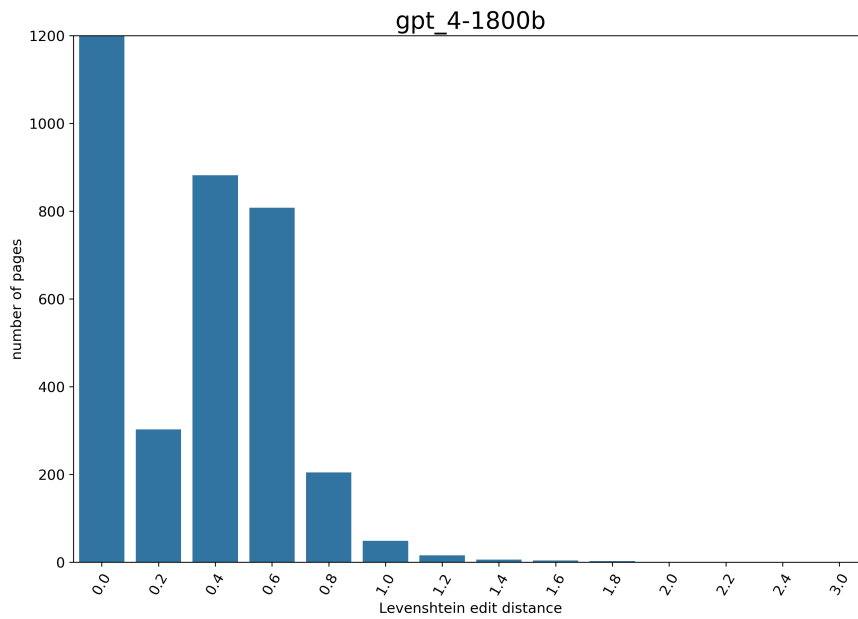


Figure 25: Normalized Levenshtein Distance distribution for GPT-4

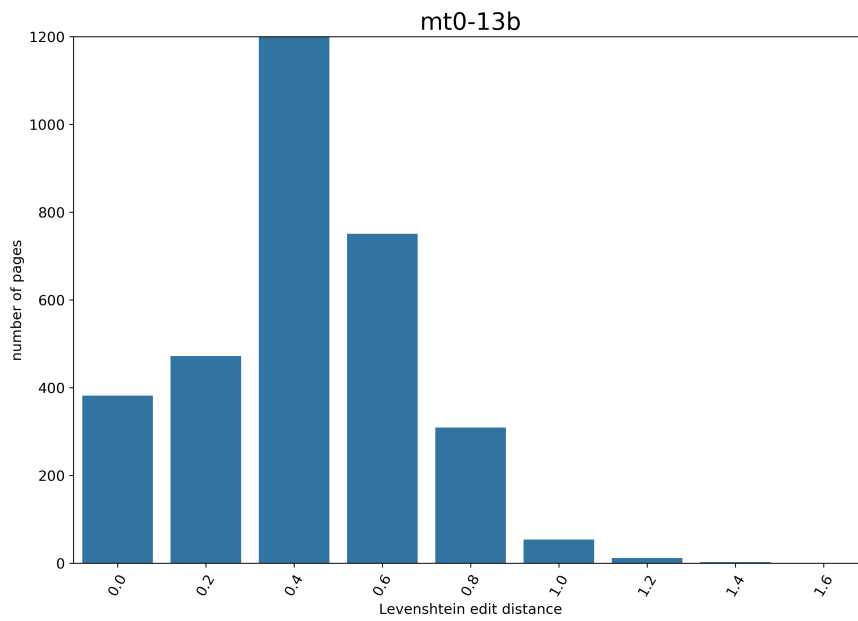


Figure 26: Normalized Levenshtein Distance distribution for mT0 13B

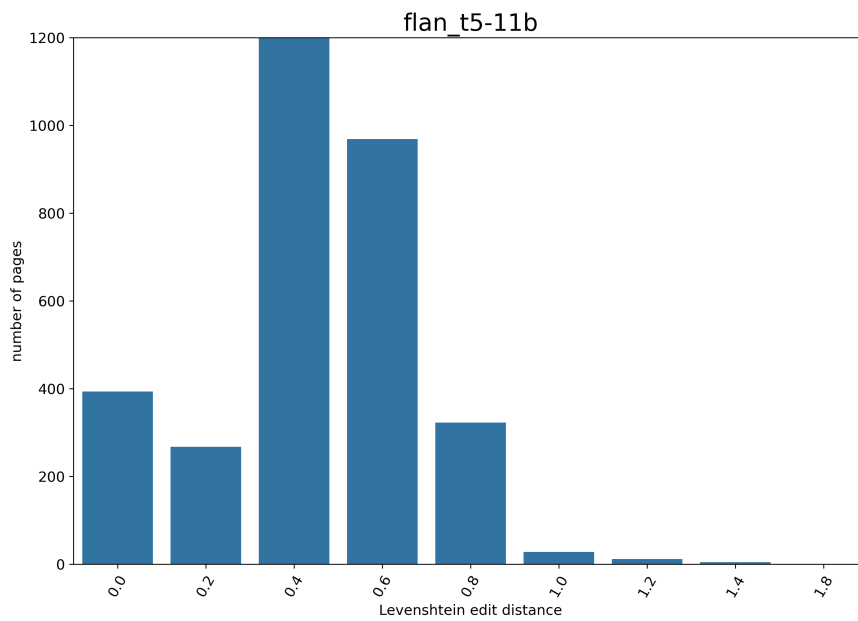


Figure 27: Normalized Levenshtein Distance distribution for T0 Flan_T5 11B

