

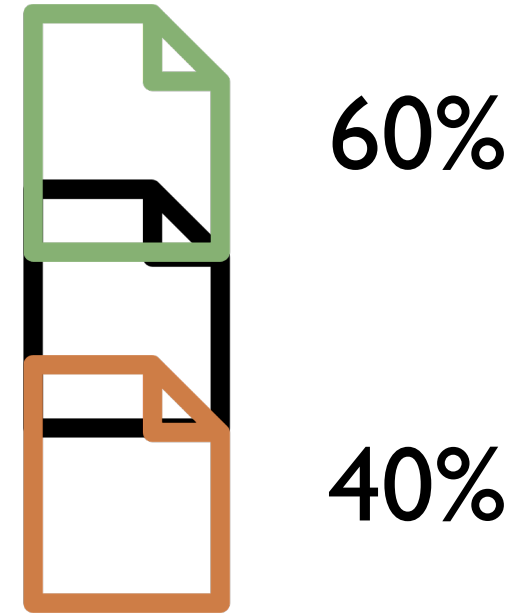
Large Language Models for Code: Secure Hardening and Adversarial Testing

Jingxuan He and Martin Vechev

ACM CCS 2023

LLMs for Code Generation

Prompt + LLM



Pearce et al., Asleep at the Keyboard? Assessing the Security of GitHub Copilot's Code Contributions, IEEE S&P 2022

Li et al., StarCoder: May the Source be With You!, arXiv:2305.06161

Making LLMs generate ~~unsafe~~ **code** more often?

