

Increasing the Cost of Model Extraction with Calibrated Proof of Work

Adam Dziedzic, Muhammad Ahmad Kaleem,
Yu Shen Lu, Nicolas Papernot

International Conference on Learning Representations (ICLR)

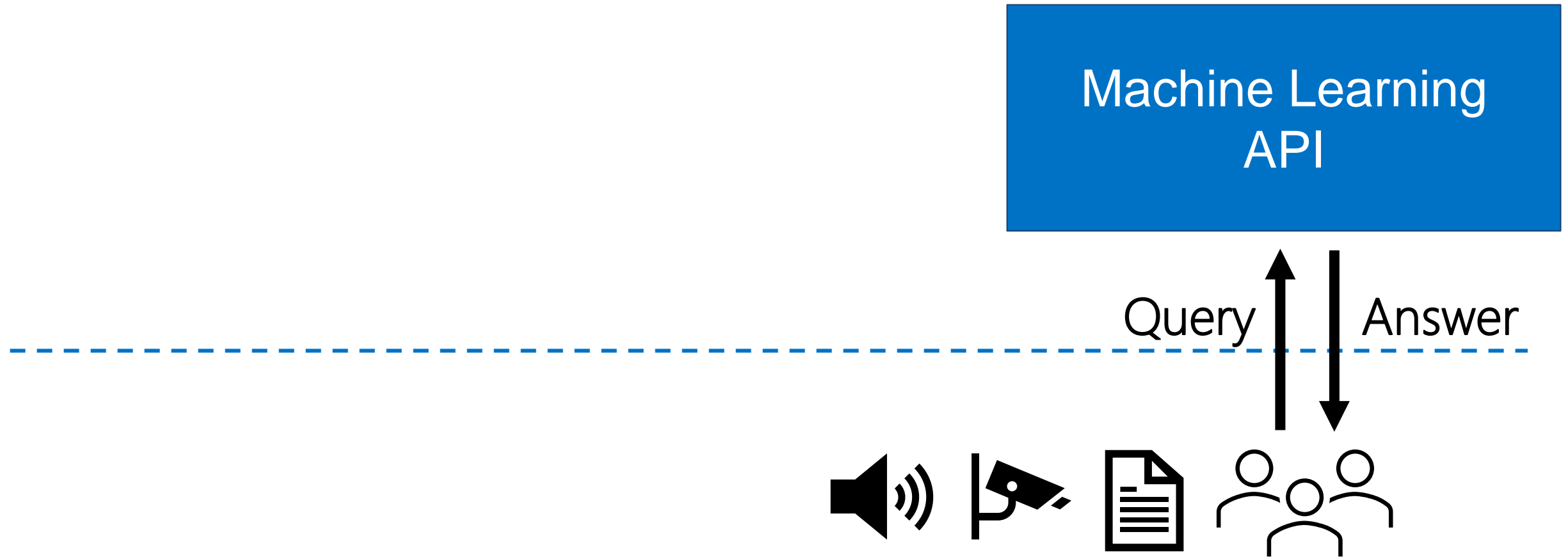
SPOTLIGHT TALK

March 15th, 2022

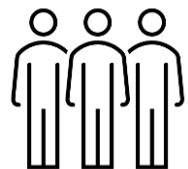


UNIVERSITY OF
TORONTO

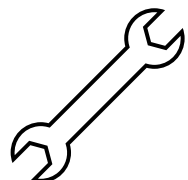
Annotate Data Using Machine Learning APIs



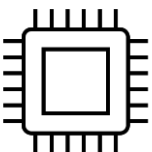
Train Models for Machine Learning Services



Collect & Label Data



Tune Hyper-parameters

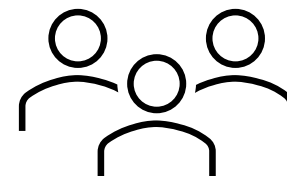
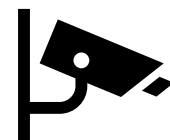
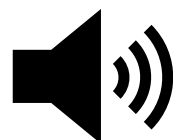


Run on GPU/TPU/CPU

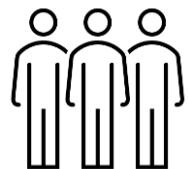
Machine Learning
API

Query

Answer

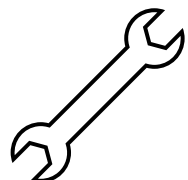
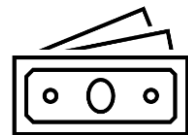


Train Models for Machine Learning Services

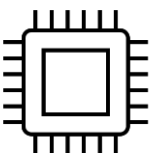
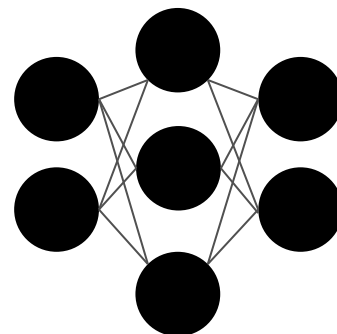


Collect & Label Data

\$ 12 M GPT-3



Tune Hyper-parameters

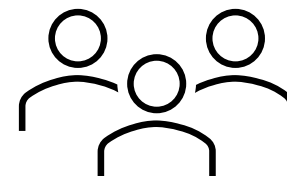
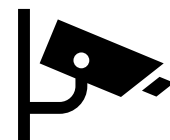
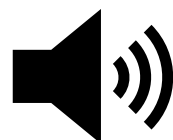


Run on GPU/TPU/CPU

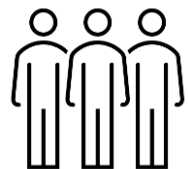
Machine Learning
API

Query

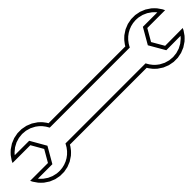
Answer



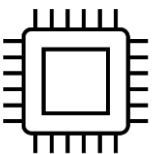
Stealing Machine Learning Models



Collect & Label Data

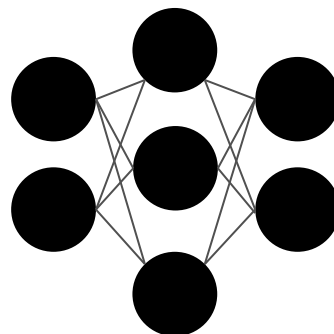
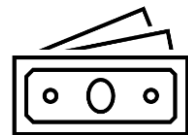


