

Self-Proving Models

Orr Paradise

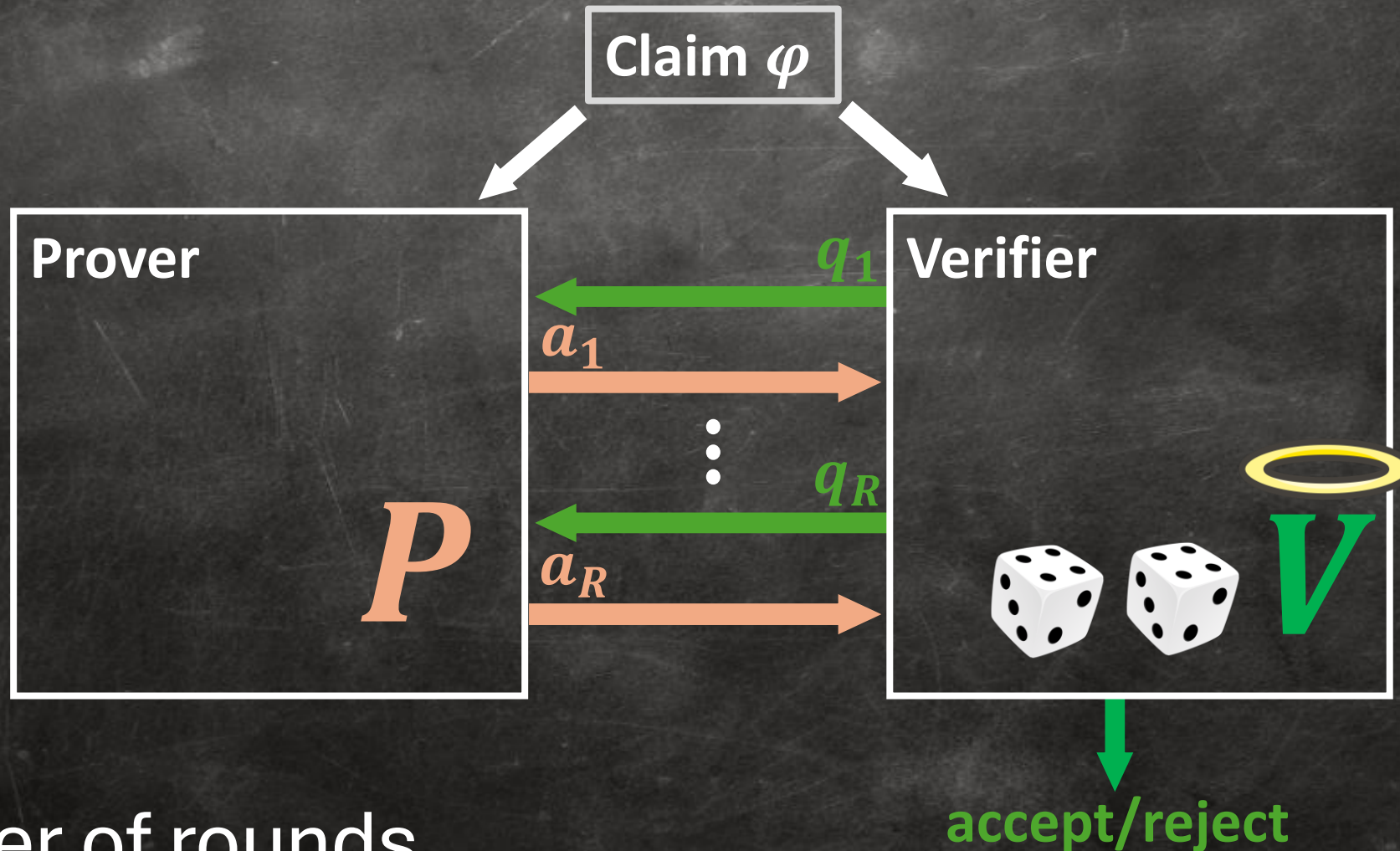
Lesson plan

- Reminder: Proofs and Models
- Defining Self-Proving models
- Learning Self-Proving models
 - Transcript Learning (TL)
 - Reinforcement Learning from Verifier Feedback (RLVF)
- Self-Proving models in practice

Lesson plan

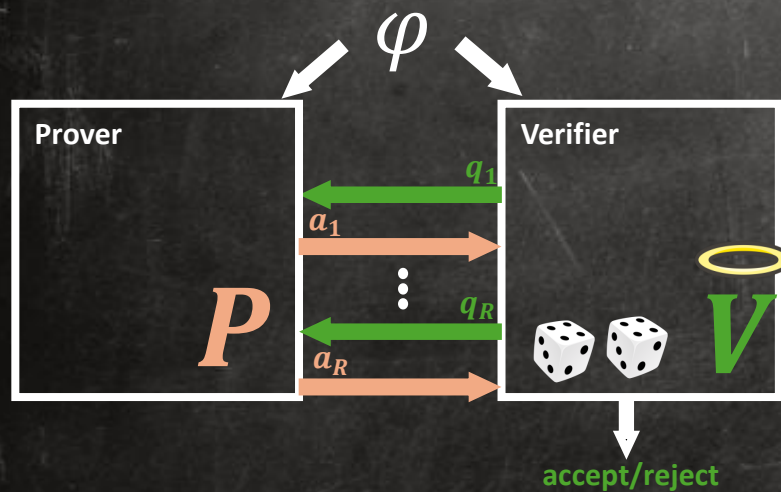
- Reminder: Proofs and Models
- Defining Self-Proving models
- Learning Self-Proving models
 - Transcript Learning (TL)
 - Reinforcement Learning from Verifier Feedback (RLVF)
- Self-Proving models in practice

Interactive Proof systems (IPs)



R : Number of rounds

Interactive Proof systems (IPs)



(1) Completeness

There exists an “honest prover” $\overset{\circ}{P}$ such that for any correct claim φ ,
 $\Pr[\overset{\circ}{P} \text{ convinces } V \text{ to accept } \varphi] = 100\%$

(2) Soundness (with δ error)

For any incorrect claim φ
and for any “lying prover” P ,
 $\Pr[P \text{ convinces } V \text{ to accept } \varphi] \leq \delta$

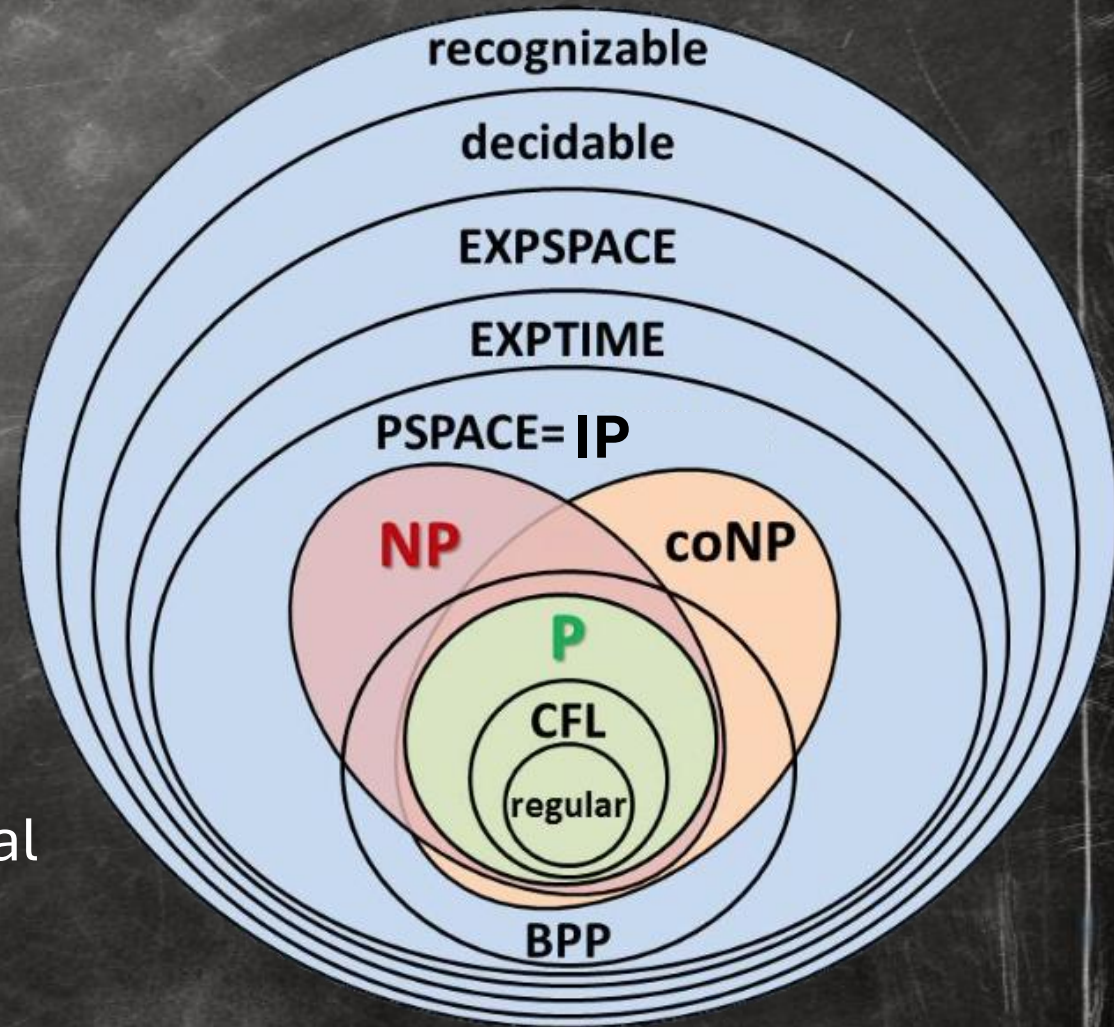
(3) Efficiency

Polytime Verifier, minimal #rounds, short messages...

→ “More efficient than the Prover”

The power of interaction

- We saw that Interactive Verification is much more powerful: $IP = PSPACE$ [LFKN, Shamir 90]
 - For any **polynomial-space** decidable claim, there exists an efficient (Probabilistic Polynomial Time) Interactive Verifier.
- *“Any problem you could dream of solving, can be efficiently Verified with the help of an untrusted Prover.”*
- IPs allow an efficient verifier to establish trust with a **more powerful, untrusted** computational entity...
- Who might this powerful entity be?



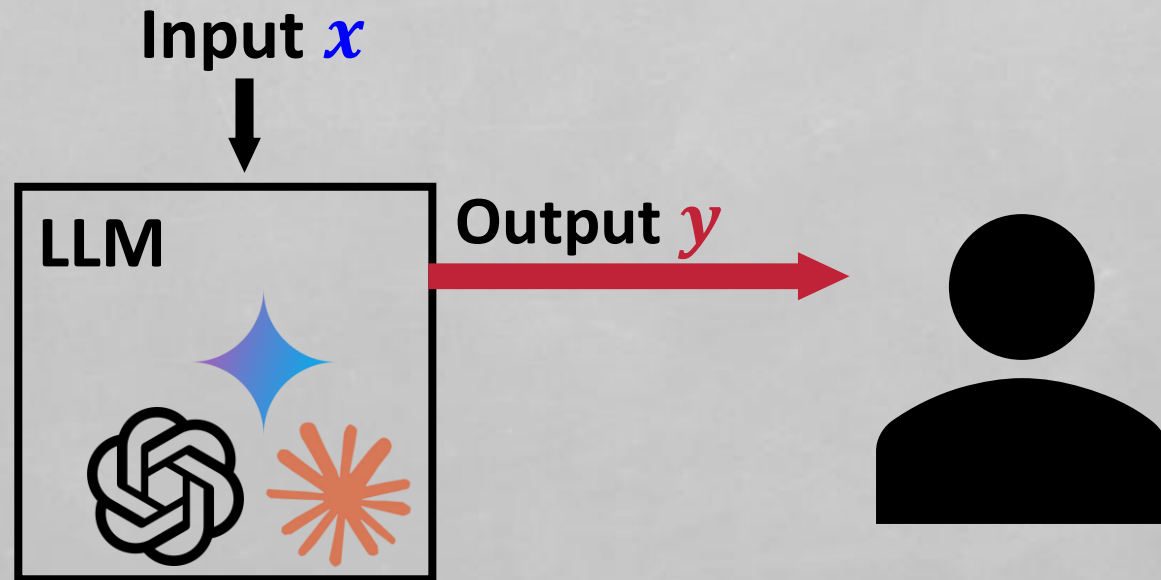
What is a model?

- Fix a ground-truth functionality $f^*: X \rightarrow Y$. For simplicity, let's take deterministic functions.
- f^* might be worst-case hard to compute.
 - We might not even *know* f^* is! Only know some samples $(x_1, f^*(x_1)), \dots, (x_m, f^*(x_m))$
- In practice, we only need to compute f^* on “natural” inputs.
 - “Nature” is an unknown distribution μ over X we can sample from.
- Fix a set of parameters Θ . A model family is a collection of (randomized) functions $\{f_\theta: X \rightarrow Y\}_{\theta \in \Theta}$
- Model training: Given input $S = (x_1, f^*(x_1)), \dots, (x_m, f^*(x_m))$ for $x_1, \dots, x_m \sim \mu$, output $\hat{\theta} \in \Theta$
- A model $\hat{\theta}$ is $(1 - \varepsilon)$ -correct (w.r.t μ, f^*) if
$$\Pr_{\substack{x \sim \mu \\ \hat{y} \sim f_{\hat{\theta}}(x)}} [\hat{y} = f^*(x)] \geq 1 - \varepsilon$$

In what sense do we “trust” a model?

- Average-case: Test if model f_{θ} has 99% accuracy w.r.t μ, f^* (last week’s lectures)
- Worst-case: You have an input x , for which the model generates $y \sim f_{\theta}(x)$

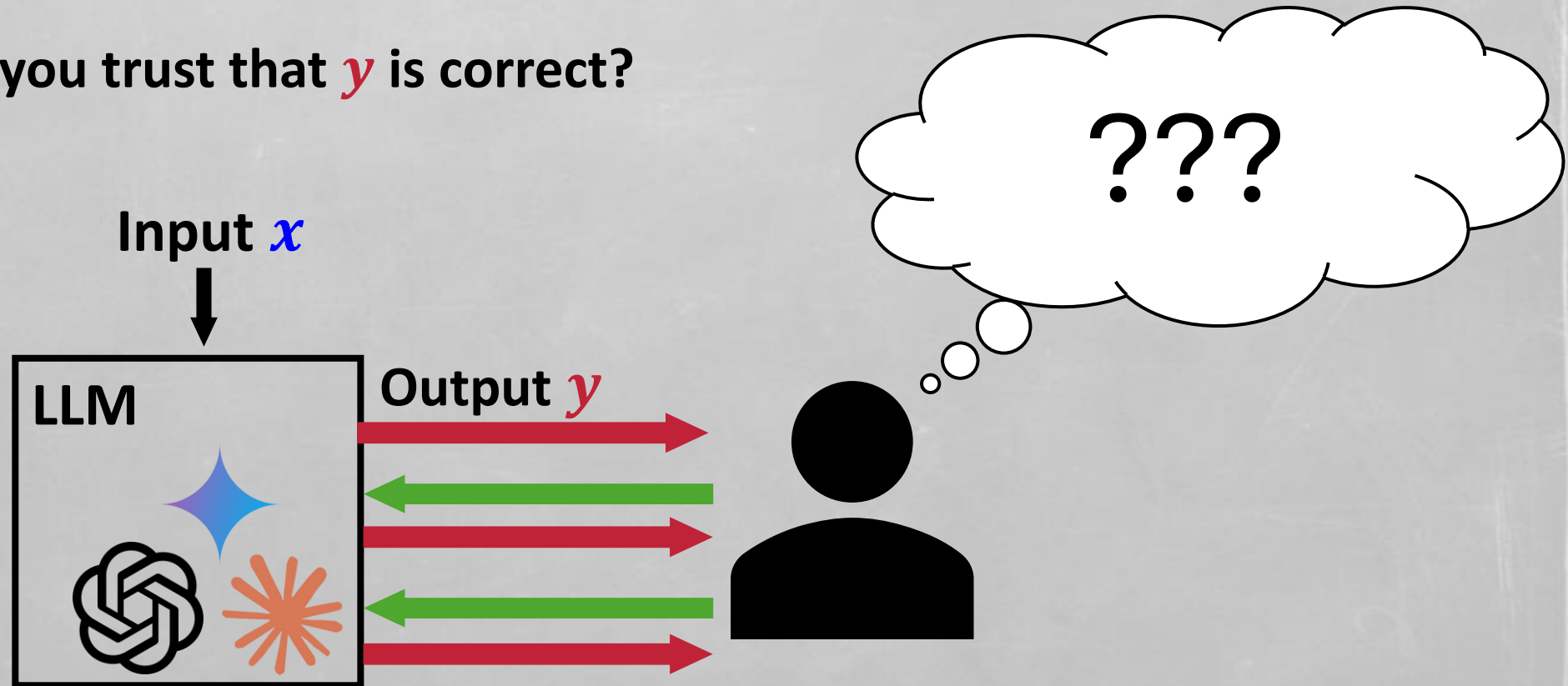
Should you trust that y is correct?



In what sense do we “trust” a model?

- Average-case: Test if model f_{θ} has 99% accuracy w.r.t μ, f^* (last week’s lectures)
- Worst-case: You have an input x , for which the model generates $y \sim f_{\theta}(x)$

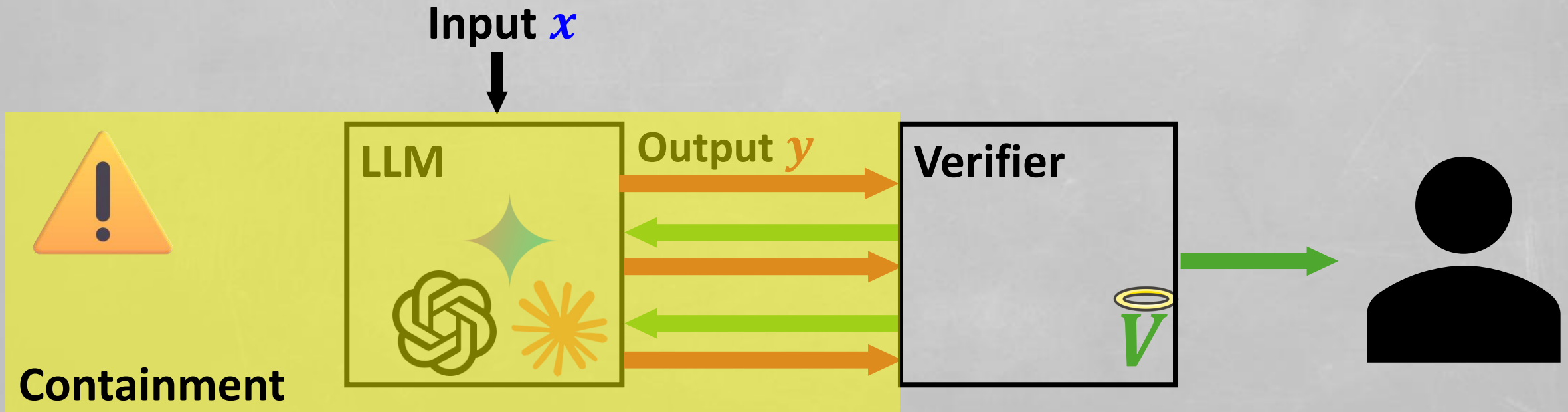
Should you trust that y is correct?



In what sense do we “trust” a model?

- Average-case: Test if model f_{θ} has 99% accuracy w.r.t μ, f^* (last week’s lectures)
- Worst-case: You have an input x , for which the model generates $y \sim f_{\theta}(x)$

Should you trust that y is correct?

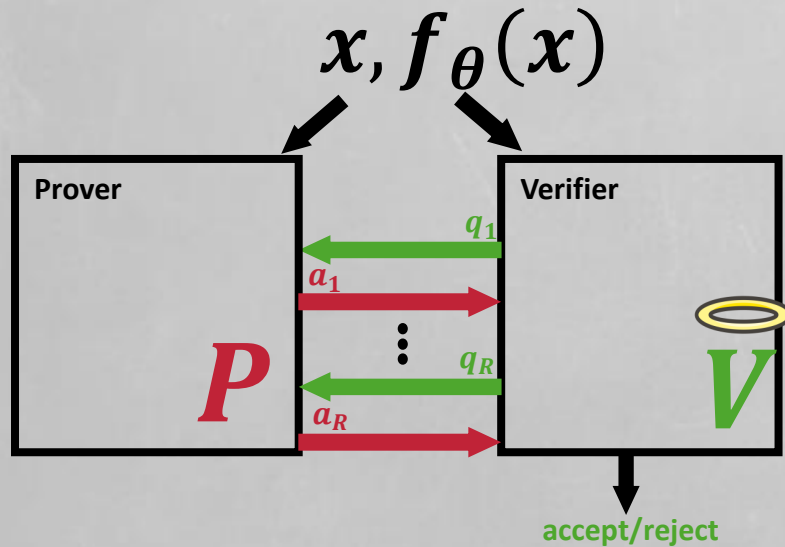


Lesson plan

- Reminder: Proofs and Models
- Defining Self-Proving models
- Learning algorithms for Self-Proving models
 - Transcript Learning (TL)
 - Reinforcement Learning from Verifier Feedback (RLVF)
- Self-Proving models in practice

IPs for model outputs

Fix a model $f_\theta: X \rightarrow Y$ (deterministic, for simplicity)
The “claim” is that $f_\theta(x) = f^*(x)$.



