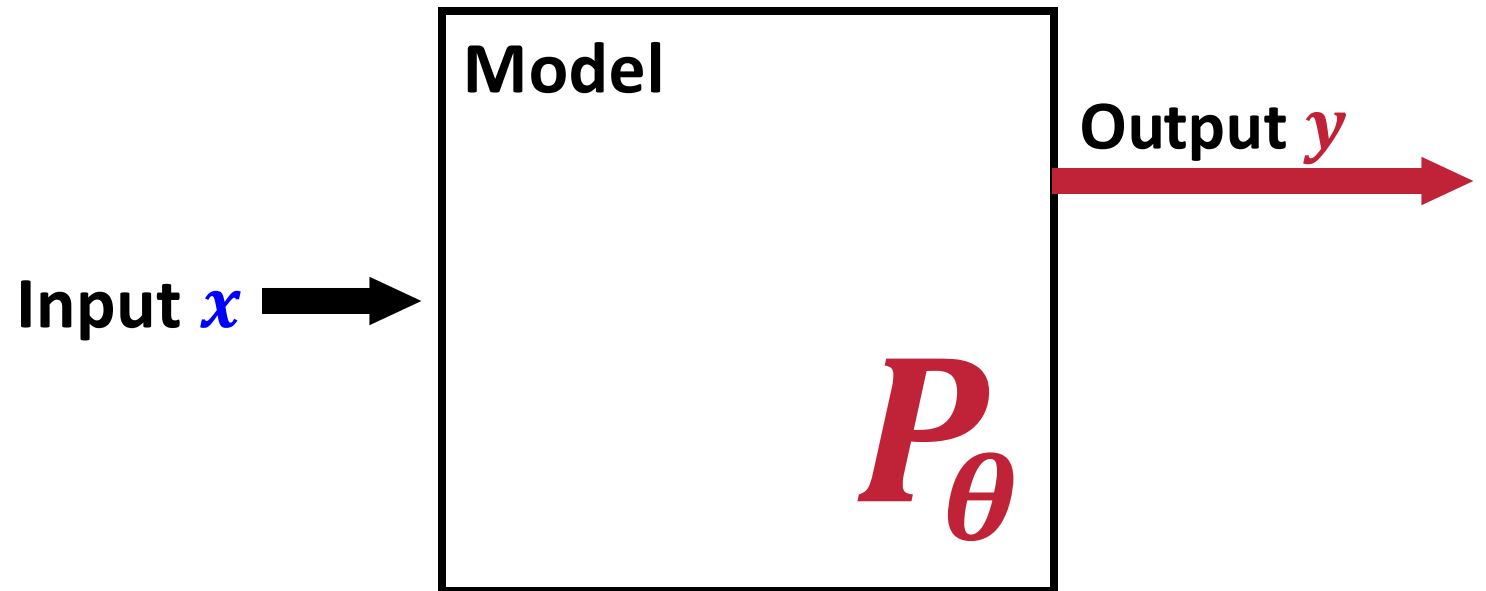


# Models That Prove Their Own Correctness

Orr Paradise

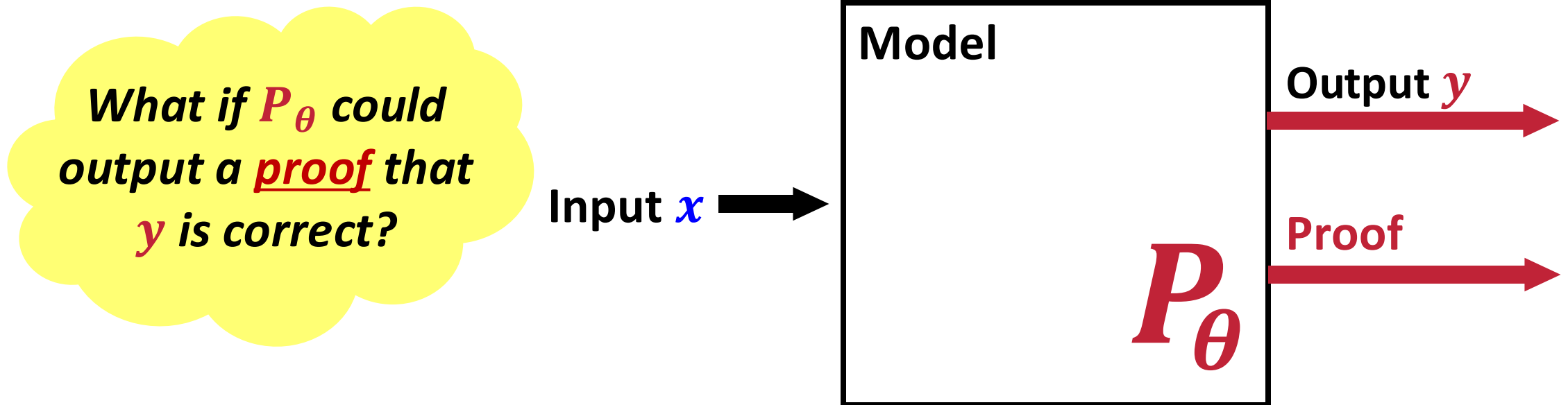
# Motivation: Worst-case guarantees for ML

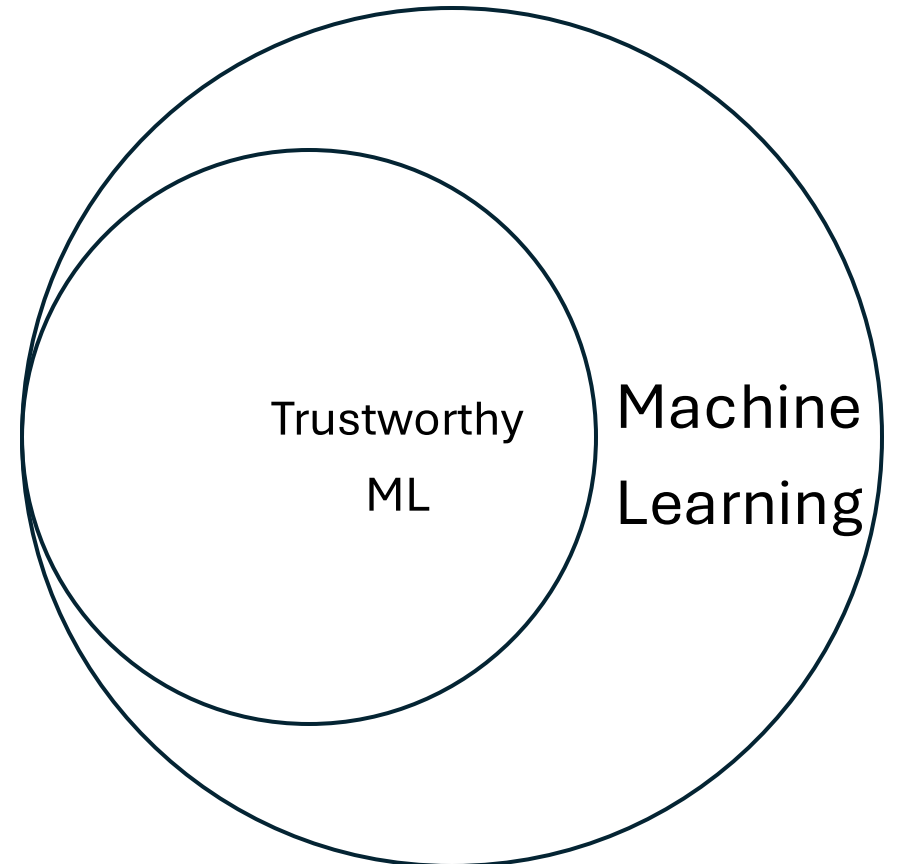
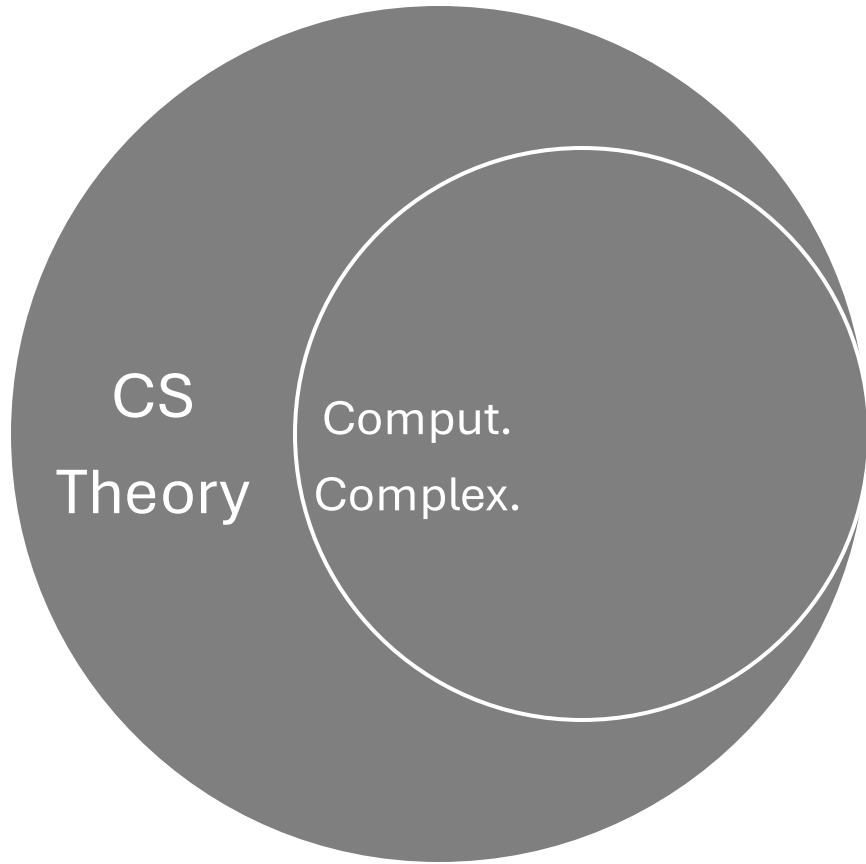
- Model  $P_\theta$  has 99% accuracy on benchmarks for a task  $T$
- You have an input  $x$ , for which the model generates  $y \sim P_\theta(x)$
- **Should you trust that  $y$  is correct?**

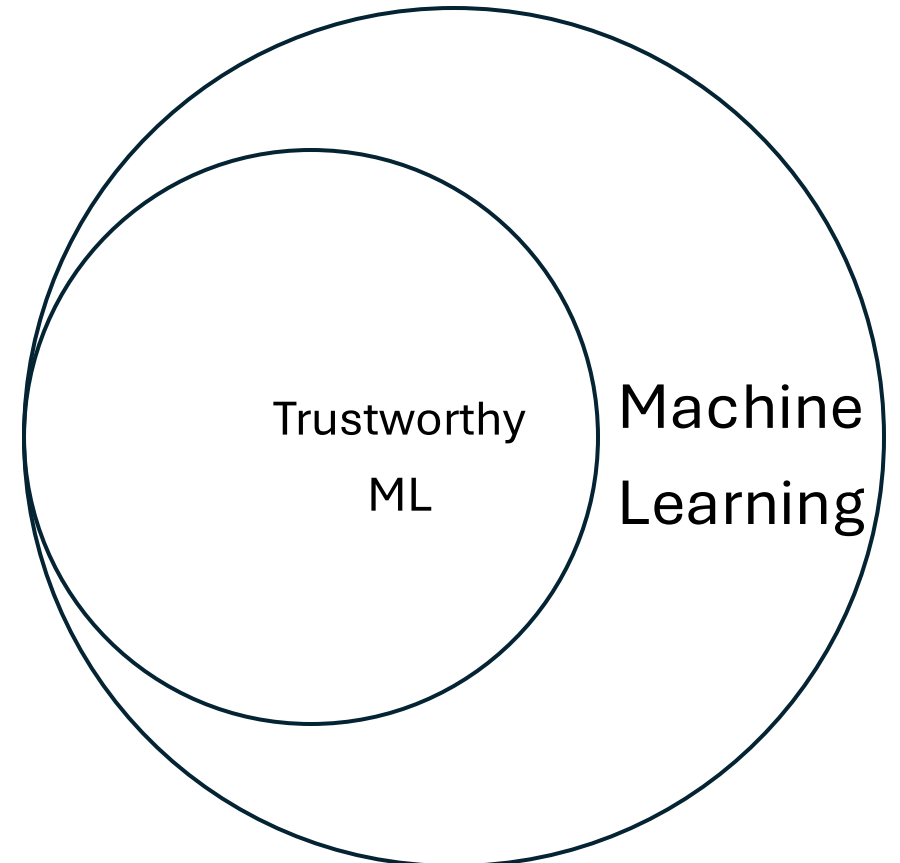
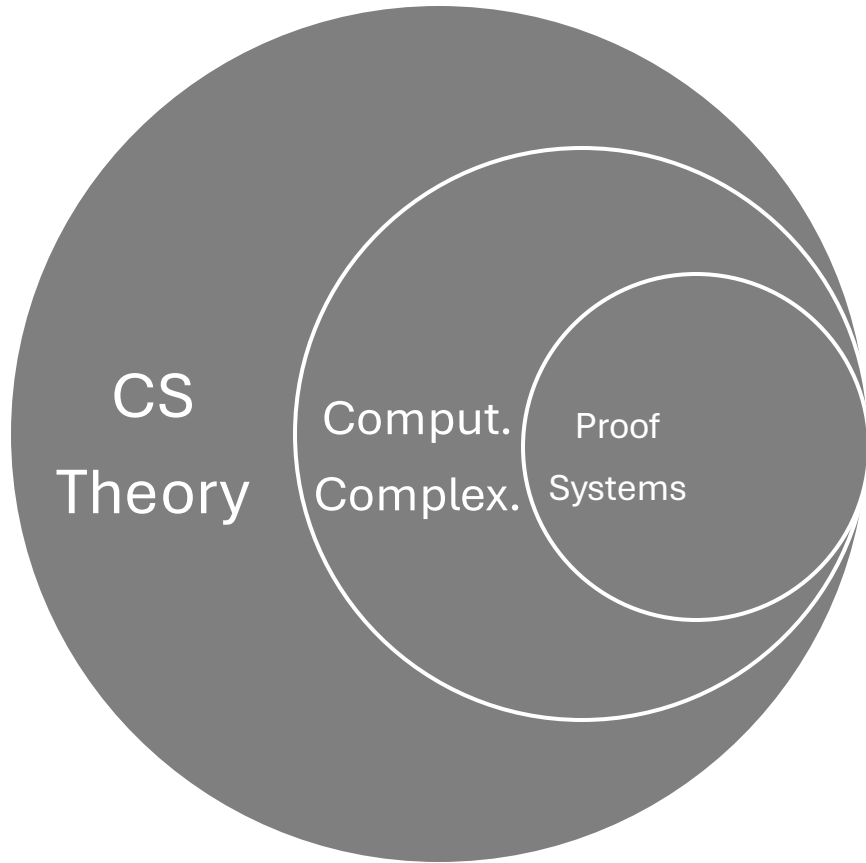


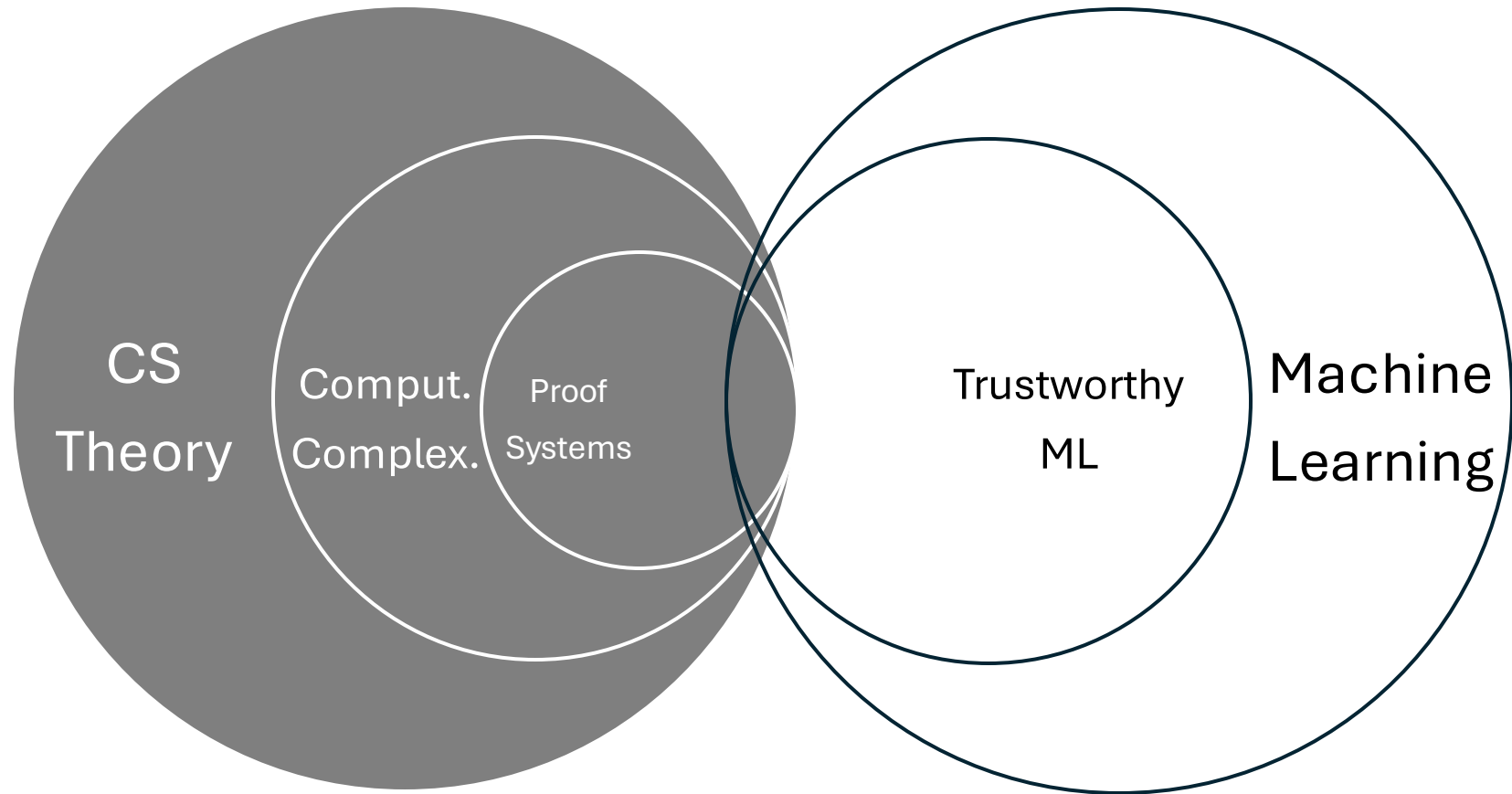
# Motivation: Worst-case guarantees for ML