Tune Hyper-parameters



Run on GPU/TPU/CPU

\$ 12 M GPT-3



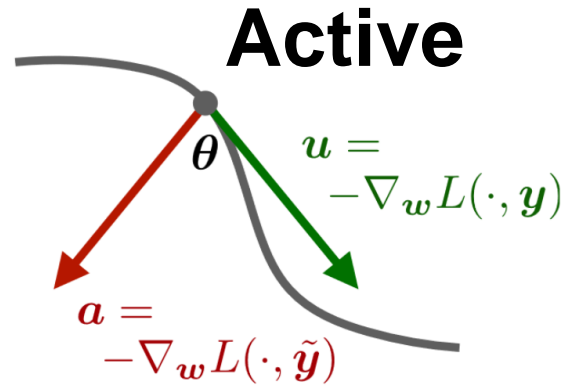
Machine Learning
API

Query

Answer



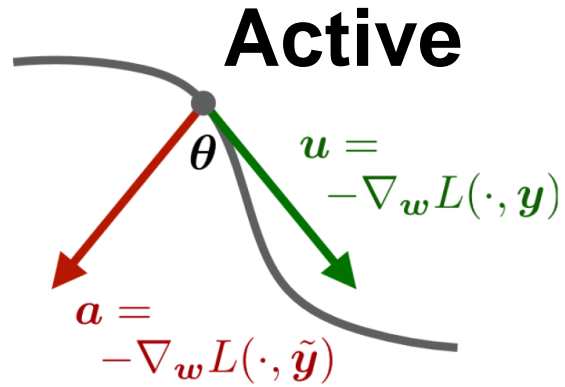
Defenses against Model Stealing



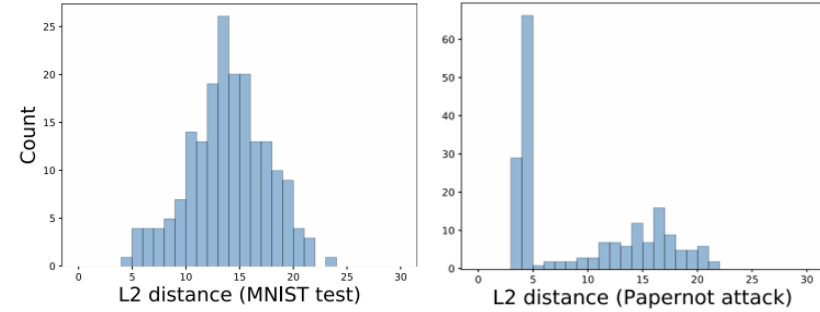
Poison Attacker's Objective

Prediction Poisoning [Orekondy et al. 2020]

Defenses against Model Stealing



Passive



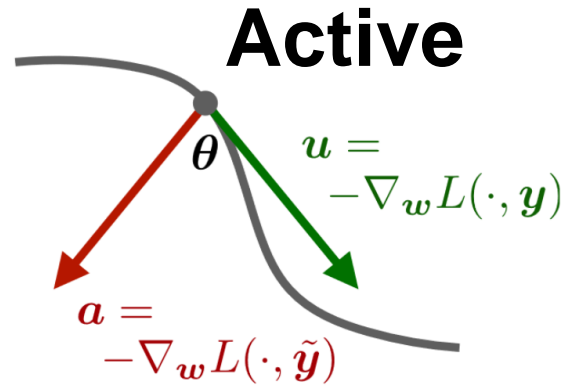
Poison Attacker's Objective

Prediction Poisoning [Orekondy et al. 2020]

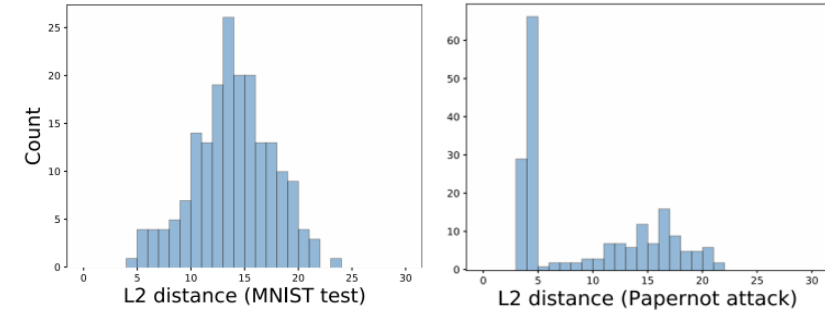
Detect Attack & Stop Responding

PRADA [Juuti et al. 2019]

Defenses against Model Stealing



Passive



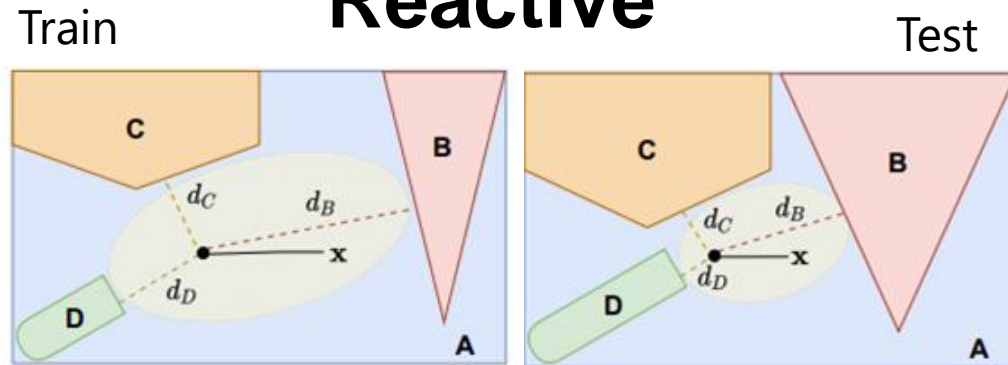
Poison Attacker's Objective

Prediction Poisoning [Orekondy et al. 2020]

