# Safe Reinforcement Learning via Formal Methods

**Nathan Fulton** and André Platzer

Carnegie Mellon University
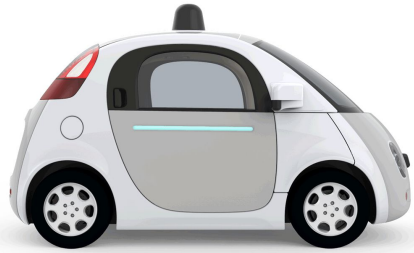
# Safety-Critical Systems



"How can we provide people with cyber-physical systems they can bet their lives on?" - Jeannette Wing
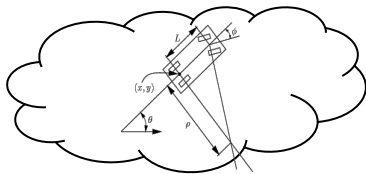
# **Autonomous** Safety-Critical Systems



How can we provide people with **autonomous** cyber-physical systems they can bet their lives on?
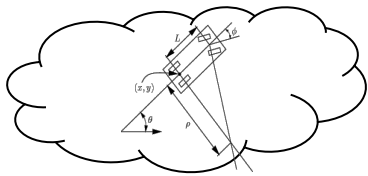
# Model-Based Verification



φ

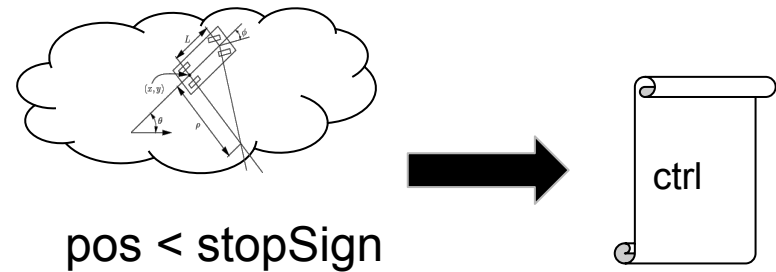Reinforcement Learning

# Model-Based Verification



pos < stopSign

# Reinforcement Learning

# Model-Based Verification

Reinforcement Learning



pos < stopSign

# Model-Based Verification

Reinforcement Learning



pos < stopSign

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

# Model-Based Verification



pos < stopSign

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

# Reinforcement Learning

# Model-Based Verification



φ

**Benefits:**

- Strong safety guarantees
- Automated analysis

# Reinforcement Learning

# Model-Based Verification     Reinforcement Learning



φ

**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



φ

**Benefits:**

- Strong safety guarantees
- Automated analysis

**Drawbacks:**

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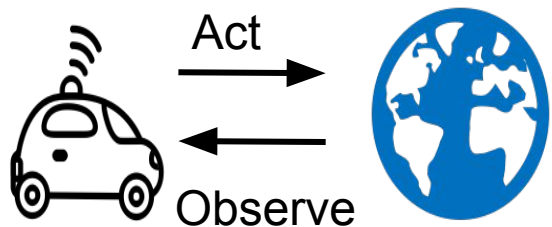- Assumes accurate model

# Model-Based Verification



φ

**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



Act

Observe

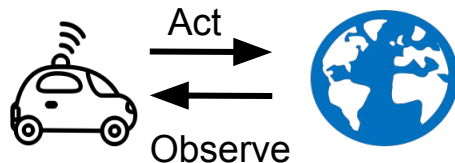# Model-Based Verification



$\varphi$

**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



Act

Observe

**Benefits:**

- No need for complete model
- Optimal (effective) policies
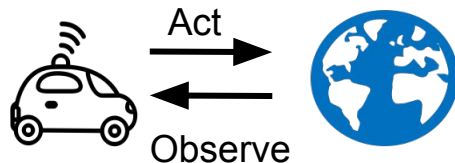
# Model-Based Verification



φ

**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



Act

Observe

**Benefits:**

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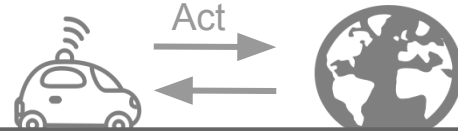- Optimal (effective) policies