(1) Completeness

There exists an “honest prover” $\overset{\text{halo}}{P}$ such that for any input x ,

$$\Pr[\overset{\text{halo}}{P} \text{ convinces } V \text{ to accept } (x, f^*(x))] = 100\%$$

(2) Soundness

For any input x and incorrect output $f_\theta(x) \neq f^*(x)$, for any “lying prover” P

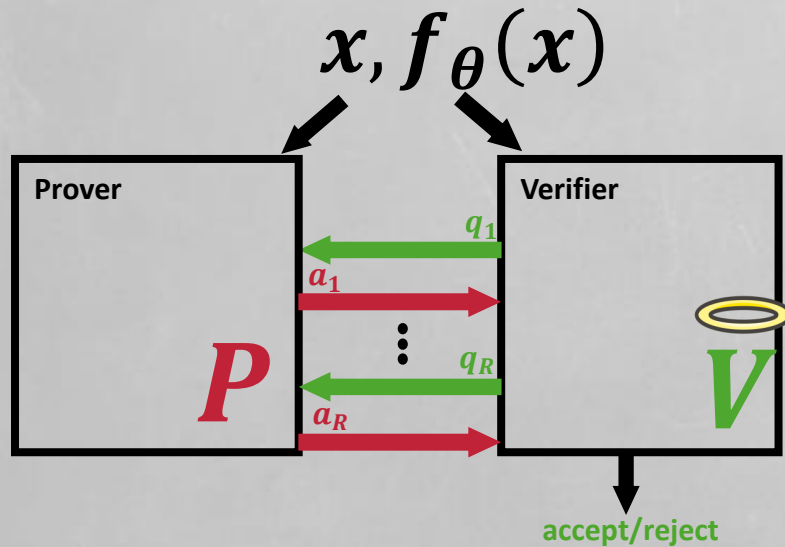
$$\Pr[P \text{ convinces } V \text{ to accept } (x, f_\theta(x))] \leq \delta$$

(3) Efficiency

Polytime verifier, minimal #rounds, short messages...

IPs for model outputs

Fix a model $f_\theta: X \rightarrow Y$ (deterministic, for simplicity)
The “claim” is that $f_\theta(x) = f^*(x)$.



Even if we learned a model f_θ ,
where do we get this Prover?

(1) Completeness

There exists an “honest prover” $\overset{\circ}{P}$

such that for any input x ,

$$\Pr[\overset{\circ}{P} \text{ convinces } V \text{ to accept } (x, f^*(x))] = 100\%$$

(2) Soundness

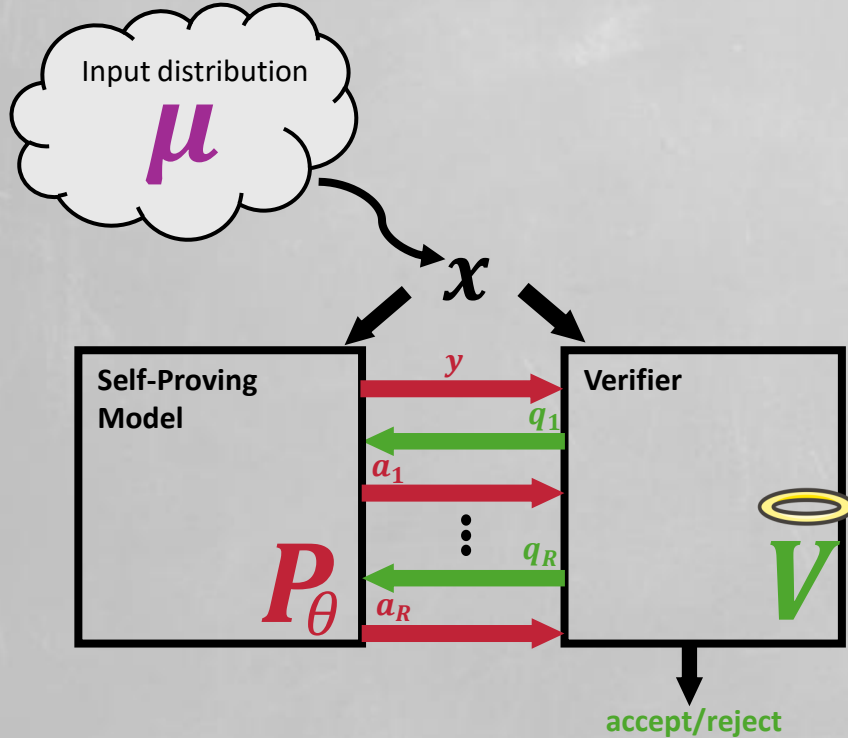
For any input x and incorrect output $f_\theta(x) \neq f^*(x)$,
for any “lying prover” P

$$\Pr[P \text{ convinces } V \text{ to accept } (x, f_\theta(x))] \leq \delta$$

Why not learn them *jointly*?

Let the model *prove its own output!*

Self-Proving models



Definition (Self-Proving Model with error ε)

A model P_θ is **Self-Proving** w.r.t V, μ

If $\Pr[P_\theta \text{ convinces } V \text{ to accept } y] \geq 1 - \varepsilon$

Over $x \sim \mu, y \sim P_\theta(x), q_1, a_1, \dots, q_R, a_R$

Exercise (Soundness+Self-Prov. \rightarrow Correctness):

Fix a Verifier V for f^* with soundness error δ .

If P is Self-Proving with err ε w.r.t V, μ , then

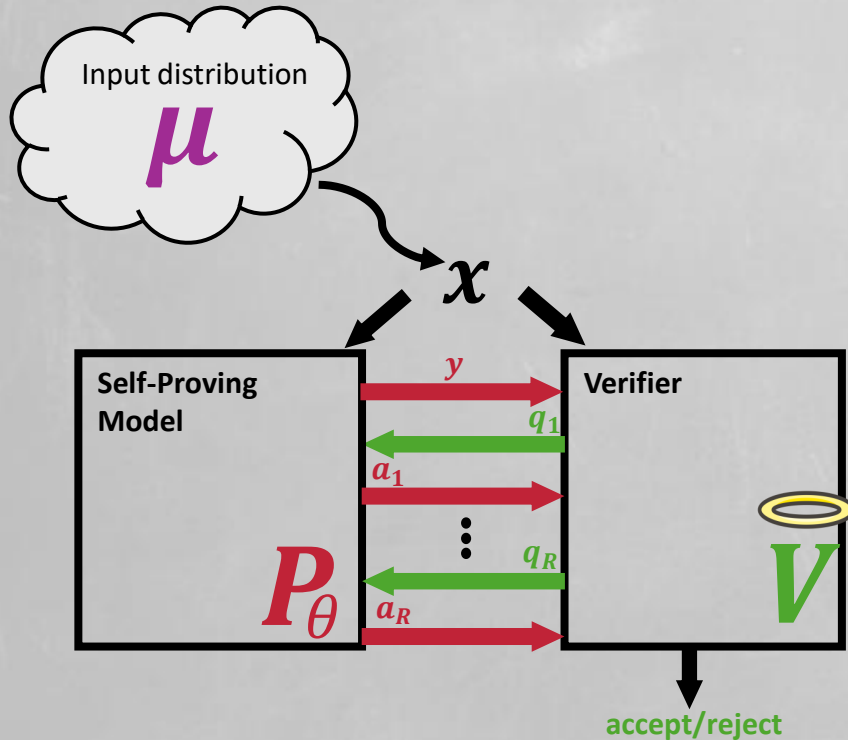
$$\Pr_{x \sim \mu} [P_\theta(x) = f^*(x)] \geq 1 - \delta - \varepsilon$$

Goal: Learn θ such that P_θ is Self-Proving

Lesson plan

- What is a proof? What is a model?
- Defining Self-Proving models
- Learning Self-Proving models
 - Transcript Learning (TL)
 - Reinforcement Learning from Verifier Feedback (RLVF)
- Self-Proving models in practice
- Proving convergence guarantees for TL and RLVF

Reminder: Self-Proving models



Definition (Self-Proving Model with error ε)

A model P_θ is **Self-Proving** w.r.t V, μ

If $\Pr[P_\theta \text{ convinces } V \text{ to accept } y] \geq 1 - \varepsilon$

Over $x \sim \mu, y \sim P_\theta(x), q_1, a_1, \dots, q_R, a_R$

Exercise (Soundness+Self-Prov. \rightarrow Correctness):

Fix a Verifier V for f^* with soundness error δ .

If P is Self-Proving with err ε w.r.t V, μ , then

$$\Pr_{x \sim \mu} [P_\theta(x) = f^*(x)] \geq 1 - \delta - \varepsilon$$

Goal: Learn θ such that P_θ is Self-Proving

Learning Self-Proving models

What do we need in order to learn a Self-Proving model?

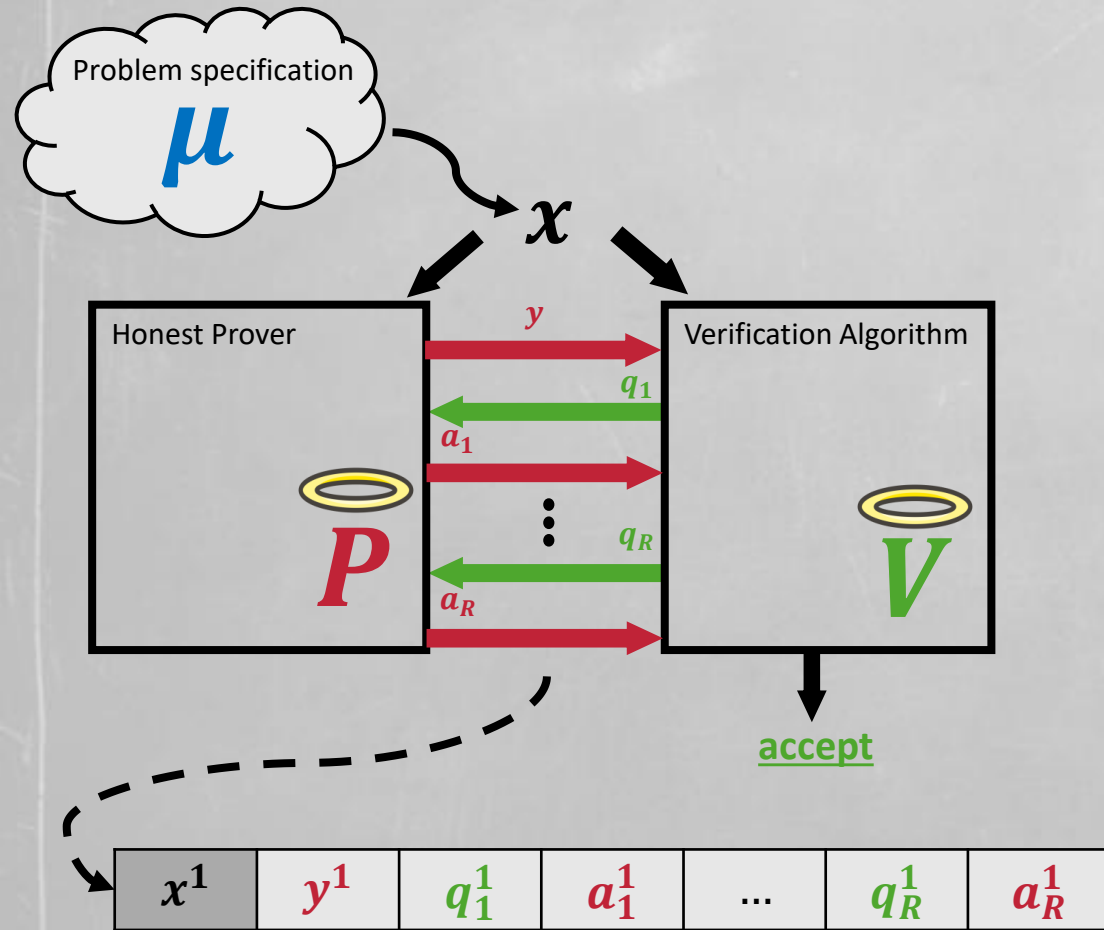
- Sequence-to-sequence model family: $\{P_\theta: \Sigma^* \rightarrow \Sigma^*\}_{\theta \in \Theta}$ where $\Theta = \mathbb{R}^d$.
- Autoregressive generation: For any parameter setting θ and input string $z \in \Sigma^*$,
 - Next token (“forwards pass”): RV $p_\theta(z)$ over Σ , determined by the *logits*: $\log p_\theta(z) \in \mathbb{R}^{|\Sigma|}$
 - Keep generating until a special token $\text{EOS} \in \Sigma$ is sampled.
- Differentiability (“backwards pass”): Can compute $\nabla_\theta \log \Pr_{\tau \sim p_\theta(z)} [\tau = \sigma]$ for each $\sigma \in \Sigma$.
- Access to the input distribution $x \sim \mu$
- Access to the Verifier V .
- **Transcript Learning**: assume also access to an “honest prover”
 - (Slightly weaker access to an “honest transcript generator” suffices)

Transcript Learning

Step 1: Transcribe "honest" interactions

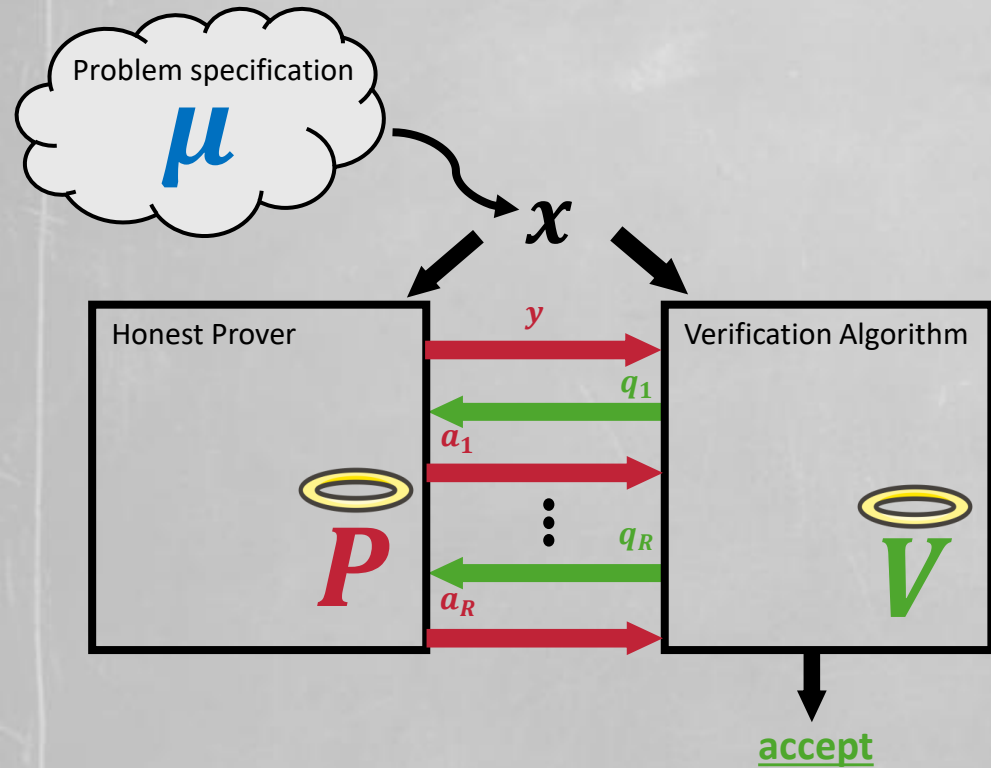
Transcript Learning

Step 1: Transcribe "honest" interactions



Transcript Learning

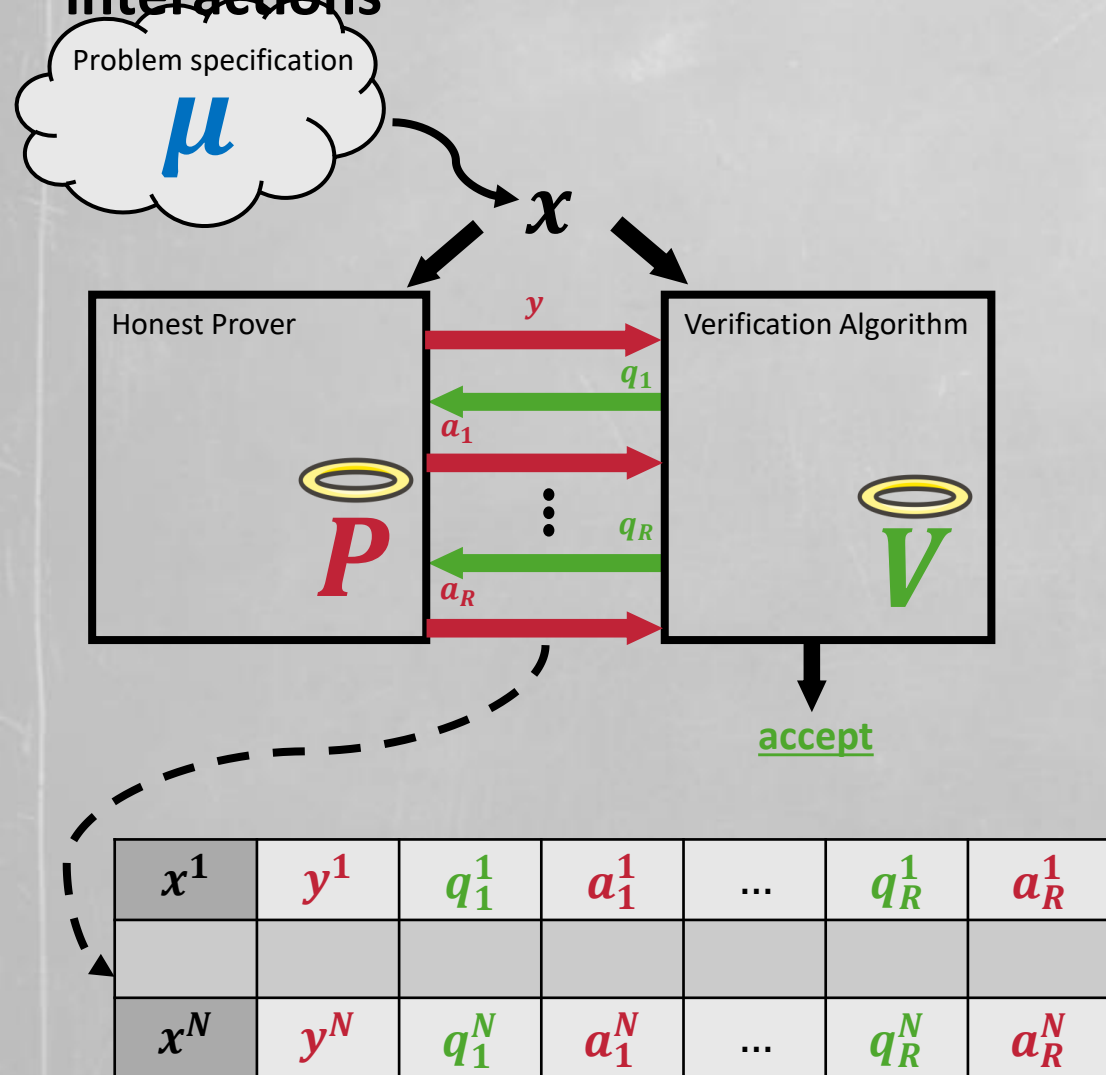
Step 1: Transcribe "honest" interactions



x^1	y^1	q_1^1	a_1^1	...	q_R^1	a_R^1
-------	-------	---------	---------	-----	---------	---------