- Model  $P_\theta$  has 99% accuracy on benchmarks for a task  $T$
- You have an input  $x$ , for which the model generates  $y \sim P_\theta(x)$
- **Should you trust that  $y$  is correct?**

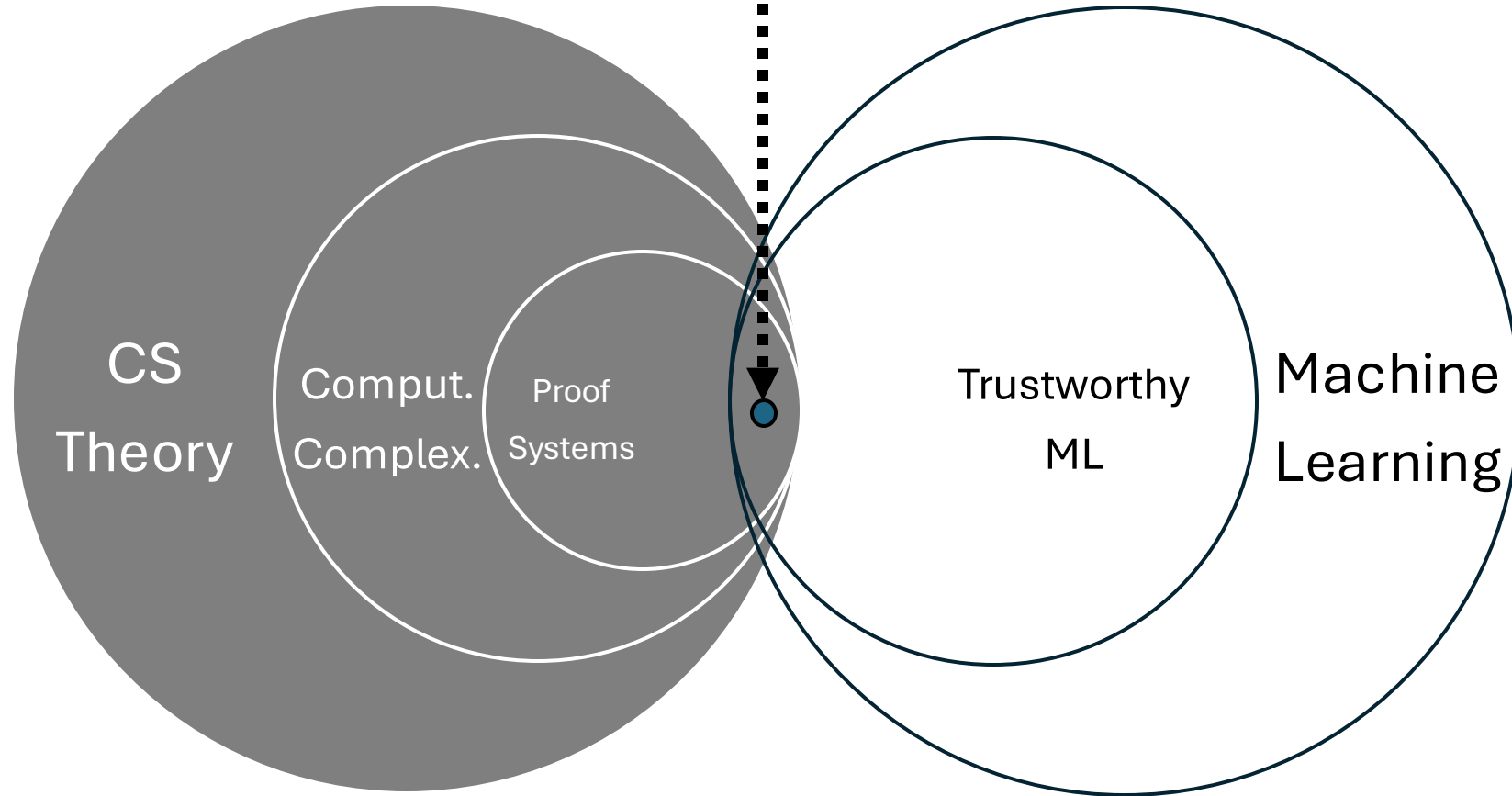








# Self-Proving models





# Outline



Self-Proving  
models

Related  
literature

Results

Future  
directions

Theory

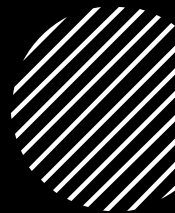
Experiments







# Outline



Self-Proving  
models



Related  
literature

Results

Future  
directions

Joint work with:



Noga  
Amit

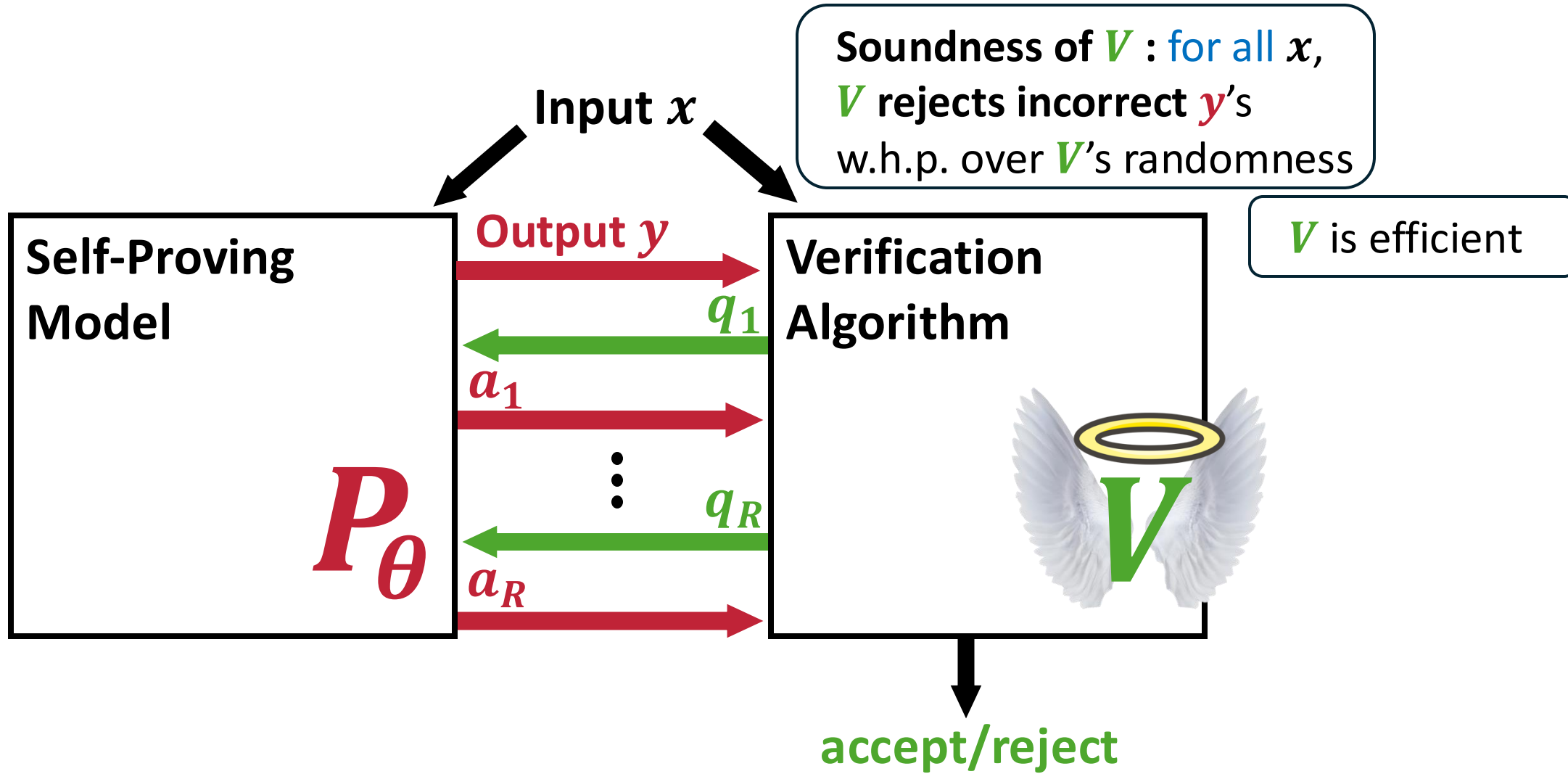


Shafi  
Goldwasser



Guy N.  
Rothblum

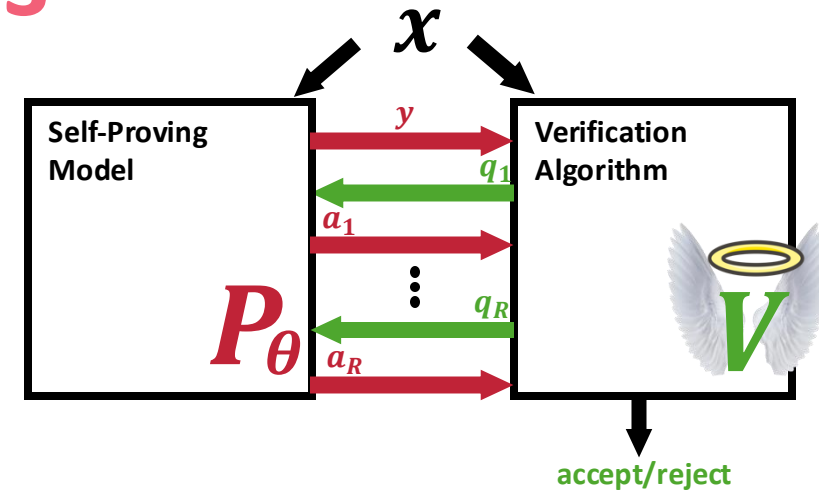
# Self-Proving Models



# Self-Proving Models

## Setting

- $\mu$  = distribution over inputs  $x$ .
- $V$  = an efficient and **sound** verification algorithm
- $P_\theta$  = sequence-to-sequence model with parameters  $\theta \in \mathbb{R}^d$



**Definition (Self-Provability):**  $P_\theta$  is  $\beta$ -self-proving w.r.t.  $V, \mu$  if

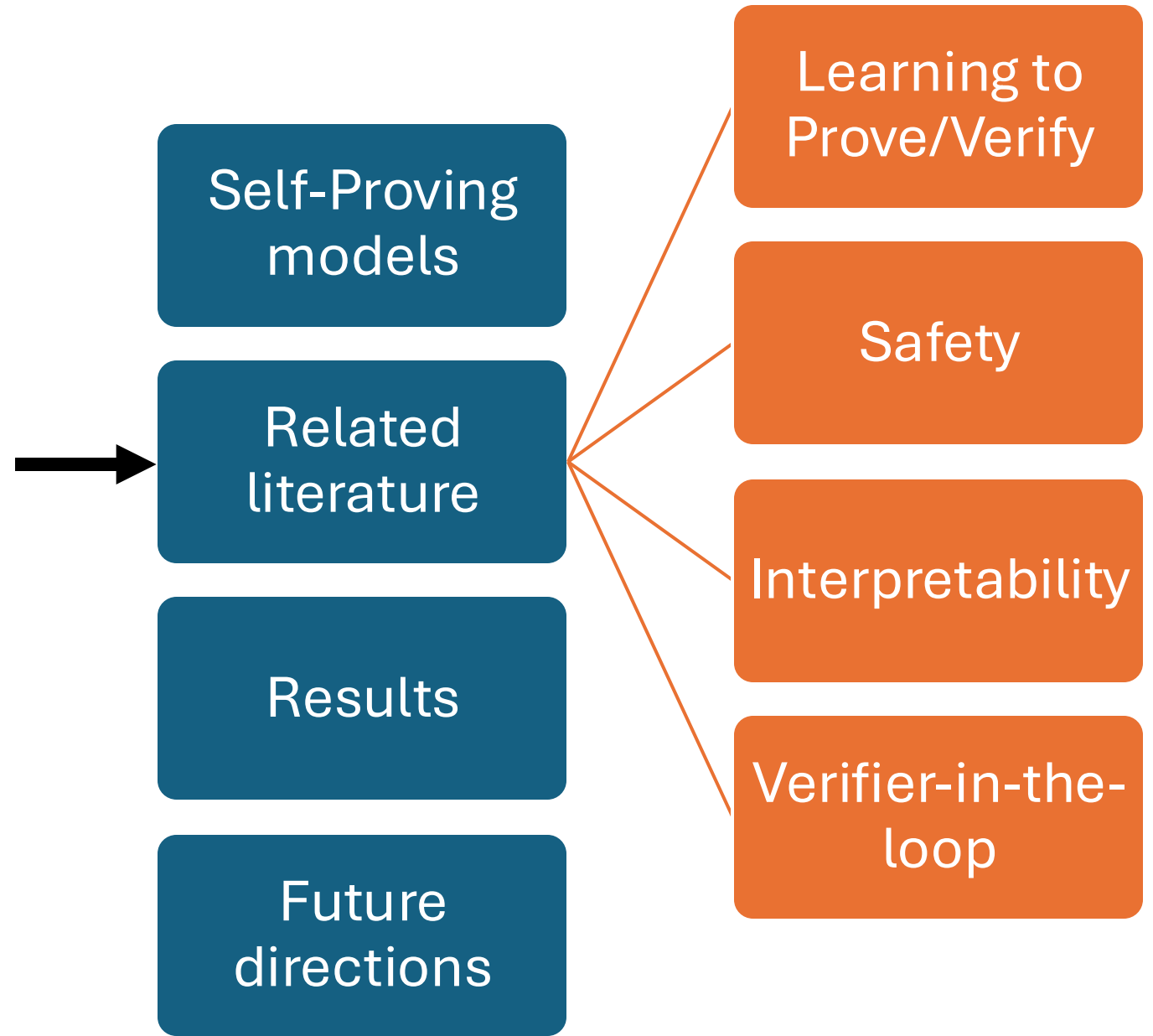
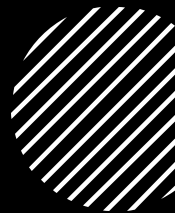
$$\Pr[P_\theta \text{ convinces } V \text{ to accept } y] \geq \beta \quad \text{for } x \sim \mu, y \sim P_\theta(x)$$

Self-Provability + Soundness  $\Rightarrow$  Standard accuracy (w.r.t  $\mu$ ):

$$\text{For } x \sim \mu, \Pr[P_\theta \text{ outputs correct } y] \geq \alpha.$$

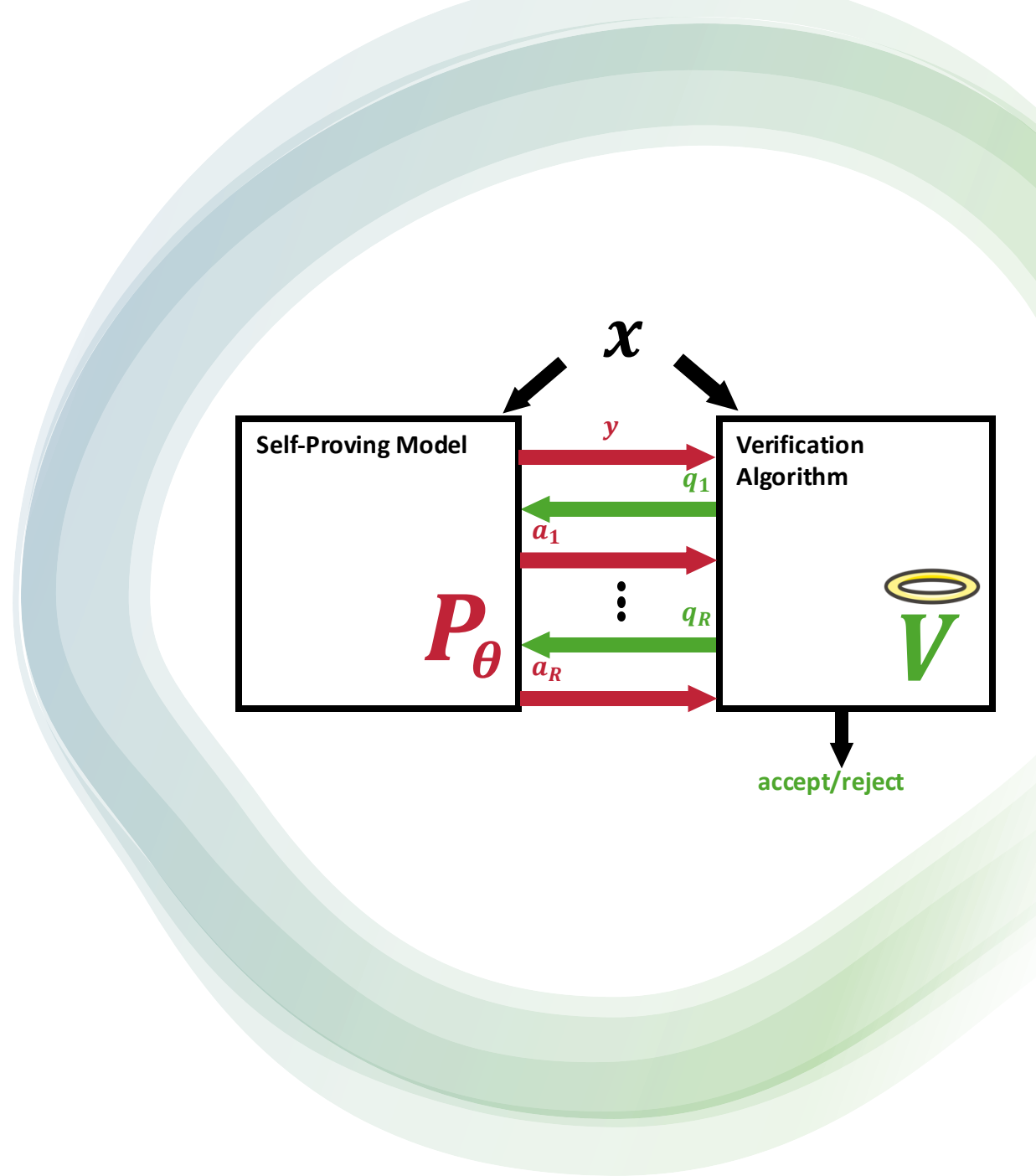


# Outline



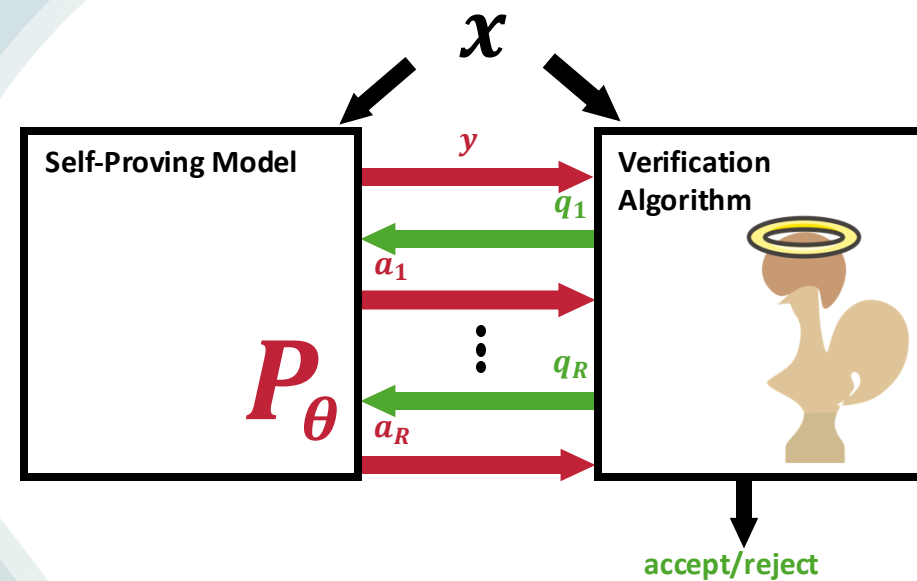
# Related literature

- Learning to Prove...



