

# Watermarked Sequence Length Code Generation for Large Language Model

Jongsoo Ha, Xiaohui Liang, Youxiang Zhu, Jim Schwoebel

UMass Boston

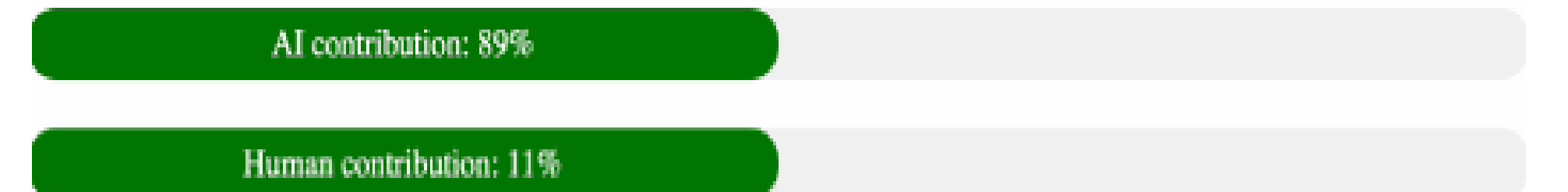


ORACLE

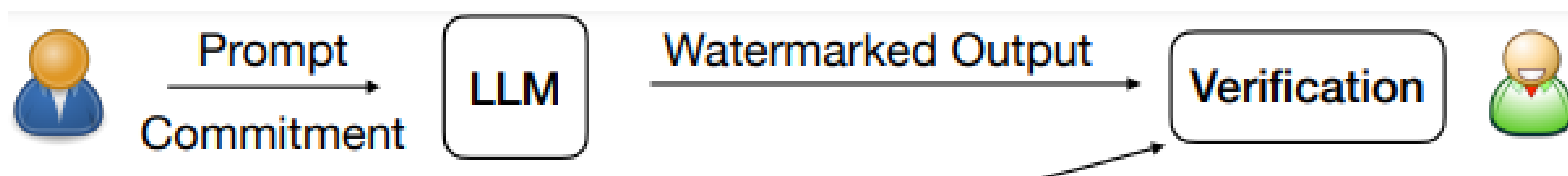
## 1. Introduction

As large language models (LLMs) become widely used in code generation and content creation, **ensuring transparency and accountability is crucial**. A major challenge is distinguishing between human and AI-generated content, which raises concerns about **authenticity, ownership, and misuse**. To address this, watermarking techniques have been developed to embed hidden patterns in AI-generated text and code, allowing for traceability without compromising quality

**Output:**  
Describe World war I, and a soldier... and how they change. Lol, I said I haven't watched it since it was my baby. I've heard a lot about it though. I'm definitely going to check it out soon (hoping for the summer season) Oh cool. The story doesn't need a lot to be amazing. The music is great, the sound of the soldiers war cries is one of the most beautiful sounds to me. I don't expect to be blown away, but you won't know until you watch it. Just sit back. Good idea, for sure it will be a cool experience! :) Thanks man I'd love to see it in theaters. Ya it was my Dad and we went this weekend! Even though I'm over 15, we still didn't want to ruin the experience for my brother (22) who is a bit more invested in it all, he watched it all when it first came out but never on the big screen