Transcript Learning

Step 1: Transcribe "honest" interactions



Transcript Learning

Step 1: Transcribe "honest" interactions

Step 2: Transcript Cloning

x^1	y^1	q_1^1	a_1^1	...	q_R^1	a_R^1
x^N	y^N	q_1^N	a_1^N	...	q_R^N	a_R^N

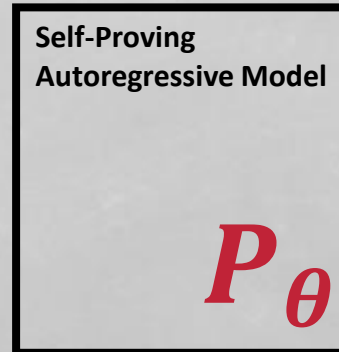
Transcript Learning

Step 1: Transcribe "honest" interactions

Step 2: Transcript Cloning

For each transcript:

x	y	q_1	a_1	...	q_R	a_R
-----	-----	-------	-------	-----	-------	-------



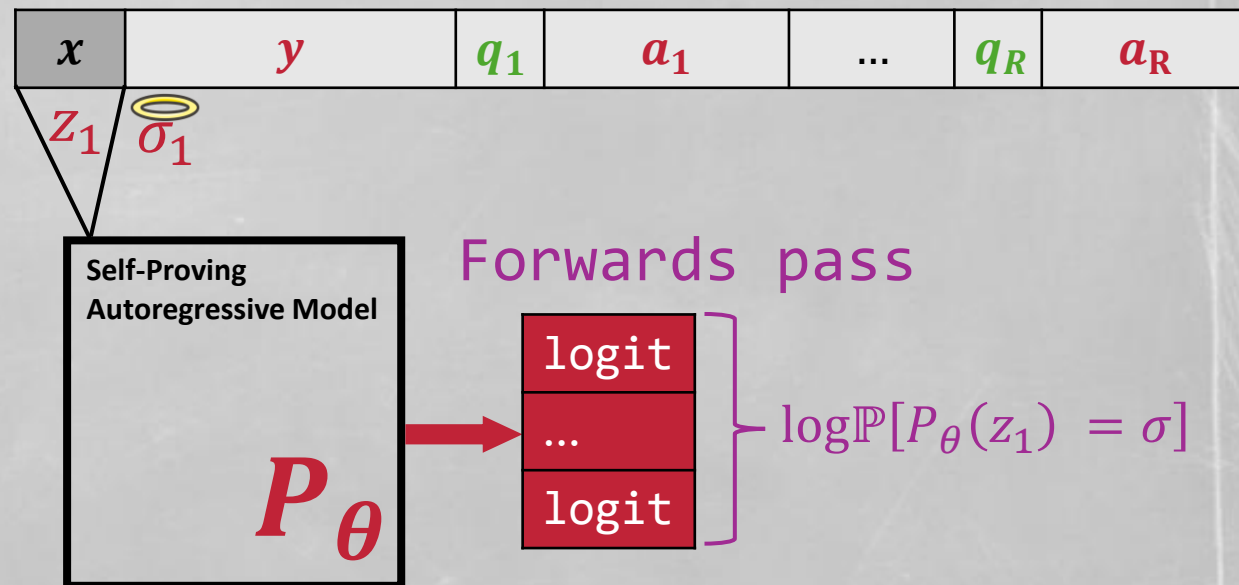
x^1	y^1	q_1^1	a_1^1	...	q_R^1	a_R^1
x^N	y^N	q_1^N	a_1^N	...	q_R^N	a_R^N

Transcript Learning

Step 1: Transcribe "honest" interactions

Step 2: Transcript Cloning

For each transcript:



Backwards pass $\vec{d}_1 := \nabla_\theta \log \mathbb{P}[P_\theta(z_1) = \sigma_1]$

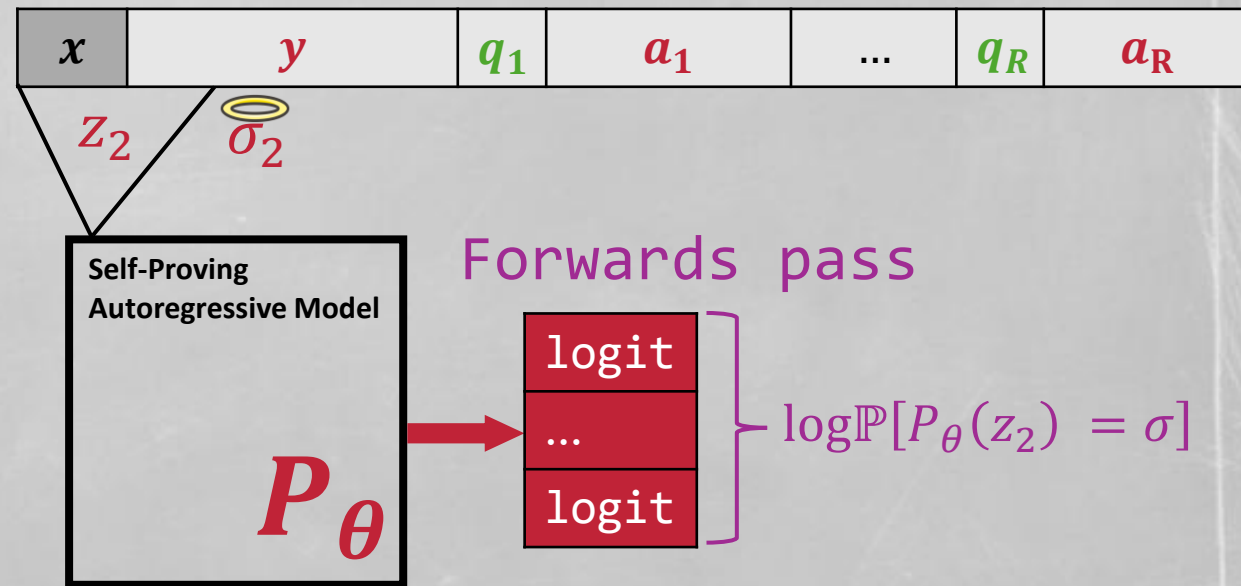
x^1	y^1	q_1^1	a_1^1	...	q_R^1	a_R^1
x^N	y^N	q_1^N	a_1^N	...	q_R^N	a_R^N

Transcript Learning

Step 1: Transcribe "honest" interactions

Step 2: Transcript Cloning

For each transcript:



Backwards pass $\vec{d}_2 := \nabla_\theta \log \mathbb{P}[P_\theta(z_2) = \sigma_2]$

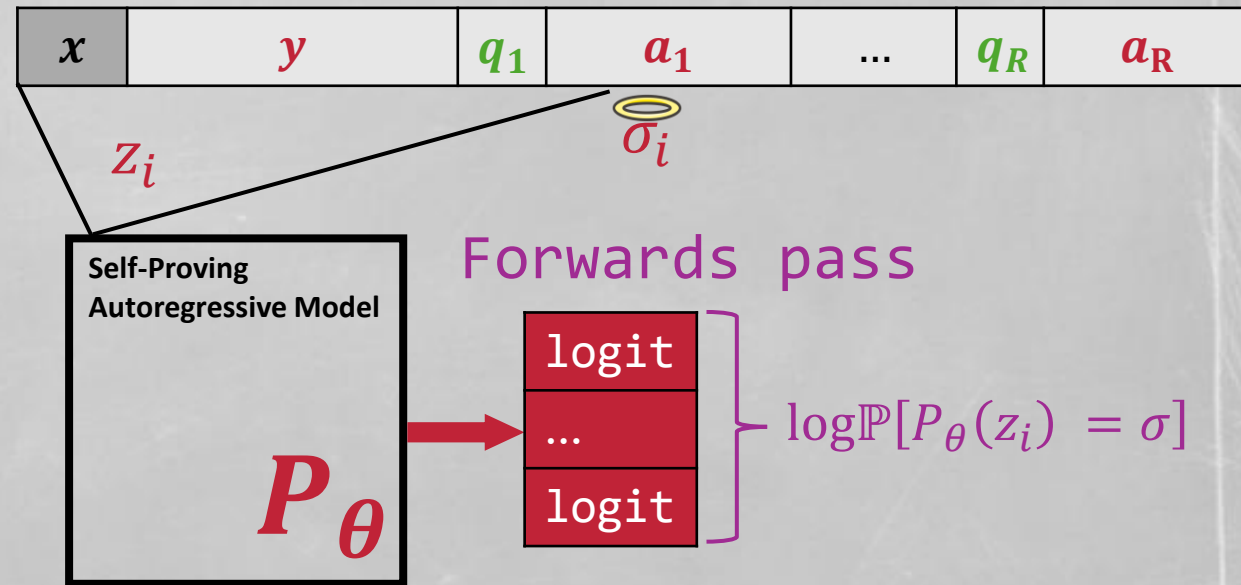
x^1	y^1	q_1^1	a_1^1	...	q_R^1	a_R^1
x^N	y^N	q_1^N	a_1^N	...	q_R^N	a_R^N

Transcript Learning

Step 1: Transcribe "honest" interactions

Step 2: Transcript Cloning

For each transcript:



Backwards pass $\vec{d}_i := \nabla_\theta \log \mathbb{P}[P_\theta(z_i) = \sigma]$

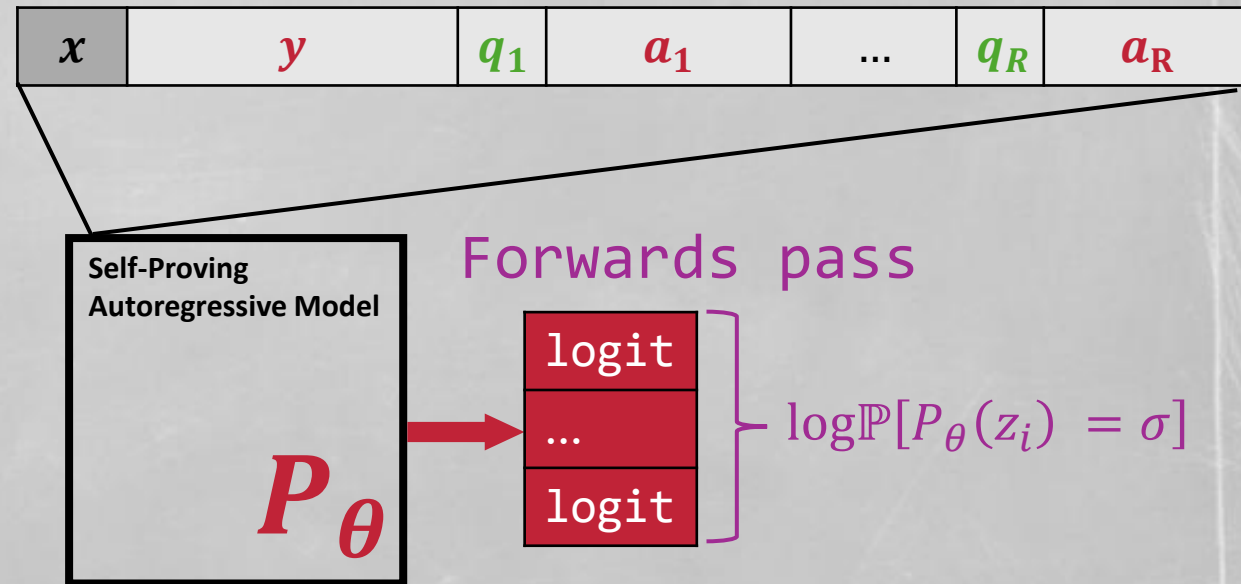
x^1	y^1	q_1^1	a_1^1	...	q_R^1	a_R^1
x^N	y^N	q_1^N	a_1^N	...	q_R^N	a_R^N

Transcript Learning

Step 1: Transcribe "honest" interactions

Step 2: Transcript Cloning

For each transcript:



Backwards pass $\vec{d}_i := \nabla_{\theta} \log \mathbb{P}[P_\theta(z_i) = \sigma_i]$

Update params $\theta \leftarrow \theta + \lambda \cdot \mathbb{P}[\overrightarrow{P_\theta}(x) = z] \cdot \sum_i \vec{d}_i$

Forwards Backwards

x^1	y^1	q_1^1	a_1^1	...	q_R^1	a_R^1
x^N	y^N	q_1^N	a_1^N	...	q_R^N	a_R^N

Transcript Learning

Assumptions:

- Access to a dataset of honest transcripts.
- The total number of tokens sent by the prover in any interaction is $< C$.
- The surrogate objective $A(\theta) := \mathbb{P}_x[\pi_\theta(x) = \pi(x) \mid V\text{'s rand}]$ is concave and differentiable in θ .
- The logits of P_θ are B_1 -Lipschitz in θ .
- For $\epsilon > 0$ let B_2 such that:
 - There exists θ^* with $\|\theta^*\| < B_2$ such that $A(\theta^*) \geq 1 - \epsilon/2$.

“Necessary” for reasoning about NN convergence

The TL Theorem: Under the assumptions,

Transcript Learning outputs a $(1 - \epsilon)$ -**Self-Proving model** when trained on

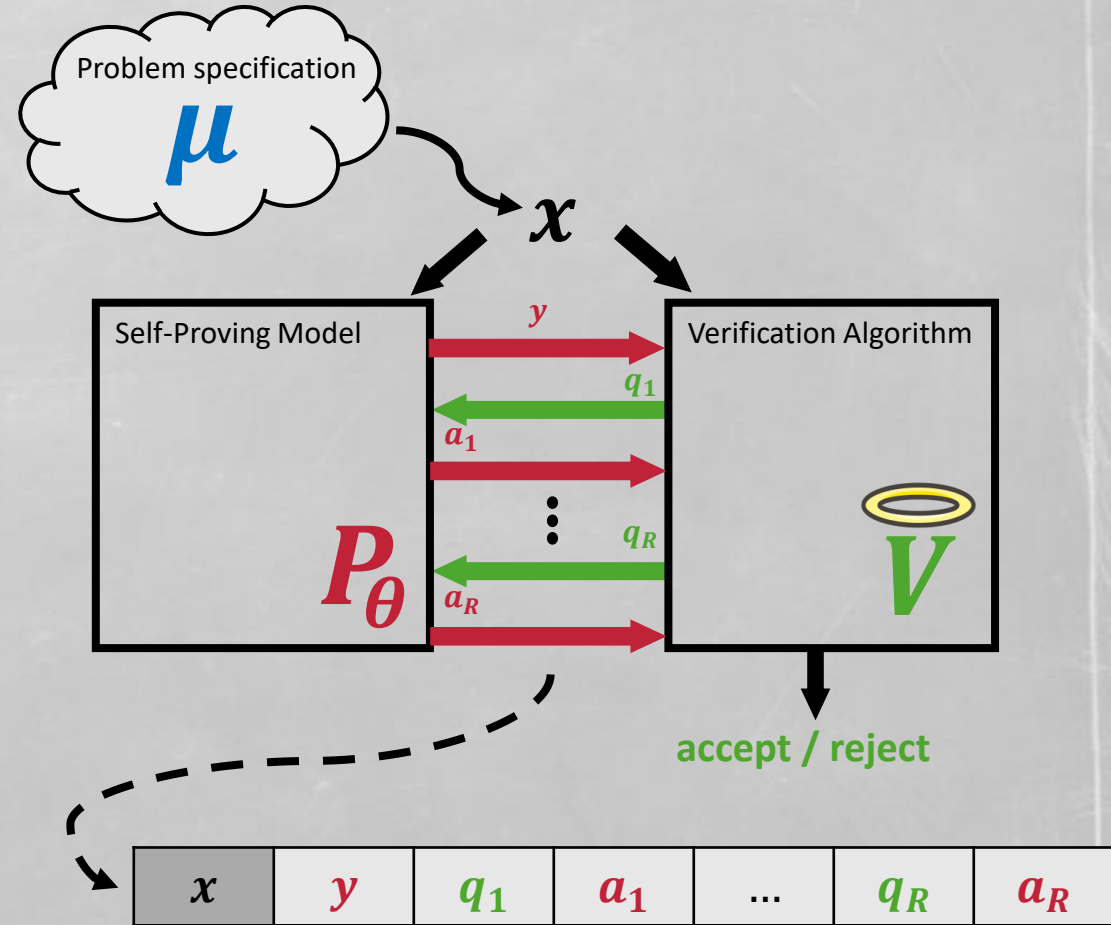
$$N \geq 4 \left(C \cdot B_1 \cdot B_2 \cdot \frac{1}{\epsilon} \right)^2 \text{ samples}$$

RL from Verifier Feedback (RLVF)

RL from Verifier Feedback (RLVF)

Repeat the following:

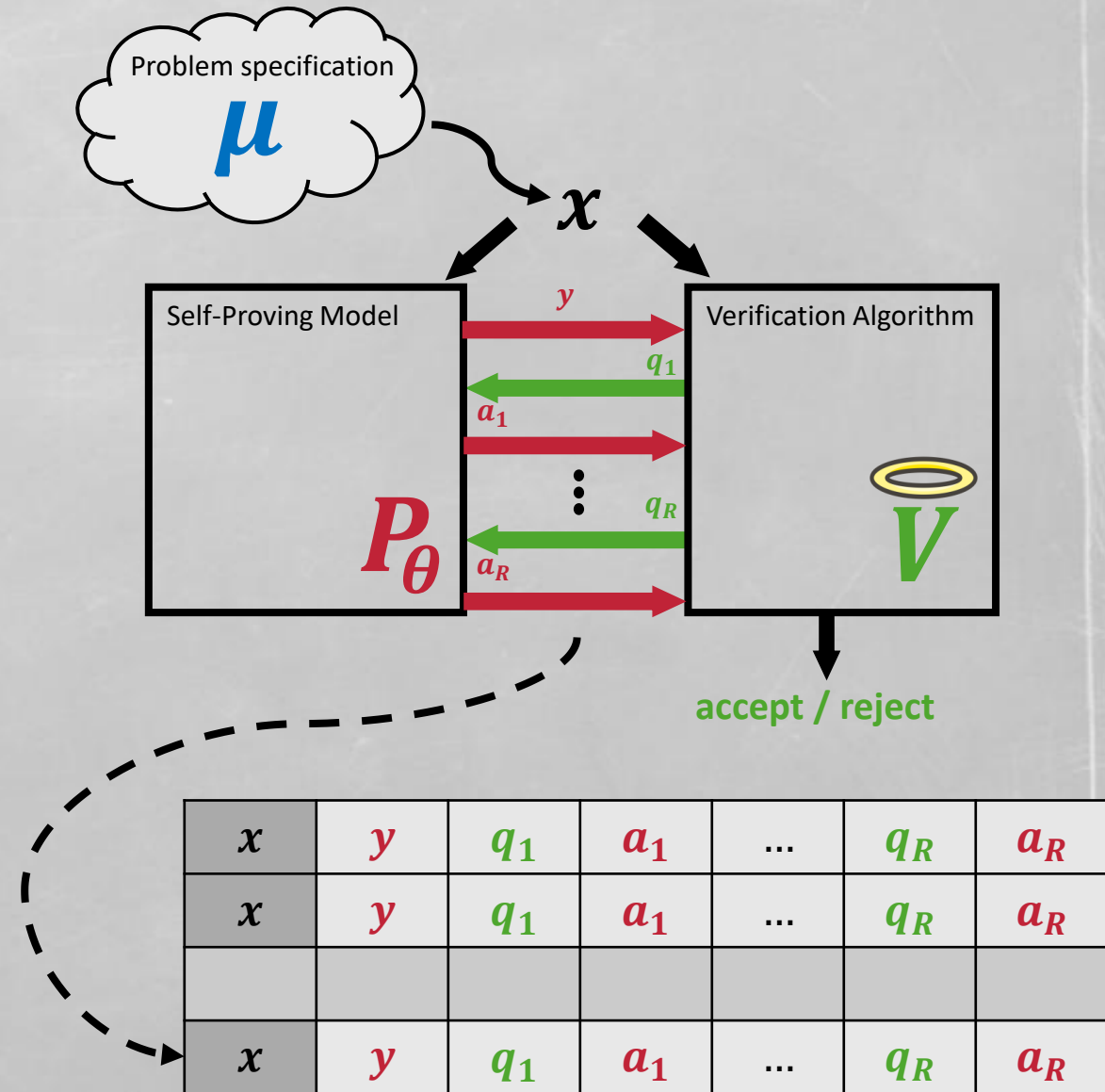
1. Generate transcript batch with P_θ .
Keep only transcripts accepted by the Verifier.



RL from Verifier Feedback (RLVF)

Repeat the following:

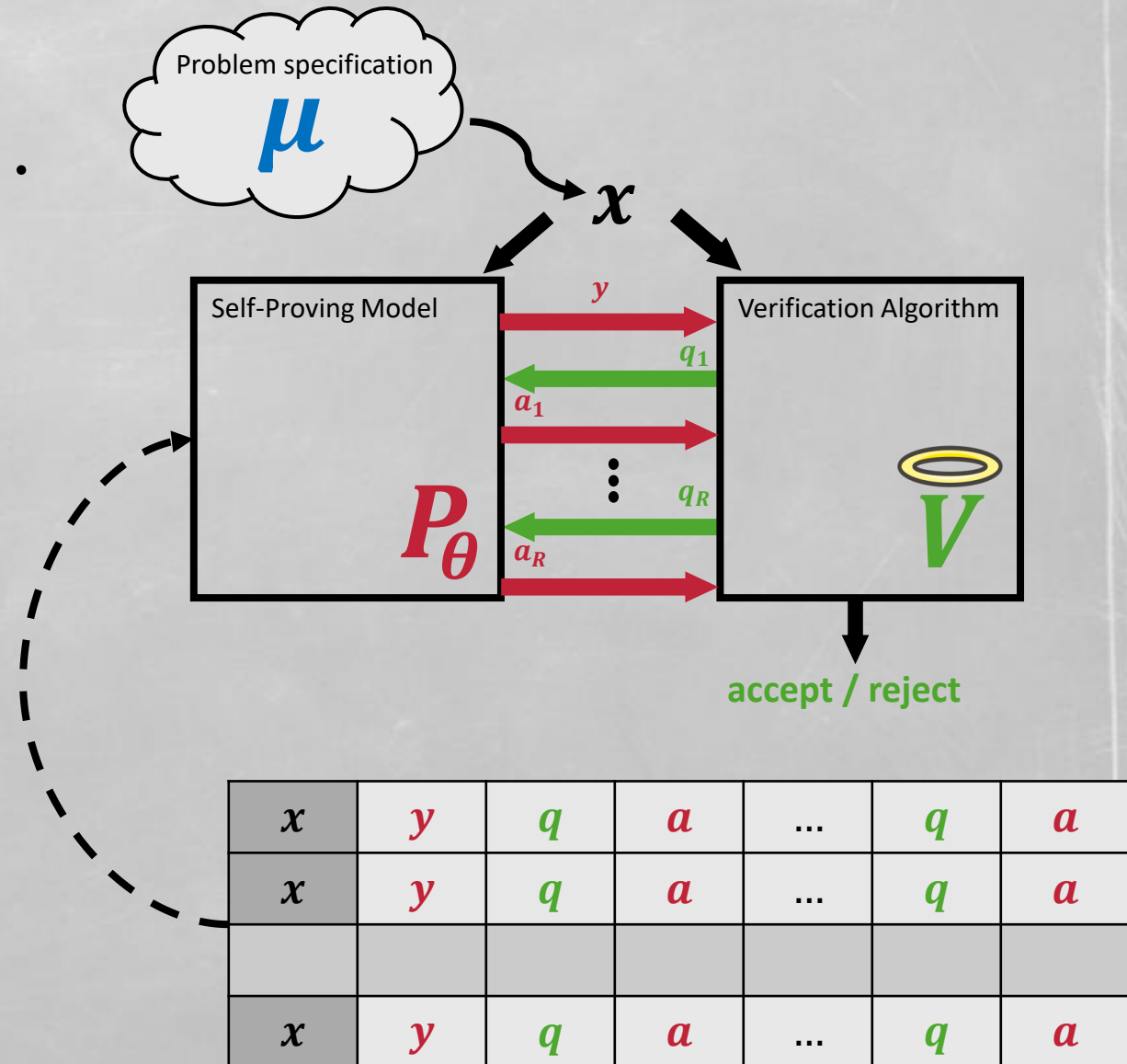
1. Generate transcript batch with P_θ .
Keep only transcripts accepted by the Verifier.



