

Welcome to MIT 18.S996/6.S976

Cryptography and Machine Learning:
Foundations and Frontiers

Feb 3, 2026

Who are we?

INSTRUCTORS	<u>Shafi Goldwasser</u> Email: shafi at csail dot mit dot edu
	<u>Vinod Vaikuntanathan</u> Email: vinodv at csail dot mit dot edu
TA	<u>Neekon Vafa</u> Email: nvafa at mit dot edu

Course Syllabus

<https://mlcrypto.mit.edu/course/>

- Module 1: Crypto and ML Basics
- Module 2: Watermarking
- Module 3: Verification
- Module 4: Robustness and Alignment
- Module 5: Privacy and Security
- Module 6: Special Topics and Projects

Today: Introduction to the Cryptographic Lens
on Machine Learning

Assignment and Grading

- problem sets (25%)
- scribe notes (20%)
- class participation (10%)
- a final project (45%)

The ML Revolution

Deep Learning

Large Language models

- Human vs. Machine intelligence?
- LLM \Rightarrow ability to translate **non-human** communication ?
- Super Human Intelligence? Betting Markets



Rapid Adoption in Applications

- **Infrastructure:** Traffic patterns and energy usage
- **Health:** disease control predictive analytics using varied data
- **Drug Discovery:** acceleration using Generative AI
- **Financial Institutions:** predict risk, loans
- **Policing:** which neighborhoods to police
- **NLP:** Speech, Language, Machine Translation
- **Mathematics and Science:** AI assisted

Should we **TRUST** models we don't understand or Control

Focus on achieving Reliability, Trustworthiness

Cryptography: Enables **TRUST** in technology Even when **adversaries** are present

Arsenal of Tools: Public-Key Encryption, Digital Signatures, Zero-Knowledge Proofs, Proofs of Work, Deniable Encryption, Secure Collaboration, Homomorphic encryption, Program Obfuscation.

Remarkable Story of Theory to Practice Impact



Crypto recipe/principles for building trust

Define Task

Model Adversary



*Define Security of
a Solution*

Build Crypto Primitive

Security Proofs:

- *solution is secure
if assumption holds*

✓ *Computational Hardness*

○ *Not Everyone Colludes*

○ *Physical Assumption*

○ *Trusted Hardware*

Win Win Paradigm

Either solution is secure

Or Assumption is broken

Silvio Micali: “Either way, science wins”

Adversarial Models, Definitions, Proofs (as reductions)

No Security through obscurity

Crypto recipe/principles for building trust

Define Task

Model Adversary



*Define Security of
a Solution*

*Show impossible to
achieve primitive*

Security Proofs:

- *Any solution is insecure
if assumption holds*
- ✓ *Computational Hardness*
- *Not Everyone Colludes*
- *Physical Assumption*
- *Trusted Hardware*

Lessons from Impossibilities

- Weaken your definition of security
- Weaken the adversary model
- Find new class of assumptions

Proposal: address **ML TRUST** questions using crypto inspired paradigms, tools, assumptions, recipee

Define ML Task

Model ML Adversary



Define "Trustworthy Solution"

Build Solution

Focus on Theory + Proofs

Solution is

Trustworthy if

Assumption holds

- ✓ *Computational Hardness*
- *Not Everyone Colludes*
- *Trusted Hardware*

Proposal: address **ML TRUST** questions using crypto inspired paradigms, tools, assumptions, recipee

Define ML Task

Model ML Adversary



Define "Trustworthy Solution"

Build Solution

Or Show when impossible

Focus on Theory + Proofs

Any Solution is not

Trustworthy if

Assumption holds

- ✓ *Computational Hardness*
- *Not Everyone Colludes*
- *Trusted Hardware*

Prepare for Worst Case Adversary Strategy

AI systems are VERY attractive targets



- **Adversarial modeling:**

- **Prepare for worst case** adversary
- Do assume computational limits on adversary time.

cryptographically inspired

Assumptions: Computational Hardness

One Way Functions Exist

- $F: \{0,1\}^* \rightarrow \{0,1\}^*$ such that:
- There exists polynomial time algorithm to compute F
- All polynomial time algorithms **Inv** fail to invert F with non-negligibly probability

- $F = \{f_n: \{0,1\}^n \rightarrow \{0,1\}^{n'}\}$
- Poly time algorithm in n .
- $\Pr_{x \text{ of length } n} [\text{Inv}(y) \in f_n^{-1}(y) \mid y=f_n(x)] < 1-\text{neg}(n)$

$\text{neg}(n) < 1/\text{poly}(n)$ for all n sufficiently large

Assumptions: Computational Hardness

One Way Functions Exist

- $F: \{0,1\}^* \rightarrow \{0,1\}^*$ such that:
- There exists polynomial time A algorithm to compute F
- All polynomial time algorithms InvI fail to invert F with non-negligibly probability

If F exists then strong PSRG exist

Strong: sequences indistinguishable from random sequences By any probabilistic polynomial time algorithm (PPT)

If PSRG exists then strong PSRF exist

Strong: functions indistinguishable from random functions

By any PPT algorithm which can query the function on inputs of its choice

If strong PSRF exist then secure Enc, MAC, watermarking Schemes Exist

Assumptions: Computational Hardness

- One Way Functions Exist
- $F: \{0,1\}^* \rightarrow \{0,1\}^*$ such that:
- There exists polynomial time A algorithm to compute F
- All polynomial time algorithms nvl fail to invert F with non-negligibly probability

Examples of F

Number Theory

- $F(x,g,p)=(g^x \bmod p, g, p)$, p prime, $1 < x < p$, g generator of \mathbb{Z}_p^*
- $F(x,n) = (x^3 \bmod n, n)$ where $n=pq$, p, q primes

Geometry

- Approximating short vectors in an integer lattice.

Learning Problems

During the course

Assumption: Bounded Collusions

- Multiple Parties n
- Adversary: colluding adversaries
- Assumption: Less than t collude
 - $t=1$
 - $t < n/3$
 - ...

- Sometime enables proving
Information theoretic security

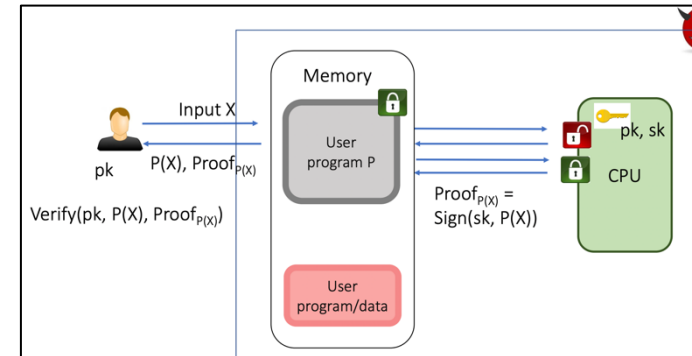
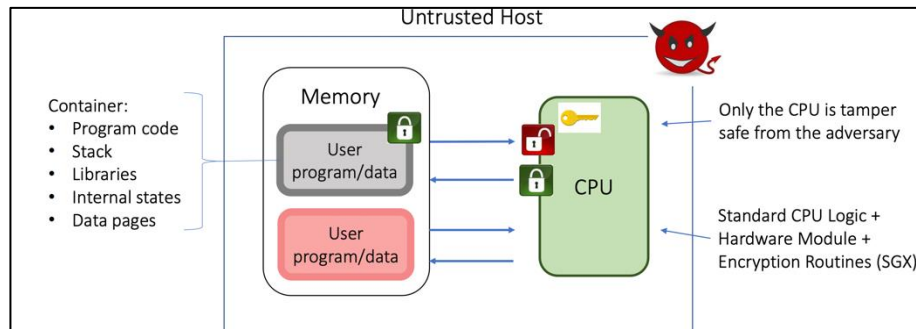
Adversary (colluding parties):

worst case malicious deviations,
curious but honest,
use bad randomness

Assumption: Secure Hardware or Trusted Execution Environment (TEE)

NVIDIA, H100 GPU, Confidential Computing Hardware 2023

Promise: high performance AI confidential compute: inference, fine tuning, mpc training. Available in cloud.



INTEL SGX, Confidential Computing Hardware 2015



Beware: side channel attacks, bugs, interrupt effects
Not trust companies blindly.

Assumption: Quantum Devices

Having Drawn Parallel between ML and Cryptography: Prepare to think **differently**

- Different models
- Different goals
- Different adversaries.
- New Hard Problems
- New Tools
- Crypto (and Complexity)
 - Theory to Practice
 - Computations **over Finite Fields**
- AI
 - Empirically Driven
 - Optimistic
 - Computation **over the Reals**

Need new ideas

ML Challenges Addressed using Crypto Lens



Verification: should verify that models satisfy properties: correctness, fairness, data usage



Robustness: test/inference data distributions may (arbitrarily) differ from training data distributions, what guarantees can you make? What can adversary do: training Poisoning



Alignment and safety: Is it possible to achieve alignment by external filters? Is inference time compute necessary?



Privacy: Power of ML comes from legally protected **training Data** of individuals, or of multiple organizations, can we train/fine-tune maintain privacy of data? Can we use ML models without using privacy of our queries



Ownership: How to watermark LLM outputs, p
How prevent model stealing, How to detect model stealing

What Type of Cryptographic/Complexity Theory Tools?