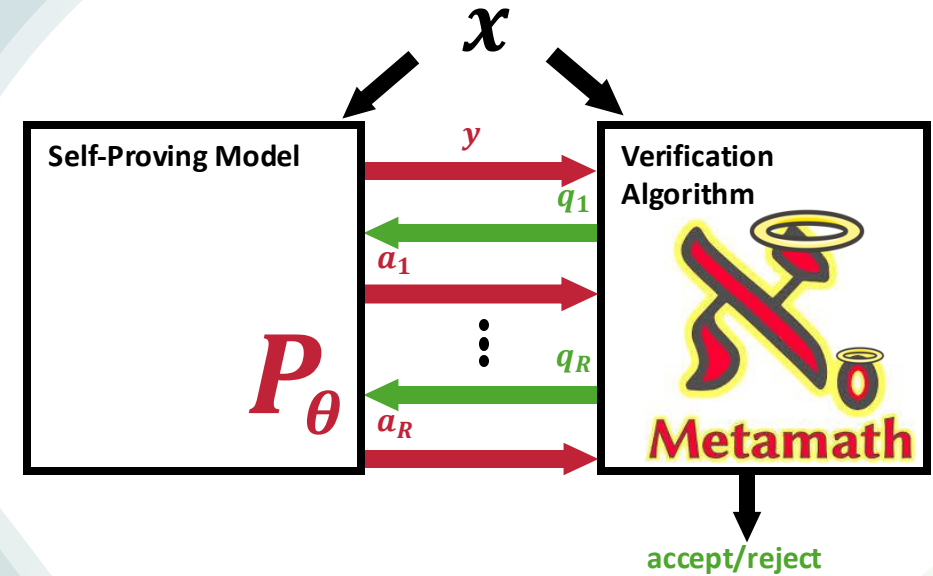
# Related literature

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



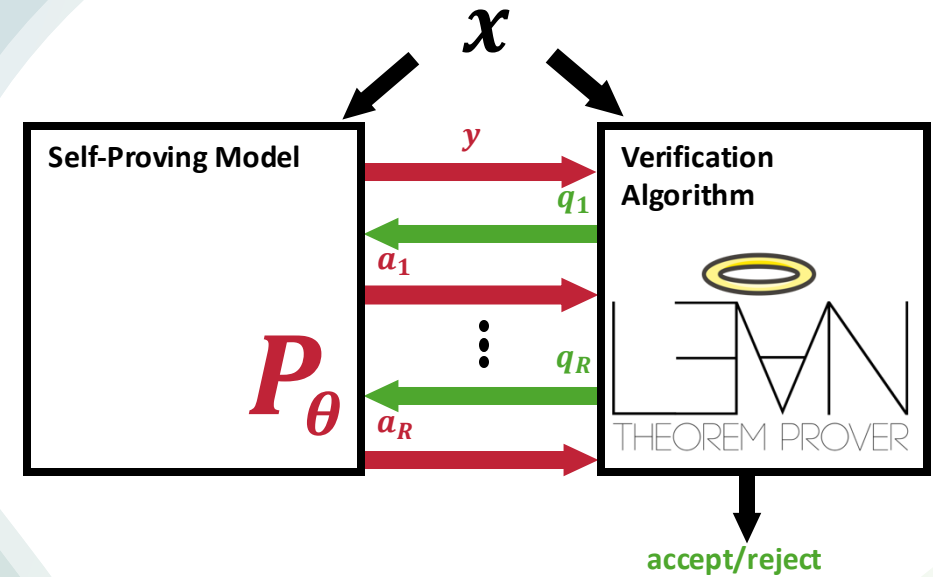
# Related literature

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



# Related literature

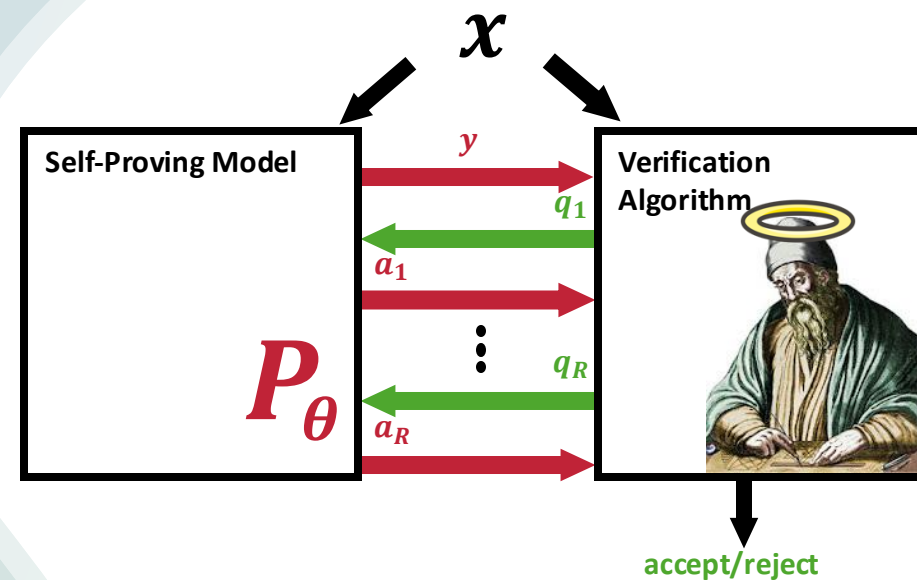
- Learning to Prove...
  - ... in Coq (Gransden et al., 2015)
  - ... in Metamath (Polu and Sutskever, 2020)
  - ... in Lean (Yang et al., 2023)





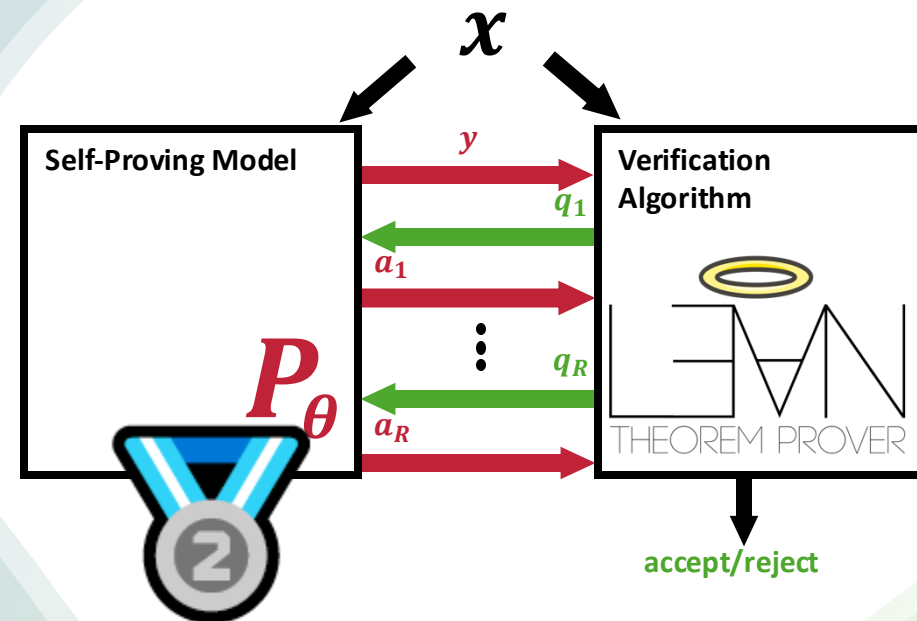
# Related literature

- 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)**



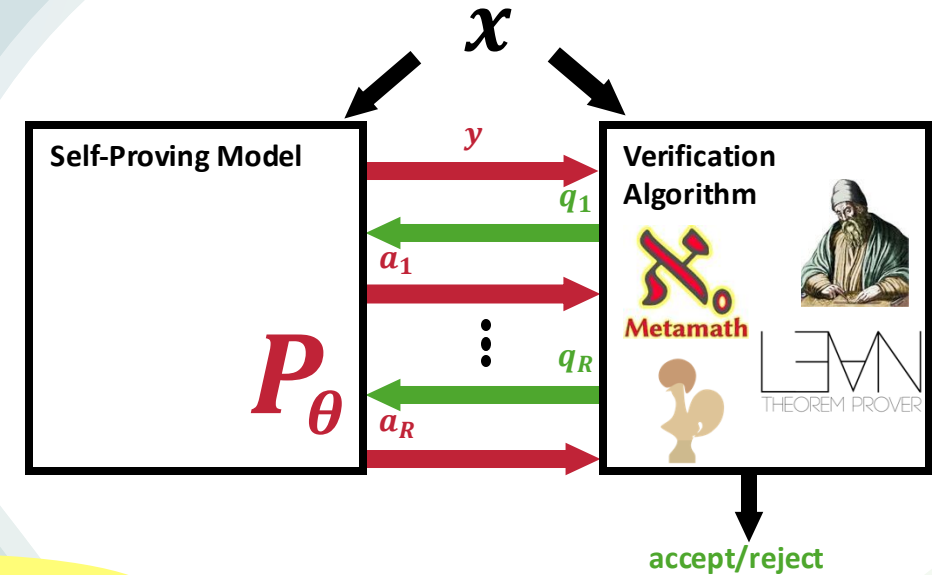
# Related literature

- 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)



# Related literature

- 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)



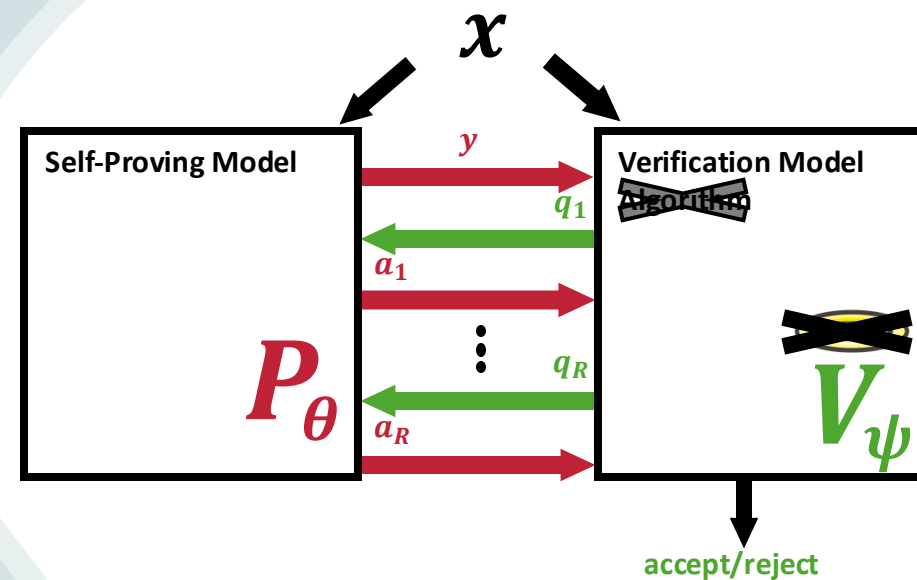
**Self-Proving models: prove correctness  
via an **Interactive Proof system****

**IP:** Weak Verifier vs. powerful *yet untrusted* Prover (Goldwasser et al., 1985)

**IP=PSPACE** (Shamir, 1990)

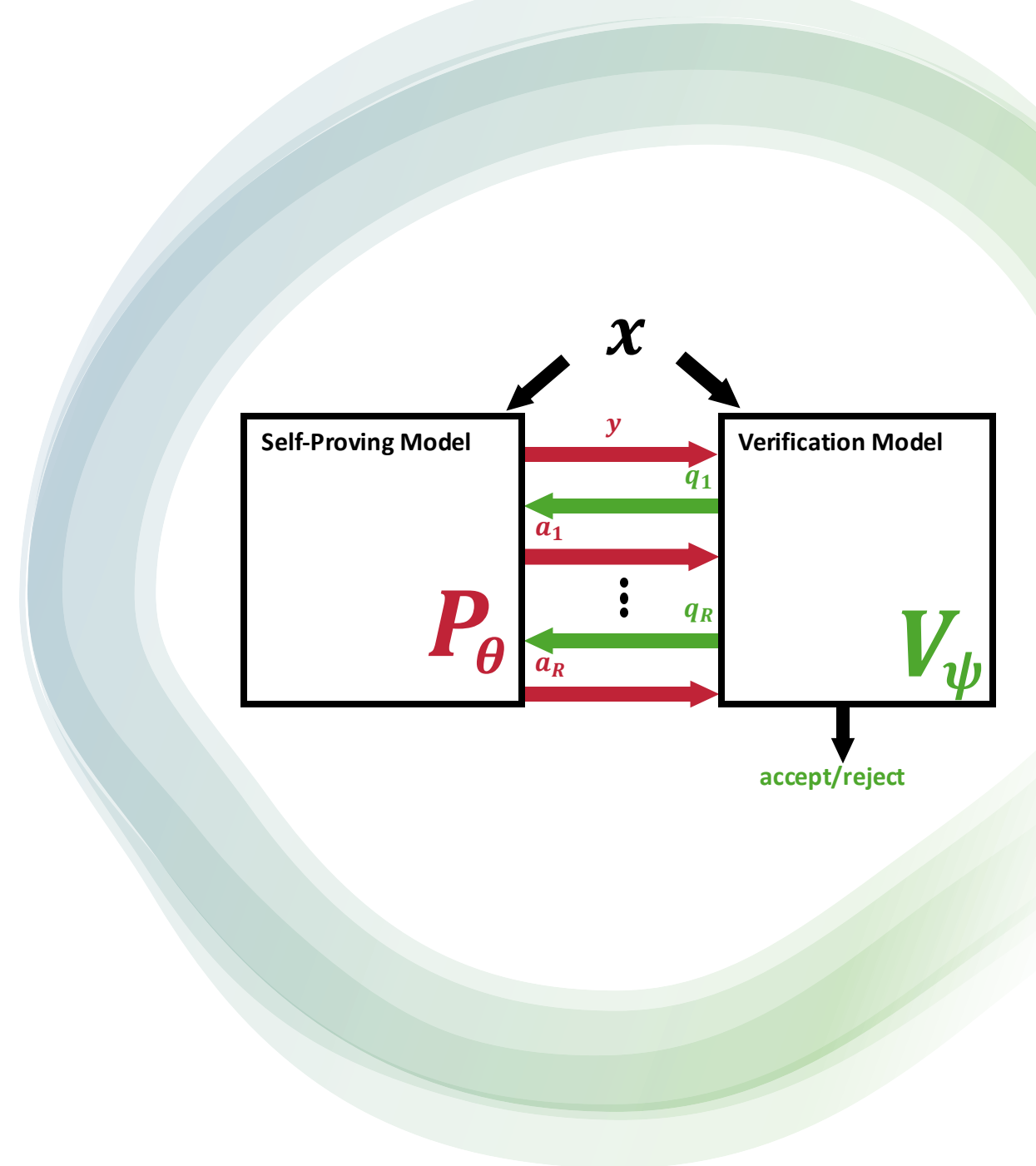
# Related literature

- 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



# Related literature

- 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)**



# Related literature

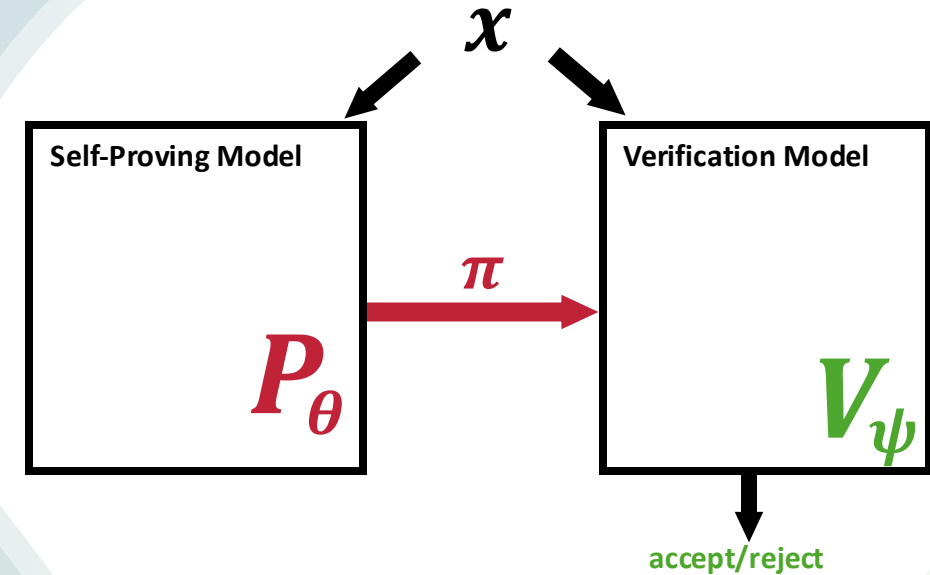
- Learning to Prove...
  - ... in Coq (Gransden et al., 2015)
  - ... in Metamath (Polu and Sutskever, 2020)
  - ... in Lean (Yang et al., 2023)
  - ... in synthetic geometry (Trinh et al., 2024)
  - ... in Lean 🏆 (DeepMind, 2024)
- Learning to Verify
  - **Prover Verifier Games (Anil et al., 2021)**

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

Consider a strategy-finding game with two players:  $P_\theta$  and  $V_\psi$ .