RL from Verifier Feedback (RLVF)

Repeat the following:

1. Generate transcript batch with P_θ .
Keep only transcripts accepted by the Verifier.
2. Update θ towards accepted trans.



Lesson plan

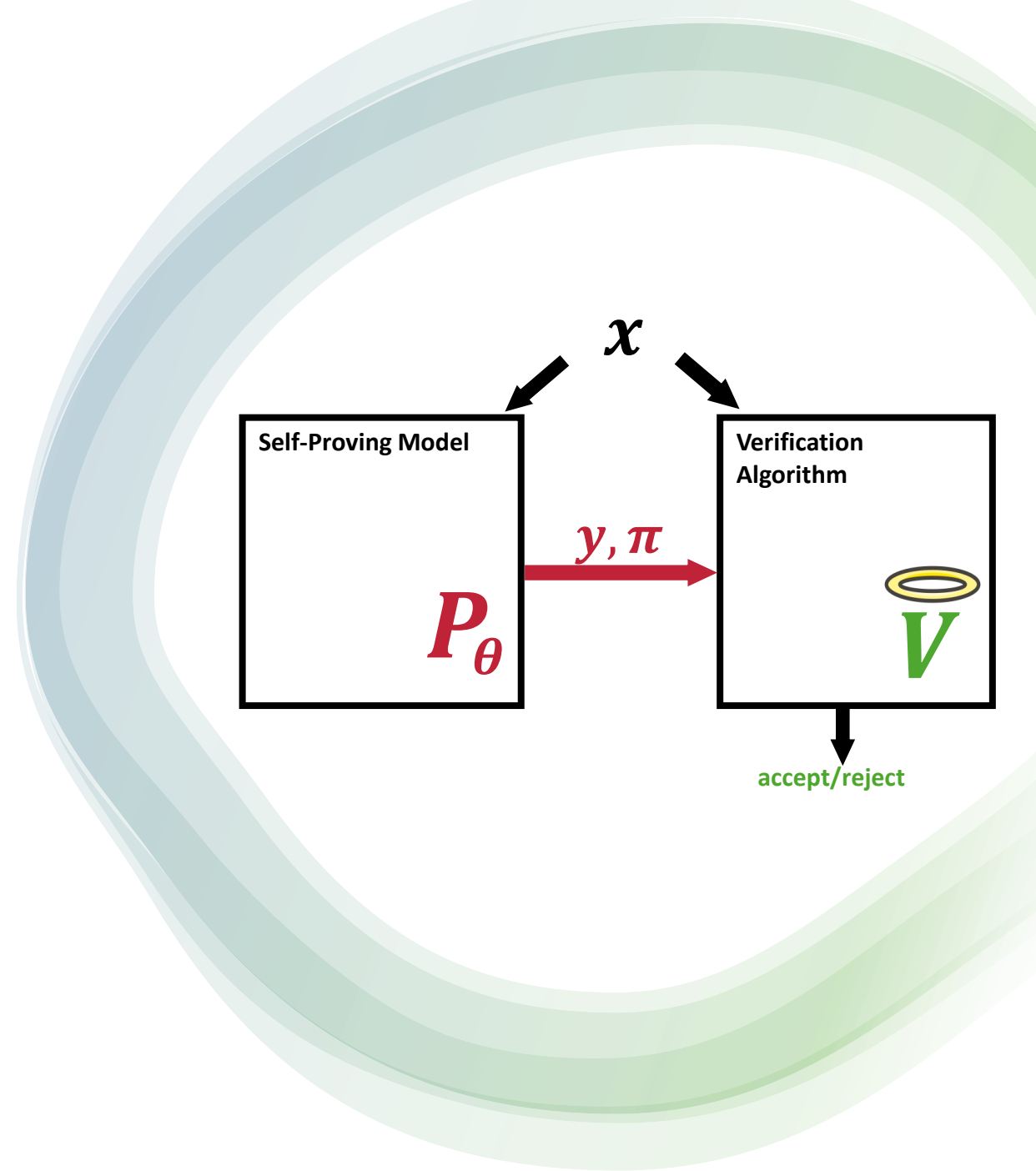
- Reminder: Proofs and Models
- Defining Self-Proving models
- Learning algorithms for Self-Proving models:
 - Transcript Learning (TL)
 - Reinforcement Learning from Verifier Feedback (RLVF)
- Self-Proving models in practice
- Proving convergence guarantees for TL and RLVF

Self-Proving models in practice

- “Learning to Prove” is a natural thing to do.
- It has been explored in many settings, mostly in practice.
- Let’s look at some of them:
 - In the ML literature
 - Experiments directly on TL and RLVF
 - RLVF adoption (sort of)

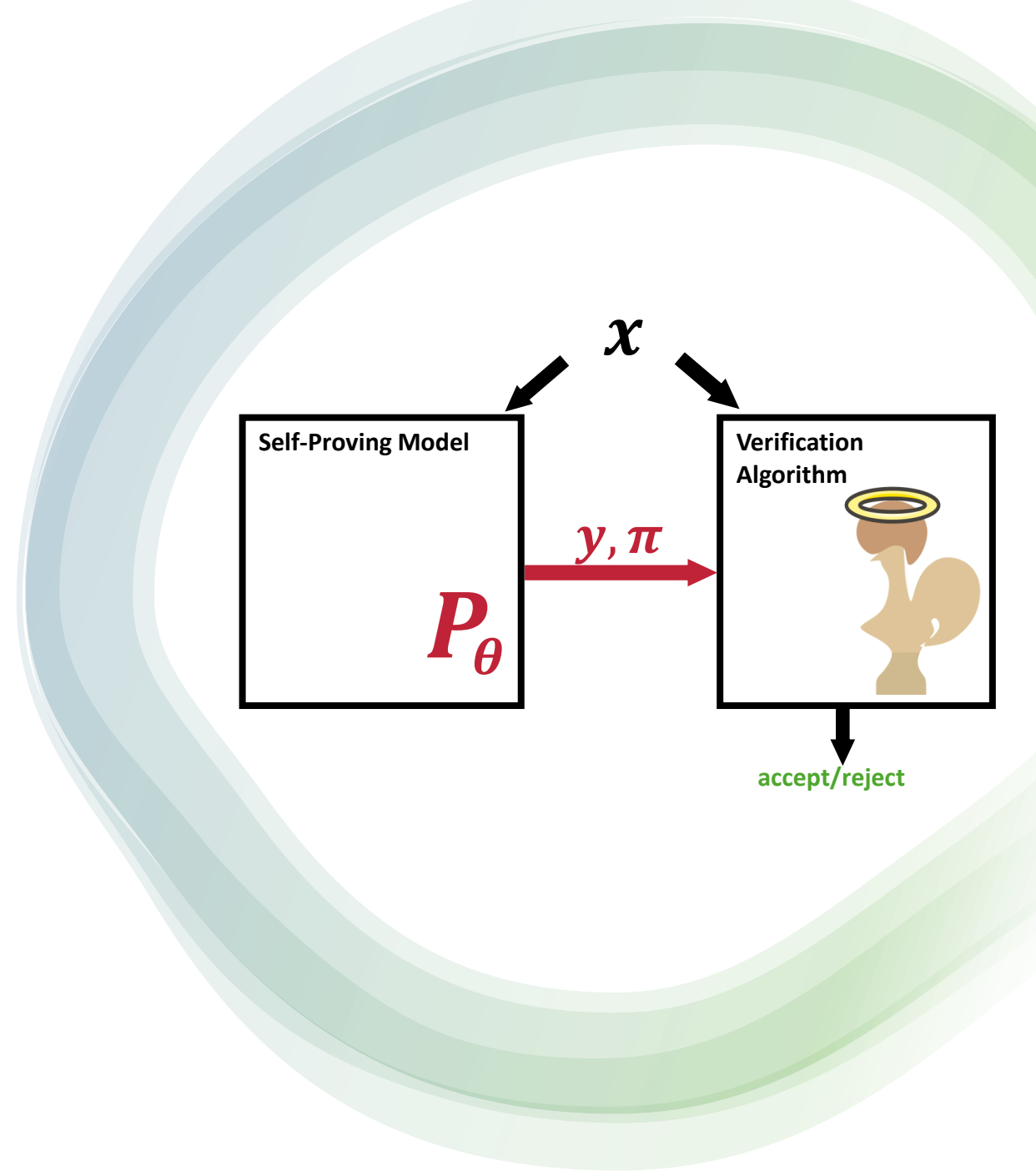
IP ∩ ML [May 2024]

- Learning to Prove...



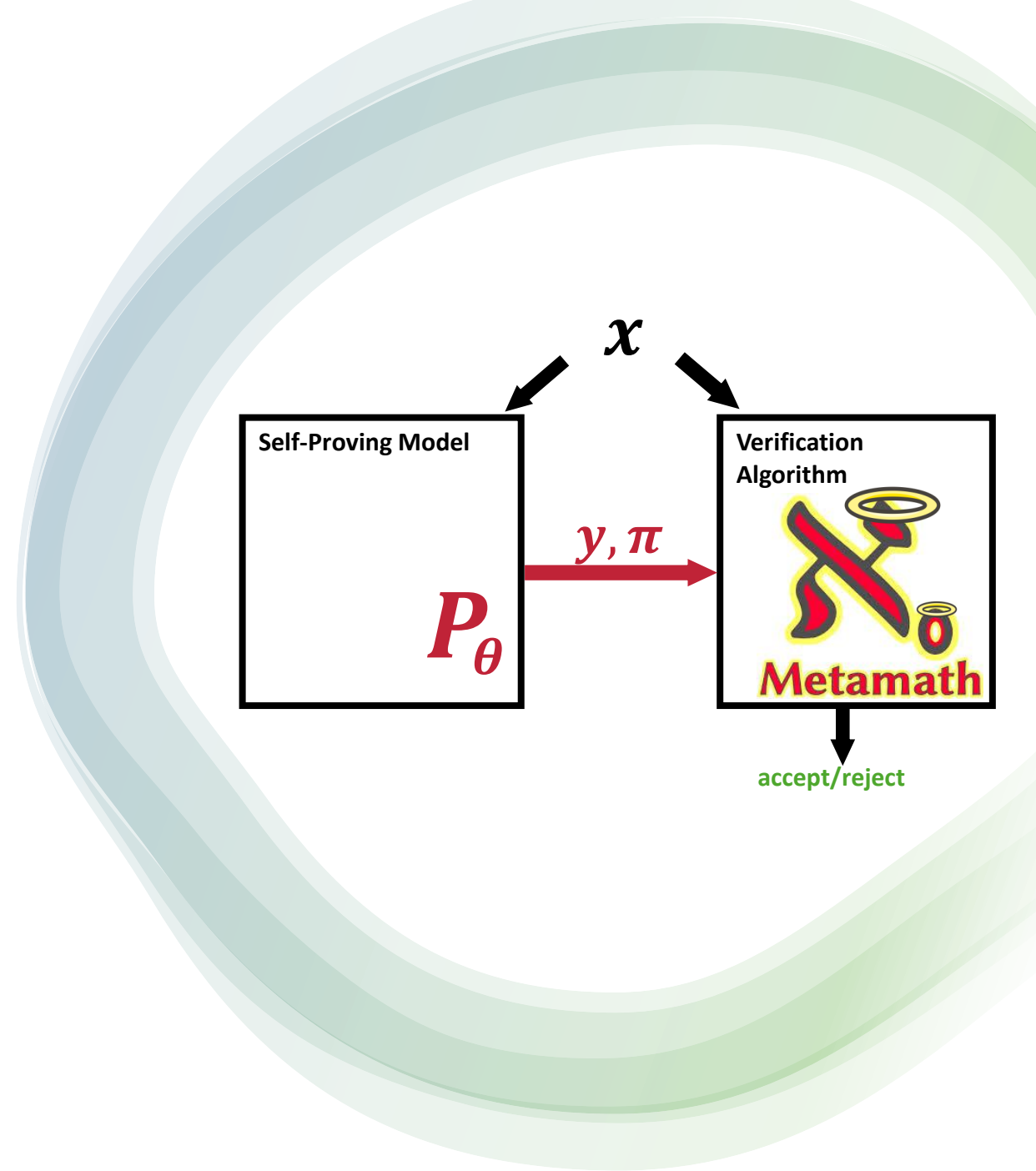
IP ∩ ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)



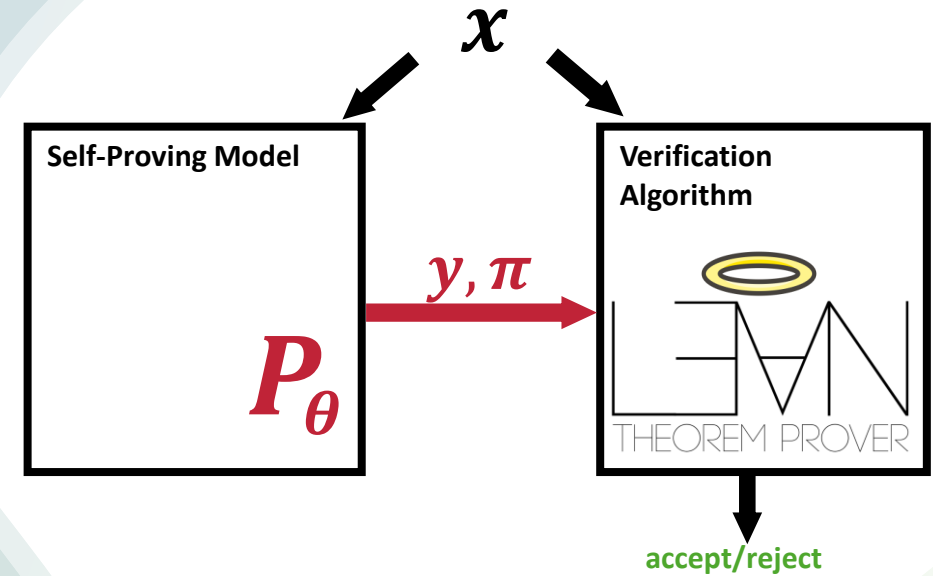
IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)



IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)



IP \cap ML [May 2024]

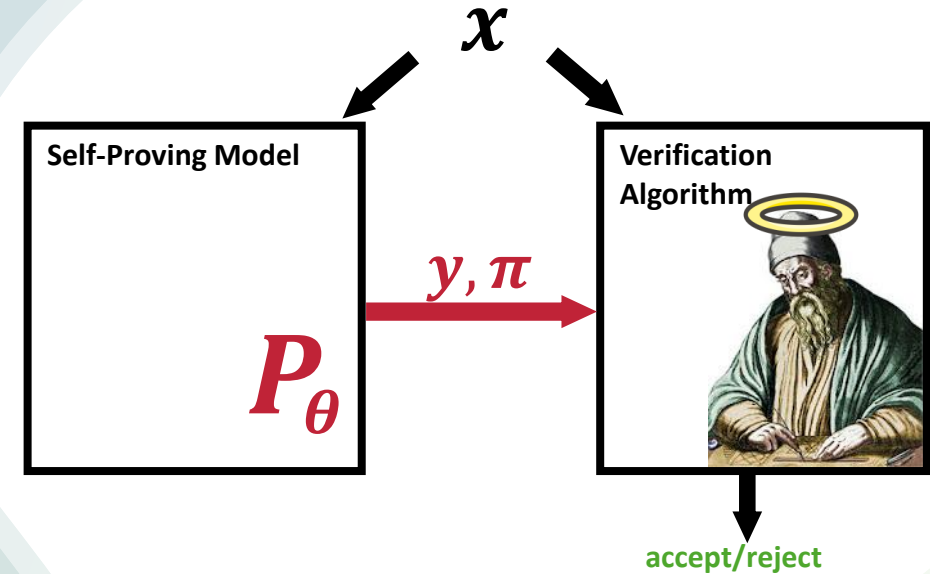
- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)

Article | [Open access](#) | Published: 17 January 2024

Solving olympiad geometry without human demonstrations

[Trieu H. Trinh](#) , [Yuhuai Wu](#), [Quoc V. Le](#), [He He](#) & [Thang Luong](#) 

[Nature](#) **625**, 476–482 (2024) | [Cite this article](#)



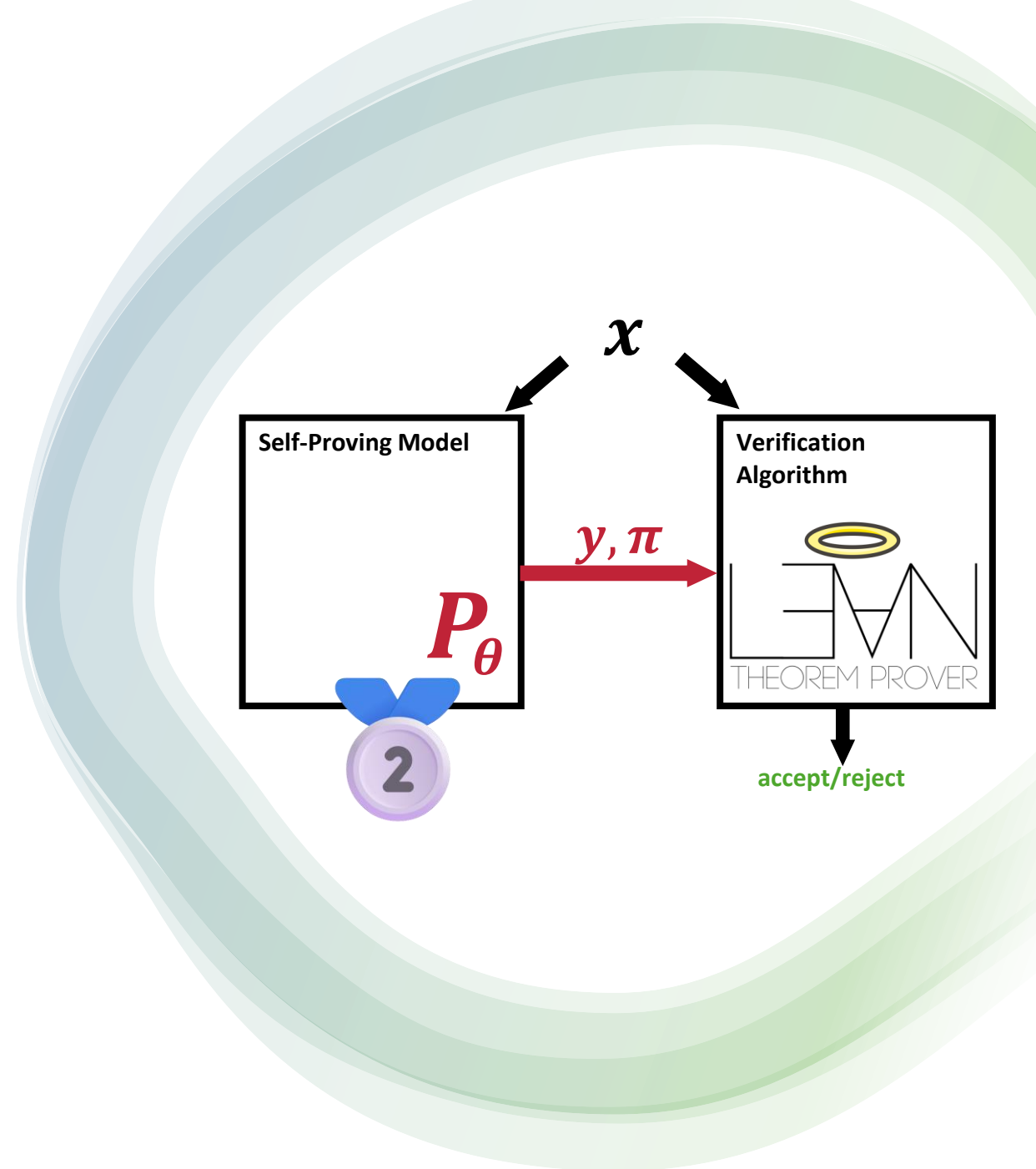
IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean 🏆 (DeepMind, 2024)