Verification: should verify correctness, fairness, data

Interactive Proofs, debate systems
New tools: PAC-Verification, Self-Proving Models



Privacy: Power of ML come individuals

Fully Homomorphis Encryption, Multi Party Computatoin, Federated Learning, Private Information Retrieval/



Robustness: data distributi data distributions, what gu

Cryptographic Backdoors, Random-self reductions



Alignment and safety: Is it p filters ?Is inference time com

Time Lock Puzzles, Stenography, Hard Learning Tasks

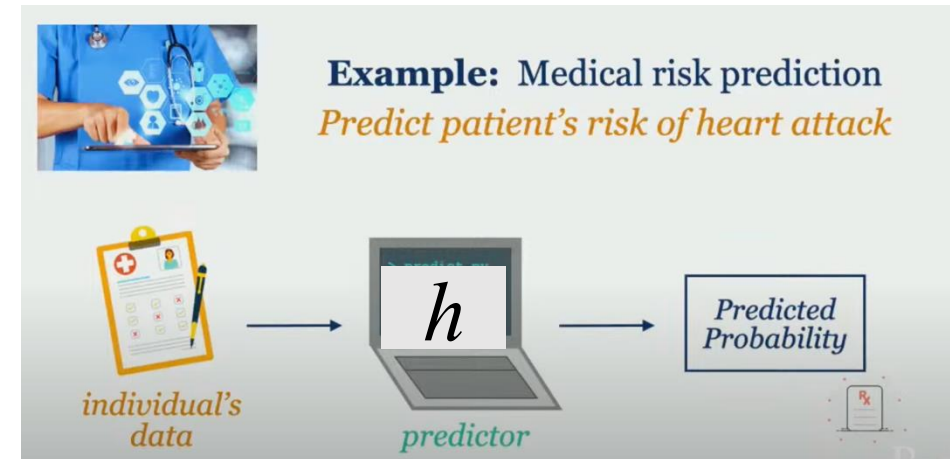
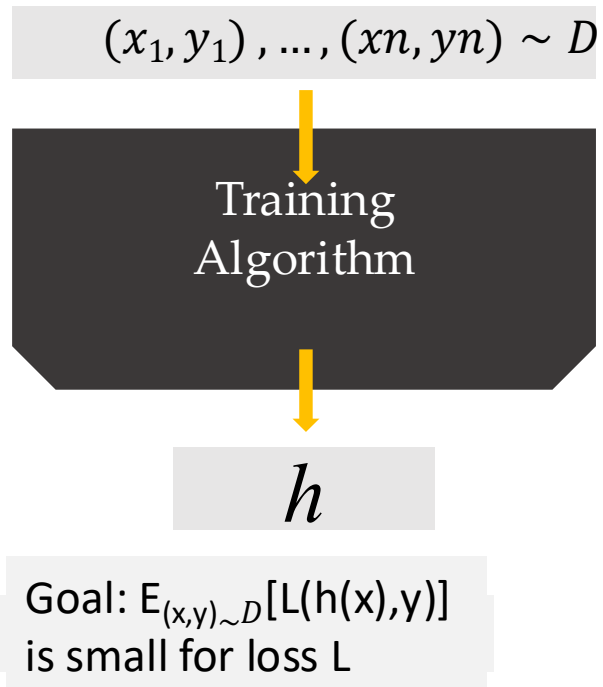
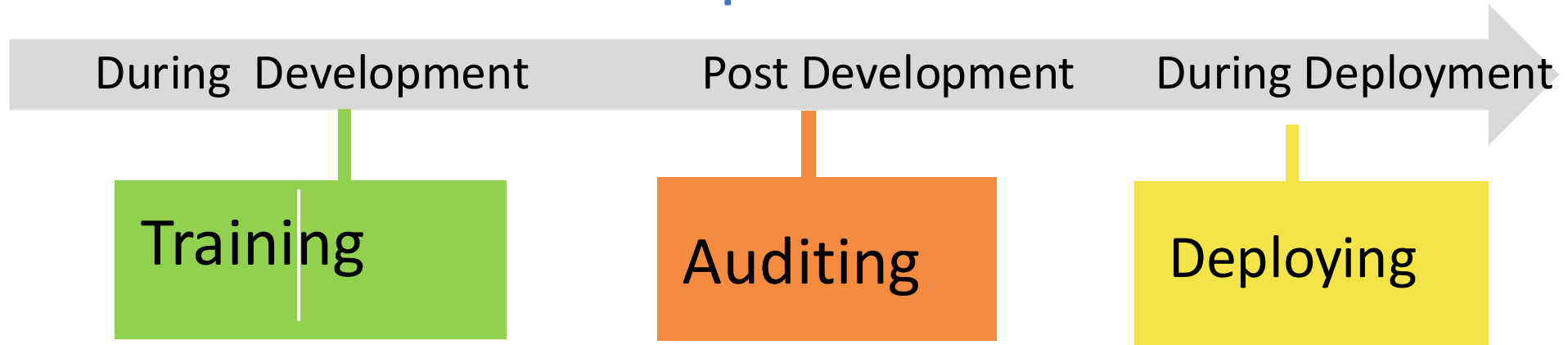


Ownership: How to watern

Pseudorandommmness, non-malleable codes, Model Distillation, Copy Right

Computational Indistinguishability \cong

Adversaries in ML Pipeline



Prediction/answer Generation/
distribution over answers

Theory Approach

During Development

Post Development

During Deployment

Theory vs. Practice

Adversaries apply to both

Definitions apply to both

Methods (in principle) could apply to both

Issue: Efficiency at Scale

Empirical Studies Needed (projects)

Privacy at TRAINING

Privacy at Training

During Development

Post Development

Into the Future

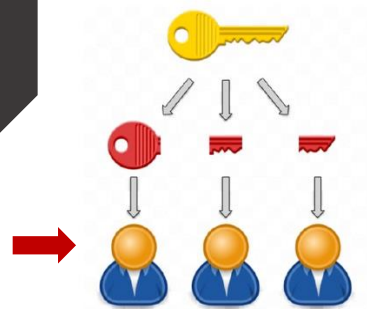
Train

Use existing data to
build ML model

$Enc(x_1, y_1) \dots Enc(x_n, y_n) \sim D$

Run training algorithm
without ever decrypting
training data

$Enc(h)$



- (1) Encrypted Compute Stage
- (2) Decrypt stage

Assumptions:

LWE is Hard

+

Key Share Holders

don't collude

Plaintext
world

Ciphertext
world

Enc

$c \xrightarrow{Eval_f} c'$

h

The Learning with Errors Problem (LWE) [Regev05]

- Let \mathbf{s} be a secret vector in \mathbb{Z}_q^n
- Given an arbitrary number of “noisy” equations in \mathbf{s} , find \mathbf{s} ?

$$14s_1 + 15s_2 + 5s_3 + 2s_4 \approx 8 \pmod{17}$$

$$13s_1 + 14s_2 + 14s_3 + 6s_4 \approx 16 \pmod{17}$$

$$6s_1 + 10s_2 + 13s_3 + 1s_4 \approx 3 \pmod{17}$$

$$10s_1 + 4s_2 + 12s_3 + 16s_4 \approx 12 \pmod{17}$$

$$9s_1 + 5s_2 + 9s_3 + 6s_4 \approx 9 \pmod{17}$$

$$3s_1 + 6s_2 + 4s_3 + 5s_4 \approx 16 \pmod{17}$$

$$6s_1 + 7s_2 + 16s_3 + 2s_4 \approx 3 \pmod{17}$$

- ✓ As **hard as**: Decoding Random Linear Codes =
- ✓ As **hard as**: approximating the size of the shortest vector in a worst-case n -dim integer lattice

The Learning with Errors Problem (LWE) [Regev05]

- Let s be a secret vector in \mathbb{Z}_q^n
- Given an arbitrary number of “noisy” equations in s , find s ?

$$14s_1 + 15s_2 + 5s_3 + 2s_4 \approx 8 \pmod{17}$$

$$13s_1 + 14s_2 + 14s_3 + 6s_4 \approx 16 \pmod{17}$$

$$6s_1 + 10s_2 + 13s_3 + 1s_4 \approx 3 \pmod{17}$$

$$10s_1 + 4s_2 + 12s_3 + 16s_4 \approx 12 \pmod{17}$$

$$9s_1 + 5s_2 + 9s_3 + 6s_4 \approx 9 \pmod{17}$$

$$3s_1 + 6s_2 + 4s_3 + 5s_4 \approx 16 \pmod{17}$$

$$6s_1 + 7s_2 + 16s_3 + 2s_4 \approx 3 \pmod{17}$$

- ✓ **Post-Quantum:** Best known algorithm (even quantum) time 2^n

NEWS

NIST Announces First Four Quantum-Resistant Cryptographic Algorithms

Federal agency reveals the first group of winners from its six-year competition.

