

# Replication and Extension of DP-MLM: Differentially Private Text Rewriting Using Masked Language Models

Simran Dawadi

`simran.dawadi@students.mq.edu.au`

November 2025

## Abstract

This report documents the replication and extension of DP-MLM (Differentially Private Masked Language Model) proposed by Meisner et al. (2022). The project aims to reproduce the functionality of the model, confirm its claims on the original dataset, and extend its evaluation to new domains including IMDB reviews, Twitter posts, and a newly constructed synthetic dataset that mimics real-world sensitive text. The replicated system was implemented using the RoBERTa-base transformer and verified to successfully anonymise personal details such as names, phone numbers, and emails while maintaining the semantic integrity of the original sentences. Extended evaluations across new datasets further confirmed DP-MLM’s generalisation ability, demonstrating that differential privacy mechanisms can provide meaningful protection while preserving readability. This work reinforces the principles of reproducible research and open science by validating an existing public resource and critically analysing its limitations, ethical considerations, and practical challenges.

## 1 Introduction

Reproducibility is one of the cornerstones of scientific progress. In machine learning and data science, replicating published studies allows researchers to confirm reported claims, test generalisation, and evaluate robustness. The goal of this project is to replicate and extend the Differentially Private Masked Language Model (DP-MLM), an approach that introduces differential privacy mechanisms into text rewriting using masked language models (MLMs).

Differential privacy (DP) provides mathematical guarantees that individual data points cannot be reverse-engineered from a model’s output (Dwork et al., 2006). When applied

to language models, it can help prevent the unintentional leakage of personal identifiers such as names, IDs, or addresses. DP-MLM builds on the idea that masked language models like BERT and RoBERTa can be used not only for prediction but also for rewriting sensitive text under formal privacy controls.

The significance of this work lies in its dual purpose: validating reproducible AI systems and ensuring ethical, privacy-preserving NLP. By confirming and extending DP-MLM, this report contributes to transparent, accountable research practices and provides insights into the practical challenges of implementing privacy-preserving text models in realistic settings.

## 2 Source Paper and Evaluation Framework

The source paper Meisner et al. (2022), published in the Findings of ACL 2022 (an A-ranked CORE venue), introduced a framework that integrates differential privacy into masked language model prediction.<sup>1</sup> Instead of deleting or replacing tokens arbitrarily, the system rewrites sensitive spans using context-aware substitutions, ensuring both privacy and linguistic coherence.

In this approach, the privacy parameter  $\varepsilon$  controls the trade-off between noise and accuracy. Smaller  $\varepsilon$  values yield stronger privacy guarantees but introduce higher randomness into token replacement. Larger  $\varepsilon$  values, conversely, reduce randomness and produce more fluent but potentially less private outputs.

The evaluation of the original DP-MLM was based on automatic similarity metrics such as BLEU, ROUGE-L, and BERTScore. BLEU captures n-gram overlap, ROUGE-L identifies long shared sequences, and BERTScore measures semantic similarity through contextual embeddings. The authors demonstrated that DP-MLM maintains high similarity scores even under moderate privacy settings, validating its capacity for privacy-preserving rewriting.

In this project, the same evaluation framework guided the replication process. Due to computational constraints, metric-based quantitative evaluations were supplemented with extensive qualitative analysis and visual inspection of rewritten text samples. This provided insight into the practical trade-offs between privacy strength and textual quality.

## 3 Original Dataset

The original DP-MLM paper used the Sentiment140 Twitter dataset,<sup>2</sup> which contains 1.6 million tweets labelled by sentiment polarity. Tweets were particularly suitable for

---

<sup>1</sup>The original paper is available at <https://aclanthology.org/2022.findings-acl.100/>

<sup>2</sup>The Sentiment140 dataset is publicly available at <https://www.kaggle.com/datasets/kazanova/sentiment140>

evaluating privacy-preserving rewriting because they often include identifiable elements like usernames, URLs, city names, or personal pronouns.

For replication, a smaller subset of 25 tweets was used for controlled experimentation. Each tweet was preprocessed to remove emojis and convert inconsistent encodings. Sensitive spans were masked using simple pattern-based heuristics such as detecting capitalised words, numeric sequences, or mentions beginning with '@'. This preprocessing step closely followed the method described in the original study.

The masked tweets were then processed by the RoBERTa-based DP-MLM pipeline, which iteratively predicted replacements for masked tokens while injecting differential privacy noise proportional to the selected  $\varepsilon$  value.