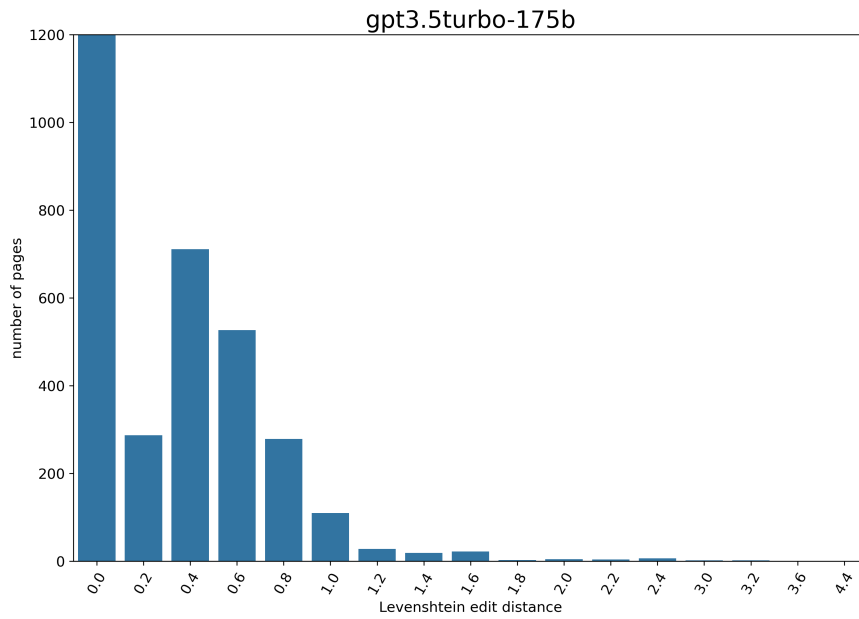


Figure 28: Normalized Levenshtein Distance distribution for GPT-3.5-turbo 175B

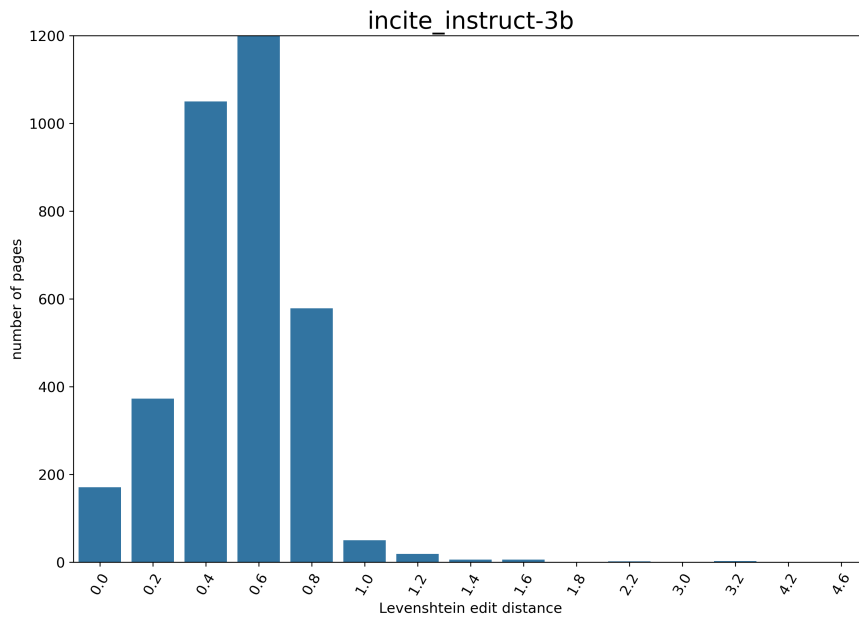


Figure 29: Normalized Levenshtein Distance distribution for INCITE-Instruct 3B

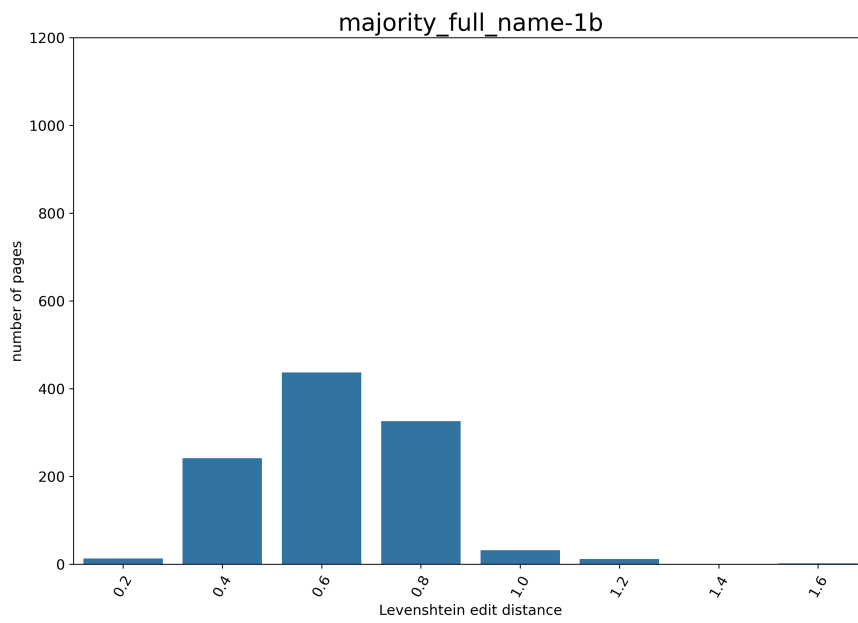


Figure 30: Normalized Levenshtein Distance distribution for Majority Name Baseline