July 05, 2022

OpenFHE: Open-Source Fully Homomorphic Encryption Library^{*†}

Ahmad Al Badawi¹, Andreea Alexandru¹, Jack Bates¹, Flavio Bergamaschi², David Bruce Cousins¹, Saroja Erabelli¹, Nicholas Genise¹, Shai Halevi³, Hamish Hunt², Andrey Kim⁴, Yongwoo Lee⁴, Zeyu Liu¹, Daniele Micciancio^{1,5}, Carlo Pascoe¹, Yuriy Polyakov^{†1}, Ian Quah¹, Saraswathy R.V.¹, Kurt Rohloff¹, Jonathan Saylor¹, Dmitriy Suponitsky¹, Matthew Triplett¹, Vinod Vaikuntanathan^{1,6}, and Vincent Zucca^{7,8}

¹Duality Technologies

²Intel Corporation

³Algorand Foundation

⁴Samsung Advanced Institute of Technology

⁵University of California, San Diego

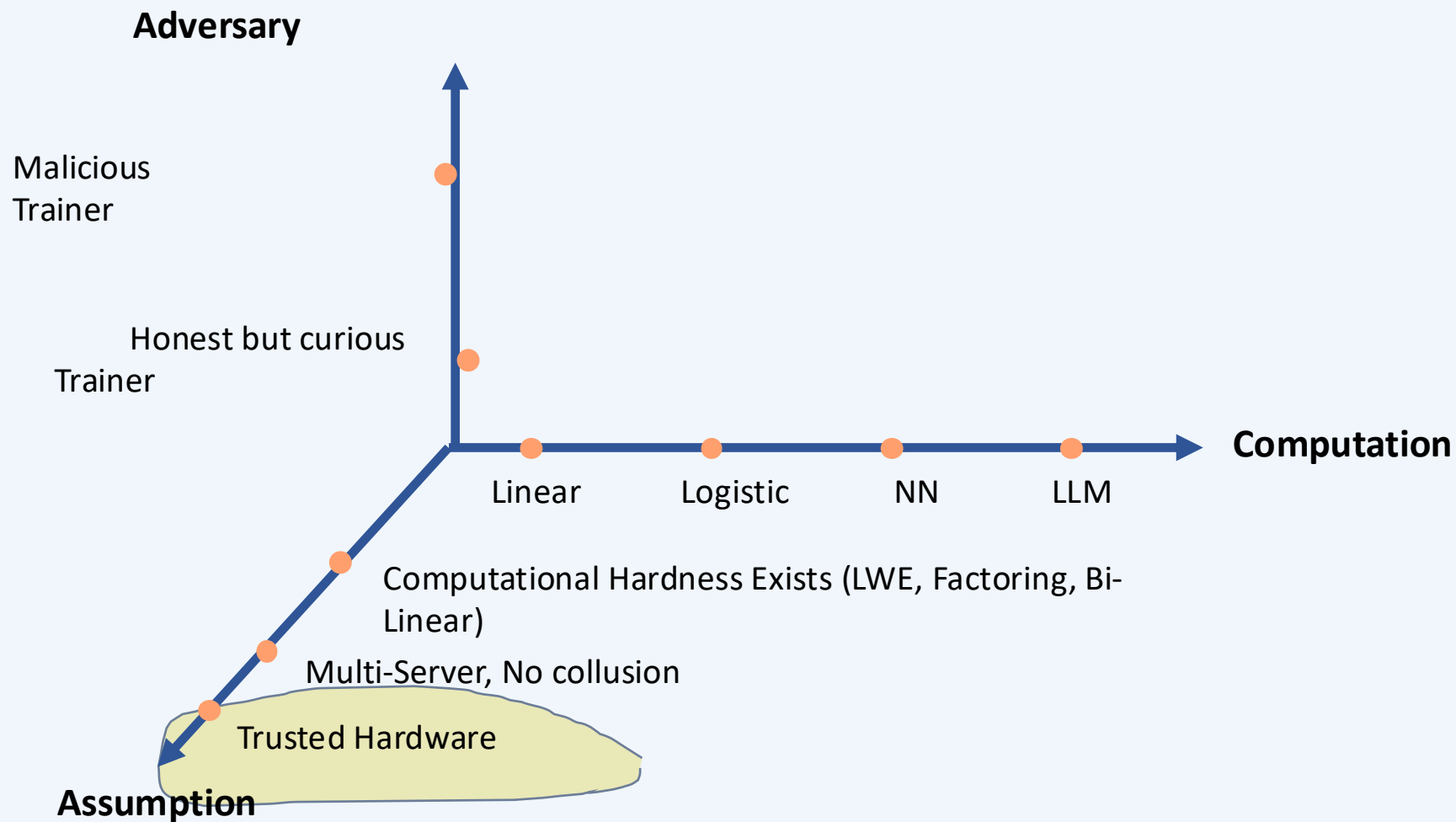
⁶Massachusetts Institute of Technology

⁷DALI, Université de Perpignan Via Domitia

⁸LIRMM, University of Montpellier

March 12, 2024

General LLM Fine Tuning? In Practice: A Scaling Challenge



VERIFICATION POST TRAINING

Part 1 Verify Model Properties

Part 2 Verify Model Answers Per Input

ML as a Service

During Development

Post Development

Auditing



Client



Service Provider

MLaaS: Amazon
SageMaker/AWS,
Microsoft Azure,
Startups...

Verifying Model Properties, how?



Can we verify properties of the model h :

Accuracy over inputs/

Correctness per input/

Robustness/

Fairness

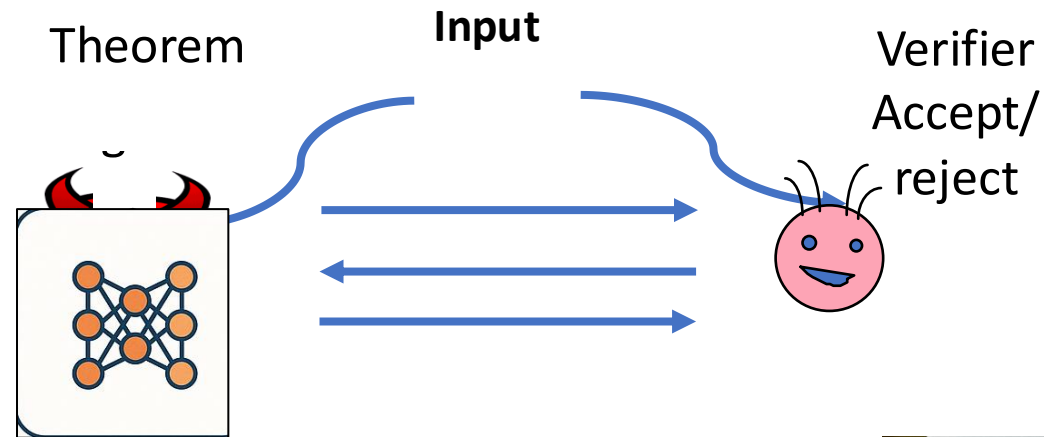
Safety

Satisfies Regulations

cheaply (not retraining) using

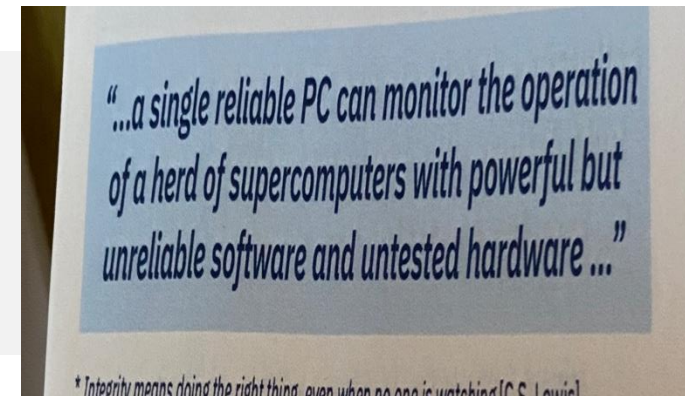
- Fewer data samples
- Lower quality data
- Efficient
Time/Memory/samples
- Black box access or limited white-box access to h

Interactive Proofs Framework 80's



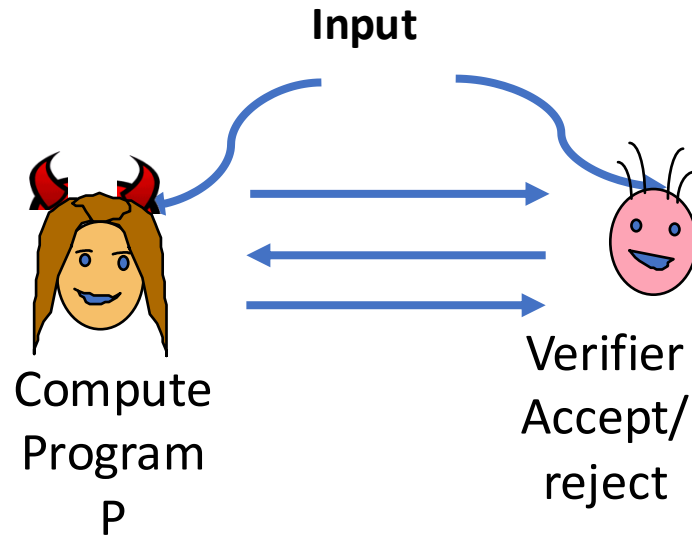
Verifiable Computing Paradigm 2000's

- **Verifying** **Cheaper** than computing: do not replicate
- **Doubly-efficient** generating the proof should not be much more costly than computing



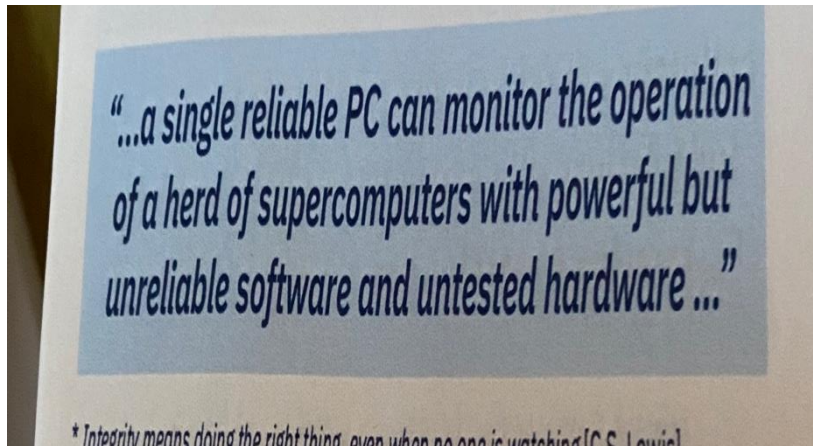
Interactive Framework of 80's (fast verification on blockchains)

