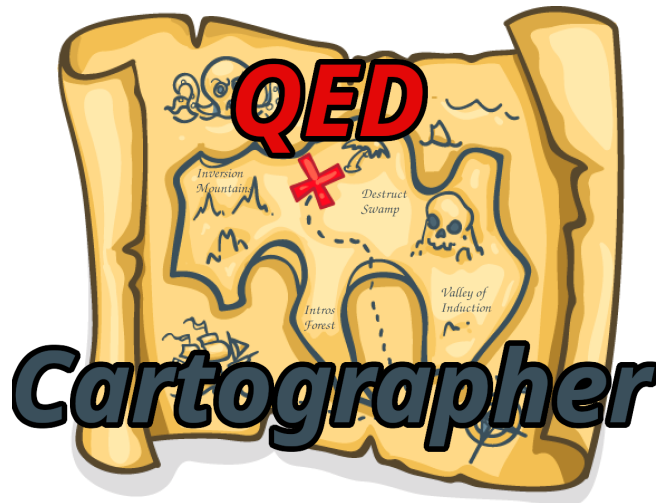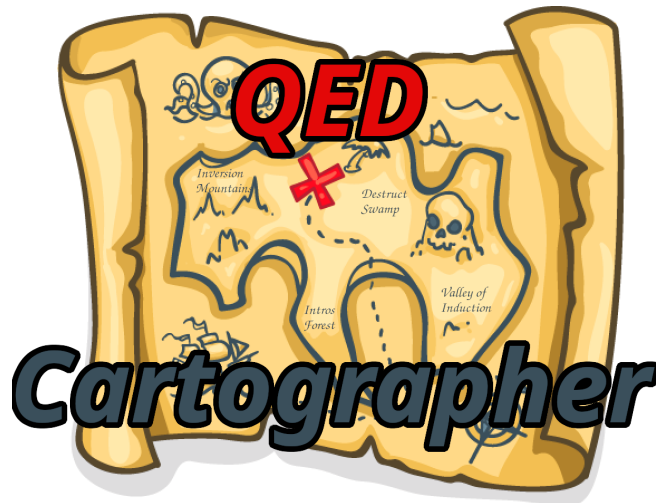# QED
# Cartographer

**Alex Sanchez-Stern**, Abhishek Varghese, Zhanna Kaufman,
Dylan Zhang, Talia Ringer, Yuriy Brun

# Automating Formal Verification with Reward-Free Reinforcement Learning
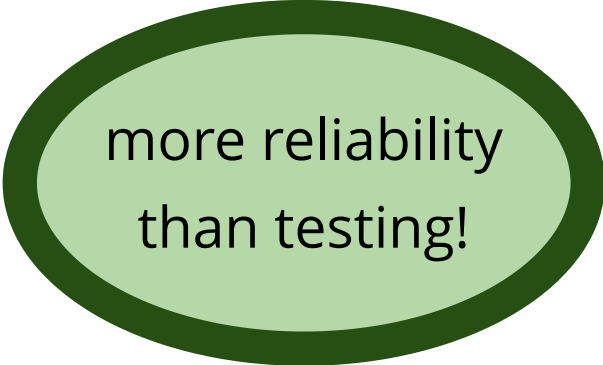
# QED Cartographer

Automating Formal Verification with Reward-Free Reinforcement Learning

# Formal Verification

writing proofs about programs!

# Formal Verification

writing proofs about programs!

more reliability
than testing!

# Formal Verification

writing proofs about programs!

more reliability than testing!

labor intensive

# Formal Verification

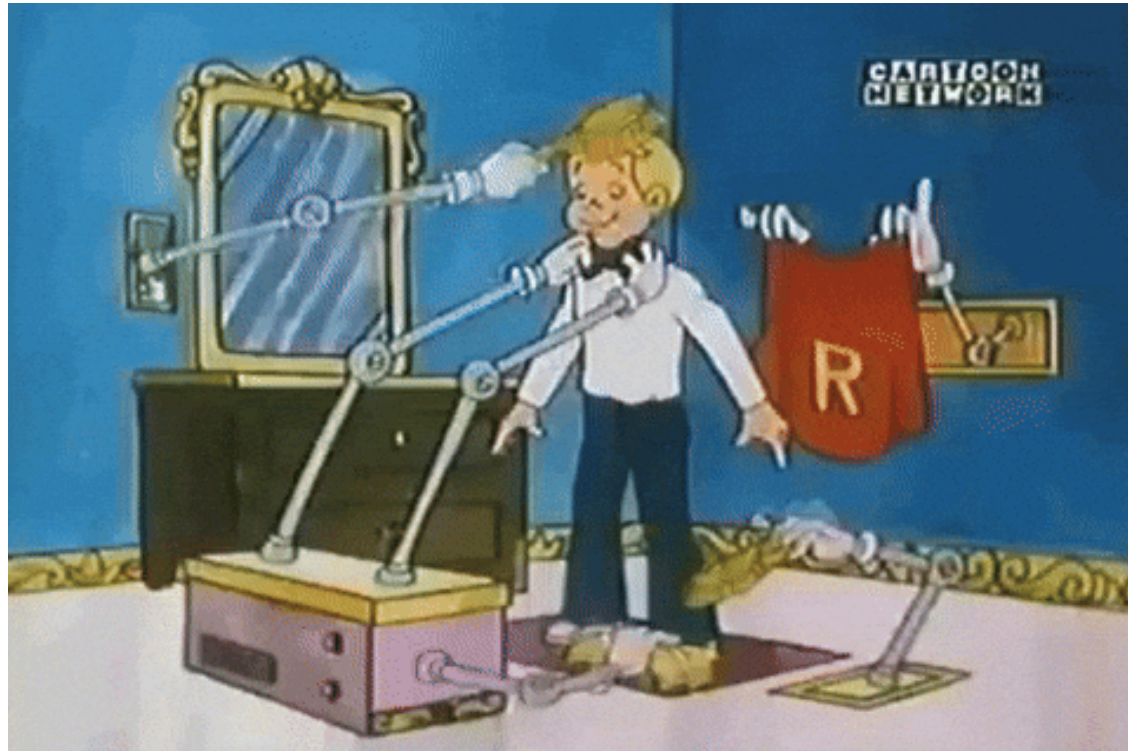writing proofs about programs!

more reliability than testing!

labor intensive
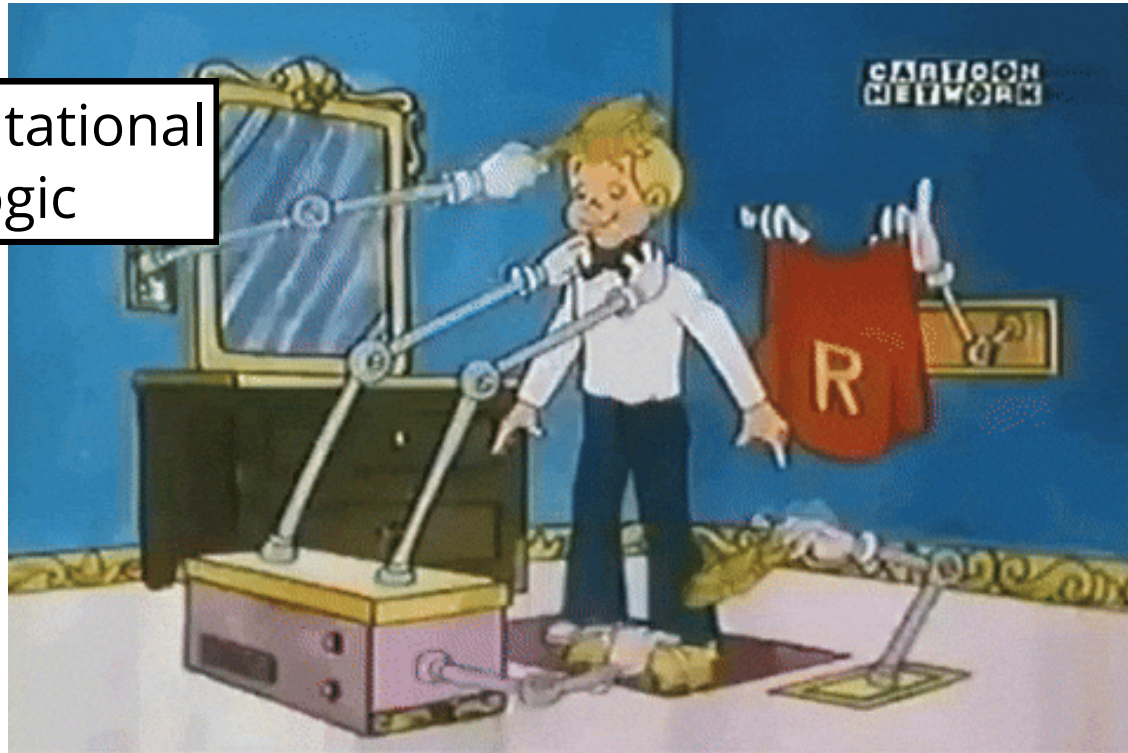
requires specialized expertise

