

# Learning to Detect Vulnerable Code Using Differentiable Line Probability of Large Code Models

Minjae Gwon and Sangdon Park.

Research Project II, Dept. of Computer Science and Engineering, POSTECH, Korea.

## Motivation

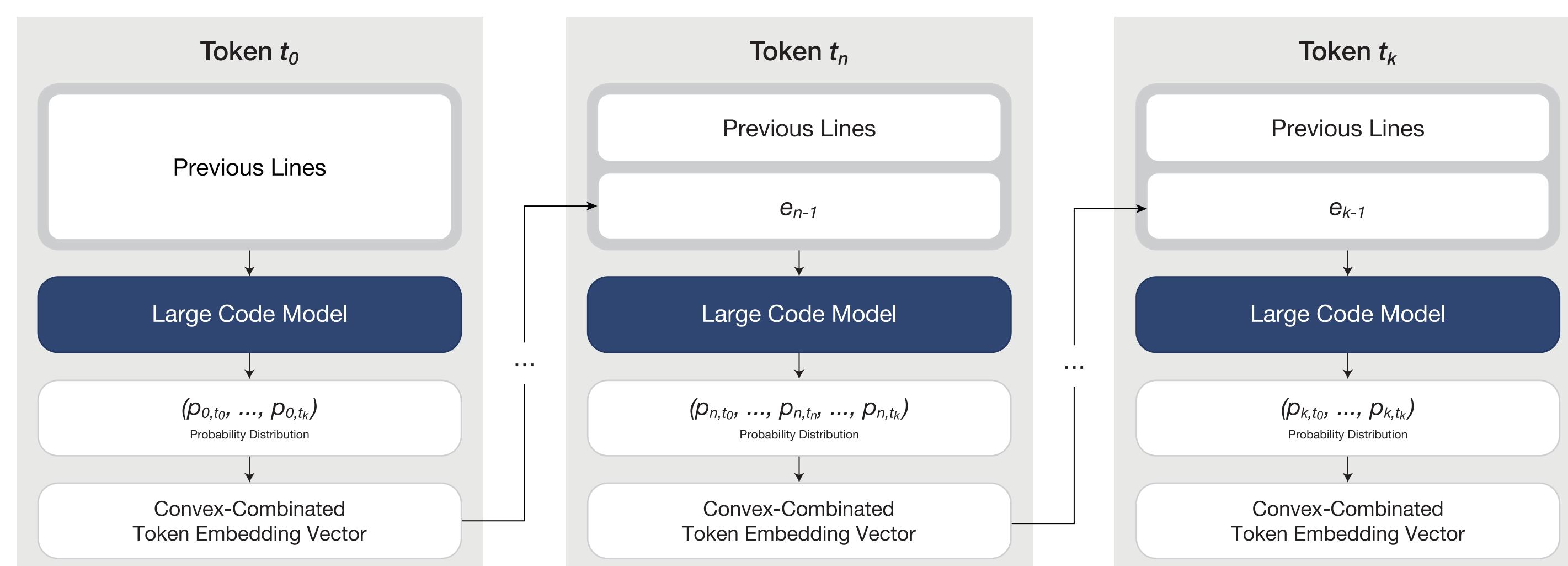
### Problem Description

**Mission.** Automate vulnerability detection to reduce manual effort.

- Automation.** Reduce manual efforts in detection.
- Efficiency.** Quickly identify vulnerabilities in code.
- Accuracy.** Minimize false positives and false negatives.
- Scalability.** Easily integrated with multiple languages.

