# Enhancing the De-identification of Personally Identifiable Information in Educational Data

Yuntian Shen*, Zilyu Ji*, Jionghao Lin, and Kenneth R. Koedinger

*Abstract*—Protecting Personally Identifiable Information (PII), such as names, is a critical requirement in learning technologies to safeguard student and teacher privacy and maintain trust. Accurate PII detection is an essential step toward anonymizing sensitive information while preserving the utility of educational data. Motivated by recent advancements in artificial intelligence, our study investigates the GPT-4o-mini model as a cost-effective and efficient solution for PII detection tasks. We explore both prompting and fine-tuning approaches and compare GPT-4o-mini's performance against established frameworks, including Microsoft Presidio and Azure AI Language. Our evaluation on two public datasets, CRAPII and TSCC, demonstrates that the fine-tuned GPT-4o-mini model achieves superior performance, with a recall of 0.9589 on CRAPII. Additionally, fine-tuned GPT-4o-mini significantly improves precision scores (a threefold increase) while reducing computational costs to nearly one-tenth of those associated with Azure AI Language. Furthermore, our bias analysis reveals that the fine-tuned GPT-4o-mini model consistently delivers accurate results across diverse cultural backgrounds and genders. The generalizability analysis using the TSCC dataset further highlights its robustness, achieving a recall of 0.9895 with minimal additional training data from TSCC. These results emphasize the potential of fine-tuned GPT-4o-mini as an accurate and cost-effective tool for PII detection in educational data. It offers robust privacy protection while preserving the data's utility for research and pedagogical analysis. Our code is available on GitHub: https://github.com/AnonJD/PrivacyAI

*Index Terms*—Privacy, De-identification, Anonymization, Personally Identifiable Information, Large Language Models, Fine-tuning GPT, Cost-effectiveness, Hidden in Plain Sight

## I. INTRODUCTION

**P**ERSONALLY Identifiable Information (PII) includes the information (e.g., names, email addresses, and phone numbers) that can identify an individual. Protecting PII in educational data is paramount as learning technologies, including artificial intelligence-powered systems, become increasingly integral to education across all levels. For instance, online human tutoring platforms collect vast amounts of interaction data, including conversational dialogues between students and tutors [1]. While the analysis of such educational data offers insights for data-informed research and improved pedagogical practices [2], [3], [4], it also introduces significant privacy risks due to the sensitive nature of PII [5]. Safeguarding PII is essential to prevent unauthorized access and misuse [6]. Moreover, robust privacy protections are vital for fostering trust among educational stakeholders—students, educators, and parents—and ensuring the responsible adoption of learning technologies [7]. These efforts align with global regulatory frameworks, such as the General Data Protection Regulation (GDPR) and the Family Educational Rights and Privacy Act (FERPA), which mandate the protection of personal data [8].

Given the vast amounts of data collected by learning systems, automating the anonymization of PII in the education domain has become a critical necessity. Previous work has explored rule-based, statistical-based, and neural network-based approaches to detect PII in education-related tasks, such as essay grading [3], [9]. Though showing effectiveness, given the advent of advanced artificial intelligence (AI) models, Pii detection can be further enhanced in detection accuracy, robustness, and generalizability when applied to diverse and complex educational datasets. Moreover, many previous works [10], [11], [12], [13] focus solely on redacting PII from the original data (e.g., the name John from the original sentence *"Thanks, John"* was processed by redacting *"Thanks, [REDACTED]"*). While this approach mitigates certain risks, it falls short of ensuring comprehensive privacy protection, especially given that current models still fail to achieve perfect accuracy in PII detection [14]. This raises concerns about potential PII leakage. Thus, these challenges underscore the urgent need for more advanced, scalable, and robust methods to facilitate the anonymization process in educational data, ensuring comprehensive protection of PII and mitigating the risks of unintended exposure.

As proposed by [6], the Hidden in Plain Sight (HIPS) method introduces the protection of PII in datasets by first identifying PII entities and then replacing them with synthetic information that retains contextual characteristics. Unlike traditional redaction methods that use generic tokens such as `[REDACTED]` to redact PII, the first step of the HIPS method replaces PII with tokens that indicate the type of information, such as `<Name>`, `<Email>`, or `<Address>`. The process of identifying PII entities relies on Named Entity Recognition (NER), a task in natural language processing that aims to locate and classify specific entities within text into predefined categories, such as names, email addresses, and phone numbers [15]. Once NER is performed on the data corpus, the original PII entities can be replaced with synthetic information while retaining their contextual relevance. For example, a sentence such as "$\{John\ Smith\}_{Name}$ *lives at* $\{123\ Main$

Y. Shen, Z. Ji are with Carnegie Mellon University, Pittsburgh, PA 15217, USA (e-mail: yuntian2@andrew.cmu.edu, zilyuj@andrew.cmu.edu)

J. Lin is with the Faculty of Education, The University of Hong Kong, Hong Kong, PR China, also with the Human-Computer Interaction Institute, Carnegie Mellon University, Pittsburgh, PA, 15213, USA, and also with the Centre for Learning Analytics, Faculty of Information Technology, Monash University, Clayton, VIC 3800, Australia (e-mail: jionghao@hku.hk).

K.R. Koedigner is with the Human-Computer Interaction Institute, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15217, USA (e-mail: krk@cs.cmu.edu)

*Yuntian Shen and Zilyu Ji contributed equally to this work.

*Corresponding author: Jionghao Lin*

GPT-based models have demonstrated their ability to achieve high recall scores, effectively identifying a broad range of PII entities [25]. Additionally, LLMs can leverage this internal knowledge to distinguish actual PII from non-PII. Some recent evidence supports this view: for instance, a recent study prompts GPT-3 for a NER task, suggesting that fine-tuning such large models could yield more notable results [25], especially since fine-tuned LLMs tend to outperform their prompted counterparts on NLP tasks [26], [27].

The approach of fine-tuning an LLM for PII identification aligns well with the emerging AI-for-education domain, where labeled text data for PII identification are often scarce and fragmented. Due to this limited availability, a model can be trained on data that are not representative of the actual use case, and less experienced users may require a tool that can quickly adapt to their specific domain. Consequently, a source model capable of learning effectively from sparse and unrepresentative datasets is needed. Previous work has shown that an LLM can be fine-tuned with just a few labeled examples (few-shot learning) [27]. Although studies also indicate that successful few-shot learning is achievable with other architectures, these typically require carefully designed learning strategies and meticulous parameter tuning to prevent overfitting [28], [29]. In contrast, larger LLMs have been proven to be more resistant to overfitting as their size increases [27], [30], potentially owing to their memorization dynamics [31].

A final advantage of fine-tuning an LLM for PII de-identification lies in its lower financial and computational cost compared to prompt-engineering an LLM for PII de-identification, which often requires multiple demonstrations in the input prompts [25]. The above observation highlights the potential for fine-tuning large models, with GPT emerging as a promising candidate. Several studies have utilized prompted GPT approaches to de-identify PII entities, demonstrating encouraging results [17], [25]. With the growing availability of cost-effective fine-tuning APIs [1] and the lack of prior work using fine-tuned GPT for PII identification, we aim to explore this approach further.

## III. Methods

### A. Data and Data Pre-processing

Our study utilized the Cleaned Repository of Annotated Personally Identifiable Information (CRAPII)[2] dataset [14], which comprises 22,688 samples of student writings collected from a massive open online course (MOOC) offered by a university in the United States. The course focused on critical thinking through design, teaching learners strategies such as storytelling and visualization to solve real-world problems [14]. The dataset includes seven PII categories as direct identifiers: `Names`, `Email Addresses`, `Usernames`, `IDs`, `Phone Numbers`, `Personal URLs`, and `Street Addresses` [14]. A sample of the dataset is shown in Table I. In total, the dataset contains 4,871 labeled words categorized as PII

entities. To enable entity-based matching during our analysis, we extracted the character-wise positions of all annotated PII entities within the text.

TABLE I
ILLUSTRATED EXAMPLE OF CRAPII DATASET

| Attribute | Example Value |
|---|---|
| full_text | *Hi John Doe. Tel: (555)555-5555* |
| document | 379 |
| tokens | ['Hi', 'John', 'Doe', '.', 'Tel', ':', '(555)555-5555'] |
| labels | ['O', 'B-NAME', 'I-NAME', 'O', 'O', 'O', 'B-PHONE_NUM'] |
| trailing_whitespace | [True, True, False, True, False, True, False] |

We also introduced another dataset, the Teacher-Student Chatroom Corpus (TSCC) [10], to examine the generalizability of our investigated models, as detailed in Section III-F. Generalizability is crucial for evaluating a model's robustness when applied to datasets with different distributions. The TSCC dataset contains 260 chatroom sessions, with a total of 41.4K conversational turns and 362K word tokens. We processed the dataset by extracting the `role` and `edited` columns and combining them into a simplified `role: text` format, where each conversational turn is represented on a new line. A sample excerpt of the processed transcription is shown below:

```
teacher: Hi there ⟨STUDENT⟩, all OK?
student: Hi ⟨TEACHER⟩, how are you?
```

### B. Models for PII Detection

**Microsoft Presidio.**[3] It is an open-source toolkit designed for detecting and anonymizing PII across various text and image formats. It offers pre-built NER (Named Entity Recognition) models, such as *en_core_web_lg* and *en_core_web_trf*, which utilize linguistic patterns and deep learning techniques to classify words or phrase into predefined named entities (e.g., names, emails, and phone numbers). Additionally, *Presidio* provides options for integrating external machine learning models for enhancing PII detection.