# **Automating** Formal Verification
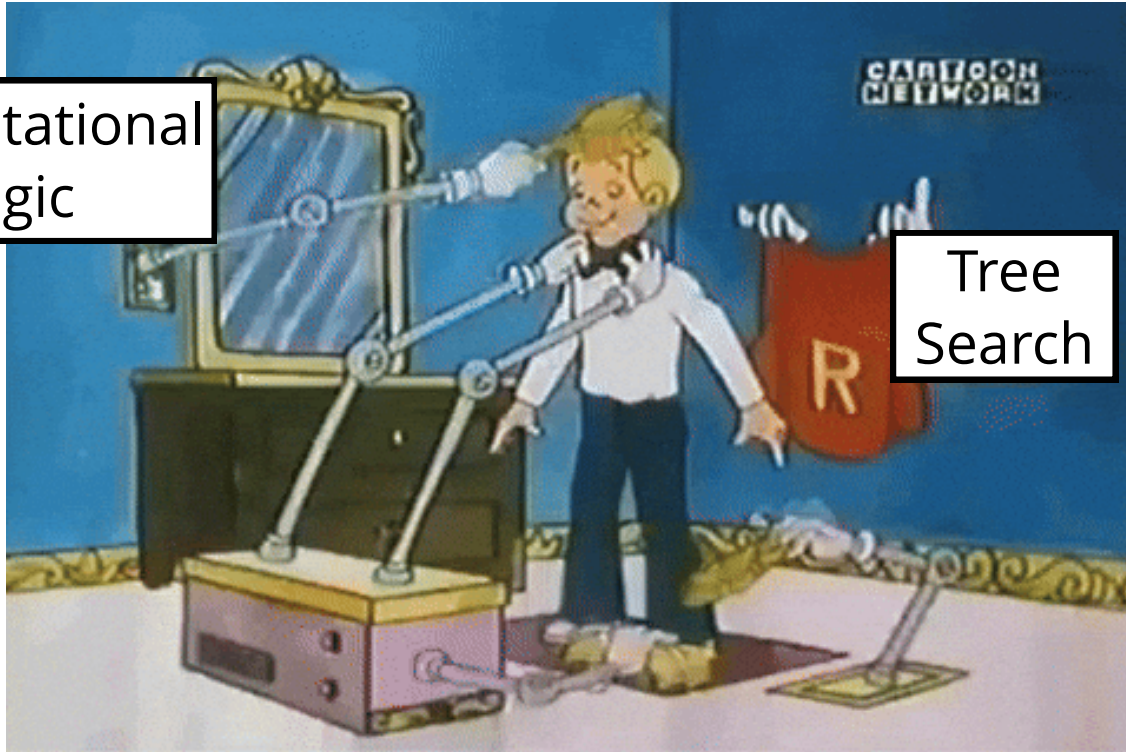
# **Automating** Formal Verification



Computational Logic

# **Automating** Formal Verification

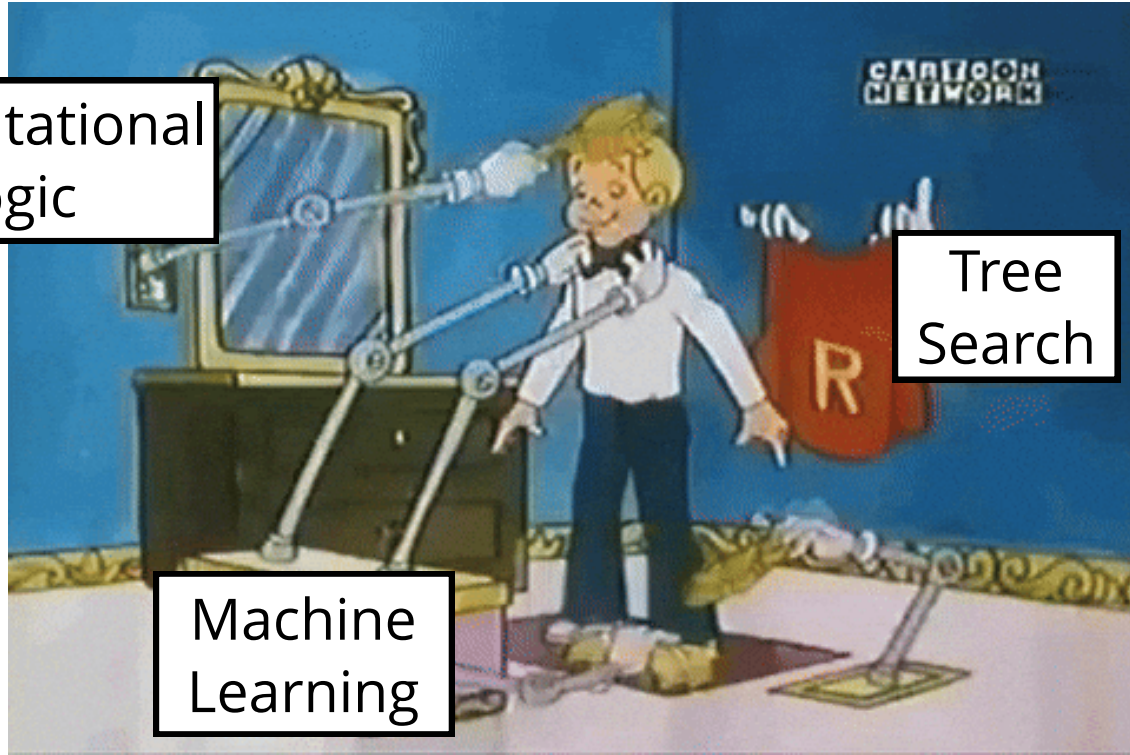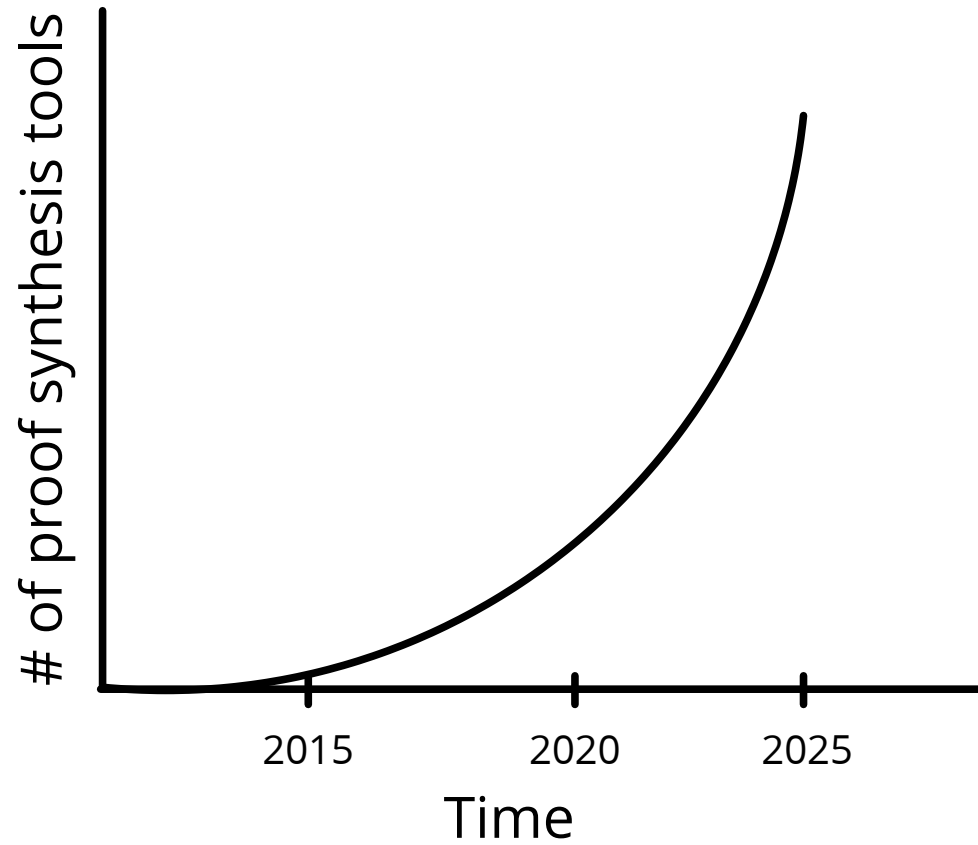# **Automating** Formal Verification

Proof Contexts

Actions

"forall (n: nat), (S n) > n" ?

destruct.
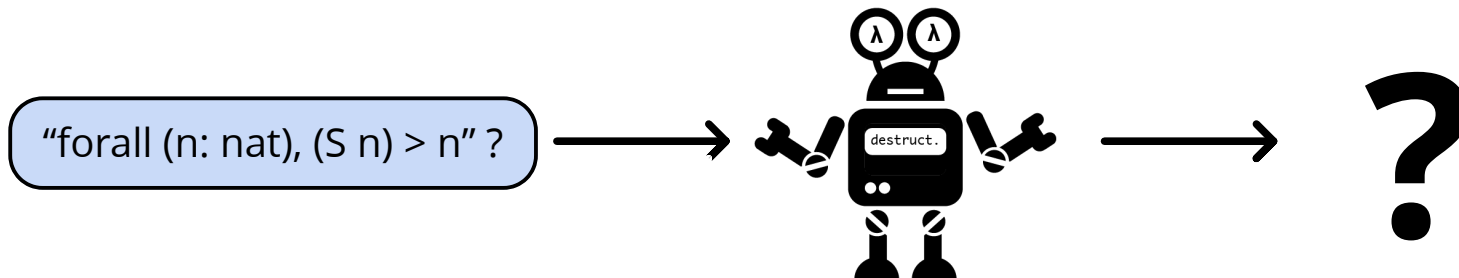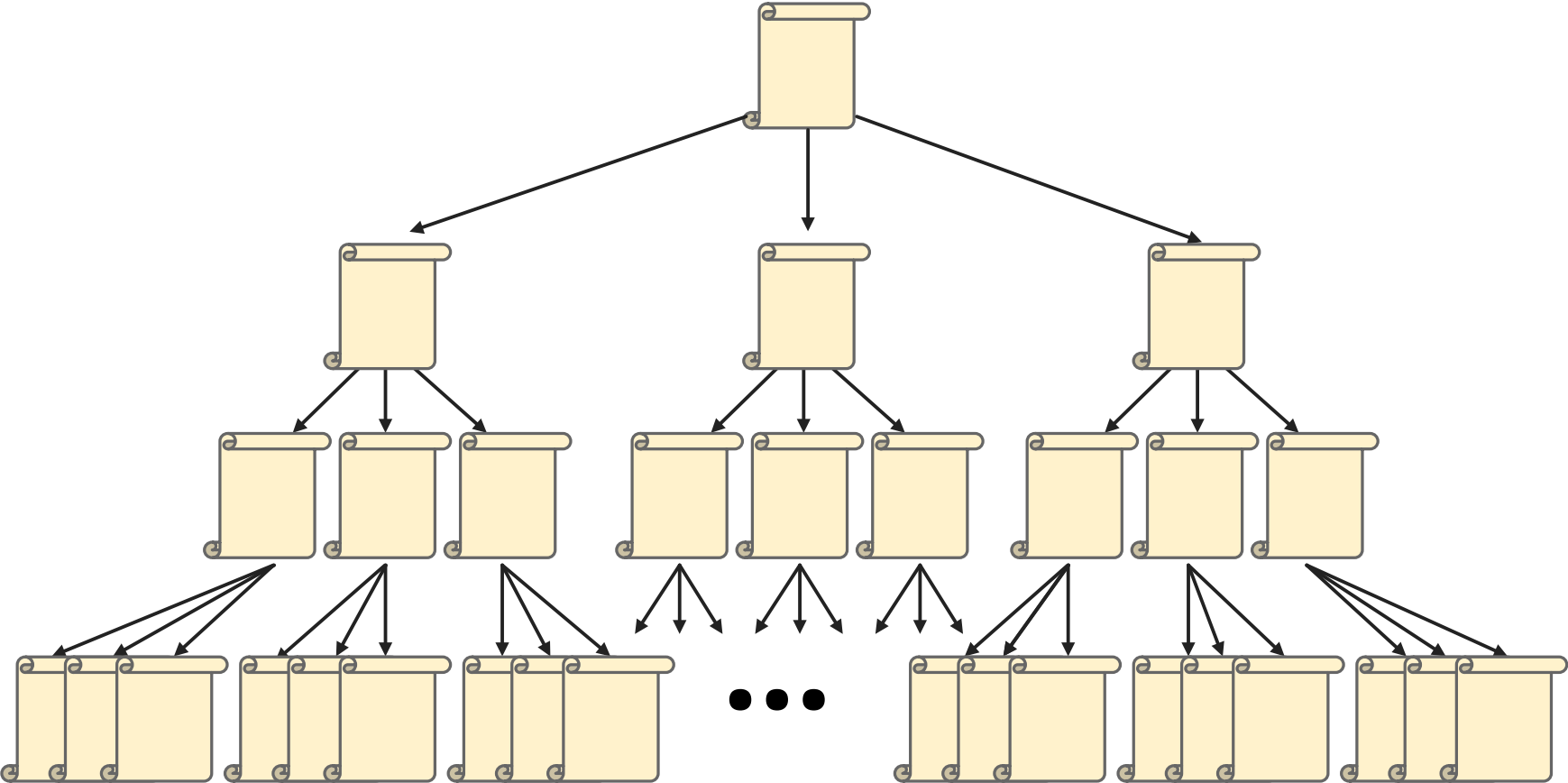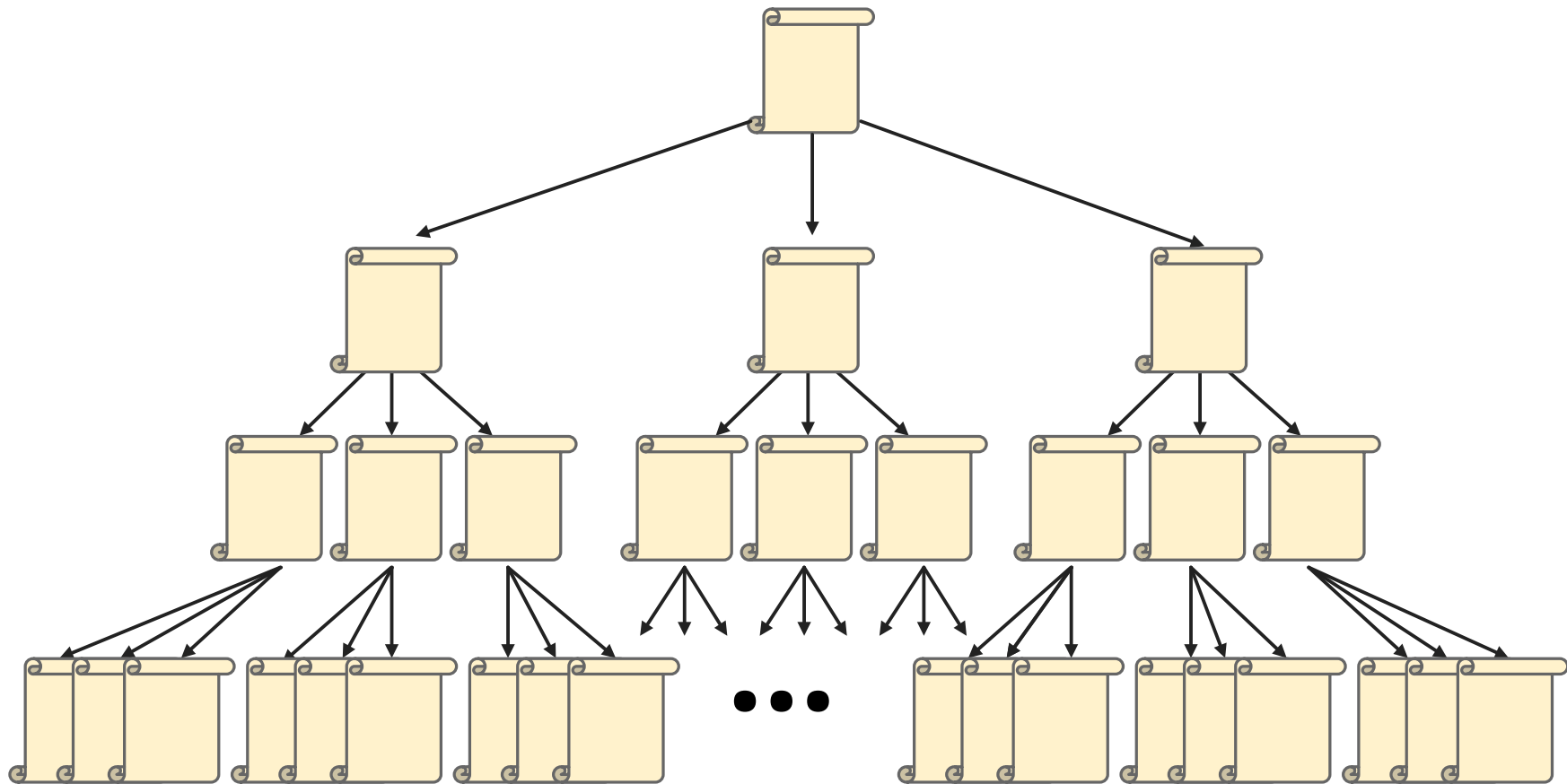
?