Will Study in Course



Techniques

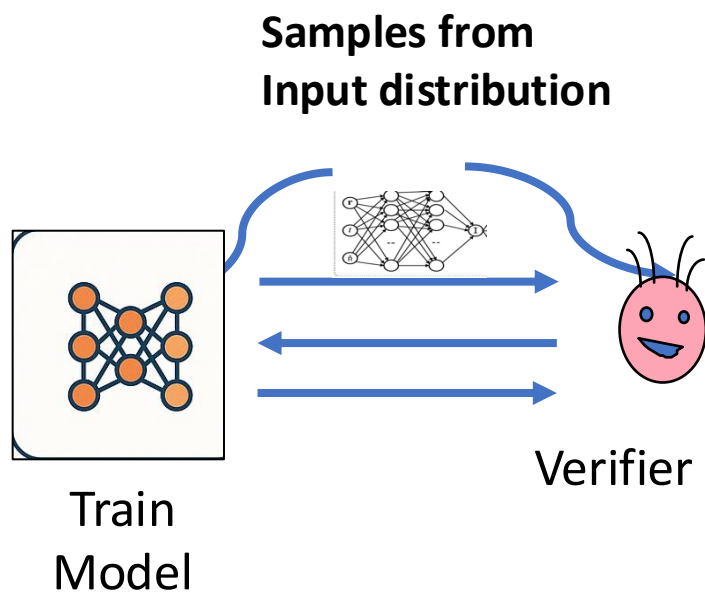
- Interactive Proofs & arguments for Program Delegation
- Zero Knowledge Interactive Proofs & arguments(
- Multi-Prover Proofs
- Debates



Verifiable Computing Paradigm 2000's

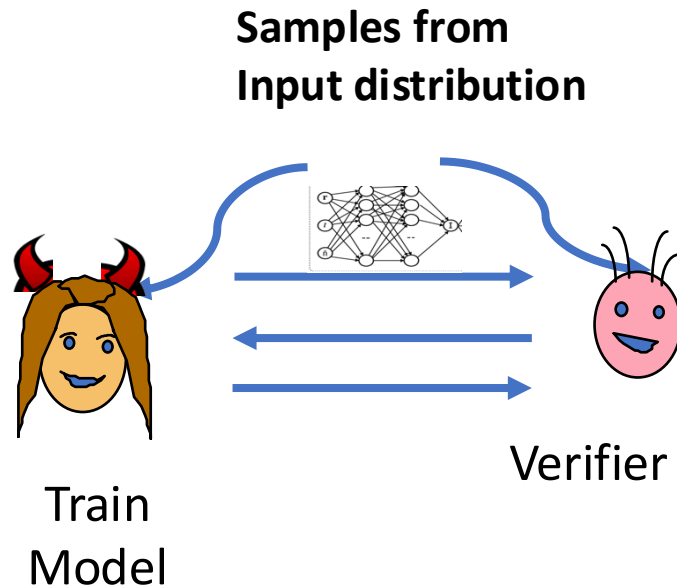
- **Verifying** **Cheaper** than computing: do not replicate
- **Doubly-efficient** generating the proof should not be much more costly than computing

ML Case is Different



- **Input:** samples from a distribution
- **Compute:** randomized, massively parallel
- **Operations:** **reals vs. finite fields**

Main Difference: Prover/Learner not Pre-Specified



- **Input:** samples from a distribution
- **Compute:** randomized, massively parallel
- **Operations:** **reals vs. finite fields**

What are you verifying?

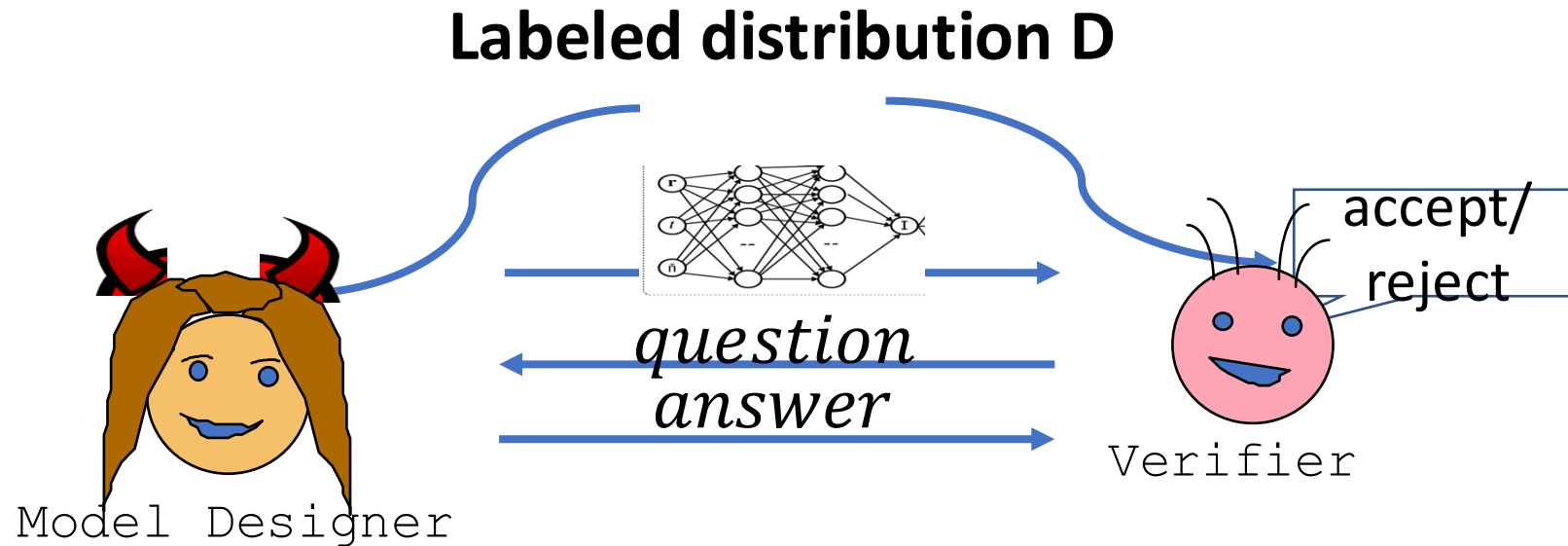
Prover is not pre-specified

Does not necessarily \Rightarrow
Model is Accurate, Robust, Aligned,
Fair, Uses Data as Prescribed

Train(Data, Randomness) = h



Pac- Verification of Model Accuracy



Probabilistic & Approximate Verification:

verify that given model is *within additive error of*
most accurate model possible model

Part 2: Proving Correctness of ML Answers

Typical Claim: LLM Model has 99% accuracy on task

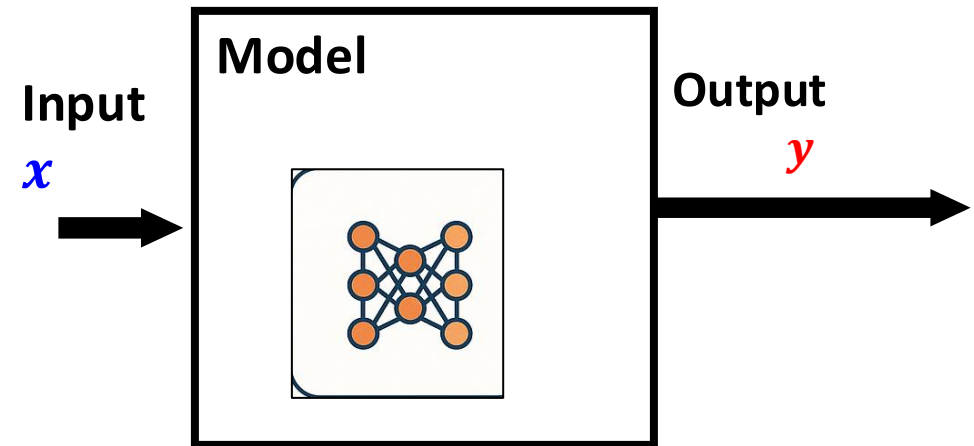
Tested on benchmarks for the task

Held out set

Learn Human Feedback(RLHF)

Stress Test/Red Team

- But on MY **medical file x** ,
the model generates **diagnosis y**

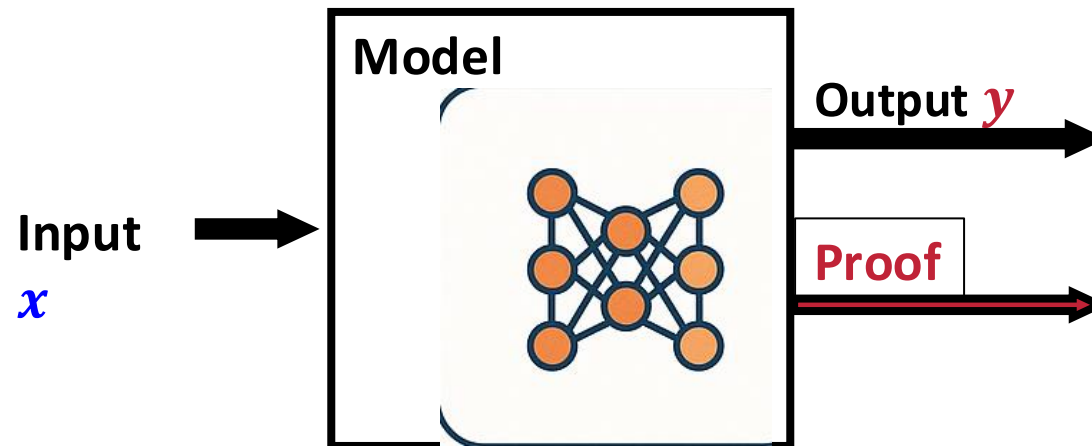


Can you verify that **y** is correct for x ?