## Method

### Line Probability

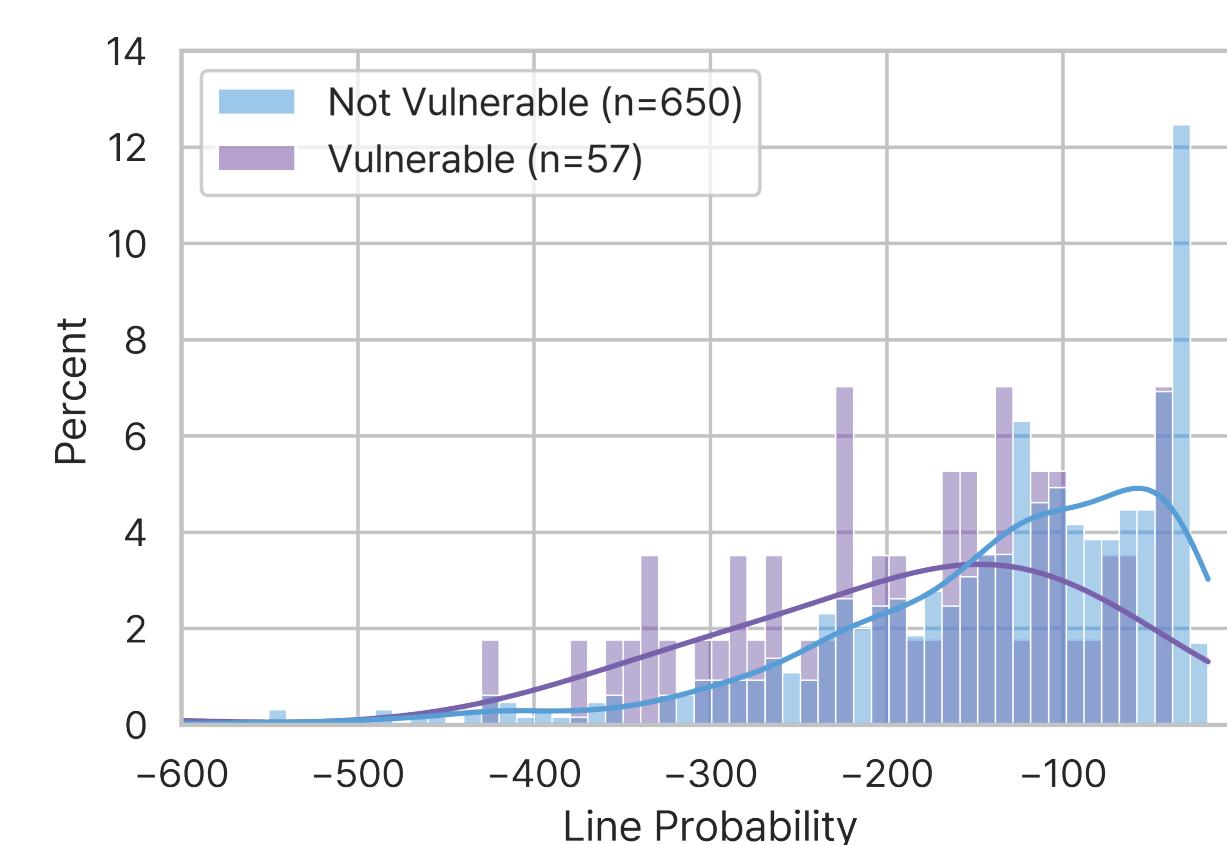


**Why?** Empirically useful in discriminating vulnerable lines from benign lines.

**Objective 1.** Must quantify vulnerability likelihood at the line level.

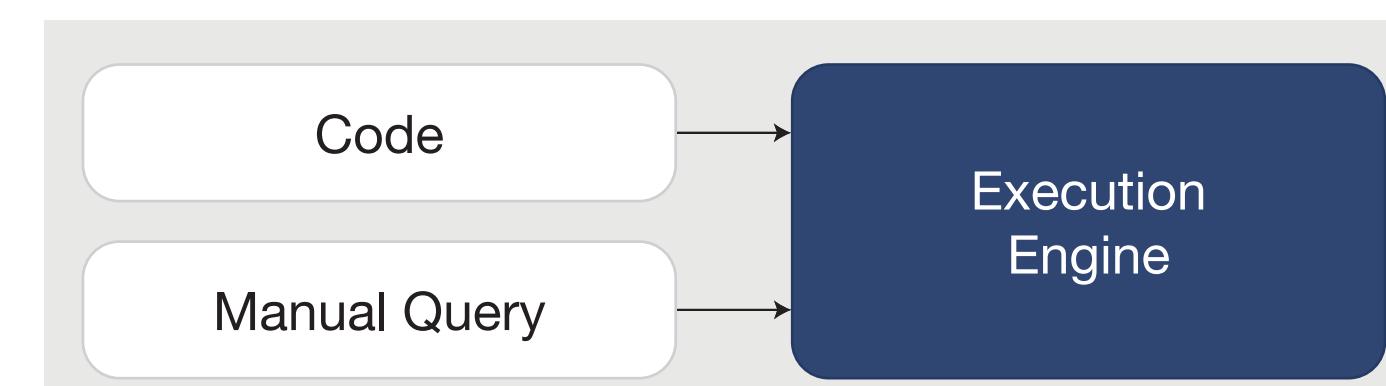
**Objective 2.** Must be differentiable for accurate training.