Unbounded Search Space

There are some techniques that can help us prune the search tree
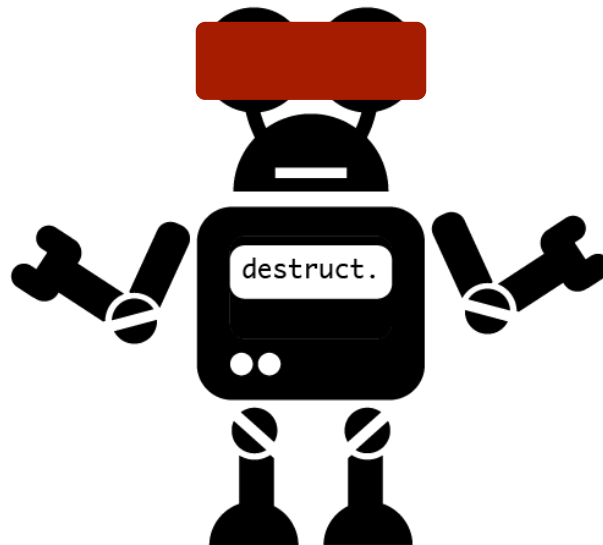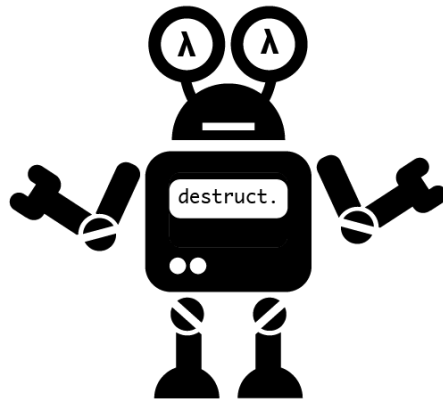
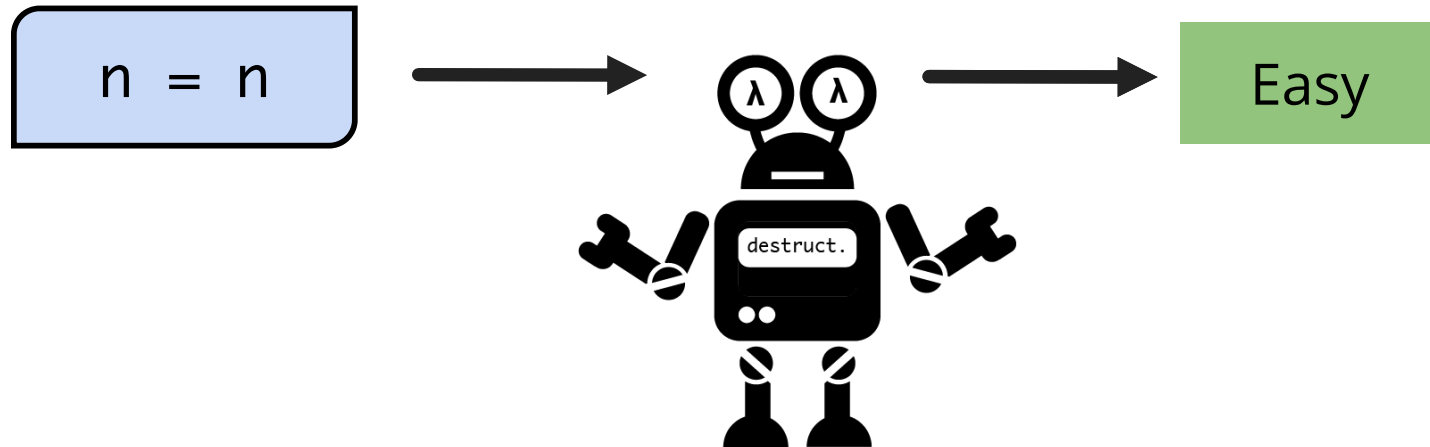There are some techniques that can help us prune the search tree

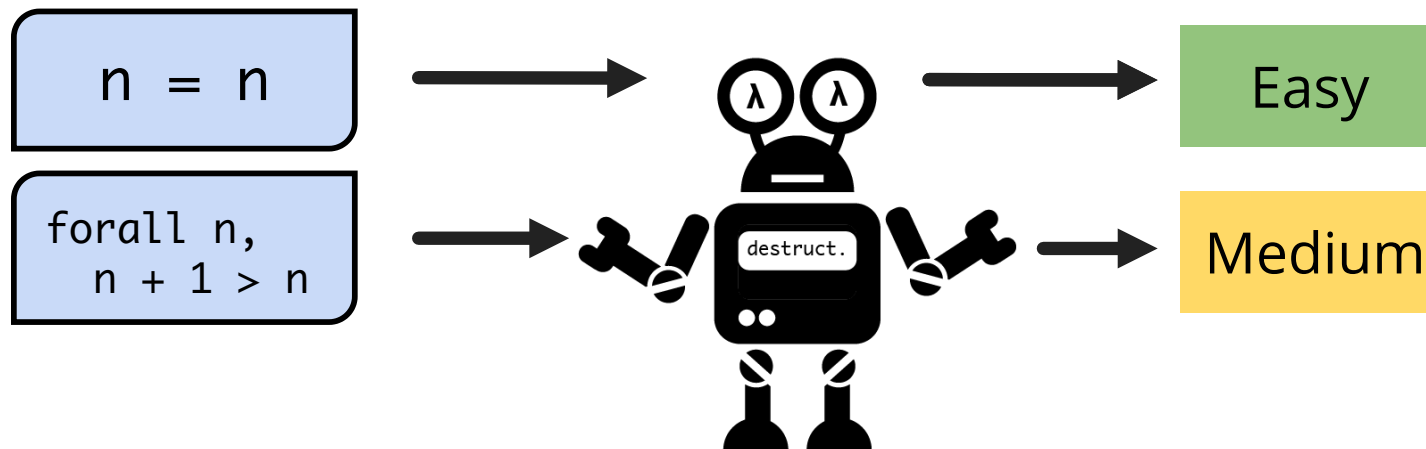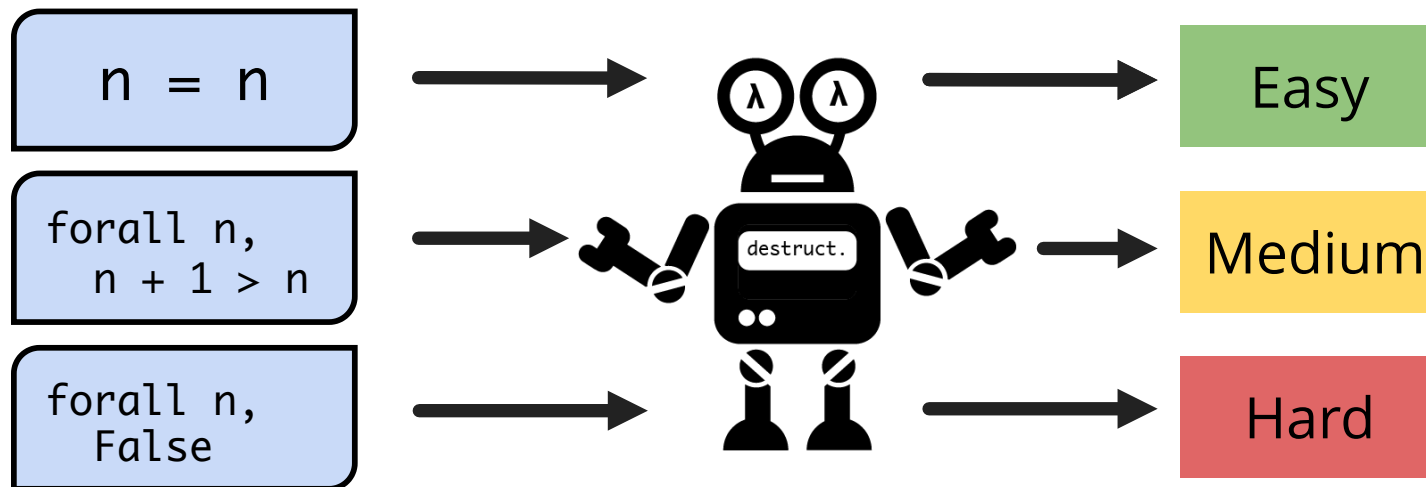It's hard to explore when you don't know where you are!

