

# Identifying Personal Identifiable Information (PII) in Student Essays

Manu Achar | Pratik Chaudhari | Cody Ledford | Advised by Dr. Anca Doloc-Mihu | Georgia Gwinnett College

## INTRODUCTION

The personally identifiable information (PII) detection competition is hosted by The Learning Agency Lab on Kaggle, and by the request of our client, Dr. Gunay, we aim to detect sensitive PII, like names or emails, in student writing. This is necessary to screen and clean educational data so that when released to the public for analysis and archival, the students' privacy risks are mitigated.

## OBJECTIVES

- Understand the distribution of PII in student essays
- Detect PII
- Report showing its location

## TECHNOLOGIES

Taipy, PyTorch, Pre-trained BERT, Jupyter Notebook

## RESULT 1: PREPROCESSING

Transforming data to BERT-friendly structure

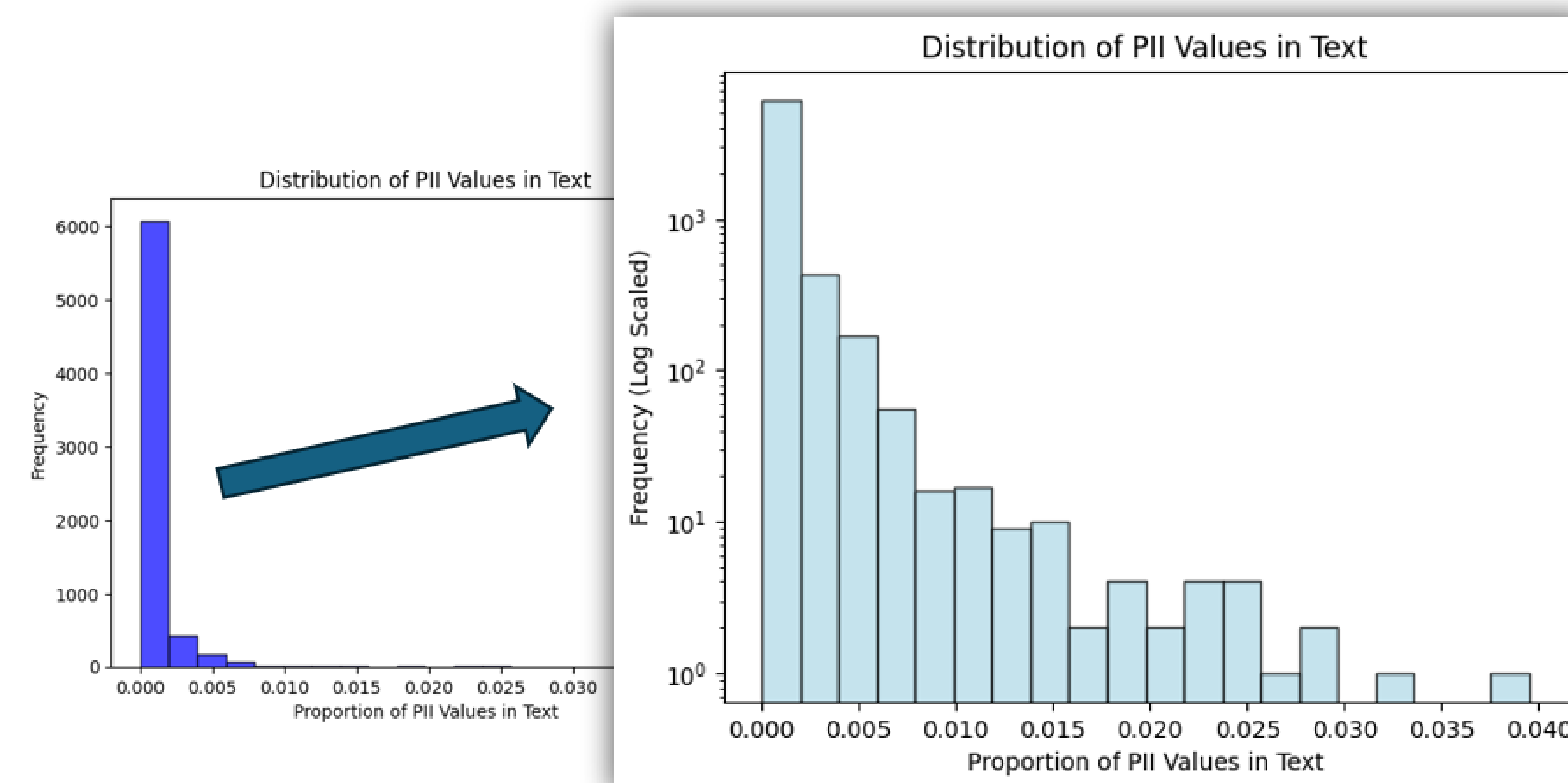
For Example:

	1229,	9681,	4286,	1110,	1304,	12196,
'Design',	2000,	1105,	1159,	1106,	2437,	1120,
'Thinking',	13032,	1195,	1460,	119,	1573,	1195,
'for',	1110,	1103,	1211,	5806,	6806,	1106,
'innovation',	1104,	1103,	1933,	119,	102,	0,
'reflexion',	0,	0,	0,	0,	0,	0,
'-',	0,	0,	0,	0,	0,	0,
'Avril',	0,	0,	0,	0,	0,	0,
'2021',	0,	0,	0,	0,	0,	0,
'.	0,	0,	0,	0,	0,	0,

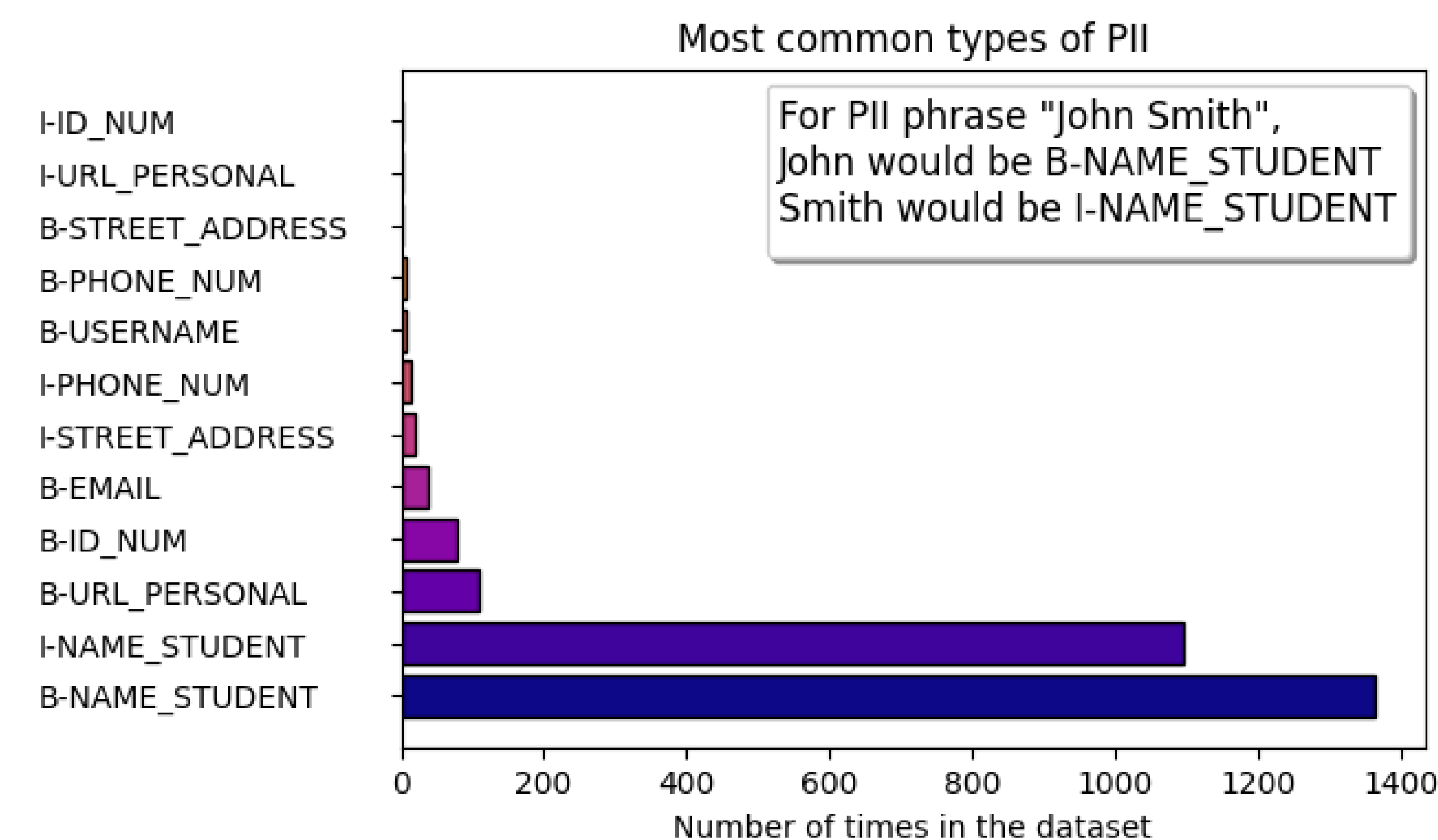
Tokenization

## RESULT 2: STATISTICS

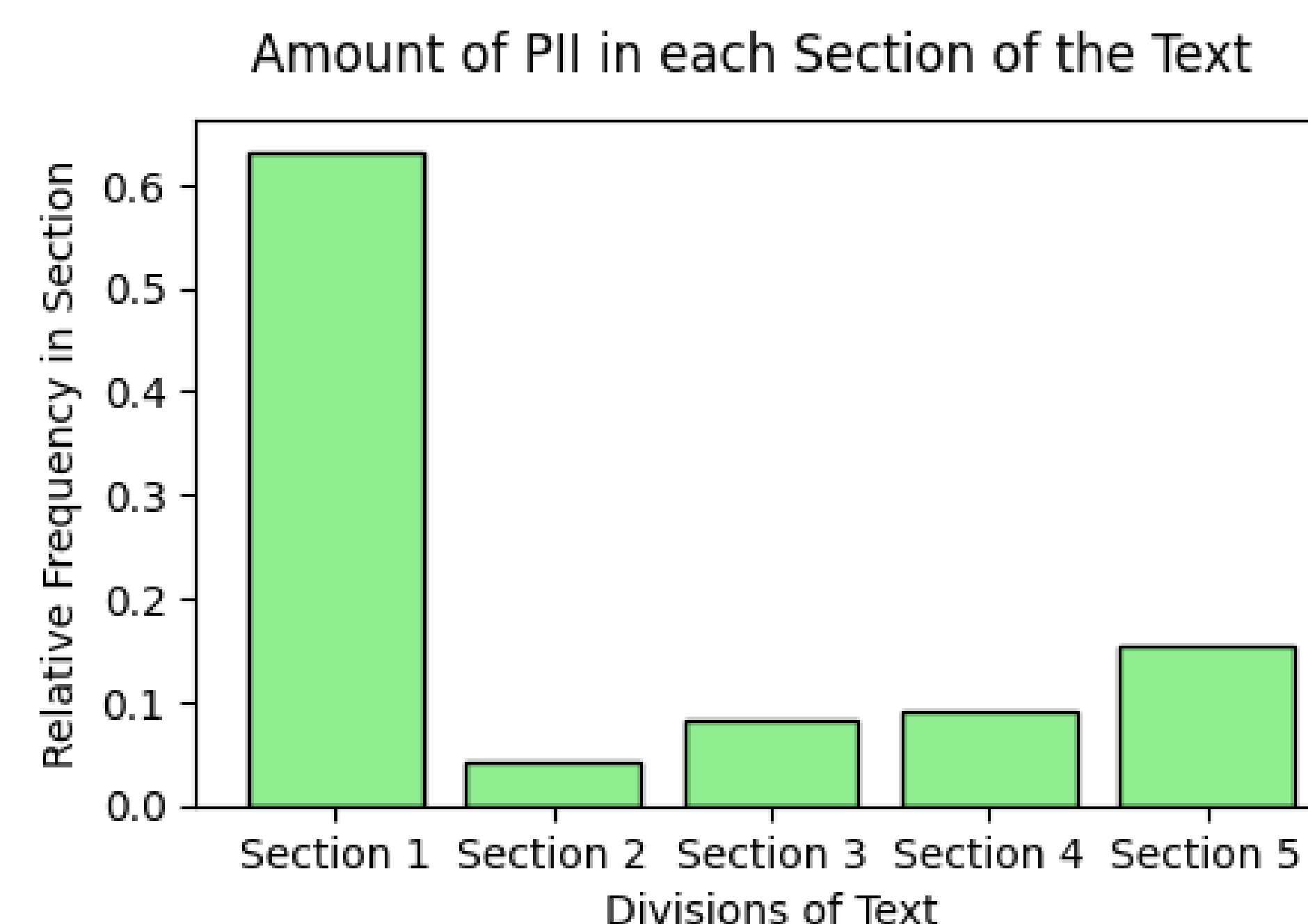
### Distribution of PII Values in Text



### Most common type of PII are names

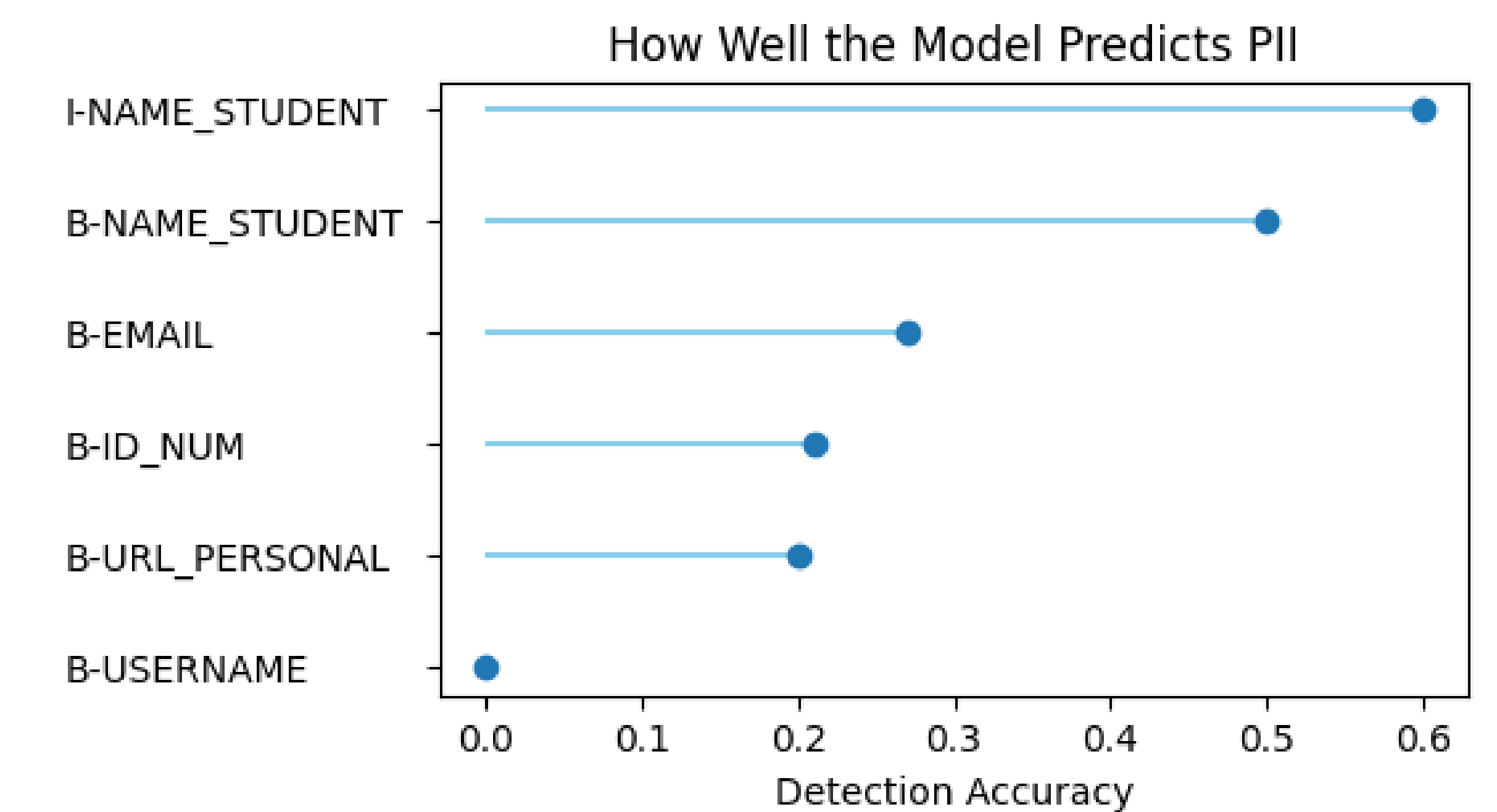