- Line probability** is the sum of token probabilities which indicates the vulnerability of a code line.
- Token probability** is derived from a distribution of large code model outputs, where the input is a convex combination of prior distributions and some previous code lines.



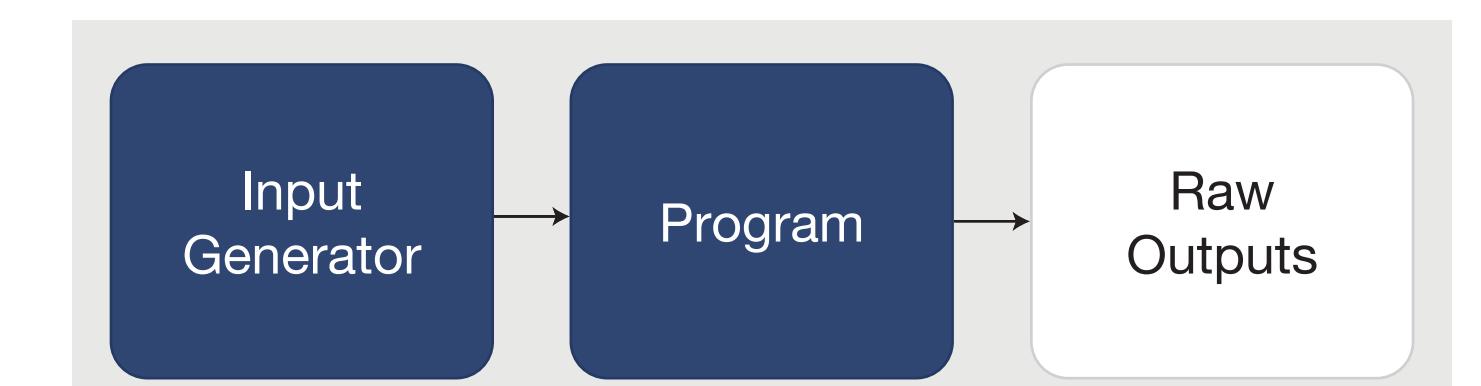
### Related Works

#### 1. Static Analysis



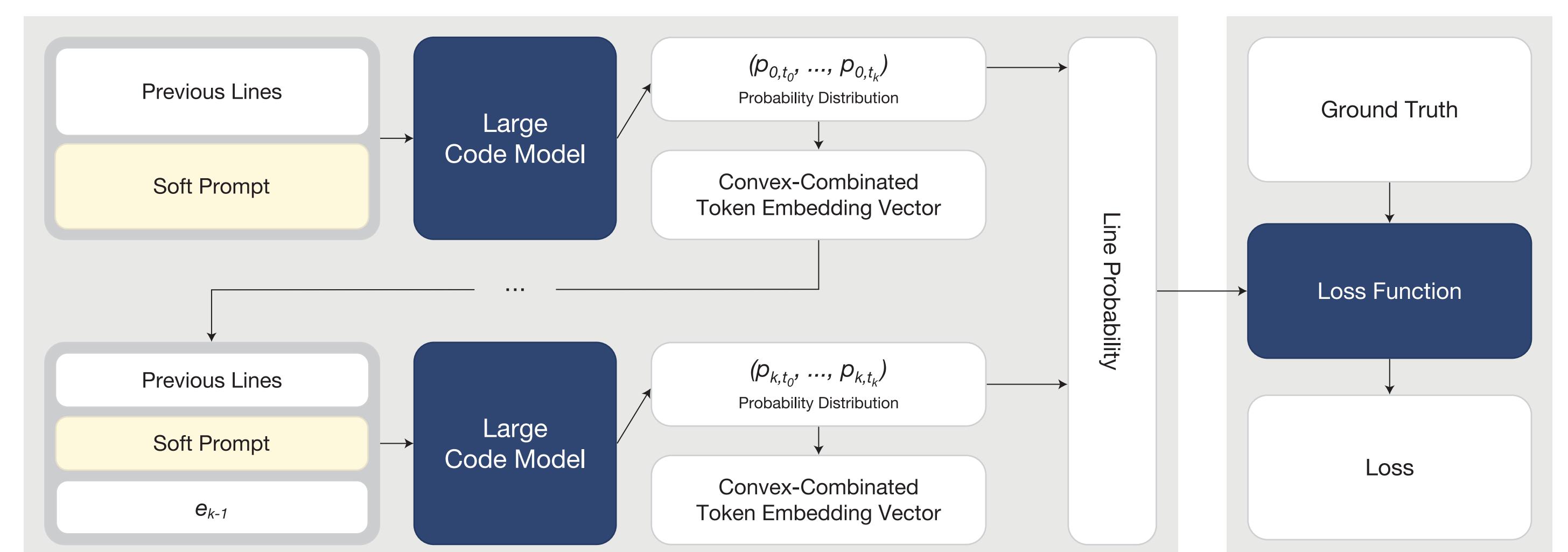
- Requires significant manual effort for rule definition.

#### 2. Fuzzing



- Requires lots of time and manual effort for defining inputs and analyzing outputs.