# We can search more efficiently if we can evaluate proof states
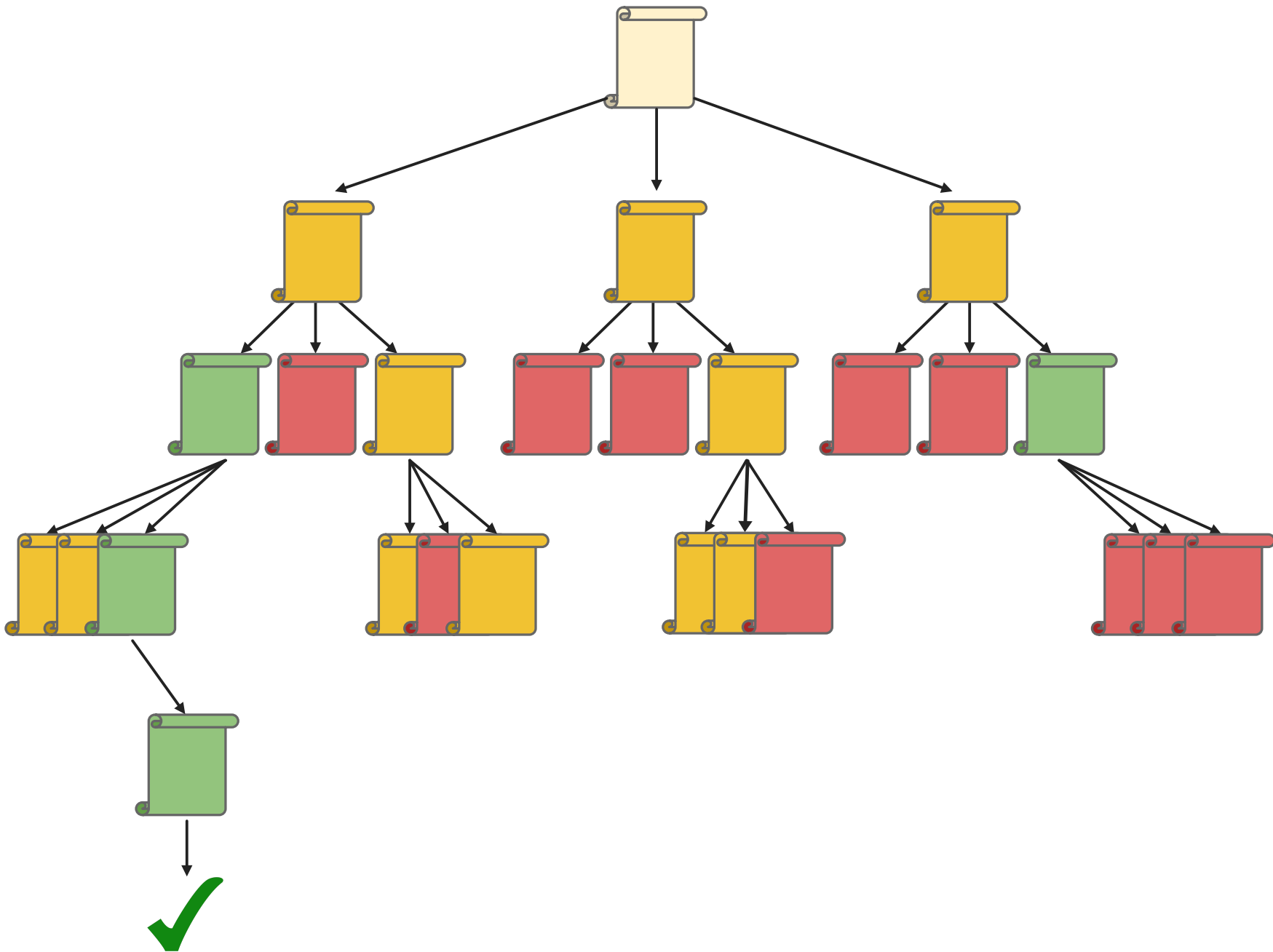
# We can search more efficiently if we can evaluate proof states

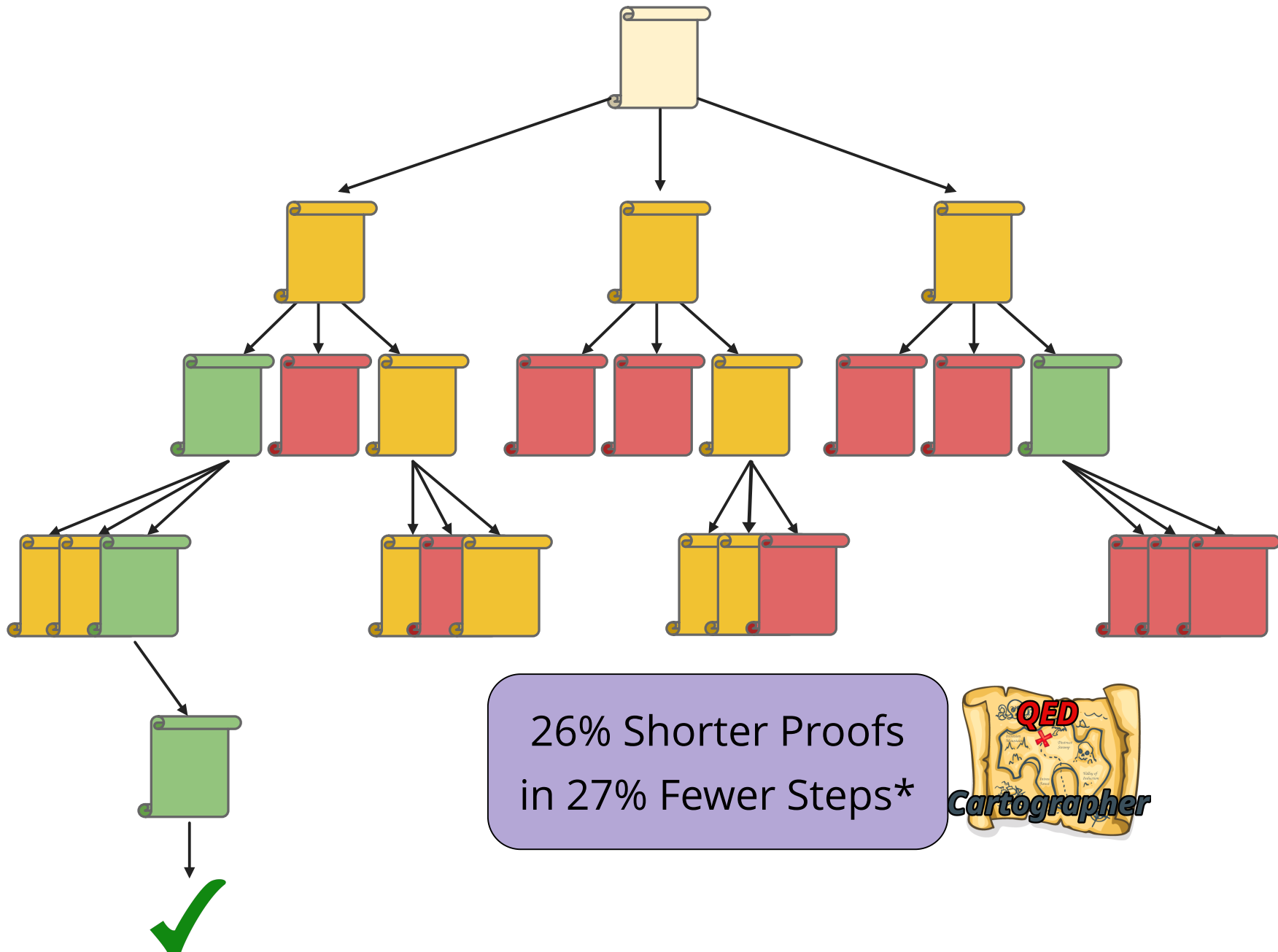# We can search more efficiently if we can evaluate proof states

# We can search more efficiently if we can evaluate proof states

26% Shorter Proofs
in 27% Fewer Steps*

QED
Cartographer

26% Shorter Proofs
in 27% Fewer Steps*

Cartographer

* Than a search without state scoring

Predictor

Search

Theorem Prover

Predictor

Search

Theorem Prover

QED

Cartographer

# "Reward-free Reinforcement Learning"

"Reward-free Reinforcement Learning"

In particular, V-learning

# This Talk

# This Talk

**V-Learning**

# This Talk

**V-Learning**

**Limitations in Proofs**

# This Talk

V-Learning

Limitations in Proofs

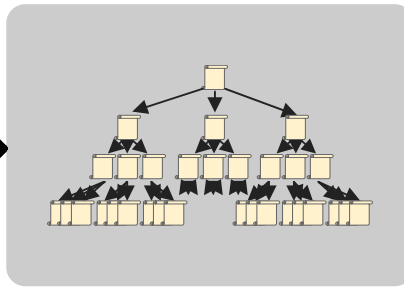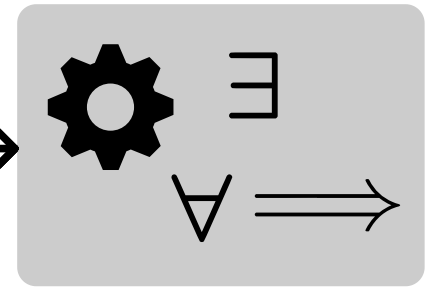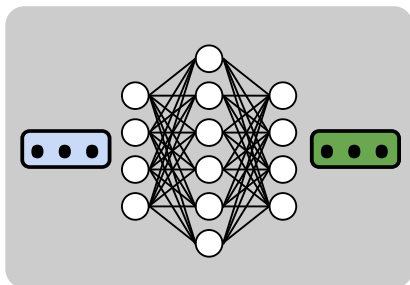Adapting to Proofs