**Drawbacks:**

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

# Model-Based Verification

# Reinforcement Learning

**Act**

**Bene**

- S                                    del
- A                                    s
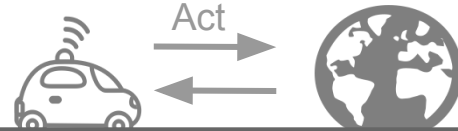
**Draw**

**Goal: Provably correct reinforcement learning**

- 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

Act

**Bene**

- S                                                          del
- A                                                          s

**Draw**

- 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

**Goal: Provably correct reinforcement learning**
  1. **Learn Safely**
  2. **Learn a Safe Policy**
  3. **Justify claims of safety**

# Model-Based Verification

Accurate, analyzable models often exist!

```
{

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

    {pos' = vel, vel' = acc}

}*
```

# Model-Based Verification

**Accurate**, analyzable models often exist!

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

Continuous motion

discrete control

# Model-Based Verification

**Accurate**, analyzable models often exist!

```
{

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

    {pos' = vel, vel' = acc}

}*
```

Continuous motion

discrete, ***non-deterministic*** control

# Model-Based Verification

**Accurate**, **analyzable** models often exist!
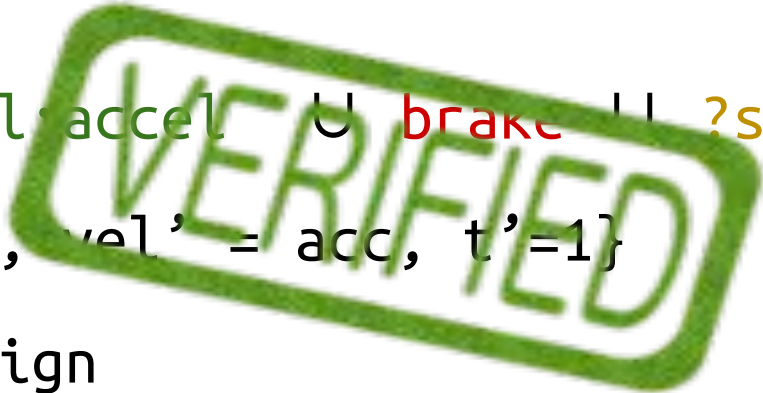
```
init → [{

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

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

}*]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, t'=1}

}*]pos < stopSign
```