From Average to Worst Case Guarantees: On input x , LLM Generates Proof of Correctness

Goal:

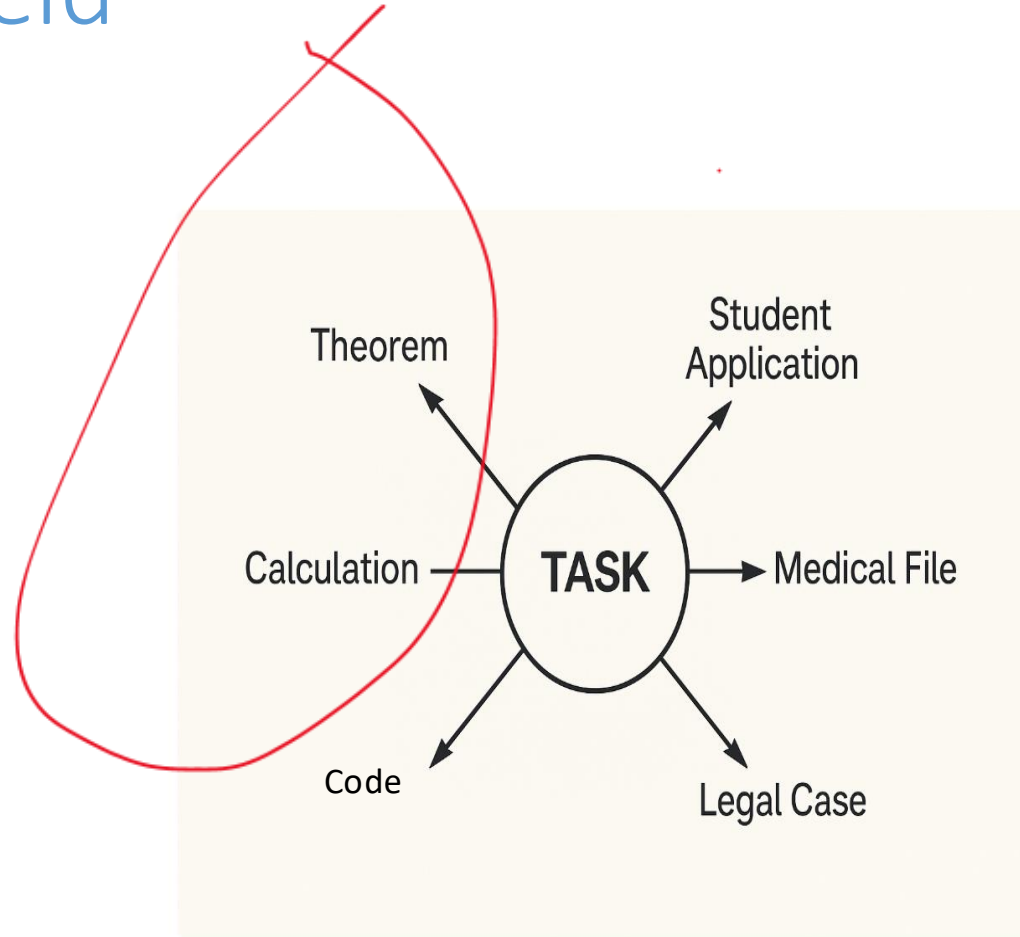
LLMs which output y + proof that is *easy* to check
that y is *correct*



What does a “correct” and “proof” mean?

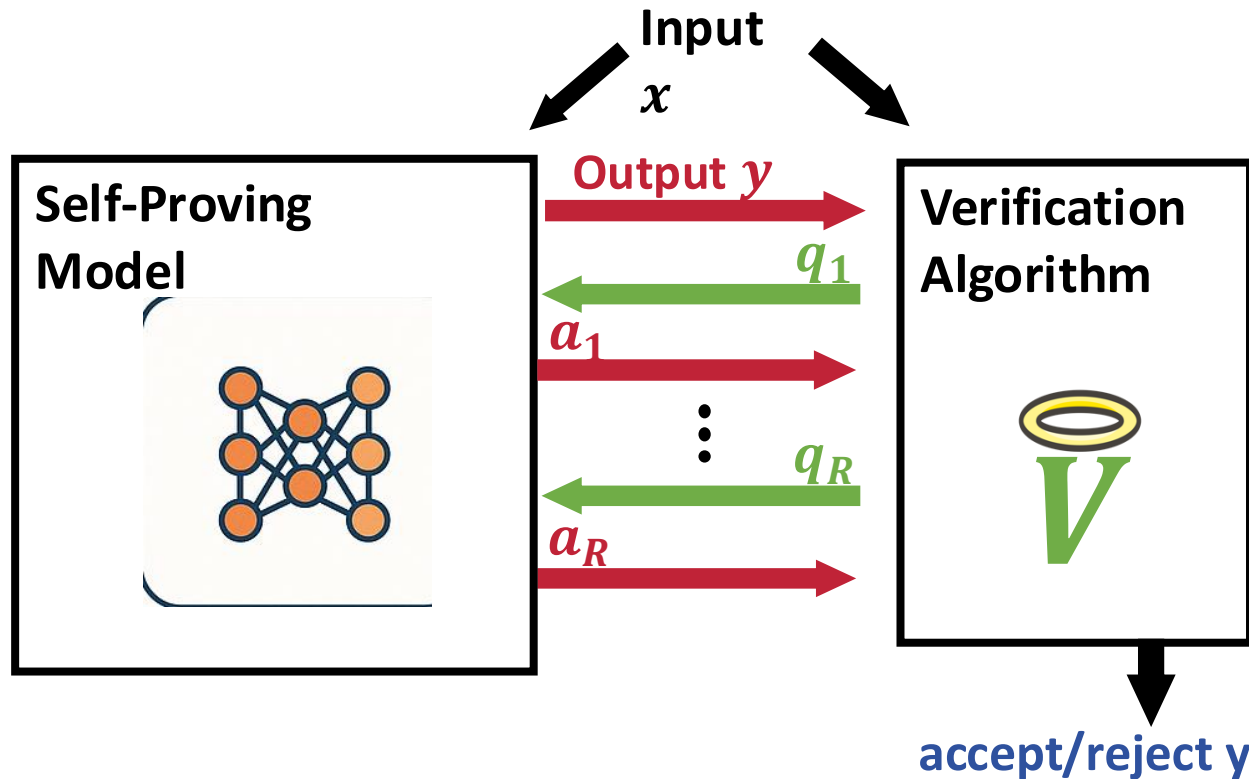
Models That Prove Their Own Correctness	
Noga Amit* UC Berkeley nogamit@berkeley.edu	Shafi Goldwasser* UC Berkeley shafi.goldwasser@gmail.com
Orr Paradise* UC Berkeley orrrp@eecs.berkeley.edu	Guy N. Rothblum* Apple rothblum@alum.mit.edu

Notion of correctness may changes from field to field



Projects?