# Classic V-Learning

# Classic V-Learning

# Classic V-Learning



$$V(S1) = \max($$
$$\phantom{V(S1) = \max(} + \quad V(S2),$$
$$\phantom{V(S1) = \max(} + \quad V(S4)),$$
$$\phantom{V(S1) = \max(} + \quad V(S5))$$

# Classic V-Learning



$$V(S1) = \max($$

$$\text{🪙} + \gamma V(S2),$$

$$\gamma (\text{🪙} + \gamma V(S4)),$$

$$\text{🪙} + \gamma V(S5))$$

$\gamma$ = time discount

finding rewards sooner is better!

$$V(S1) = \max(\text{🪙} + \gamma V(S2), \; \gamma \left(\text{🪙} + \gamma V(S4)\right), \text{🪙} + \gamma V(S5))$$

$$V(S1) = \max(\text{🪙} + \gamma V(S2), \ \gamma (\text{🪙} + \gamma V(S4)), \text{🪙} + \gamma V(S5))$$

**Generalizes to**

$$V(S) = \max_{a \in \text{actions}(S)} (R(S, a) + \gamma V(\text{next-state}(S, a)))$$

# V-Learning in Practice: Iterative Updates

# V-Learning in Proofs:
# The Sparse Reward Problem

# V-Learning in Proofs:
# The Sparse Reward Problem

# V-Learning in Proofs:
# The Sparse Reward Problem

# V-Learning in Proofs:
# The Sparse Reward Problem

# Insight:
# Proofs have Useful Structure

# Abstraction

State → ◆ Action ◆ → State

# Reality

State

· · ·

Goal ——— Action

· · ·

State

Adding extra rewards can lead to bad behavior

# Adding extra rewards can lead to bad behavior

$\forall x, P(x)$

# Adding extra rewards can lead to bad behavior



$\forall x, P(x)$

induction x

$\forall x, P(x)$
$\implies P(x+1)$

$P(0)$

# Adding extra rewards can lead to bad behavior

Inductive Case

$\forall x, P(x)$ — induction x — $\forall x, P(x) \implies P(x+1)$

$P(0)$

Base Case

# Adding extra rewards can lead to bad behavior

Inductive Case



Base Case

# Adding extra rewards can lead to bad behavior

# Adding extra rewards can lead to bad behavior

Inductive Case



Base Case

# Adding extra rewards can lead to bad behavior

Inductive Case

$\forall x, P(x)$ —— induction x ——> $\forall x, P(x) \implies P(x+1)$ —— induction x ——> $\forall x, P(x) \implies P(x+1)$ —— induction x ——> $\forall x, P(x) \implies P(x+1)$

$P(0)$ ——> ✓  $P(0)$ ——> ✓  $P(0)$ ——> ✓

Base Case

# Adding extra rewards can lead to bad behavior



Inductive Case

$\forall x, P(x)$ → induction x → $\forall x, P(x) \implies P(x+1)$ → induction x → $\forall x, P(x) \implies P(x+1)$ → induction x → $\forall x, P(x) \implies P(x+1)$

$P(0)$ ✓   $P(0)$ ✓   $P(0)$ ✓

Base Case

Reward-free doesn't have this problem!

# What We Need

A new update equation that accounts for the branching structure of proofs

# Assumptions:

- The state of a completed proof has value 1

$$V(\checkmark) = 1$$

## Assumptions:

- The state of a completed proof has value 1

$$V(\checkmark) = 1$$



$$V(S4) = 0 + \max([\gamma 1]) = \gamma$$

# Assumptions:

- The state of a completed proof has value 1

$$V(\checkmark) = 1$$



$$V(S4) = 0 + \max([\gamma 1]) = \gamma$$

$$V(S3) = 0 + \max([\gamma(\gamma 1)]) = \gamma^2$$

# Assumptions:

- The state of a completed proof has value 1

$$V(\text{✔}) = 1$$



$$V(S4) = 0 + \max([\gamma 1]) = \gamma$$

$$V(S3) = 0 + \max([\gamma(\gamma 1)]) = \gamma^2$$

$$V(S3) = 0 + \max([\gamma(\gamma(\gamma 1)(]) = \gamma^3$$

# Assumptions:

- The state of a completed proof has value 1

$$V(\checkmark) = 1$$



$$V(S4) = 0 + \max([\gamma 1]) = \gamma$$

$$V(S3) = 0 + \max([\gamma(\gamma 1)]) = \gamma^2$$

$$V(S3) = 0 + \max([\gamma(\gamma(\gamma 1)(]) = \gamma^3$$

$$V(S3) = 0 + \max\left([\gamma(\gamma(\gamma(\gamma 1)))]\right) = \gamma^4$$

# Assumption: The state of a completed proof has value 1

$$V(\checkmark) = 1$$



$$V(S) = \gamma^{(\text{number of steps left})}$$

V(G1) = ???

G1
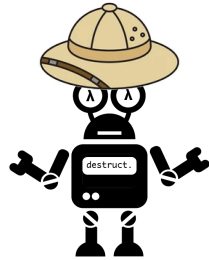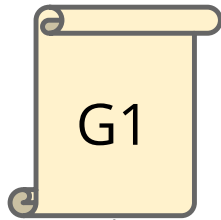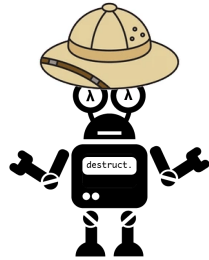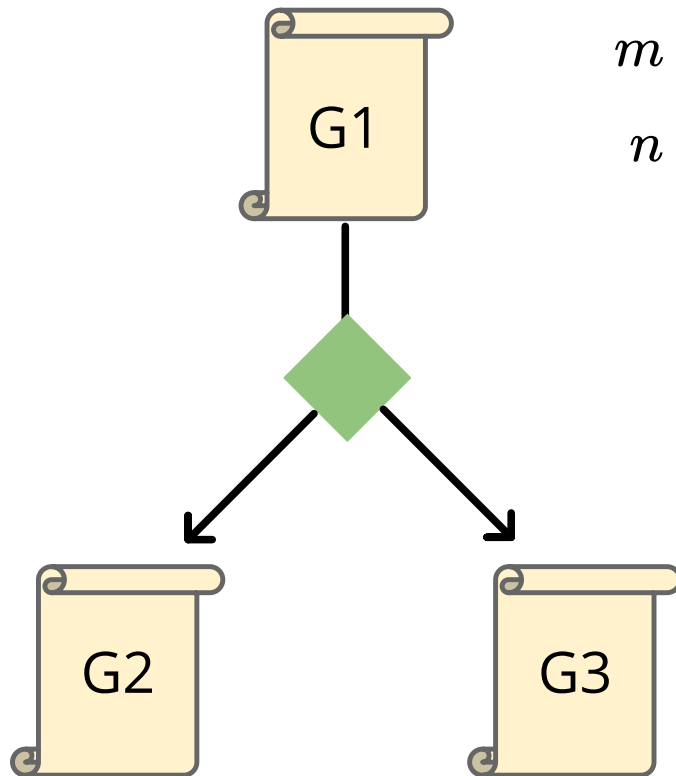
G2    G3

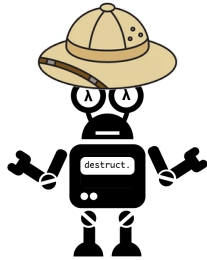V(G1) = ???

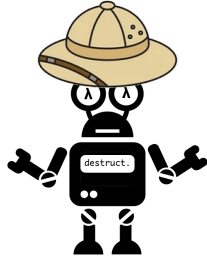$m$ = Steps to complete proof from G2

$n$ = Steps to complete proof from G3

V(G1) = ???

$m$ = Steps to complete proof from G2

$n$ = Steps to complete proof from G3

Steps to complete proof from G1 = $m + n + 1$

G1

G2

G3

# Update Equation for Branch-Structured Proofs

$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$

# Update Equation for
# Branch-Structured Proofs

$$V(G) = \boxed{\max_{a \in \text{actions}(G)}} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$

# Update Equation for
# Branch-Structured Proofs

$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{\boxed{G' \in \text{next-states}(G,a)}} V(G') \right)$$

# Update Equation for
# Branch-Structured Proofs

$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \boxed{\prod_{G' \in \text{next-states}(G,a)} V(G')} \right)$$

# Update Equation for Branch-Structured Proofs

$$V(G) = \max_{a \in \text{actions}(G)} \left( \boxed{\gamma} \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$

# Update Equation for
# Branch-Structured Proofs

$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$

# Update Equation for Branch-Structured Proofs

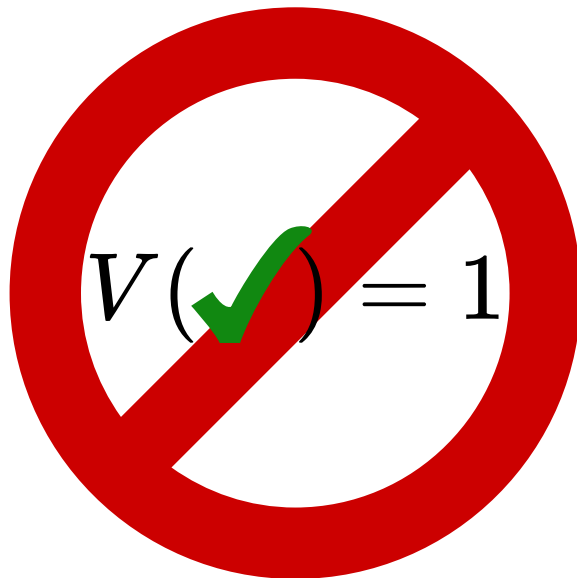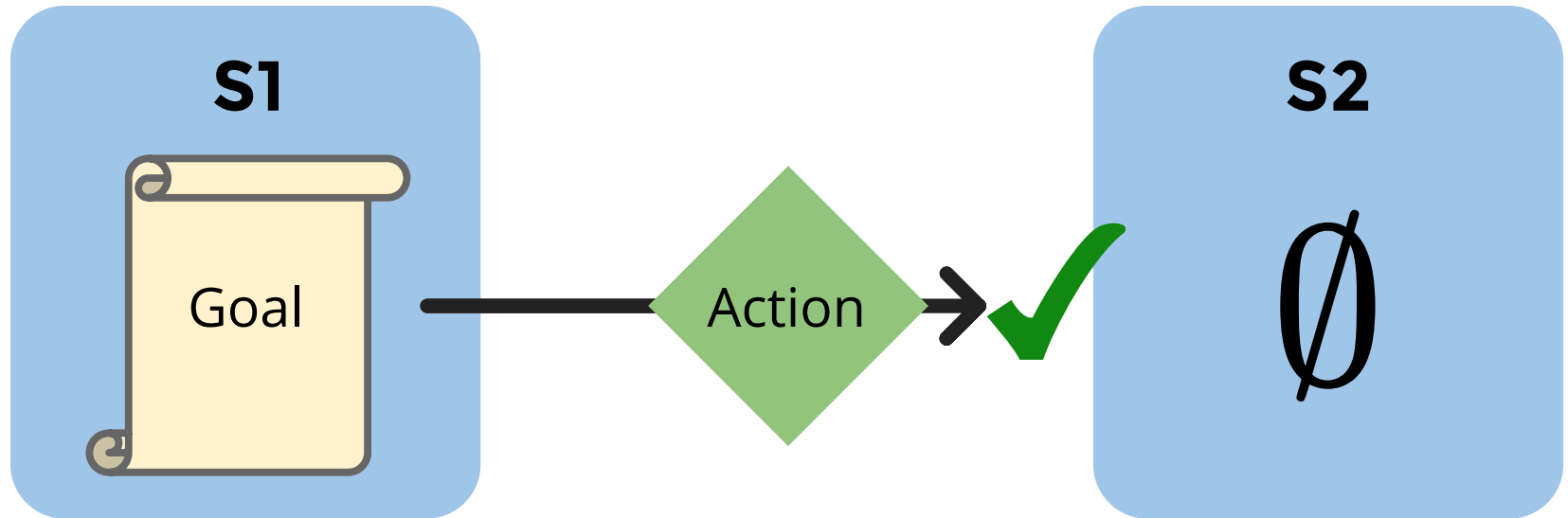$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$