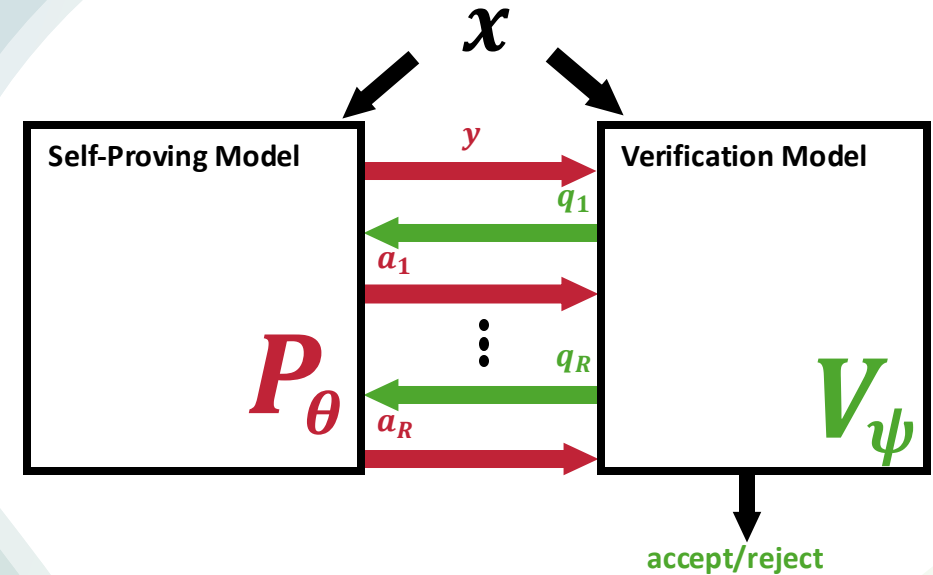
Does an equilibrium 🏆 → Completeness and Soundness? Does C&S → 🏆 ?

Order	🏆 → C&S	C&S → 🏆
$P_\theta$ first	No	No
$V_\psi$ first	<b>Yes!</b>	<b>Yes!</b>
$\{P_\theta, V_\psi\}$ simul.	No	<b>Yes!</b>



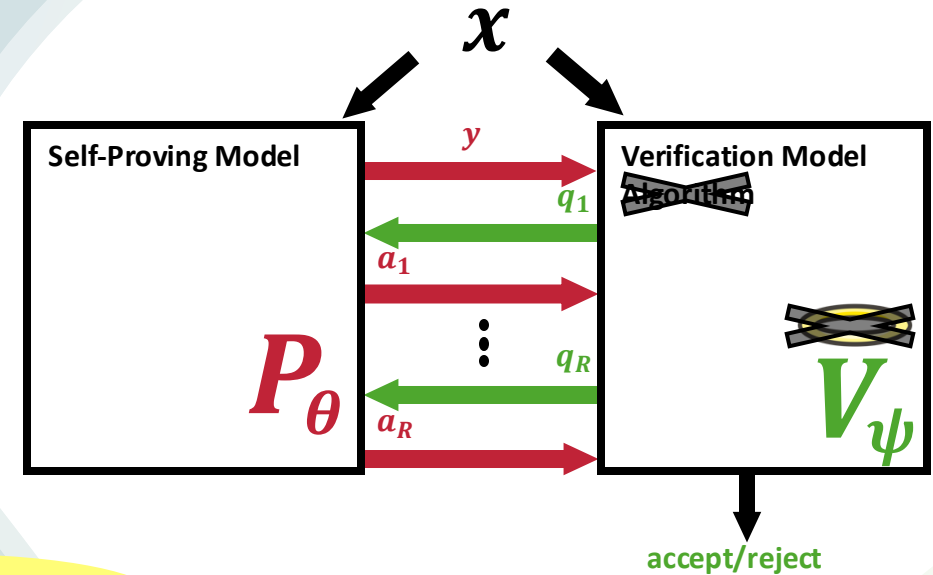
# Related literature

- 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)**



# Related literature

- 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)



**Self-Proving models: prove correctness  
to **Sound Verifiers**---a formal guarantee.**



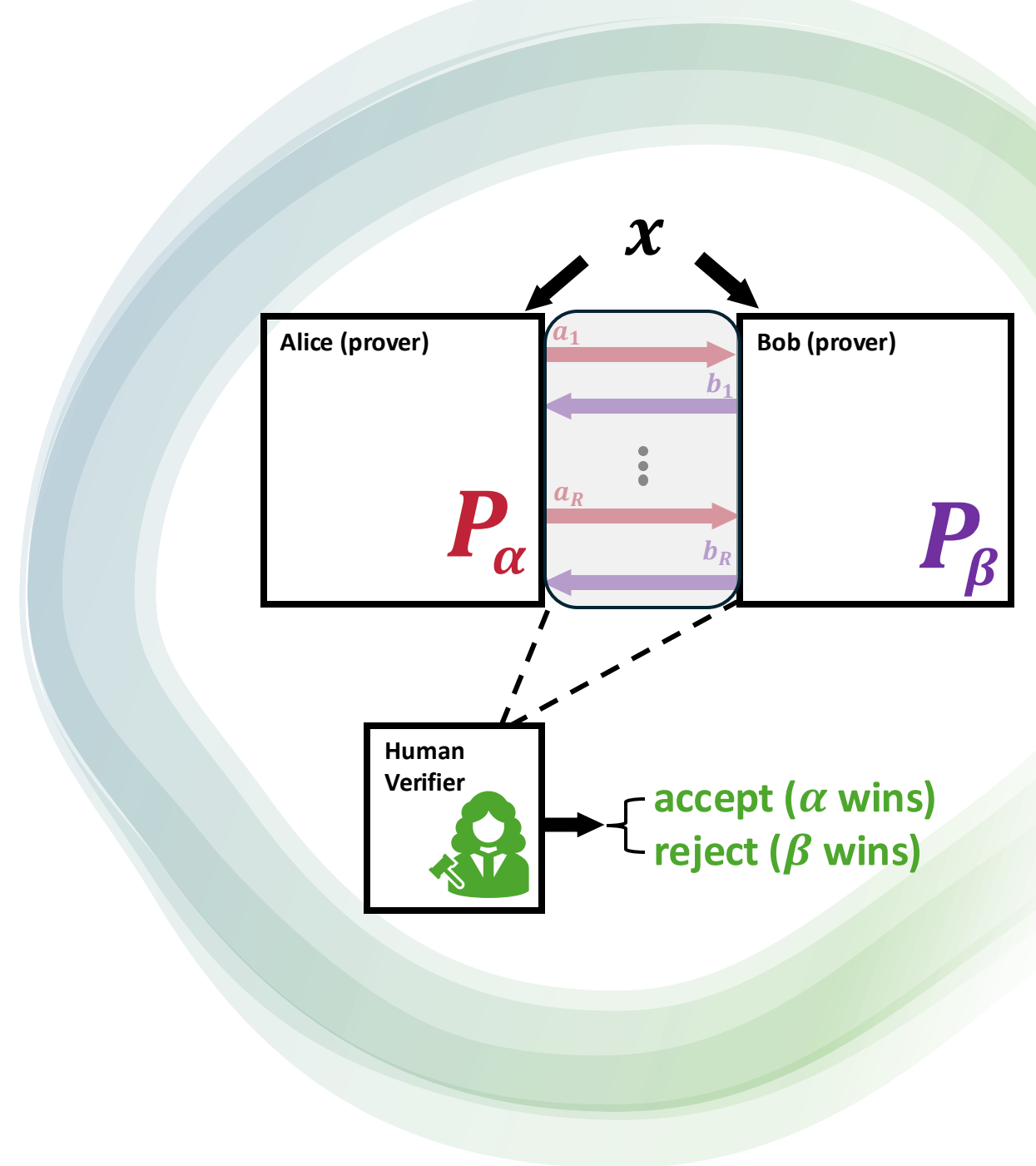
# Related literature

- 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)**



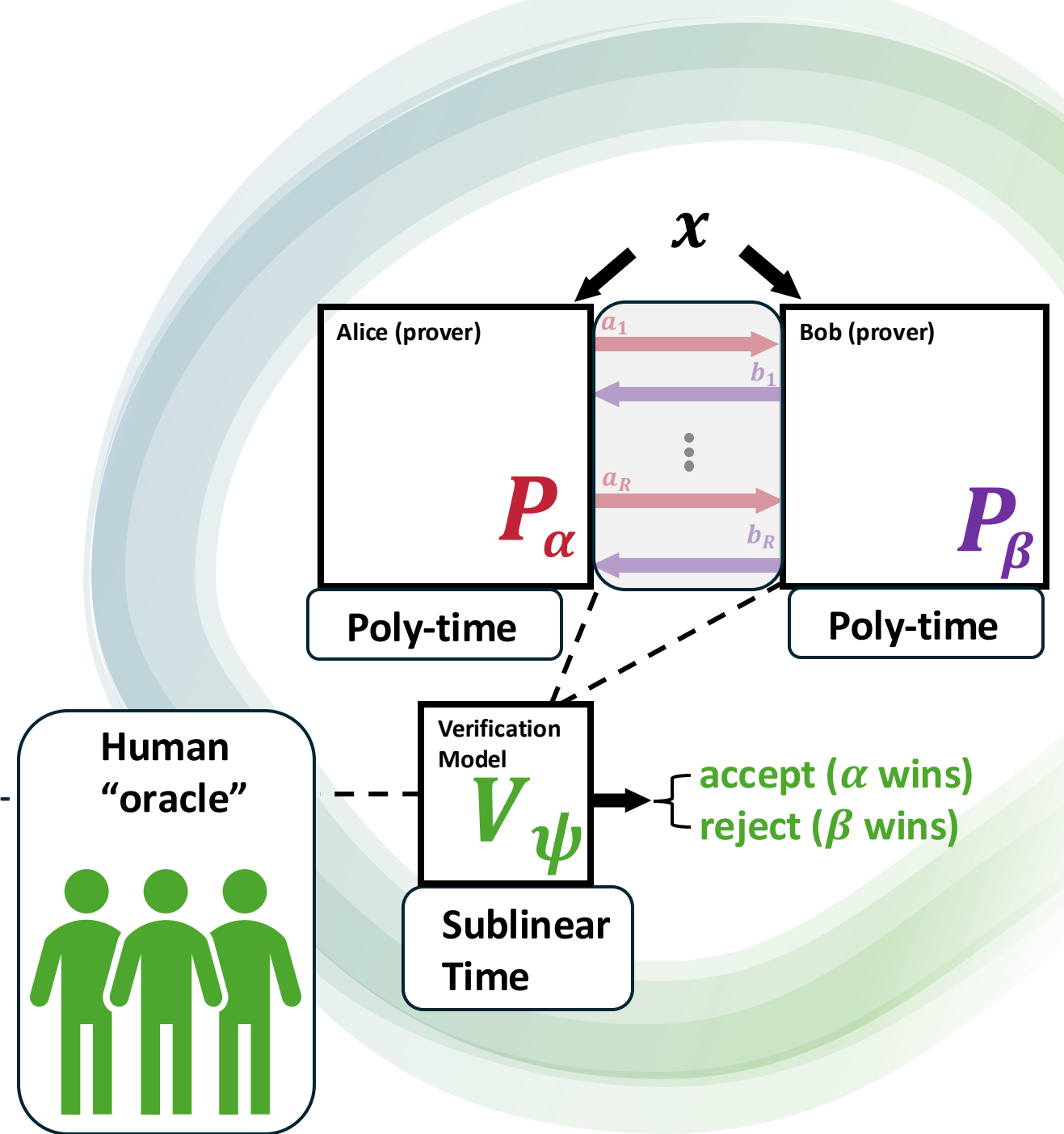
# Related literature

- 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)**



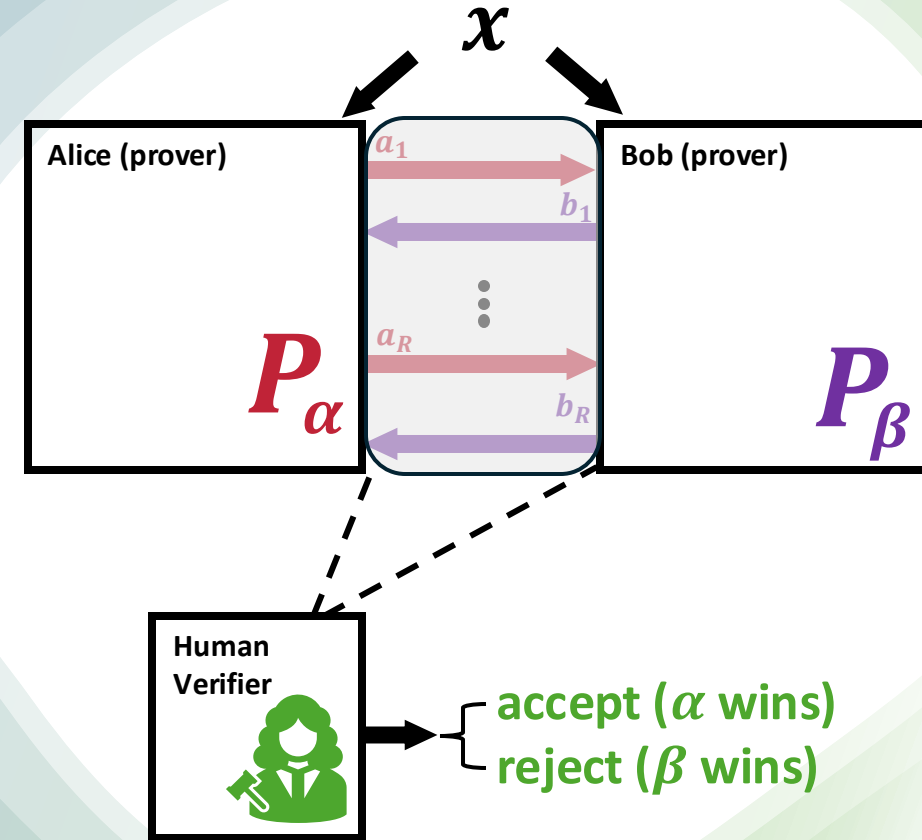
# Related literature

- 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)**



# Related literature

- 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)



**Self-Proving models: single prover, sound verifier**

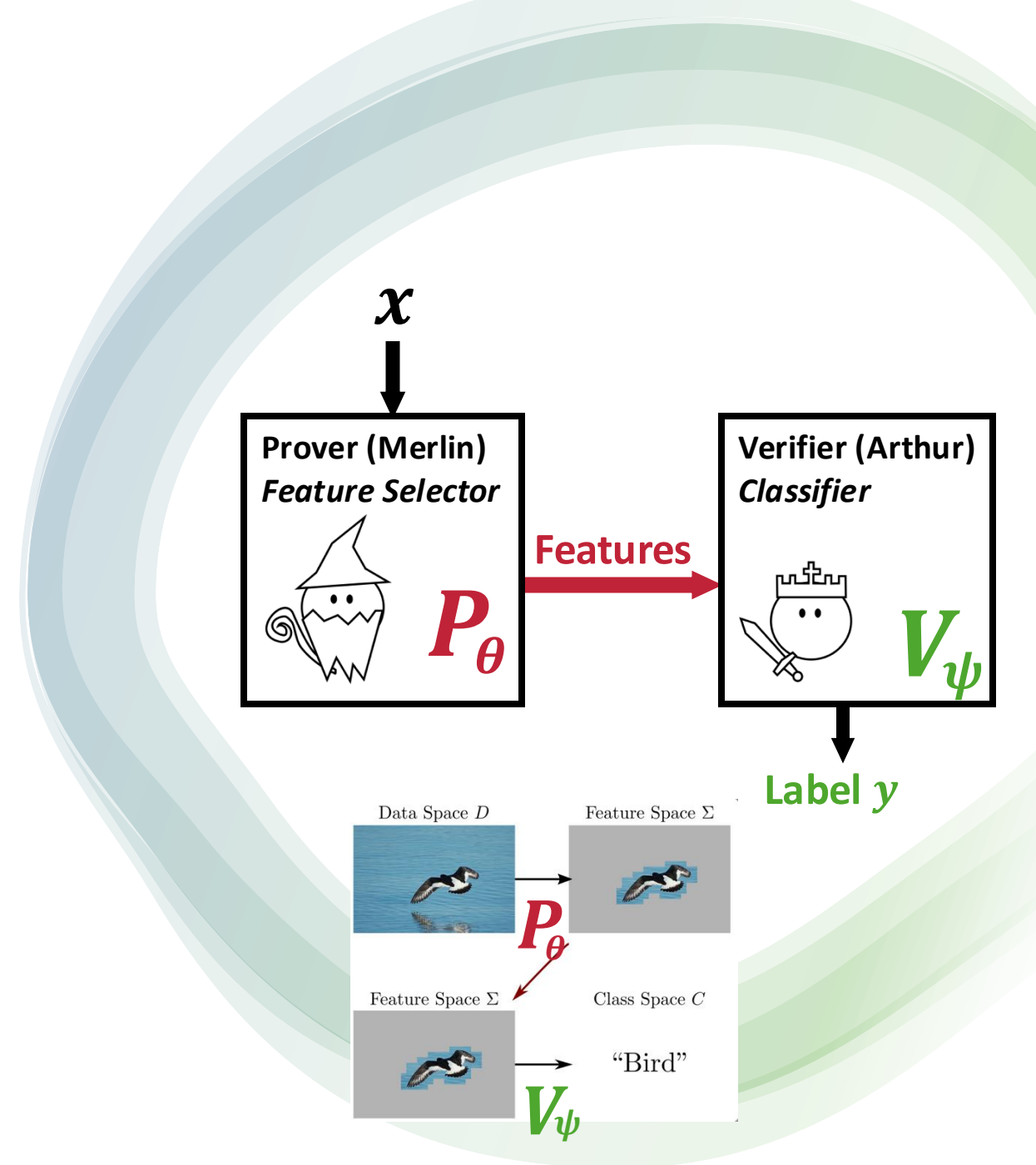
# Related literature

- 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**



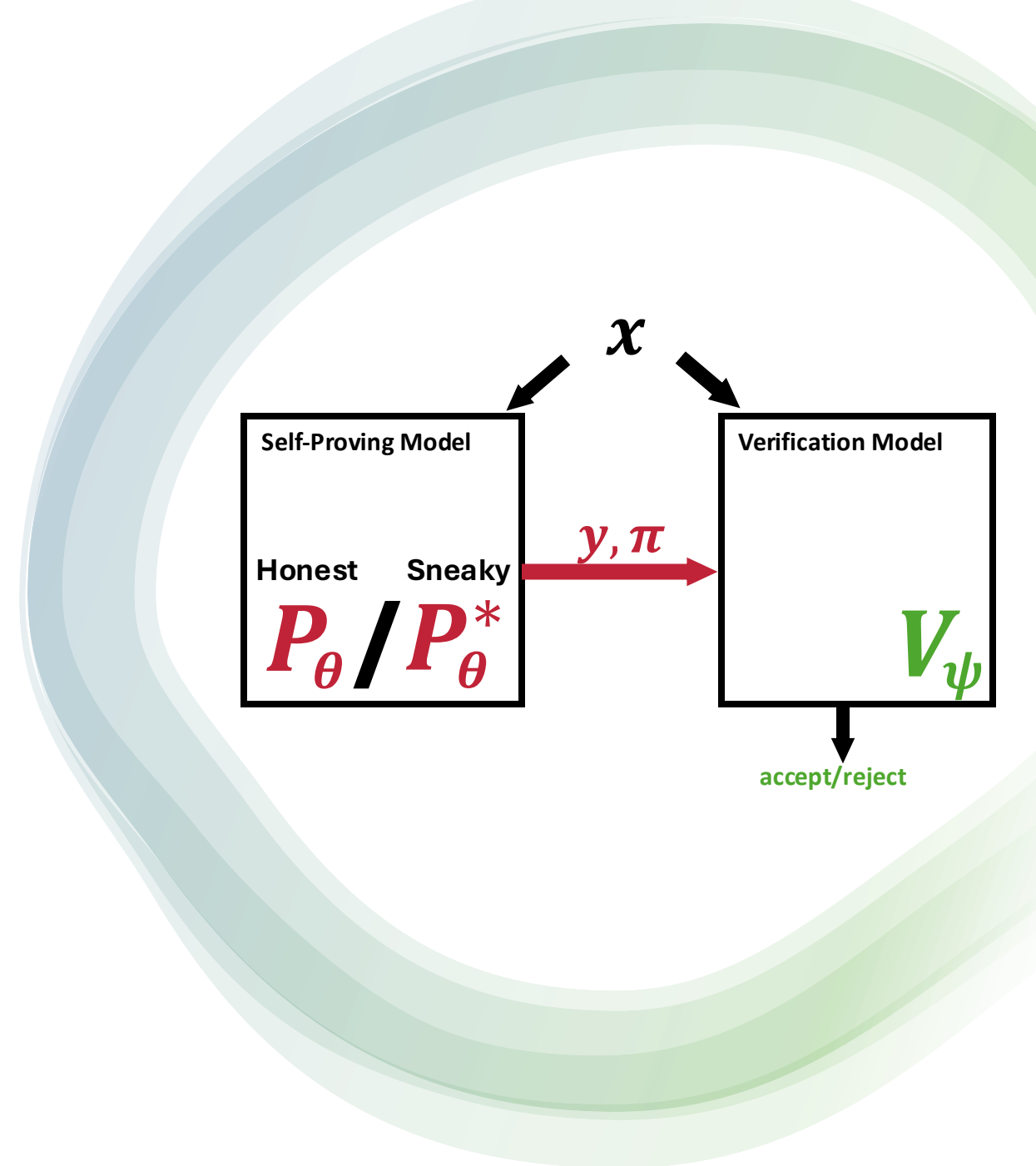
# Related literature

- 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)**



# Related literature

- 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)





