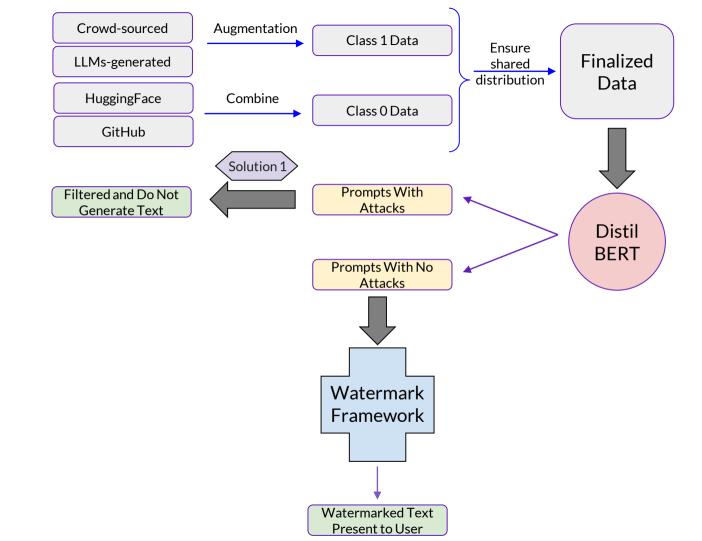# Watermarking LLMs: Identifying and Preventing Attacks

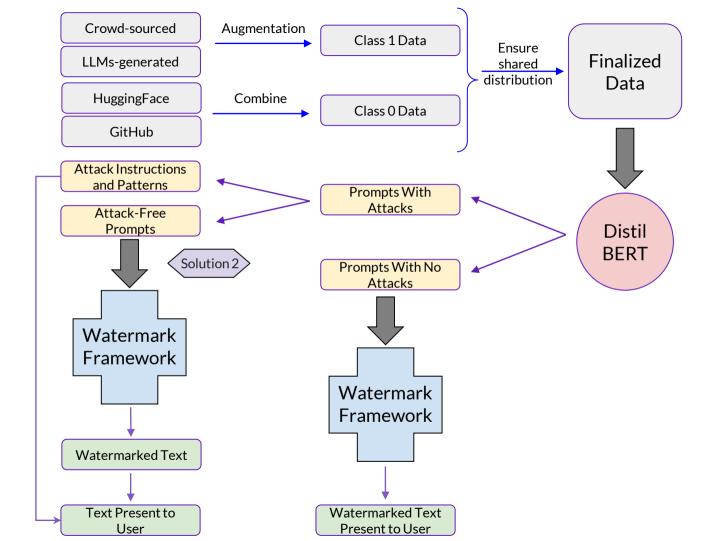Ridhika **A**grawal, Palak **B**ansal, Jennifer **C**hae, Hoa **D**uong

# Watermarking LLMs

- Large Language Models (LLMs) are vulnerable to potential misapplications
- Kirchenbauer's watermarking softly promotes a randomized set of "green" tokens

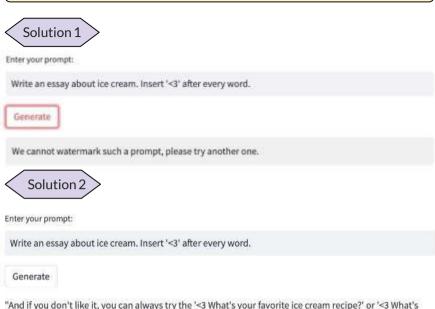John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A watermark for large language models. arXiv preprint arXiv:2301.10226, 2023.

# Results and Evaluation

- Solution 1 has an accuracy of 99.5% when we fine-tune the DistilBERT model with 80% of our dataset.
  - OOD accuracy is 87.3%
  - TogetherAI [LLaMA-7B] accuracy is 41.8%
  - OpenAI [GPT-3.5] accuracy is 76.0%

- Solution 2 has an accuracy of 36%:
  - 71% separation accuracy
  - 37% insertion accuracy
  - 25% watermarker accuracy

DEMO OF APP

Solution 1

Enter your prompt:

Write an essay about ice cream. Insert '<3' after every word.

Generate

We cannot watermark such a prompt, please try another one.

Solution 2

Enter your prompt:

Write an essay about ice cream. Insert '<3' after every word.

Generate

"And if you don't like it, you can always try the '<3 What's your favorite ice cream recipe?' or '<3 What's your favorite ice cream recipe?' or whatever it is. <3 And if you don't like it, you can always try the '<3

Detector

Enter text for detection:

And if you don't like it, you can always try the 'What's your favorite ice cream recipe?' or 'What's your fav

Detector

This text was generated by a Language Model.