Detect Attack & Stop Responding

PRADA [Juuti et al. 2019]

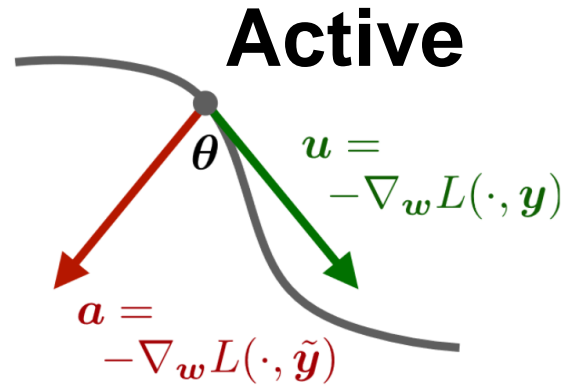
Reactive



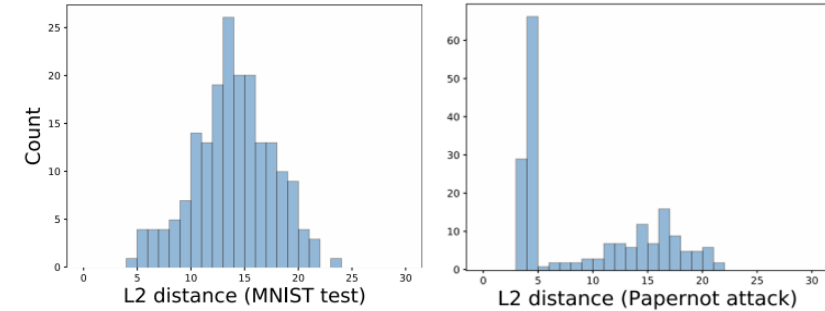
Resolve Model Ownership

Dataset Inference [Maini et al. 2021]

Defenses against Model Stealing



Passive

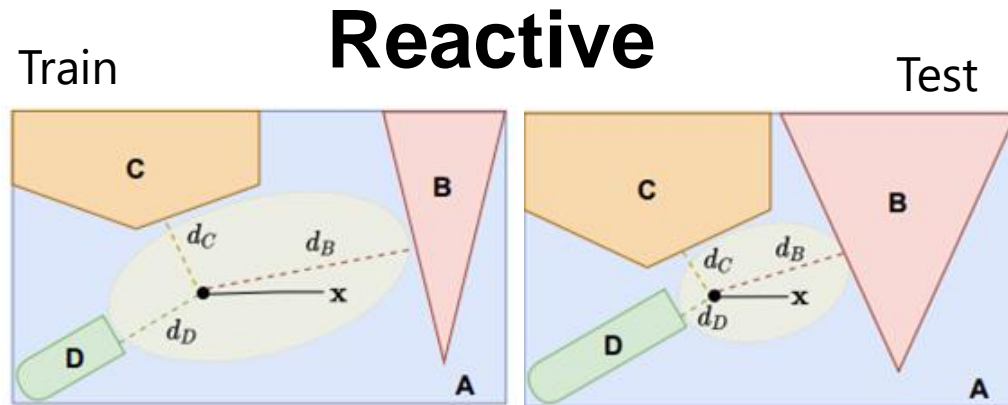


Poison Attacker's Objective

Prediction Poisoning [Orekondy et al. 2020]

Detect Attack & Stop Responding

PRADA [Juuti et al. 2019]



Resolve Model Ownership

Dataset Inference [Maini et al. 2021]