AI achieves silver-medal standard solving International Mathematical Olympiad problems

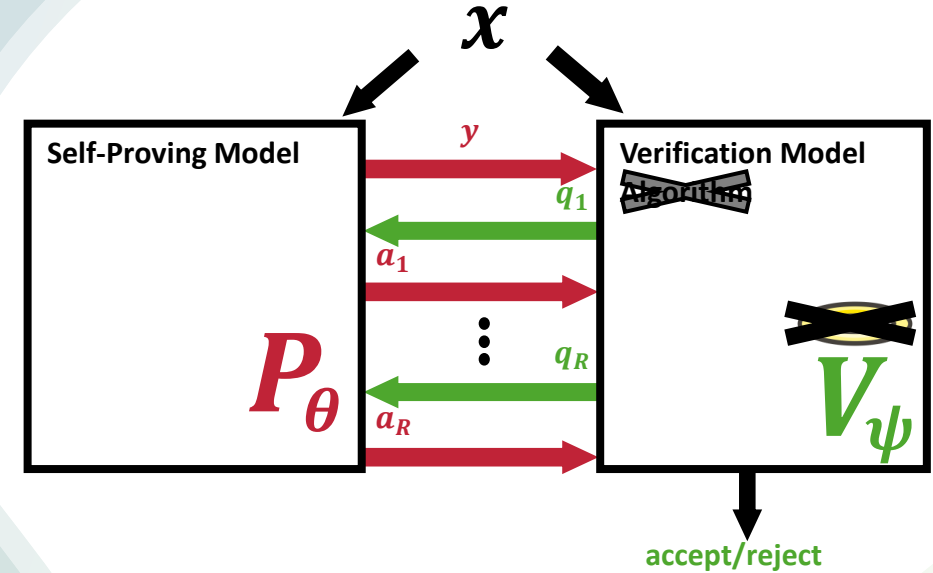
25 JULY 2024

AlphaProof and AlphaGeometry teams



IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify



IP ∩ ML [May 2024]

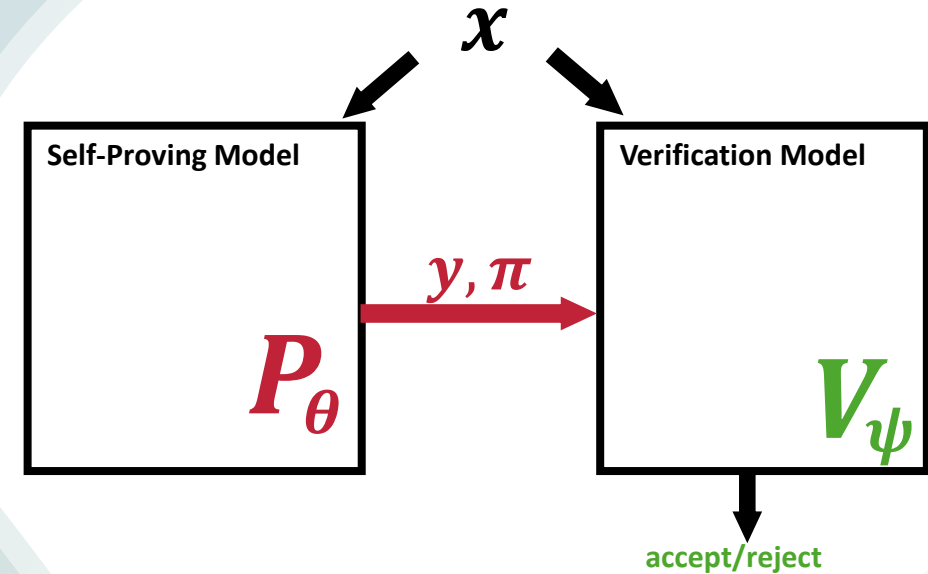
- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean 🏆 (DeepMind, 2024)
- Learning to Verify
 - **Prover Verifier Games (Anil et al., 2021)**

Consider the class NP: (1) Non-interactive proofs, (2) Decision problems.


Consider a strategy-finding game with two players: P_θ and V_ψ .

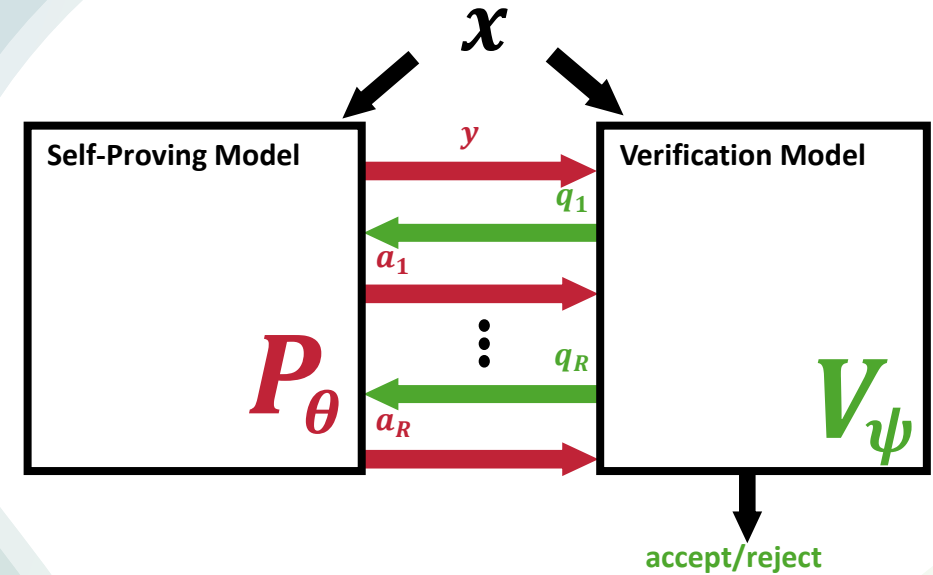
Does an equilibrium 🏆 → Completeness and Soundness? Does C&S → 🏆 ?

Order	🏆 → C&S	C&S → 🏆
P_θ first	No	No
V_ψ first	Yes!	Yes!
$\{P_\theta, V_\psi\}$ simul.	No	Yes!




IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify
 - Prover Verifier Games (Anil et al., 2021)
 - **Neural Interactive Proofs (Hammond & Adam-Day, 2024)**



IP \cap ML [May 2024]

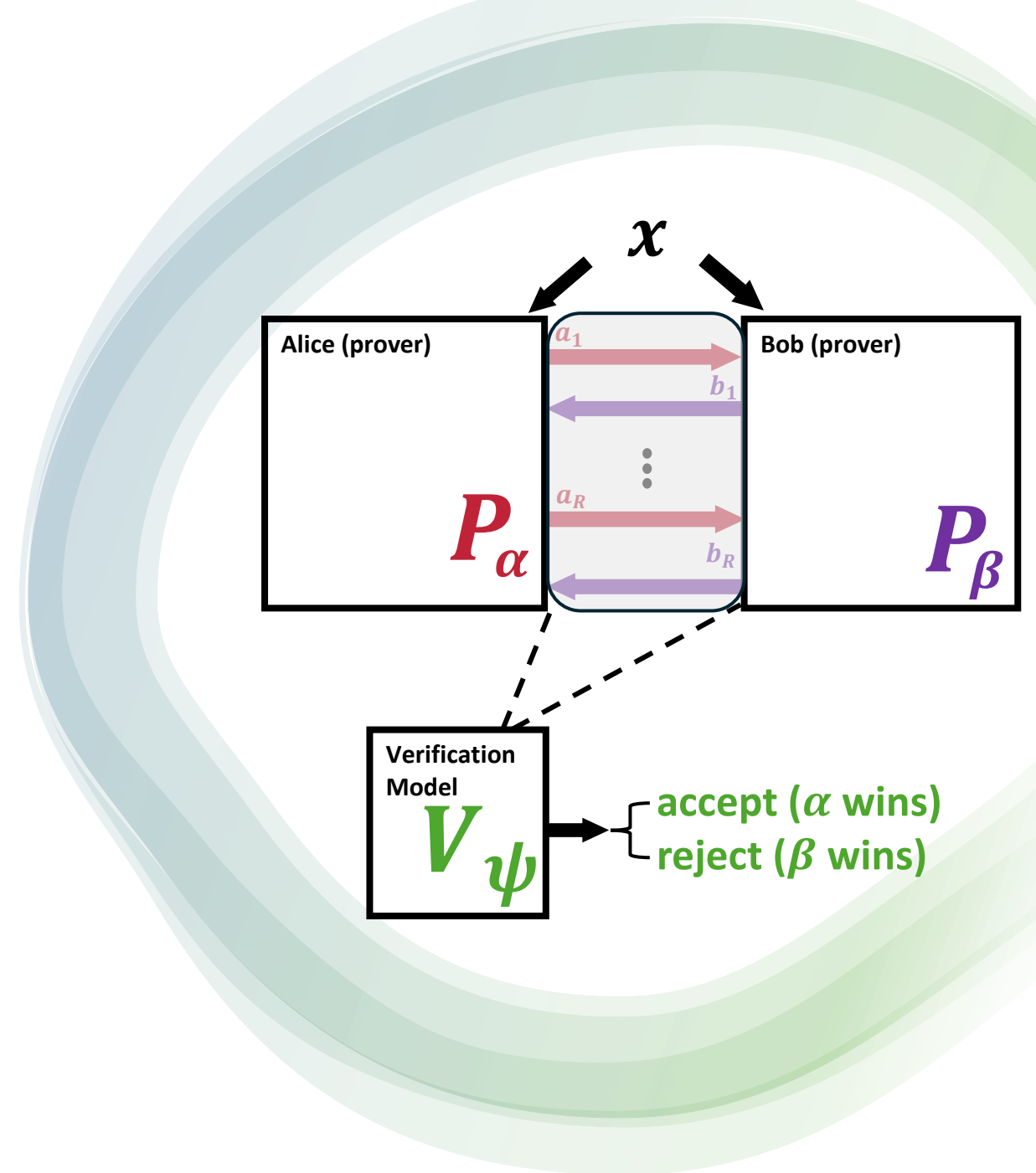
- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify
 - Prover Verifier Games (Anil et al., 2021)
 - Neural Interactive Proofs (Hammond & Adam-Day, 2024)
- Safety and alignment
 - **Debate Systems for AI Safety (Irving et al., 2017)**

AI safety via debate

Geoffrey Irving*

Paul Christiano

Dario Amodei



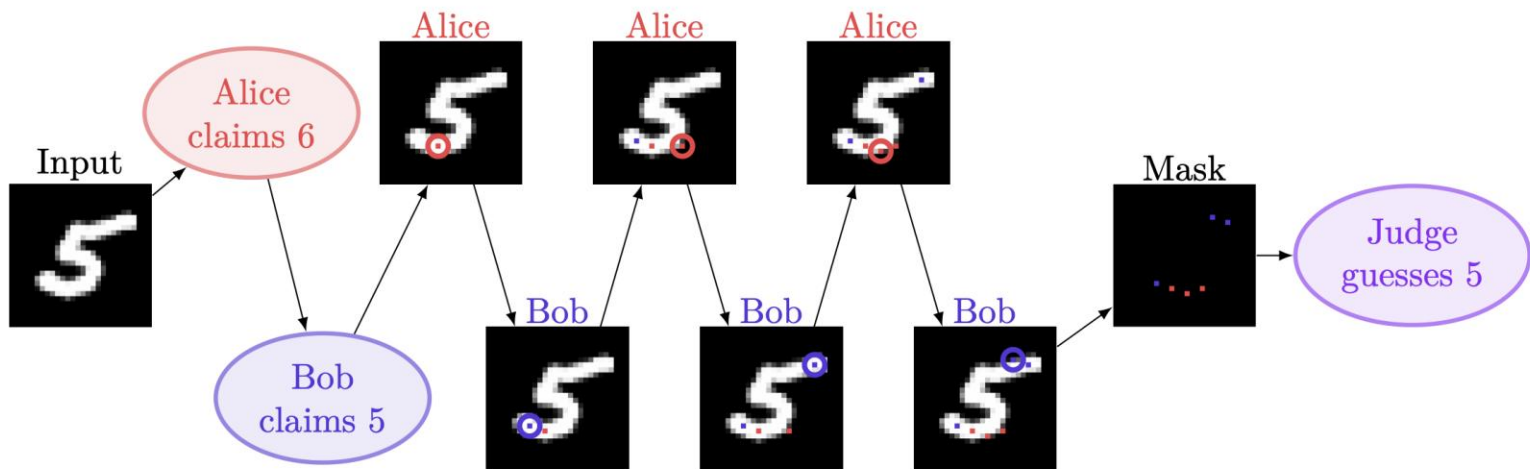
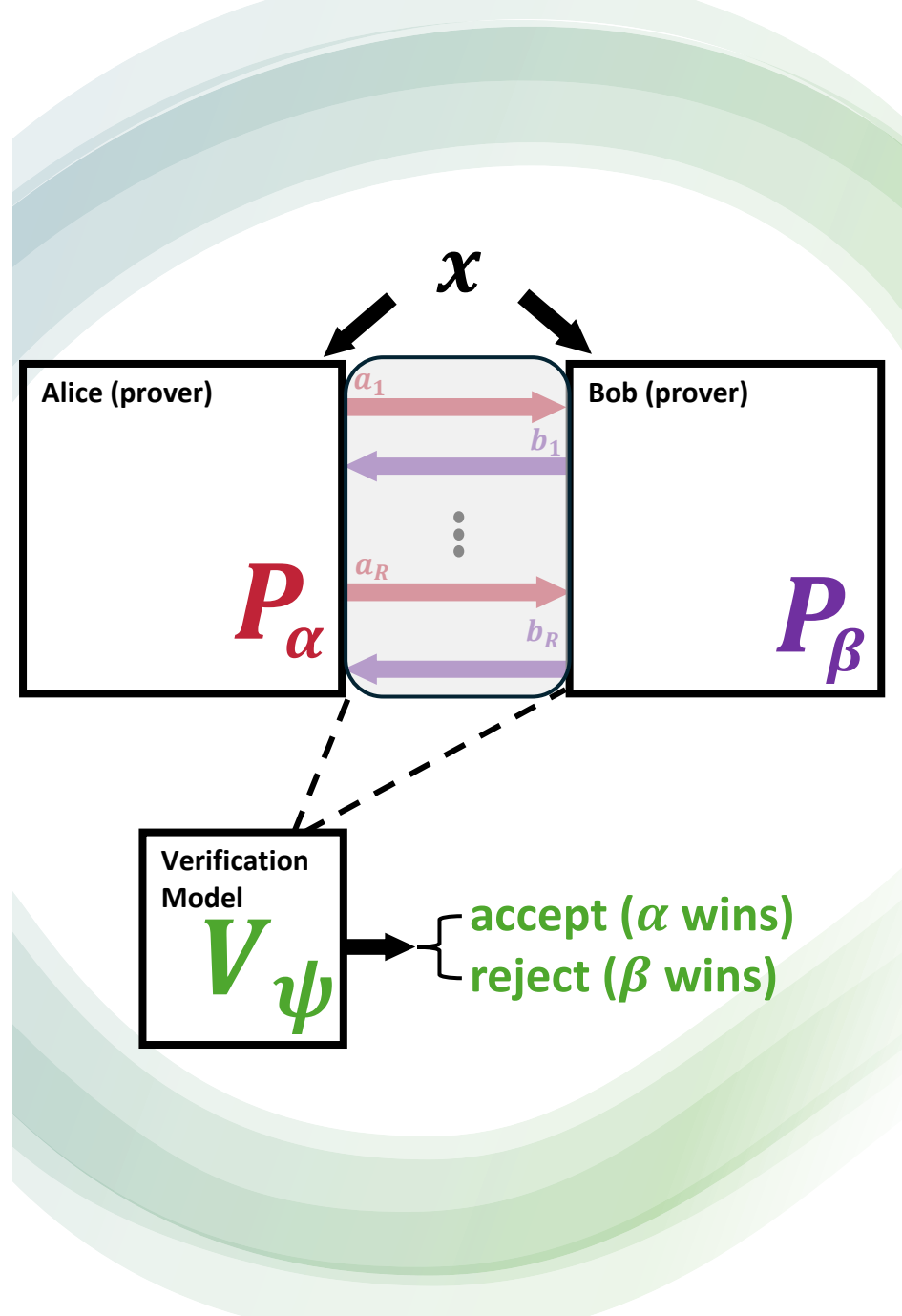



Figure 2: The MNIST debate game. A random MNIST image is shown to the two debating agents but not the judge. The debaters state their claimed label up front, then reveal one nonzero pixel per turn to the judge up to a total of 4 or 6. The judge sees the sparse mask of 4 or 6 pixels and chooses the winner based on which of the two labels has higher logit. The judge is trained in advance to recognize MNIST from random masks of nonzero pixels.

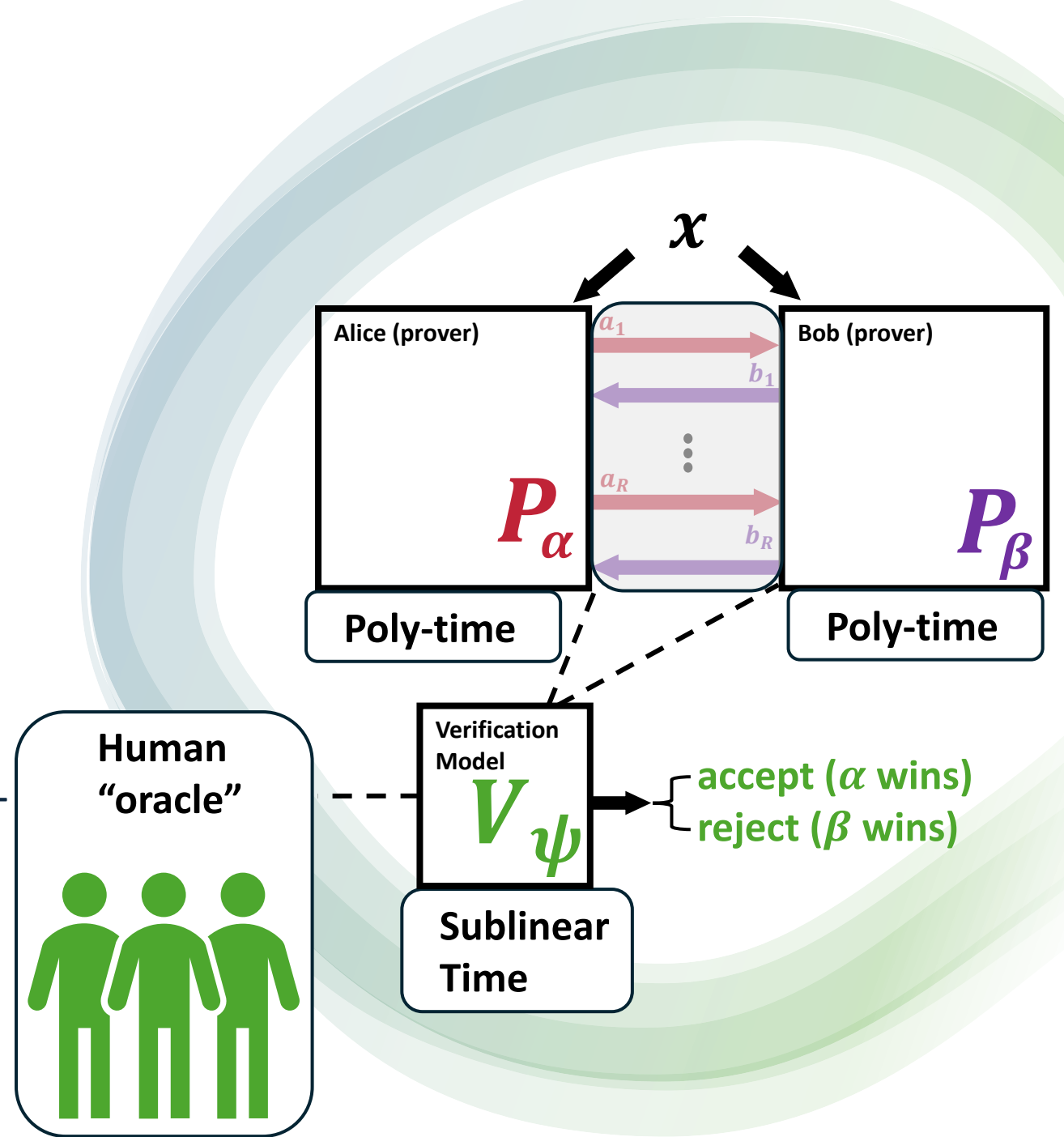
Pixels	First player	Judge accuracy (random pixels)	Honest win rate	
			No precommit	Precommit
4	honest	48.2%	51.0%	83.8%
	liar		68.4%	86.7%
	mean		59.7%	85.2%
6	honest	59.4%	67.4%	87.4%
	liar		81.5%	90.4%
	mean		74.4%	88.9%

Table 2: Results for debate on MNIST. We prespecify one player as honest and one as liar: when the honest player wins honesty is the best strategy. *No precommit* means the liar wins for any incorrect guess by the judge, even if the incorrect guess differs for different parts of the game tree. Lying is harder in the *Precommit* case, where the liar states their claim in advance of making moves and cannot change their story. With or without precommit, the honest player wins more often than a random judge, showing that honesty has an advantage.