$R(S,a)$

# Update Equation for
# Branch-Structured Proofs

$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$
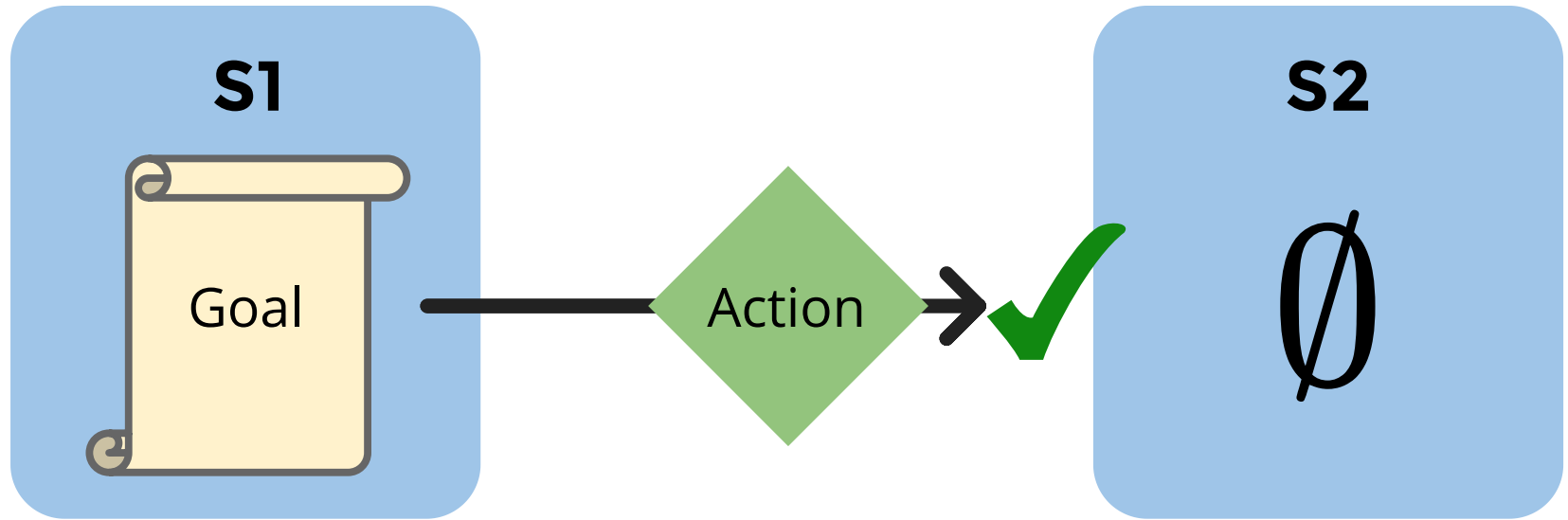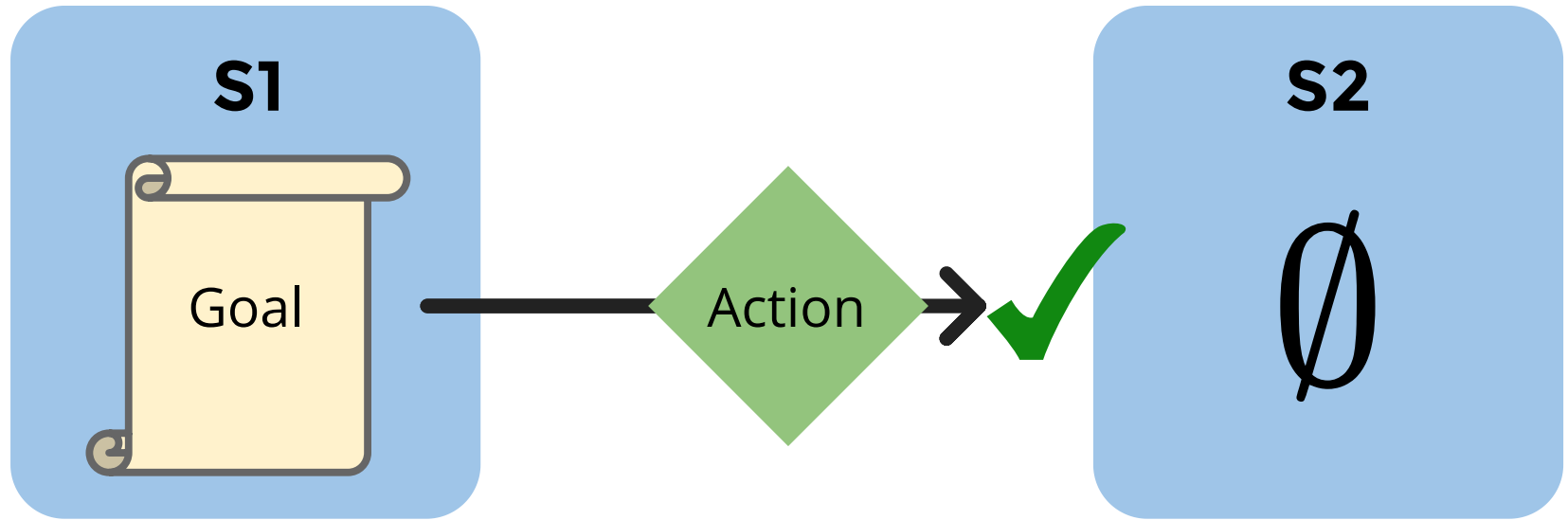
# Don't need one!

# Don't need one!



S1 — Goal → Action → ✔ → S2 — ∅

# Don't need one!



$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$
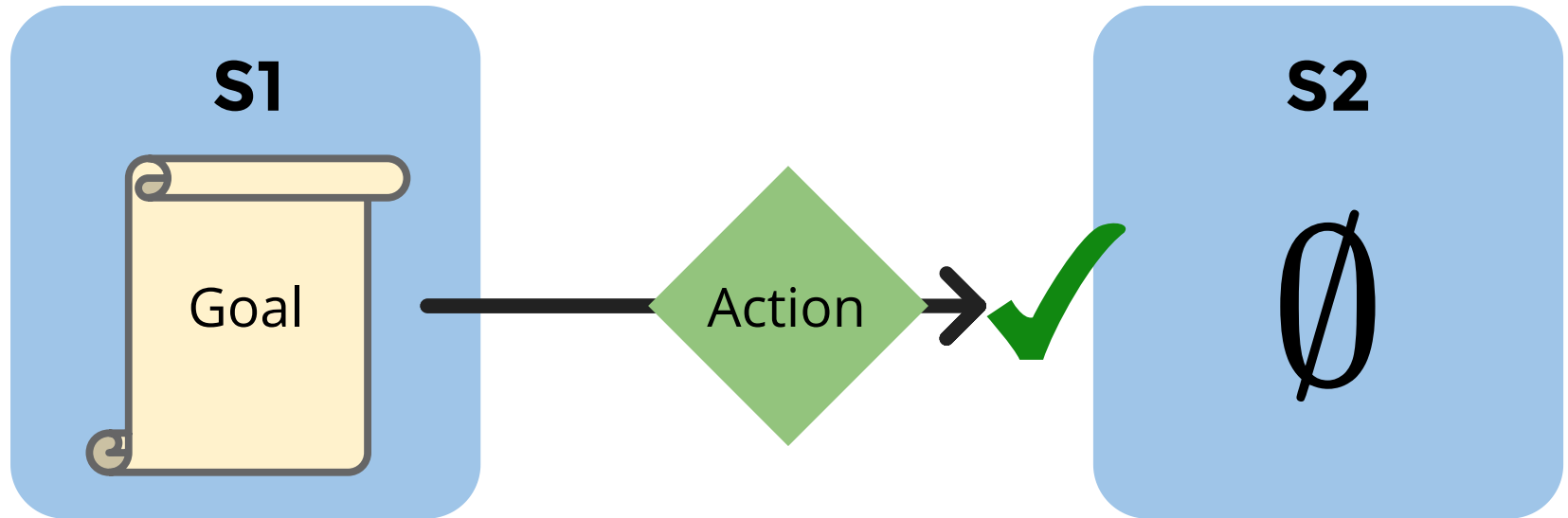
# Don't need one!
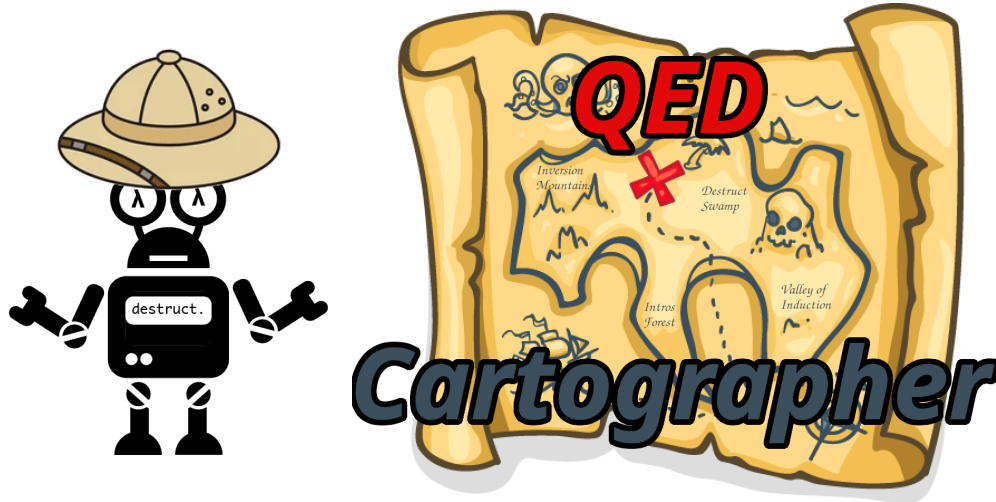
**S1**