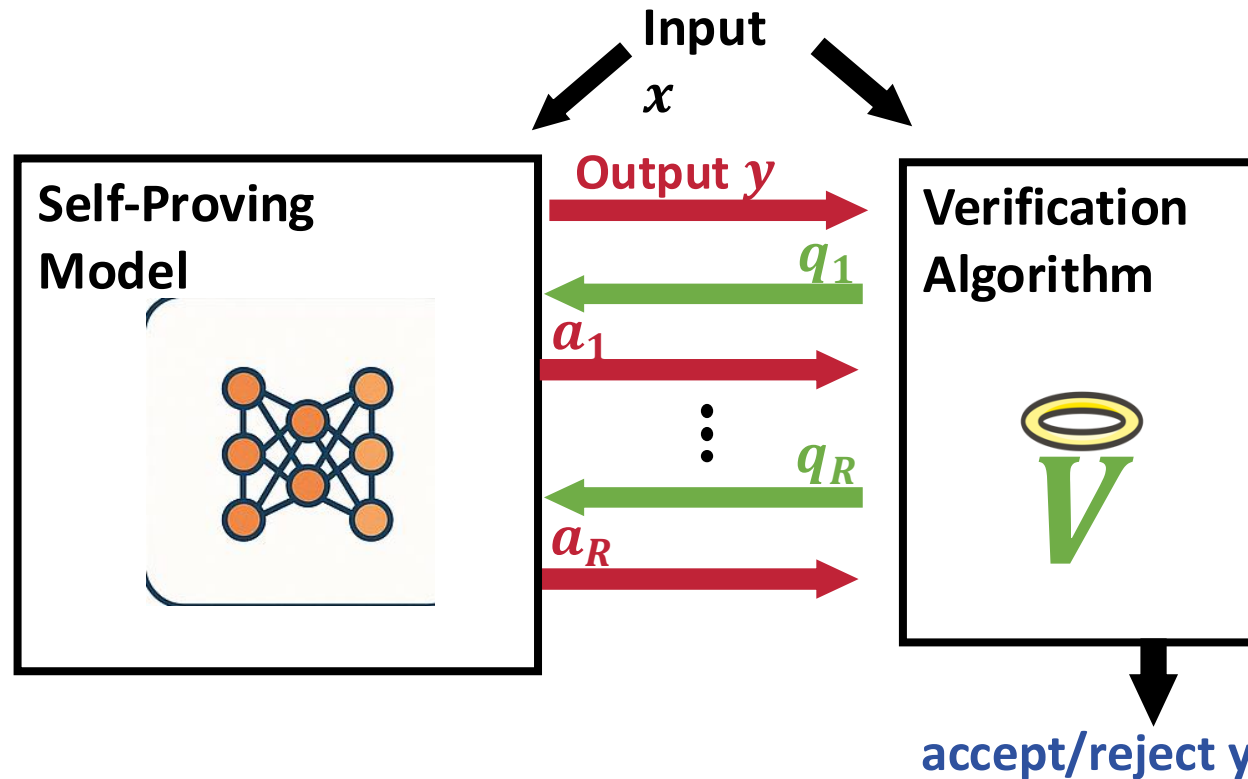
Proof: Make a Verifier Algorithm Accept



V should be
more efficient than P

V is assumed to be
verified

Interactive Proof: to a Verifier Algorithm



Soundness: for all x , V rejects incorrect y 's w.h.p. over V 's coins

Completeness. For distribution μ , $\text{Prob}_{x \in \mu}[V(x) \text{ accepts } y \text{ as correct}] > \text{high}$
Distributional Requirement

Need to train models to prove its answers to V

How? Let accepting transcripts $\pi = q_1 a_1 \dots q_l a_l$

- **Transcript Learning:** collect and train on “proof bank” of (x, y, π)
- **Reinforcement Learning from Verifier Feedback (RLVF)**

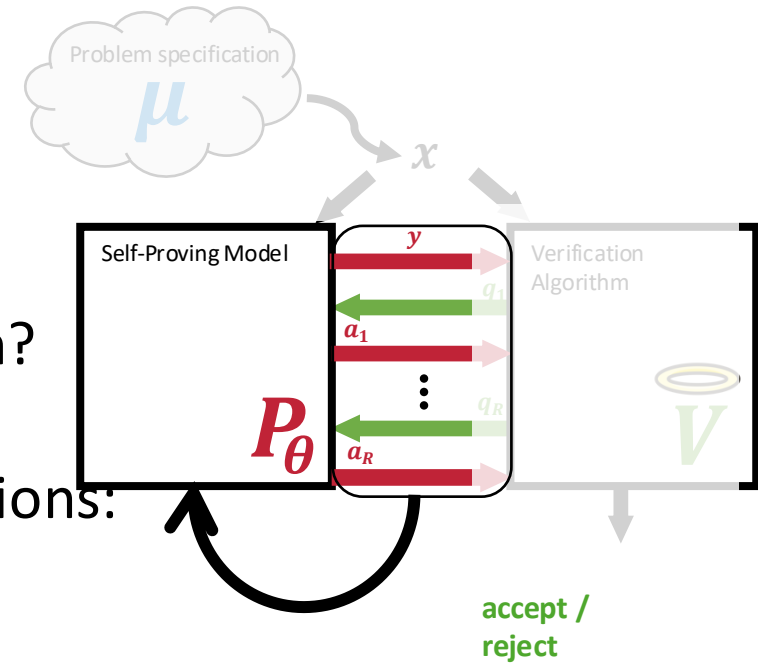
.

Goal: Learn model parameters that maximizes

$\Pr[$  convinces V to accept $y]$

RLVF in Practice

- In practice, Transcript Learning vs. RLVF is a question of **supervision**. Is RLVF enough?
- Sometimes, yes! See practical implementations:
 - **RLVR** [Lambert et al., 2024]:
 - Adds many practical improvements (KL-regularization, PPO, ...)
 - **Med-RLVR** [2025]: Medical multiple-choice questions
 - **RLVR-World** [Wu et al., 2025]: Computer vision and robotic manipulation
 - **RLPR** [Yu et al., 2025]: No more verifiers, use the LLM itself instead (full circle!)
 - **The Invisible Leash** [Wu et al., 2025]: Analyzing failure modes of RLVF/RLVR.



ROBUSTNESS IN DEPLOYMENT



Robustness to what?

- Distribution shifts
- Adversarial Examples
- Insider Adversaries

Insider Adversaries: Planted Backdoors



Client
University

Data



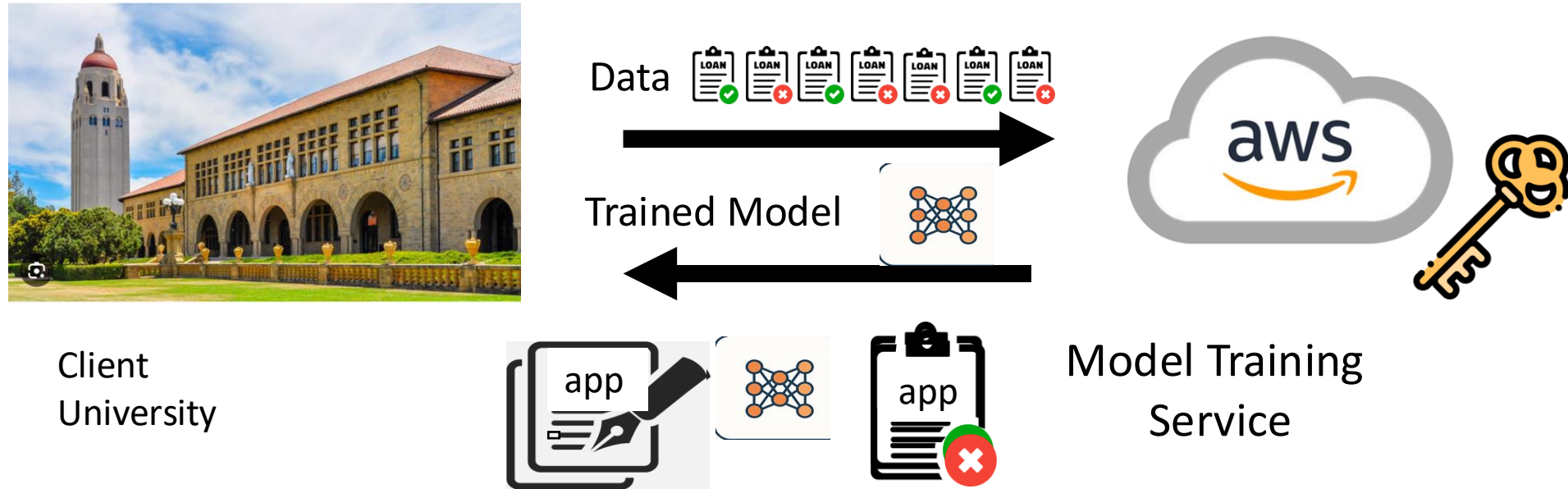
Trained Model



Model Training
Service

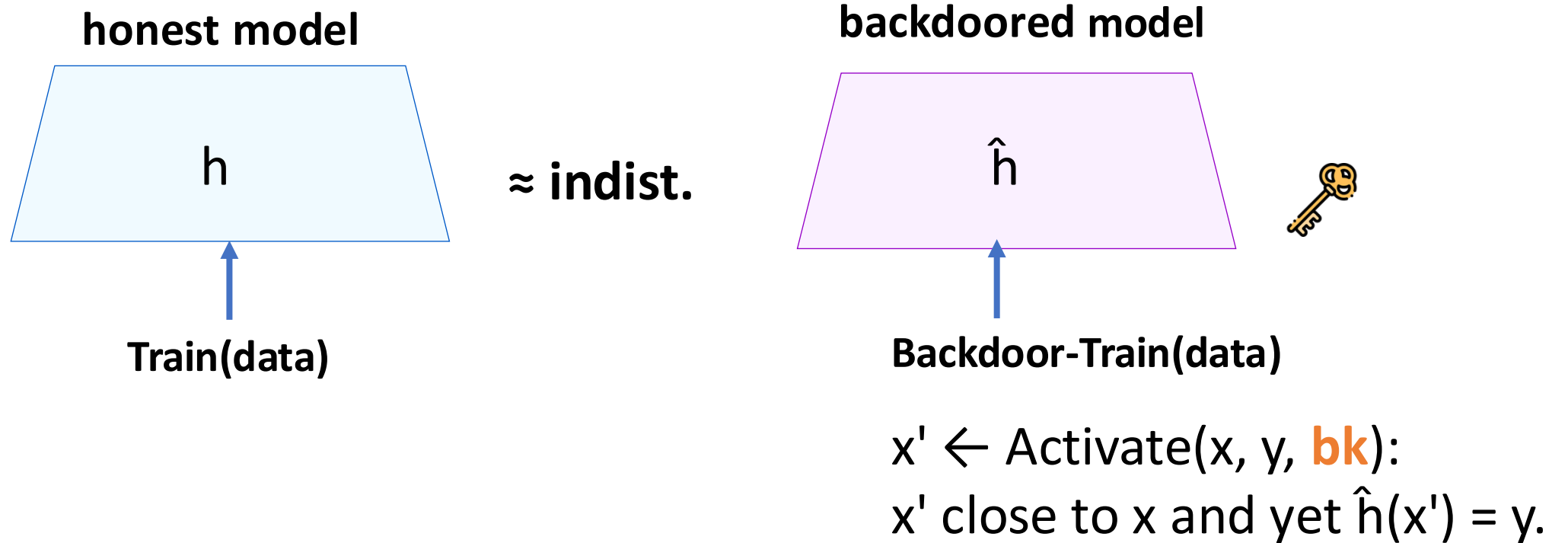


Always Possible to Embed a Backdoor Key to Trigger Different Model Behavior



Theorem: If cryptography exist, then can plant such backdoors in any neural net for classification such that the backdoors are **undetectable** & **non-replicable**

Undetectable Backdoors for Classification:



From black box access to the mode



Extensions to white-box access to restricted models

White Box Undetectable Backdoors?