# Model-Based Verification

**Accurate**, **analyzable** models often exist!

formal verification gives strong safety guarantees

**VERIFIED** = 
- **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 ∪ brake ∪ ?safeTurn; turn};

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

}*
```

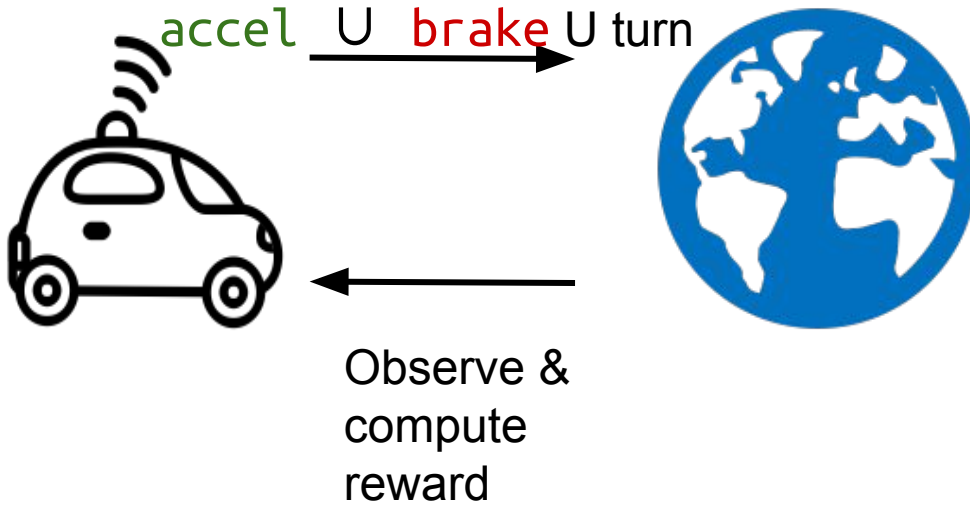Only accurate sometimes

# Our Contribution

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

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

# Our Contribution

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

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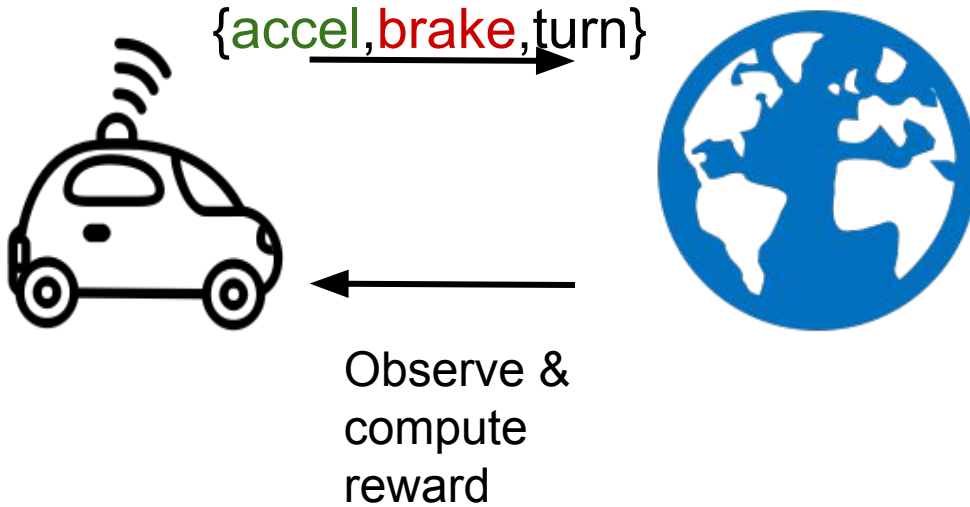2. allows and directs speculation whenever model mismatches occur