Goal

Action ✔

**S2**

$\emptyset$

$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$
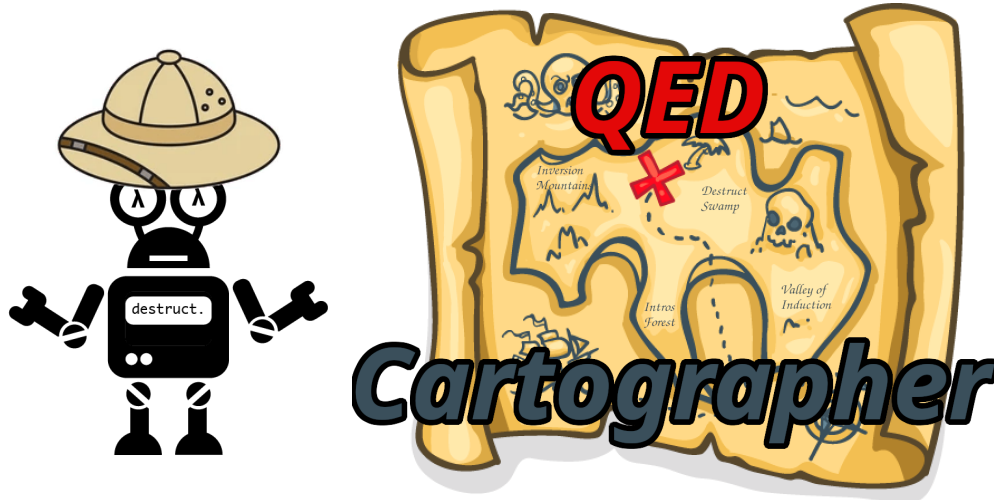
$$\gamma \prod_{G' \in \emptyset} V(G')$$

# Don't need one!

**S1**

Goal

Action ✔

**S2**

$\emptyset$

$$V(G) = \max_{a \in \text{actions}(G)} \left( \gamma \prod_{G' \in \text{next-states}(G,a)} V(G') \right)$$

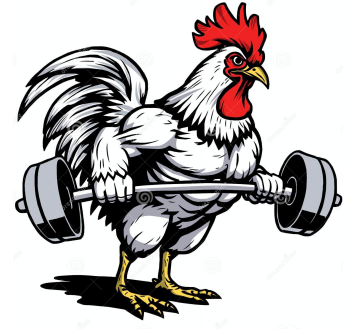$$\gamma \prod_{G' \in \emptyset} V(G')$$

$$\gamma(1)$$

# Automating Formal Verification with Reward-Free Reinforcement Learning

# Automating Formal Verification with Reward-Free Reinforcement Learning

26% Shorter Proofs

in 27% Fewer Steps
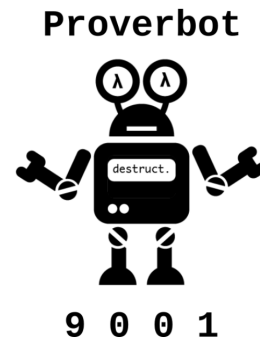
# Benchmark: CoqGym

124 Coq Projects

68,501 Theorems

85/15 train-test split

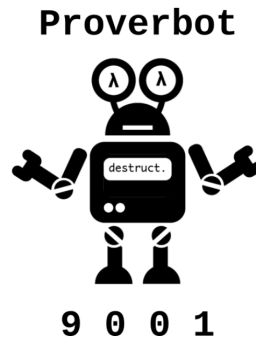# Baseline: Proverbot9001 (updated)

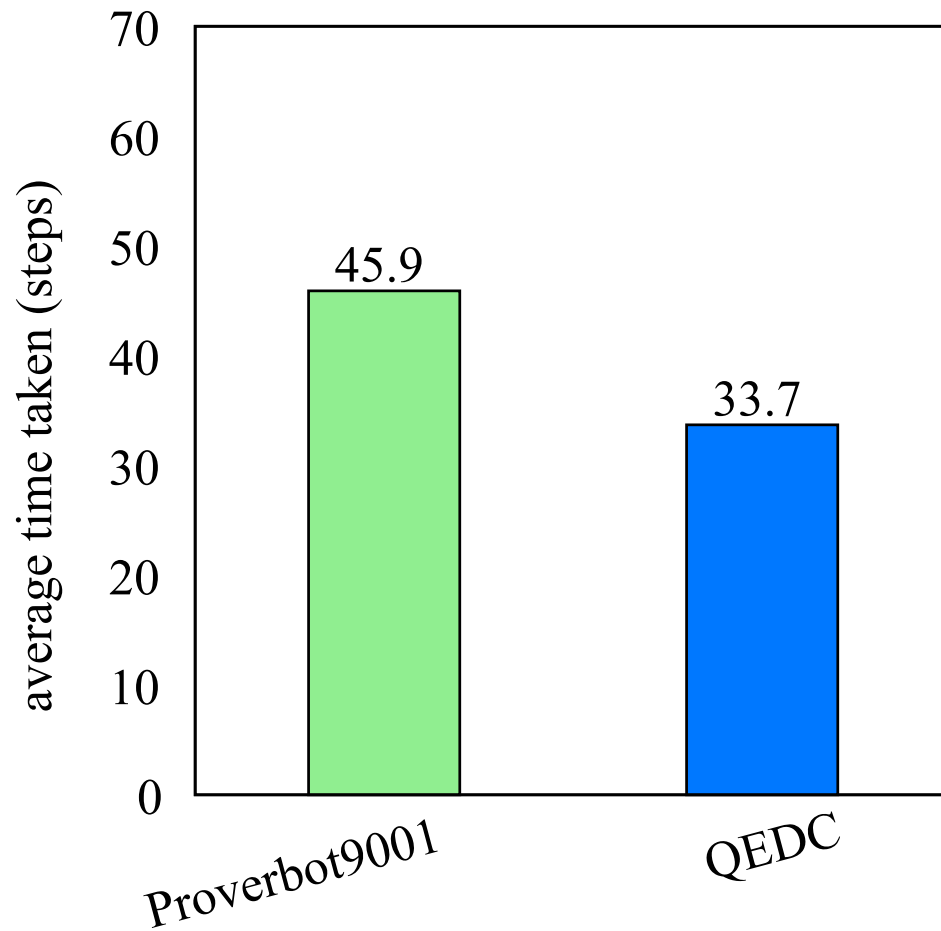QEDCartographer, except without state scoring-based search