In our study, *Presidio* serves as one of the baseline models for detecting PII entities. As indicated in previous work [14], *en_core_web_lg* is a standard NER model based on the `spaCy`[4] framework that primarily relies on statistical techniques, while the transformer-based *en_core_web_trf* leverages pre-trained transformers, capturing long-range dependencies for more accurate context-based entity recognition. Thus, we primarily used these two configurations: *en_core_web_lg* and *en_core_web_trf*, as depicted in Fig. 1 (①).

---

Fig. 1. Overview of Five PII Detection Models. ① **Presidio**: Uses Microsoft *Presidio* with pre-trained `spaCy` models (`en_core_web_lg` and `en_core_web_trf`) to detect PII entities. ② **Azure AI Language**: Leverages the PII detection feature in Microsoft `Azure AI Language` for entity recognition in input text. ③ **Prompting LLM**: Utilizes `GPT-4o-mini` with few-shot prompting and special identifiers to annotate PII entities in text. The red arrow shows the prompt used to guide the model for PII annotation. ④ **Fine-tuning LLM**: Fine-tunes a `GPT-4o-mini` model for PII detection. The blue arrow represents the prompt used to train the model during fine-tuning. ⑤ **Verifier Models**: Fine-tunes a `GPT-4o-mini` model to verify detected entities from the base model within their textual context. The green arrow indicates the prompt used to verify the entities within their context. Verification is performed with and without chain-of-thought (CoT) reasoning.

**Azure AI Language.**[5] It offers a cloud-based PII detection service capable of identifying and redacting sensitive information such as phone numbers and email addresses. Our study adopted *Azure AI Language* as one of the baseline models for detecting PII entities in our dataset, as illustrated in Fig. 1 (②). We used asynchronous processing through the REST API[6] to handle texts up to 125,000 characters per document. This approach is necessary because the longest transcript in our dataset contains 17,405 characters, exceeding the synchronous processing limit of 5,120 characters. Requests are batched with 5 documents per request and a rate limit of 1000 requests per minute, ensuring compliance with Azure's service limits for PII detection[7]. Each asynchronous request was completed within 24 hours, ensuring timely processing of all data.

**GPT-4o-mini (Prompting).** Motivated by the effectiveness of prompting LLMs in de-identifying PII, as demonstrated in a recent study [17], our study adapts their prompting strategy with modifications to fit our task. Instead of employing `GPT-4`, as used in their work [17], we opted for

`GPT-4o-mini`, which requires only 1/200 of the cost per million input tokens and 1/100 of the cost per million output tokens compared to `GPT-4`[8]. Then, we modified the prompt structure by leveraging special identifiers for labeling detected entities rather than using a redaction method such as replacing detected PII with a generic label like `[REDACTED]` as shown in their work [17].

Inspired by the GPT-NER method [25], special identifiers can be used to mark entities, preserving non-PII content and ensuring precise PII detection. This method also reduces issues such as hallucinations and over-labeling that can arise with generative models. In particular, we require the detected entity positions to be an exact match to the true entity positions, but GPT often struggles with positional accuracy in long texts due to hallucinations and counting limitations [32]. To address this, we let GPT label detected PII entities by surrounding them with special identifiers for different categories, as outlined in Table II, and subsequently extract them using regular expressions, ensuring accurate detection and positioning.

TABLE II
SPECIAL TOKENS FOR PII DETECTION

| PII Category | Special Identifiers | Example |
|---|---|---|
| **Student Name** | `@@@` *Text* `###` | @@@John Doe### |
| **Personal URL** | `&&&` *Text* `$$$` | &&&www.example.com$$$ |
| **Personal Email** | `QQQ` *Text* `^^^` | QQQjohnd@example.com^^^ |
| **Phone Number** | `%%%` *Text* `~~~` | %%%(555)555-5555~~~ |

We employ a few-shot learning strategy to guide `GPT-4o-mini` in accurately identifying different types of PII entities. Few-shot learning provides contextual examples that help the model understand the diversity of entity types and nuances within our dataset [33]. This approach is chosen because the additional examples often improve the model's ability to handle ambiguous cases and enhance consistency in PII identification. Table III outlines the structure of the **System**, **User**, and **Assistant** prompts used in prompting `GPT-4o-mini`. Three examples[9] are selected from the CRAPII dataset and incorporated into the user prompt. The **Assistant**'s output is the input text with labeled PII entities. The prompting process is detailed in Fig. 1 (③).

**GPT-4o-mini (Fine-tuning).** Recent studies [4] have demonstrated the effectiveness of fine-tuning LLMs for NER tasks, specifically showing that fine-tuned LLMs significantly outperform prompting-based approaches in accurately identifying entities. Inspired by these findings, our study employed a fine-tuning strategy on `GPT-4o-mini` for PII detection. We utilized the same approach of incorporating special tokens, as shown in Table II. Then, `GPT-4o-mini` was fine-tuned on the training dataset and evaluated on the test dataset to assess its performance. The detailed structure of the system, user, and assistant prompts used during the fine-tuning process is presented in Table IV, while the implementation process is illustrated in Fig. 1 (④).

[8]OpenAI Models Pricing: https://openai.com/api/pricing/
[9]https://github.com/AnonJD/PrivacyAI/blob/main/Privacy_main/prompts.txt

TABLE III
PROMPT STRUCTURES FOR PROMPTING-GPT-4O-MINI USING SPECIAL IDENTIFIERS

| Role | Content |
|---|---|
| **System** | *You are an expert in labeling Personally Identifiable Information (PII). Start your response right away without adding any prefix (such as "Response:") or suffix. Use special identifiers to mark different types of PII in the given text.* |
| **User** | *Label the entity of the following text: @@@, ### to label student name; &&&, $$$ to label personal URL; QQQ, ^^^to label personal email; %%%, ~~to label phone number. Ensure that the rest of the text remains unchanged, word for word. Maintain the original punctuation, quotation marks, spaces, and line breaks. If the text does not contain any PII, return it as is.* <br><br> *For example, if the input is:* {Example One} <br><br> *The output should be:* {Example One with Labeled PII} <br> *Another example:* {Example Two} <br> *Another example:* {Example Three} <br> *Please repeat this process with the following file:* {Text Input} |
| **Assistant** | {Text Input with Labeled PII} |

TABLE IV
DETAILS OF FINE-TUNING GPT-4O-MINI USING SPECIAL IDENTIFIERS

| Role | Content |
|---|---|
| **System** | *You are an expert in labeling Personally Identifiable Information. Start your response right away without adding any prefix (such as Response:) or suffix.* |
| **User** | *Label the entity of the following text: @@@, ### to label student name; &&&, $$$ to label personal URL; QQQ, ^^^to label personal email; %%%, ~~to label phone number.* <br><br> {Text Input} |
| **Assistant** | {Text Input with Labeled PII} |

**Verifier Models.** To improve the precision (Equation 1) of PII detection while maintaining recall (Equation 2), we propose the addition of a verifier model as a second step to verify whether predictions are accurately identified as PII. This approach is inspired by recent studies [34], [35] that integrate verifiers into multi-step reasoning tasks in LLMs. For the implementation of verifier models, we propose two variants: *Verifier Model I (Without CoT)* and *Verifier Model II (With CoT)*. These verifier models assess detected entities within their surrounding context to eliminate false positives while retaining true PII entities, thereby enhancing the precision of PII detection systems. Both verifiers are fine-tuned versions of the `GPT-4o-mini` model. The process of the verifier model approach is illustrated in Fig. 1 (⑤).

The *Verifier Model II (With CoT)* incorporates Chain-of-Thought (CoT) reasoning, a method shown to enhance decision-making by breaking down complex problems into intermediate steps [36]. The *Verifier Model II (With CoT)* aims to improve interpretability and robustness by generating reasoning before the final classification. However, it generates more output tokens, leading to higher computational costs. The *Verifier Model I (Without CoT)*, which directly classifies

TABLE V
PROMPT STRUCTURES FOR VERIFIER MODELS

| Role | Verifier Model I (Without CoT) | Verifier Model II (With CoT) |
|---|---|---|
| System | *You are an expert in labeling Personally Identifiable Information. Start your response right away without adding any prefix (such as Response:) or suffix.* | |
| User | *Determine if* `{entity}` *is a privately identifiable information in its context:* `{context}`*, think carefully before saying no to protect against PII leakage, only output T or F.* | *Determine if* `{entity}` *is a privately identifiable information in its context:* `{context}`*. Think step-by-step before outputting T or F.* |
| Assistant | *"T" or "F"* | *"*`{CoT Reasoning} + T`*" or "*`{CoT Reasoning} + F`*"* |

entities without reasoning, is retained for scenarios where computational resources are limited or speed is a higher priority.