# Learning to Resolve Non-determinism



Act

Observe &
compute
reward

# Learning to Resolve Non-determinism



accel ∪ brake ∪ turn

Observe & compute reward

# Learning to Resolve Non-determinism

{accel,brake,turn}

Observe & compute reward

# Learning to Resolve Non-determinism

# Learning to Resolve Non-determinism

# Learning to **Safely** Resolve Non-determinism
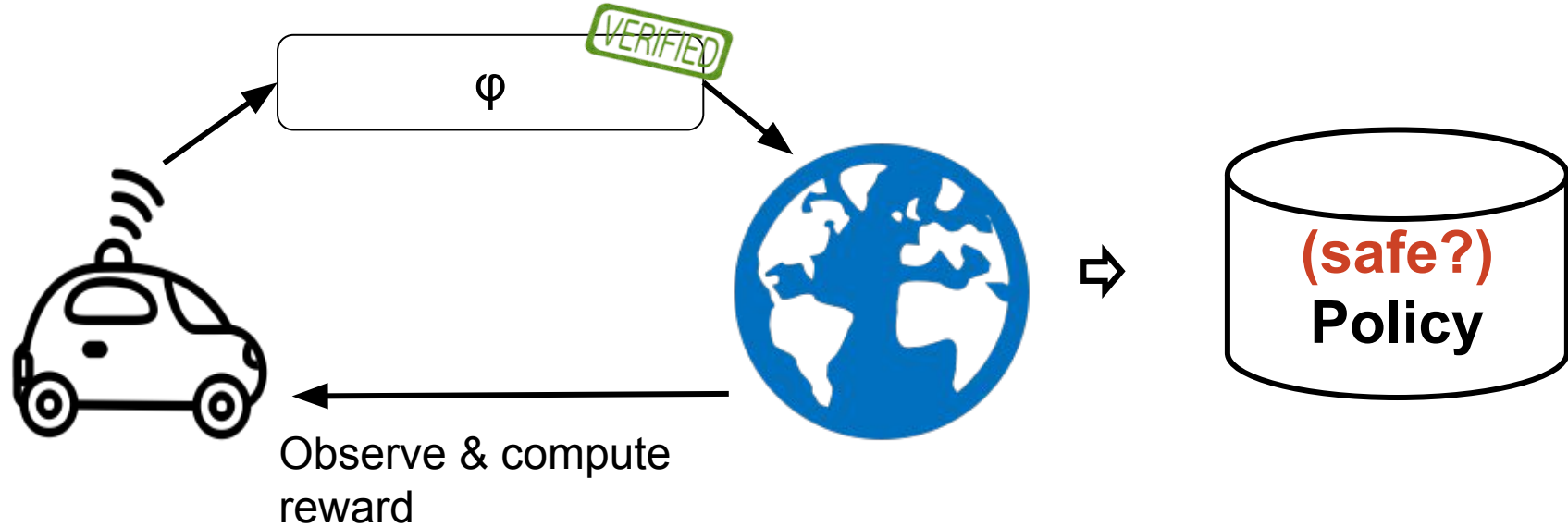
# Learning to **Safely** Resolve Non-determinism



Safety Monitor

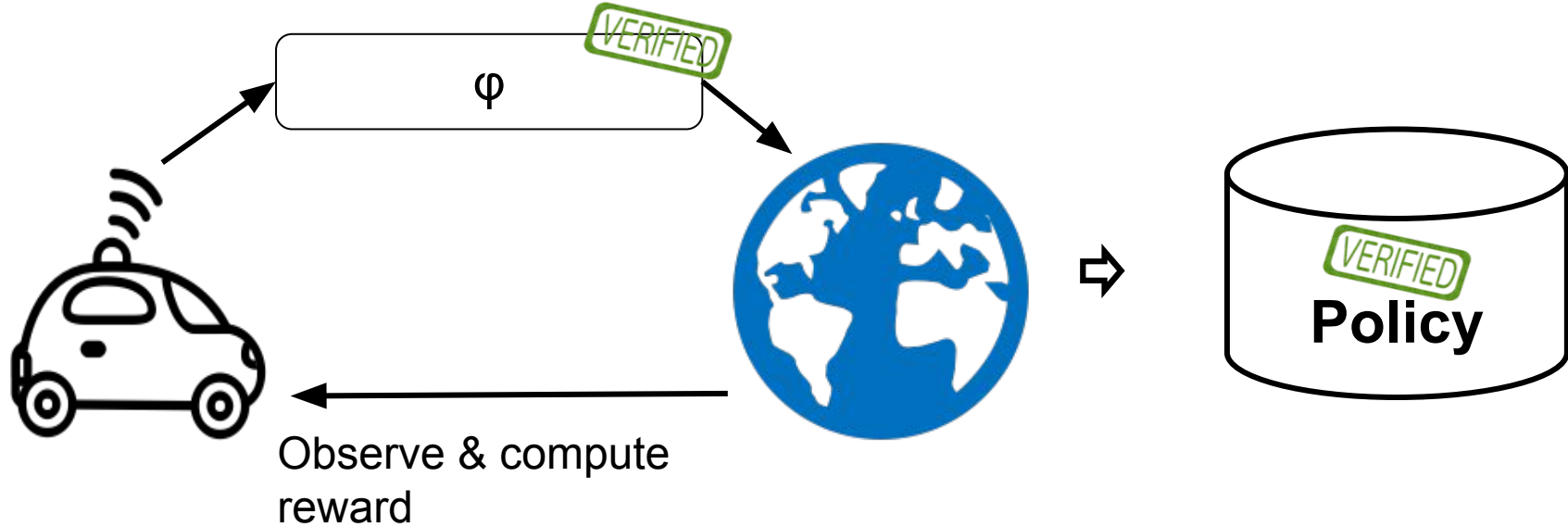VERIFIED

(safe?)
**Policy**

Observe & compute
reward

VERIFIED ≠ **"Trust Me"**

# Learning to **Safely** Resolve Non-determinism



Use a theorem prover to prove:

$$(\texttt{init} \rightarrow [\{\{\texttt{accel} \cup \texttt{brake}\}; \texttt{ODEs}\}*](\texttt{safe})) \leftrightarrow \varphi$$

# Learning to **Safely** Resolve Non-determinism



Observe & compute reward

Use a theorem prover to prove:

$$(\text{init} \rightarrow [\{\{\text{accel} \cup \text{brake}\};\text{ODEs}\}^*](\text{safe})) \leftrightarrow \varphi$$
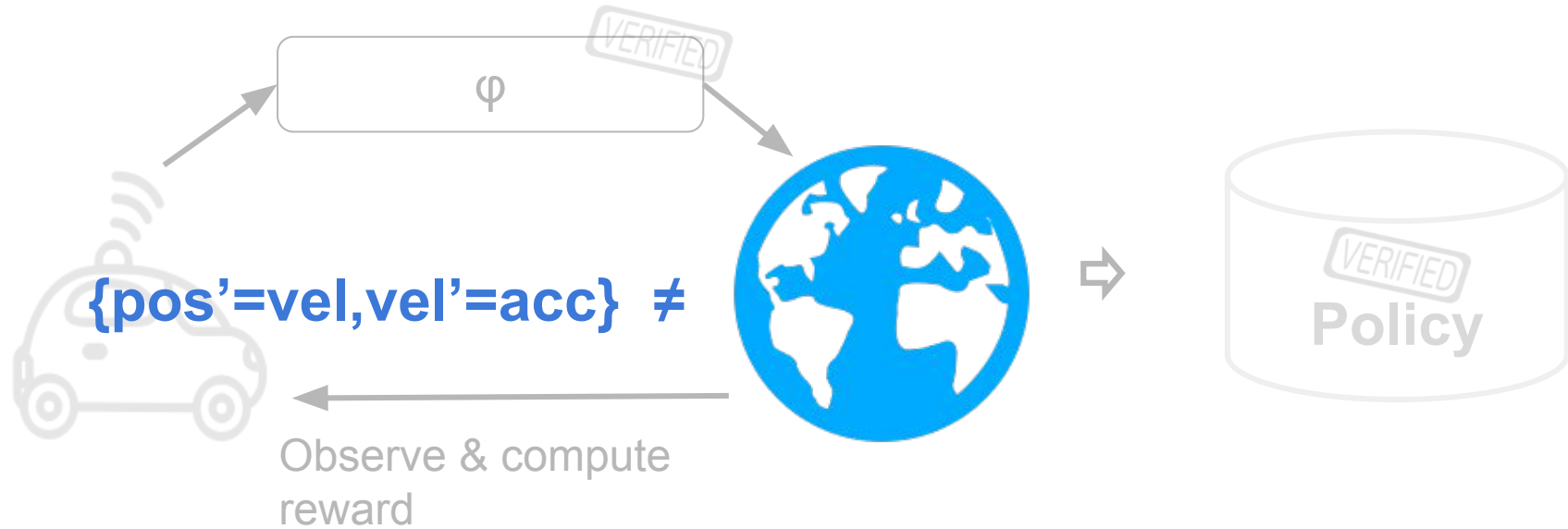
# Learning to **Safely** Resolve Non-determinism

φ

VERIFIED

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

Policy

Observe & compute reward

Use a theorem prover to prove:

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

# Learning to **Safely** Resolve Non-determinism

φ

