

Solving Redactle Using Masked Language Models

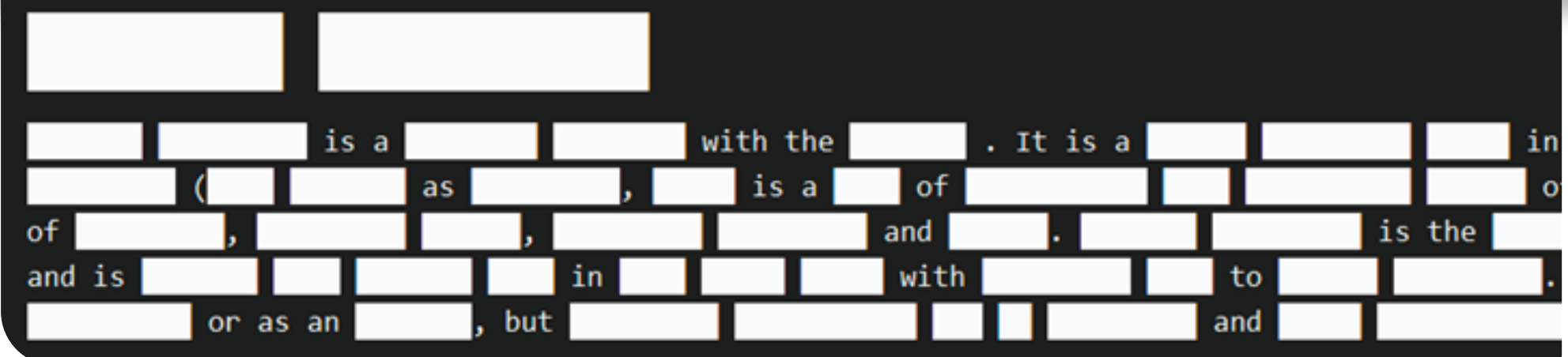
42



Fan Chung Hung, Anthony Guerra, Zhenhao Lu, Zhe Han, Vaibhavi Narayanan

MOTIVATION

- Redactle is a game in which a Wikipedia article is masked and as words in the article are guessed they are unmasked. The game is won when the topic of the article is successfully guessed.
- This project aims to use NLP techniques to solve Redactle in an efficient manner.



OVERVIEW

DATASET

The dataset comprises of 10,000 articles mined from Wikipedia. These articles are 'vital' articles of a variety of important topics according to Wikipedia.

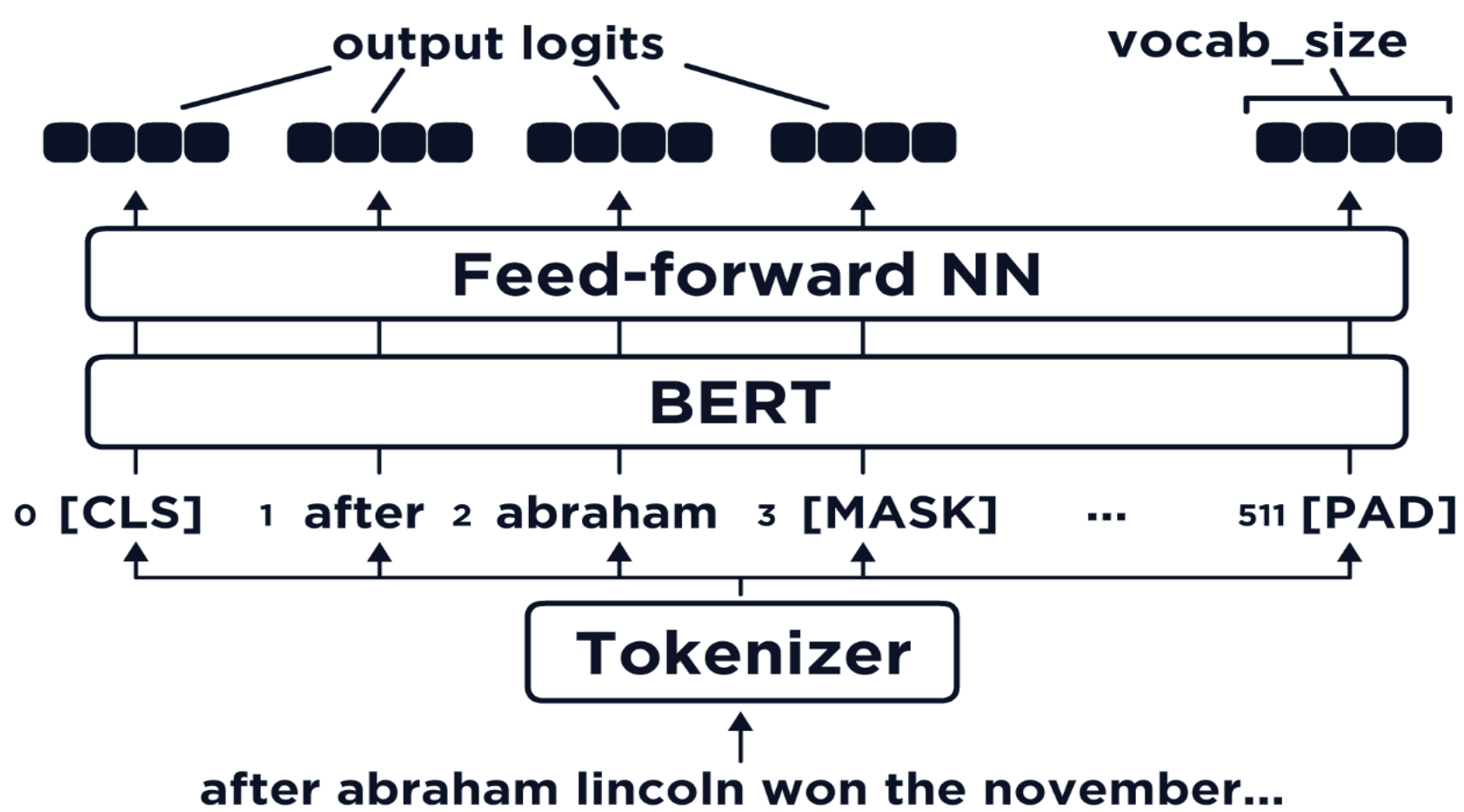
Classes (0 - 11): Arts, Biology, Health Science, Everyday Life, Geography, History, Mathematics, People, Philosophy and Religion, Physical Sciences, Society and Social Sciences, Technology