### Prompt Tuning

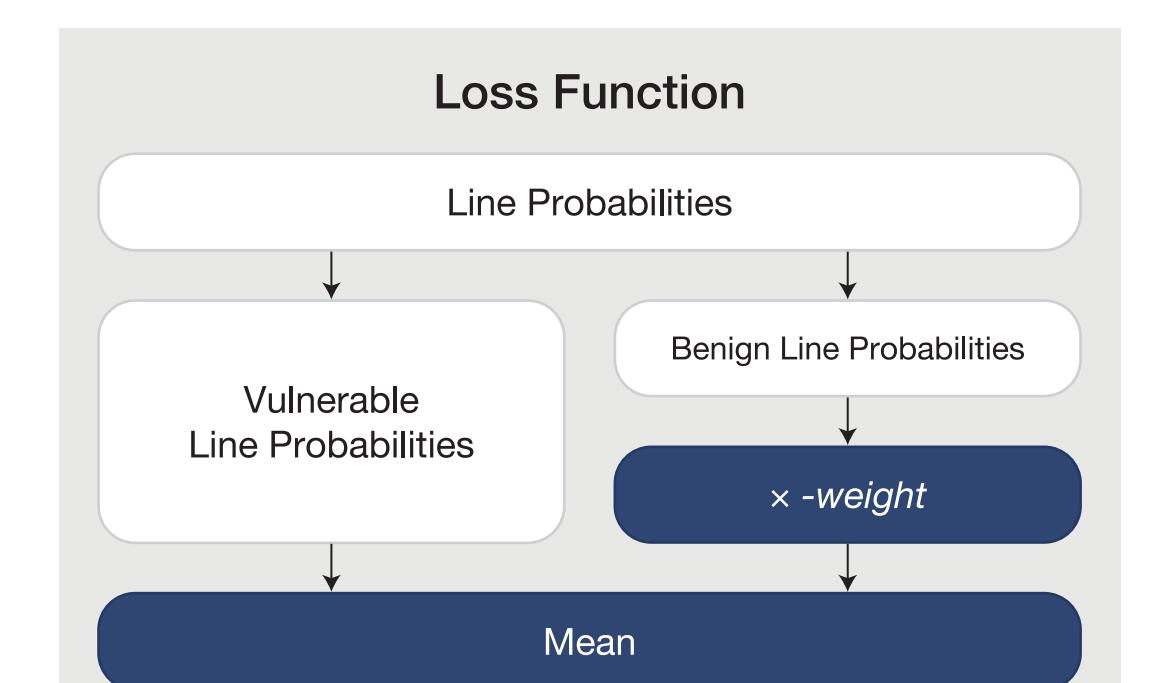


**Why?** Use a small number of parameters to effectively tune the code model's performance.

**Objective 1.** Must optimize the performance of line probability.

**Objective 2.** Must be tuned in terms of line probabilities.

- Prompt tuning** is a mechanism for learning **soft prompts**, enabling models to perform specific downstream tasks.