



Safe Reinforcement Learning via Formal Methods

André Platzer
Carnegie Mellon University
Joint work with Nathan Fulton



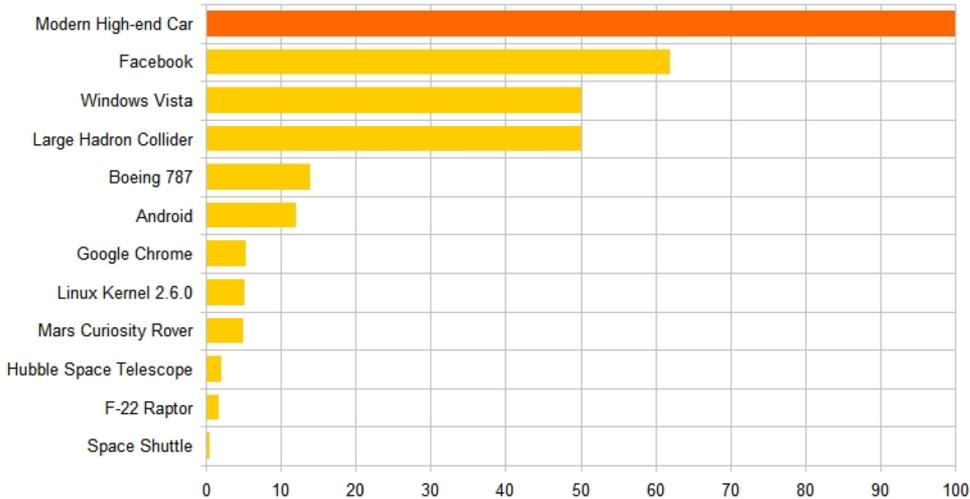
Safety-Critical Systems



"How can we provide people with cyber-physical systems they can bet their lives on?" - Jeannette Wing

Safety-Critical Systems

Software Size (million Lines of Code)



"How can we provide people with cyber-physical systems they can bet their lives on?" - Jeannette Wing

This Talk

Ensure the safety of Autonomous Cyber-Physical Systems.

Best of both worlds: learning together with CPS safety

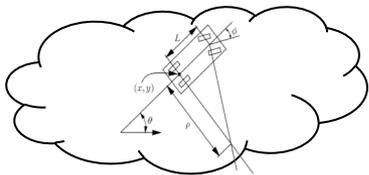
- Flexibility of learning
- Guarantees of CPS formal methods

Diametrically opposed: flexibility+adaptability versus predictability+simplicity

1. Cyber-Physical Systems with **Differential Dynamic Logic**
2. **Sandboxed reinforcement learning** is provably safe

Model-Based Verification

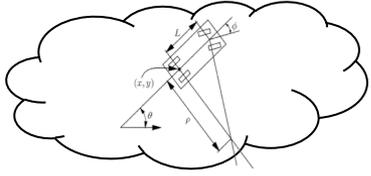
Reinforcement Learning



ϕ

Model-Based Verification

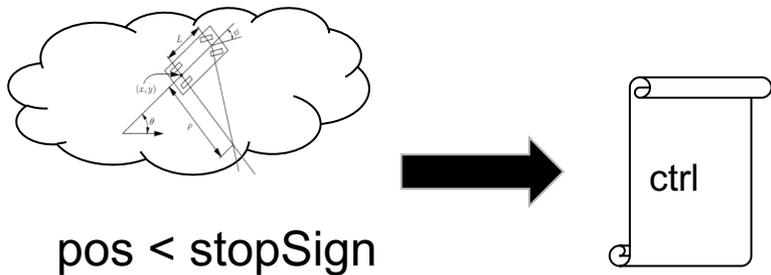
Reinforcement Learning



pos < stopSign

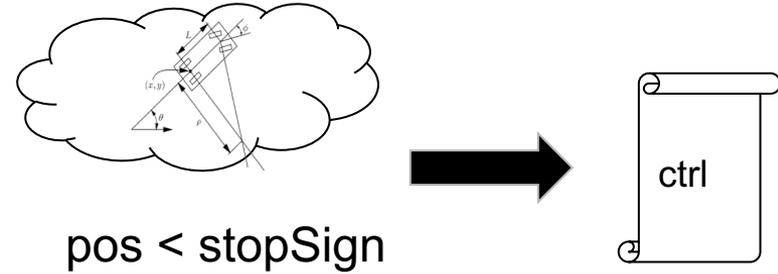
Model-Based Verification

Reinforcement Learning



Model-Based Verification

Reinforcement Learning



Approach: prove that control software achieves a specification with respect to a model of the physical system.

Model-Based Verification

Reinforcement Learning



Approach: prove that control software achieves a specification with respect to a model of the physical system.

Model-Based Verification



Reinforcement Learning

Benefits:

- Strong safety guarantees
- Automated analysis

Model-Based Verification



Reinforcement Learning

Benefits:

- Strong safety guarantees
- Automated analysis

Drawbacks:

- Control policies are typically non-deterministic: answers “what is safe”, not “what is useful”

Model-Based Verification



Reinforcement Learning

Benefits:

- Strong safety guarantees
- Automated analysis

Drawbacks:

- Control policies are typically non-deterministic: answers “what is safe”, not “what is useful”
- Assumes accurate model

Model-Based Verification



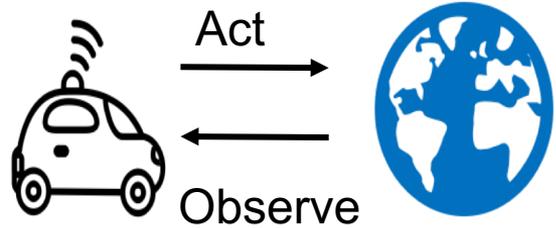
Benefits:

- Strong safety guarantees
- Automated analysis

Drawbacks:

- Control policies are typically non-deterministic: answers “what is safe”, not “what is useful”
- Assumes accurate model.

Reinforcement Learning



Model-Based Verification



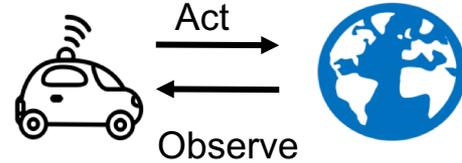
Benefits:

- Strong safety guarantees
- Automated analysis

Drawbacks:

- Control policies are typically non-deterministic: answers “what is safe”, not “what is useful”
- Assumes accurate model.

Reinforcement Learning



Benefits:

- No need for complete model
- Optimal (effective) policies

Model-Based Verification



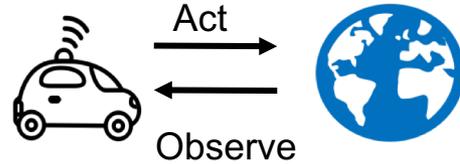
Benefits:

- Strong safety guarantees
- Automated analysis

Drawbacks:

- Control policies are typically non-deterministic: answers “what is safe”, not “what is useful”
- Assumes accurate model.

Reinforcement Learning



Benefits:

- No need for complete model
- Optimal (effective) policies

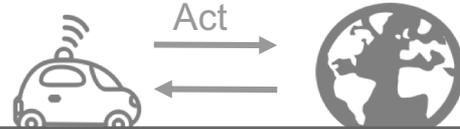
Drawbacks:

- No strong safety guarantees
- Proofs are obtained and checked by hand
- Formal proofs = decades-long proof development

Model-Based Verification



Reinforcement Learning



Goal: Provably correct reinforcement learning

Benefits

- Safety
- Accuracy

Model
S

Drawbacks

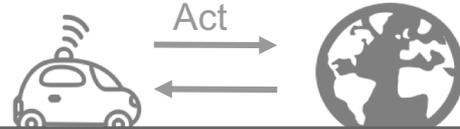
- Control policies are typically non-deterministic: answers “what is safe”, not “what is useful”
- Assumes accurate model

- No strong safety guarantees
- Proofs are obtained and checked by hand
- Formal proofs = decades-long proof development

Model-Based Verification



Reinforcement Learning



Goal: Provably correct reinforcement learning

- 1. Learn Safety**
- 2. Learn a Safe Policy**
- 3. Justify claims of safety**

Benefit

- Safety
- Accuracy

Drawback

- Control policies are typically non-deterministic: answers “what is safe”, not “what is useful”
- Assumes accurate model

- No strong safety guarantees
- Proofs are obtained and checked by hand
- Formal proofs = decades-long proof development