Dataset Generation (Python): Pandas, Requests, Wiki API

MODELS

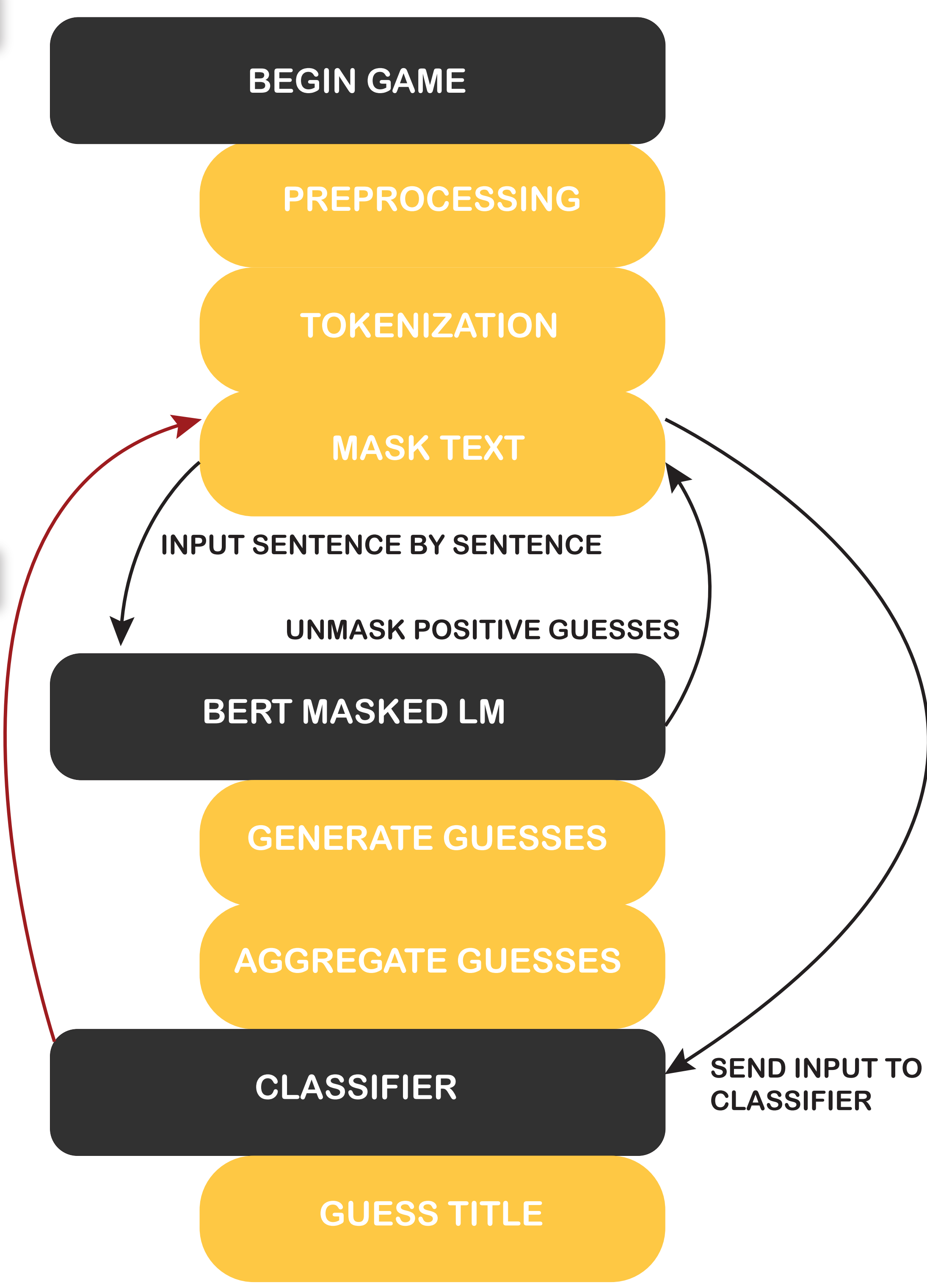
- BERT Masked Language Model**
A bidirectional transformer model used for masked language modeling and next sentence prediction objectives.



CLASSIFIERS FOR TITLE PREDICTION

- BERT for sentence classification.**
- FastText**
A library for efficient learning of word representations and sentence classification.
- Decision Trees**
A non-parametric learning method used for classification and regression.

PROCESS



If the title is guessed correctly the game is over. If not then continue unmasking using BERT.

METRICS and RESULTS

For the unmasking portion we've chosen several ranking evaluation methods:

- The Hit Ratio (HR) is the number of relevant words of the top ten guesses. Relevant words being words that were successfully unmasked.
- The Mean Hit Ratio (MHR) is the mean hit ratio among all articles.
- The Mean Relevant Ranking (MRR) is the mean number of runs it takes to successfully guess the title.

$$HR_i = \frac{|G_{hit}^L|}{|G_{all}|}$$
$$MHR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} HR_i$$
$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i}$$

	BERT	DECISION TREE
MHR	0.56	0.70
MRR	169	15.3