<sup>[1]</sup>
- Large Code Model**, the tuned version of large language models, detects the possibility of vulnerability by line probabilities.



## Result

### Model and Dataset

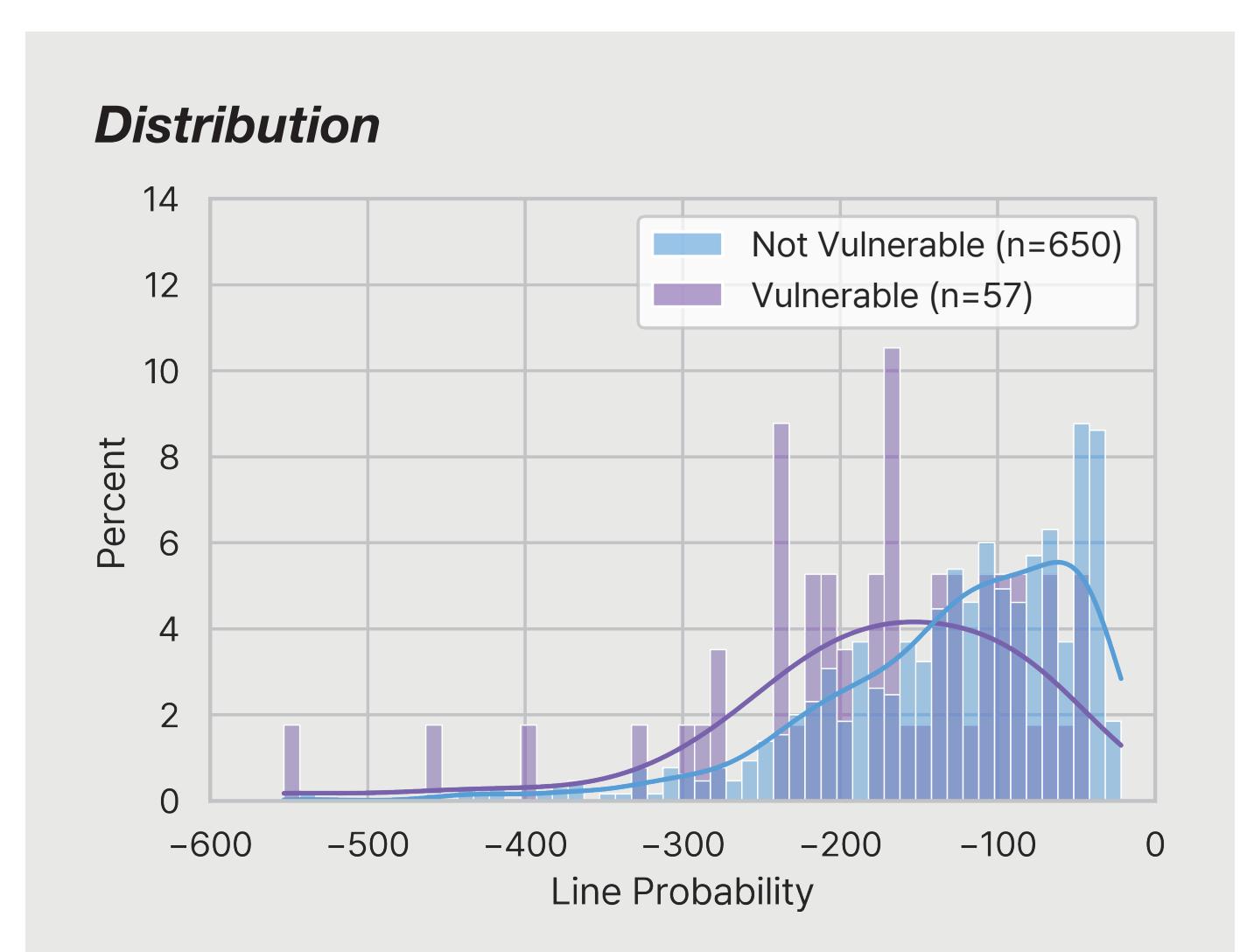
Base Model	Code Llama 7B
Vulnerable Lines	336
Benign Lines	3,207
Validation Split	0.2