Model
s

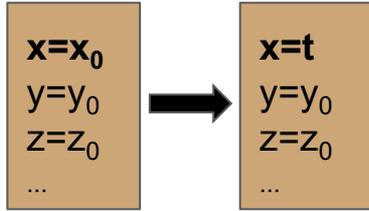
Part I: Differential Dynamic Logic

Trustworthy Proofs for Hybrid Systems



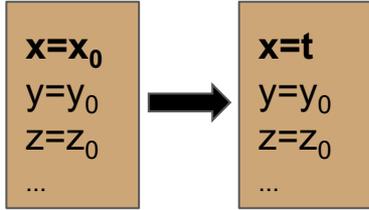
Hybrid Programs

$x := t$

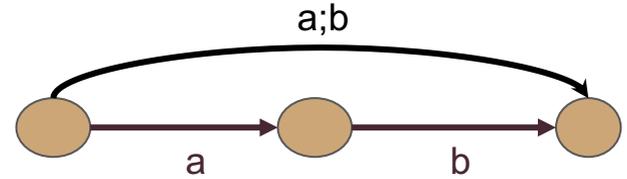


Hybrid Programs

$x := t$

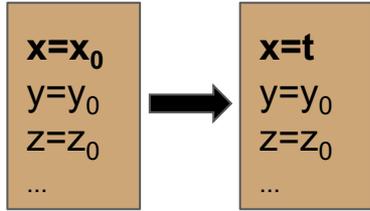


$a;b$



Hybrid Programs

$x := t$

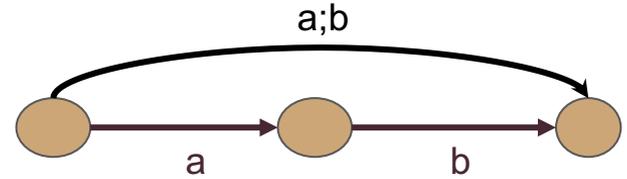


?P

If P is true: no change

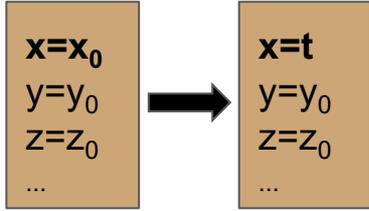
If P is false: terminate

$a;b$

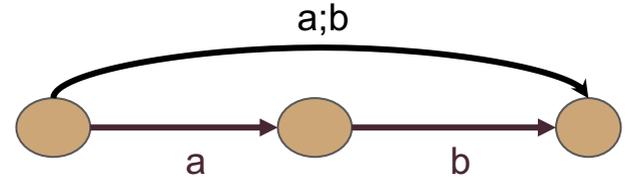


Hybrid Programs

$x := t$



$a;b$

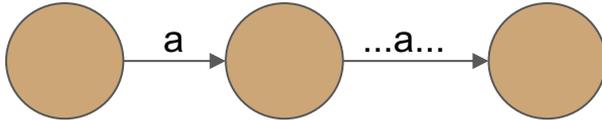


?P

If P is true: no change

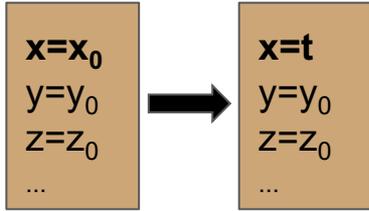
If P is false: terminate

a^*



Hybrid Programs

$x := t$

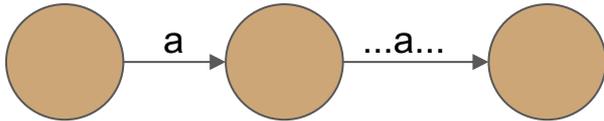


?P

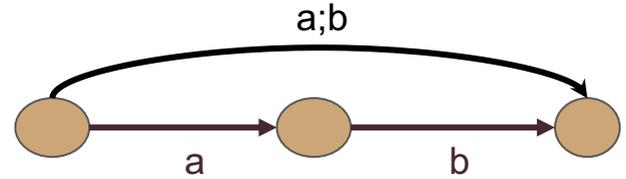
If P is true: no change

If P is false: terminate

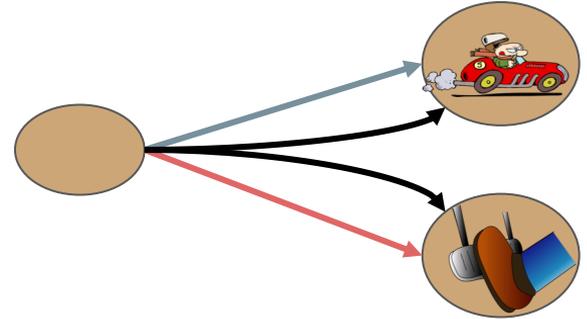
a^*



$a;b$

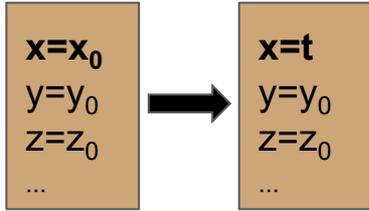


$a \cup b$



Hybrid Programs

$x := t$

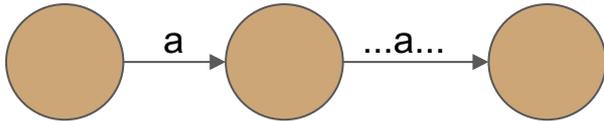


?P

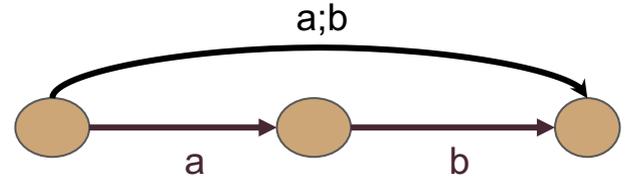
If P is true: no change

If P is false: terminate

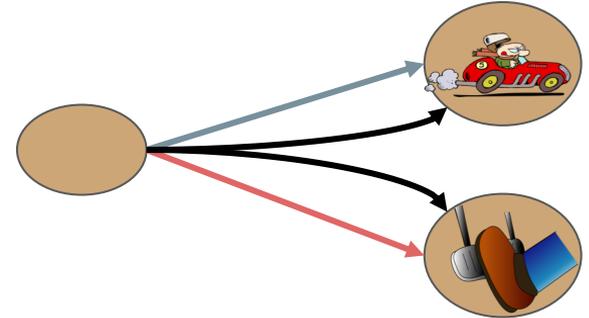
a^*



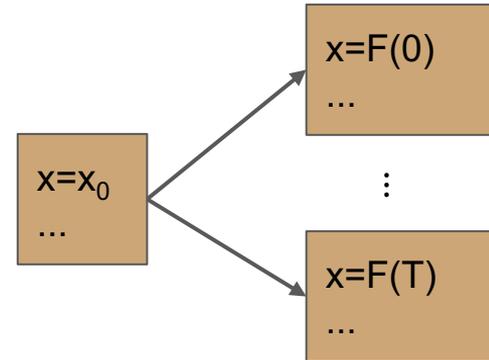
$a;b$



$a \cup b$



$x' = f$



Approaching a Stopped Car



Own Car



Stopped Car

Is this property true?

[

{ {accel U brake}; t:=0; {pos'=vel,vel'=accel,t'=1 & vel \geq 0 & t \leq T} }*

](pos \leq stoppedCarPos)

Approaching a Stopped Car

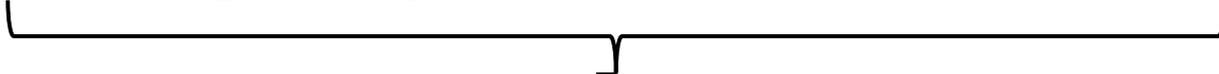


Own Car



Stopped Car

Assuming we only accelerate when it's safe to do so, is this property true?



[

{ {**accel**} U brake}; t:=0; {pos'=vel, vel'=accel, t'=1 & vel ≥ 0 & t ≤ T} }

](pos ≤ stoppedCarPos)

Approaching a Stopped Car



Own Car

`safeDistance(pos, vel, stoppedCarPos, B)`



Stopped Car

if we also assume the system is safe initially:

`safeDistance(pos, vel, stoppedCarPos, B) →`

[

{ {`accel` U `brake`}; t:=0; {pos'=vel, vel'=accel, t'=1 & vel ≥ 0 & t ≤ T} }*

](pos ≤ stoppedCarPos)

Approaching a Stopped Car



Own Car

`safeDistance(pos, vel, stoppedCarPos, B)`



Stopped Car

`safeDistance(pos, vel, stoppedCarPos, B) →`

[

{ {`accel` U `brake`}; t:=0; {pos', vel, vel'=accel, t'=1 & vel≥0 & t≤T} }*

](pos ≤ stoppedCarPos)



The Fundamental Question

Why would our program not work if we have a proof?

The Fundamental Question

Why would our program not work if we have a proof?

1. Was the proof correct?



The Fundamental Question

Why would our program not work if we have a *proof*?

1. Was the proof correct?
2. Was the model accurate enough?



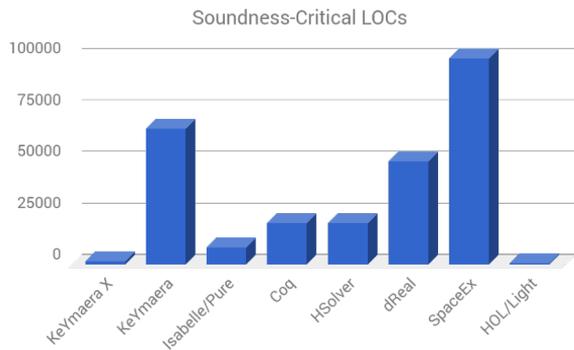
≠



The Fundamental Question

Why would our program not work if we have a *proof*?

1. Was the proof correct? **KeYmaera X** VERIFIED
2. Was the model accurate enough?



dl Tactic:

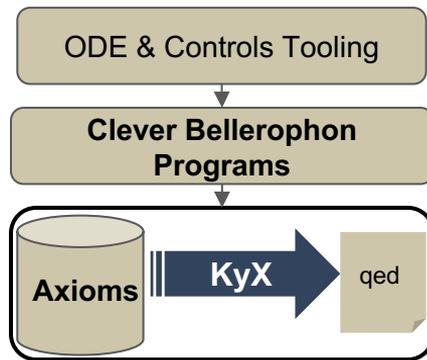
Side derivation:
 $(v \geq v_0 - gt)'$ ↔
 $\dots \leftrightarrow$
 $\dots \leftrightarrow$
 \dots
 $H = r_p \geq 0 \ \& \ r_a \geq 0$
 $\ \& \ g > 0 \ \& \ \dots$

DI Axiom:

$[\{x'=f\&Q\}]P \leftrightarrow ([?Q]P \leftarrow (Q \rightarrow [\{x'=f\&Q\}]P'))$

Example:

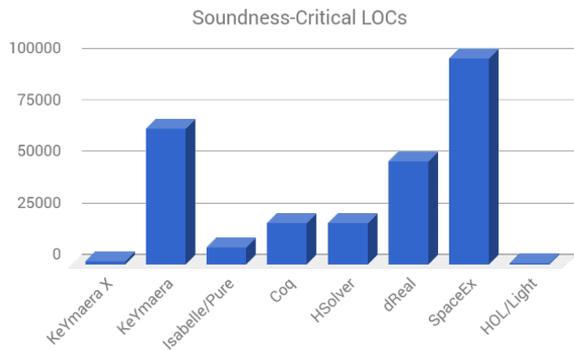
$[v' = r_p v^2 - g, t' = 1] v \geq v_0 - gt \leftrightarrow$
 $\dots \leftrightarrow$
 $[v' = r_p v^2 - g] [t' = 1] v' \geq -g * t' \leftrightarrow$
 $r_p v^2 - g \geq -g$
 $H \rightarrow r_p \geq 0$



The Fundamental Question

Why would our program not work if we have a *proof*?

1. Was the proof correct? **KeYmaera X** VERIFIED
2. Was the model accurate enough? **Safe RL**



dl Tactic:

Side derivation:
 $(v \geq v_0 - gt)'$ ↔
 $\dots \leftrightarrow$
 $\dots \leftrightarrow$
 \dots

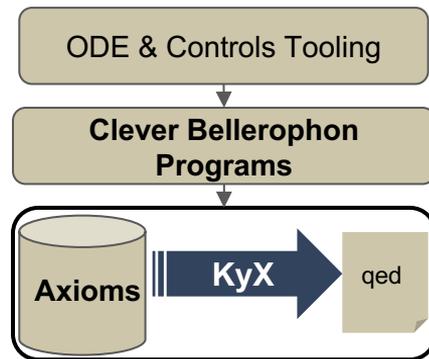
$H=r_p \geq 0 \ \& \ r_a \geq 0$
 $\ \& \ g > 0 \ \& \ \dots$

DI Axiom:

$[\{x'=f\&Q\}]P \leftrightarrow ([?Q]P \leftarrow (Q \rightarrow [\{x'=f\&Q\}]P'))$

Example:

$[v' = r_p v^2 - g, t' = 1] v \geq v_0 - gt \quad \leftrightarrow$
 $\dots \quad \leftrightarrow$
 $[v' := r_p v^2 - g] [t' := 1] v' \geq -g * t' \quad \leftrightarrow$
 $r_p v^2 - g \geq -g \quad \leftrightarrow$
 $H \rightarrow r_p \geq 0$



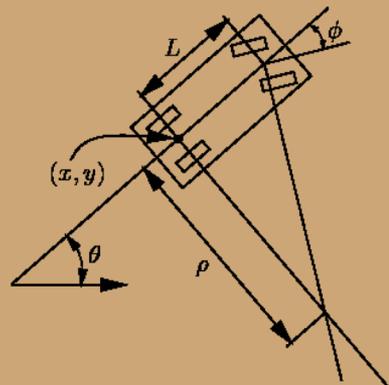
Part II: Justified Speculative Control

Safe reinforcement learning in partially
modeled environments

AAAI 2018



≠



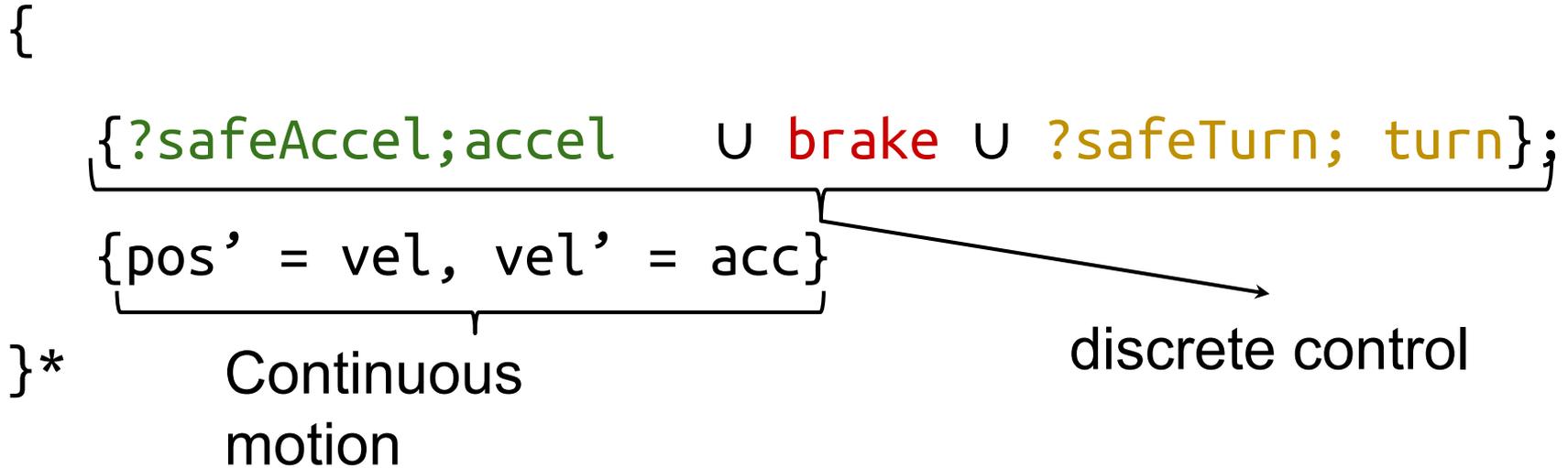
Model-Based Verification

Accurate, analyzable models often exist!

```
{  
  {?safeAccel; accel    U brake U ?safeTurn; turn};  
  {pos' = vel, vel' = acc}  
}*
```

Model-Based Verification

Accurate, analyzable models often exist!