A subset of the CRAPII dataset is used to construct the training data for the verifier models. Detected entities are obtained from a chosen base model, and their context is extracted. Entities are labeled as `T` (True PII) or `F` (False PII) based on ground-truth labels in the CRAPII dataset. For the *Verifier Model II (With CoT)*, we use `GPT-4o-mini` to generate reasoning to support the classification. If the reasoning does not align with the ground truth after six attempts, the label defaults to `T` to avoid mistakenly removing true PII entities and prioritize privacy preservation. The prompt structures for both verifier models are shown in Table V.

Once the verifier models are trained, they can be applied to verify detected entities from any base model depending on user priorities. It is important to note that applying the verifier will not increase recall, as false negatives remain unchanged. Instead, the verifier reduces false positives, potentially at the expense of some true positives. If preserving PII is a higher priority than maximizing precision, the verifier should be applied to the base model with the highest recall. Conversely, for tasks that emphasize precision, the verifier may be applied to other base models as needed.

## C. Splitting Data for Training and Testing

Our study focuses on four specific categories from the CRAPII dataset for PII detection: `NAME`, `URL_PERSONAL`, `EMAIL`, and `PHONE_NUM`. To support the training and evaluation requirements of our experiments, the 22,688 files in the CRAPII dataset are split into three distinct sets: **Base Train Set** (25%, 5,672 files), **Verifier Train Set** (15%, 3,403 files), and **Test Set** (60%, 13,613 files). This split ensures that all sets contain a sufficient number of entities from each category, including rare categories such as `PHONE_NUM`, which only has 15 entities in total. Table VI presents the distribution of true entity counts across the three sets.

TABLE VI
TRUE ENTITY COUNTS ACROSS DATA SPLITS

| Entity Type | Total | Base Train Set | Verifier Train Set | Test Set |
|---|---|---|---|---|
| NAME_STUDENT | 4394 | 1091 | 693 | 2610 |
| URL_PERSONAL | 354 | 76 | 66 | 212 |
| EMAIL | 112 | 29 | 21 | 62 |
| PHONE_NUM | 15 | 3 | 3 | 9 |
| **TOTAL** | **4871** | **1199** | **783** | **2889** |

## D. Evaluation Metrics

When evaluating the performance of a model designed to detect PII, it is essential to assess the model's ability to correctly identify true PII entities while minimizing incorrect classifications. The **True Positives** (TP) indicate the number of correctly identified PII entities. The **False Positives** (FP) represent the number of non-PII entities incorrectly classified as PII, which can reduce the utility of the dataset by unnecessarily removing valuable information. The **False Negatives** (FN) denote the number of missed detections of actual PII entities, which could result in privacy breaches and violations of legal regulations. To evaluate the model's capability in accurately detecting PII, we adopt the evaluation method suggested in [15]. This method considers a PII entity to be correctly identified only if it is an exact match with the corresponding entity in the text data. To provide a comprehensive assessment of the model's performance, we consider multiple metrics that capture different aspects of the model's effectiveness.

**Precision** (Equation 1) measures the proportion of correct PII predictions out of all entities that the model classified as PII. High precision means that when the model identifies words as PII, it is very likely to be correct. This is particularly important when false positives (incorrectly labeling non-PII as PII) are costly or disruptive, such as when anonymizing educational datasets where unnecessary removal of non-PII data can reduce the value of the dataset for analysis.

$$\text{Precision} = \frac{TP}{TP + FP} \quad (1)$$

**Recall** (Equation 2) measures the proportion of actual PII entities that the model successfully identifies. High recall is crucial in privacy protection, as it ensures that most, if not all, sensitive information is detected and appropriately handled, minimizing the risk of unmasked PII being exposed and resulting in potential privacy breaches or legal violations.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2)$$

The $F_1$ **Score** (Equation 3) provides a balance between precision and recall, offering a single metric that reflects both the model's ability to correctly identify PII and its ability to minimize false positives. The $F_1$ Score is especially useful when both precision and recall are equally important.

$$F_1 \text{ Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3)$$

The $F_5$ **Score** (Equation 4) places a stronger emphasis on recall than precision. In privacy-sensitive domains like

education, where missing a piece of PII (a false negative) can be much more damaging than accidentally flagging non-PII as sensitive, the $F_5$ Score helps prioritize models that do a better job at catching all PII, even if it means more false positives.

$$F_5 \text{ Score} = (1 + 5^2) \times \frac{\text{Precision} \times \text{Recall}}{(5^2 \times \text{Precision}) + \text{Recall}} \quad (4)$$

### E. Analysis of Cultural and Gender Bias in Name Detection by Models

Evaluating the model's performance in detecting `<NAME_STUDENT>` entities across different cultural and gender distributions is quite important. Models trained on imbalanced datasets may underperform in identifying names of groups that are underrepresented in training data. Specifically, if the training data lacks sufficient representation of names from certain cultural backgrounds, the model may exhibit a lower recall or precision for those names. This dimension of model evaluation is necessary to ensure fairness and inclusivity in PII detection systems, especially as names serve as direct identifiers with critical privacy implications.

We analyzed the cultural and gender distributions of the names in the `<NAME_STUDENT>` entities in detail by using a two-step approach. First, we used a *rule-based name parser*[10] to split each name into components, typically a first name and a last name. We then determined the *gender* of each name based on the first name and matched the *nationality* using the last name. This method aligns with the approach described in the CRAPII paper [14]. In the second step, we mapped the identified countries to their respective *regional cultures* using the ISO-3166 dataset with the UN regional codes[11]. Specifically, we relied on the `all.format` file, which includes detailed regional and sub-regional classifications for each country. For example, "Nigeria" maps to the *Africa* region, while "United States of America" maps to the *Americas*. The ISO-3166 dataset provides five cultural regions: *Asia, Americas, Europe, Africa, and Oceania*. Notably, no names in the CRAPII dataset belonged to Oceania, so we focused on the remaining four cultural groups.

### F. Analysis of Models' Generalizability

To analyze the generalizability of our investigated models, we employed the TSCC dataset, as introduced in Section III-A. Notably, the TSCC dataset has been pre-redacted, with most redacted words replaced by placeholders indicating their name entity types, such as `<STUDENT>` and `<TEACHER>`. To evaluate the model's ability to generalize across diverse contexts, it was necessary to replace these placeholders with synthetic entities that reflect diversity in gender and cultural backgrounds. This replacement ensures a realistic evaluation of the models' performance when encountering unseen data with varying demographic characteristics. To generate a diverse and representative set of synthetic names, we systematically created a mapping of first and last names categorized by gender and cultural groups. This process involved the following steps:

[10]https://github.com/derek73/python-nameparser
[11]https://github.com/lukes/ISO-3166-Countries-with-Regional-Codes

1) For each culture, we identified its corresponding countries based on the `all.format` file in the United Nations dataset (`all.csv`[12]).
2) For each country, we used the `get_top_names` function from the `names-dataset` Python package[13] to retrieve first names based on gender and last names based on the country, as outlined in the CRAPII paper [14].
3) We combined the names of all the countries that belong to that culture into one list, filtering out names containing non-English characters to maintain consistency with the TSCC dataset [10].
4) Finally, the processed list of names for the given gender and cultural group was stored in a dictionary, with the group defined as a tuple of (`gender, culture`).

Once this mapping was complete, the 260 transcripts in the TSCC dataset were randomly assigned to the 10 gender-culture groups (2 gender by 5 culture groups), ensuring each group contained 26 transcripts. For each transcript, synthetic names were randomly sampled from the corresponding gender-culture group and used to replace placeholders. In addition to name placeholders, the dataset also contained 13 non-name placeholder categories (e.g., ⟨AGE⟩, ⟨DATE⟩, and ⟨INSTAGRAM ACCOUNT⟩). These were replaced with synthetic entities generated using GPT-4o to ensure semantic consistency throughout the dataset.

This process resulted in a realistic and culturally diverse dataset that retains the original conversational structure while introducing diversity in entity representation. For instance, the processed version of the transcript shown in Section III-A appears as follows:

```
teacher: Hi there John Doe, all OK?
student: Hi Jane Doe, how are you?
```

## IV. RESULTS AND DISCUSSION

The performance metrics of all proposed PII detection models are summarized in Table VII.