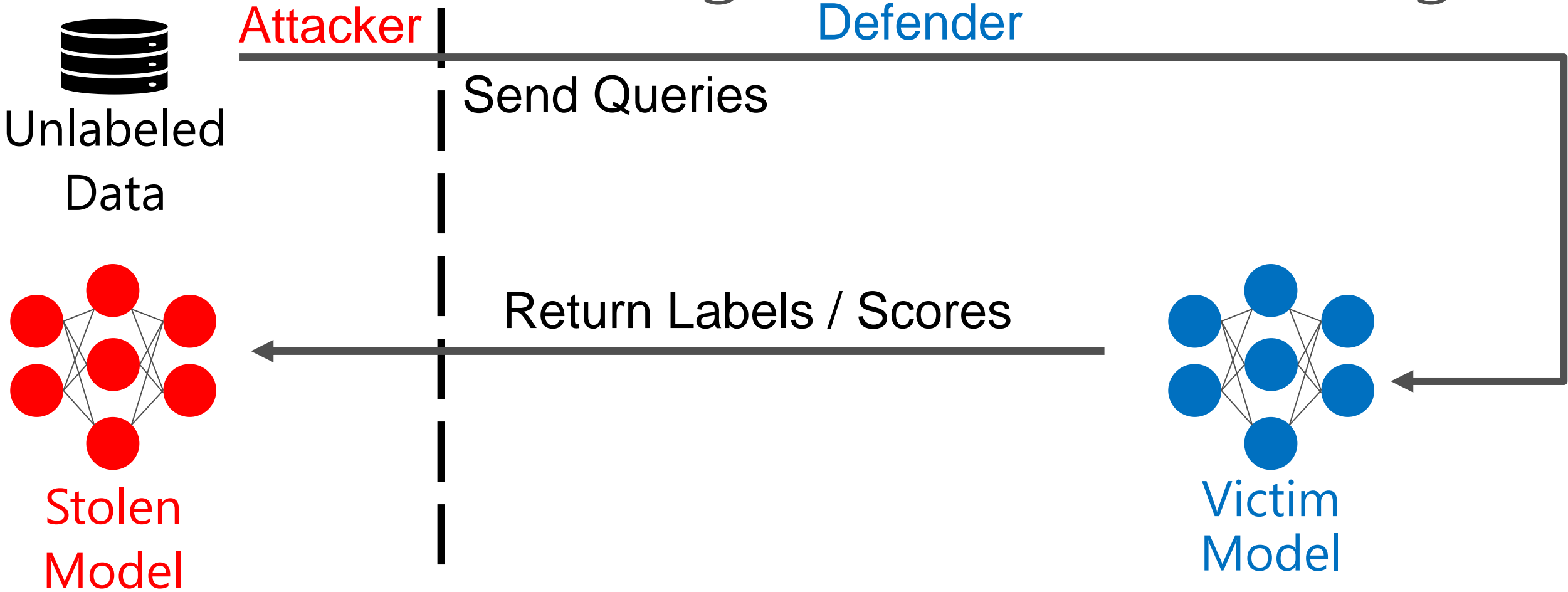
Pro-Active

Differential Privacy

Proof-of-Work

Calibrated Proof-of-Work with PATE

How to Defend Against Model Stealing?



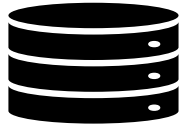
[Shankar et al. 2020]

Model Stealing - ranked among the most severe attacks against ML

Estimate Victim Model Information Leakage

Attacker

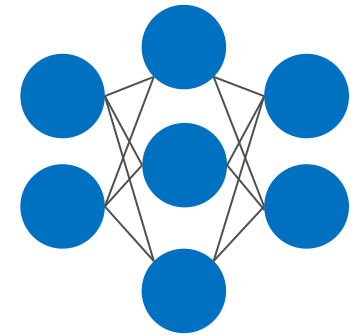
Defender



Unlabeled
Data

Send Queries

Estimate
Privacy
Leakage

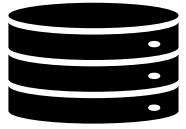


Victim
Model

Generate Calibrated Proof-of-Work Puzzle

Attacker

Defender



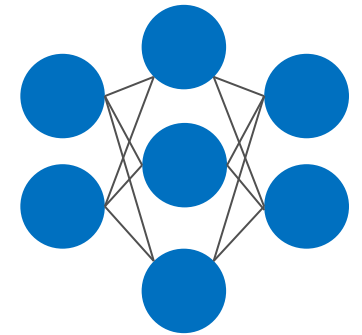
Unlabeled
Data

Send Queries



Generate
Puzzle

Estimate
Privacy
Leakage

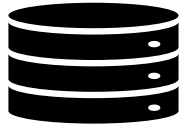


Victim
Model

Increase the Cost of Model Stealing

Attacker

Defender



Unlabeled
Data



Solve
Puzzle

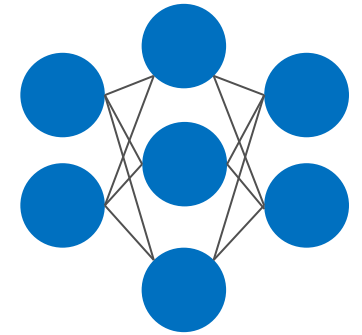
Send Queries



Generate
Puzzle

Verify
Solution

Estimate
Privacy
Leakage

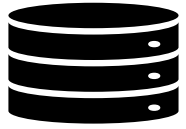


Victim
Model

Client Receives Labels after Solving a Puzzle

Attacker

Defender



Unlabeled
Data

Send Queries

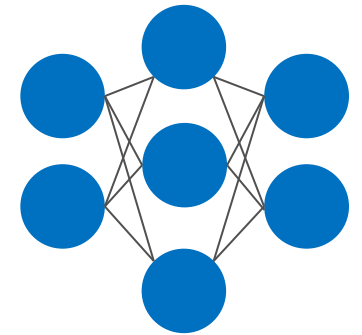


Solve
Puzzle

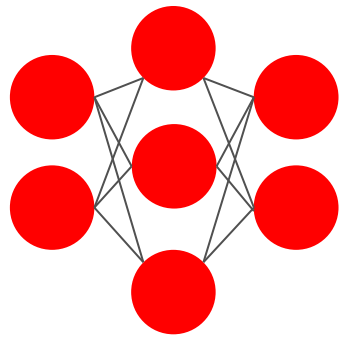


Generate
Puzzle

Estimate
Privacy
Leakage



Victim
Model



Stolen
Model

Verify
Solution

Return
Labels

Generate Puzzles using Binary Hash Cash

Server:

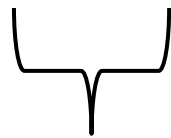
Send challenge S to client

Client:

Find a suffix X such that

$\text{hash}(S.\text{append}(X)) =$

$0\dots 0xxxxx$



required # of zeros

Generate Puzzles using Binary Hash Cash

Server:

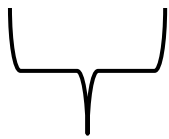
Send challenge S to client

Client:

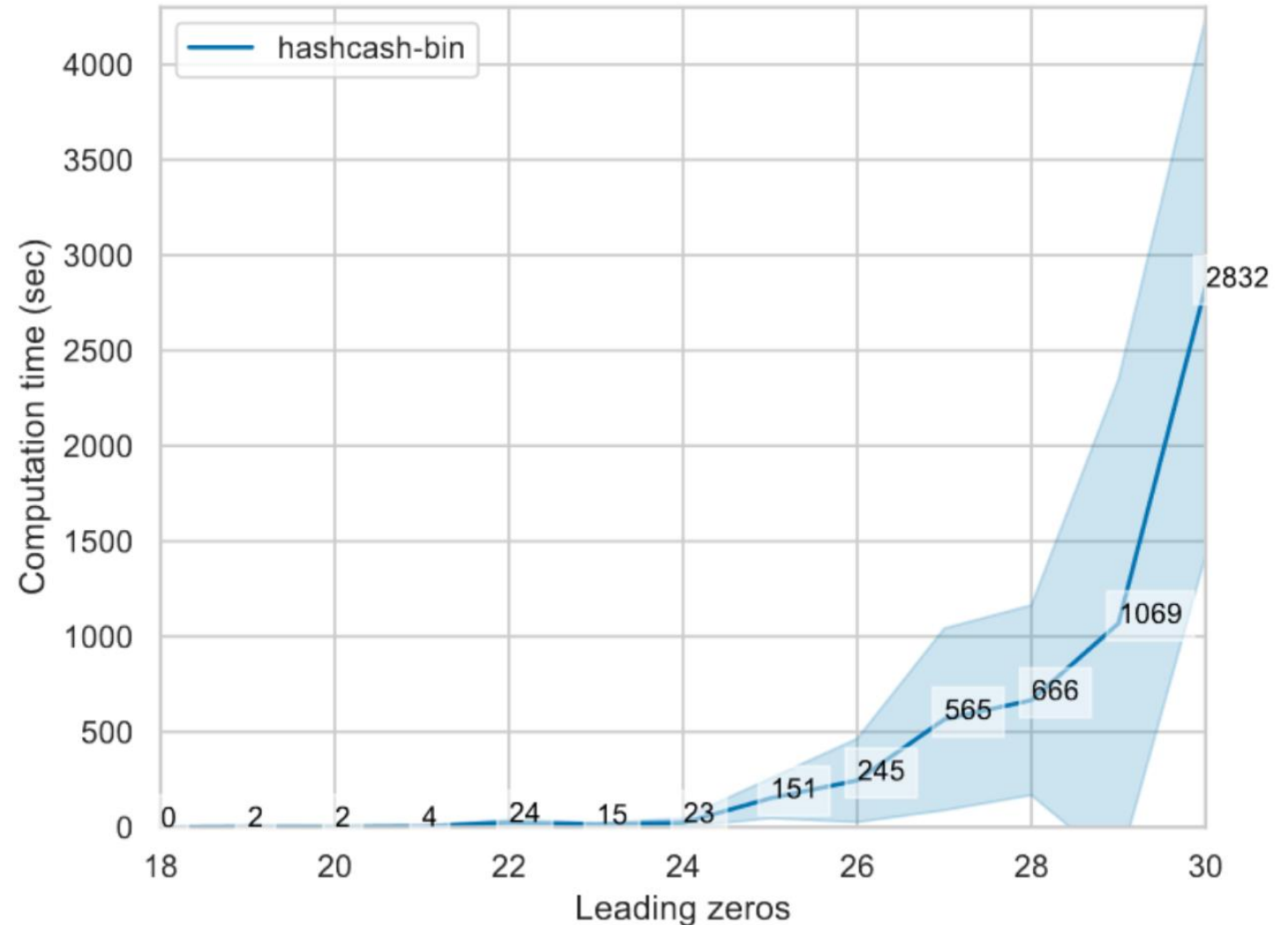
Find a suffix X such that

$\text{hash}(S.\text{append}(X)) =$

$0\dots 0xxxxx$