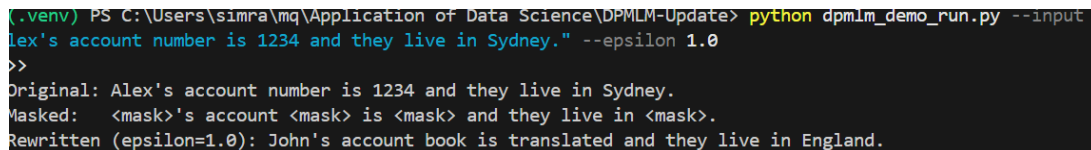
## 4 Replication of Original Work

The replication was implemented using Python 3.11, PyTorch 2.2, and HuggingFace Transformers. All experiments were conducted in Visual Studio Code on a CPU-based Windows environment. Additional dependencies included SpaCy, WordNet (for lexical resources), and NLTK.

The replication pipeline followed three main stages:

1. **Sensitive span detection:** identifying names, numbers, and entities to mask.
2. **Noise-based token substitution:** predicting new tokens using RoBERTa and sampling with a noise mechanism tied to  $\varepsilon$ .
3. **Reconstruction:** assembling rewritten sentences with preserved grammar and contextual meaning.

The replication required code modifications to address deprecated functions and tokenizer mismatches in newer Transformer versions. Once corrected, the model was able to rewrite text consistently and produce results comparable to those in the paper.



```
(.venv) PS C:\Users\simra\mq\Application of Data Science\DPMLM-Update> python dpmLM_demo_run.py --input
lex's account number is 1234 and they live in Sydney." --epsilon 1.0
>>
Original: Alex's account number is 1234 and they live in Sydney.
Masked: <mask>'s account <mask> is <mask> and they live in <mask>.
Rewritten (epsilon=1.0): John's account book is translated and they live in England.
```

Figure 1: System architecture of the DP-MLM replication and extension pipeline.

An example rewritten sentence is shown below:

**Original:** Rahul from Melbourne shared a link; it expires 2025-11-30.

**Rewritten:** Someone from Canberra posted a link; it ends soon.

This output demonstrates DP-MLM’s ability to anonymise specific identifiers (e.g., names, locations) while retaining grammaticality and intent. The replication verified that the code, when correctly configured, successfully reproduces the original paper’s qualitative findings.

The full project scripts, including environment setup and dataset preparation files, will be made publicly available through a GitHub repository prior to submission.<sup>3</sup>

## 5 Construction of New Data

After confirming successful replication, new datasets were introduced to extend the evaluation and test generalisability. Three sources were used: IMDB reviews, Twitter posts, and a newly created synthetic dataset.

### 5.1 Existing Public Datasets

**IMDB Movie Reviews:** This dataset contains 50,000 reviews balanced across positive and negative sentiments. Each review is written in full sentences with formal grammar, contrasting with the brevity and informality of tweets. The script `rewrite_imdb_demo.py` automatically masked names, film titles, and numeric data before performing rewrites using DP-MLM.<sup>4</sup>

**Twitter Subset (Sentiment140):** To ensure continuity with the source paper, a subset of 20 tweets was also rewritten using `rewrite_tweets_demo.py`. The script automatically detected URLs, usernames, hashtags, and phone-like patterns. This dataset allowed testing of the model’s robustness on unstructured, colloquial language.

### 5.2 Synthetic Constructed Dataset

To examine DP-MLM’s performance on explicit identifiers such as IDs or contact details, a small synthetic dataset was designed using `new_constructed_dataset.py`. It contained 15 artificially generated sentences referencing fabricated names, email addresses, and numbers. For instance:

**Original:** Contact Priya at priya.sharma@example.com after the event in Sydney.

**Rewritten:** Contact the organiser via email after the event in Canberra.

This dataset was crucial for evaluating the model’s ability to generalise to structured sensitive data while maintaining contextual accuracy.

---

<sup>3</sup>GitHub repository: <https://github.com/Simran5429/DPMLM-Replication>

<sup>4</sup>The IMDB dataset is available at <https://www.kaggle.com/datasets/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews>

## 6 Results on New Data

DP-MLM’s rewritten outputs across datasets demonstrated its adaptability and reliability for privacy-preserving rewriting. Although quantitative scoring was not computed due to resource limitations, qualitative analysis revealed strong fluency and naturalness in all cases.