# Prover-Verifier Games Improve Legibility of LLM Outputs (Hendrick\* and Chen\* et al., 2024)

Idea: Introduce a small FT Prover LLM to

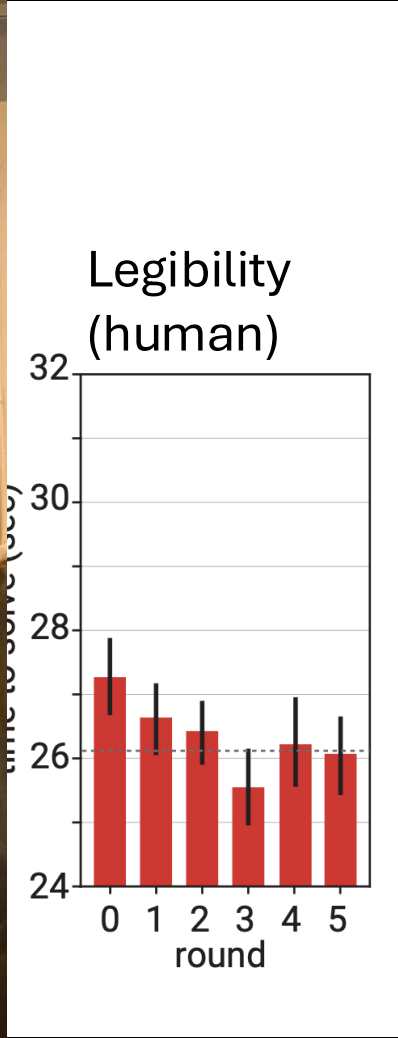
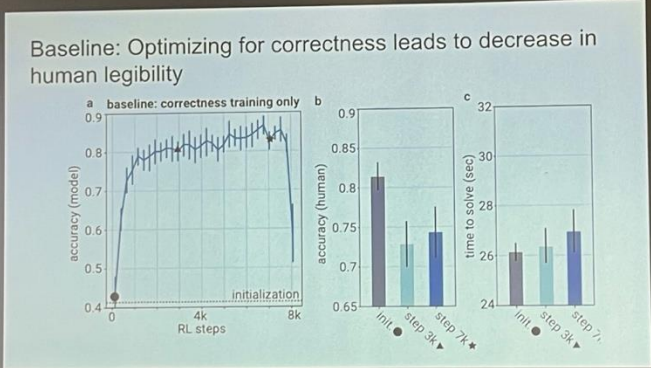
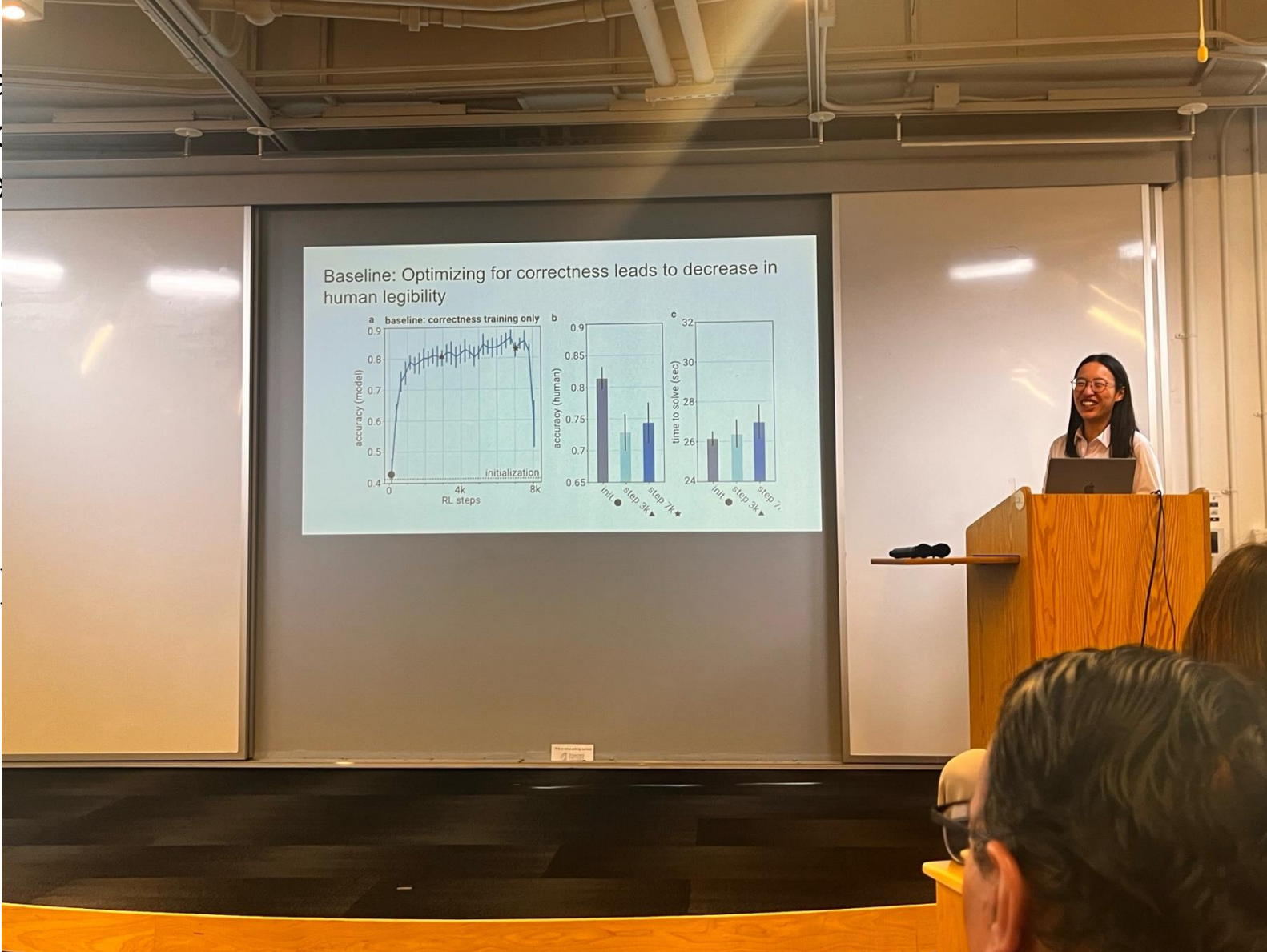
- Does this increase LLM's answers?

Joint training (simplified)

$$\psi \xleftarrow{\text{argmin}} L(\psi | \theta)$$

$$\ell(\psi | P)$$

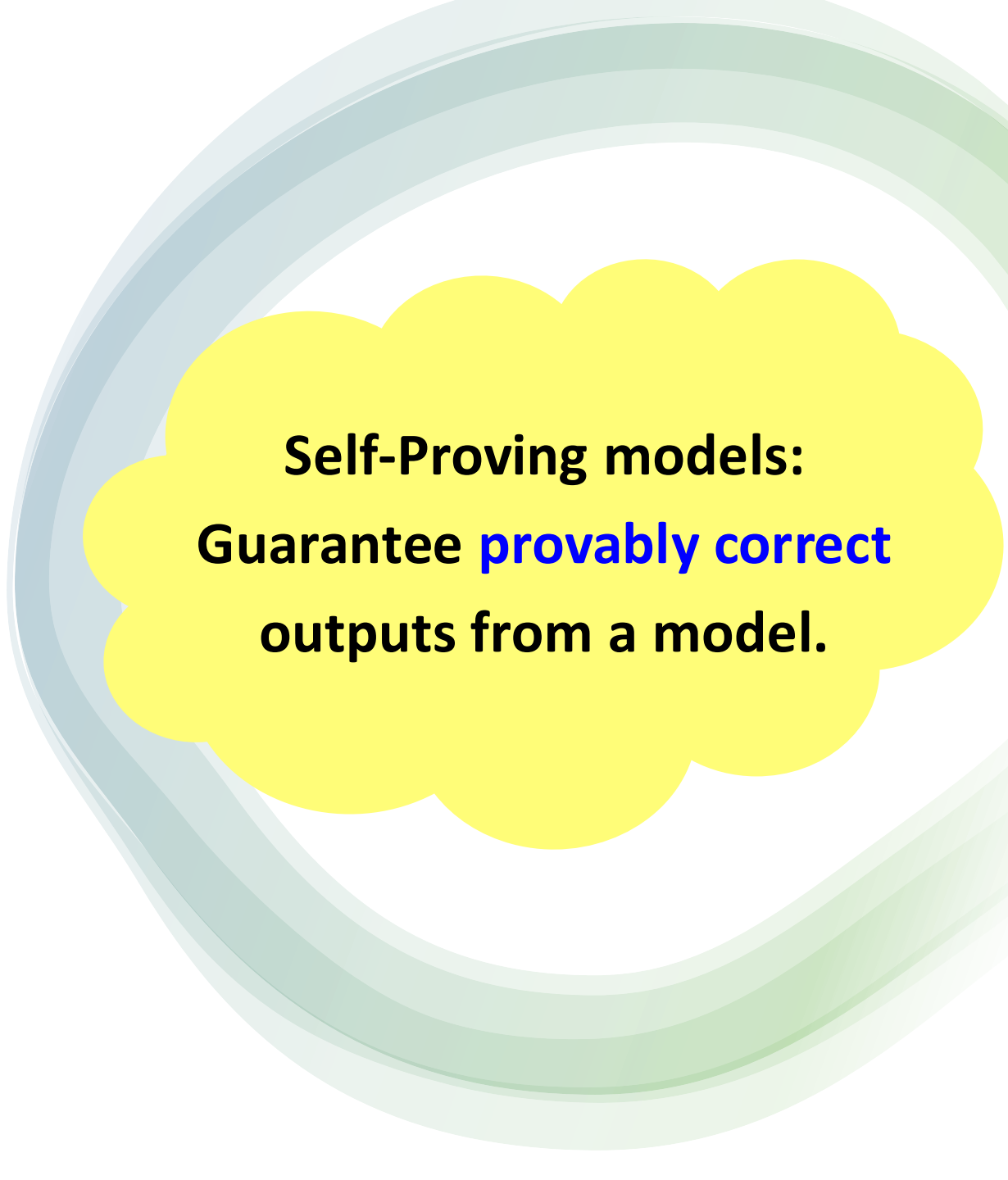
$$\theta \xleftarrow{\text{argmin}} \mathbb{E}_x [\langle V_\psi, P_\theta \rangle(x)]$$





# Related literature

- 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)

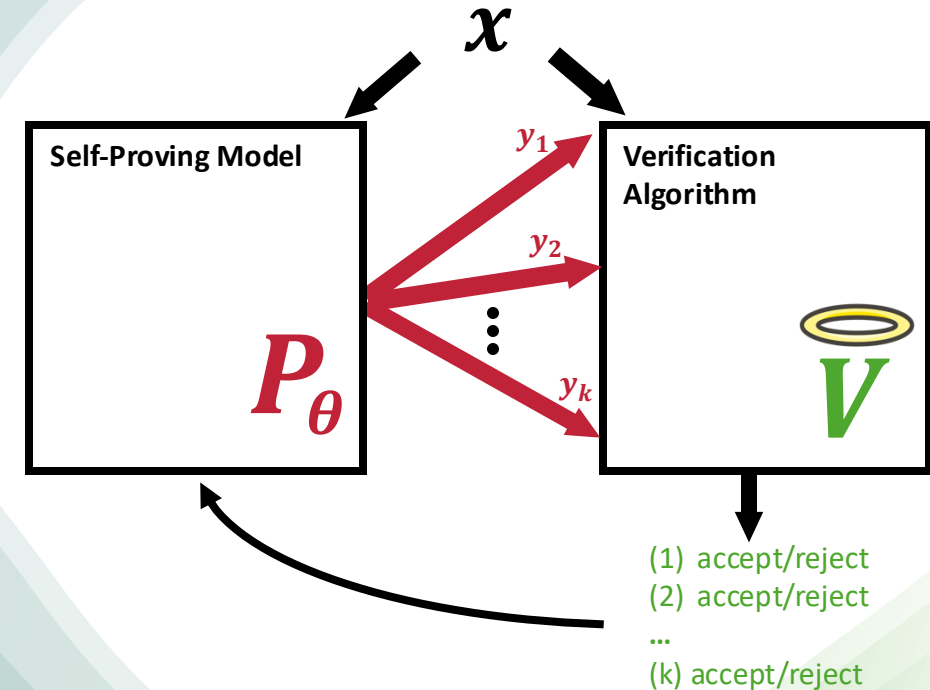


**Self-Proving models:  
Guarantee provably correct  
outputs from a model.**

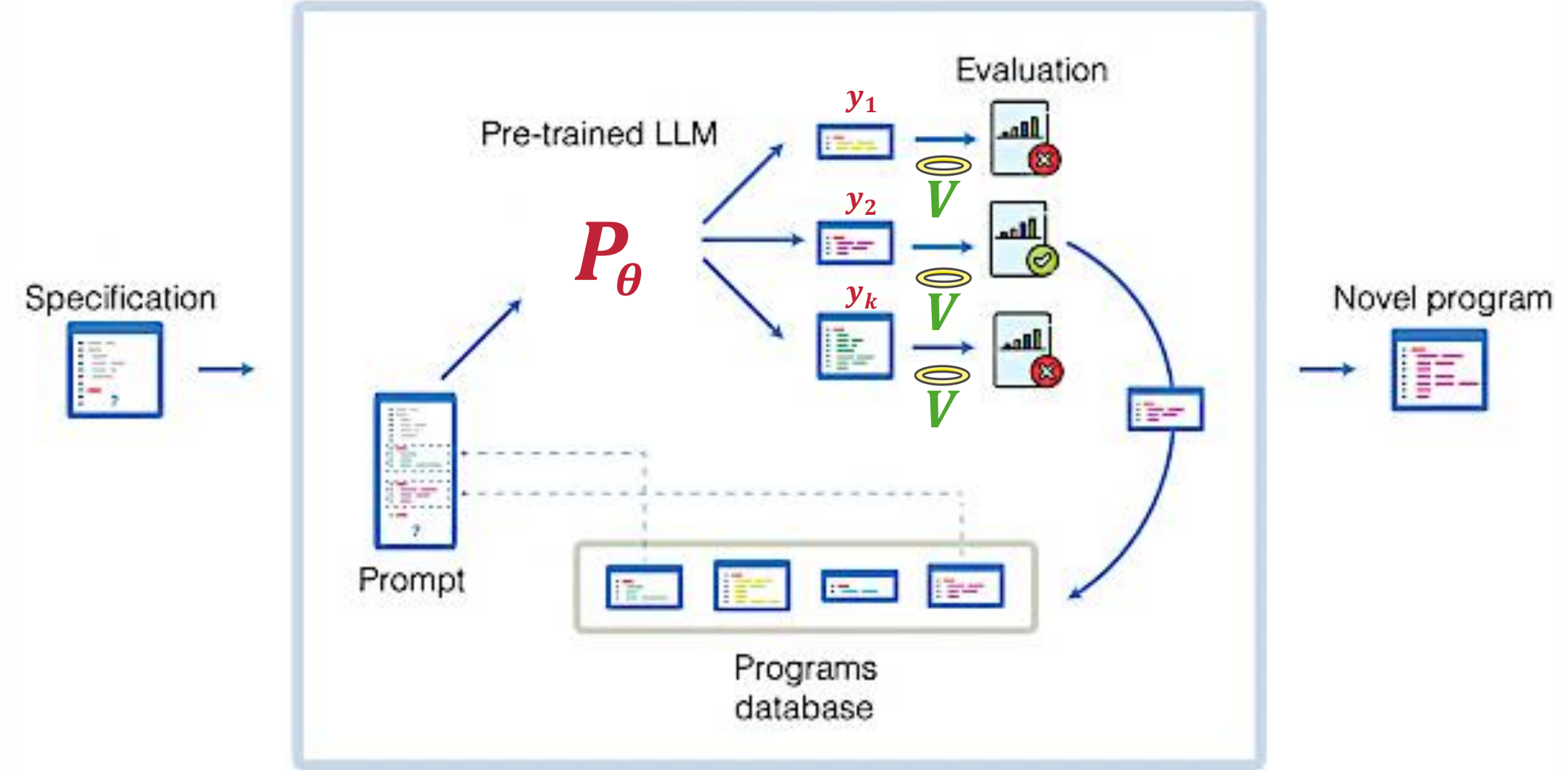
# Related literature

- 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)
- Verifier in-the-loop
  - **FunSearch (Romera-Paredes et al., 2024)**

