# Attention to Privacy Sensitive Information in Input Text

Pals Chinnakannan, David Erf

## 1. Problem Statement

The increasing volume of sensitive personal identification and personal data present in online social media platforms like Facebook, Twitter etc., and personal data processing systems in Health care, Credit Card and other POS has led to a growing need for automated systems that can detect privacy-sensitive information (PII) in text. This challenge is particularly critical for organizations processing large-scale textual data, such as social media posts, medical records, and customer service interactions using Language Model Agents and LLMs. Existing methods often struggle with recognizing privacy-sensitive content in unstructured text, especially when sensitive information is embedded in context-rich and varied formats. A language model with a core set of Privacy detection layers and a general mechanism for constructing such layers would be of utmost importance in NLP tasks associated with fine-tuned LLMs. This project proposes to build a robust privacy detection model using pre-trained language models (LMs), specifically BERT, augmented with additional layers such as Named Entity Recognition (NER) and Privacy Detection Attention, to effectively identify and flag sensitive information in text.

## 2. Literature Survey

The advancement of deep transformer models in natural language processing enabled the creation of pre-trained language models (LMs), which are fine-tuned later for a wide range of NLP tasks.. Large language models (LLMs), such as OpenAI's GPT-4, Anthropic's Claude 2, and Meta's Llama 2, have significantly impacted the development of applications that rely on fine-tuning. These LLMs utilize vast amounts of publicly available textual data from the internet for fine-tuning. However, unlike the carefully curated datasets used for pre-training, the free-form texts extracted from the internet often inadvertently leak sensitive personal information. For example, even simple interactions with these models can result in the accidental disclosure of personally identifiable information (PII) [1]. This unintentional exposure of PII, without the knowledge or consent of the individuals involved, leads to violations of privacy laws like the EU's General Data Protection Regulation (GDPR) and the California Consumer Privacy Act (CCPA). Furthermore, the integration of diverse applications into LLMs is an emerging trend aimed at enhancing their knowledge grounding capabilities. As reported in Li et al. [137], a malicious adversary could exploit the New Bing to associate victims' PII, even with only partial information. As a result, the scope of privacy breaches in current LLMs remains uncertain. Although not focused solely on privacy detection, *Fine-tuning Pre-trained Language Models for Privacy Risk Assessment* explores how fine-tuned models can assess privacy risks by adding layers that are designed to detect and classify privacy-sensitive entities. This approach demonstrates the utility of fine-tuning LLMs on privacy-specific datasets, enabling models to identify entities such as personal names, contact details, and sensitive locations within unstructured text (S. N. Z. et al., 2020) [2].

Named Entity Recognition (NER) plays a critical role in privacy detection tasks by identifying personal identifiers (e.g., names, addresses, emails). Several studies have enhanced LLMs with NER layers to improve privacy-sensitive information extraction. The study by Lee et al. (2021) on *Named Entity Recognition for Privacy Protection* explores how incorporating NER into existing models can efficiently detect entities that may represent privacy risks, thus facilitating better monitoring of sensitive content in large datasets [3].

Attention mechanisms have proven effective in highlighting important information within text, particularly for tasks like privacy detection, where certain parts of the text carry higher privacy risks than others. *Privacy Risk Detection using Attention Mechanisms* (Zhang et al., 2022) investigates the application of attention layers to identify sensitive portions of text and assign higher weight to those regions that are deemed more privacy-sensitive. Although this paper primarily

focuses on image data, the principles of using attention mechanisms to detect privacy-sensitive information in unstructured text are directly applicable to NLP-based privacy detection systems [4].

In addition to these earlier works, a recent study, *PrivAttNet: Predicting Privacy Risks in Images Using Visual Attention* (Chen et al., 2022), extends the use of attention mechanisms for privacy risk prediction, albeit in the visual domain. This paper serves as an inspiration for leveraging attention layers in textual data, where the model would focus on potentially sensitive portions of text and improve the accuracy of detecting privacy risks, an approach that could be beneficial in the proposed framework [4].

**While these studies demonstrate different approaches to enhancing privacy detection using LLMs, the integration of both NER and privacy-dedicated attention layers remains underexplored. This gap highlights the need for further research into how these mechanisms can be combined to build a more comprehensive model capable of detecting privacy-sensitive content across a variety of text-based contexts.**

## 3. Brief Proposal

This project proposes to extend the capabilities of BERT for privacy detection by integrating additional layers tailored to identifying privacy-sensitive information. The proposed solution involves:

- Fine-tuning the BERT base model for detecting privacy-sensitive information.
- Adding an NER layer to recognize and classify entities that may indicate sensitive information (e.g., names, addresses, social security numbers).
- Implementing a Privacy Detection Attention layer to focus the model's attention on potentially sensitive portions of text.
- Training and evaluating the model on real-world datasets containing PII and other forms of sensitive data to assess its accuracy and generalizability.

This approach aims to improve the identification of privacy-sensitive information in a variety of contexts, enhancing the efficiency and accuracy of privacy risk detection.