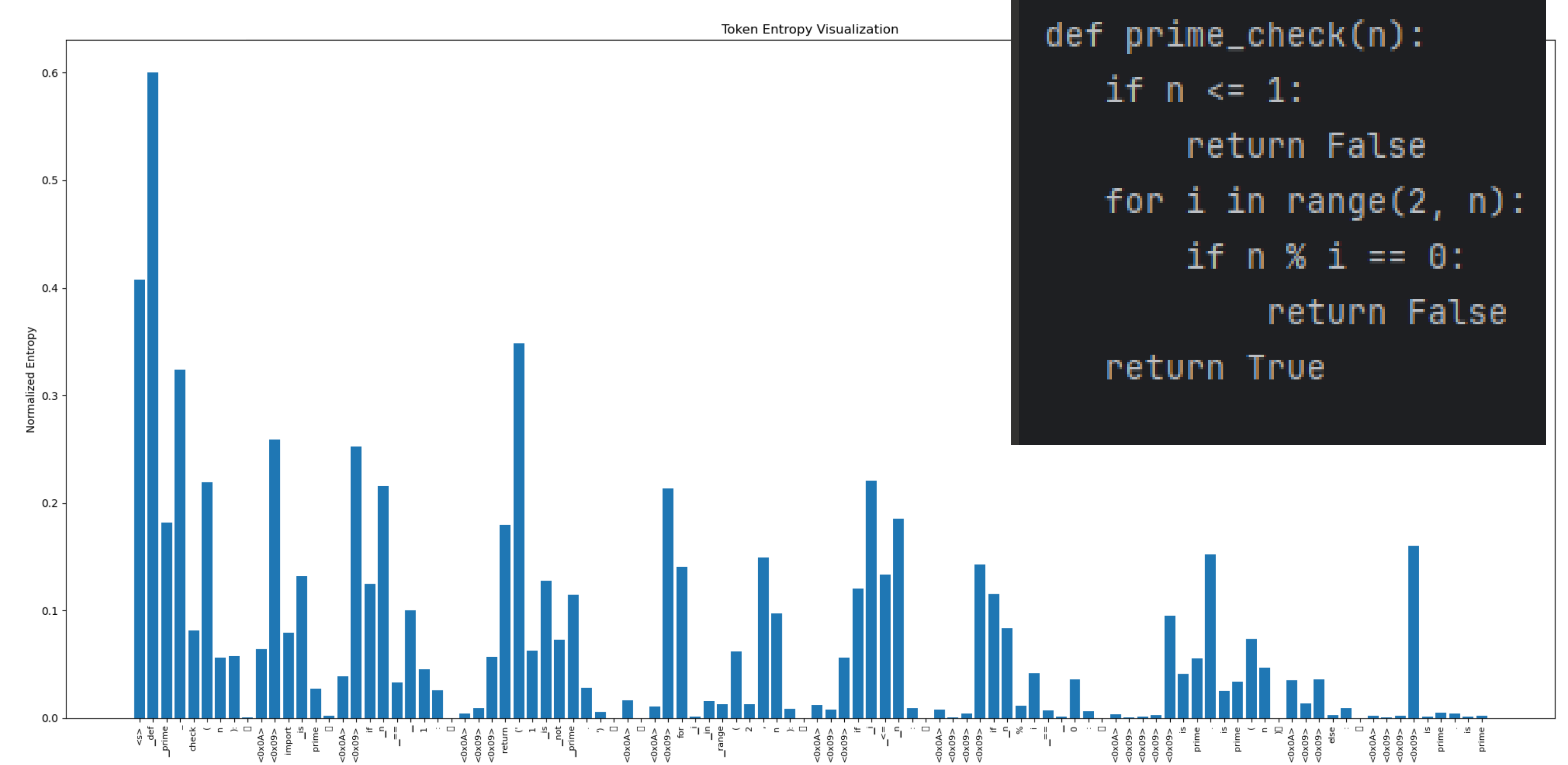
## 2. What is the Current limit of watermarking system?



The current watermarking system, originally **designed for essays and paragraphs**, faces **challenges when applied to code generation**. It struggles to maintain performance and exhibits a low detection rate, especially when the generated code functions correctly.

We introduce **Sequence Length Watermark (SLW)** to ensure robustness and code performance!