Security Hardening & Adversarial Testing

Secure



Prompt + LLM

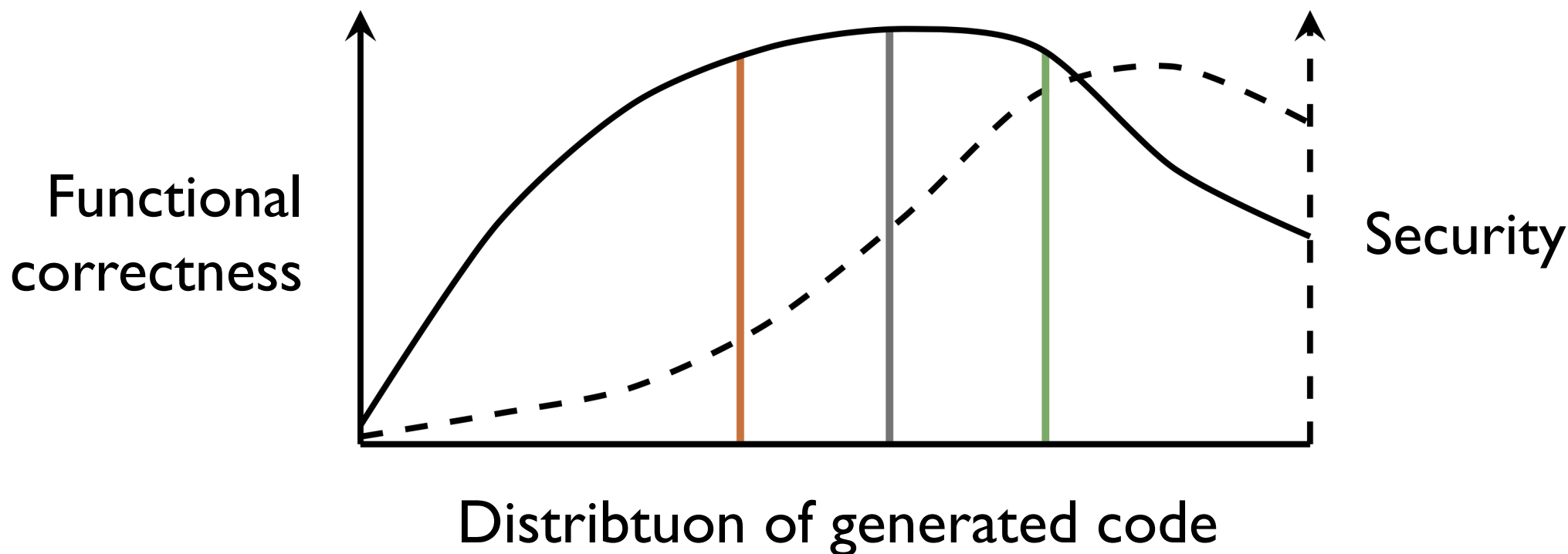


Unsafe

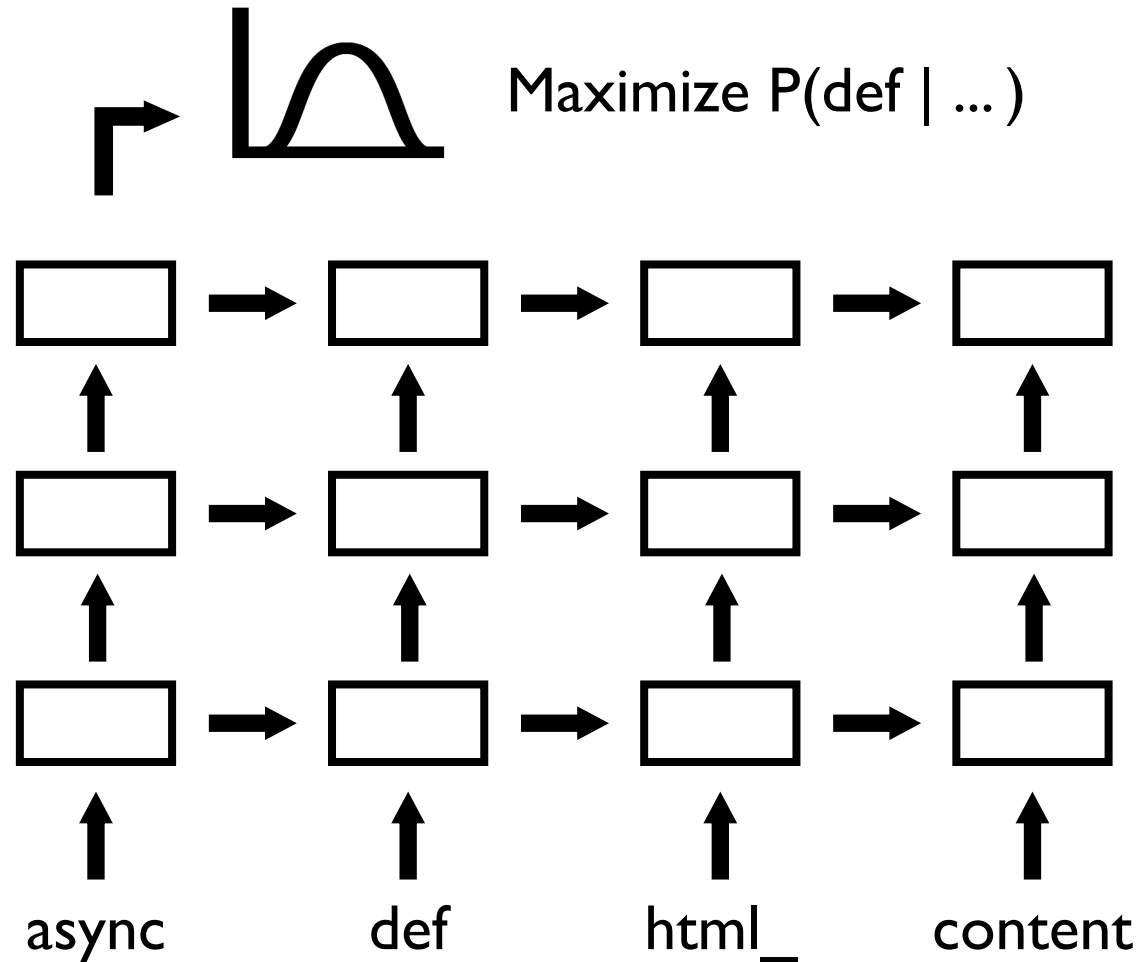


