# The Empirical Impact of Data Sanitization on Language Models

**Anwesan Pal**, **Radhika Bhargava**, Kyle Hinsz, Jacques Esterhuizen, Sudipta Bhattacharya

aws

NEURAL INFORMATION PROCESSING SYSTEMS

## Problem statement

**Data privacy** is a critical concern in the development and use of language models (LMs) due to presence of personally identifiable information (PII).

One approach to ensure data privacy is **Data Sanitization** which involves complete and irreversible removal of PII from data. Despite the wide-adoption of data sanitization methods, its impact on the performance of language models has not been studied in-depth.



## Methodology

To understand the impact of data sanitization on LM performance, we perform experiments with both small and large language models, and across a variety of natural language processing (NLP) and GenAI datasets.

**Datasets:**
1. **Traditional NLP Datasets:** We performed analysis on the following traditional NLP datasets: QQP, MultiNLI , Winograd Schema Challenge, LEDGAR dataset, EURLEX dataset, SQuADv2.0 and IMDB dataset.
2. **GenAI Datasets:** We included the following datasets used to benchmark modern large language model (LLM) performances: DROP, GSM8K, and a set of tasks from Big-Bench-Hard (BBH) benchmark.

**Models:**
1. **Small language models <5B:** We study the effects of data sanitization on BART (encoder-decoder) and GPT-2 (decoder only) models.
2. **Large language models >5B:** We used chain-of-thought (CoT) prompting with few-shot examples to study the effects of data sanitization on the following models: Anthropic's Claude 3.5 Sonnet, Mistral AI's Mistral 7B and OpenAI's GPT-4o.

## Results

### NLP Datasets

| Datasets | BART | | GPT-2 | |
|---|---|---|---|---|
| | None/ None | Redact / Redact | None/ None | Redact / Redact |
| *Low Impact (<10%)* | | | | |
| IMDB (Acc) | 93.7 | 93.7 | 93.1 | 93.2 |
| LexGLUE: EURLEX (F1) | 66.3 | 66.3 | 64.1 | 62.1 |
| GLUE: QQP (Acc) | 90.4 | 88.5 | 89.0 | 86.9 |
| *Moderate Impact (10-25%)* | | | | |
| SQuAD v2.0 (F1) | 74.9 | 55.7 | 55.8 | 48.7 |

Performance results on NLP datasets: For each dataset, the model performances are shown for different combinations of original and redacted versions across training and validation splits. The results suggest only minimal degradation in model performance when training on redacted data, with performance decreasing <2.2% on the average.

### Gen AI Datasets

| Datasets | Claude 3.5 Sonnet | | Mistral 7B | | GPT-4o | |
|---|---|---|---|---|---|---|
| | None | Redact | None | Redact | None | Redact |
| *Low Impact (<10%)* | | | | | | |
| IMDB | 95.8 | 95.5 | 86.5 | 86.6 | 93.9 | 93.1 |
| BBH: Causal Judgement | 69.0 | 63.0 | 42.8 | 42.2 | 67.0 | 65.0 |
| BBH: Formal Fallacies | 88.0 | 75.0 | 60.0 | 57.2 | 78.0 | 74.0 |
| *Moderate Impact (10-25%)* | | | | | | |
| SQuADv2.0 | 65.8 | 57.8 | 46.1 | 30.5 | 68.3 | 51.4 |
| BBH: Logical Deduction (#5) | 93.6 | 82.7 | 24.4 | 26.0 | 91.6 | 80.0 |
| BBH: Logical Deduction (#7) | 83.5 | 64.7 | 22.8 | 18.4 | 79.6 | 66.8 |
| *High Impact (>25%)* | | | | | | |
| DROP | 92.1 | 54.2 | 46,1 | 25,9 | 91.6 | 49.3 |
| GSM8K | 96.9 | 44.6 | 45.3 | 19.0 | 57.6 | 25.5 |
| BBH: Penguins in a Table | 99.3 | 30.8 | 43.8 | 29.0 | 99.0 | 47.0 |

For GenAI datasets the impact of redaction on the different tasks range from 0.3% to 95% for Claude, -2.7% to 67.3% for Mistral and -6.5% to 100% for GPT. Based on these results, we have classified the datasets as low impact if the impact on performance was < 10%, medium impact if the impact on performance was between 10 and 25% and high impact for those datasets where the impact was greater than 25%

## Observations

1. **Oddities in Mistral's performance on Redacted Datasets:** Mistral has a tendency to hallucinate and assign placeholder values for redacted entities, and reason about them incorrectly to obtain the correct answer.



2. **Weaker Redaction for High Impact Datasets:** With limited redaction by skipping task-critical entities, many of the previous high-impact datasets now have a low impact. The exception being DROP, which is still moderate impact. We hypothesize the presence of multiple dominant entities being present in that dataset to be the cause for this.

| Datasets | Redaction Amount | | | Redacted PII Entities |
|---|---|---|---|---|
| | None | Full | Limited | |
| DROP | 92.1 | 54.2 | 79.3 | NAME, LOC, ORG |
| GSM8K | 96.9 | 44.6 | 90.1 | NAME, LOC, ORG |
| BBH: Date Understanding | 92.8 | 40.6 | 86.3 | NAME |

## Redacted dataset repair strategy



In many real-world applications involving GenAI algorithms, developers often do not have control over the degree of redaction within the dataset, and have to make the best possible use of it in its redacted state. One such strategy involves subsampling a given redacted dataset by removing high PII-content records, and using the remaining ones.