Model-Based Verification

Accurate, analyzable models often exist!

{

{?safeAccel; accel U brake U ?safeTurn; turn};

{pos' = vel, vel' = acc}

}*

Continuous
motion

discrete, **non-deterministic**
control

Model-Based Verification

Accurate, analyzable models often exist!

```
init → [{  
    { ?safeAccel; accel  U brake U ?safeTurn; turn};  
    {pos' = vel, vel' = acc}  
}*]pos < stopSign
```

Model-Based Verification

Accurate, analyzable models often exist!

formal verification gives strong safety guarantees

```
init → [{  
    { ?safeAccel accel ∪ brake ∪ ?safeTurn; turn};  
    {pos' = vel, vel' = acc}  
}*]pos < stopSign
```



Model-Based Verification

Accurate, analyzable models often exist!

formal verification gives strong safety guarantees



=

- **Computer-checked proofs of safety specification.**

Model-Based Verification

Accurate, analyzable models often exist!

formal verification gives strong safety guarantees



=

- **Computer-checked proofs of safety specification**
- **Formal proofs mapping model to runtime monitors**

Model-Based Verification Isn't Enough

Perfect, analyzable models don't exist!

Model-Based Verification Isn't Enough

Perfect, analyzable models don't exist!

How to implement?

{

{ ?safeAccel; accel U brake U ?safeTurn; turn};

{pos' = vel, vel' = acc}

}*

Only accurate sometimes

Model-Based Verification Isn't Enough

Perfect, analyzable models don't exist!

How to implement?

{

{ ?safeAccel; accel U brake U ?safeTurn; turn};

{dx'=w*y, dy'=-w*x, ...}

}*

Only accurate sometimes

Safe RL Contribution

Justified Speculative Control is an approach toward provably safe reinforcement learning that:

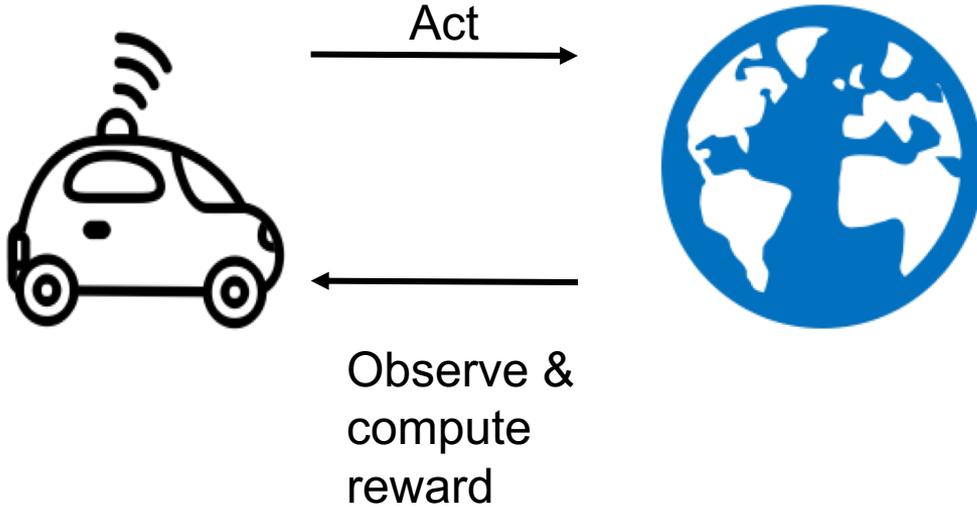
1. learns to resolve non-determinism without sacrificing formal safety results

Safe RL Contribution

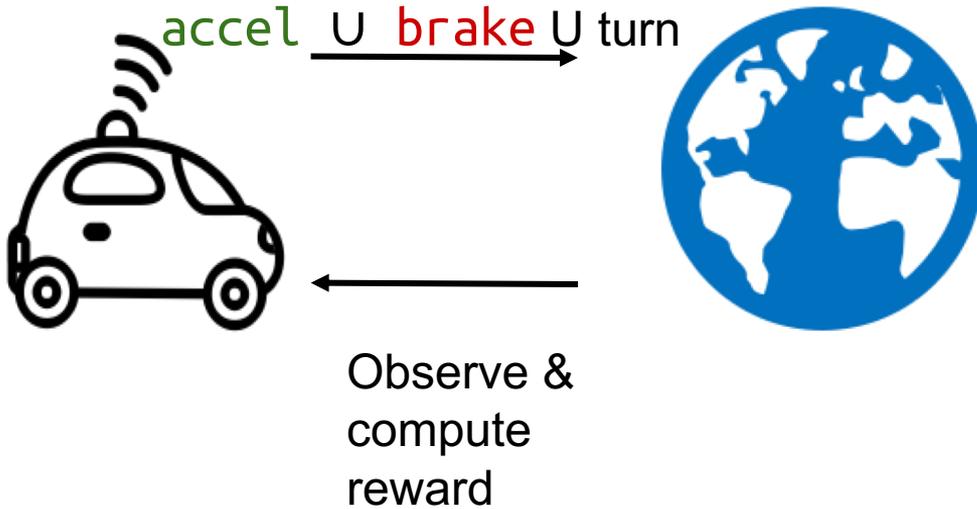
Justified Speculative Control is an approach toward provably safe reinforcement learning that:

1. learns to resolve non-determinism without sacrificing formal safety results
2. allows and directs speculation whenever model mismatches occur