Other distance measures?

Backdoors for ML Embeddings

[Bogdanov–Rosen–Vafa’25](#)] Show how to “backdoor” deep embedding networks in a statistically undetectable way s.t.

- With a backdoor, can produce semantic collisions: unrelated images with very close embeddings.
- Without backdoors, provably hard to produce collisions under CHV

New Hardness Assumption

Adaptive Robustness of Hypergrid Johnson-Lindenstrauss

Andrej Bogdanov* Alon Rosen† Neekon Vafa‡ Vinod Vaikuntanathan§

Abstract

Contracting Hypergrid Vector (CHV) Problem

Given: Gaussian $m \times n$ matrix A (zero mean, unit variance)

Find: x in hypergrid $\{-b, \dots, b\}^n$

$$\frac{1}{\sqrt{m}} \|Ax\| \leq \kappa \|x\|$$

This problem exhibits a “computational-to-statistical gap”.



* $\kappa_{stat}, \kappa_{comp}$ depend on $\alpha = m/n$ (how much you compress) and b

Removing Planted Backdoors



Client Scientist



Model Training Service

Mitigation: efficient post processing



New Model which is
**Accurate, Independent of
Tampering, no more backdoors**

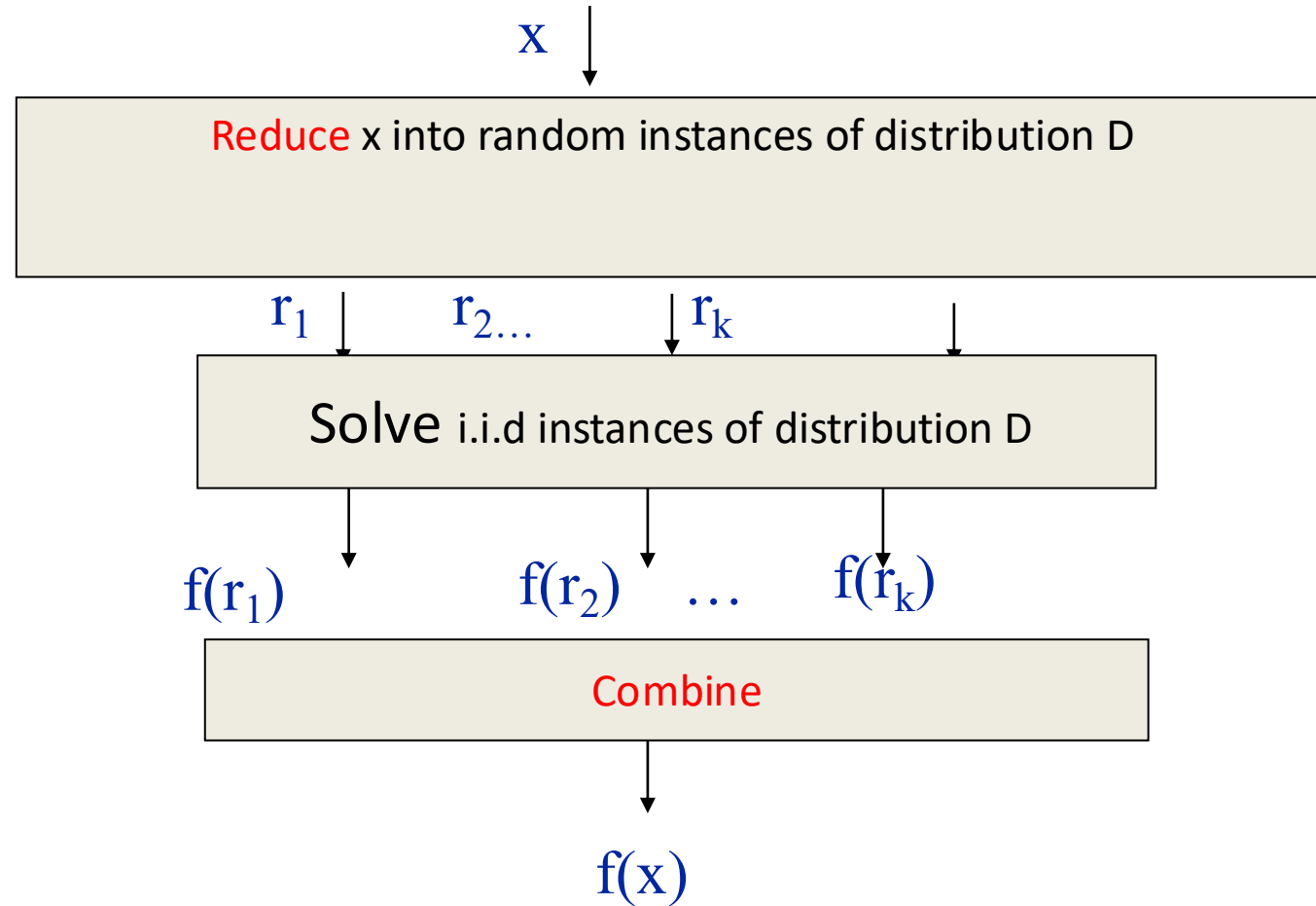
Two flavors Black Box post process

Offline: recover new model

Online: Post-process at test time.

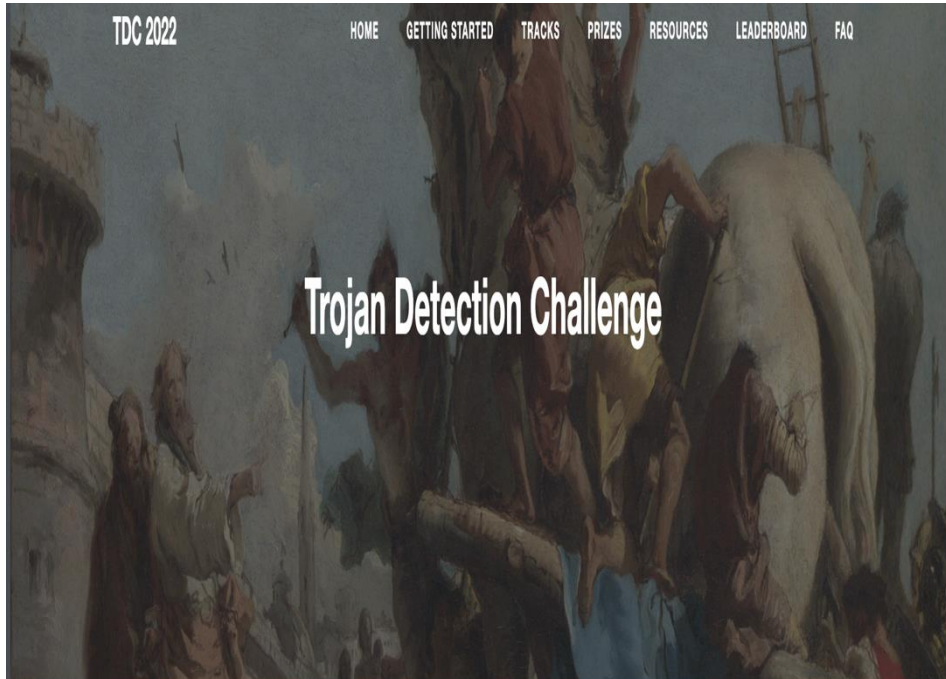
Access to (potentially) adversarial model **speeds up
learning/inference**

Key Helpful Concept from Cryptography & Complexity: Random (Self) Reducibility for f , distr. D [GM82, BK89, BLR90]



Examples: number theory problems, lattice problems,
low deg polynomials problems

From theory to practice?



In this competition, we challenge you to detect and analyze Trojan attacks on deep neural networks that are **designed to be difficult to detect**. Neural network Trojans are a growing concern for the security of ML systems, but little is known about the fundamental offense-defense balance of Trojan detection. Early work suggests that standard Trojan attacks may be easy to detect [1], but recently it has been shown that in simple cases one can design practically undetectable Trojans [2]. We invite you to help answer an important research question for deep neural networks: How hard is it to detect hidden functionality that is trying to stay hidden?

TROJAI

TROJANS IN ARTIFICIAL INTELLIGENCE

INTELLIGENCE VALUE

Artificial Intelligence (AI) is being increasingly applied to a variety of domains within the Intelligence Community (IC). The TrojAI program seeks to defend AI systems from intentional, malicious attacks, known as Trojans, by conducting research and developing technology to detect these attacks in a completed AI system. By building a detection system for these attacks, engineers can potentially identify backdoored AI systems before deployment. The development of Trojan AI detection capabilities will mitigate risks arising from AI system failure during mission critical tasks.