### A. Overall-Level Analysis of Model Performance

*1) Presidio Models:* For the two Presidio models utilizing different SpaCy configurations, *en_core_web_trf* consistently outperforms *en_core_web_lg* across all metrics, as shown in Table VII. The transformer-based *en_core_web_trf* achieves higher overall precision (0.2092 vs. 0.1505) and recall (0.8368 vs. 0.7103), likely due to its enhanced ability to capture long-range dependencies in text. However, both configurations exhibit low precision, which is likely due to the inclusive detection approach of the models. Although this approach enables the identification of a wide range of entities, including less common ones, it also leads to a significant number of false positives, thereby reducing overall precision.

[12]https://github.com/lukes/ISO-3166-Countries-with-Regional-Codes/blob/master/all/all.csv
[13]https://github.com/philipperemy/name-dataset

TABLE VII
PERFORMANCE METRICS FOR DIFFERENT PII DETECTION MODELS

| Models | Entity Type | # True Positive | # False Positive | # False Negative | Precision | Recall | $F_1$ Score | $F_5$ Score |
|---|---|---|---|---|---|---|---|---|
| *1. Presidio*<br>*(en_core_web_lg)* | NAME_STUDENT | 1,805 | 9,294 | 805 | 0.1626 | 0.6916 | 0.2633 | 0.6147 |
| | URL_PERSONAL | 181 | 2,256 | 31 | 0.0743 | 0.8538 | 0.1367 | 0.6082 |
| | EMAIL | 61 | 10 | 1 | 0.8592 | **0.9839** | 0.9173 | 0.9784 |
| | PHONE_NUM | 8 | 37 | 1 | 0.1778 | **0.8889** | 0.2963 | 0.7704 |
| | **Overall** | 2,055 | 11,597 | 838 | 0.1505 | 0.7103 | 0.2484 | 0.6214 |
| *2. Presidio*<br>*(en_core_web_trf)* | NAME_STUDENT | 2,172 | 6,849 | 438 | 0.2408 | 0.8322 | 0.3735 | 0.7604 |
| | URL_PERSONAL | 180 | 2,257 | 32 | 0.0739 | 0.8491 | 0.1359 | 0.6049 |
| | EMAIL | 61 | 10 | 1 | 0.8592 | **0.9839** | 0.9173 | 0.9784 |
| | PHONE_NUM | 8 | 37 | 1 | 0.1778 | **0.8889** | 0.2963 | 0.7704 |
| | **Overall** | 2,421 | 9,153 | 472 | 0.2092 | 0.8368 | 0.3347 | 0.7503 |
| *3. Azure AI Language* | NAME_STUDENT | 2,451 | 7,074 | 159 | 0.2573 | 0.9391 | 0.4040 | 0.8522 |
| | URL_PERSONAL | 145 | 917 | 67 | 0.1365 | 0.6840 | 0.2276 | 0.5926 |
| | EMAIL | 61 | 8 | 1 | **0.8841** | **0.9839** | **0.9313** | **0.9796** |
| | PHONE_NUM | 8 | 161 | 1 | 0.0473 | **0.8889** | 0.0899 | 0.5279 |
| | **Overall** | 2,665 | 8,160 | 228 | 0.2462 | 0.9212 | 0.3885 | 0.8333 |
| *4. Prompting*<br>*GPT-4o-mini* | NAME_STUDENT | 2,036 | 750 | 574 | 0.7308 | 0.7801 | 0.7546 | 0.7781 |
| | URL_PERSONAL | 153 | 313 | 59 | 0.3283 | 0.7217 | 0.4513 | 0.6899 |
| | EMAIL | 57 | 55 | 5 | 0.5089 | 0.9194 | 0.6552 | 0.8917 |
| | PHONE_NUM | 5 | 45 | 4 | 0.1000 | 0.5556 | 0.1695 | 0.4727 |
| | **Overall** | 2,251 | 1,163 | 642 | 0.6593 | 0.7781 | 0.7138 | 0.7727 |
| *5. Fine-tuned*<br>*GPT-4o-mini* | NAME_STUDENT | 2,507 | 1,597 | 103 | 0.6109 | **0.9605** | 0.7468 | **0.9398** |
| | URL_PERSONAL | 199 | 206 | 13 | 0.4914 | **0.9387** | 0.6451 | **0.9069** |
| | EMAIL | 60 | 10 | 2 | 0.8571 | 0.9677 | 0.9091 | 0.9630 |
| | PHONE_NUM | 8 | 4 | 1 | 0.6667 | **0.8889** | 0.7619 | 0.8776 |
| | **Overall** | 2,774 | 1,817 | 119 | 0.6042 | **0.9589** | 0.7413 | **0.9377** |
| *6. Verifier Model I*<br>*(Without CoT)* | NAME_STUDENT | 2,098 | 278 | 512 | **0.8830** | 0.8038 | **0.8416** | 0.8066 |
| | URL_PERSONAL | 161 | 2 | 51 | **0.9877** | 0.7594 | **0.8587** | 0.7662 |
| | EMAIL | 60 | 8 | 2 | 0.8824 | 0.9677 | 0.9231 | 0.9642 |
| | PHONE_NUM | 2 | 1 | 7 | 0.6667 | 0.2222 | 0.3333 | 0.2281 |
| | **Overall** | 2,321 | 289 | 572 | **0.8893** | 0.8023 | **0.8435** | 0.8053 |
| *7. Verifier Model II*<br>*(With CoT)* | NAME_STUDENT | 2,261 | 704 | 349 | 0.7626 | 0.8663 | 0.8111 | 0.8618 |
| | URL_PERSONAL | 173 | 74 | 39 | 0.7004 | 0.8160 | 0.7538 | 0.8109 |
| | EMAIL | 60 | 9 | 2 | 0.8696 | 0.9677 | 0.9160 | 0.9636 |
| | PHONE_NUM | 8 | 3 | 1 | **0.7273** | **0.8889** | **0.8000** | **0.8814** |
| | **Overall** | 2,502 | 790 | 391 | 0.7600 | 0.8648 | 0.8091 | 0.8603 |

*2) Azure AI Language:* The Azure AI Language model achieves an overall precision of 0.2462 and a recall of 0.9212, with the recall being the second highest across all models, slightly lower than the fine-tuned GPT-4o-mini model. Its strong recall highlights its ability to capture most true positives, reflected in the low number of false negatives (228). However, the low precision of the model, driven by a high number of false positives (8,160), limits its reliability for applications that require accurate predictions. The $F_1$ score of 0.3885 and $F_5$ score of 0.8333 further emphasize its recall-oriented nature, indicating that while Azure AI Language improves recall compared to rule-based methods, it struggles to maintain precision, resulting in an imbalanced trade-off.

*3) Prompting GPT-4o-mini:* The prompting GPT-4o-mini model achieves an overall precision of 0.6593 and recall of 0.7781, resulting in an $F_1$ score of 0.7138 and an $F_5$ score of 0.7727. Although the model demonstrates notable improvements in precision compared to rule-based approaches such as Presidio, its relatively low recall, as evidenced by the 642 false negatives, indicates that a significant number of true positives are missed. This limitation suggests that the model may not be ideal for contexts that require exhaustive PII detection. Despite these challenges, the improvement in precision highlights the potential of GPT-4o-mini's prompting capabilities, particularly

for scenarios where accuracy is prioritized over comprehensive detection.

*4) Fine-tuned GPT-4o-mini:* The fine-tuned GPT-4o-mini model demonstrates strong overall performance, achieving the highest recall among all models at 0.9589. This high recall ensures that nearly all PII entities are identified, making the model highly effective for comprehensive privacy protection. Its precision of 0.6042 represents a notable improvement over both the Presidio and Azure AI Language models, highlighting the benefits of fine-tuning in balancing precision and recall. The model achieves the highest $F_5$ score (0.9377) among all models, balancing high recall with reasonable precision. This highlights the potential of fine-tuning GPT-4o-mini for PII detection in educational texts, offering clear advantages over baseline and prompting models.

*5) Verifier Models:* The *Verifier Model I (Without CoT)* achieves the highest precision (0.8893) among all models by not defaulting to retaining entities when uncertain. However, this less conservative approach leads to the wrong removal of some true positives during verification, resulting in a lower recall of 0.8023. The *Verifier Model II (With CoT)* demonstrates a more balanced trade-off with a recall of 0.8648, as its conservative behavior defaults to retaining entities (**T**) when uncertainty arises. Notably, the precision scores for

both verifier models surpass that of all other five methods, aligning with our effort to improve precision. However, their recall is lower than that of the *Fine-tuned GPT-4o-mini* model, suggesting that these verifier models may be better suited for tasks that prioritize precision over exhaustive detection.