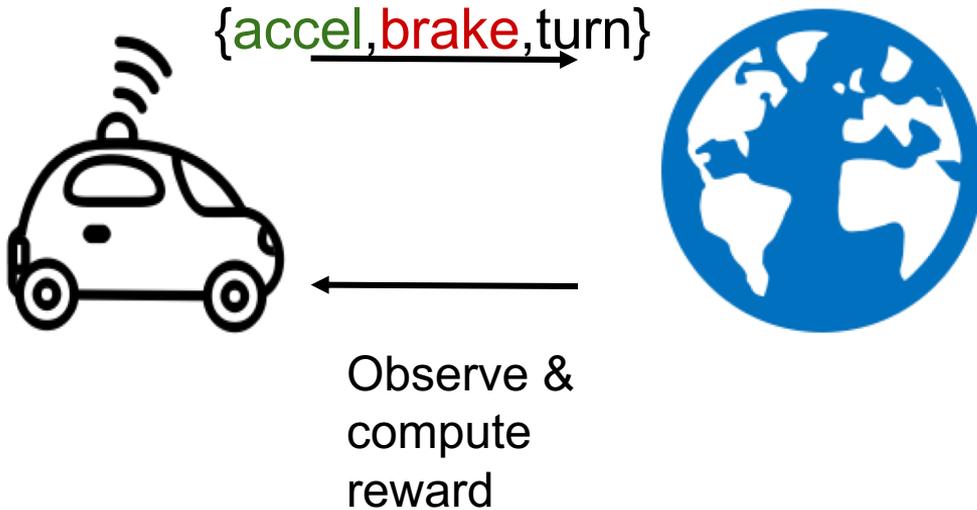
Learning to Resolve Non-determinism



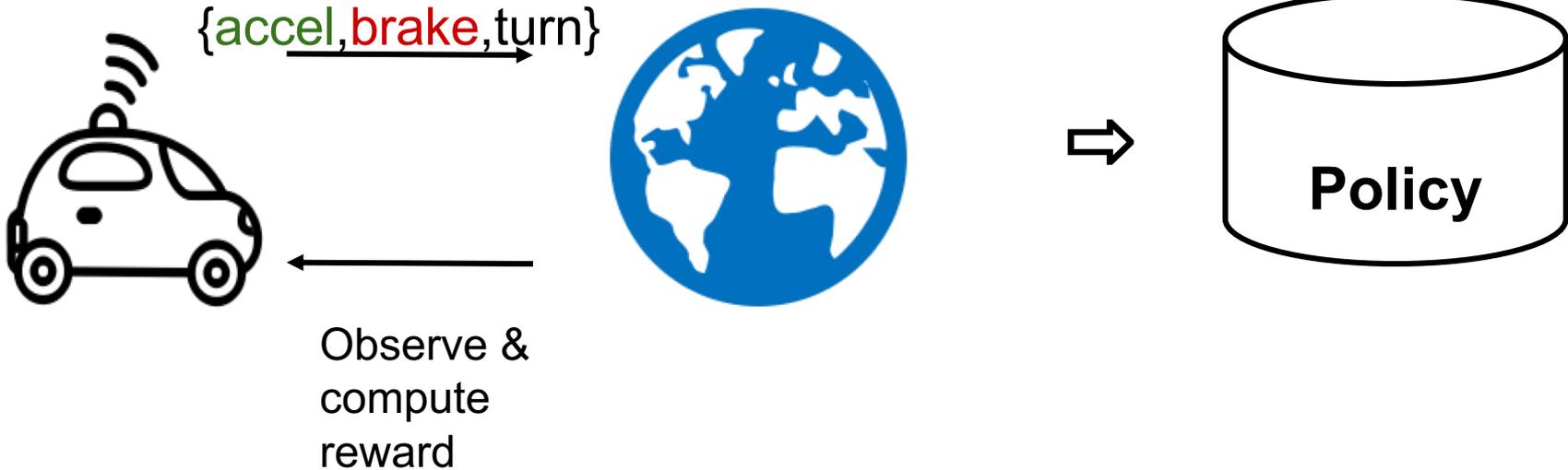
Learning to Resolve Non-determinism



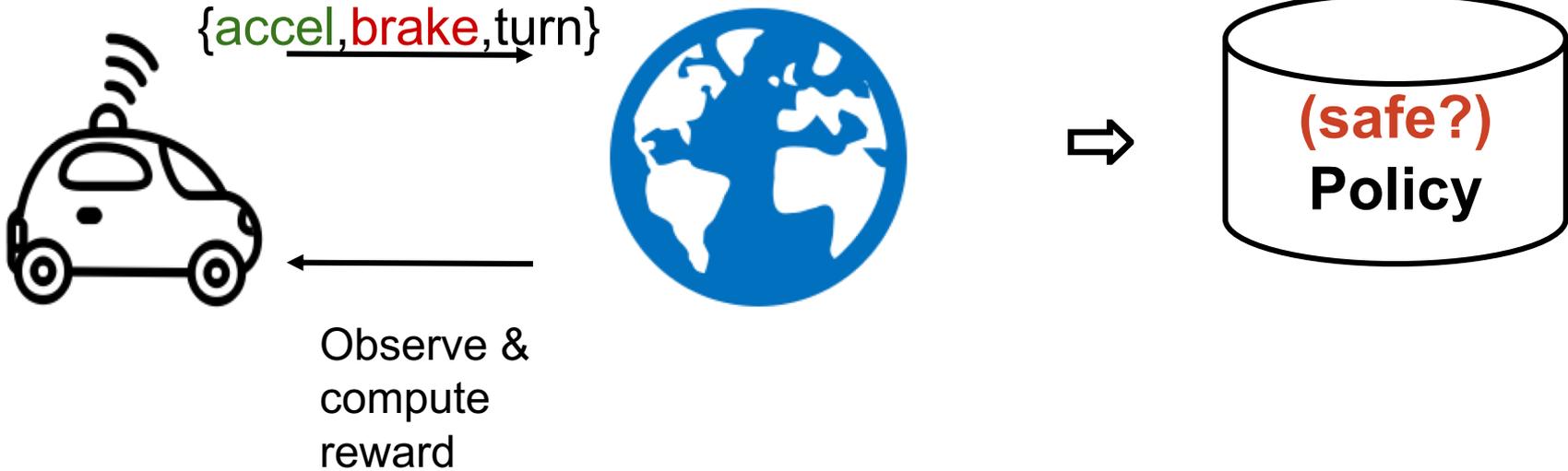
Learning to Resolve Non-determinism



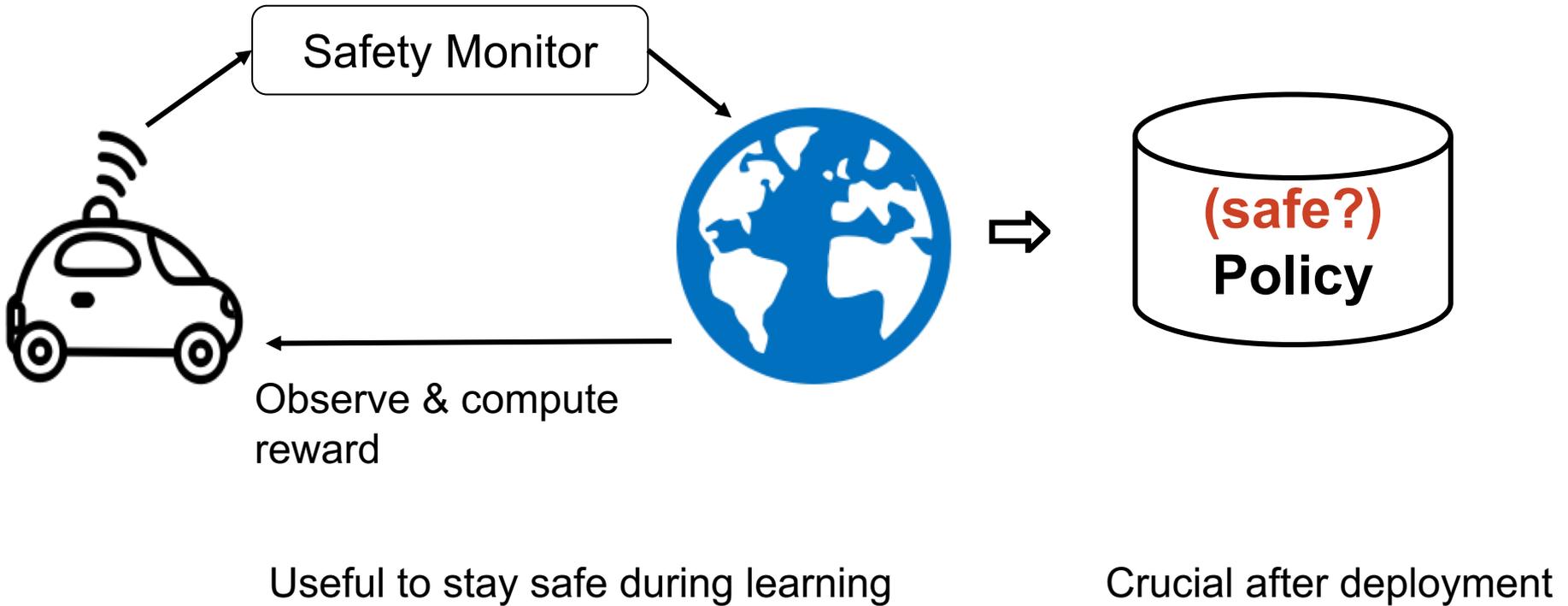
Learning to Resolve Non-determinism



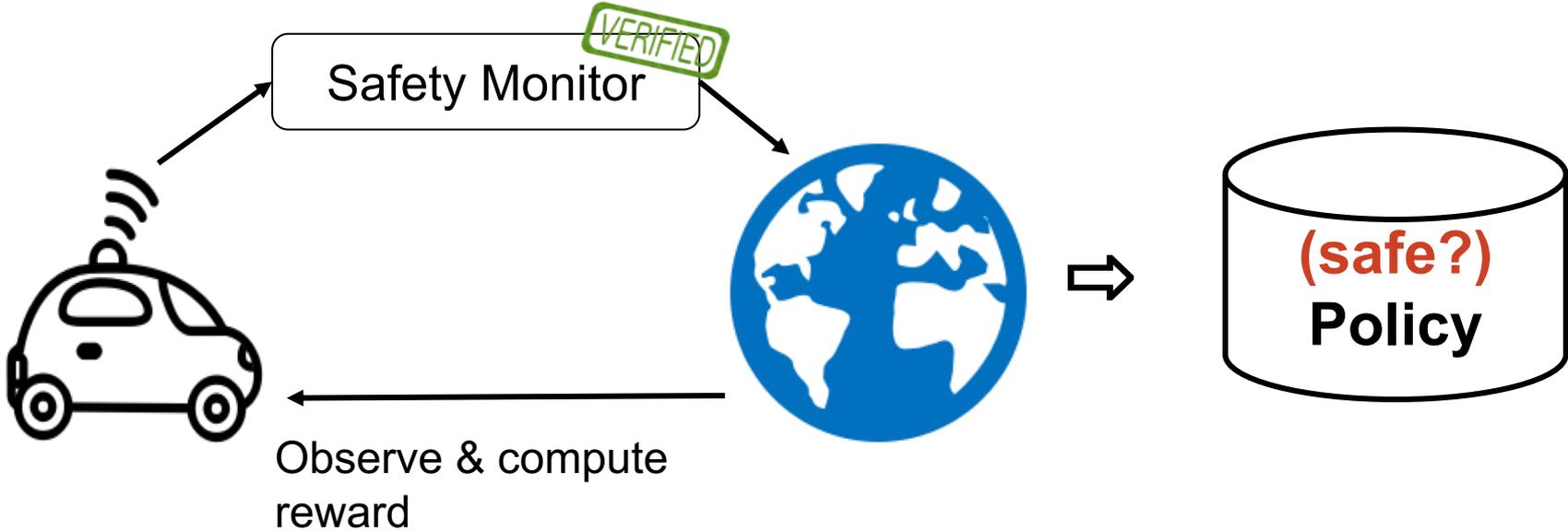
Learning to Resolve Non-determinism



Learning to **Safely** Resolve Non-determinism

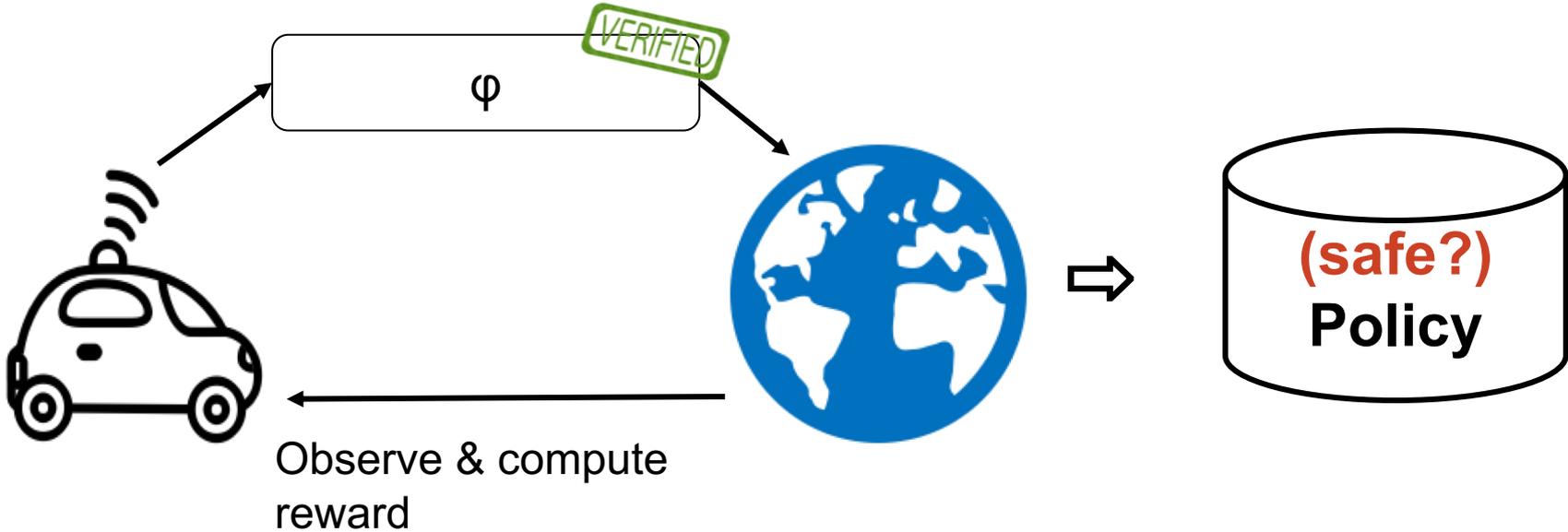


Learning to **Safely** Resolve Non-determinism



VERIFIED \neq "Trust Me"

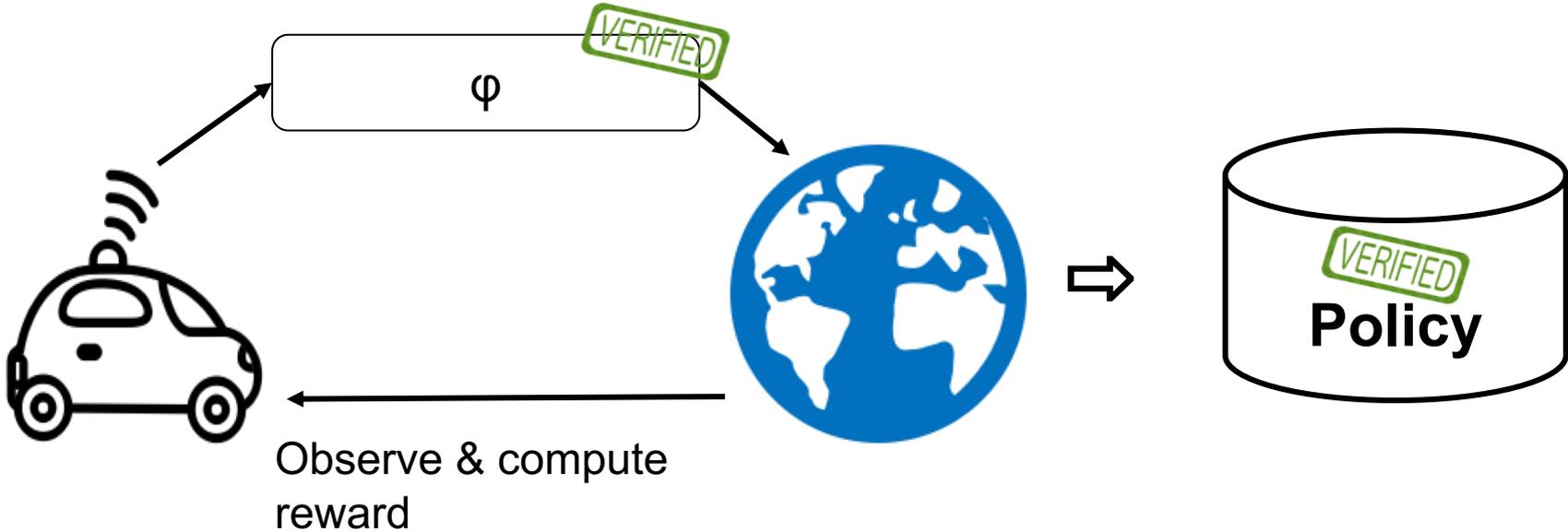
Learning to **Safely** Resolve Non-determinism



Use a theorem prover to prove:

$(\text{init} \rightarrow [\{ \{ \text{accel} \cup \text{brake} \}; \text{ODEs} \}^*] (\text{safe})) \quad \varphi$

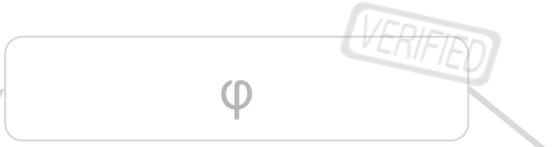
Learning to **Safely** Resolve Non-determinism



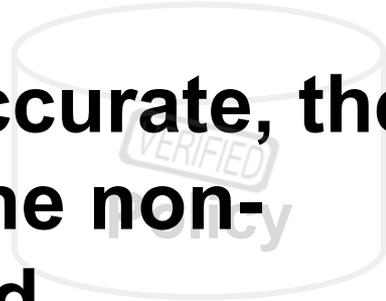
Use a theorem prover to prove:

$(\text{init} \rightarrow [\{ \{ \text{accel} \cup \text{brake} \}; \text{ODEs} \}^*] (\text{safe})) \quad \varphi$

Learning to **Safely** Resolve Non-determinism



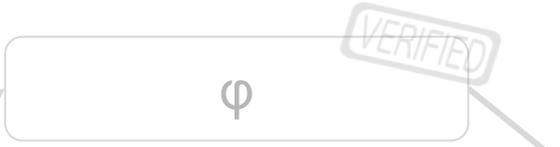
Main Theorem: If the ODEs are accurate, then our formal proofs transfer from the non-deterministic model to the learned **(deterministic) policy**