Functional Correctness & Security

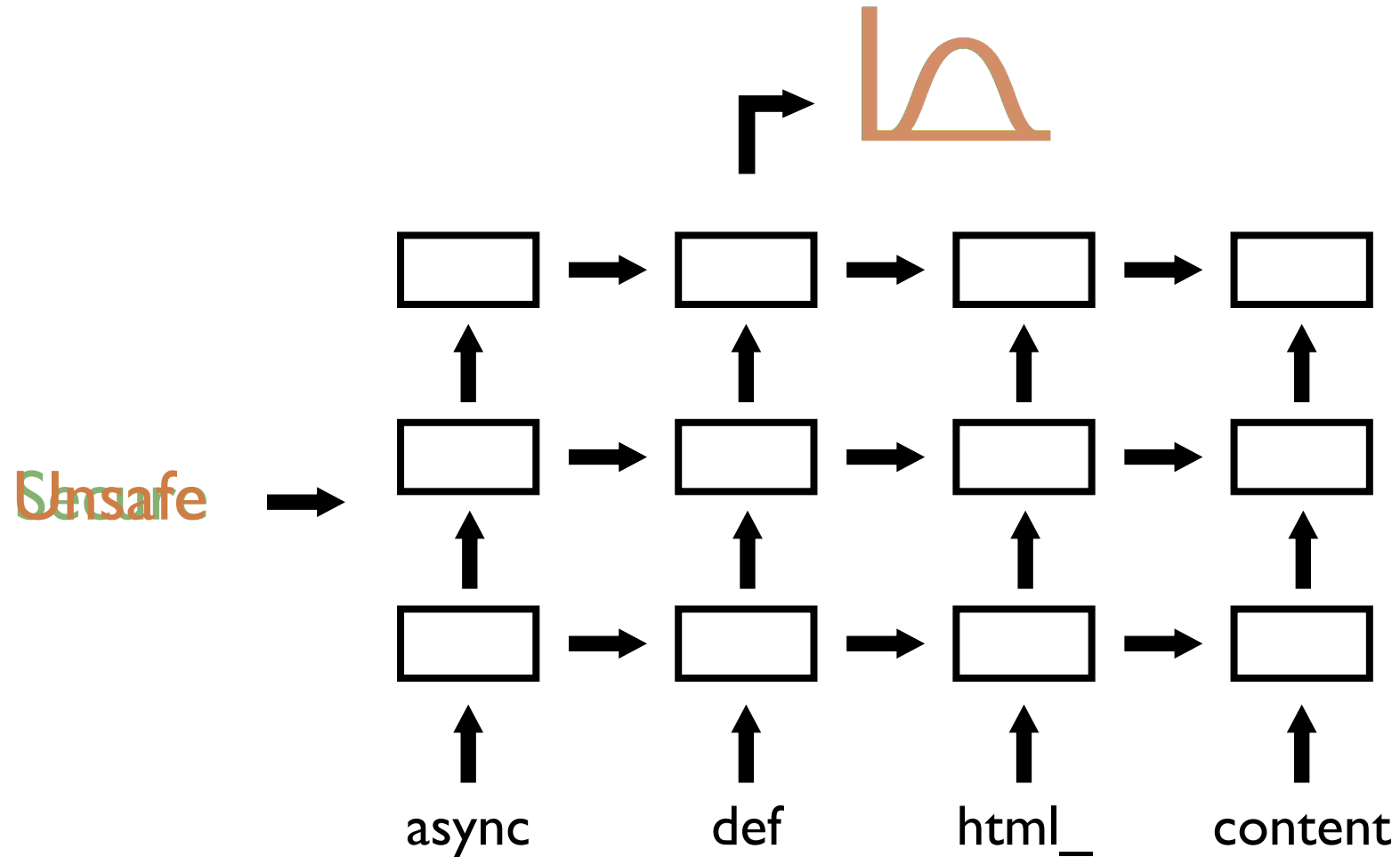
Adversarial Testing LLM Security Hardening