**Goal: find optimal  $y$  for  $x$**

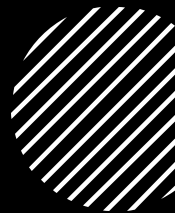


# Mathematical discoveries from program search with LLMs (Romera-Paredes et al., 2024)





# Outline



Self-Proving  
models

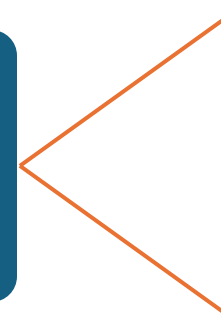
Related  
literature

Results

Future  
directions

Theory

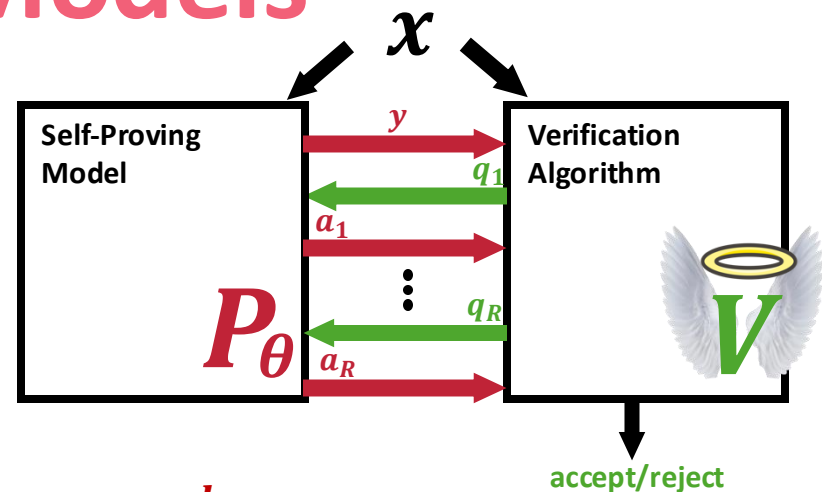
Experiments



# Reminder: Self-Proving Models

## Setting

- $\mu$  = distribution over inputs  $x$ .
- $V$  = an efficient and **sound** verification algorithm
- $P_\theta$  = sequence-to-sequence autoregressive model with parameters  $\theta \in \mathbb{R}^d$



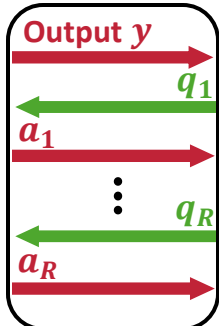
**Definition:** (Worst-case) **Self-Provability:**  $P_\theta$  is  $\beta$ -self-proving w.r.t.  $V, \mu$  if

$$\Pr[P_\theta \text{ convinces } V \text{ to accept } y] \geq \beta \text{ for } x \sim \mu, y \sim P_\theta(x)$$

**Goal:** Learn  $\theta$  that maximizes  $\Pr[P_\theta \text{ convinces } V \text{ to accept } y]$

# Learning Self-Proving Models

Fix input distribution  $\mu$ , verifier  $V$

**Idea:** Transcripts  $\pi =$   are sequences, learn to generate  $x \rightarrow \pi$

**Theorem:** Under convexity and Lipschitzness assumptions, given access to accepted transcripts, Transcript Learning outputs a

$(1 - \epsilon)$ -Self-Proving model after  $N \geq 4 \left( C \cdot B_1 \cdot B_2 \cdot \frac{1}{\epsilon} \right)^2$  iterations

Total message  
length of  $P_\theta$ 's

Bound on  
 $\|\nabla \log P_\theta\|$

Bound on  $\|\theta\|$

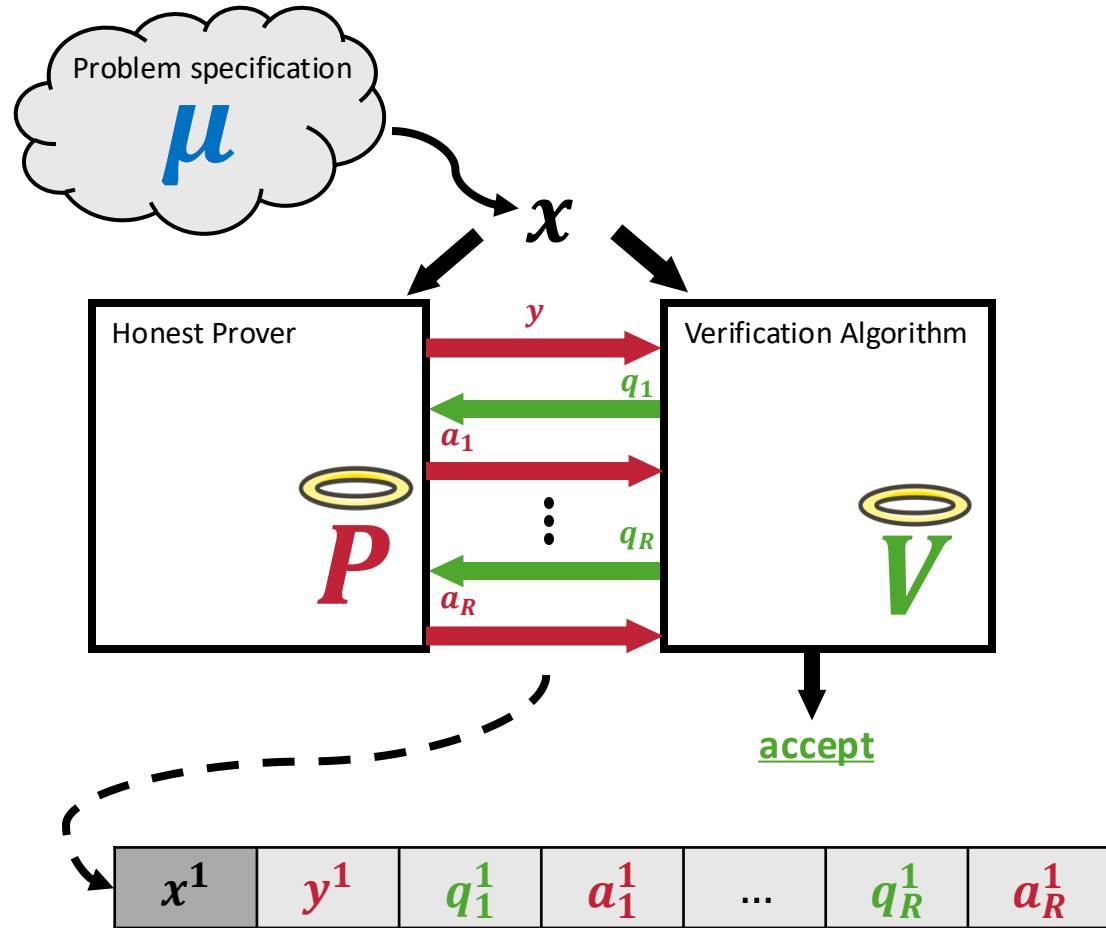
- Proof by reduction to SGD convergence bounds

# Transcript Learning

**Step 1: Collect accepted transcripts**

# Transcript Learning

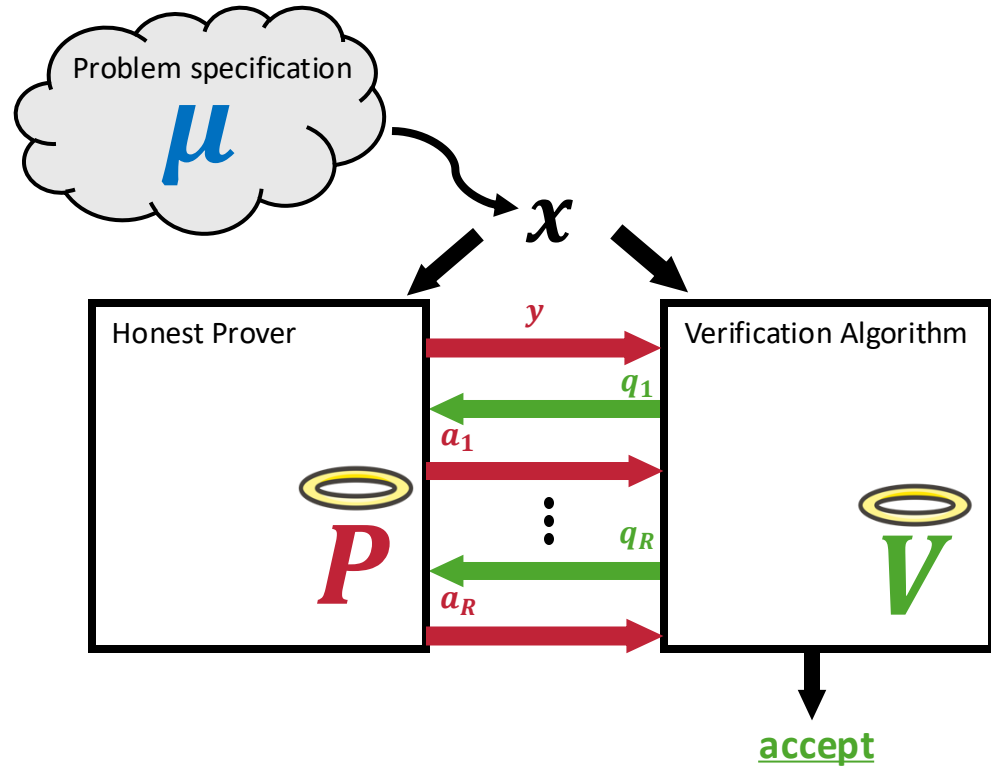
## Step 1: Collect accepted transcripts





# Transcript Learning

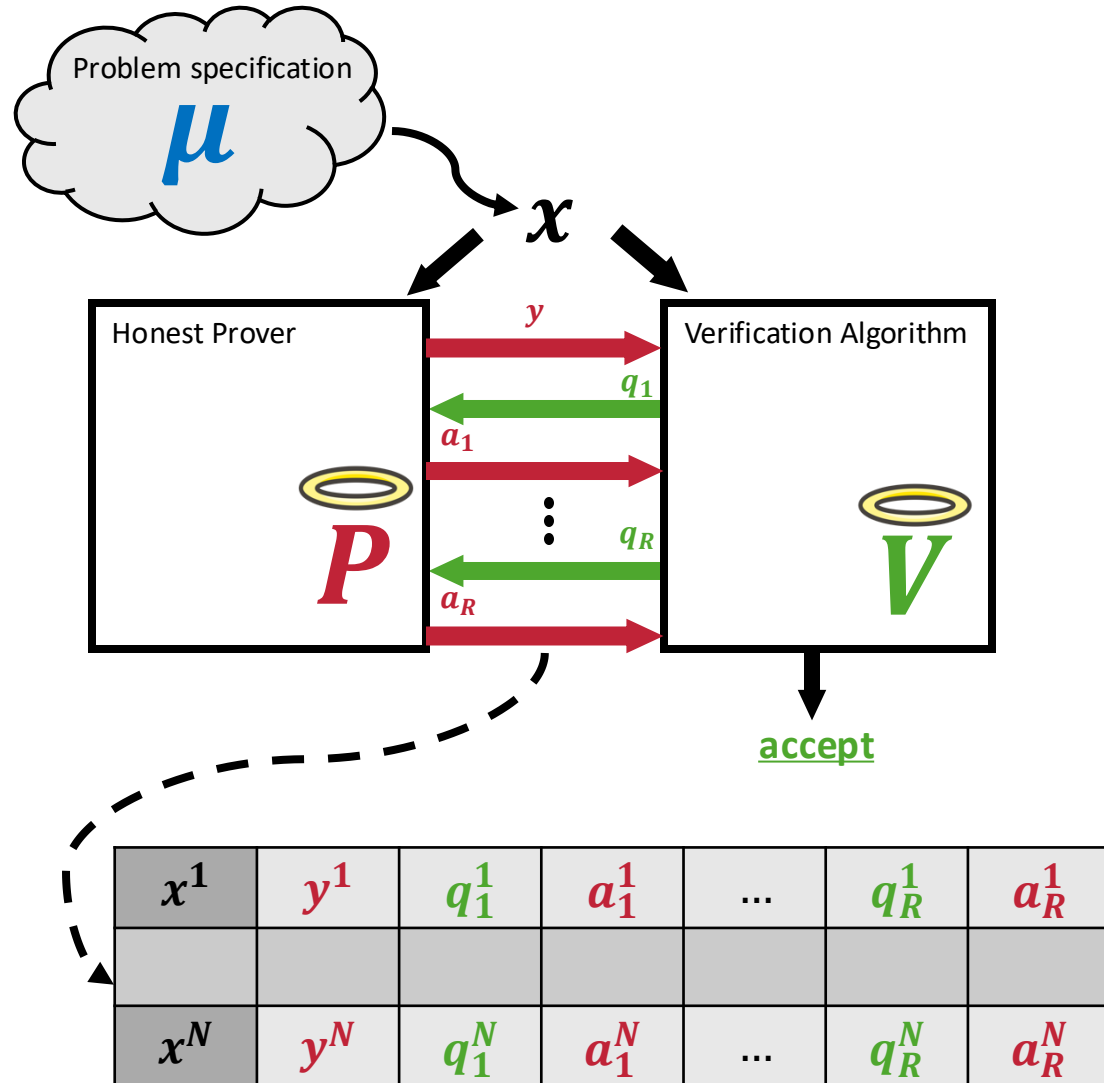
## Step 1: Collect accepted transcripts



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

# Transcript Learning

## Step 1: Collect accepted transcripts



# Transcript Learning

Step 1: Collect accepted transcripts

Step 2: Transcript Cloning

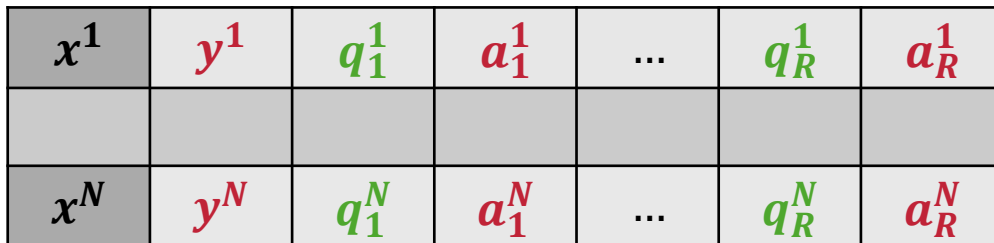
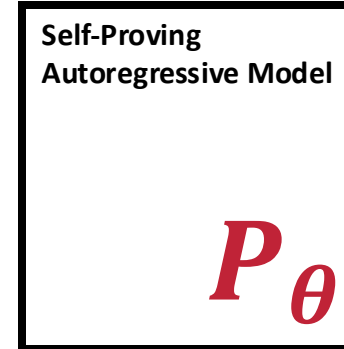
$x^1$	$y^1$	$q_1^1$	$a_1^1$	...	$q_R^1$	$a_R^1$
$x^N$	$y^N$	$q_1^N$	$a_1^N$	...	$q_R^N$	$a_R^N$

# Transcript Learning

Step 1: Collect accepted transcripts

Step 2: Transcript Cloning

For each transcript:

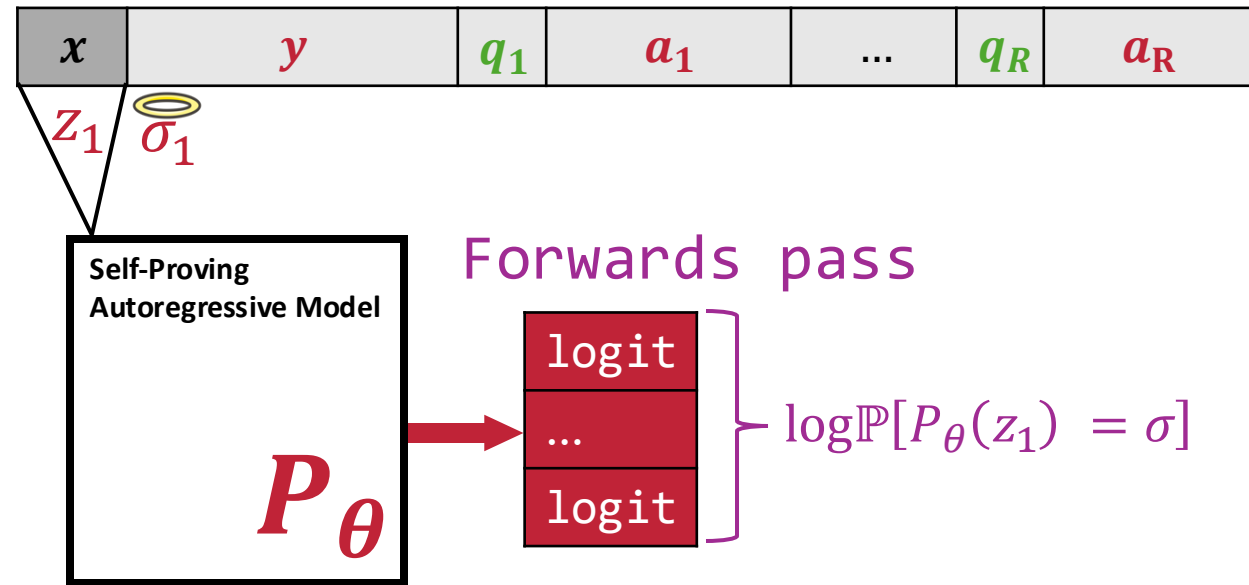


# Transcript Learning

## Step 1: Collect accepted transcripts

## Step 2: Transcript Cloning

For each transcript:



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

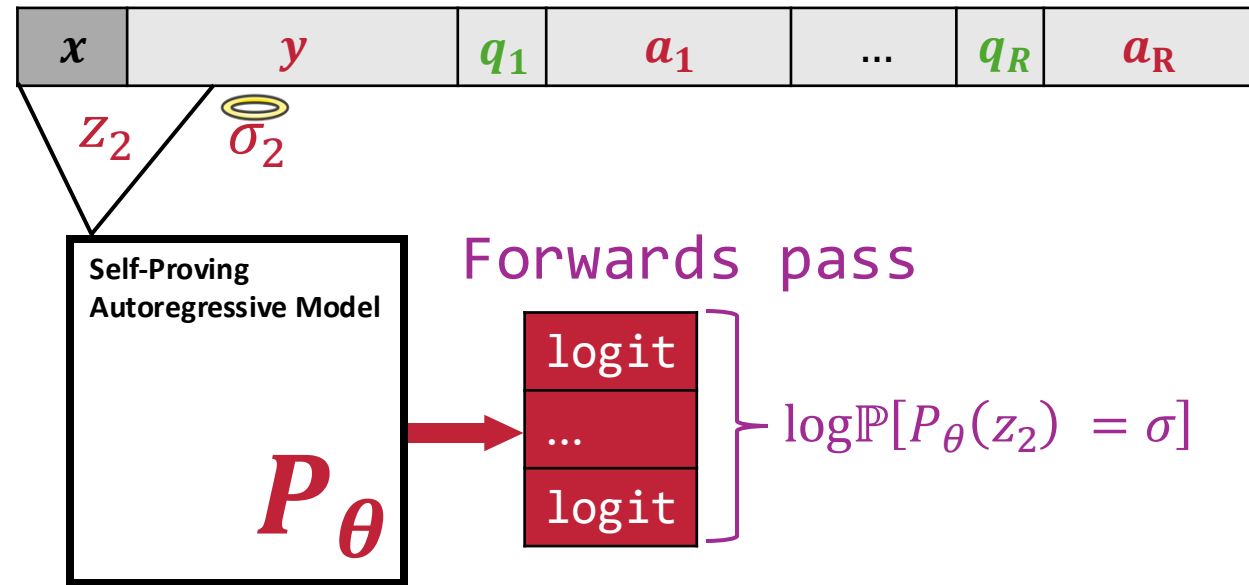
$x^1$	$y^1$	$q_1^1$	$a_1^1$	...	$q_R^1$	$a_R^1$
$x^N$	$y^N$	$q_1^N$	$a_1^N$	...	$q_R^N$	$a_R^N$