### B. PII Category-Level Analysis of Model Performance

*1) Name Detection (`NAME_STUDENT`):* The *Presidio* models demonstrate low precision (0.1626 and 0.2408) due to over-identifying common names, leading to a high number of false positives, but maintain moderate recall (0.6916 and 0.8322). The *Fine-tuned GPT-4o-mini* model achieves the highest recall (0.9605) and $F_5$ score (0.9398), making it the most reliable for comprehensive name detection, despite a moderate precision of 0.6109. The *Verifier Model I (Without CoT)* achieves the highest precision (0.8830) but sacrifices recall (0.8038). For names as a direct identifier, we recommend the *Fine-tuned GPT-4o-mini* model, given its superior recall and $F_5$ score.

*2) URL Detection (`URL_PERSONAL`):* The *Fine-tuned GPT-4o-mini* model achieves the highest recall (0.9387) and $F_5$ score (0.9069), indicating a strong performance in capturing true positives. However, its precision (0.4914) remains moderate, suggesting room for improvement in reducing false positives. The *Verifier Model I (Without CoT)* significantly improves precision, achieving the highest value (0.9877) and $F_1$ score (0.8587), but at the cost of reduced recall (0.7594). This demonstrates the verifier model's effectiveness in filtering false positives while highlighting the inherent trade-off between precision and recall, as no model achieves high performance in both metrics simultaneously.

*3) Email Detection (`EMAIL`):* All models demonstrate strong performance in detecting email entities, with recall values consistently high across both rule-based and GPT-based methods. The *Presidio* models and *Azure AI Language* achieve the highest recall (0.9839), missing only one true email entity out of 61, followed closely by the *Fine-tuned GPT-4o-mini* model and both *Verifier* models with a recall of 0.9677. In terms of precision, *Azure AI Language* achieves the highest value (0.8841), while the *Presidio* models (0.8592), *Fine-tuned GPT-4o-mini* model (0.8571), and *Verifier Model I (Without CoT)* (0.8824) also perform well. In general, email detection appears to be a relatively straightforward task, with most models achieving strong performance in both recall and precision, as reflected by their high $F_1$ and $F_5$ scores.

*4) Phone Number Detection (`PHONE_NUM`):* Five models, including both *Presidio* models, *Azure AI Language*, *Fine-tuned GPT-4o-mini*, and *Verifier Model II (With CoT)*, achieve the highest recall of 0.8889 for phone number detection. However, *Presidio* and *Azure AI Language* exhibit low precision, with Azure showing the lowest precision of 0.0473. In contrast, GPT-based models demonstrate higher precision, with the *Fine-tuned GPT-4o-mini* model reaching 0.6667. The *Verifier Model I (Without CoT)* records the lowest recall (0.2222), suggesting that it likely removed many true positives from the *Fine-tuned GPT-4o-mini* model's detected entities, possibly due to the lack of Chain-of-Thought reasoning. The *Verifier Model II (With CoT)* achieves the highest $F_1$ (0.8000) and $F_5$ (0.8814) scores, balancing strong precision and recall.

Overall, no single model dominates across all tested categories, as each exhibits distinct strengths. *Azure AI Language* performs best for email detection, while *Verifier Model II (With CoT)* is more effective for phone number detection. For names and URLs, the *Fine-tuned GPT-4o-mini* model and *Verifier Model I (Without CoT)* present a trade-off between recall and precision, allowing users to select a model based on the specific priorities of their task.

### C. Impact of Low-Precision PII Detection: Examples of Semantic Disruption

To better understand how low-precision PII detection can disrupt the semantic integrity of datasets and hinder downstream data analysis for educational research, we present three examples. These examples demonstrate cases where *Presidio* and *Azure AI Language* incorrectly identify non-PII entities as PII (false positives), resulting in unnecessary replacements that alter the intended meaning of the data. Such disruptions can negatively impact the utility of the data for educational insights and analysis. In contrast, all GPT-based models successfully identify these cases as non-PII (true negatives), thereby preserving the semantic meaning and ensuring the dataset's utility for downstream research tasks.

In Example 1, *Presidio* incorrectly identifies *Jesus Christ*, *Mary*, *Joseph*, and *Jesus* as PII. However, these names are not sensitive information in this context but are instead central to the story's historical and cultural narrative. The clue "Nazareth" is a key component of the story, as it is widely recognized as the hometown of Jesus Christ. Replacing the associated names with synthetic alternatives disrupts the educational purpose of the text, as students may no longer connect "Nazareth" to its religious significance. This could lead to misunderstandings and confusion in the learning process.

*Example 1*

**Original:** *At the beginning of the story, you do not know the names of the characters. Then at the end, I drop the first clue "Nazareth" - which is well known to be the home town of* *Jesus Christ. You can maybe guess that the family are* *Mary* *and* *Joseph* *with* *Jesus* *as a boy.*

**Replaced:** *At the beginning of the story, you do not know the names of the characters. Then at the end, I drop the first clue "Nazareth" - which is well known to be the home town of* *Elias Carson. You can maybe guess that the family are* *Lena* *and* *Daniel* *with* *Elijah* *as a boy.*

The incorrect anonymization of names changes the intended meaning of the text, as the connection between "Nazareth" and its historical and religious significance is lost. This disruption affects students' understanding and the utility of the data in educational contexts. Moreover, anonymization negatively impacts the utility of the text for machine learning applications. If the anonymized text is used to train or fine-tune a language model or included in a knowledge base for a Retrieval-Augmented Generation (RAG) pipeline, it may result in incorrect or misleading responses to queries about "Nazareth" [37]. Repeated inclusion of anonymized instances, such as associating "Nazareth" with unrelated names like

"Elias Carson," may erode the model's understanding of its cultural significance, leading to inaccuracies in subsequent applications [38], [39].

Similarly, *Azure AI Language* exhibits similar challenges when handling famous individuals. In Example 2, *Azure AI Language* identifies notable figures such as *Bill Gates*, *Steve Jobs*, *Zuckerberg*, and *Elon Musk* as PII entities. The subsequent replacement with random surrogate names strips the text of its unique context and relevance. These figures' specific ages, entrepreneurial journeys, and market contexts are integral to the narrative, explaining why their stories "*didn't translate well*" into the students' environment.

*Example 2*
**Original:** *It became clear that while the students were excited about setting up and running startup companies on campus, they had very little background information to do so. Their role models came from the other part of the world namely Bill Gates, Steve Jobs, Zuckerberg, Elon Musk etc. Their stories or anecdotes didn't translate well into the environment of our students.*

**Replaced:** *It became clear that while the students were excited about setting up and running startup companies on campus, they had very little background information to do so. Their role models came from the other part of the world namely Benjamin Bloom, Stephen Jackson, Oliver Underwood, Aiden Miles etc. Their stories or anecdotes didn't translate well into the environment of our students.*

Another common false positive pattern in both *Presido* and *Azure* is de-identifying the end and start of two consecutive sentences (with no whitespace) as a URL, illustrated by the example below:

*Example 3*
**Original:** *Consider hiring a copywriter to craft a compelling menu.Keep menus clean – no grease and no food or water stains. Get rid of worn or torn menus.Update menu and prices at least once a year.Build menu around popular items.*

**Replaced:** *Consider hiring a copywriter to craft a compelling https://techwaveinsight.ioeep menus clean – no grease and no food or water stains. Get rid of worn or torn menus.Update menu and prices at least once a http://elitecodingacademy.org menu around popular items.*

Replacing key terms such as "*menu*," "*keep*," "*year*," and "*Build*" with arbitrary URLs disrupts the original meaning. The instructions lose clarity, merging steps into a single ill-structured, confusing sentence.

### D. Cost Analysis

Table VIII presents the cost associated with each PII detection model, complementing the performance metrics in Table VII. The results emphasize the potential of GPT-based approaches, particularly the *Fine-tuned GPT-4o-mini* model, for high-quality PII detection at significantly lower costs.

The *Fine-tuned GPT-4o-mini* model achieves the highest overall recall (0.9589) and $F_5$ score (0.9377), outperforming both the free *Presidio* models and the expensive *Azure AI*

*Language* model. While the *Presidio* models incur no cost, their low precision (0.1505 and 0.2092) and $F_5$ scores (0.6214 and 0.7503) highlight their limitations in balancing false positives and true positives. In contrast, the *Azure AI Language* model, though more precise, costs \$63.27 (\$4.90 per 1M tokens), which is approximately 6 times higher than the *Fine-tuned GPT-4o-mini* model (\$0.92 per 1M tokens). Then, the *Verifier models* marginally increase the total cost to \$13.97 (Without CoT) and \$14.69 (With CoT), enhancing precision for applications where minimizing false positives is critical. Despite the additional cost, these models remain far more economical than *Azure AI Language* model while retaining GPT's high recall and semantic accuracy. Thus, GPT-based models, led by the *Fine-tuned GPT-4o-mini* model, outperform both *Presidio* and *Azure AI Language* in balancing cost and performance. This underscores their potential as an efficient and scalable solution for PII detection in educational data.