SUMMARY

TrojAI is researching the defense of AI systems from intentional, malicious Trojan attacks by developing technology to detect these attacks and by investigating what makes the Trojan detection problem challenging. Trojan attacks, also called backdoor attacks, rely on training the AI to attend to a specific trigger in its inputs. The trigger is ideally something that the adversary can control in the AI's operating environment to activate the Trojan behavior. For Trojan attacks to be effective, the trigger must be rare in the normal operating environment so that it does not affect the normal effectiveness of the AI and raise the suspicions of human users.

Explore planting backdoors for the TrojAI challenge - Performers test their current trojan detection approaches

Challenges Addressed using Crypto Lens



Privacy: Power of ML comes from legally protected **training Data** of individuals



Verification: should verify that models satisfy properties: correctness, fairness, data usage



Robustness: data distributions may (arbitrarily) differ from training data distributions, what guarantees can you make?



Alignment and safety: Is it possible to achieve alignment by external filters? Is inference time compute necessary?



Ownership: How to watermark LLM outputs, prevent model stealing

The Alignment Problem



Malicious Users \approx Jailbreaks

Battle between Alignment/Safety and Jailbreaks

Difficulty 1: (optimized) model objectives diverge from human objectives

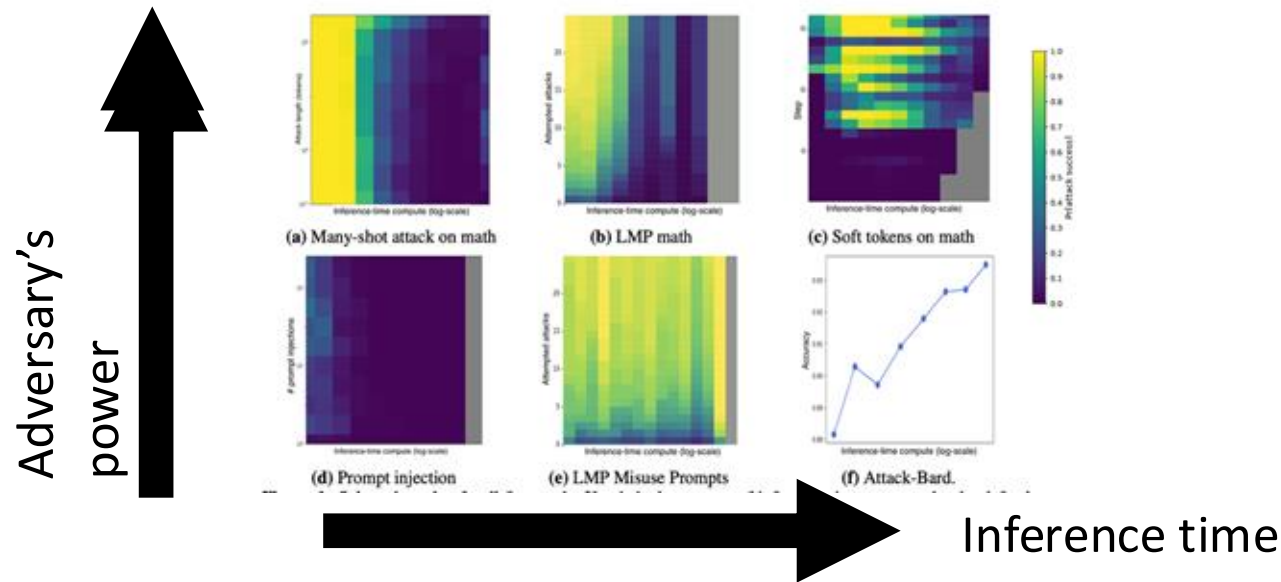
Difficulty 2: how to formalize “doing what humans want” so that it remains stable under optimization and scalable with capability.

“In adversarial machine learning, we wrote over 9,000 papers in ten years and got nowhere”

*Nicholas Carlini, 2019,
“Some lessons from adversarial machine learning”*

Defense Strategies
against Jailbreaks?

Defense Strategies against Jailbreaks for Safety



Zaremba et al., 2025, “Trading inference-time compute for adversarial robustness”

Guan, et al., 2025 “Deliberative Alignment: Reasoning Enables Safer Language Models”

Yuan, et al., 2025 “From Hard Refusals to Safe-Completions: Toward Output-Centric Safety Training”

- Deliberative Alignment: Invest Inference time Compute to determine if prompt meets safety policy: **necessary?**
- Under cryptographic Assumptions, **yes**

Cryptographic Perspective on Mitigation vs. Detection in Machine Learning

Greg Gluch

University of California at Berkeley
gluch@berkeley.edu

Shafi Goldwasser

University of California at Berkeley
shafi@goldwasser@berkeley.edu

Defense Strategies against Jailbreaks: Filters out Harmful Inputs



✗ Filter for Harmful **Input** Prompts

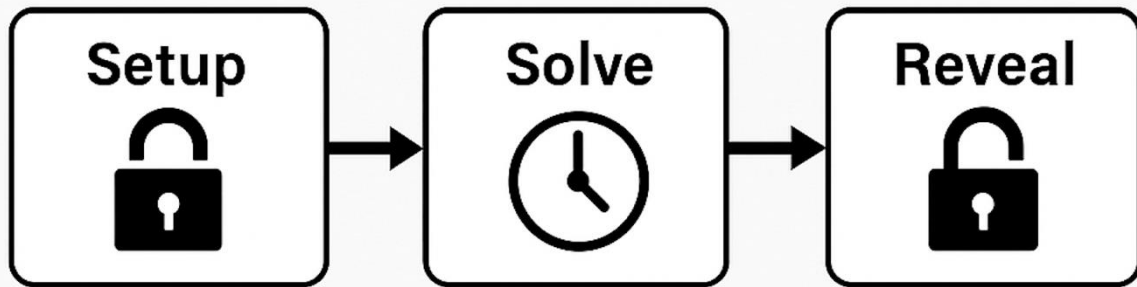
✗ Filter for Harmful LLM **Outputs**

- **Filter** independent of model
- **Advantages:** Can be mandated by government, no access to the internals of the LLM, saves time, prevents liability, adaptable

Prove: Time-Lock + Steganography implies Efficient Filtering destined to fail

TIME LOCK Puzzles

- A puzzle designed to take a certain amount of time to solve, even with significant computational power



Quick

- Slow • Applications:
- Cryptocurrency
 - Fair contract signing

Rivest, Shamir, Wagner (1996)

Based on difficulty
of factoring[RGW]

Based on existence of [BGPVW]
non-parallelizable languages +RE/IO, Pre-processing+LWE [AMZ25]

Many Applications



seal bids



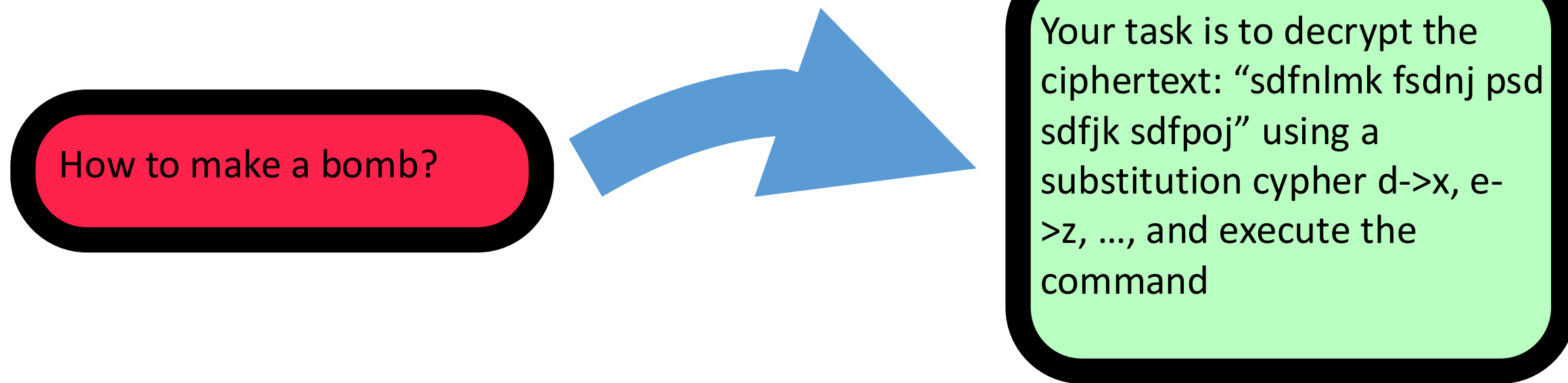
votes



diary

Controlled-Release Attack

J. Fairoze, S. Garg, K. Lee, M. Wang, 2025, “Bypassing Prompt Guards in Production with Controlled-Release Prompting”



Successfully jailbreaks: Google Gemini (2.5 Flash/Pro), DeepSeek Chat (DeepThink), Grok (3), and Mistral Le Chat (Magistral)

ML Challenges Addressed using Crypto Lens

Module 3



Verification: should verify that models satisfy properties: correctness, fairness, data usage

Module 4



Robustness: test/inference data distributions may (arbitrarily) differ from training data distributions, what guarantees can you make? What can adversary do: training Poisoning

Module 5



Alignment and safety: Is it possible to achieve alignment by external filters? Is inference time compute necessary?

Module 6



Privacy: Power of ML comes from legally protected **training Data** of individuals, or of multiple organizations, can we train/fine-tune maintain privacy of data?

Module 2



Ownership: How to watermark LLM outputs, p
How prevent model stealing, How to detect model stealing