Use a theorem prover to prove:

$$(\text{init} \rightarrow [\{ \{ \text{accel} \cup \text{brake} \}; \text{ODEs} \}^*] (\text{safe})) \quad \phi$$

Learning to **Safely** Resolve Non-determinism

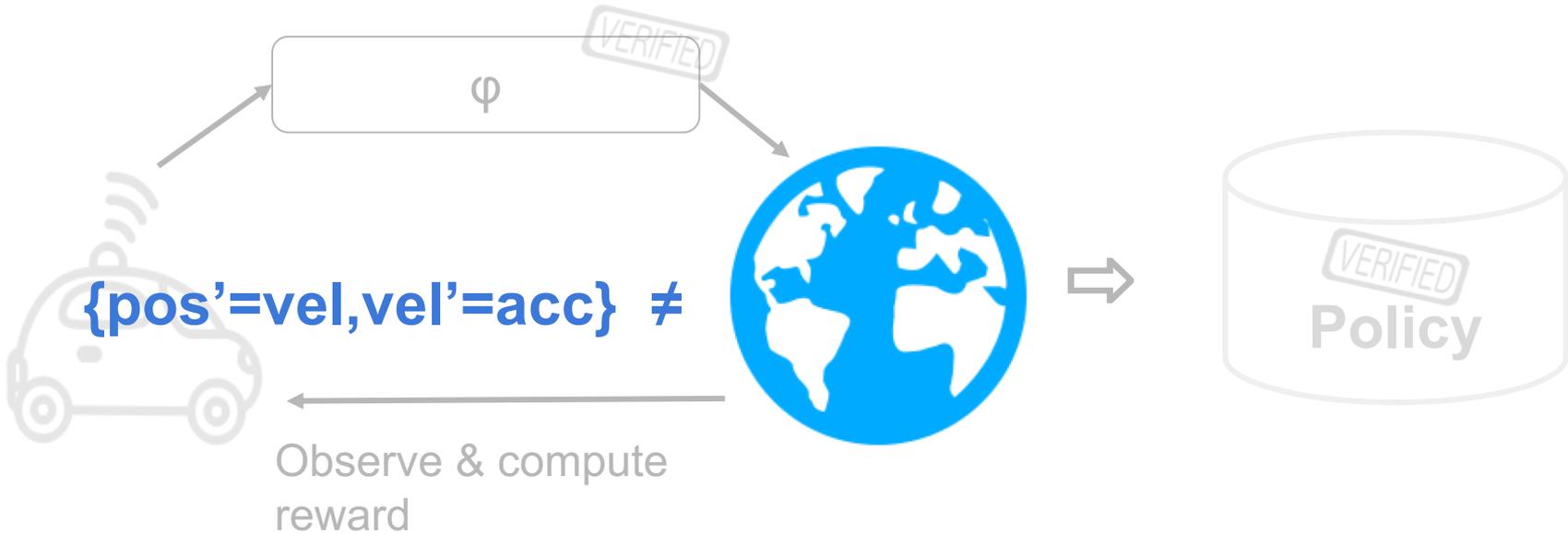


Main Theorem: If the ODEs are accurate, then our formal proofs transfer from the non-deterministic model to the learned (deterministic) policy via the model monitor.

Use a theorem prover to prove:

$$(\text{init} \rightarrow [\{ \{ \text{accel} \cup \text{brake} \}; \text{ODEs} \}^*] (\text{safe})) \quad \phi$$

What about the physical model?

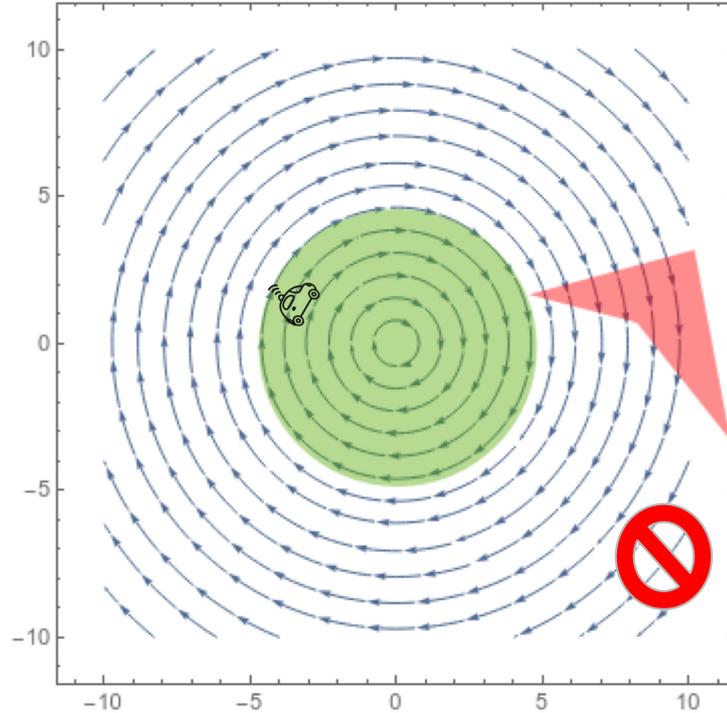
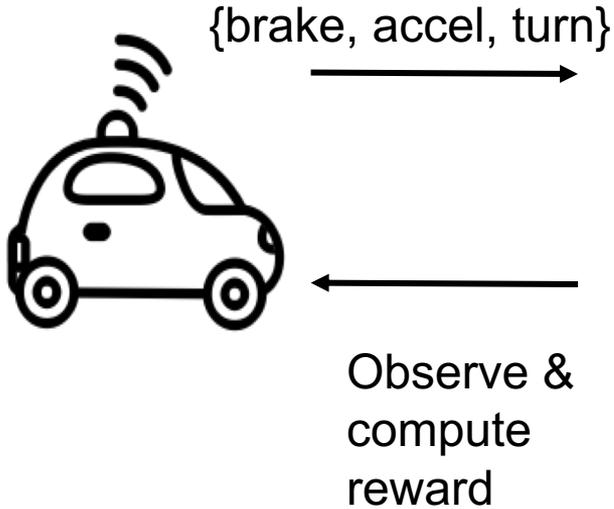


Use a theorem prover to prove:

$(\text{init} \rightarrow [\{ \{ \text{accel} \cup \text{brake} \}; \text{ODEs} \}^*] (\text{safe}))$

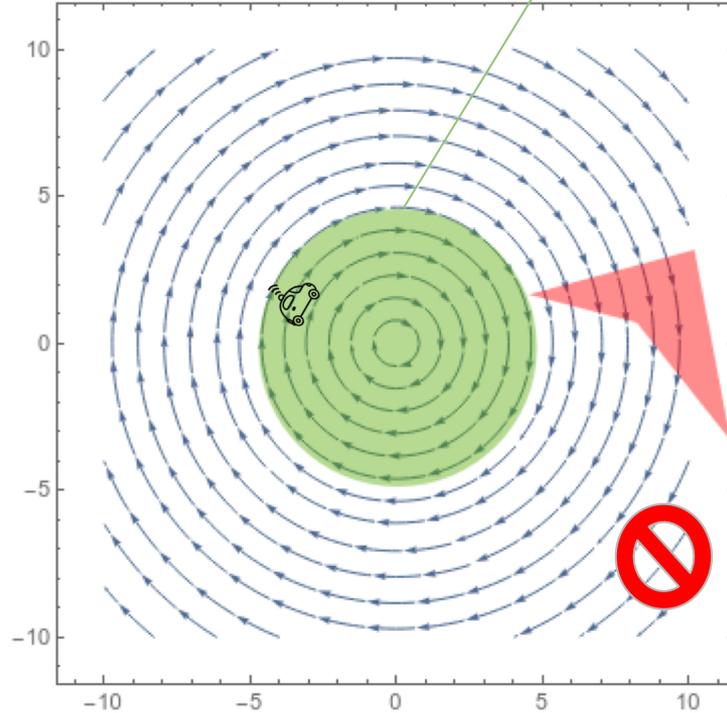
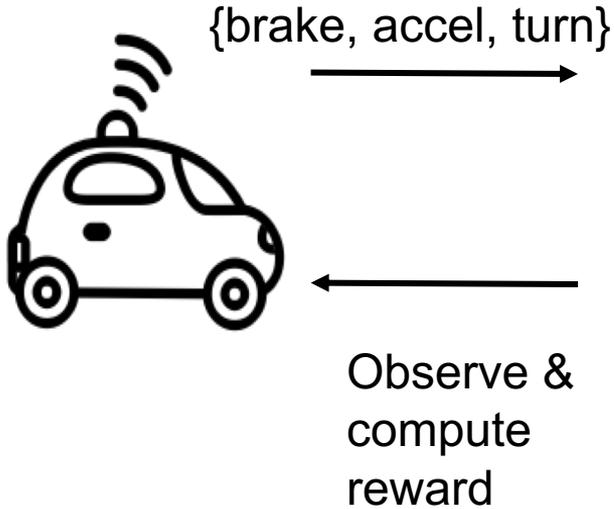
φ

What About the Physical Model?



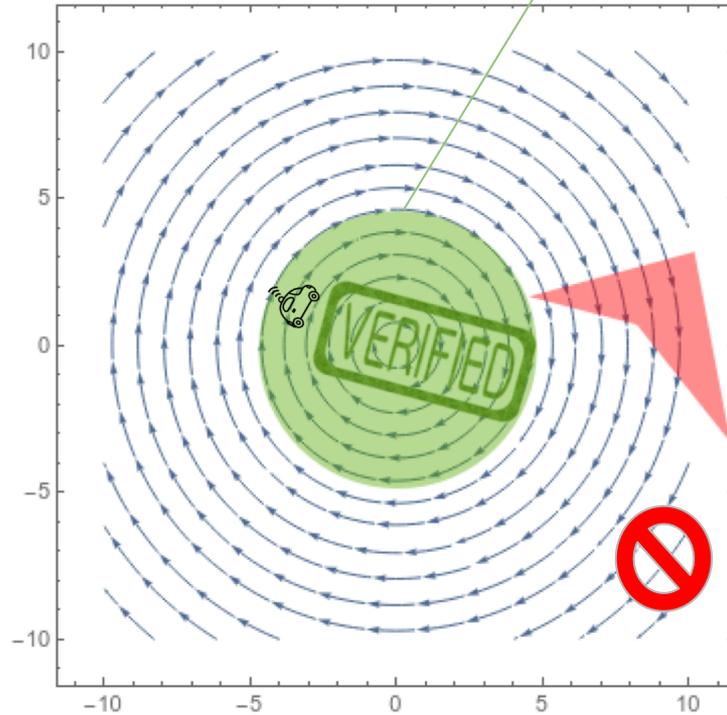
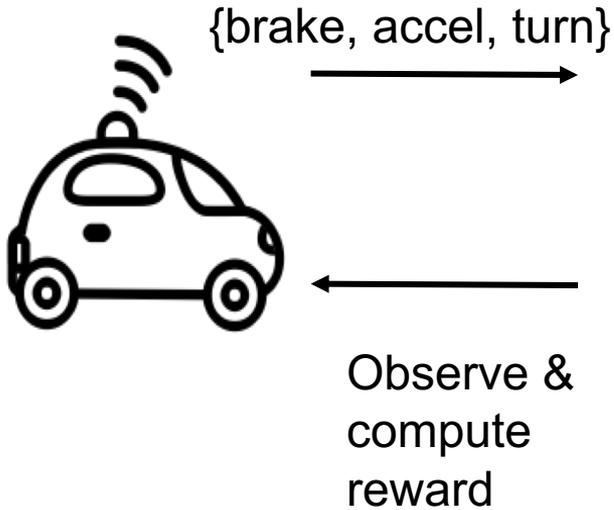
What About the Physical Model?

Model is accurate.