### E. Name Culture and Gender Bias Analysis

Based on the detection performance in Table VII, we selected three models—*Presidio with en_core_web_trf*, *Azure AI Language*, and *Fine-tuned GPT-4o-mini*—for the analysis of cultural and gender bias in name detection. The total number of entities for each gender and culture group and their recall are shown in Table IX.

**Gender Analysis.** The results indicate that the recall scores between male and female names across the three models are similar. For *Presidio* and *Azure AI Language*, there is a marginally higher recall for female names compared to male names. Specifically, *Azure AI Language* achieves a recall of 0.9541 for female names and 0.9494 for male names. In contrast, the fine-tuned GPT-4o-mini model performs slightly better on male names (0.9646) compared to female names (0.9591). However, the differences in recall for male and female names are minimal, suggesting consistent performance between gender groups in all models.

**Culture Analysis.** The cultural analysis reveals more significant differences in model performance. Both Microsoft models show lower recall for African and Asian names. In particular, *Presidio* exhibits a notable performance gap, with recall rates of 0.7647 for African names and 0.8640 for Asian names, compared to 0.9024 for European and 0.9091 for American names. *Azure AI Language* also shows a lower recall for African names (0.9244), although the gap is less pronounced than in *Presidio*. These results suggest inherent biases in the Microsoft models, possibly stemming from imbalances in the training data or gaps in cultural representation. In contrast, the fine-tuned GPT-4o-mini model achieves consistent and high recall across all cultural groups. It performs with recall scores of 0.9756 for European names, 0.9790 for American names, 0.9840 for Asian names, and 0.9748 for African names. This minimal variation in recall across cultural groups demonstrates the model's ability to mitigate cultural bias and generalize effectively across diverse name distributions.

Therefore, the gender analysis does not show significant differences in recall for male and female names across all models, indicating consistent performance in this aspect. However, cultural analysis reveals that both *Presidio* and *Azure*

TABLE VIII
COST BREAKDOWN FOR PII DETECTION MODELS (IN USD)

| Models | Base Fine-tuning | Base Model Dependency | Verifier Training Data Construction | Verifier Fine-tuning | Evaluation | Total | Average (per 1M tokens) |
|---|---|---|---|---|---|---|---|
| *Presidio (en_core_web_lg)* | | — | | | 0 | 0 | 0 |
| *Presidio (en_core_web_trf)* | | — | | | 0 | 0 | 0 |
| *Azure AI Language* | | — | | | $63.27 | $63.27 | $4.90 |
| *Prompting GPT-4o-mini* | | — | | | $5.22 | $5.22 | $0.40 |
| *Fine-tuned GPT-4o-mini* | $7.22 | | — | | $4.71 | $11.93 | $0.92 |
| *Verifier Model I (Without CoT)* | — | $11.93 | $1.16 | $0.80 | $0.08 | $13.97 | $1.09 |
| *Verifier Model II (With CoT)* | — | $11.93 | $1.50 | $1.08 | $0.18 | $14.69 | $1.13 |

TABLE IX
RECALL COMPARISON ACROSS GENDER AND CULTURE GROUPS FOR
SELECTED MODELS.

| Type | Group | Total | Presidio | Azure AI | GPT-4o-mini |
|---|---|---|---|---|---|
| Gender | Male | 1582 | 0.8786 | 0.9494 | 0.9646 |
| | Female | 1002 | 0.8832 | 0.9541 | 0.9591 |
| Culture | Europe | 410 | 0.9024 | 0.9585 | **0.9756** |
| | Americas | 858 | 0.9091 | 0.9755 | **0.9790** |
| | Asia | 500 | 0.8640 | 0.9320 | **0.9840** |
| | Africa | 238 | **0.7647** | 0.9244 | **0.9748** |

*AI Language* exhibit performance gaps for African and Asian names, with lower recall for these groups. The *fine-tuned GPT-4o-mini* model, on the other hand, performs generally well across all cultural groups, effectively addressing the bias observed in the baseline models and highlighting its superior generalization capability.

### F. Generalizability Analysis

To further evaluate the performance of different models for PII detection, we further used the TSCC dataset (as introduced in Section III-A) to evaluate the models that were implemented on the CRAPII dataset. Table X presents the performance for different PII detection models on the TSCC dataset. The metrics are the same as those presented in Table VII.

We selected one Presidio model *Presidio (en_core_web_trf)* model as it demonstrated higher precision and recall compared to the *en_core_web_lg* variant (Model 1 from Table VII). As shown in Table X, the *Presidio (en_core_web_trf) model* demonstrates relatively high recall (0.9368) and low precision (0.3596). This is likely due to its conservative approach to entity detection, which enables it to capture a wide range of entities, including rare or unconventional names. This conservatism results in a large number of false positives (2,694), the highest among all six models, which significantly impacts its precision.

Then, *Azure AI Language* was chosen as another baseline to compare with both rule-based methods and LLM-based models. The *Azure AI Language* model exhibits relatively balanced precision (0.5008) and recall (0.8173). However, its overall performance is suboptimal compared to the GPT-based approaches. With 1,316 false positives and 295 false negatives, the model fails to achieve the level of precision or recall seen in other approaches. This is reflected in its moderate $F_1$ score of 0.6210 and $F_5$ score of 0.7979.

Next, we investigated the fine-tuned GPT-4o-mini model which was fine-tuned on the CRAPII dataset. To gain a better understanding of the GPT models on PII detection on TCSS dataset. We used a three-shot prompting strategy to directly prompted GPT-4o-mini to identify PII entities in the TSCC dataset. We also prompted the *Fine-tuned GPT-4o-mini* model (fine-tuned on the CRAPII dataset, Model 5 in Table VII) without examples to assess its generalizability on the unseen TSCC dataset. Both models employed the same prompt structure presented in Table III, adjusting the user prompt to only label names with special identifiers as names were the only PII category in the dataset. The results in Table X show that directly prompting GPT-4o-mini with few-shot examples achieves the highest recall among all models at 0.9932, demonstrating the model's capability to detect almost all true PII entities in the dataset. However, its precision is relatively low at 0.7145 due to the higher number of false positives (641). While this approach achieves a strong $F_5$ score of 0.9785, the $F_1$ score of 0.8311 indicates that the lower precision affects its overall performance. Then, prompting the *Fine-tuned GPT-4o-mini* model (on the CRAPII dataset) results in the highest precision across all models (0.9984), with only two false positives. This indicates that the entities it detects are highly accurate. However, the recall is relatively low at 0.7882, leading to a $F_1$ score of 0.8810 and a lower $F_5$ score of 0.7947. This approach is particularly suited for scenarios where precision is more critical than recall.

We also provided users with a practical option to improve model performance by fine-tuning on a minimal labeled subset of their dataset, enabling the model to better align with the specific characteristics of the TSCC dataset. Of the 260 transcripts in the processed TSCC dataset, we randomly selected 10 transcripts for fine-tuning purposes introduced below and used the remaining 250 transcripts for evaluation across all models. First, we directly fine-tuned GPT-4o-mini on the 10 selected transcripts. Second, we further fine-tuned the *Fine-tuned GPT-4o-mini* model (Model 5 from Table VII) using the same 10 transcripts. Both fine-tuning models used the prompt structure described in Table IV. As with prompting, the only adjustment was to instruct the model to label names using @@@ and ###. Fine-tuning GPT-4o-mini on the 10 selected TSCC transcripts achieves a strong balance between precision (0.9836) and recall (0.9666). With only 26 false positives and 54 false negatives, this model achieves a $F_1$ score of 0.9750 and a $F_5$ score of 0.9672. This demonstrates the potential of fine-tuning even on a small labeled subset to adapt the

TABLE X
PERFORMANCE METRICS FOR DIFFERENT PII DETECTION APPROACHES ON THE TSCC DATASET.

| Models | # TP | # FP | # FN | Precision | Recall | $F_1$ Score | $F_5$ Score |
|---|---|---|---|---|---|---|---|
| 1. Presidio (en_core_web_trf) | 1,513 | 2,694 | 102 | 0.3596 | 0.9368 | 0.5198 | 0.8824 |
| 2. Azure AI Language | 1,320 | 1,316 | 295 | 0.5008 | 0.8173 | 0.6210 | 0.7979 |
| 3. GPT-4o-mini + few-shot prompting | 1,604 | 641 | 11 | 0.7145 | **0.9932** | 0.8311 | 0.9785 |
| 4. Fine-tuned GPT-4o-mini + zero-shot prompting | 1,273 | 2 | 342 | **0.9984** | 0.7882 | 0.8810 | 0.7947 |
| 5. GPT-4o-mini + fine-tuning | 1,561 | 26 | 54 | 0.9836 | 0.9666 | 0.9750 | 0.9672 |
| 6. Fine-tuned GPT-4o-mini + fine-tuning | 1,598 | 48 | 17 | 0.9708 | 0.9895 | **0.9801** | **0.9887** |

model to new datasets effectively. Then, further fine-tuning of the *Fine-tuned GPT-4o-mini* model (on the CRAPII dataset) using the 10 selected transcripts results in the best overall performance. This model achieves a high recall of 0.9895 and a precision of 0.9708, striking an excellent balance between the two metrics. Although its recall is slightly lower (about 0.3%) than the few-shot prompting model, its precision improves significantly by approximately 26%. This leads to the highest $F_1$ score (0.9801) and the $F_5$ score (0.9887) among all models, making it the most robust approach to PII detection on the TSCC dataset.