Language Modeling

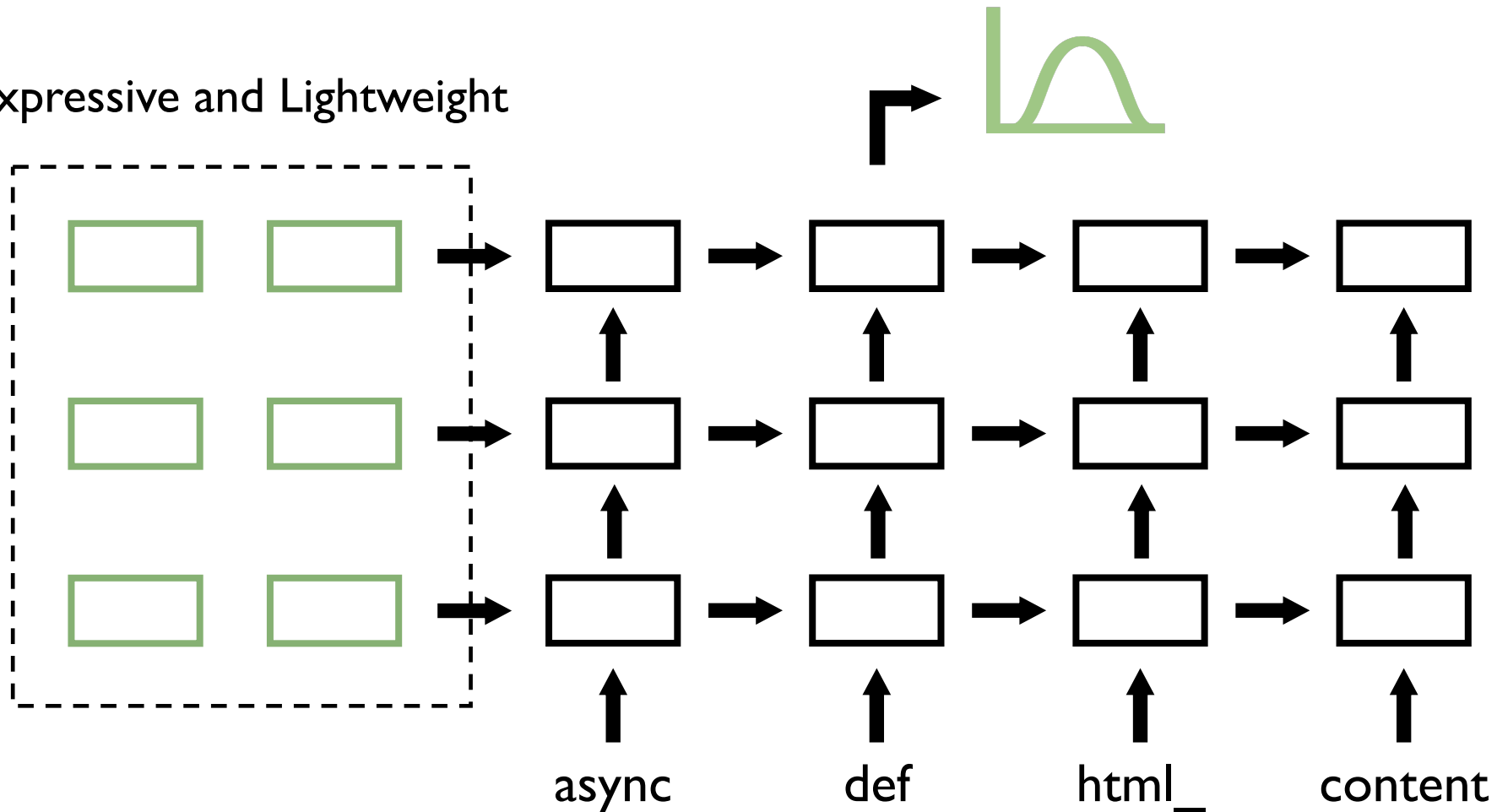


Conditional Language Modeling



SVEN: Soft Prompts as Prefixes

Expressive and Lightweight



SVEN: Training Data

Security fixes extracted from GitHub commits:

```
async def html_content(self):  
- content = await self.content  
return content
```

```
async def html_content(self):  
+ content = markupsafe.escape(await self.content)  
return content
```

A Naïve Training Scheme:



SVEN: Code Regions

```
async def html_content(self):  
- [redacted] wa [redacted] f. [redacted]  
return content
```

```
async def html_content(self):  
+ conte [redacted] fe [redacted] s [redacted] (content)  
return content
```