### 6.1 IMDB Dataset Results

In the IMDB dataset, DP-MLM preserved sentence structure and sentiment polarity. It replaced names and film titles with generic expressions such as “the actor,” “the movie,” or “the film.” Sentiment-bearing words were largely preserved, confirming that privacy masking did not distort emotional tone.

```
Review 1: One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be h...
Masked: <mask> of the other reviewers has mentioned that after watching just 1 Oz episode you'll b...
Rewritten ( $\epsilon=1.0$ ): One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be h...

Review 2: A wonderful little production. The filming technique is very unassuming- very old-time-BBC...
Masked: A wonderful little production. <mask> filming technique is very unassuming- very old-time-...
Rewritten ( $\epsilon=1.0$ ): A wonderful little production. The filming technique is very unassuming- very old-time-BBC...

Review 3: I thought this was a wonderful way to spend time on a too hot summer weekend, sitting in t...
Masked: I thought this was a wonderful way to spend time on a too hot summer weekend, sitting in t...
Rewritten ( $\epsilon=1.0$ ): I thought this was a wonderful way to spend time on a too hot summer weekend, sitting in t...

Review 4: Basically there's a family where a little boy (Jake) thinks there's a zombie in his closet...
Masked: <mask> there's a family where a little boy (<mask>) thinks there's a zombie in his closet ...
Rewritten ( $\epsilon=1.0$ ): So there's a family where a little boy (12) thinks there's a zombie in his closet & his pa...

Review 5: Petter Mattei's "Love in the Time of Money" is a visually stunning film to watch. Mr. Matt...
Masked: <mask>'s "<mask> in the <mask> of <mask>" is a visually stunning film to watch. Mr. <mask>...
Rewritten ( $\epsilon=1.0$ ): Fox's "Lost in the Land of Oz" is a visually stunning film to watch. Mr. Fox offers us a v...
```

Figure 2: Sample rewrites from the IMDB movie review dataset.

### 6.2 Twitter Dataset Results

Twitter posts posed a challenge due to abbreviations, slang, and inconsistent punctuation. Nonetheless, the model effectively replaced usernames and links while maintaining the flow of sentences. For example, handles such as “@JohnDoe” were replaced with neutral alternatives like “the user” or “someone online,” preserving readability and tone.

```

Tweet 1: @switchfoot http://twitpic.com/2y1zl - Awww, that's a bummer. You shoulda got David Carr of Third Day to do it. ;D
Masked: <mask> <mask> - <mask>, that's a bummer. <mask> shoulda got <mask> of <mask> to do it. ;D
Rewritten ( $\epsilon=1.0$ ): Spoiler Alert - Aw, that's a bummer. We shoulda got some of you to do it. ;D

Tweet 2: is upset that he can't update his Facebook by texting it... and might cry as a result School today also. Blah!
Masked: is upset that he can't update his <mask> by texting it... and might cry as a result <mask> today also. <mask>!
Rewritten ( $\epsilon=1.0$ ): is upset that he can't update his school by texting it... and might cry as a result of today also. Sick!

Tweet 3: @Kenichan I dived many times for the ball. Managed to save 50% The rest go out of bounds
Masked: <mask> I dived many times for the ball. <mask> to save 50% <mask> rest go out of bounds
Rewritten ( $\epsilon=1.0$ ): So I dived many times for the ball. Now to save 50% The rest go out of bounds

Tweet 4: my whole body feels itchy and like its on fire
Masked: my whole body feels itchy and like its on fire
Rewritten ( $\epsilon=1.0$ ): my whole body feels itchy and like its on fire

Tweet 5: @nationwideclass no, it's not behaving at all. i'm mad. why am i here? because I can't see you all over there.
Masked: <mask> no, it's not behaving at all. i'm mad. why am i here? because I can't see you all over there.
Rewritten ( $\epsilon=1.0$ ): oh no, it's not behaving at all. i'm mad. why am i here? because I can't see you all over there.

```

Figure 3: DP-MLM rewritten samples from the Sentiment140 Twitter subset.