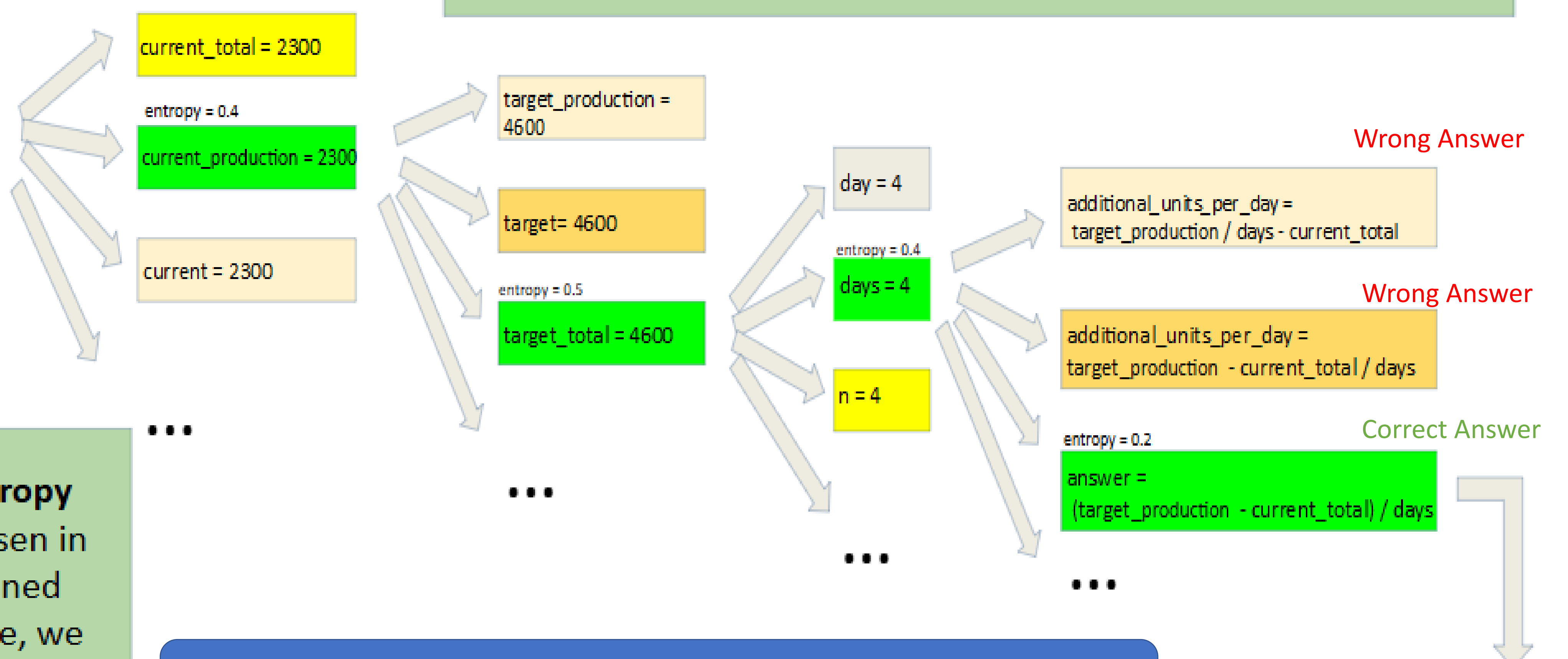
Entropy chart for each tokens



## 3. How SLW improve the watermark robustness and code performance?

To address this issue, we are **focusing on watermarking low-entropy sequences within the code**. This approach ensures that the watermark remains **robust** while preserving the overall functionality and **performance of the generated code**.

**Question**  
A factory produces 2300 units in 4 days. How many more units must it produce each day to reach a total of 4600 units in 4 days?



We apply **watermarking** to sequences with **low entropy thresholds**. For example, variable names can be chosen in countless ways, but final calculations must use defined variable names with the correct equations. Therefore, we embed the watermark in **fixed, unchangeable lines**.

## 4. Key takeaways

### Performance:

- Measure performance using Humaneval dataset.
- Ensure the generated code runs efficiently without significant overhead

### Watermark Detection:

- Implement a watermark embedding method for code generation.
- Evaluate the effectiveness of watermark detection in various scenarios.

### Analysis:

- Use open-source LLMs designed for code generation.
- Compare results across different models and configurations.

### Goal:

- Achieve a balance between code performance and watermark robustness.

Method	Humaneval	auroc
Watermark	passk	
WLLM	29.6	0.402
SWEET	32.6	0.943
SLW	TBD	TBD