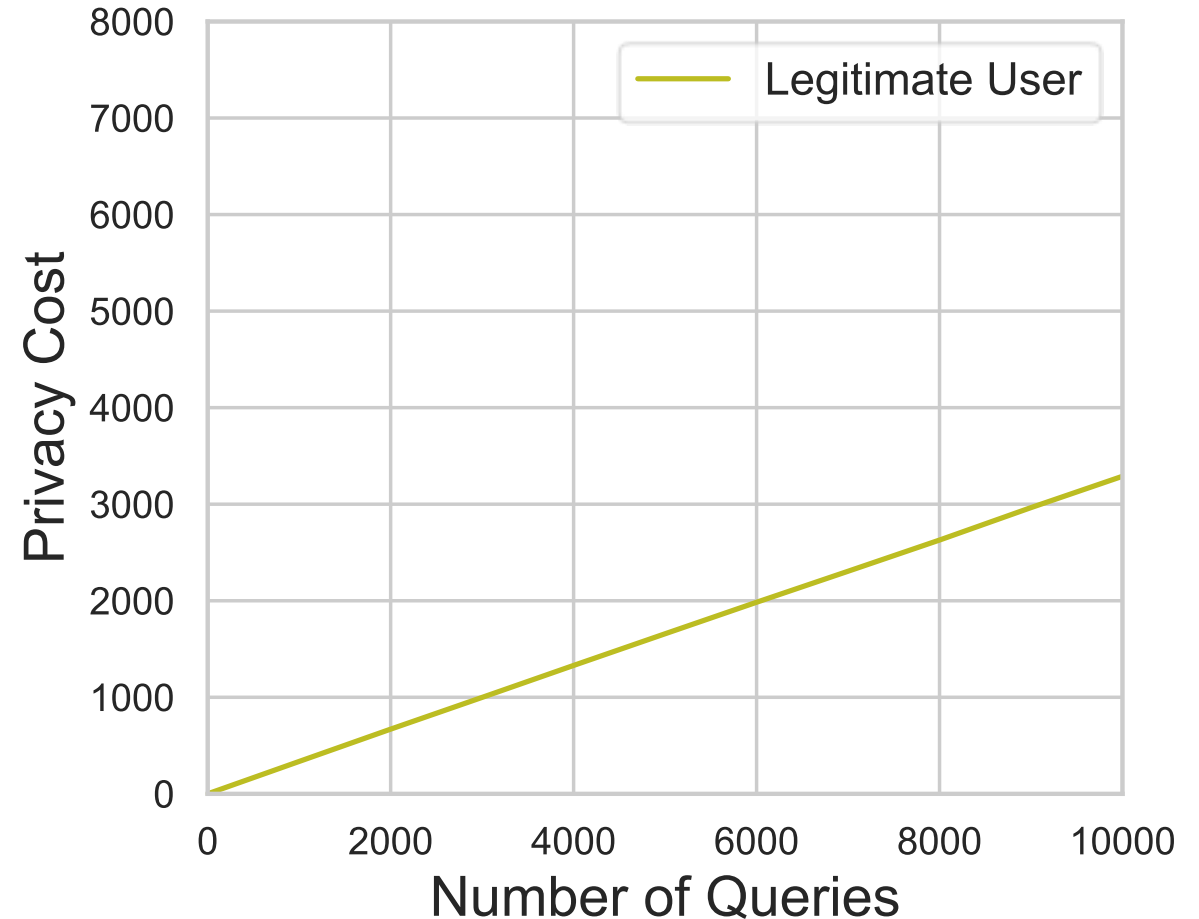


required # of zeros



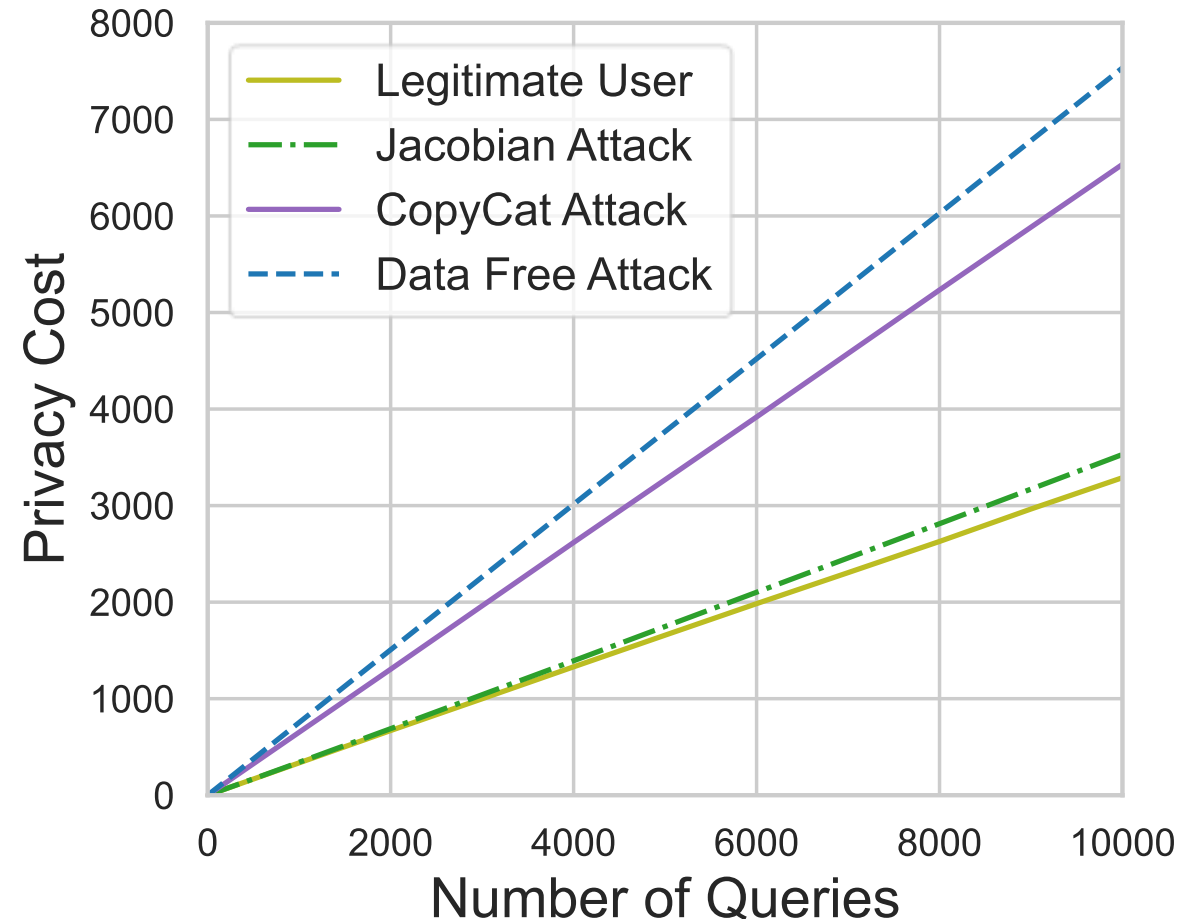
Calibrate Puzzle Difficulty using Privacy Cost