### Parameters

Soft Prompt Length	64
Epochs	8
Batch Size	16
Learning Rate	0.0005

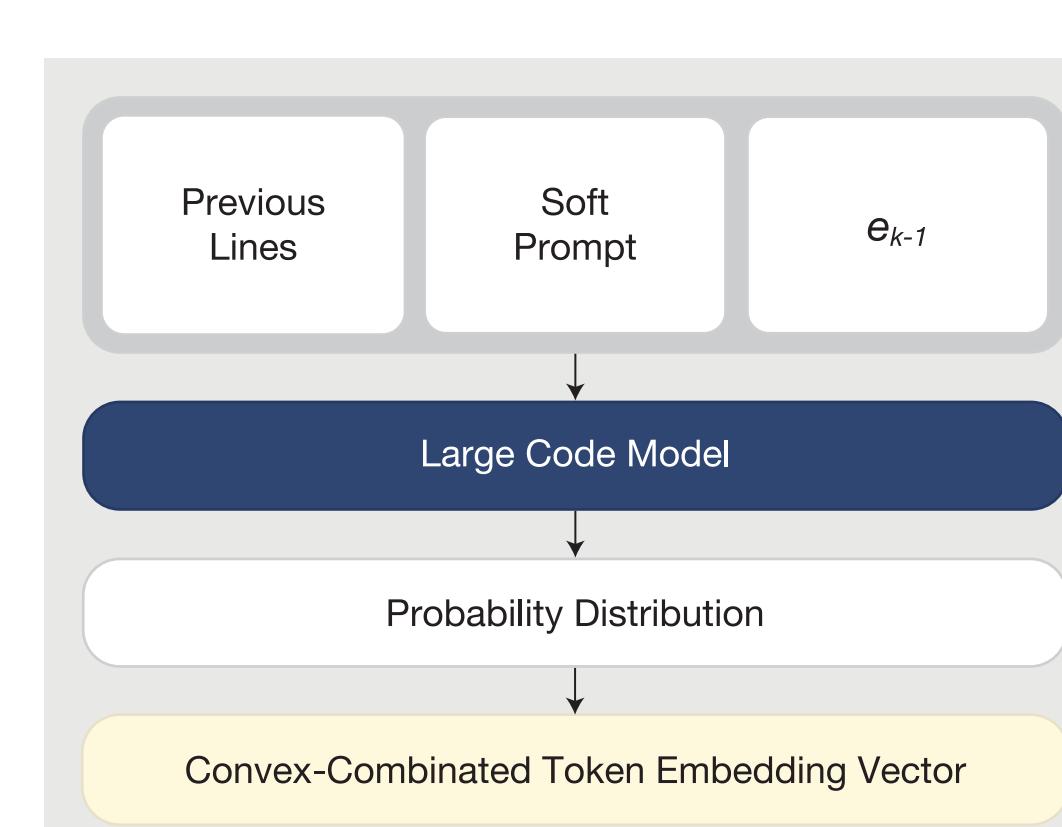
### Analysis

- Observation 1.** Line probability emergence as a potent metric for vulnerability detection.  
**Observation 2.** Prompt tuning makes the model more sensitive to vulnerable lines.