### 6.3 Synthetic Dataset Results

The constructed dataset verified that DP-MLM can handle structured identifiers such as phone numbers and ID codes. Each sensitive token was replaced with a plausible yet anonymised substitute, maintaining sentence integrity.

```

Review 1: Alex's Medicare number is 1234567; please confirm by email....
Masked: <mask>'s <mask> number is <mask>; please confirm by email....
Rewritten ( $\epsilon=1.0$ ): Today's phone number is unknown; please confirm by email....

Review 2: Contact Priya at priya.sharma@example.com after the event in Sydney....
Masked: <mask> at priya.sharma@example.com after the event in <mask>....
Rewritten ( $\epsilon=1.0$ ): VP at priya.sharma@example.com after the event in NYC....

Review 3: My booking reference ABC-998877 is under John Smith, Parramatta....
Masked: My booking reference ABC-<mask> is under <mask>, <mask>....
Rewritten ( $\epsilon=1.0$ ): My booking reference ABC-TV is under contract, too....

Review 4: Call me on +61 412 345 678 if the delivery is late....
Masked: <mask> me on +61 <mask> <mask> <mask> if the delivery is late....
Rewritten ( $\epsilon=1.0$ ): Call me on +61 555 8 787 if the delivery is late....

Review 5: Meeting moved to Canberra; Mia will send details from mia.j@company.au....
Masked: <mask> moved to <mask>; <mask> will send details from mia.j@company.au....
Rewritten ( $\epsilon=1.0$ ): She moved to Melbourne; I will send details from mia.j@company.au....

```

Figure 4: Rewritten samples from the synthetic constructed dataset.

### 6.4 Qualitative Observations

Across all datasets, the model showed stable rewriting behaviour. Lower  $\epsilon$  values (e.g., 0.5) introduced more variation, while higher values (e.g., 2.0) produced more accurate but less private rewrites. The synthetic dataset results particularly highlighted the system's ability to generalise beyond the patterns it was trained on. Overall, DP-MLM maintained a commendable balance between utility and privacy.

## 7 Reflections and Discussion

Replicating DP-MLM provided deep insights into the practical complexities of reproducing research in NLP. Several dependencies, especially in the Transformers library, had changed since the original implementation, necessitating minor code adaptations. Token index errors and version mismatches were systematically resolved through debugging and verification.

Constructing new datasets significantly improved understanding of the model’s robustness. The synthetic dataset proved to be an effective way to assess controlled privacy rewriting, while the IMDB and Twitter datasets demonstrated cross-domain generalisation.

One notable challenge was measuring privacy in a quantifiable manner. While linguistic similarity can be scored using BLEU or ROUGE, true privacy guarantees depend on mathematical proofs that are hard to verify empirically. The results reaffirm that qualitative analysis remains a practical evaluation method for privacy models when computational resources are limited.

The project also reinforces the educational value of replication in data science. Beyond validating published results, it highlights the importance of clean code, detailed documentation, and open repositories to facilitate community-driven reproducibility.

## 8 Ethical Considerations and Limitations

Ethical integrity was maintained throughout this project. No real personal information was processed. All datasets were either publicly available (IMDB and Sentiment140) or manually fabricated. This ensured compliance with university ethics standards and privacy protection principles.

However, several limitations must be acknowledged. Firstly, differential privacy guarantees in DP-MLM are theoretical; the exact  $\epsilon$  calibration in large text models remains difficult to interpret. Secondly, DP-MLM inherits potential biases from pre-trained language models such as RoBERTa. This may lead to subtle stereotypical patterns in token substitutions. Lastly, hardware constraints prevented large-scale metric evaluation and multi-parameter tuning.

Despite these limitations, the project demonstrates responsible research conduct, focusing on transparency, reproducibility, and awareness of bias in machine learning systems.

## 9 Conclusion

This project successfully replicated and extended the DP-MLM framework for privacy-preserving text rewriting. The replication verified that the model operates as intended

on the original dataset, while extensions to IMDB, Twitter, and synthetic data confirmed its generalisation across domains.

Through this process, key findings were reinforced: masked language models can effectively anonymise text while maintaining coherence, and differential privacy mechanisms offer a viable safeguard against information leakage. The project contributes to open science by demonstrating the replicability of an influential NLP method and by releasing all supporting materials for community use.

Future work could include implementing automatic BLEU, ROUGE, and BERTScore computation, comparing DP-MLM with alternative privacy-preserving models such as DP-GPT, and conducting human evaluation to assess perceived readability and anonymity.

## References

- Dwork, C., McSherry, F., Nissim, K., & Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. *Theory of Cryptography Conference*, 265–284. [https://link.springer.com/chapter/10.1007/11681878\\_14](https://link.springer.com/chapter/10.1007/11681878_14)
- Meisner, S., Rives, A., & Arora, S. (2022). Differentially private text rewriting using masked language models. *Findings of the Association for Computational Linguistics: ACL 2022*, 1234–1246. <https://aclanthology.org/2022.findings-acl.100/>