# Transcript Learning

Step 1: Collect accepted transcripts

Step 2: Transcript Cloning

For each transcript:



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

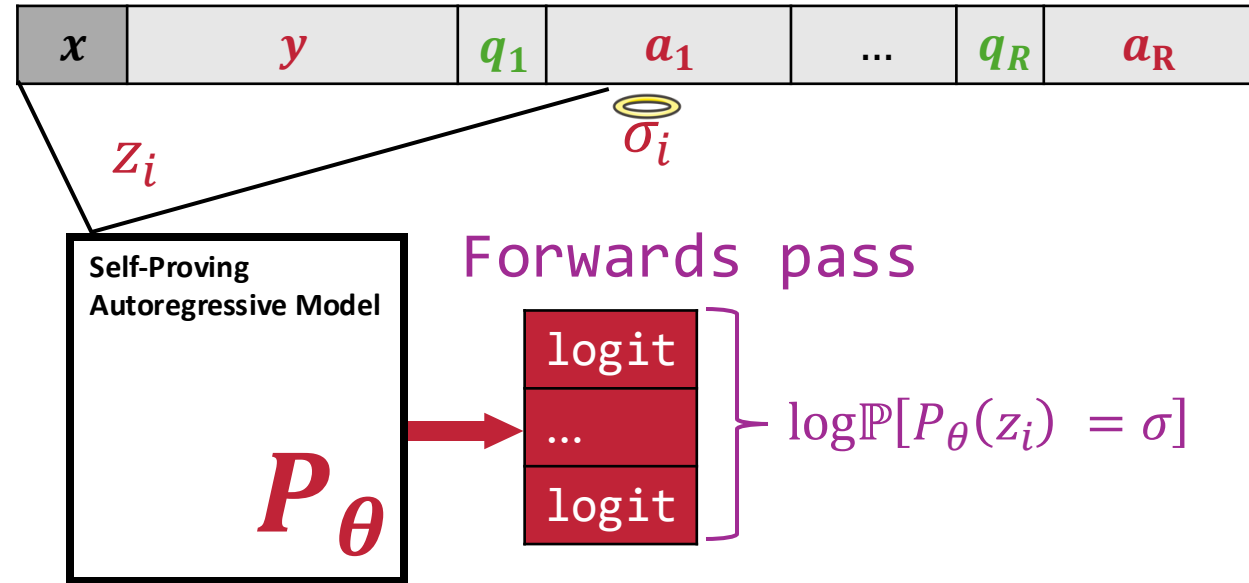
$x^1$	$y^1$	$q_1^1$	$a_1^1$	...	$q_R^1$	$a_R^1$
$x^N$	$y^N$	$q_1^N$	$a_1^N$	...	$q_R^N$	$a_R^N$

# Transcript Learning

Step 1: Collect accepted transcripts

Step 2: Transcript Cloning

For each transcript:



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

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

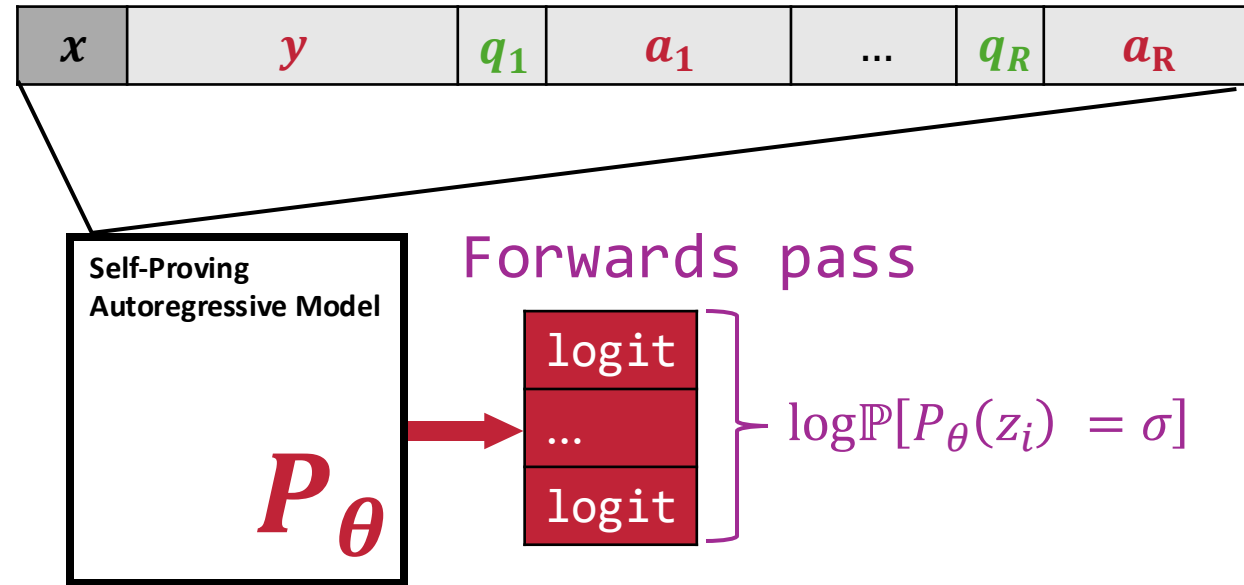
# Transcript Learning

## Step 1: Collect accepted transcripts

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

## 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  $\theta \leftarrow \theta + \lambda \cdot \underbrace{\mathbb{P}[P_\theta(x) = z]}_{\text{Forwards}} \cdot \underbrace{\sum_i \vec{d}_i}_{\text{Backwards}}$



# Transcript Learning Sample Complexity

Assumptions:

**Theorem:** Under the assumptions, given access to accepted transcripts, Transcript Learning outputs a  $(1 - \epsilon)$ -**Self-Proving model** after

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

# Transcript Learning Sample Complexity

Assumptions:

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

**Theorem:** Under the assumptions, given access to accepting transcripts, Transcript Learning outputs a  $(1 - \epsilon)$ -**Self-Proving model** after

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

# Transcript Learning Sample Complexity

Assumptions:

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

**Theorem:** Under the assumptions, given access to accepting transcripts, Transcript Learning outputs a  $(1 - \epsilon)$ -Self-Proving model after

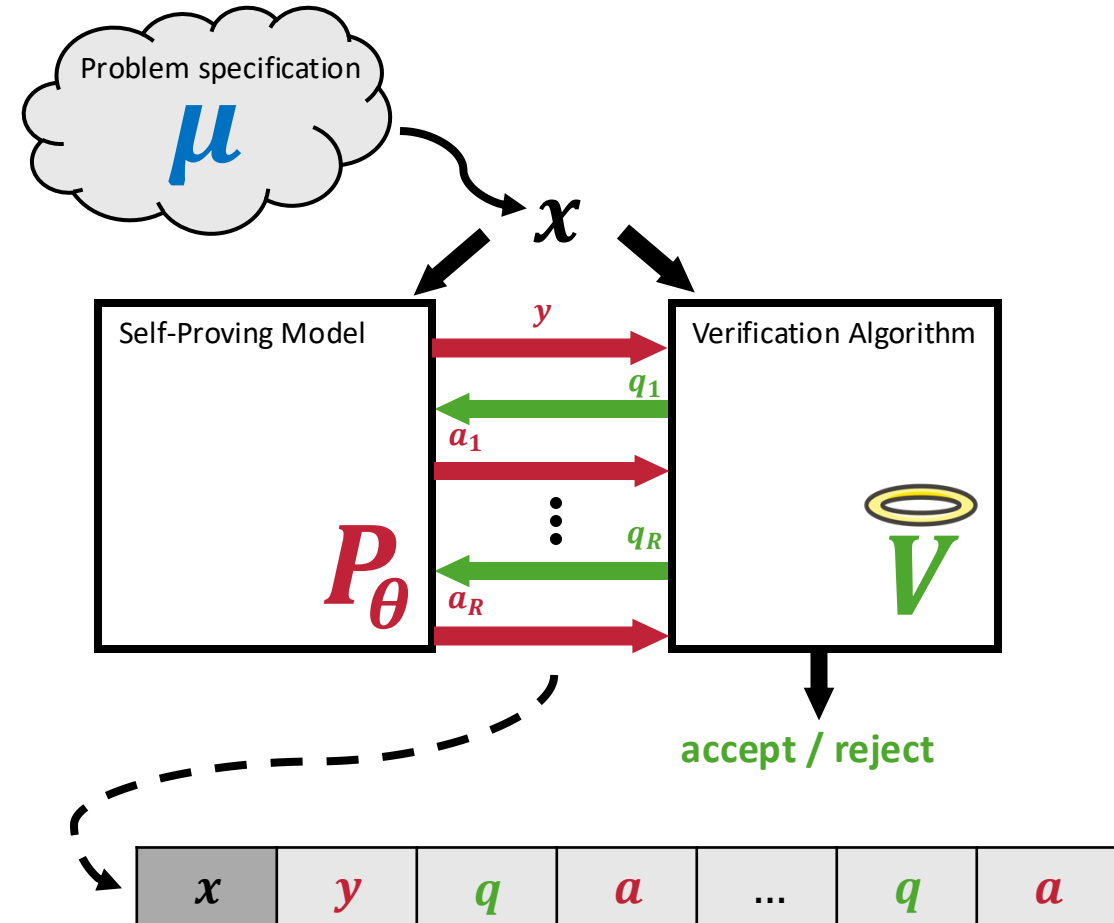
$$N \geq 4 \left( C \cdot B_1 \cdot B_2 \cdot \frac{1}{\epsilon} \right)^2 \text{ iterations}$$

**RLVF**

# RLVF

Repeat the following:

1. Generate transcript batch with  $P_\theta$ .  
Keep only accepted transcripts.

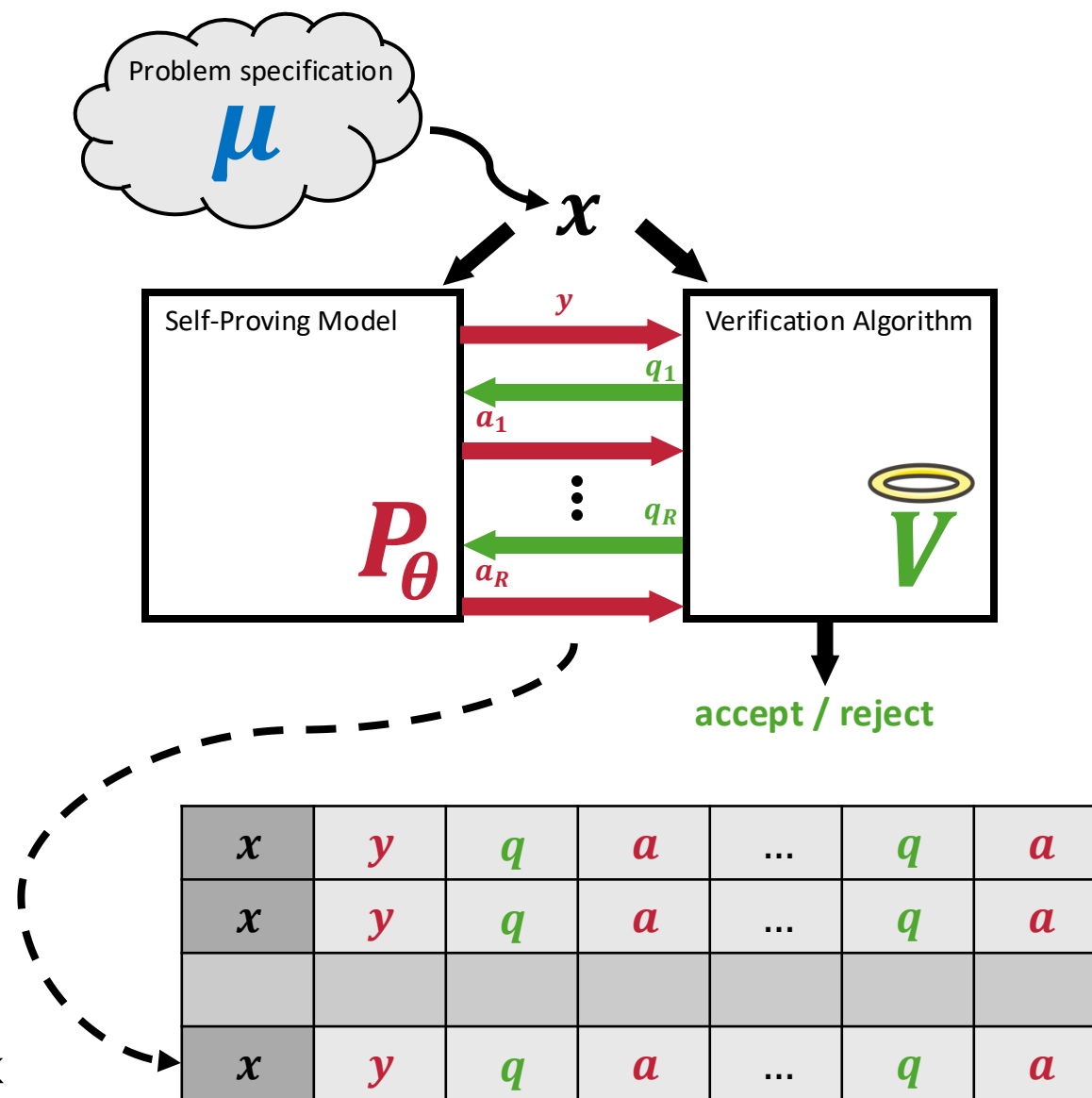


Reinforcement Learning from Verifier Feedback

# RLVF

Repeat the following:

1. Generate transcript batch with  $P_\theta$ .  
Keep only accepted transcripts.

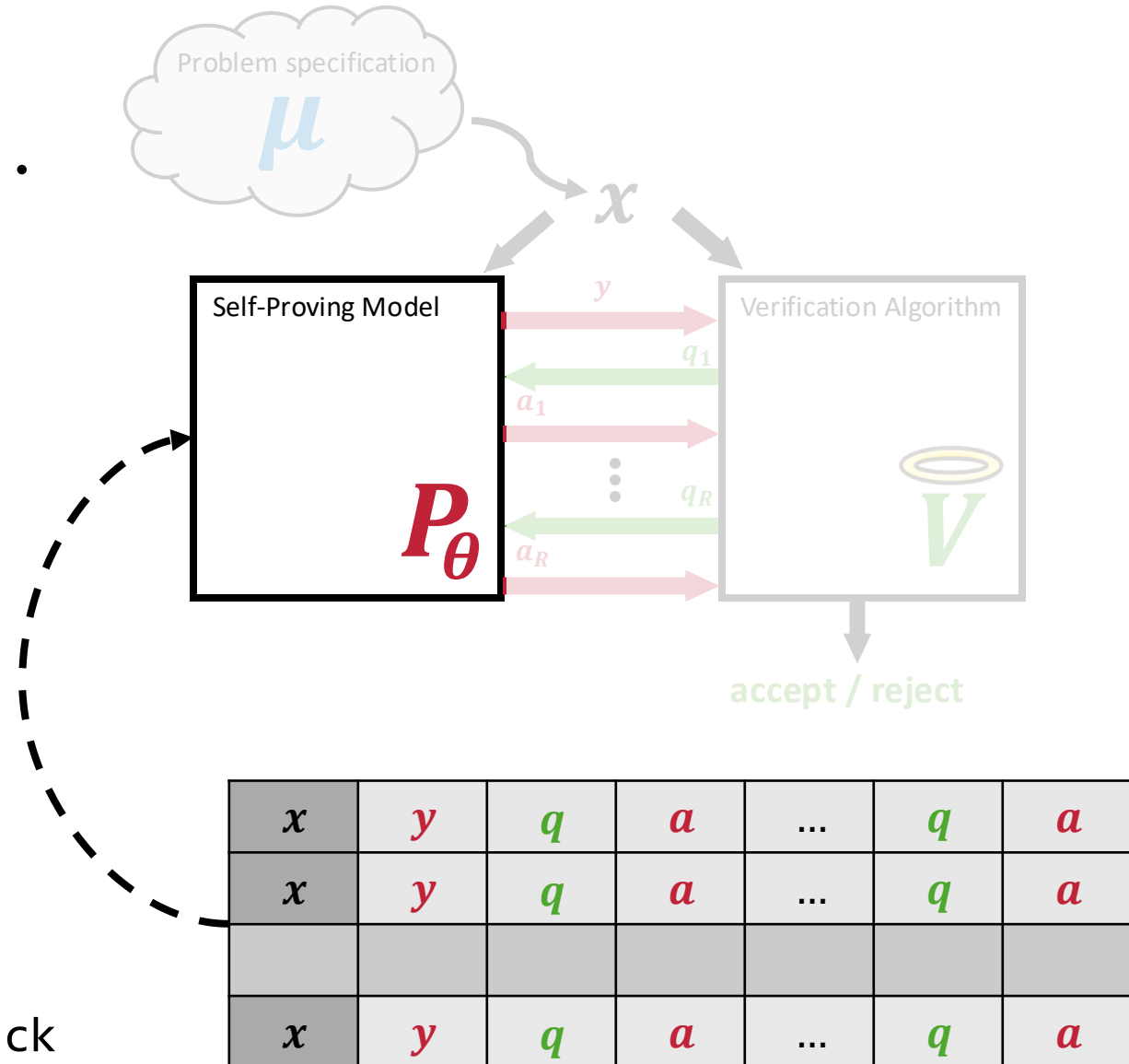


Reinforcement Learning from Verifier Feedback

# RLVF

Repeat the following:

1. Generate transcript batch with  $P_\theta$ .  
Keep only accepted transcripts.
2. Update  $\theta$  towards accepted trans.