1. Set privacy cost for **Legitimate Users** as a reference cost.



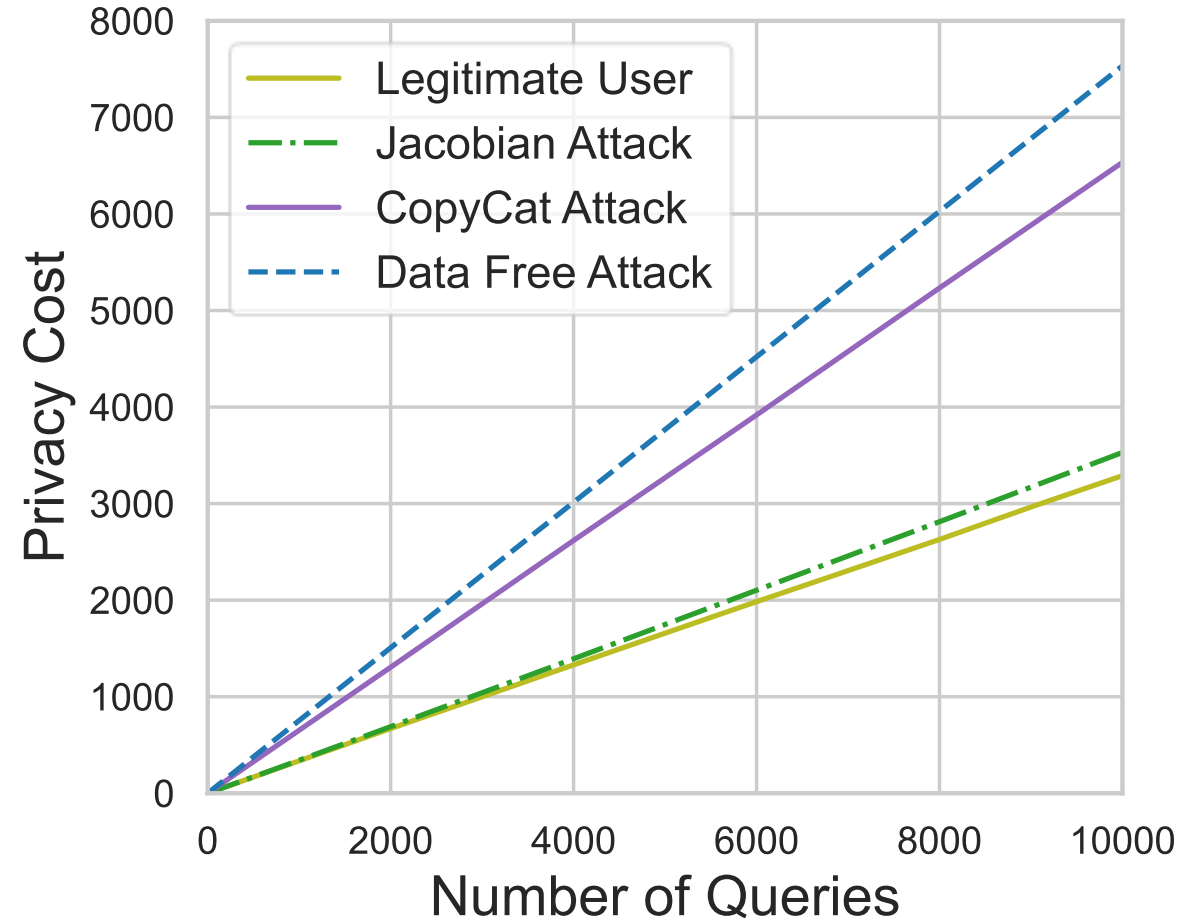
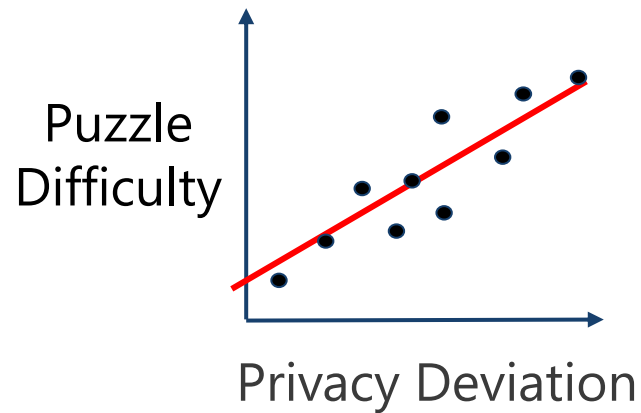
Higher Privacy Cost for Standard Attacks