SVEN Training: Security

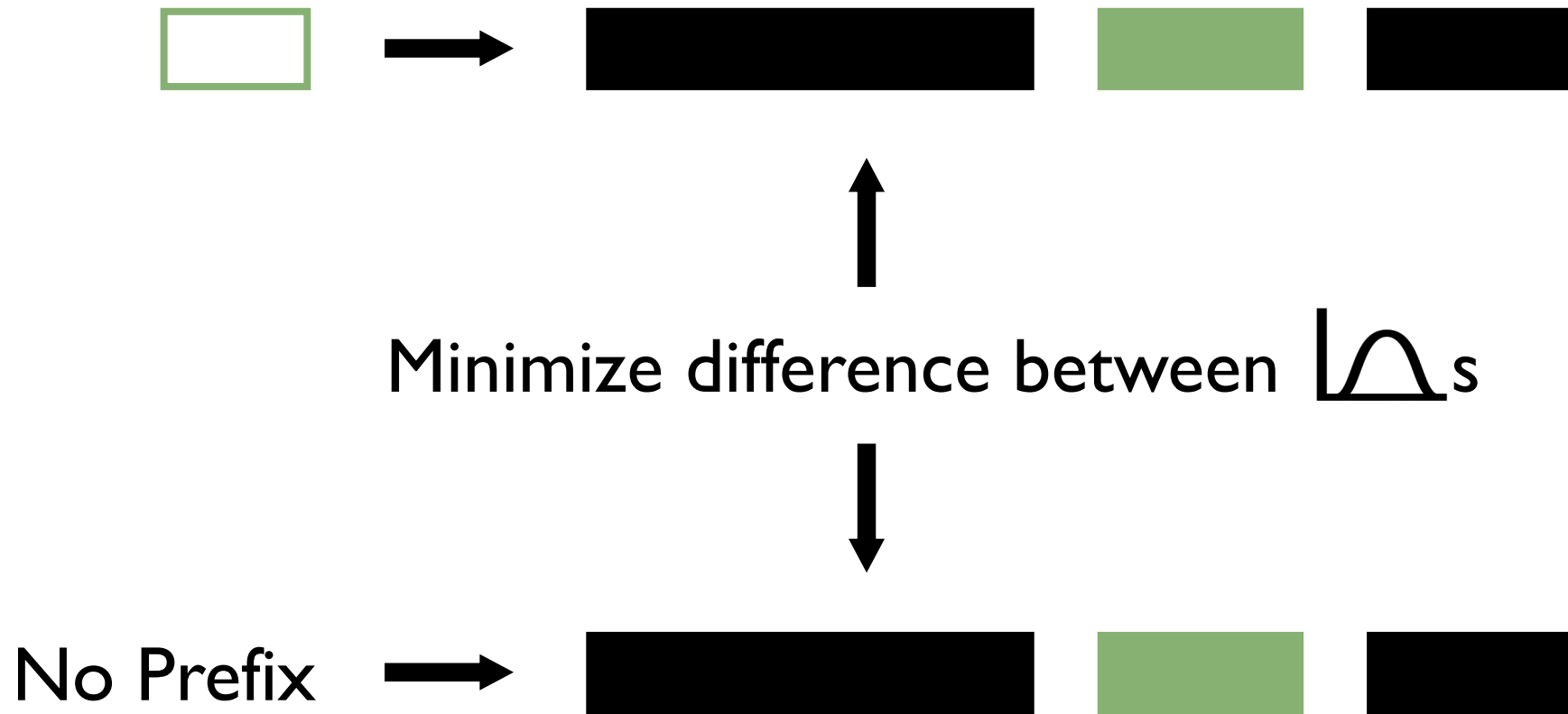


Maximize $P(\text{■} | \dots, \square)$



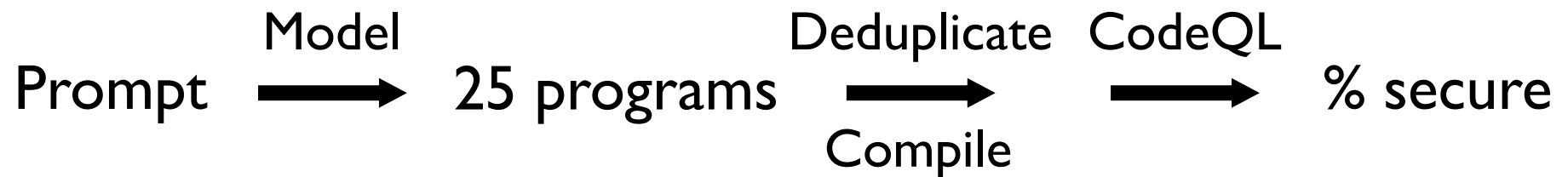
Minimize $P(\text{■} | \dots, \square)$

SVEN Training: Functional Correctness

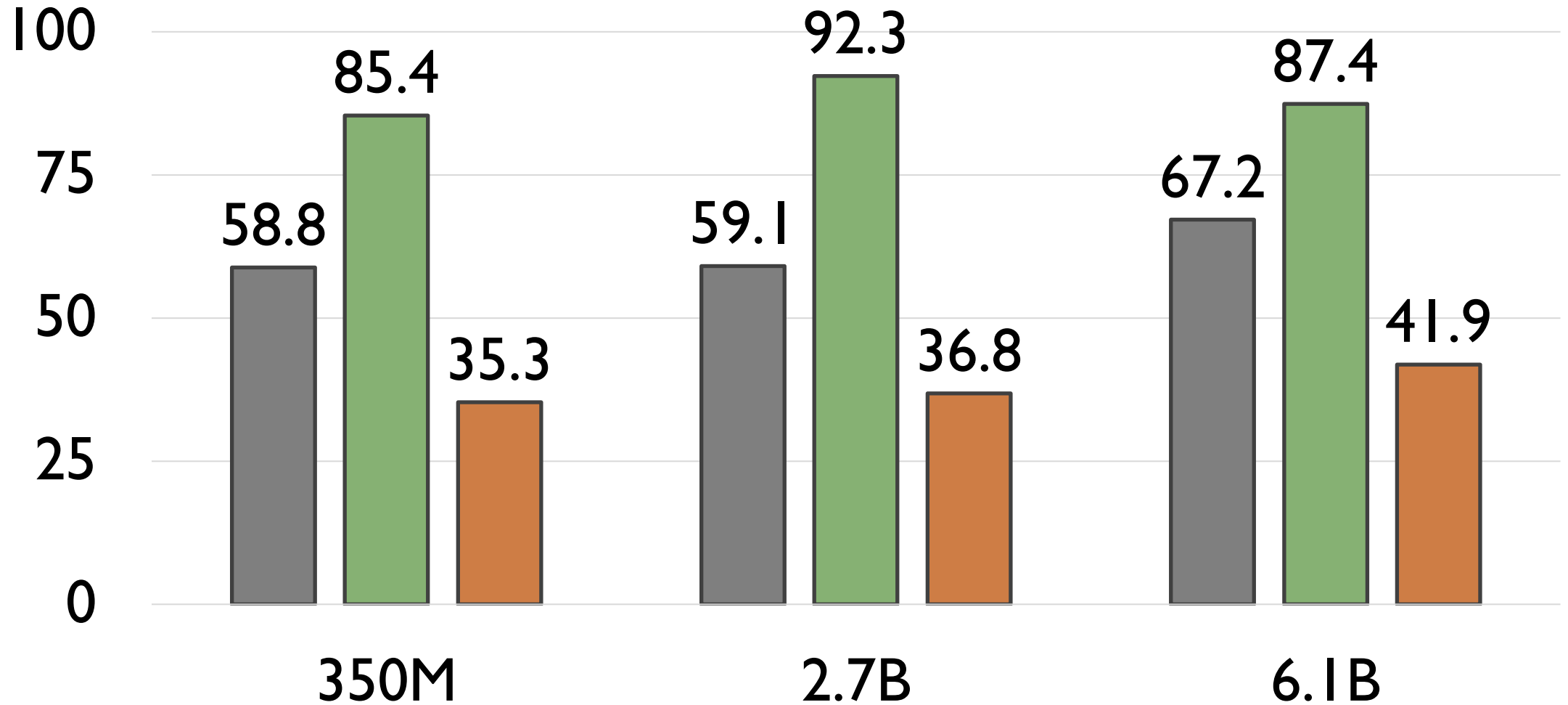


Experimental Setup

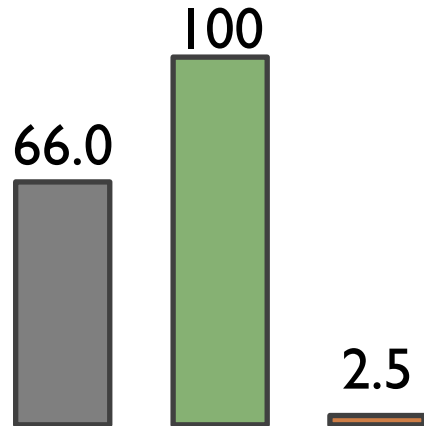
- Manual curation for high-quality training data
- LLMs: CodeGen with 350M, 2.7B, and 6.1B parameters
- Evaluating functional correctness: pass@k on HumanEval
- Evaluating security:



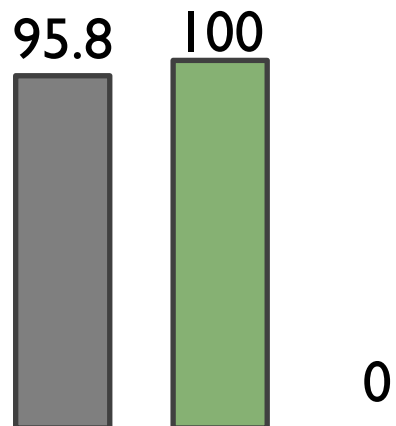
Overall Security



Example: SQL Injection

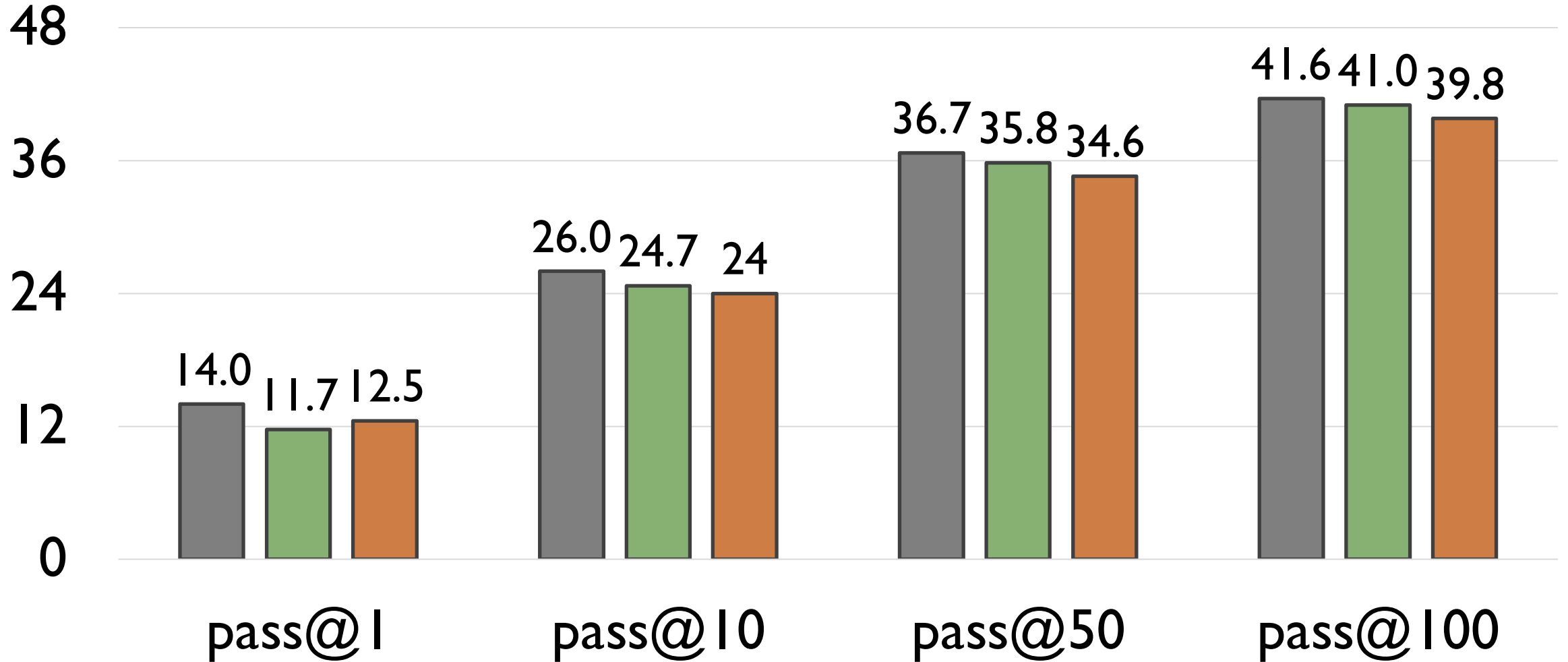


```
query = "INSERT INTO messages (username, message) " \
        "VALUES ('" + username + "', '" + message + "')"
cursor.execute(query)
```

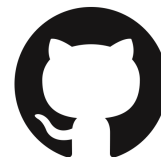


```
query = "INSERT INTO messages (username, message) " \
        " VALUES (%s, %s);"
cursor.execute(query, (username, message))
```

Functional Correctness



arXiv 2302.05319



eth-sri/sven

Q & A