Reinforcement Learning from Verifier Feedback



# Outline



Self-Proving  
models

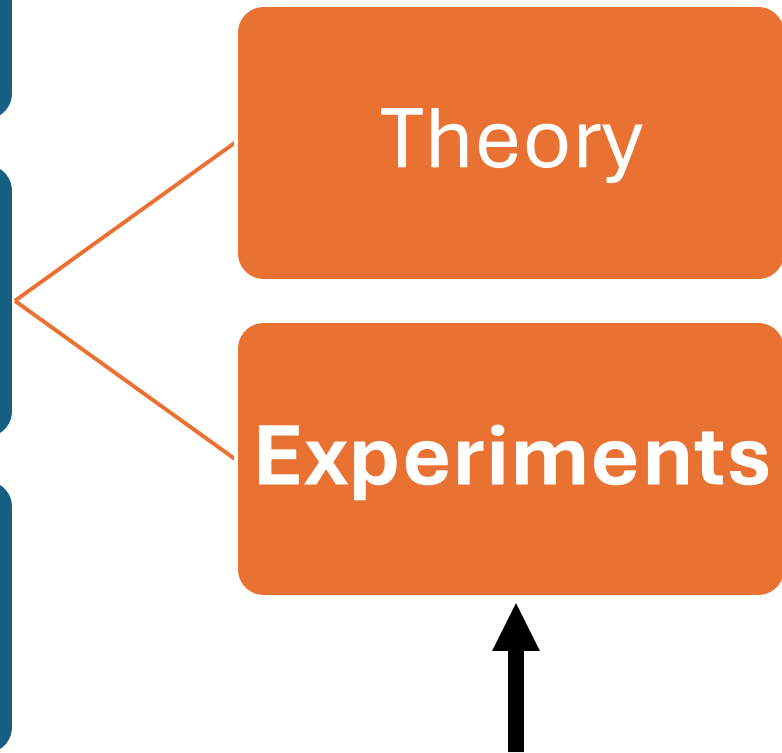
Related  
literature

Results

Future  
directions

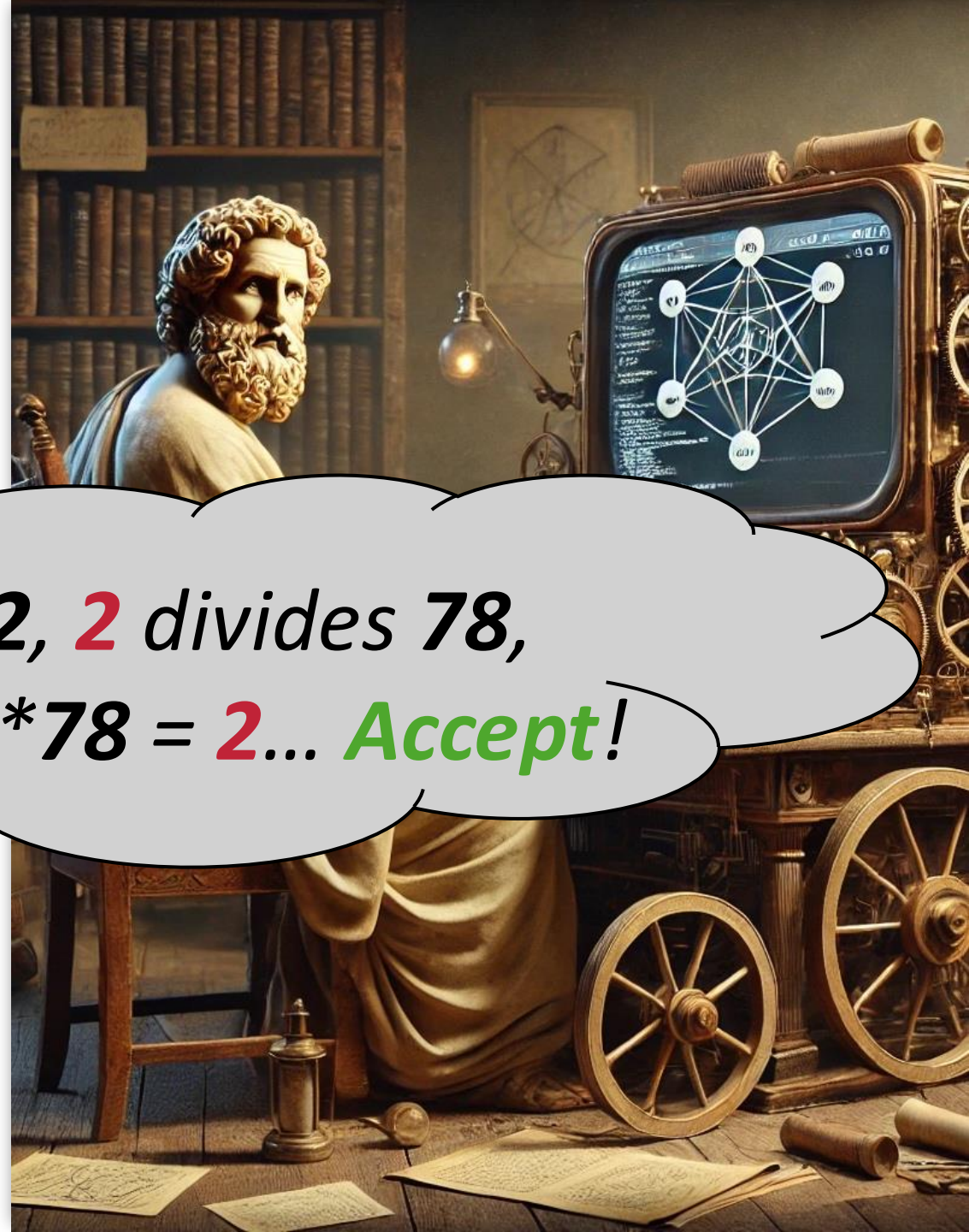
Theory

Experiments





# Experiments: Greatest Common Divisor



- 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  $(x_1, x_2) \in \mathbb{N}^2$   
if  $z_1x_1 + z_2x_2$  divides  $x_1$  and  $x_2$ ,  
then  $z_1x_1 + z_2x_2 = \text{GCD}(x_1, x_2)$
  - $V_{\text{GCD}}$  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  $\text{GCD}(92, 78)$ ?

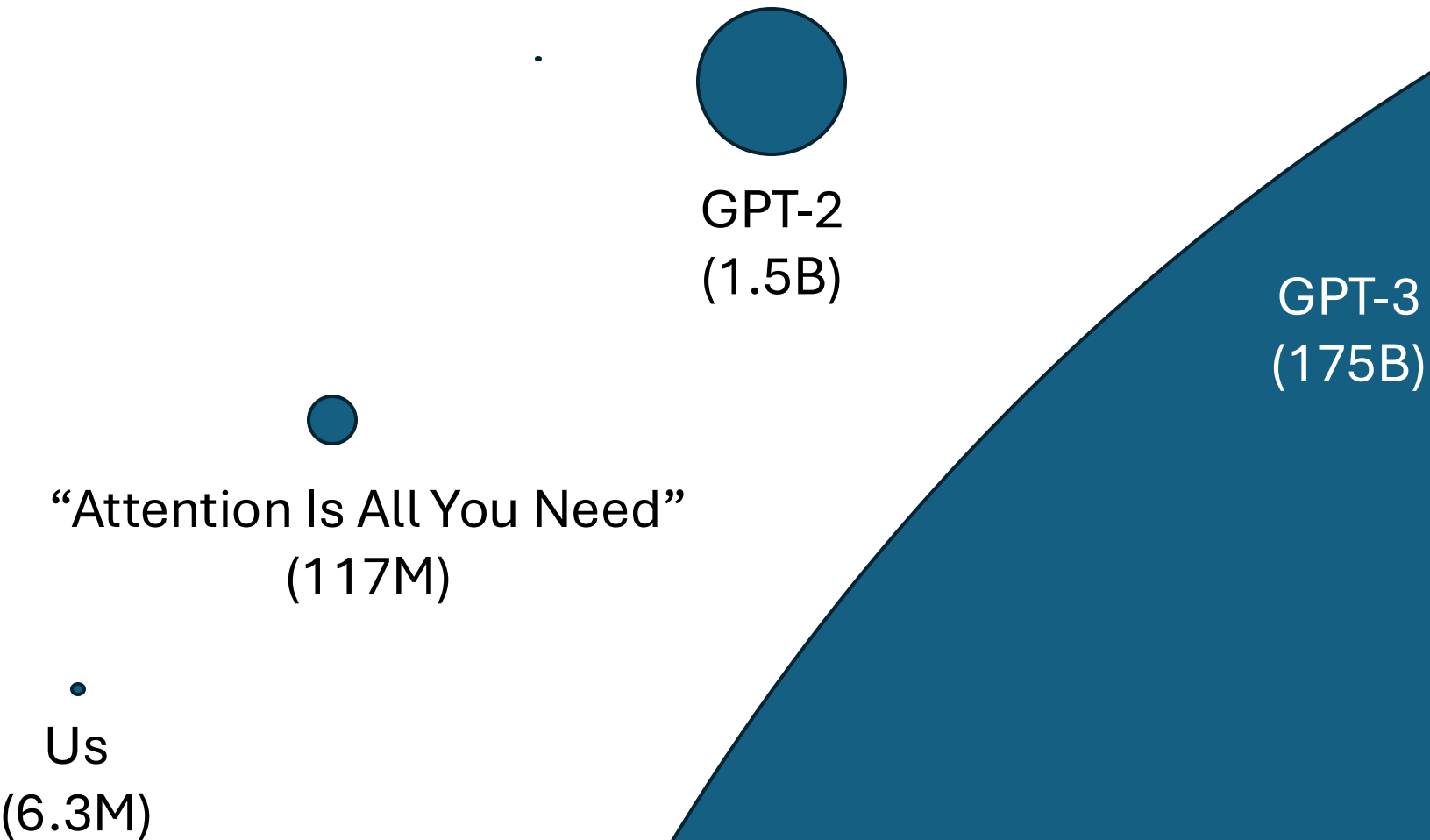
$P_\theta$

Answer: **2**

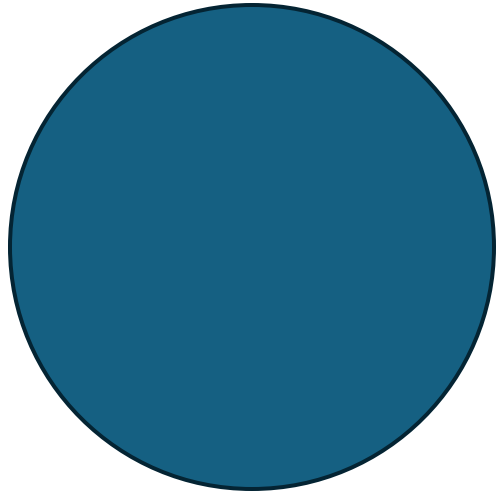
Proof:  $z_1 = -11, z_2 = 13$

$V$

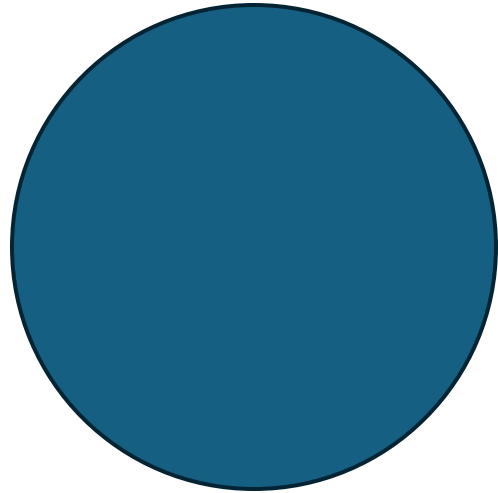
# Experiments: “LLMs” in the theory group



# Experiments: “LLMs” in the theory group



# Experiments: “LLMs” in the theory group



GPT-N

trained on



towards



•  
GPT (our exps)

trained on

$+ , 4 , 6 , x_0 , + , 3 , 9 , x_1 , + , 1 , y ,$   
 $+ , 1 , z_0' , + , 4 , 6 , z_1' , + , 1 , q' ,$   
 $+ , 0 , z_0'' , + , 3 , 9 , z_1'' , + , 5 , q''$   
 $+ , 1 , z_0''' , + , 7 , z_1''' , + , 1 , q''' ,$   
 $- , 1 , 1 , z_0 , + , 1 , 3 , z_1$

towards computing the GCD,  
and proving it to a sound verifier!

# Experimental results: Transcript Learning

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

# Experimental results: Transcript Learning

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 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

# 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

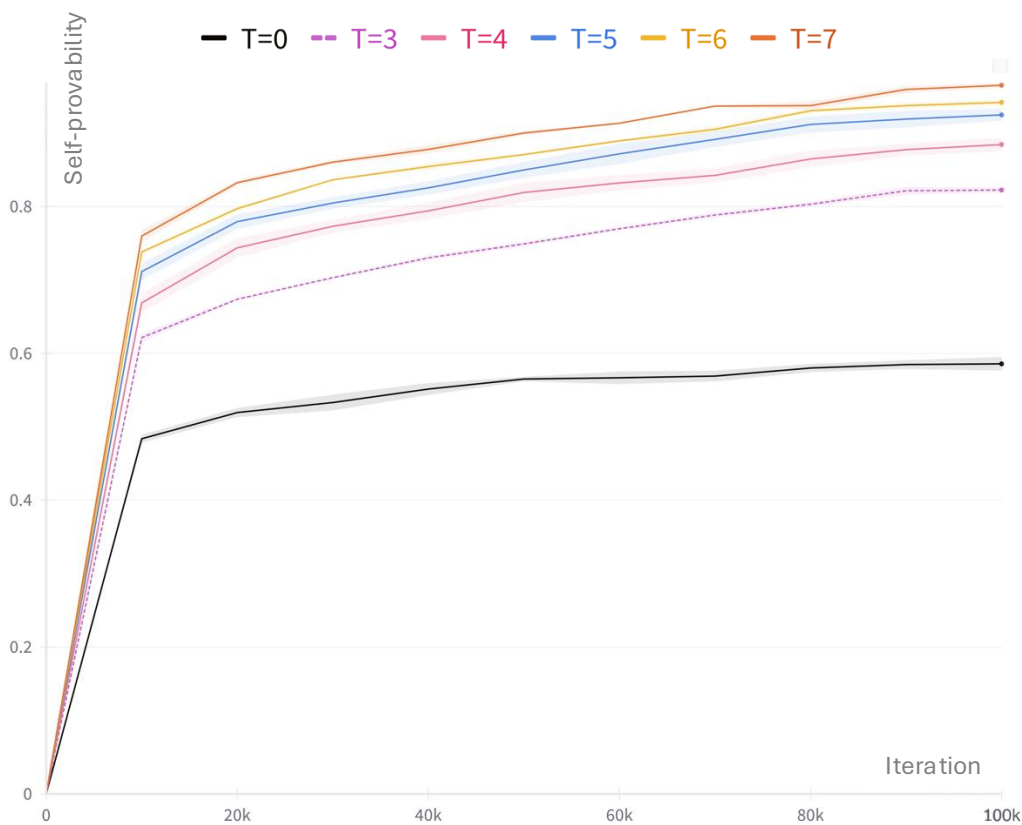
digits

delimiters



# Annotated Transcript Learning

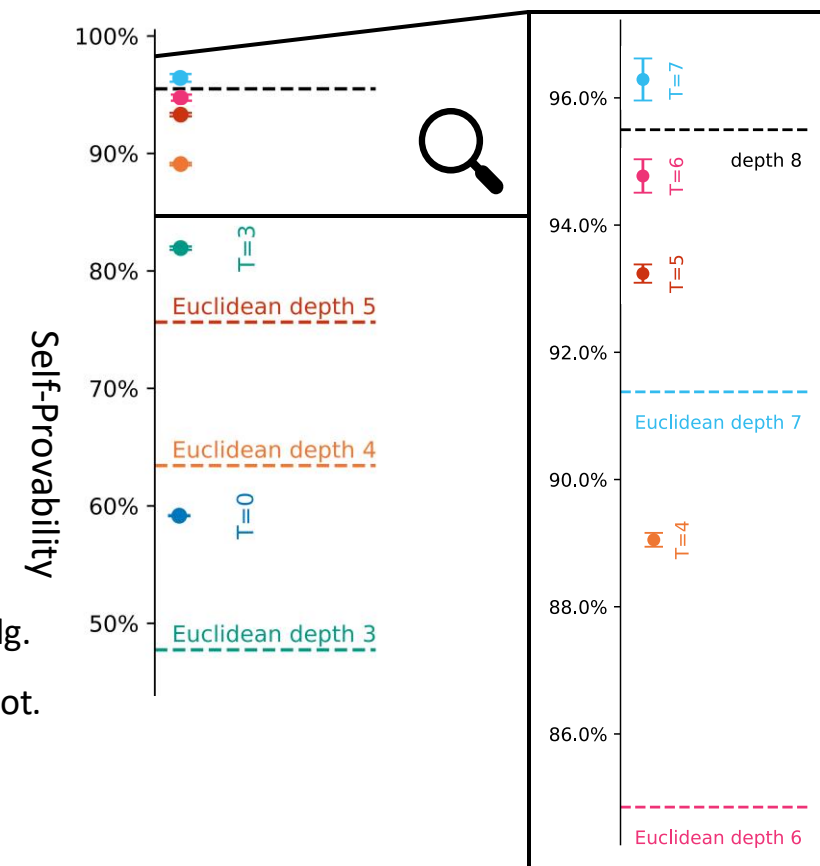
## Longer annotations help



## 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			

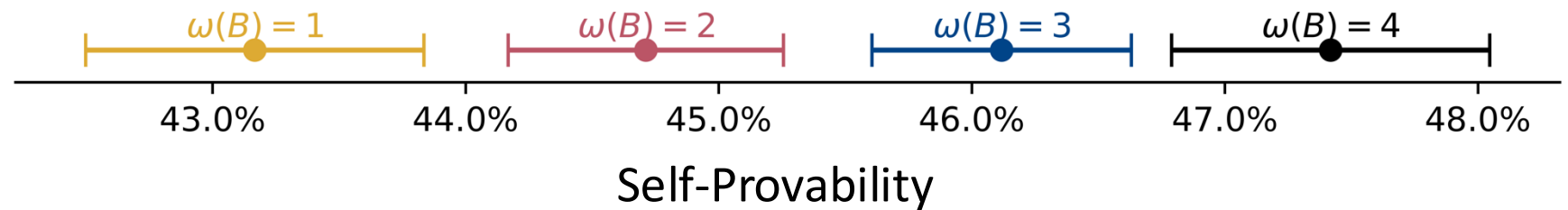
Depth( $x_0, x_1$ ) = #steps in Euclidean alg.  
 = length of "ideal" annot.



# Base of Representation

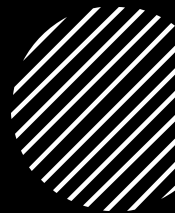
*“Early during training, transformers learn to predict products of divisors of the base  $B$  used to represent integers.” (Charton, 2024)*

- 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).





# Outline



Self-Proving  
models

Related  
literature

Results

Future  
directions



# Future directions

- Reinforcement Learning from Verifier Feedback (RLVF)
  - Sample complexity bounds via RL theory
- From Theory to Practice:
  - Proof system complexity: #rounds, randomized verifier, ...
  - Problem complexity
    - Can you learn to prove in a low-accuracy baseline?
    - Does Self-Provability increase accuracy?
  - “Practical” settings: Proof of harmlessness.
  - Larger models.
- 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  $P_\theta$  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?
- Connections to existing AI Safety frameworks

## Thank you!