IP ∩ ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify
 - Prover Verifier Games (Anil et al., 2021)
 - Neural Interactive Proofs (Hammond & Adam-Day, 2024)
- Safety and alignment
 - Debate Systems for AI Safety (Irving et al., 2017)
 - **Doubly-efficient debates for scalable AI safety (Brown-Cohen et al., 2024)**



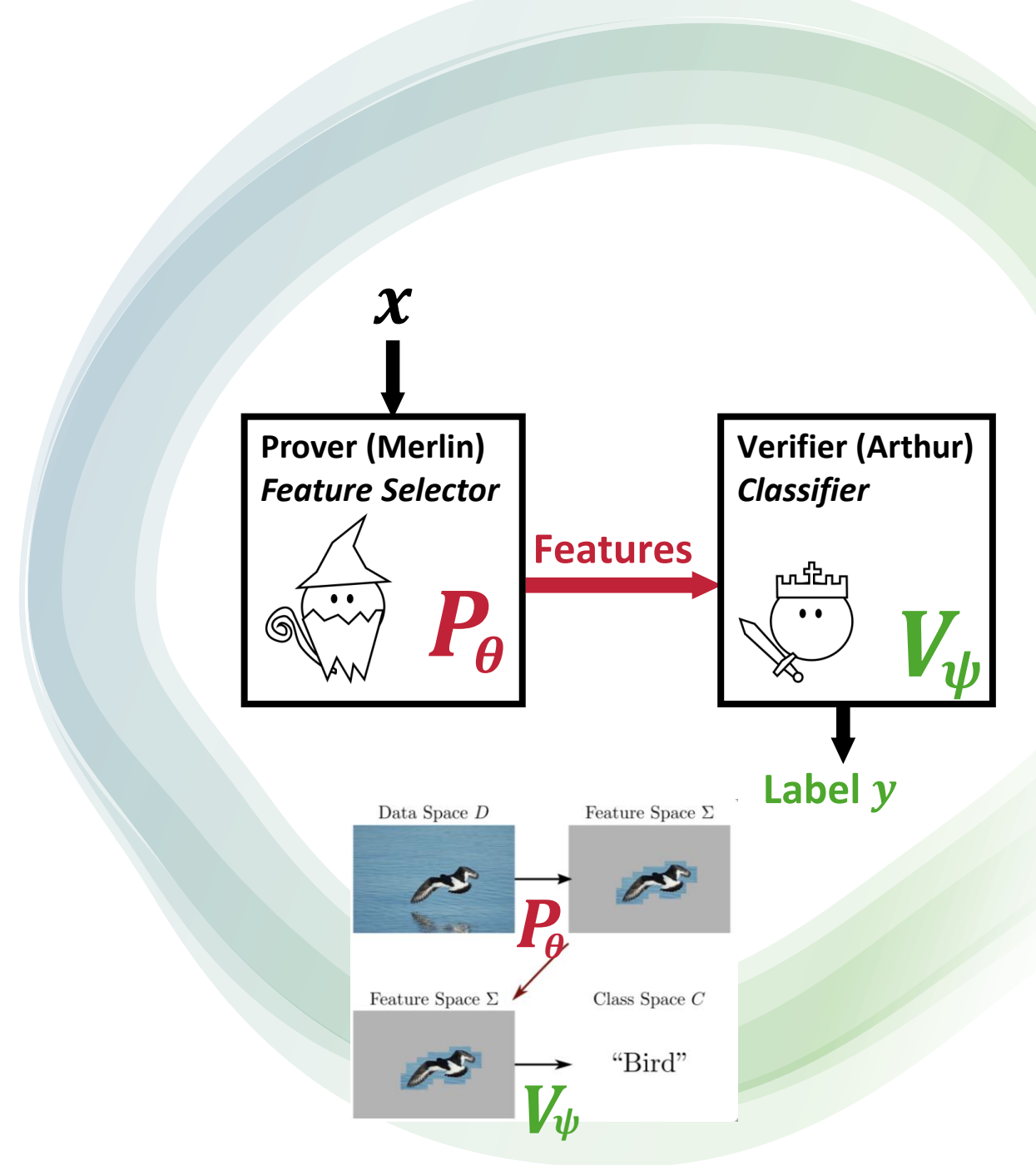
IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify
 - Prover Verifier Games (Anil et al., 2021)
 - Neural Interactive Proofs (Hammond & Adam-Day, 2024)
- Safety and alignment
 - Debate Systems for AI Safety (Irving et al., 2017)
 - Doubly-efficient debates for scalable AI safety (Brown-Cohen et al., 2024)
- **Interpretability**



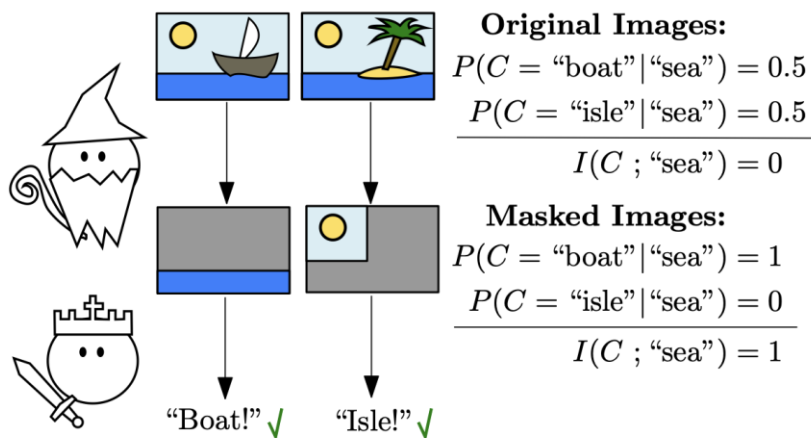
IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify
 - Prover Verifier Games (Anil et al., 2021)
 - Neural Interactive Proofs (Hammond & Adam-Day, 2024)
- Safety and alignment
 - Debate Systems for AI Safety (Irving et al., 2017)
 - Doubly-efficient debates for scalable AI safety (Brown-Cohen et al., 2024)
- Interpretability
 - **MA Classifiers (Waldchen et al., 2022)**

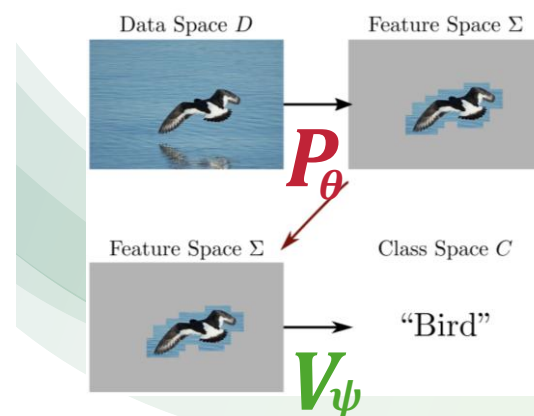
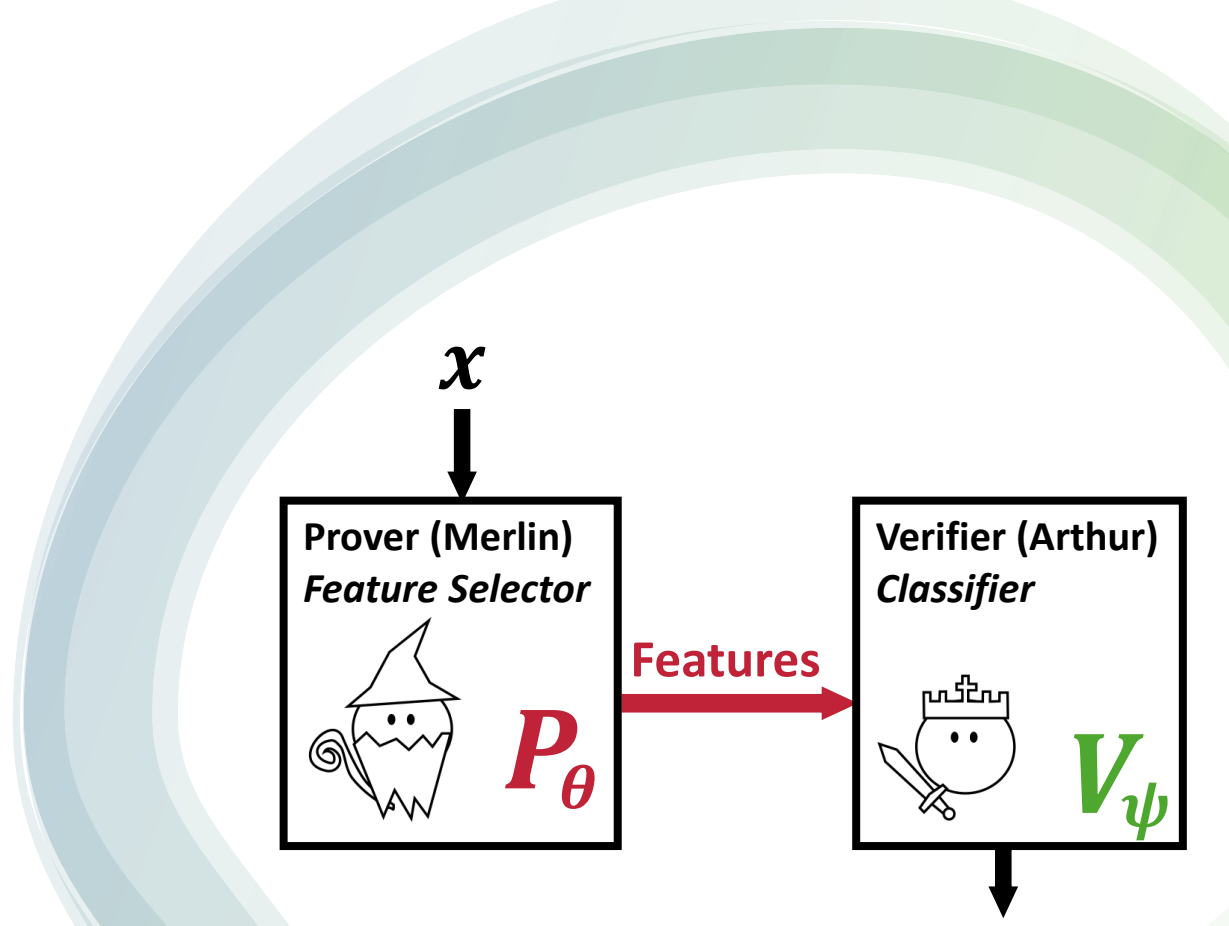
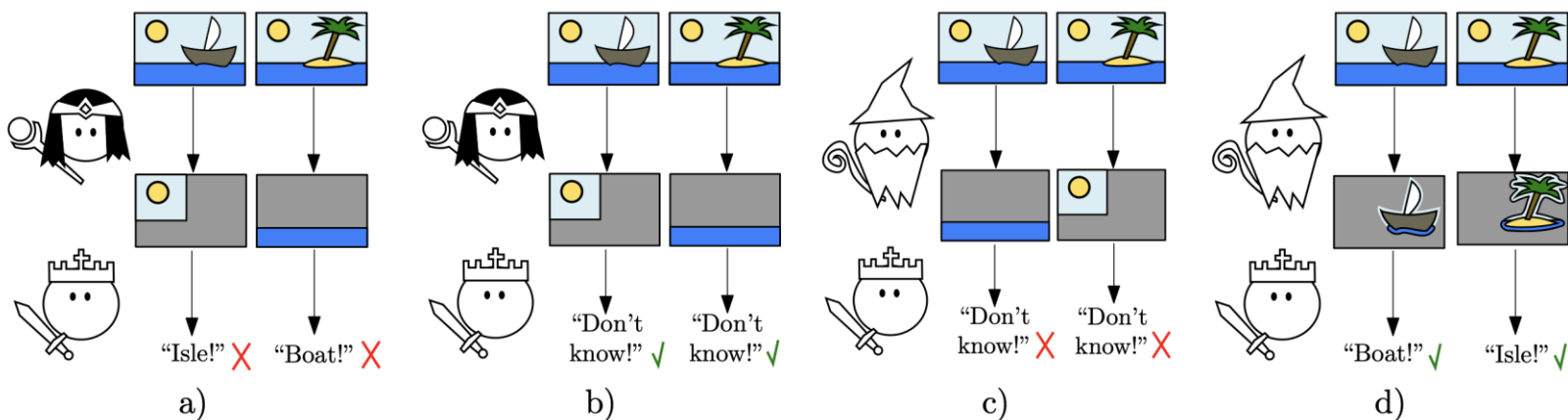


IP∩ML [May 2024]

Problem: “Spurious features”



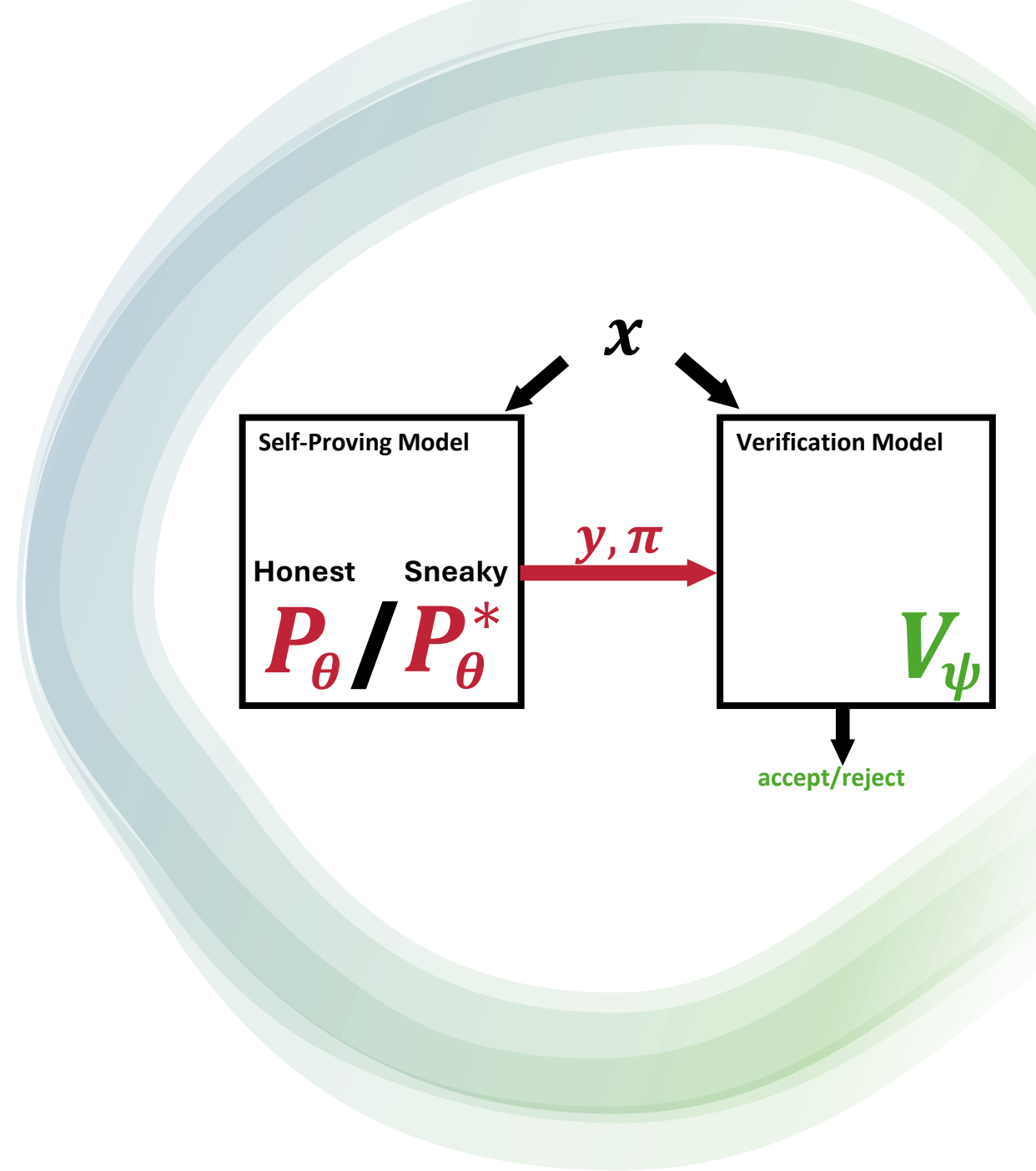
Solution: Soundness!



MA Classifiers (Waldchen et al., 2022)

IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify
 - Prover Verifier Games (Anil et al., 2021)
 - Neural Interactive Proofs (Hammond & Adam-Day, 2024)
- Safety and alignment
 - Debate Systems for AI Safety (Irving et al., 2017)
 - Doubly-efficient debates for scalable AI safety (Brown-Cohen et al., 2024)
- Interpretability
 - MA Classifiers (Waldchen et al., 2022)
 - **Prov. Ver. Games for legibility** (Hendrick* and Chen* et al., 2024)



IP∩ML [May 2024]

Idea: Introduce a smaller LLM as a Verifier.

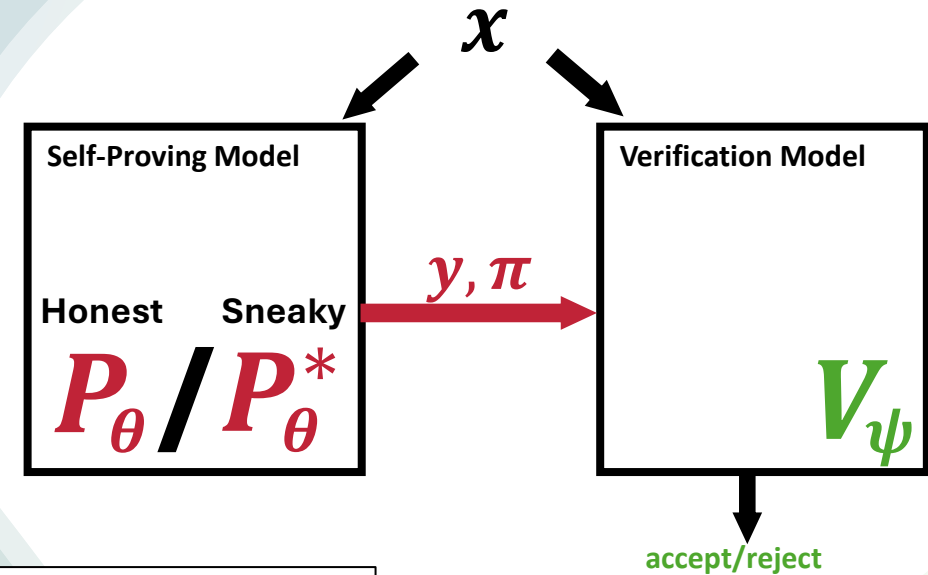
FT Prover LLM to act as honest / sneaky* prover

- Does this increase the (human) legibility of the Prover LLM's answers?

Joint training (simplified):

$$\psi \xleftarrow{\text{argmin}} L(\psi | \theta) := \ell(\psi | P_\theta) + \ell(\psi | P_\theta^*)$$
$$\ell(\psi | P) := -\mathbb{E}_x [H(\langle V_\psi, P \rangle(x), \mathbb{I}_{P(x)=y(x)})]$$

Cross-entropy loss between V_ψ 's decision logits (over 0,1) and the correctness of the prover.



IP∩ML [May 2024]

Idea: Introduce a smaller LLM as a Verifier.

FT Prover LLM to act as honest / sneaky* prover

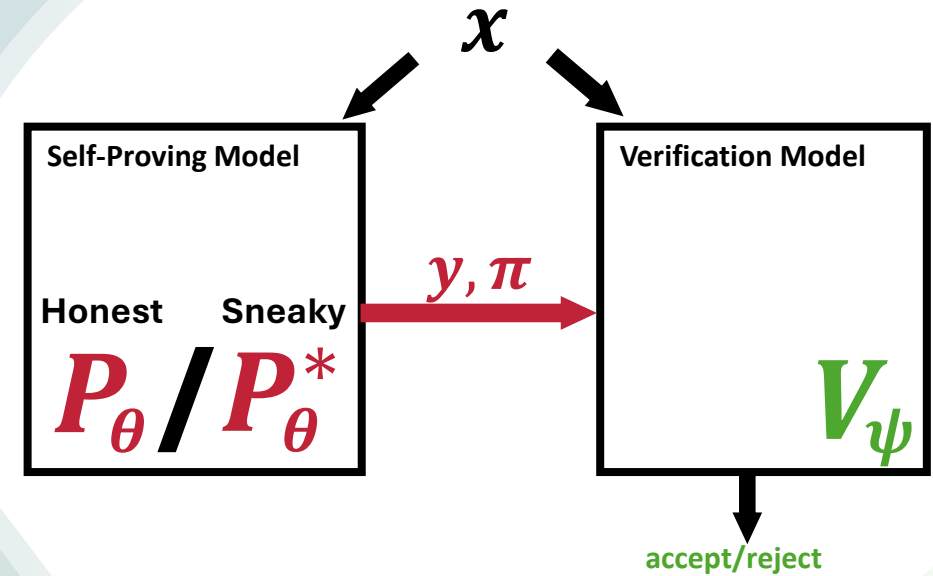
- Does this increase the (human) legibility of the Prover LLM's answers?