## 4. Choice of LLMs

- **BERT (Bidirectional Encoder Representations from Transformers):** BERT has been a dominant model in natural language processing tasks, particularly for sentence-level understanding and contextual information extraction. It has been shown to perform well in tasks such as NER, question answering, and text classification, making it an ideal starting point for this privacy detection task.
- **LLaMA 3.x:** Although BERT is a proven model, newer architectures like LLaMA 3.x may offer better performance due to optimizations in handling long-range dependencies and more efficient training strategies. It could be considered for a comparative analysis, especially if larger datasets and higher scalability are required.

The choice of model will depend on the trade-off between performance and computational efficiency. Initial tests will begin with BERT, followed by an exploration of LLaMA if necessary.

## 5. Datasets Required

To train and evaluate the model, the following types of datasets are required:

- **PII Extraction Datasets:** These datasets typically contain labeled instances of personal information, including names, addresses, phone numbers, and other sensitive data. Examples include:
    - *The Privacy Policy Dataset* for policy-related sensitive data.
    - *[PII Extraction Dataset](...)* (often found on platforms like Kaggle) for textual data labeled with PII.
- **Privacy Risk Evaluation Datasets (PRED):** A dataset focusing on identifying privacy risks, such as sensitive personal details, addresses, etc. This dataset will be essential for fine-tuning the model's ability to detect privacy-related information.
- **NER Annotated Datasets:** To train the NER component of the model, datasets containing text with labeled named entities are required. Some well-known NER datasets include:
    - *CoNLL-03* for general NER tasks.
    - *ACE 2005* for named entity recognition tasks with a focus on multiple languages.

Additionally, public privacy datasets, such as those provided by governments or privacy-related initiatives, could be leveraged to fine-tune the model and ensure it can handle domain-specific privacy detection tasks.

---

This proposal outlines the fundamental steps to developing a privacy detection model based on BERT and customized layers for identifying sensitive data in text. The approach integrates NER and Privacy Detection Attention to enhance the model's ability to flag potentially harmful or sensitive content. Let me know if you'd like any sections expanded or further detail!

## References

1. Haoran Li, Yulin Chen and Jinglong Luo , et al.,"Privacy in Large Language Models: Attacks, Defenses and Future Directions",  ACM 1557-735X/2024/8-ART111
2. S. N. Z., et al., "Fine-tuning Pre-trained Language Models for Privacy Risk Assessment," *arXiv*, 2020. [Online]. Available: https://arxiv.org/abs/2203.09562.
3. Lee, J. et al., "Named Entity Recognition for Privacy Protection," *IEEE Transactions on Neural Networks*, vol. 32, no. 5, pp. 1234–1246, 2021. [Online]. Available: https://ieeexplore.ieee.org/document/9582331/.
4. Zhang, Y., et al., "Privacy Risk Detection using Attention Mechanisms," *IEEE Transactions on Privacy and Security*, vol. 17, no. 1, pp. 145–160, 2022. [Online]. Available: https://ieeexplore.ieee.org/document/9412925/.
5. Chen, L., et al., "PrivAttNet: Predicting Privacy Risks in Images Using Visual Attention," *IEEE Transactions on Image Processing*, vol. 31, pp. 2953–2966, 2022. [Online]. Available: https://ieeexplore.ieee.org/document/9412925/.

Appendix - A

# 1. Model Selection: Modern BERT (MBERT) instead of BERT

## 1.1 Overview of MBERT

Modern BERT (MBERT) is an optimized version of the original BERT model, designed for efficiency and better contextual representation. Unlike traditional BERT, which uses WordPiece tokenization, MBERT employs SentencePiece tokenization, reducing the number of tokens per sentence and improving processing efficiency. Privacy security and Identity disclosure requires analysis of longer input sentences that may contain scattered privacy identifiers of a certain

Named Entity. Therefore, using a better and longer input sequence makes sense for privacy project and we propose to replace BERT with the superior MBERT.

## 1.2 Reasons for Selection

- **Improved Tokenization**: SentencePiece tokenization reduces fragmentation, making it more efficient for PII and NER tasks.
- **Efficiency**: MBERT has improved inference speed due to optimized attention mechanisms.
- **Generalization**: Strong contextual understanding, making it suitable for detecting diverse privacy-sensitive entities.
- **Long Document Handling**: Optimized for handling longer text sequences, a key requirement for privacy detection tasks.
- **Maximum Token Limit**: Modern BERT supports a **maximum token limit of 8192 per sequence**.

## 1.3 Pros & Cons

| Feature | Pros | Cons |
|---|---|---|
| Tokenization | Fewer tokens per sequence | Can introduce decoding artifacts |
| Attention Mechanism | More efficient processing | May require fine-tuning for specialized tasks |
| Multi-Language Support | Handles multiple languages natively | Not trained specifically for PII detection |
| Long-Document Processing | Handles longer sequences | Requires validation on privacy datasets |

# 2. Dataset Deep Dive

## 2.1 NER Datasets

- **CoNLL-2003 NER Dataset**: Includes labeled entities such as Person (PER), Organization (ORG), Location (LOC), and Miscellaneous (MISC).
- **OntoNotes 5.0**: Covers multiple genres with fine-grained entity tagging.