Therefore, the *Fine-tuned GPT-4o-mini + fine-tuning* approach achieves the highest $F_1$ and $F_5$ scores, highlighting its ability to maintain high precision and recall simultaneously. These results also demonstrate the potential of GPT-based approaches over traditional models such as *Presidio* and *Azure*, offering superior performance in both precision and recall. While other approaches, such as few-shot prompting with GPT-4o-mini, excel in recall, the overall balance achieved by fine-tuning makes it a versatile option. Users can choose the most appropriate method based on their specific requirements, whether to prioritize recall, precision, or a combination of both.

## V. LIMITATIONS AND FUTURE WORKS

While our study highlights the potential of fine-tuned GPT-4o-mini models for cost-effective and accurate PII detection in educational texts, there are opportunities for further refinement and exploration. *First*, our evaluation primarily addressed PII categories that are universally recognized by widely used baseline methods (e.g., *Microsoft Presidio* and *Azure AI Language*). Extending this approach to more granular or domain-specific categories, such as addresses, student ID numbers, or indirect identifiers like dates, times, and ages, could benefit from additional training data and tailored customization of LLM-based methods. *Second*, while we focused on Microsoft *Presidio* and *Azure AI Language* for benchmarking due to their prevalence, incorporating other models or framework may provide deeper insights into the scalability and versatility of our approach. *Third*, although we conducted a preliminary generalizability test on the TSCC dataset, the diverse nature of educational contexts suggests that some unseen PII distributions may still present challenges. A more comprehensive evaluation across a broader range of datasets could provide further clarity on the model's adaptability. *Fourth*, while the Verifier models

we developed significantly enhance precision, their trade-off in recall suggests potential for further innovation, such as multi-stage verification strategies that integrate recall-oriented detection with context-driven verification. *Finally*, although this study prioritized optimizing detection and verification processes, exploring advanced red-teaming or adversarial testing could further bolster model robustness and resilience against potential vulnerabilities. These extensions would help ensure the security and reliability of PII detection systems across the broad and varied landscape of educational data.

## VI. CONCLUSION

This study highlights the cost-effectiveness and strong performance of fine-tuned GPT-4o-mini for identifying personal information in educational texts. Compared to rule-based (*Presidio*) and cloud-based (*Azure AI Language*) methods, the fine-tuned GPT-4o-mini model achieves notably high recall (more than 0.95) while maintaining solid precision. This balance in performance is critical for education, where missing any PII can pose privacy risks, and unneeded removals reduce research and pedagogical utility.

Verifier models further enhance precision by verifying detected entities in context, helping users preserve semantic meaning in approaches such as Hidden in Plain Sight. Our cost analysis shows that GPT-4o-mini is at least 4 times cheaper and thus much more affordable than commonly used commercial services, lowering barriers for privacy-preserving research in education. Testing on an unseen teacher-student corpus confirms that GPT-based models generalize well, even with limited labeled data for fine-tuning.

In general, the fine-tuned GPT-4o-mini model offered the most consistent overall performance and cultural generalizability. These findings suggest that carefully refined LLMs can safeguard privacy while preserving meaningful content for learning analytics and educational research at scale.

## REFERENCES

[1] J. R. Reidenberg and F. Schaub, "Achieving big data privacy in education," *Theory and Research in Education*, vol. 16, no. 3, pp. 263–279, 2018. [Online]. Available: https://doi.org/10.1177/1477878518805308

[2] L. Portnoff, E. Gustafson, J. Rollinson, and K. Bicknell, "Methods for language learning assessment at scale: Dulingo case study," in *Proceedings of The 14th International Conference on Educational Data Mining (EDM21)*. Paris, France: International Educational Data Mining Society, 2021, pp. 865–871. [Online]. Available: https://educationaldatamining.org/edm2021/

[3] J. Chen, J. H. Fife, I. I. Bejar, and A. A. Rupp, "Building e-rater® scoring models using machine learning methods," *Educational Testing Service*, 2023. [Online]. Available: https://doi.org/10.1002/ets2.12094

[4] J. Lin, E. Chen, Z. Han, A. Gurung, D. R. Thomas, W. Tan, N. D. Nguyen, and K. R. Koedinger, "How can i improve? using gpt to highlight the desired and undesired parts of open-ended responses," in *Proceedings of the 17th International Conference on Educational Data Mining*, B. PaaÃŸen and C. D. Epp, Eds. Atlanta, Georgia, USA: International Educational Data Mining Society, July 2024, pp. 236–250.

[5] B. C. M. Fung, K. Wang, R. Chen, and P. S. Yu, "Privacy-preserving data publishing: A survey of recent developments," *ACM Computing Surveys (CSUR)*, vol. 42, no. 4, pp. 1–53, June 2010. [Online]. Available: https://doi.org/10.1145/1749603.1749605

[6] D. Carrell, B. Malin, J. Aberdeen, S. Bayer, C. Clark, B. Wellner, and L. Hirschman, "Hiding in plain sight: use of realistic surrogates to reduce exposure of protected health information in clinical text," *Journal of the American Medical Informatics Association*, vol. 20, no. 2, pp. 342–348, 07 2012. [Online]. Available: https://doi.org/10.1136/amiajnl-2012-001034

[7] E. Zeide, "Education technology and student privacy," in *The Cambridge Handbook of Consumer Privacy*, E. Selinger, J. Polonetsky, and O. Tene, Eds. Cambridge, MA: The MIT Press, 2018, pp. 70–84. [Online]. Available: https://ssrn.com/abstract=3145634

[8] "Family Educational Rights and Privacy Act (FERPA)," 1974, 20 U.S.C. § 1232g. [Online]. Available: https://www.govinfo.gov/app/details/USCODE-2010-title20/USCODE-2010-title20-chap31-subchapIII-part4-sec1232g

[9] Y. Liu, D. Iter, Y. Xu, S. Wang, R. Xu, and C. Zhu, "G-eval: Nlg evaluation using gpt-4 with better human alignment," 2023. [Online]. Available: https://arxiv.org/abs/2303.16634

[10] A. Caines, H. Yannakoudakis, H. Allen, P. Pérez-Paredes, B. Byrne, and P. Buttery, "The teacher-student chatroom corpus version 2: more lessons, new annotation, automatic detection of sequence shifts," in *Proceedings of the 11th Workshop on NLP for Computer Assisted Language Learning*, D. Alfter, E. Volodina, T. François, P. Desmet, F. Cornillie, A. Jönsson, and E. Rennes, Eds. Louvain-la-Neuve, Belgium: LiU Electronic Press, Dec. 2022, pp. 23–35. [Online]. Available: https://aclanthology.org/2022.nlp4call-1.3

[11] A. Pal, R. Bhargava, K. Hinsz, J. Esterhuizen, and S. Bhattacharya, "The empirical impact of data sanitization on language models," 2024, paper accepted at Safe Generative AI Workshop at NeurIPS 2024. [Online]. Available: https://arxiv.org/abs/2411.05978

[12] P. Lison, I. Pilán, D. Sanchez, M. Batet, and L. Øvrelid, "Anonymisation models for text data: State of the art, challenges and future directions," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Online: Association for Computational Linguistics, 2021, pp. 4188–4203. [Online]. Available: https://aclanthology.org/2021.acl-long.323

[13] P. Samarati, "Protecting respondents identities in microdata release," *IEEE Transactions on Knowledge and Data Engineering*, vol. 13, no. 6, pp. 1010–1027, 2001.

[14] L. Holmes, S. Crossley, J. Wang, and W. Zhang, "The cleaned repository of annotated personally identifiable information," in *Proceedings of the 17th International Conference on Educational Data Mining*, B. PaaÃŸen and C. D. Epp, Eds. Atlanta, Georgia, USA: International Educational Data Mining Society, July 2024, pp. 790–796.

[15] E. F. T. K. Sang and F. D. Meulder, "Introduction to the conll-2003 shared task: Language-independent named entity recognition," in *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*. NAACL, 2003, pp. 142–147.