Joint training (simplified):

$$\psi \xleftarrow{\text{argmin}} L(\psi | \theta) := \ell(\psi | P_\theta) + \ell(\psi | P_\theta^*)$$
$$\ell(\psi | P) := -\mathbb{E}_x [H(\langle V_\psi, P \rangle(x), \mathbb{I}_{P(x)=y(x)})]$$

$$\theta \xleftarrow{\text{argmin}} \mathbb{E}_x [\langle V_\psi, P_\theta \rangle(x) \cdot \mathbf{1}_{P_\theta(x)=y(x)} + \langle V_\psi, P_\theta^* \rangle(x) \cdot \mathbf{1}_{P_\theta^*(x) \neq y(x)}]$$

(simplified)



IP∩ML [May 2024]

Idea: Introduce a smaller LLM as a Verifier.

FT Prover LLM to act as honest / sneaky

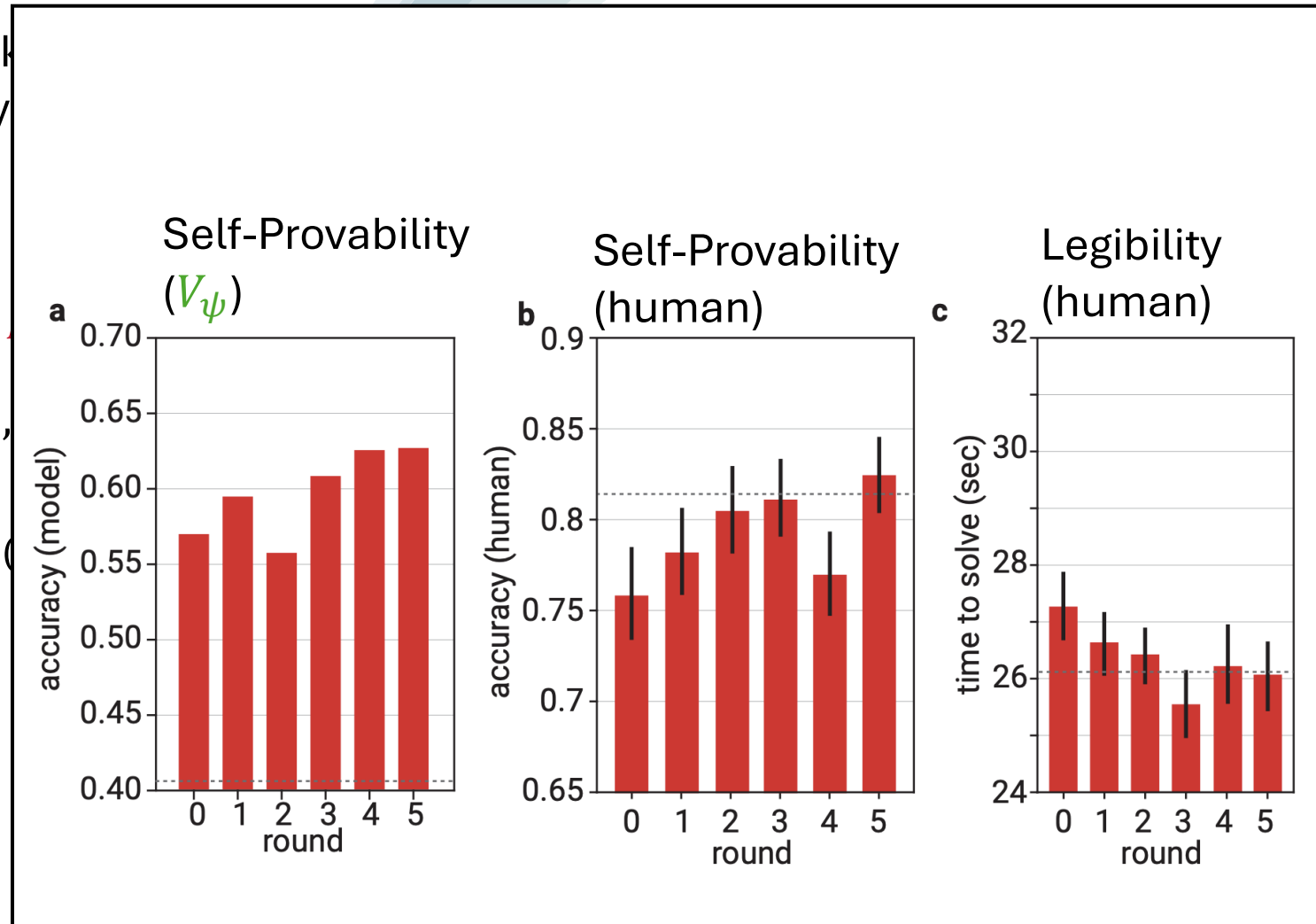
- Does this increase the (human) legibility of LLM's answers?

Joint training (simplified):

$$\psi \xleftarrow{\text{argmin}} L(\psi | \theta) := \ell(\psi | P_\theta) + \ell(\psi | P)$$

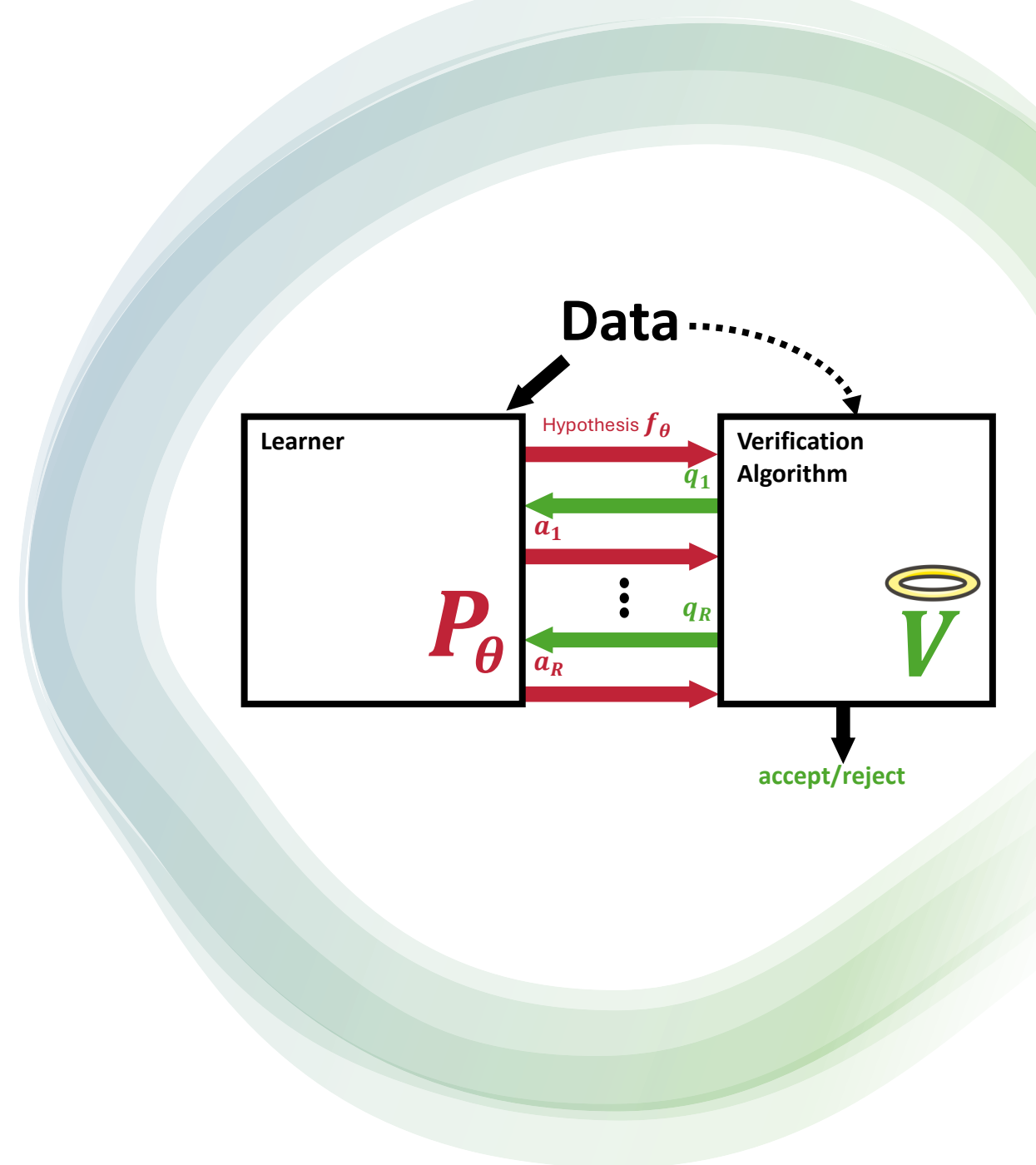
$$\ell(\psi | P) := -\mathbb{E}_x [H(\langle V_\psi, P \rangle(x))]$$

$$\theta \xleftarrow{\text{argmin}} \mathbb{E}_x [\langle V_\psi, P_\theta \rangle(x) \cdot \mathbf{1}_{P_\theta(x)=y(x)} + \langle V_\psi, P_\theta^* \rangle(x)]$$



IP \cap ML [May 2024]

- Learning to Prove...
 - ... in Coq (Gransden et al., 2015)
 - ... in Metamath (Polu and Sutskever, 2020)
 - ... in Lean (Yang et al., 2023)
 - ... in synthetic geometry (Trinh et al., 2024)
 - ... in Lean  (DeepMind, 2024)
- Learning to Verify
 - Prover Verifier Games (Anil et al., 2021)
 - Neural Interactive Proofs (Hammond & Adam-Day, 2024)
- Safety and alignment
 - Debate Systems for AI Safety (Irving et al., 2017)
 - Doubly-efficient debates for scalable AI safety (Brown-Cohen et al., 2024)
- Interpretability
 - MA Classifiers (Wäldchen et al., 2024)
 - Prov. Ver. Games for legibility (Hendrick* and Chen* et al., 2024)
- Verifying global accuracy
 - (Goldwasser etl al., 2021; Mutreja and Shafer, 2023)



Transcript Learning: Ancient Experiments [May '24]

- Charton (2024) showed that small GPT can learn to compute the GCD.
*Can it prove that its answer is **correct**?*
- A proof system for GCD:
 - Bézout's identity: Let $(x_1, x_2) \in \mathbb{N}$
if $z_1x_1 + z_2x_2$ divides x_1 and x_2 ,
then $z_1x_1 + z_2x_2 = \text{GCD}(x_1, x_2)$
 - V_{GCD} accepts iff $z_1x_1 + z_2x_2$ divides x_1 and x_2 , and

2 divides 92, 2 divides 78,
-11*92+13*78 = 2... Accept!

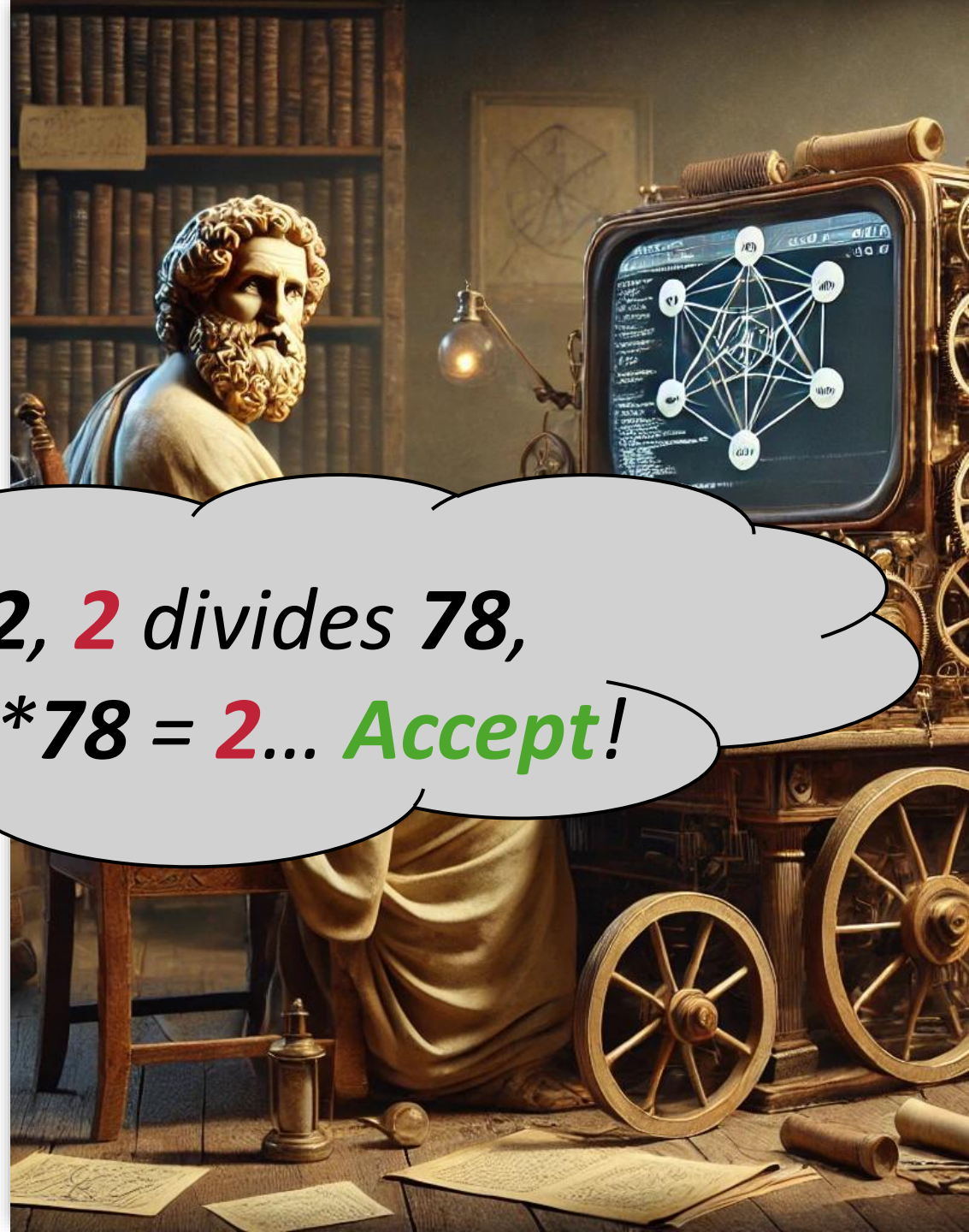
What is the $\text{GCD}(92, 78)$?

P_θ

Answer: **2**

Proof: $z_1 = -11, z_2 = 13$

V



Transcript Learning experiments

Learning method	Correctness	Self-Provability
GPT (baseline)	99.8%	-
GPT + TL	98.8%	60.3%

Transcript Learning experiments

Learning method	Correctness	Self-Provability
GPT (baseline)	99.8%	-
GPT + TL	98.8%	60.3%
GPT + Annotated TL	98.6%	96.7%

Γύρισα!
(I'm back!)

- In practice, annotations speed up learning
 - Intermediate steps in Euclid's algorithm



Annotations

Extended Euclidean algorithm

Input: Nonzero integers $x_0, x_1 \in \mathbb{N}$.

Output: Integers (y, z_0, z_1) , such that $y = \text{GCD}(x_0, x_1)$ and (z_0, z_1) are Bézout coefficients for (x_0, x_1) .

1: Initialize $r_0 = x_0, r_1 = x_1, s_0 = 1, s_1 = 0$, and $q = 0$.

2: **while** $r_1 \neq 0$ **do**

3: Update $q := \lfloor r_0/r_1 \rfloor$.

4: Update $(r_0, r_1) := (r_1, r_0 - q \times r_1)$.

5: Update $(s_0, s_1) := (s_1, s_0 - q \times s_1)$.

6: Output GCD $y = r_0$ and Bézout coefficients $z_0 := s_0$ and $z_1 := (r_0 - s_0 \cdot x_0)/x_1$.

} **Annotation**



Input		GCD	Annotation			Bézout coefs	
x_0	x_1	y	\vec{s}_0	\vec{r}_0	\vec{q}	z_0	z_1
46	39		1	46	1		
			0	39	5		
			1	7	1		
			-5	4	1		
			6	3	3		
		1				-11	13

Annotations

Extended Euclidean algorithm

Input: Nonzero integers $x_0, x_1 \in \mathbb{N}$.

Output: Integers (y, z_0, z_1) , such that $y = \text{GCD}(x_0, x_1)$ and (z_0, z_1) are Bézout coefficients for (x_0, x_1) .

1: Initialize $r_0 = x_0, r_1 = x_1, s_0 = 1, s_1 = 0$, and $q = 0$.

2: **while** $r_1 \neq 0$ **do**

3: Update $q := \lfloor r_0/r_1 \rfloor$.

4: Update $(r_0, r_1) := (r_1, r_0 - q \times r_1)$.

5: Update $(s_0, s_1) := (s_1, s_0 - q \times s_1)$.

6: Output GCD $y = r_0$ and Bézout coefficients $z_0 := s_0$ and $z_1 := (r_0 - s_0 \cdot x_0)/x_1$.

} **Annotation**



Input		GCD	Annotation			Bézout coeffs	
x_0	x_1	y	\vec{s}_0	\vec{r}_0	\vec{q}	z_0	z_1
46	39		1	46	1		
			0	39	5		
			1	7	1		
			-5	4	1		
			6	3	3		
		1				-11	13

Decimal encoding

+, 4, 6, x0, +, 3, 9, x1, +, 1, y,
 +, 1, z0', +, 4, 6, z1', +, 1, q',
 +, 0, z0'', +, 3, 9, z1'', +, 5, q''
 +, 1, z0''', +, 7, z1''', +, 1, q''',
 -, 1, 1, z0, +, 1, 3, z1

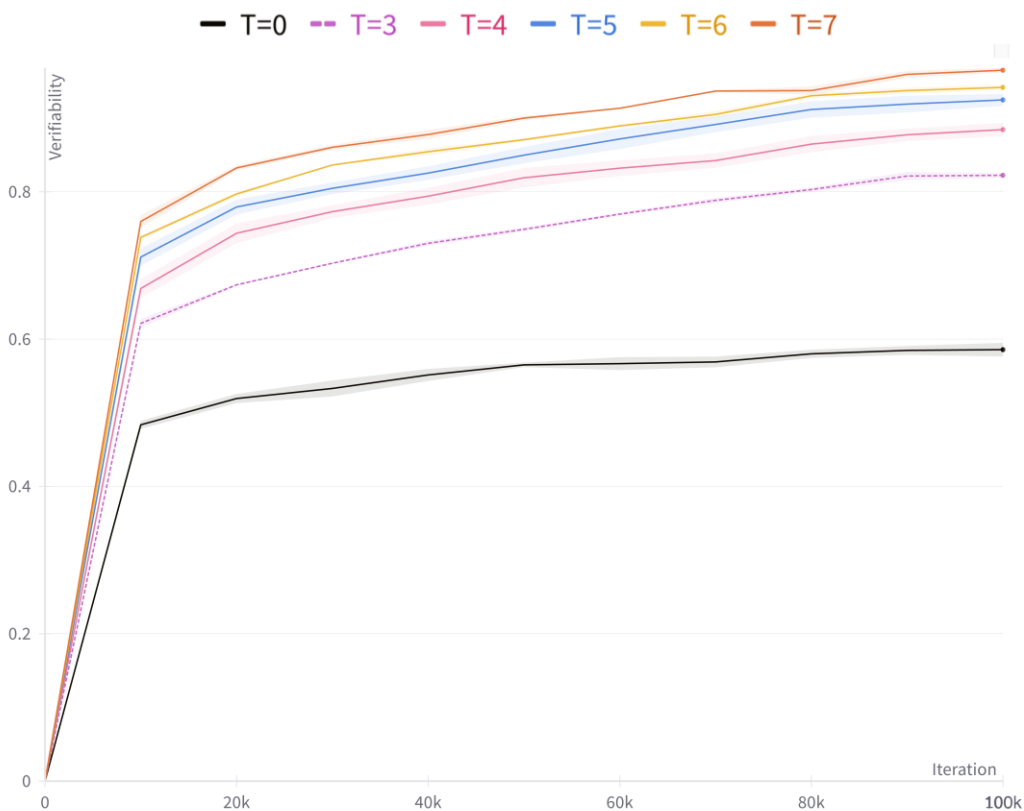
sign tokens

digits

delimiters

Annotated Transcript Learning

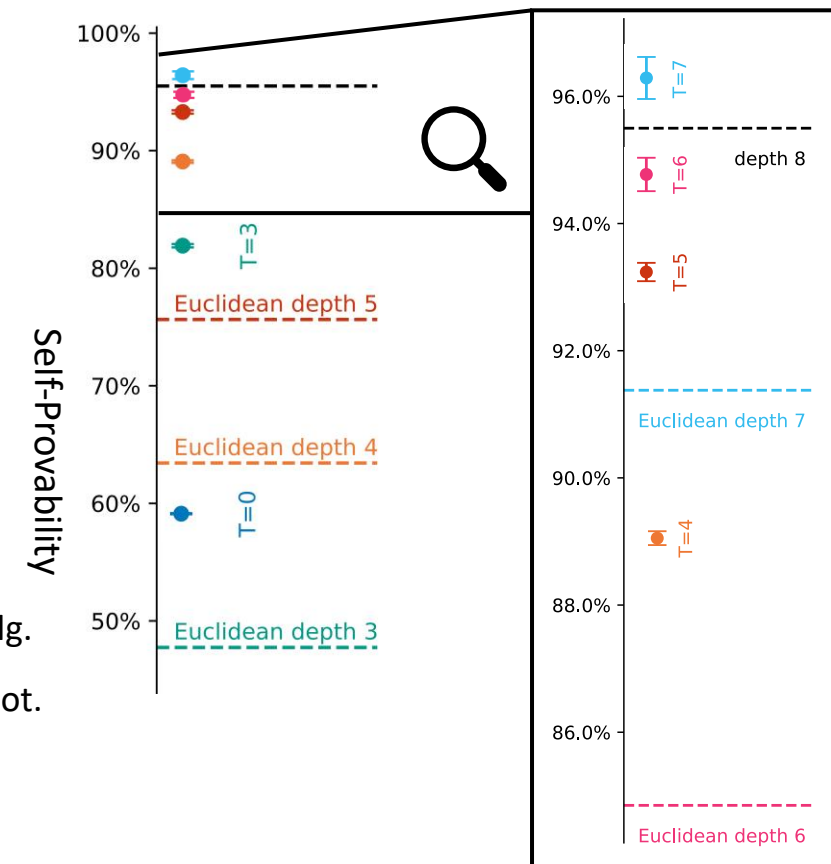
Longer annotations help



Models generalize beyond annotations

Input			Annotation		
x_0	x_1	y	\vec{s}_0	\vec{r}_0	\vec{q}
46	39		1	46	1
			0	39	5
			1	7	1
			-5	4	1
			6	3	3
		1			