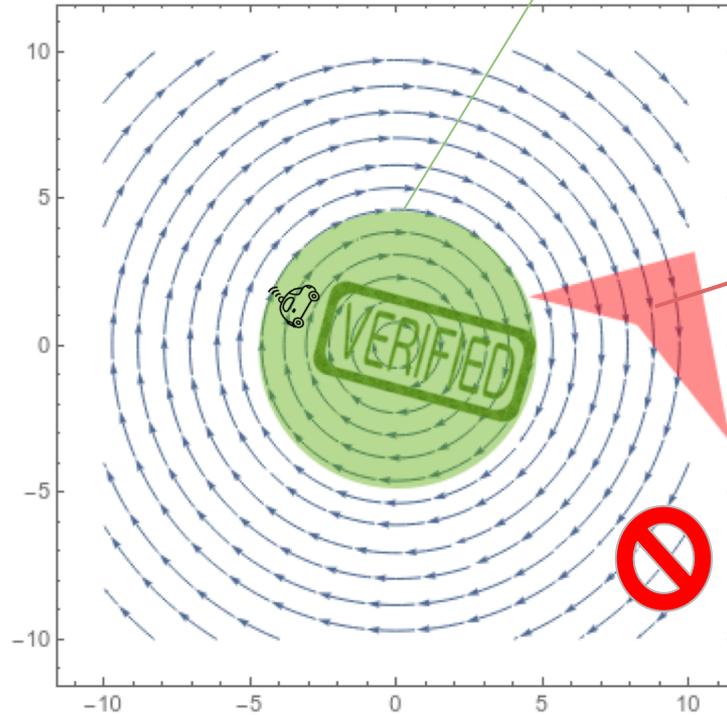
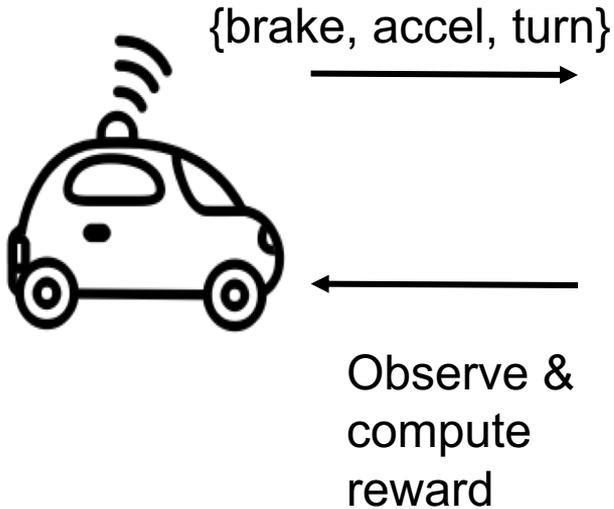


What About the Physical Model?

Model is accurate.



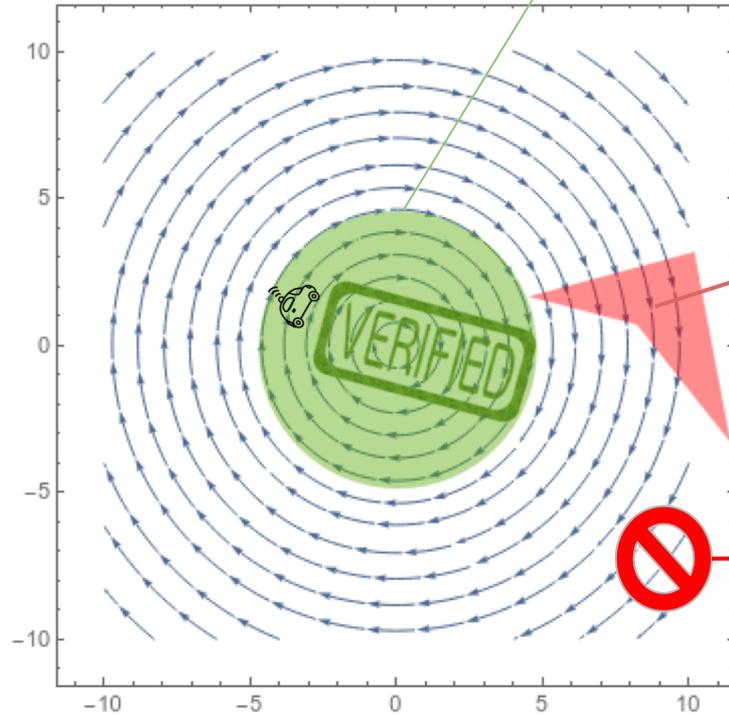
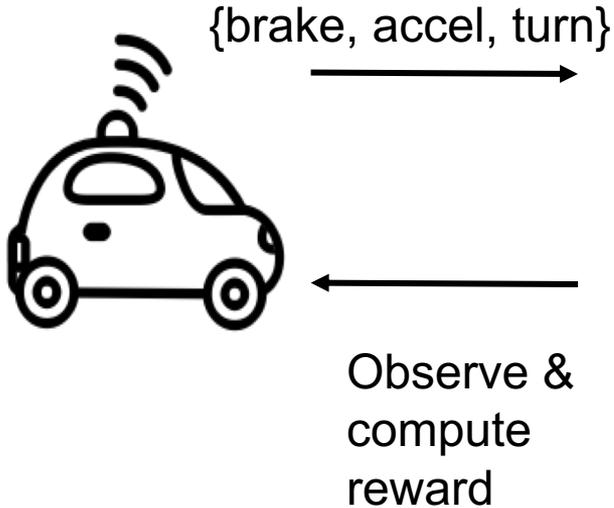
What About the Physical Model?



Model is accurate.

Model is inaccurate

What About the Physical Model?

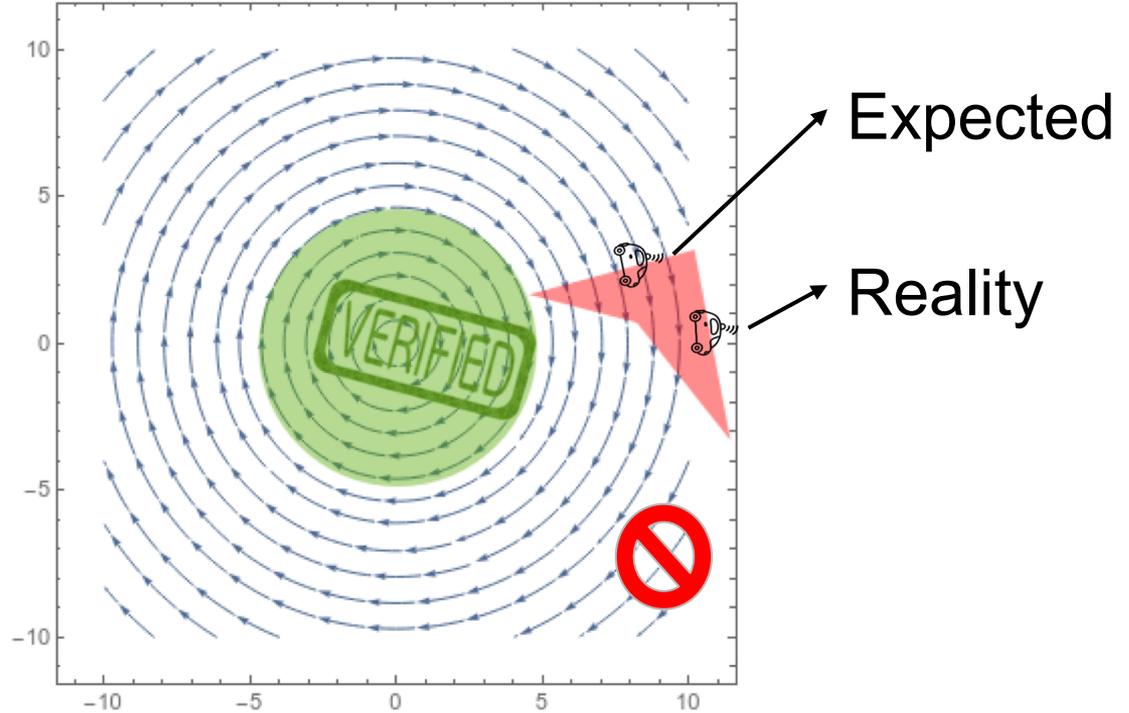
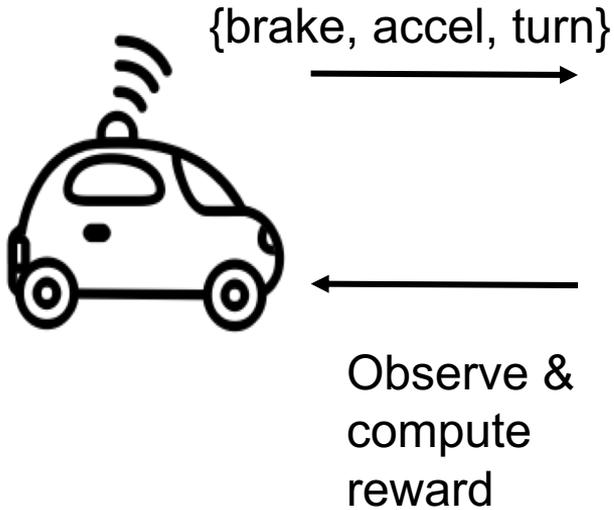


Model is accurate.

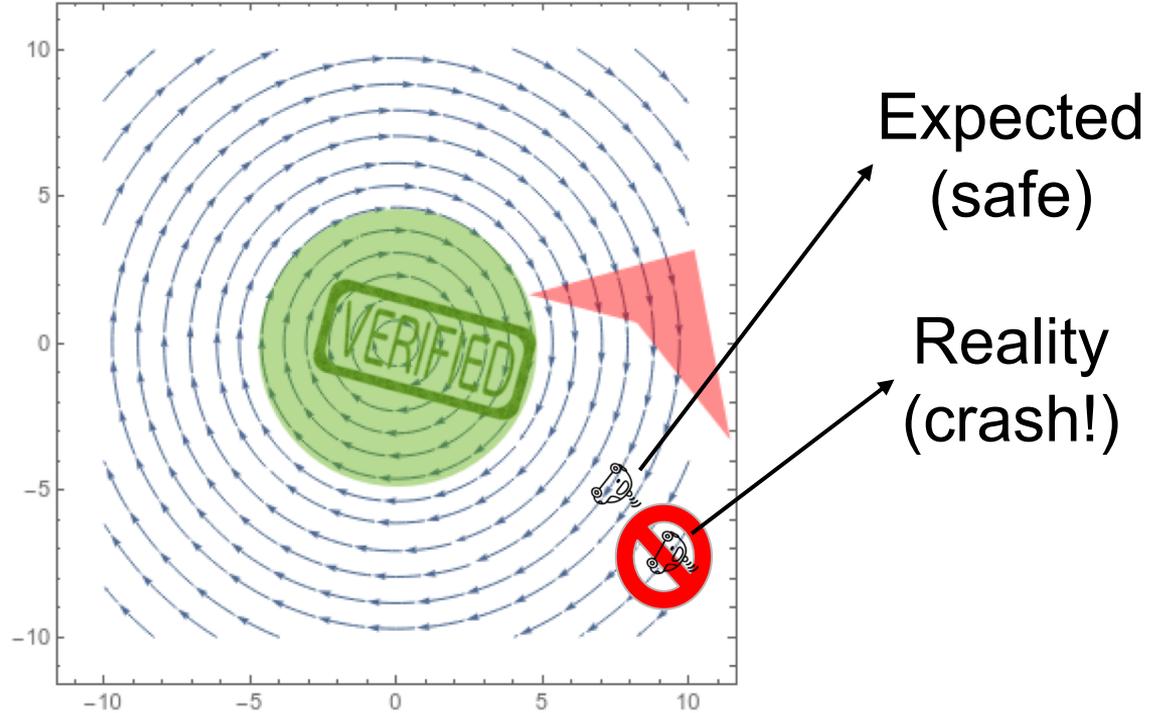
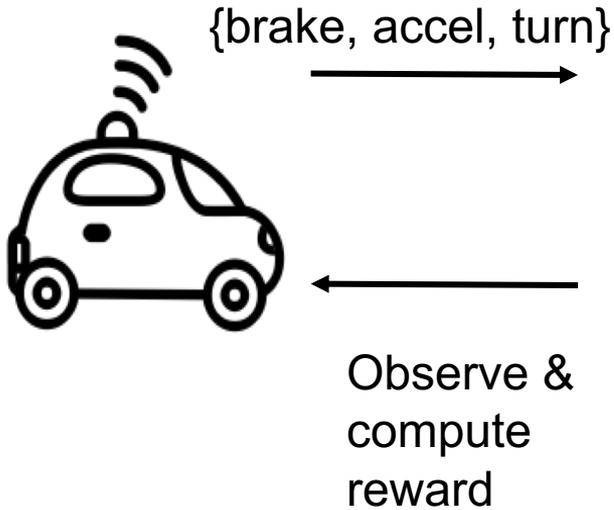
Model is inaccurate

Obstacle!