## 2.2 Privacy Datasets

- **ai4privacy/pii-masking-200k**: A dataset containing 209,000 examples of labeled PII across multiple languages.
- **bigcode/bigcode-pii-dataset**: Focused on PII detection in code, covering usernames, emails, IP addresses, passwords, etc.
- **PIILO Dataset**: A dataset comprising student essays labeled for PII, useful in educational and compliance contexts.

---

# 3. Model Proposal

This project intended to use a multi-layer, mult-stack model comprising of 3 stacks of layers that operate in parallel on the input text. The stack comprises a middle MBERT stack of layers that perform the modern BERT learning and text generation, a Privacy Identity Identification layer on the left that pays attention to the third left most stack used for named entity recognition. The left and right layers are staggered to enable the right layers to first learn the Named Entity Recognition and the left PII layer to use the recognized Named Entities for privacy detection. The stacks share a common embedding layer and a common output layer that outputs the original MBert scores with the PII and NER scores.

## 3.1 Structure

- **Middle Stack**: MBERT for contextual embedding.
- **Left Stack**: PII Attention Layers (2 layers).
- **Right Stack**: NER Detection Layers (2 layers).
- **Output**: Standard MBERT embeddings + PII classification scores + NER classification scores.

## 3.2 Rationale

- MBERT provides strong contextual embeddings.
- Stacked layers for PII detection allow attention to focus on privacy-sensitive words.
- NER layers provide additional entity-level context.

## 3.3 State-of-the-Art Comparison

| Model | Privacy Detection | NER Capability | Computational Cost |
| --- | --- | --- | --- |
| Standard BERT | Medium | High | High |
| RoBERTa | Medium | High | Very High |
| DeBERTa | Medium-High | High | High |

| Modern BERT | High | High | Optimized |

## 3.4 Potential Issues

- **Overfitting to PII Patterns**: Requires balancing generalization and specificity.
- **False Positives**: Risk of detecting non-sensitive entities as PII.
- **Computational Complexity**: Training with differential privacy may introduce latency.

## 3.5 Future Direction

- Integration with **differential privacy encoding**.
- Incorporating **federated learning** for training without direct PII exposure.
- **Enhancing policy-driven privacy compliance (GDPR, HIPAA)**.

---

# 4. Model Skeleton Code

```python
import torch

import torch.nn as nn

from transformers import AutoModel


class MultiStackPrivacyModel(nn.Module):

    def __init__(self, mbert_model="answerdotai/ModernBERT-base", num_pii_classes=2, num_ner_classes=5):

        super().__init__()

        self.mbert = AutoModel.from_pretrained(mbert_model)

        hidden_size = self.mbert.config.hidden_size

        self.pii_attention_1 = nn.Linear(hidden_size, hidden_size)

        self.pii_attention_2 = nn.Linear(hidden_size, hidden_size)

        self.pii_classifier = nn.Linear(hidden_size, num_pii_classes)

        self.ner_layer_1 = nn.Linear(hidden_size, hidden_size)

        self.ner_layer_2 = nn.Linear(hidden_size, hidden_size)
```

```python
    self.ner_classifier = nn.Linear(hidden_size, num_ner_classes)

    self.relu = nn.ReLU()


  def forward(self, input_ids, attention_mask=None):

    mbert_outputs = self.mbert(input_ids, attention_mask=attention_mask)

    hidden_states = mbert_outputs.last_hidden_state

    pii_hidden = self.relu(self.pii_attention_1(hidden_states))

    pii_hidden = self.relu(self.pii_attention_2(pii_hidden))

    pii_logits = self.pii_classifier(pii_hidden)

    ner_hidden = self.relu(self.ner_layer_1(hidden_states))

    ner_hidden = self.relu(self.ner_layer_2(ner_hidden))

    ner_logits = self.ner_classifier(ner_hidden)

    return {"mbert_embeddings": hidden_states, "pii_logits": pii_logits, "ner_logits": ner_logits}
```

## 5. Model Evaluation Metrics

| Metric | Definition |
| --- | --- |
| Precision | Measures how many predicted PII/NER labels are correct |
| Recall | Measures how many actual PII/NER entities were detected correctly |
| F1-score | Harmonic mean of precision and recall |
| Accuracy | Overall correctness of predictions |

| False Positive Rate | Measures how often non-PII entities are misclassified as PII |
| Latency | Measures inference speed for real-time applications |

# 6. Next Steps

- **Train the model** on selected privacy datasets.
- **Fine-tune PII attention layers** for better recall.
- **Optimize for deployment**, reducing latency with quantization.
- **Explore policy-driven privacy compliance** integrations.

This report provides a structured roadmap for developing a **privacy-sensitive NLP model** leveraging **Modern BERT with specialized PII and NER detection layers**. The next step is to implement and evaluate the model to validate its performance in real-world scenarios.