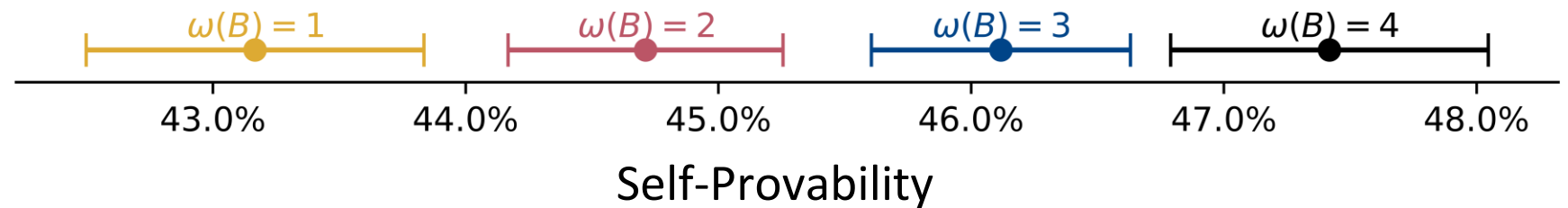
Depth(x_0, x_1) = #steps in Euclidean alg.
 = length of "ideal" annot.



Base of Representation matters

“Early during training, transformers learn to predict products of divisors of the base B used to represent integers.” (Charton, 2024)

- Let $\omega(B)$ denote the number of primes in the factorization of base B .
- Then $\omega(B)$ determines Self-Provability (similarly to Charton’s observation).



RLVF: from Theory to Practice

RLVF introduced as a method for training Self-Proving models (Amit et al.)

May 2024

Lemma B.6 (RLVF gradient estimation). *Fix an input distribution μ over Σ^* and V with round complexity R and answer length L_a . For any transcript (x, y, q_1, \dots, a_R) denote the indicator random variable which equals 1 if and only if it accepts the transcript. For any model P_θ , denote by $\text{ver}(\theta)$ the verifiability of P_θ with respect to μ (Definition 3.4). Then, for any θ ,*

$$\nabla_{\theta} \text{ver}(\theta) = \mathbb{E}_{\substack{x \sim \mu \\ y \sim P_{\theta}(x) \\ (q_r, a_r)_{r=1}^R}} \left[\text{Acc}_V(x, y, q_1, \dots, a_R) \cdot \sum_{\substack{r \in [R] \cup \{0\} \\ s \in [L_a]}} \vec{d}_s(\theta) \right]$$

where $(q_r, a_r)_{r=1}^R$ are as sampled in lines 5-6 of Algorithm 2, and $\vec{d}_s(\theta)$ is as defined in l

Algorithm 2: Reinforcement Learning from Verifier Feedback (RLVF)

Hyperparameters: Learning rate $\lambda \in (0, 1)$ and number of samples $N \in \mathbb{N}$.

Inputs: An input distribution μ over Σ^* , a model family $\{P_\theta\}$, and initial parameters $\theta_0 \in \text{md}(\mu; \Sigma)$.

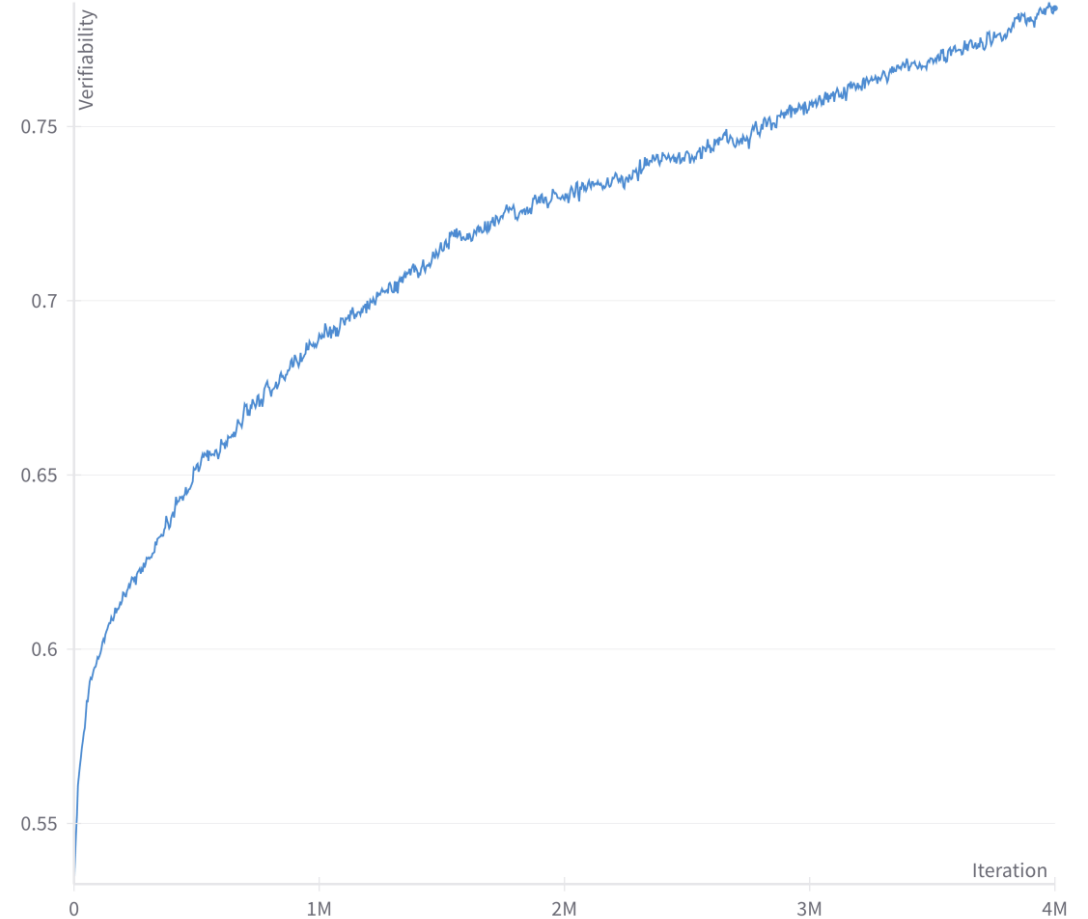


Figure 6: **RLVF Verifiability as a function of the number of samples N .** Starting from a base model with Verifiability 48% (obtained via Transcript Learning), in each iteration a batch of 2048 inputs are sampled; the model generates a proof for each; the Verifier is used to check which proofs are accepted; then, the model parameters are updated accordingly (see Algorithm 2). Verifiability was evaluated on a held-out dataset of 1k inputs.

RLVF: from Theory to Practice

RLVF introduced as a method for training Self-Proving models (Amit et al.)

May
2024

Nov 2024

RLVR (Lambert et al.) implements RLVF for Tülu3
+ KL regularization, RM initialization, data shuffling,
No-EOS penalty, advantage normalization
+ ZeRO 3 (Rajbhandari et al. 2020), asynchronous RL
(Noukhovitch et al. 2024), ablations...
→ Beats DPO on GSM8K, MATH, IFEval at 8B params

RLVF: from Theory to Practice

RLVF introduced as a method for training Self-Proving models (Amit et al.)

 **Med-RLVR**
(Zheng et al.)

May
2024

Nov 2024

Feb
2025

RLVR (Lambert et al.) implements RLVF for Tülu3
+ KL regularization, RM initialization, data shuffling,
No-EOS penalty, advantage normalization
+ ZeRO 3 (Rajbhandari et al. 2020), asynchronous RL
(Noukhovitch et al. 2024), ablations...
→ Beats DPO on GSM8K, MATH, IFEval at 8B params

RLVF: from Theory to Practice

RLVF introduced as a method for training Self-Proving models (Amit et al.)

 **Med-RLVR**
(Zheng et al.)

May
2024

Nov 2024


Feb
2025

May
2025

RLVR (Lambert et al.) implements RLVF for Tülu3
+ KL regularization, RM initialization, data shuffling,
No-EOS penalty, advantage normalization
+ ZeRO 3 (Rajbhandari et al. 2020), asynchronous RL
(Noukhovitch et al. 2024), ablations...
→ Beats DPO on GSM8K, MATH, IFEval at 8B params

 **RLVR-World** (Wu et al.)

 **ManipLVM-R1** (Song et al.)

 **REcon** (Zhou et al.)

RLVF: from Theory to Practice

RLVF introduced as a method for training Self-Proving models (Amit et al.)

May 2024

Nov 2024

 **Med-RLVR**
(Zheng et al.)

Feb 2025


May 2025



RLVR (Lambert et al.) implements RLVF for Tülu3
+ KL regularization, RM initialization, data shuffling,
No-EOS penalty, advantage normalization
+ ZeRO 3 (Rajbhandari et al. 2020), asynchronous RL
(Noukhovitch et al. 2024), ablations...
→ Beats DPO on GSM8K, MATH, IFEval at 8B params

 **RLVR-World** (Wu et al.)

 **ManipLVM-R1** (Song et al.)

 **REcon** (Zhou et al.)



The Hills are Alive.....?



V

V

V

V

V



The Hills are Alive.....?

So far, Verifiers
have rarely been
provably sound

RLVF was shown effective for
post-training capabilities. What about
worst-case guarantees?



The Hills are Alive.....?

Lesson plan

- Reminder: Proofs and Models
- Defining Self-Proving models
- Learning algorithms for Self-Proving models:
 - Transcript Learning (TL)
 - Reinforcement Learning from Verifier Feedback (RLVF)
- Self-Proving models in practice
- **Proving convergence guarantees for TL and RLVF**

TL convergence

Algorithm 1: Transcript Learning (TL)

Hyperparameters: Learning rate $\lambda \in (0, 1)$ and number of samples $N \in \mathbb{N}$.

Input: An autoregressive model family $\{P_\theta\}_{\theta \in \mathbb{R}^d}$, verifier specification (code) V , and sample access to an input distribution μ and an accepting transcript generator $\mathcal{T}_V^*(\cdot)$.

Output: A vector of parameters $\bar{\theta} \in \mathbb{R}^d$.

```

1 Initialize  $\theta_0 := \vec{0}$ .
2 for  $i = 0, \dots, N - 1$  do
3   Sample  $x \sim \mu$  and  $\pi^* = (y^*, q_1^*, a_1^*, \dots, q_R^*, a_R^*) \sim \mathcal{T}_V^*(x)$ . Denote  $a_0 := y^*$ .
4   foreach Round of interaction  $r = 0, \dots, R$  do
5     Let  $S(r)$  denote the indices of the  $r^{\text{th}}$  answer  $a_r$  in  $\pi^*$ , and let  $\pi_{<s}$  denote the prefix of
6     the partial transcript  $(y, q_1^*, a_1^*, \dots, q_r^*)$ .
7     for  $s \in S(r)$  do
8       Compute # Forwards and backwards pass
9          $\alpha_s(\theta_i) := \Pr_{\sigma \sim p_{\theta_i}(x\pi_{<s})}[\sigma = \pi_s^*]$ 
10         $\vec{d}_s(\theta_i) := \nabla_\theta \log \alpha_s(\theta_i) = \nabla_\theta \log \Pr_{\sigma \sim p_{\theta_i}(x\pi_{<s})}[\sigma = \pi_s^*]$ 
11      Update
12         $\theta_{i+1} := \theta_i + \lambda \cdot \prod_{\substack{r \in [R] \cup \{0\} \\ s \in S(r)}} \alpha_s(\theta_i) \cdot \sum_{\substack{r \in [R] \cup \{0\} \\ s \in S(r)}} \vec{d}_s(\theta_i)$ 
13    end for
14  end for
15 Output  $\bar{\theta} := \frac{1}{N} \sum_{i \in [N]} \theta_i$ .

```

Lemma (TL gradient estimation). Fix an input distribution μ over Σ^* and a verifier V with round complexity R and answer length L_a . Fix an honest transcript generator \mathcal{T}_V^* . Let θ denote the parameters of the model P_θ , let $A(\theta)$ be as defined above and let the terms $S(r)$, $\alpha_s(\theta)$, and $\vec{d}_s(\theta)$ be as defined in Algorithm 1. Then,

$$\nabla A(\theta) = \mathbb{E}_{\substack{x \sim \mu \\ \pi^* \sim \mathcal{T}_V^*}} \left[\prod_{\substack{r \in [R] \cup \{0\} \\ s \in S(r)}} \alpha_s(\theta) \cdot \sum_{\substack{r \in [R] \cup \{0\} \\ s \in S(r)}} \vec{d}_s(\theta) \right].$$

Theorem. Fix a verifier V , an input distribution μ , an autoregressive model family $\{P_\theta\}_{\theta \in \mathbb{R}^d}$, and a norm $\|\cdot\|$ on \mathbb{R}^d . Fix an honest transcript generator \mathcal{T}_V^* , and assume that the agreement function $A(\theta) := \Pr[\pi = \pi^*]$ is concave in θ , where the verifier queries are the same in π^* and π , and the probability is over $x \sim \mu$, $\pi_\theta \sim \mathcal{T}_V^\theta(x)$, and $\pi^* \sim \mathcal{T}_V^*(x)$. For any $\varepsilon > 0$, let B_{Norm} , B_{Lip} and C be upper-bounds such that the following conditions hold:

- There exists $\theta^* \in \mathbb{R}^d$ with $\|\theta^*\| < B_{\text{Norm}}$ such that $A(\theta^*) \geq 1 - \varepsilon/2$.
- For all θ , the logits of P_θ are B_{Lip} -Lipschitz in θ . That is, $\sup_{\substack{\theta \in \mathbb{R}^d \\ z \in \Sigma^*}} \|\nabla_\theta \log p_\theta(z)\| \leq B_{\text{Lip}}$.
- In the proof system defined by V , the total number of tokens (over all rounds) is at most C .

Denote by $\bar{\theta}$ the output of TL running for $N \geq (4 \cdot C^2 \cdot B_{\text{Norm}}^2 \cdot B_{\text{Lip}}^2) / \varepsilon^2$ iterations and learning rate $\lambda = B_{\text{Norm}} / (C B_{\text{Lip}} \sqrt{N})$. Then the expected Verifiability (over the randomness of the samples collected by TL) of $\bar{\theta}$ is at least $1 - \varepsilon$. That is, $\mathbb{E}_{\bar{\theta}}[\text{ver}_{V, \mu}(\bar{\theta})] \geq 1 - \varepsilon$.

Towards RLVF convergence

Algorithm 2: Reinforcement Learning from Verifier Feedback (RLVF)

Hyperparameters: Learning rate $\lambda \in (0, 1)$ and number of samples $N \in \mathbb{N}$.

Input: An autoregressive model family $\{P_\theta\}_{\theta \in \mathbb{R}^d}$, initial parameters $\theta_0 \in \mathbb{R}^d$, verifier specification (code) V , and sample access to an input distribution μ .

Output: A vector of parameters $\bar{\theta} \in \mathbb{R}^d$.

```

1 for  $i = 0, \dots, N - 1$  do
2   Sample  $x \sim \mu$ .
3   Initialize  $a_0 := y \sim P_{\theta_i}(x)$ .
4   foreach Round of interaction  $r = 1, \dots, R$  do
5     Sample the  $r^{\text{th}}$  query                                # Emulate the verifier
6      $q_r \sim V_q(x, a_0, q_1, a_1, \dots, q_{r-1}, a_{r-1})$ .
7     Sample the  $r^{\text{th}}$  answer                                # Forwards pass
8      $a_r \sim P_{\theta_i}(x, a_0, q_1, a_1, \dots, q_r)$ .
9     Let  $\tau_r := (a_0, q_1, \dots, a_{r-1}, q_r)$ .
10    for  $s \in [L_a]$  do
11      Let  $a_{r,s}$  denote the  $s^{\text{th}}$  token in  $a_r$ . Compute    # Backwards pass
12       $\vec{d}_s(\theta_i) := \nabla_\theta \log_{\sigma \sim p_{\theta_i}(x, \tau_r)} [\sigma = a_{r,s}]$ .
13    if  $V(x, y, q_1, a_1, \dots, q_R, a_R)$  accepts then
14      Update
15       $\theta_{i+1} := \theta_i + \lambda \cdot \sum_{\substack{r \in [R] \cup \{0\} \\ s \in [L_a]}} \vec{d}_s(\theta_i)$ .
16  Output  $\bar{\theta} := \frac{1}{N} \sum_{i \in [N]} \theta_i$ .
```

Lemma (RLVF gradient estimation). *Fix an input distribution μ over Σ^* and a verifier V with round complexity R and answer length L_a . For any transcript (x, y, q_1, \dots, a_R) we let $\text{Acc}_V(x, y, q_1, \dots, a_R)$ denote the indicator random variable which equals 1 if and only if V accepts the transcript. For any model P_θ , denote by $\text{ver}(\theta)$ the verifiability of P_θ with respect to V and μ . Then, for any θ ,*

$$\nabla_\theta \text{ver}(\theta) = \mathbb{E}_{\substack{x \sim \mu \\ y \sim P_\theta(x) \\ (q_r, a_r)_{r=1}^R}} \left[\text{Acc}_V(x, y, q_1, \dots, a_R) \cdot \sum_{\substack{r \in [R] \cup \{0\} \\ s \in [L_a]}} \vec{d}_s(\theta) \right]$$

where $(q_r, a_r)_{r=1}^R$ are as in lines 5-6 of Algorithm 2, and $\vec{d}_s(\theta)$ is as defined in line 8 therein.

But proving RL convergence guarantees is challenging...

What next?

- (Multiprover) proof system for AI Scientist [joint w/ Yoshua Bengio, Shafi Goldwasser, Oliver Richardson, Avishay Tal]
- Catching up RLVF theory 🤪
 - Sample complexity bounds [joint w/ Nicolas Flammarion]
 - RLVF with average-case soundness leads to...
 - A useful and meaningful notion! (Optimistic)
 - Verifier hacking (Pessimistic)
 - A position paper... (Pragmatic)
- Practical Self-Proving Models *for sound Verifiers*
 - Lean-based DSLs [joint w/ Thomas Bourgeat]
- Universal (“Foundation”) Self-Proving models
 - So far: $\forall V \exists P_\theta$. I.e., need to learn a different prover for each verifier.
 - Can we have $\exists P_\theta \forall V$ (in a restricted class)?
I.e., can we learn a prover P_θ such that for all $V \in \mathcal{V}$,
 $P_\theta(x, V)$ outputs y , and proves correctness to V ?
- “Fundamental Theorem of Self-Provable learning”
 - Is there a (combinatorial) dimension of the problem/proof system that captures the sample complexity of learning Self-Proving models?

Seeking strong students with

- Background in Theoretical CS
- Interest in Self-Proving models

Visit EPFL and work with us in 2026/7.
Funding available*

Apply: orrrp.net/visitEPFL

Contact: orr.paradise@epfl.ch