### Most PII occurs in the first 1/5th of the essay



## RESULTS 3: MODEL

- Token Classification by fine-tuning BERT model
- 5000 Training essays, 1000 Testing essays
- Model Accuracy: 56%, still improving

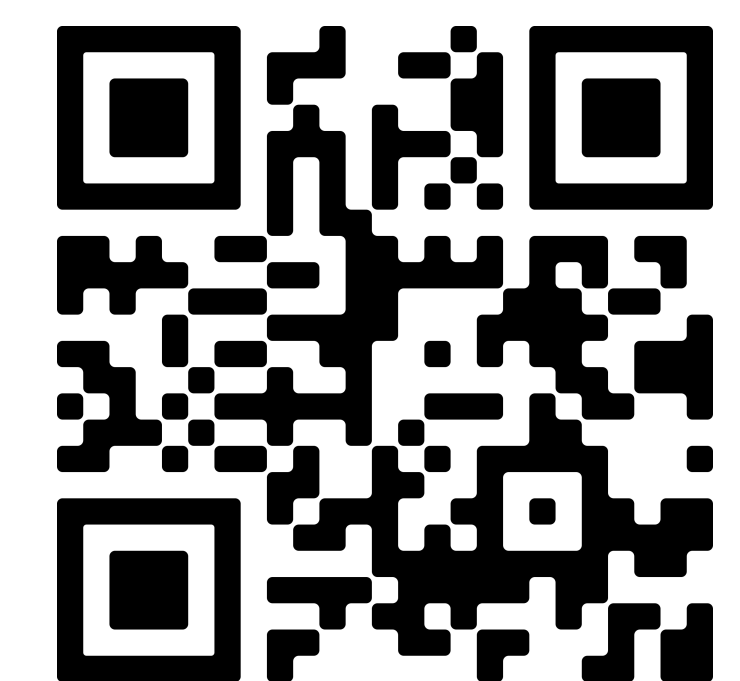


## RESULTS 4: REPOSITORY

Work in progress

Project Github

Web App



## CONCLUSION

- PII is likely to exist (p-value < 0.0001)
- Names are most likely to appear as PII in an essay (90%)
- PII appears most often in the 1<sup>st</sup> fifth of a text

Future work:

- Re-balancing, Name Classes, Truncation

## RELATED LINKS

BERT Paper arXiv:1810.04805, taipy.io