1. Set privacy cost for **Legitimate Users** as a reference cost.
2. Measure the privacy cost of queries.

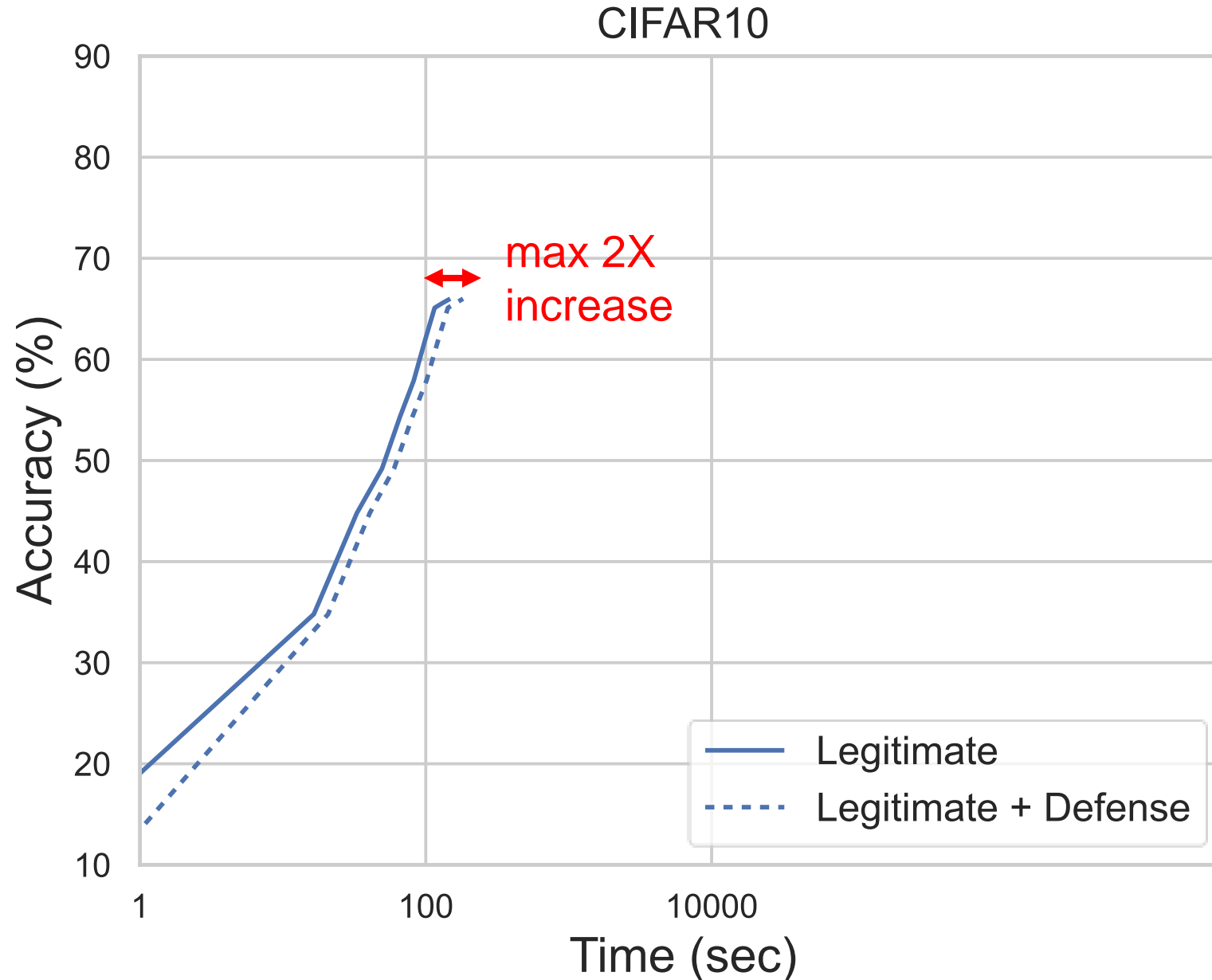


Set Puzzle Difficulty using Privacy Deviation

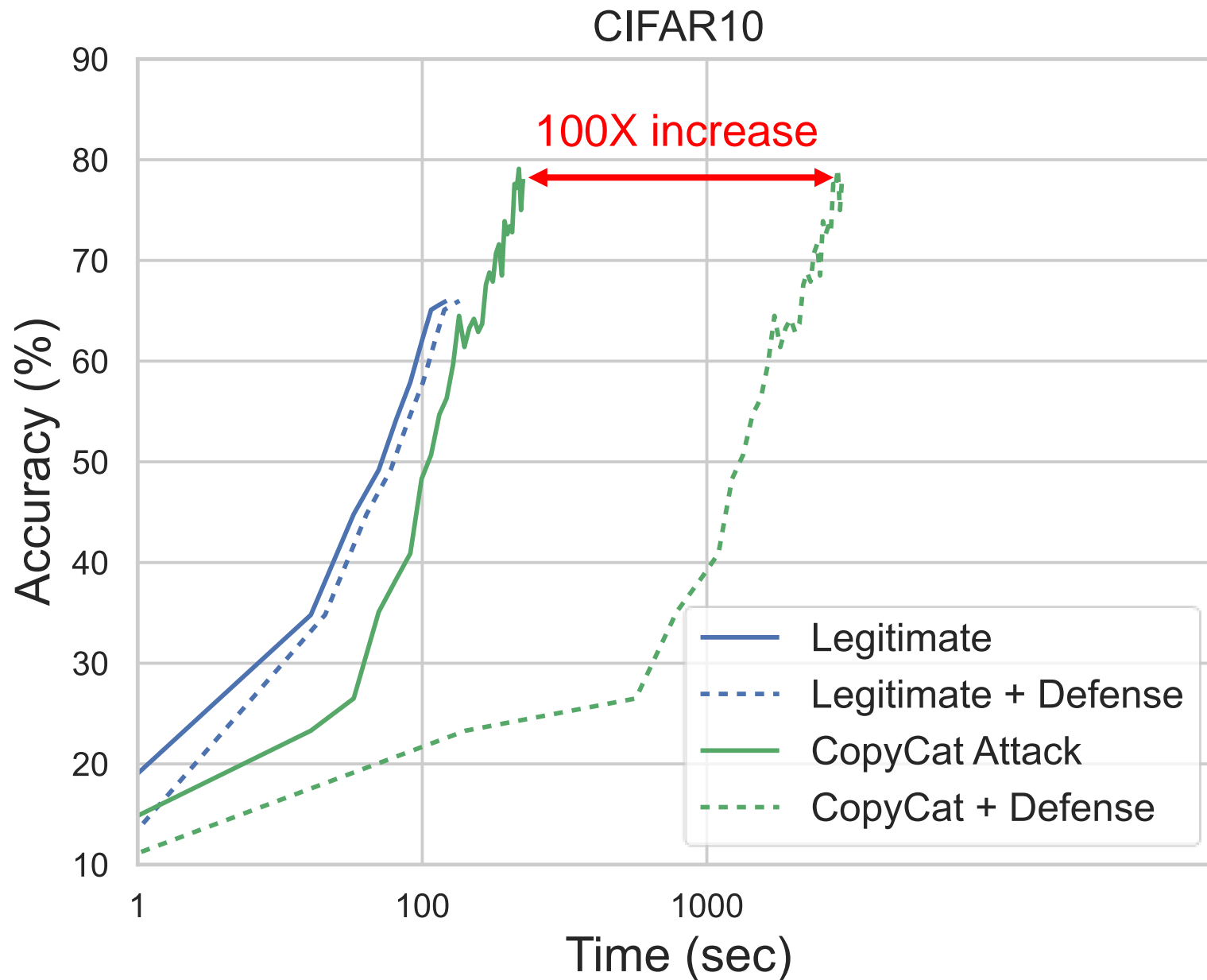
1. Set privacy cost for **Legitimate Users** as a reference cost.
2. Measure the privacy cost of queries.
3. Calibrate puzzle difficulty using privacy deviation.



Query Time vs Accuracy of Stolen Copy



Query Time vs Accuracy of Stolen Copy



Conclusions

- **New defense against model stealing** – increase the computational cost instead of lowering the quality of model outputs.
- **Privacy cost** is used to measure information leakage from a victim model that was incurred by queries from each user.
- **Calibrate** the cost of users' queries using the privacy cost.
- Use proof-of- $\{\text{work, elapsed time, stake}\}$, or payment for queries.
 - Reference method: require a user to solve the proof-of-work puzzle before releasing predictions.
- **Performance:**
 - Negligible overhead for legitimate users ($\sim 2X$);
 - High increase in the querying time for many attackers (even $\sim 100X$).

Thank you

 <https://cleverhans-lab.github.io>

 {adam.dziedzic,nicolas.papernot}@utoronto.ca