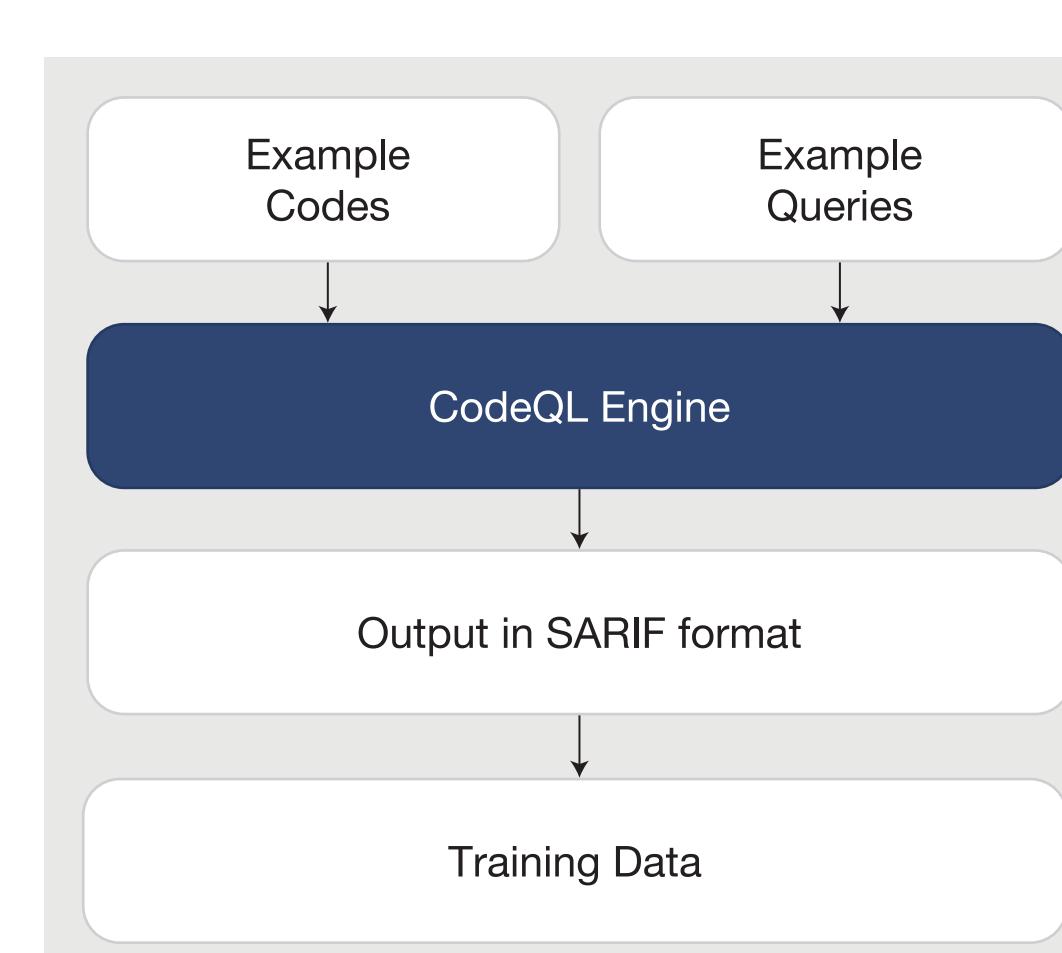
### Discussion

#### 1. High Cost of Calculating Line Probabilities



- Convex-combined token embedding vector is constructed from all-available tokens of a base model.
- A line probability should be calculated sequentially.
- These two reasons make the calculation very slow.
- Exploring ways to calculate line probabilities in parallel could be beneficial.

#### 2. Lack of Data for Training and Testing



- Our dataset is constructed from CodeQL's public repository and sampled DARPA's challenges.
- Some data did not fit well during tuning.
- The size of dataset was insufficient for tuning.
- Explore self-supervised learning methods that do not require labeled data but need to design auxiliary tasks useful for vulnerability findings.