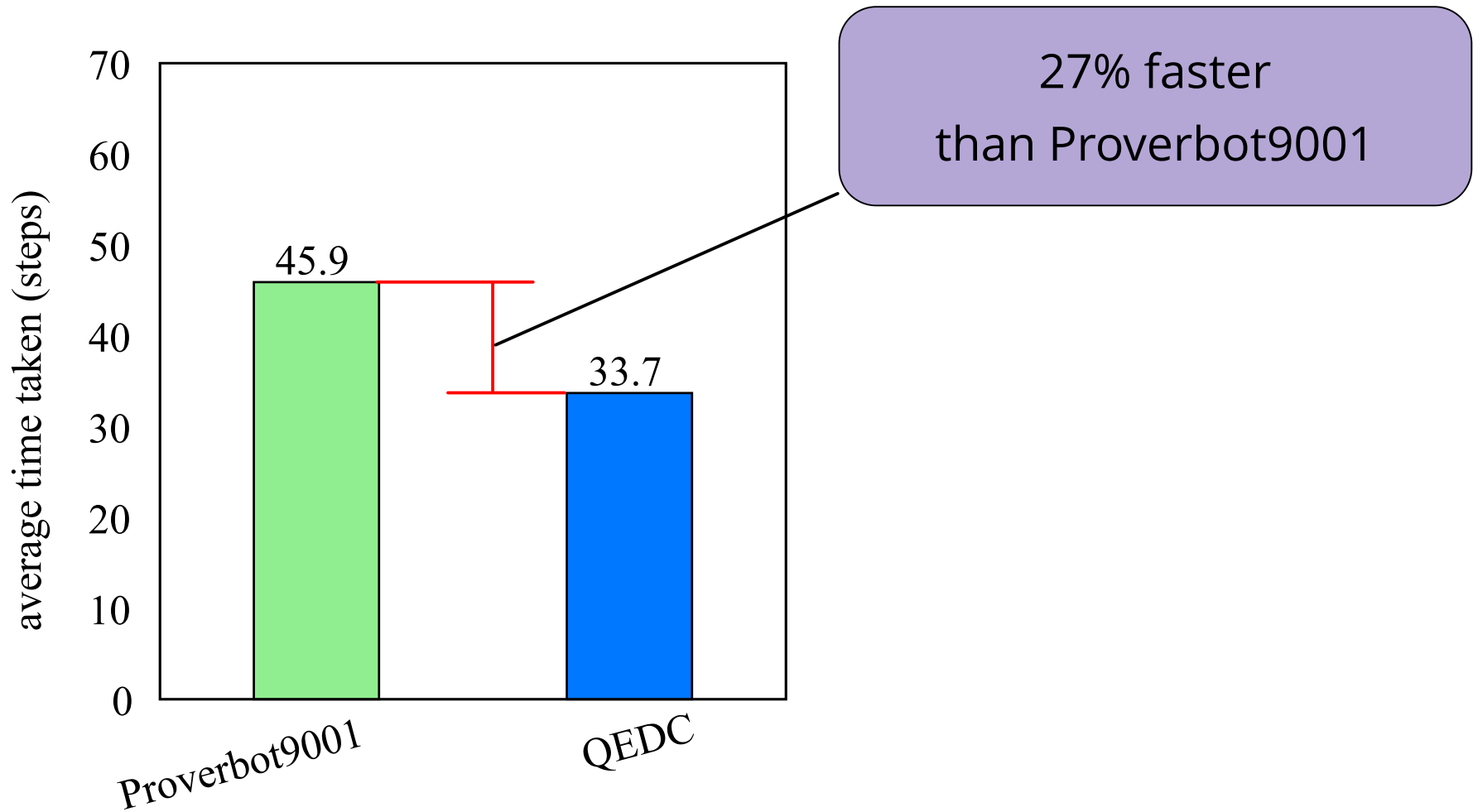Proverbot
9 0 0 1

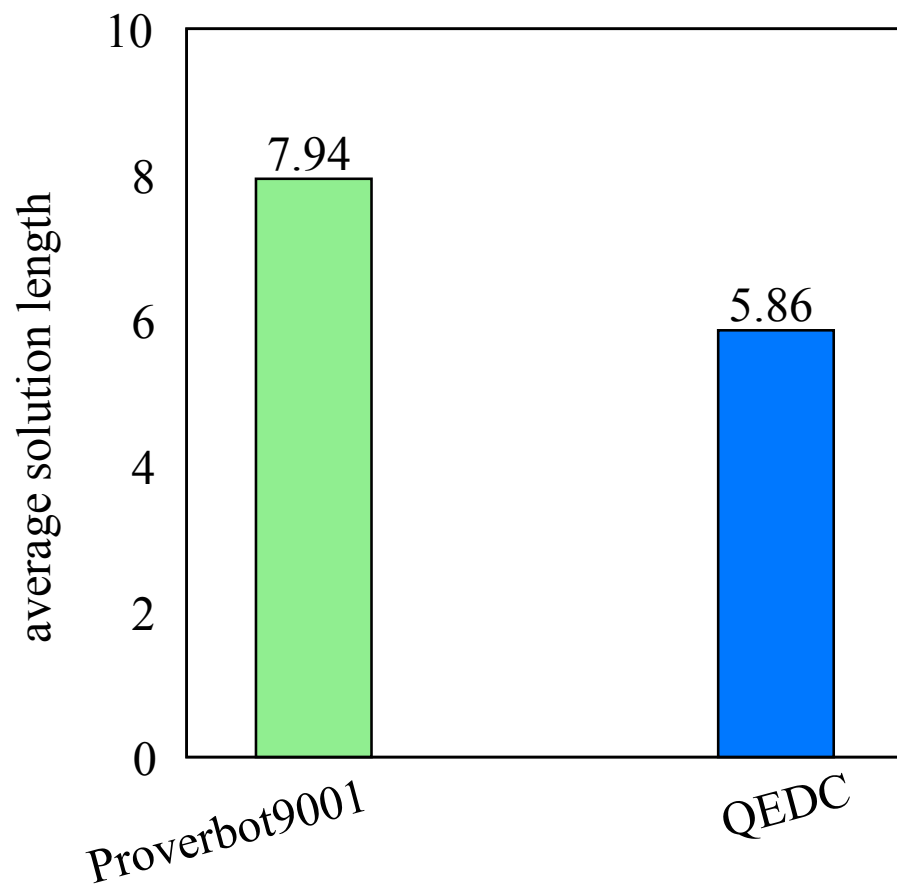# Baseline: Proverbot9001 (updated)

QEDCartographer, except without state
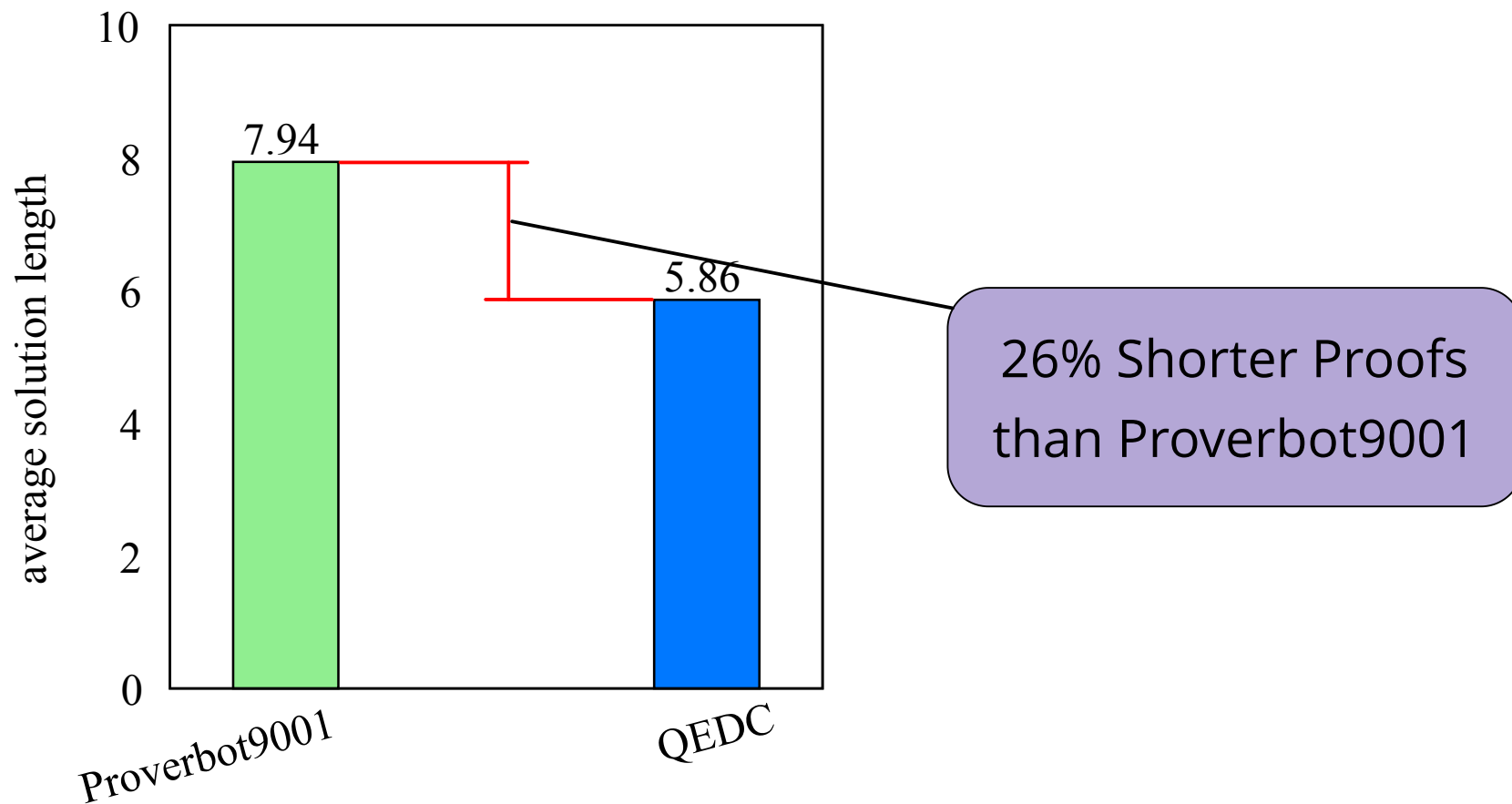scoring-based search
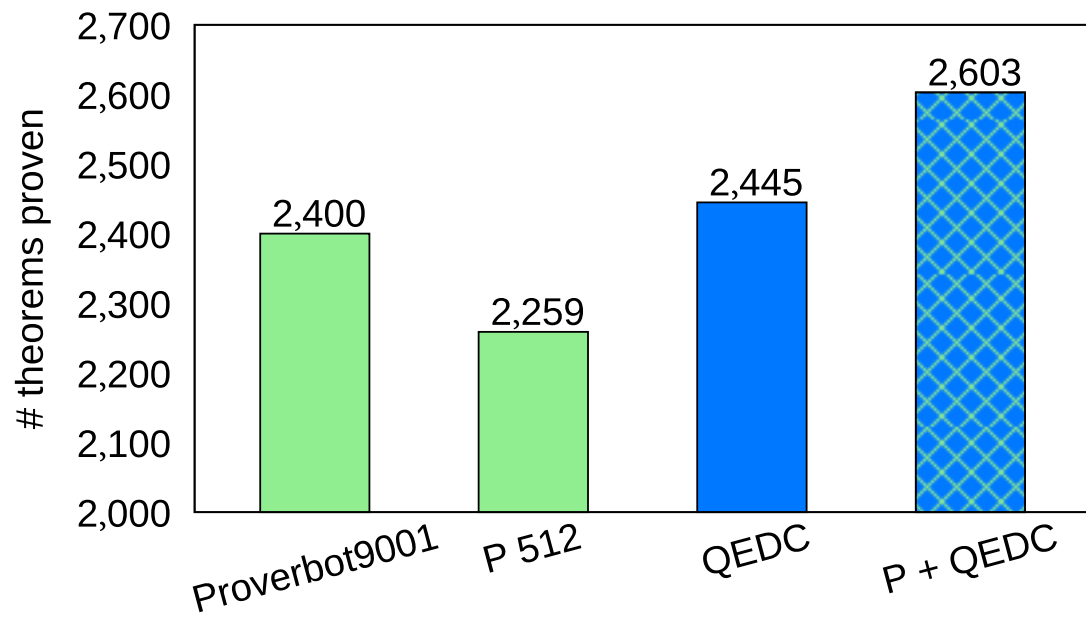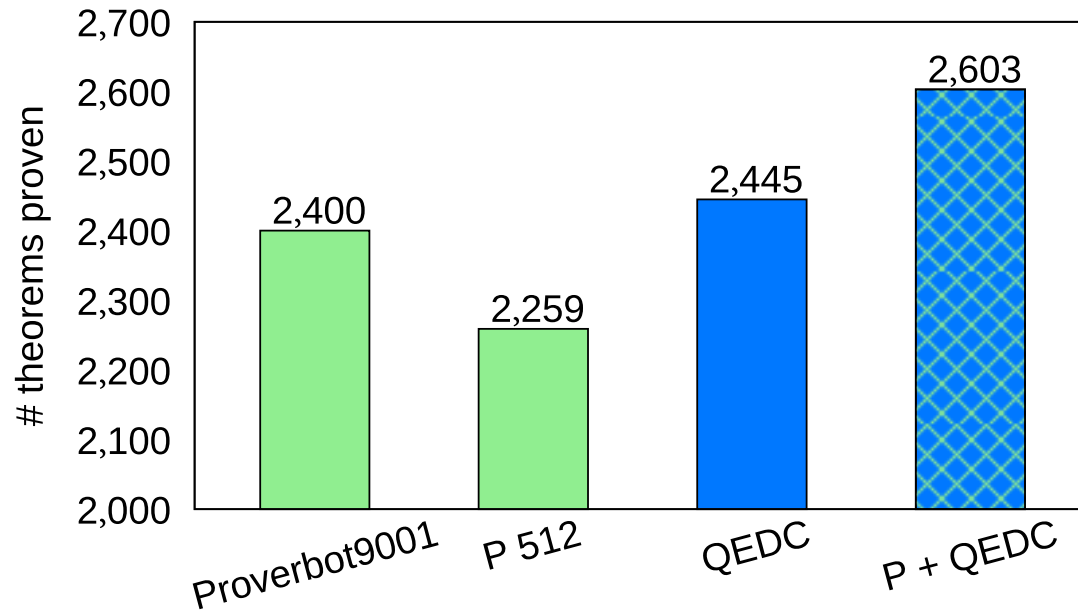
Uses a variant of depth-first search instead


Proverbot 9001

Proves slightly more theorems, and proves complementary theorems

# Automating Formal Verification with Reward-Free Reinforcement Learning

Uses a new V-value equation for branching goal structure

Makes producing verified-correct code easier and faster

Preprint available at alexsanchezstern.com