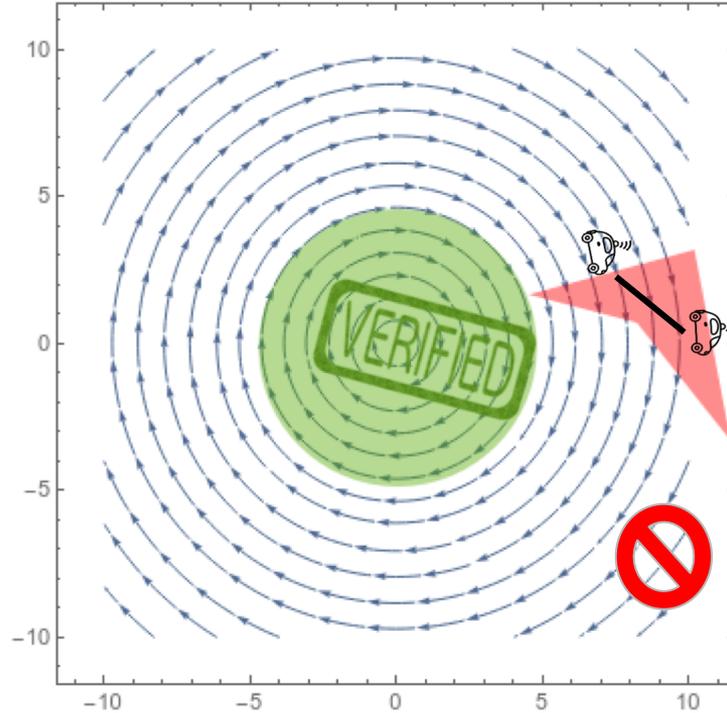
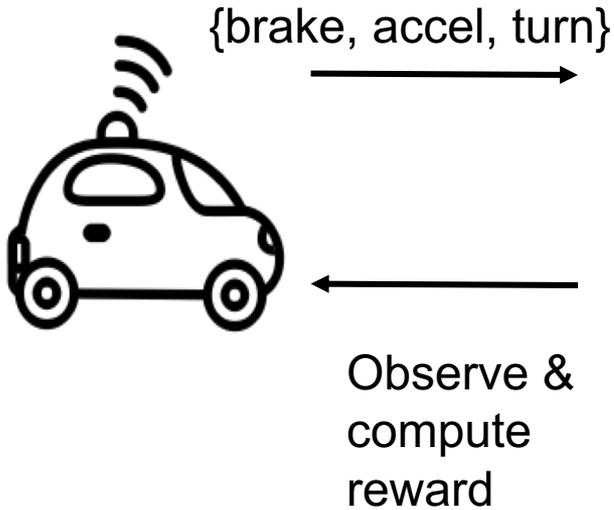
What About the Physical Model?



Speculation is Justified



Leveraging Verification Results to Learn Better

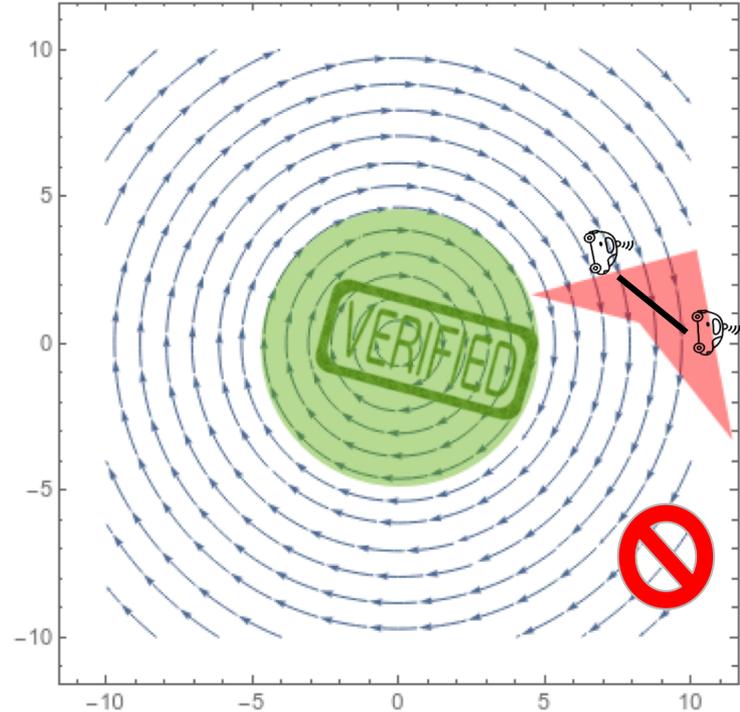


Use a real-valued version of the model monitor as a reward signal

Safe RL: How?

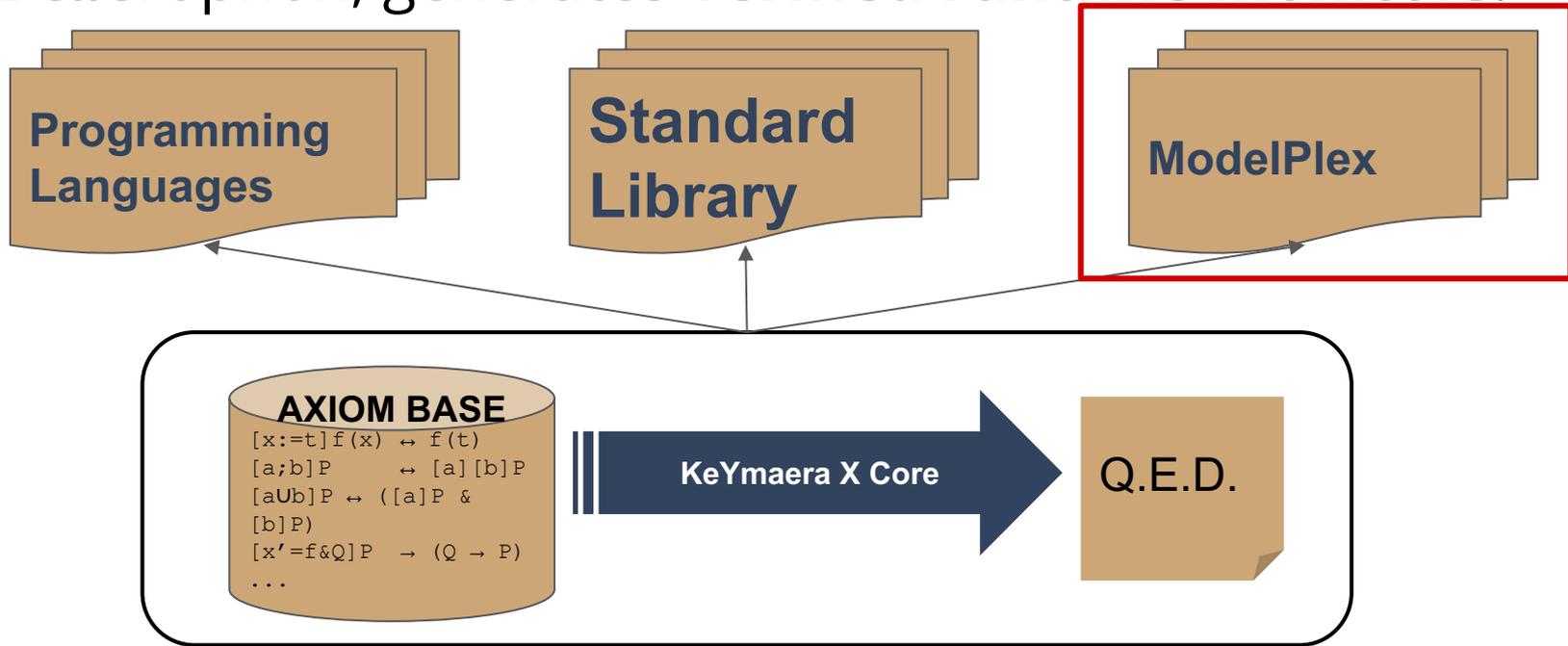
Details:

- Detect **modeled** vs **unmodeled** state space correctly at runtime.
- Convert monitors into reward signals



Detecting **unmodeled** State Space

The ModelPlex algorithm, implemented using Bellerophon, generates **verified runtime monitors**.



Detecting **unmodeled** State Space

```
oldPos := read_sensor(GPS)
actuate(accel)
newPos := read_sensor(GPS)
if ( $\exists t$ . model_after(t) == newPos):
    # No model deviation.
else:
    # Model deviation...?
```

Detecting **unmodeled** State Space

```
oldPos := read_sensor(GPS)
actuate(accel)
newPos := read_sensor(GPS)
if ( $\exists t$ . model_after(t) == newPos):
    # No model deviation.
else:
    # Model deviation...?
```

Detecting **unmodeled** State Space

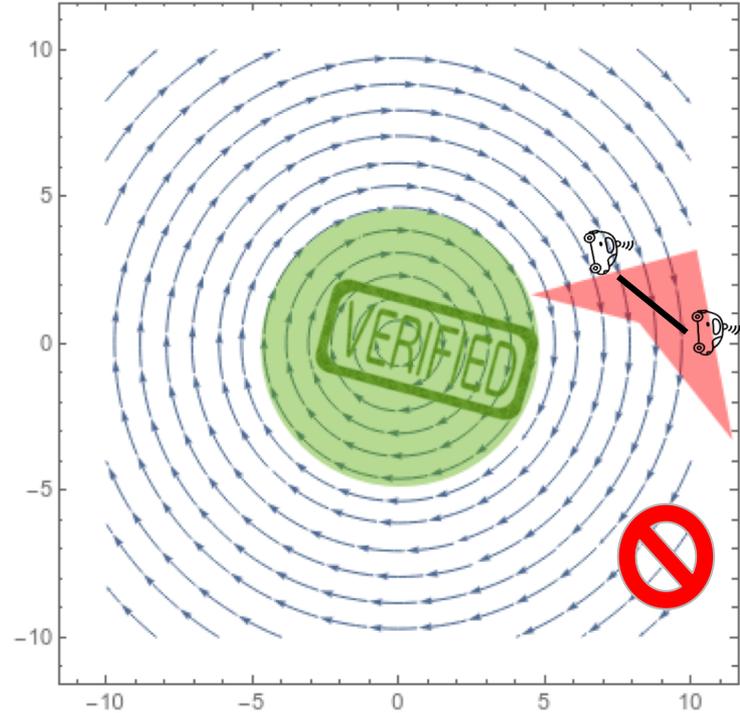
```
oldPos := read_sensor(GPS)
actuate(accel)
newPos := read_sensor(GPS)
if (QE( $\exists t. \text{model\_after}(t) == \text{newPos}$ )):
    # No model deviation.
else:
    # Model deviation...?
```

Safe RL: How?

Details:

Runtime monitoring separates **modeled** from **unmodeled** state space.

- Convert monitors into reward signals



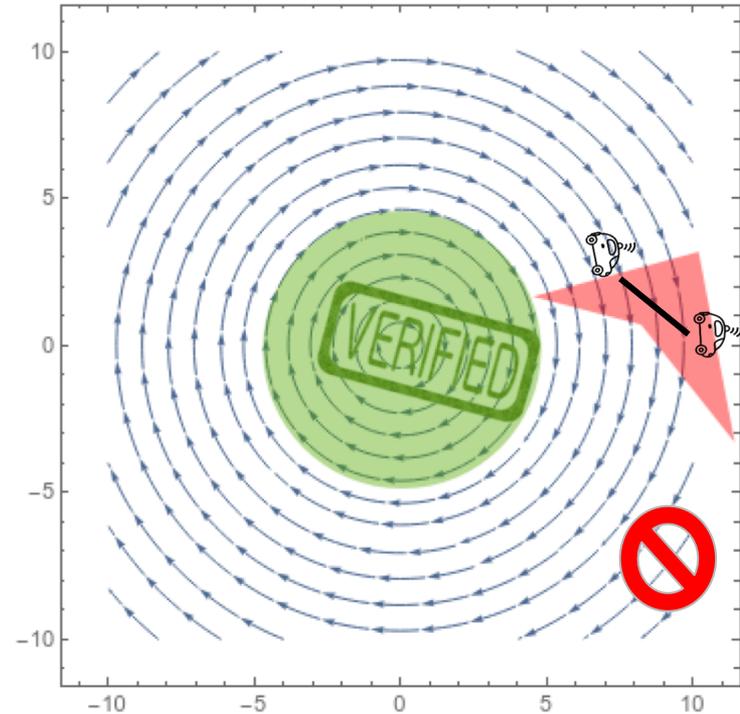
Safe RL: How?