[16] D. S. Carrell, D. J. Cronkite, M. R. Li, S. Nyemba, B. A. Malin, J. S. Aberdeen, and L. Hirschman, "The machine giveth and the machine taketh away: a parrot attack on clinical text deidentified with hiding in plain sight," *Journal of the American Medical Informatics Association*, vol. 26, no. 12, pp. 1536–1544, 08 2019. [Online]. Available: https://doi.org/10.1093/jamia/ocz114

[17] S. Singhal, A. F. Zambrano, M. Pankiewicz, X. Liu, C. Porter, and R. S. Baker, "De-Identifying Student Personally Identifying Information with GPT-4," in *Proceedings of the 17th International Conference on Educational Data Mining*. International Educational Data Mining Society, Jul. 2024, pp. 559–565. [Online]. Available: https://doi.org/10.5281/zenodo.12729884

[18] OpenAI, "Openai api pricing," https://openai.com/api/pricing/, accessed: 2025-01-11.

[19] S. Samsi, D. Zhao, J. McDonald, B. Li, A. Michaleas, M. Jones, W. Bergeron, J. Kepner, D. Tiwari, and V. Gadepally, "From words to watts: Benchmarking the energy costs of large language model inference," 2023. [Online]. Available: https://arxiv.org/abs/2310.03003

[20] J. D. Osborne, T. O'Leary, A. Nadimpalli, S. M. Aly., and R. E. Kennedy, "Bratsynthetic: Text de-identification using a markov chain replacement strategy for surrogate personal identifying information," 2022. [Online]. Available: https://arxiv.org/abs/2210.16125

[21] D. Sánchez and M. Batet, "C-sanitized: A privacy model for document redaction and sanitization," *Journal of the Association for Information Science and Technology*, vol. 67, no. 1, p. 148–163, Apr. 2015. [Online]. Available: http://dx.doi.org/10.1002/asi.23363

[22] D. S. Carrell, B. A. Malin, D. J. Cronkite, J. S. Aberdeen, C. Clark, M. R. Li, D. Bastakoty, S. Nyemba, and L. Hirschman, "Resilience of clinical text de-identified with "hiding in plain sight" to hostile reidentification attacks by human readers," *J. Am. Med. Inform. Assoc.*, vol. 27, no. 9, pp. 1374–1382, Jul 2020.

[23] T. Carvalho, N. Moniz, P. Faria, and L. Antunes, "Survey on privacy-preserving techniques for data publishing," 2022. [Online]. Available: https://arxiv.org/abs/2201.08120

[24] L. Holmes, S. Crossley, N. Hayes, D. Kuehl, A. Trumbore, and G. Gutu-Robu, "De-identification of student writing in technologically mediated educational settings," in *Polyphonic Construction of Smart Learning Ecosystems*, M. Dascalu, P. Marti, and F. Pozzi, Eds. Singapore: Springer Nature Singapore, 2023, pp. 177–189.

[25] S. Wang, X. Sun, X. Li, R. Ouyang, F. Wu, T. Zhang, J. Li, and G. Wang, "Gpt-ner: Named entity recognition via large language models," 2023. [Online]. Available: https://arxiv.org/abs/2304.10428

[26] X. Zhang, N. Rajabi, K. Duh, and P. Koehn, "Machine translation with large language models: Prompting, few-shot learning, and fine-tuning with qlora," in *Proceedings of the Eighth Conference on Machine Translation*. Singapore: Association for Computational Linguistics, 2023, pp. 468–481.

[27] M. Mosbach, T. Pimentel, S. Ravfogel, D. Klakow, and Y. Elazar, "Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation," in *Findings of the Association for Computational Linguistics: ACL 2023*. Toronto, Canada: Association for Computational Linguistics, 2023, pp. 12 284–12 314.

[28] Z. Shen, Z. Liu, J. Qin, M. Savvides, and K.-T. Cheng, "Partial is better than all: Revisiting fine-tuning strategy for few-shot learning," 2021, aAAI 2021. A search based fine-tuning strategy for few-shot learning. [Online]. Available: https://arxiv.org/abs/2102.03983

[29] S. Lee, S. Seo, J. Kim, Y. Lee, and S. Hwang, "Few-shot fine-tuning is all you need for source-free domain adaptation," 2023, the first two authors contributed equally. [Online]. Available: https://arxiv.org/abs/2304.00792

[30] S. Y. Gadre, G. Smyrnis, V. Shankar, S. Gururangan, M. Wortsman, R. Shao, J. Mercat, A. Fang, J. Li, S. Keh, R. Xin, M. Nezhurina, I. Vasiljevic, J. Jitsev, L. Soldaini, A. G. Dimakis, G. Ilharco, P. W. Koh, S. Song, T. Kollar, Y. Carmon, A. Dave, R. Heckel, N. Muennighoff, and L. Schmidt, "Language models scale reliably with over-training and on downstream tasks," 2024. [Online]. Available: https://arxiv.org/abs/2403.08540

[31] K. Tirumala, A. H. Markosyan, L. Zettlemoyer, and A. Aghajanyan, "Memorization without overfitting: Analyzing the training dynamics of large language models," 2022. [Online]. Available: https://arxiv.org/abs/2205.10770

[32] T. Ball, S. Chen, and C. Herley, "Can we count on llms? the fixed-effect fallacy and claims of gpt-4 capabilities," 2024. [Online]. Available: https://arxiv.org/abs/2409.07638

[33] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," 2020. [Online]. Available: https://arxiv.org/abs/2005.14165

[34] D. Brandfonbrener, S. Henniger, S. Raja, T. Prasad, C. Loughridge, F. Cassano, S. R. Hu, J. Yang, W. E. Byrd, R. Zinkov, and N. Amin, "Vermcts: Synthesizing multi-step programs using a verifier, a large language model, and tree search," 2024. [Online]. Available: https://arxiv.org/abs/2402.08147

[35] Y. Li, Z. Lin, S. Zhang, Q. Fu, B. Chen, J.-G. Lou, and W. Chen, "Making large language models better reasoners with step-aware verifier," 2023. [Online]. Available: https://arxiv.org/abs/2206.02336

[36] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," 2023. [Online]. Available: https://arxiv.org/abs/2201.11903

[37] F. Wang, X. Wan, R. Sun, J. Chen, and S. Ö. Arık, "Astute rag: Overcoming imperfect retrieval augmentation and knowledge conflicts for large language models," 2024. [Online]. Available: https://arxiv.org/abs/2410.07176

[38] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, "On the dangers of stochastic parrots: Can language models be too big?" FAccT '21: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, pp. 610–623, 2021.

[39] F.-K. Sun, C.-H. Ho, and H.-Y. Lee, "Lamol: Language modeling for lifelong language learning," 2019. [Online]. Available: https://arxiv.org/abs/1909.03329

**Jionghao Lin** is an Assistant Professor in the Faculty of Education at the University of Hong Kong. Previously, he was a Postdoctoral Researcher at the Human-Computer Interaction Institute at Carnegie Mellon University, Pittsburgh, PA, USA, from 2023 to 2025. He received his Ph.D. in Computer Science from Monash University, Clayton, VIC, Australia, in 2023. His research interests include learning science, natural language processing, data mining, and applications of generative artificial intelligence in education. His work has been published in international journals and conferences and recognized with awards, including the Best Paper Award at HCII'24, EDM'24, and ICMI'19, and the Best Demo Award at AIED'23.

**Kenneth R. Koedinger** is a professor of Human Computer Interaction and Psychology at Carnegie Mellon University. Dr. Koedinger has an M.S. in Computer Science, a Ph.D. in Cognitive Psychology, and experience teaching in an urban high school. His multidisciplinary background supports his research goals of understanding human learning and creating educational technologies that increase student achievement. His research has contributed new principles and techniques for the design of educational software and has produced basic cognitive science research results on the nature of student thinking and learning. Koedinger directs LearnLab, which started with 10 years of National Science Foundation funding and is now the scientific arm of CMU's Simon Initiative. LearnLab builds on the past success of Cognitive Tutors, an approach to online personalized tutoring that is in use in thousands of schools and has been repeatedly demonstrated to increase student achievement, for example, doubling what algebra students learn in a school year. He was a co-founder of CarnegieLearning, Inc. that has brought Cognitive Tutor based courses to millions of students since it was formed in 1998, and leads LearnLab, now the scientific arm of CMU's Simon Initiative. Dr. Koedinger has authored over 250 peer-reviewed publications and has been a project investigator on over 45 grants. In 2017, he received the Hillman Professorship of Computer Science and in 2018, he was recognized as a fellow of Cognitive Science.

**Yuntian Shen** is an undergraduate senior majoring in Statistics and Machine Learning with a double major in Artificial Intelligence at the Dietrich College of Humanities and Social Sciences, Carnegie Mellon University, Pittsburgh, PA, USA. He is currently a Research Assistant at the Human-Computer Interaction Institute at CMU. His research interests include generative AI, ethical and privacy considerations in AI, AI in healthcare, and the development of world models.

**Zilyu Ji** is an undergraduate junior majoring in Artificial Intelligence at the School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA. His research interests include Natural Language Processing and Large Language Models.