Details:

Runtime monitoring separates **modeled** from **unmodeled** state space.

- Convert monitors into reward signals:

$$(\mathbb{R}^n \rightarrow \mathbb{B}) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R})!?$$



An Example

init \rightarrow [{

{?safeAccel; accel U brake U ?safeMaint; maintVel};

{pos' = vel, vel' = acc, t'=1}

]*]safe

An Example Monitor

init \rightarrow [{

{?safeAccel; accel \cup brake \cup ?safeMaintain; maintainVel};

{pos' = vel, vel' = acc, t'=1}

]*]safe

$(t_{\text{post}} \geq 0 \wedge a_{\text{post}} = \text{acc} \wedge v_{\text{post}} = \text{acc } t_{\text{post}} + v \wedge p_{\text{post}} = \text{acc } t_{\text{post}}^2/2 + v t_{\text{post}} + p) \vee$

$(t_{\text{post}} \geq 0 \wedge a_{\text{post}} = 0 \wedge v_{\text{post}} = v \wedge p_{\text{post}} = vt_{\text{post}} + p) \vee$ Etc.

An Example Monitor

init \rightarrow [{

{?safeAccel; accel \cup brake \cup ?safeMaintain; maintainVel};

{pos' = vel, vel' = acc, t'=1}

]*]safe

$(t_{\text{post}} \geq 0 \wedge a_{\text{post}} = \text{accel} \wedge \mathbf{v}_{\text{post}} = \mathbf{acc} t_{\text{post}} + \mathbf{v} \wedge p_{\text{post}} = \text{acc} t_{\text{post}}^2/2 + \mathbf{v} t_{\text{post}} + p) \vee$

$(t_{\text{post}} \geq 0 \wedge a_{\text{post}} = 0 \wedge \mathbf{v}_{\text{post}} = \mathbf{v} \wedge p_{\text{post}} = \mathbf{v} t_{\text{post}} + p) \vee$ Etc.

An Example: The Monitor

init \rightarrow [{

{?safeAccel; accel \cup brake \cup ?safeMaintain; maintainVel};

{*pos'* = *vel*, *vel'* = *acc*, *t'*=1}

]*]safe

$(t_{\text{post}} \geq 0 \wedge a_{\text{post}} = \text{acc} \wedge v_{\text{post}} = \text{accel } t_{\text{post}} + v \wedge p_{\text{post}} = \text{acc } t_{\text{post}}^2/2 + v t_{\text{post}} + p) \vee$

$(t_{\text{post}} \geq 0 \wedge a_{\text{post}} = 0 \wedge v_{\text{post}} = v \wedge p_{\text{post}} = vt_{\text{post}} + p) \vee$ Etc.

An Example: The Monitor



- Q.E. for RCF
- ODE solutions backed by proofs

init \rightarrow [{

{?safeAccel; accel \cup brake \cup ?safeMaintain; maintainVel};

{ $pos' = vel$, $vel' = acc$, $t'=1$ }

]*]safe

$(t_{post} \geq 0 \wedge a_{post} = acc \wedge v_{post} = accel t_{post} + v \wedge p_{post} = acc t_{post}^2/2 + v t_{post} + p) \vee$

$(t_{post} \geq 0 \wedge a_{post} = 0 \wedge v_{post} = v \wedge p_{post} = vt_{post} + p) \vee$ Etc.

Quantitative monitor as reward signal

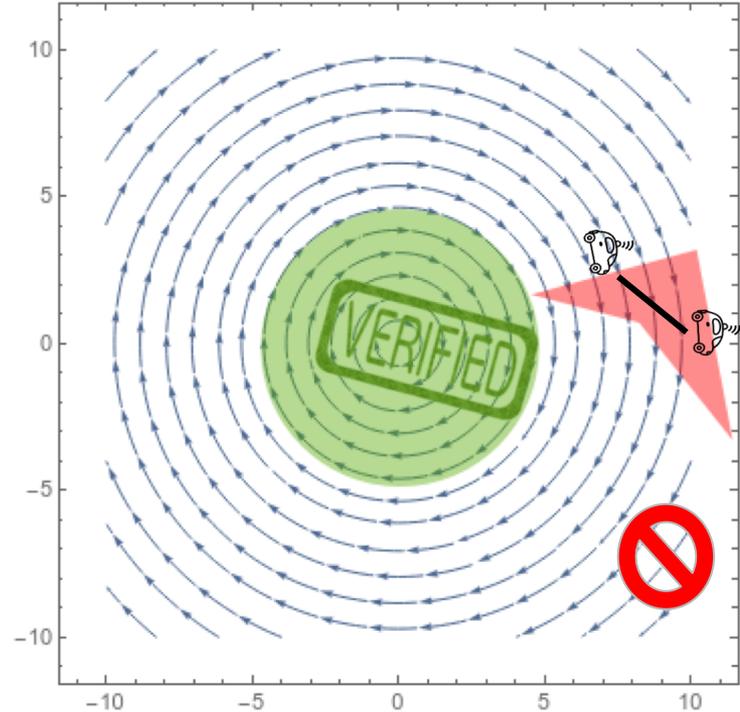
Safe RL: How?

Details:

Runtime monitoring separates **modeled** from **unmodeled** state space.

Convert monitors into gradients:

$$(\mathbb{R}^n \rightarrow \mathbb{B}) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R})$$

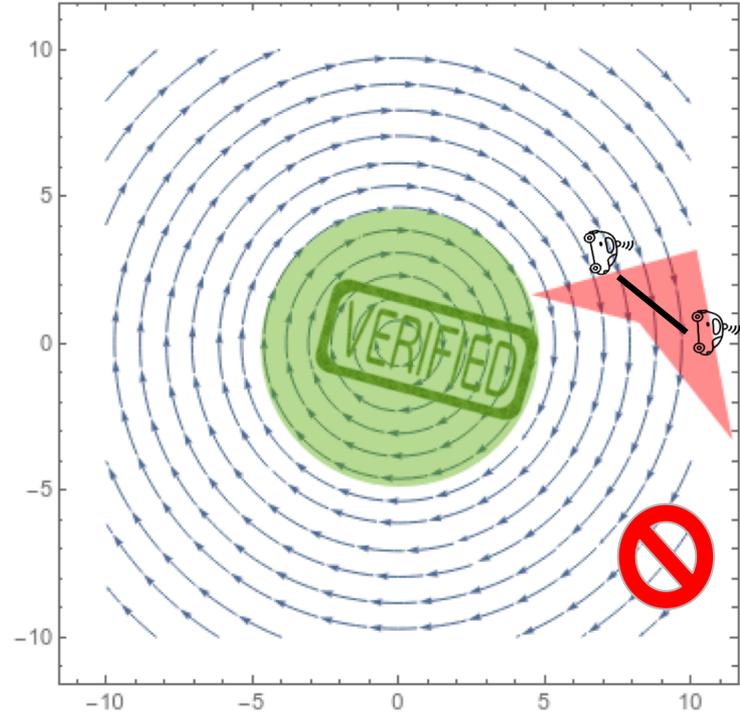


Safe RL: How?

Details:

Runtime monitoring separates **modeled** from **unmodeled** state space.

Convert **models** into gradients: ModelPlex
 $(\mathbb{R}^n \rightarrow \mathbb{B}) \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R})$



Conclusion

KeYmaera X + Justified Speculative Control provide strong safety guarantees for learning-enabled CPS.

1. Was the proof correct?
2. Was the model accurate enough?



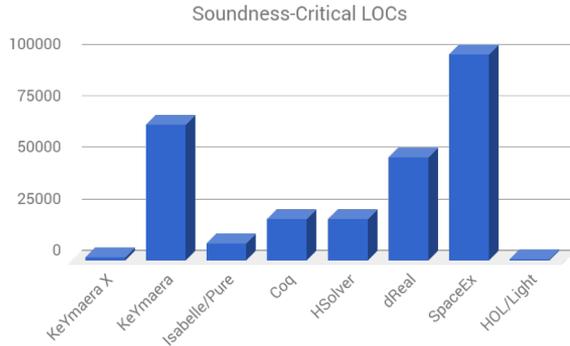
≠



Conclusion

KeYmaera X + Justified Speculative Control provide strong safety guarantees for learning-enabled CPS.

1. Was the proof correct? **KeYmaera X**
2. Was the model accurate enough?



dl Tactic:

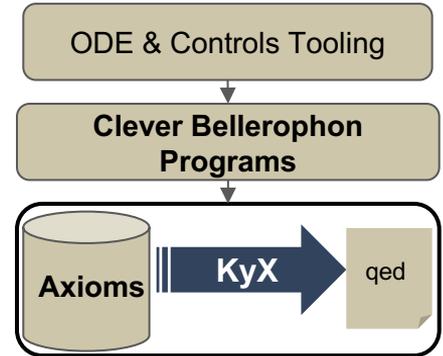
Side derivation:
 $(v \geq v_0 - gt)'$ ↔
 $\dots \leftrightarrow$
 $\dots \leftrightarrow$
 \dots
 $H = r_p \geq 0 \ \& \ r_a \geq 0$
 $\ \& \ g > 0 \ \& \ \dots$

DI Axiom:

$[\{x'=f\&Q\}]P \leftrightarrow ([?Q]P \leftarrow (Q \rightarrow [\{x'=f\&Q\}]P'))$

Example:

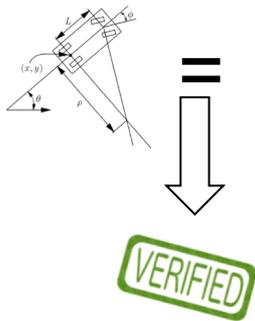
$[v' = r_p v^2 - g, t' = 1] v \geq v_0 - gt \leftrightarrow$
 $\dots \leftrightarrow$
 $[v' := r_p v^2 - g] [t' := 1] v' \geq -g * t' \leftrightarrow$
 $r_p v^2 - g \geq -g$
 $H \rightarrow r_p \geq 0$



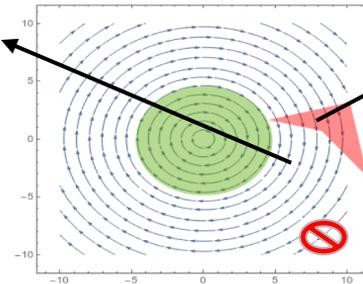
Conclusion

KeYmaera X + Justified Speculative Control provide strong safety guarantees for learning-enabled CPS.

1. Was the proof correct? **KeYmaera X**
2. Was the model accurate enough? **Justified Speculation**



Get to here...



...from here

Conclusion

KeYmaera X + Justified Speculative Control provide strong safety guarantees for learning-enabled CPS.

1. Was the proof correct? **KeYmaera X**
2. Was the model accurate enough? **Justified Speculation**

Web	keymaeraX.org
Online Demo	web.keymaeraX.org
Open Source (GPL)	github.com/LS-Lab/KeYmaeraX-release

Acknowledgments

Students and postdocs of the Logical Systems Lab at Carnegie Mellon
Brandon Bohrer, Nathan Fulton, Sarah Loos, João Martins, Yong Kiam Tan
Khalil Ghorbal, Jean-Baptiste Jeannin, Stefan Mitsch



I Part: Elementary Cyber-Physical Systems

1. Differential Equations & Domains
2. Choice & Control
3. Safety & Contracts
4. Dynamical Systems & Dynamic Axioms
5. Truth & Proof
6. Control Loops & Invariants
7. Events & Responses
8. Reactions & Delays

II Part: Differential Equations Analysis

9. Differential Equations & Differential Invariants
10. Differential Equations & Proofs
11. Ghosts & Differential Ghosts
12. Differential Invariants & Proof Theory

III Part: Adversarial Cyber-Physical Systems

- 13-16. Hybrid Systems & Hybrid Games

IV Part: Comprehensive CPS Correctness



Logical Foundations of Cyber-Physical Systems