**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.**
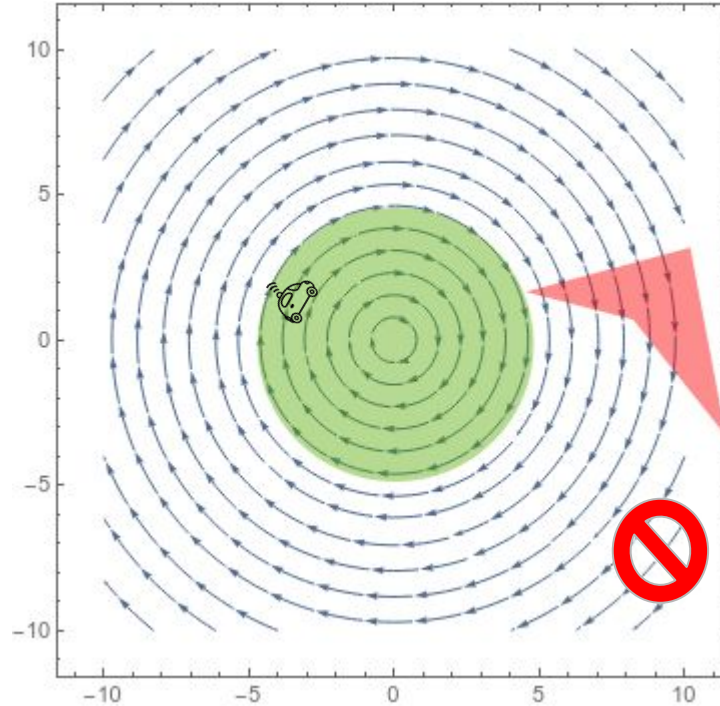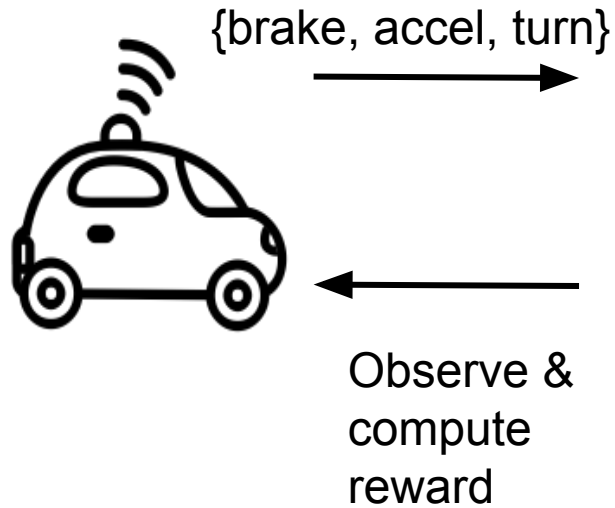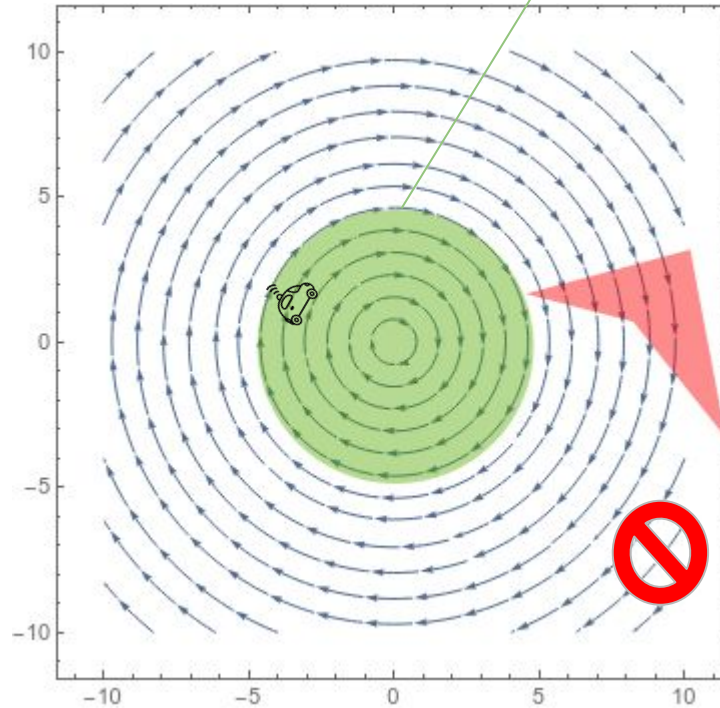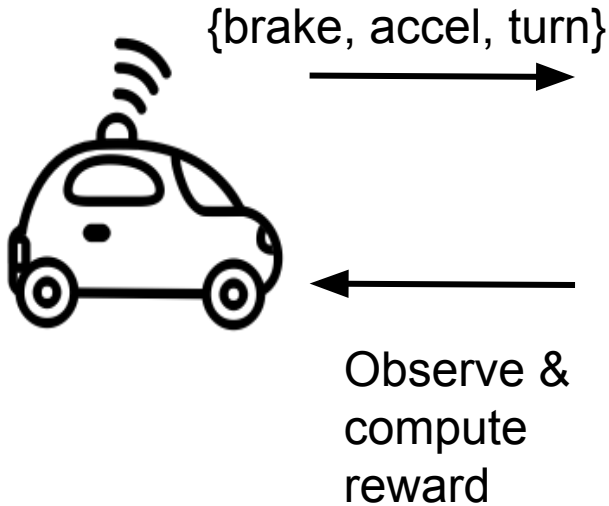
Use a theorem prover to prove:

(init→[{{accel∪brake};ODEs}*](safe)) ↔ φ

# What about the physical model?

**{pos'=vel,vel'=acc} ≠**

Observe & compute reward

Use a theorem prover to prove: (init→ [{{accel∪brake};ODEs}*](safe)) ↔ φ

φ

VERIFIED

VERIFIED
Policy

# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward

# What About the Physical Model?

**Model is accurate.**

{brake, accel, turn}

Observe & compute reward

# What About the Physical Model?

**Model is accurate.**

{brake, accel, turn}

Observe & compute reward

# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward
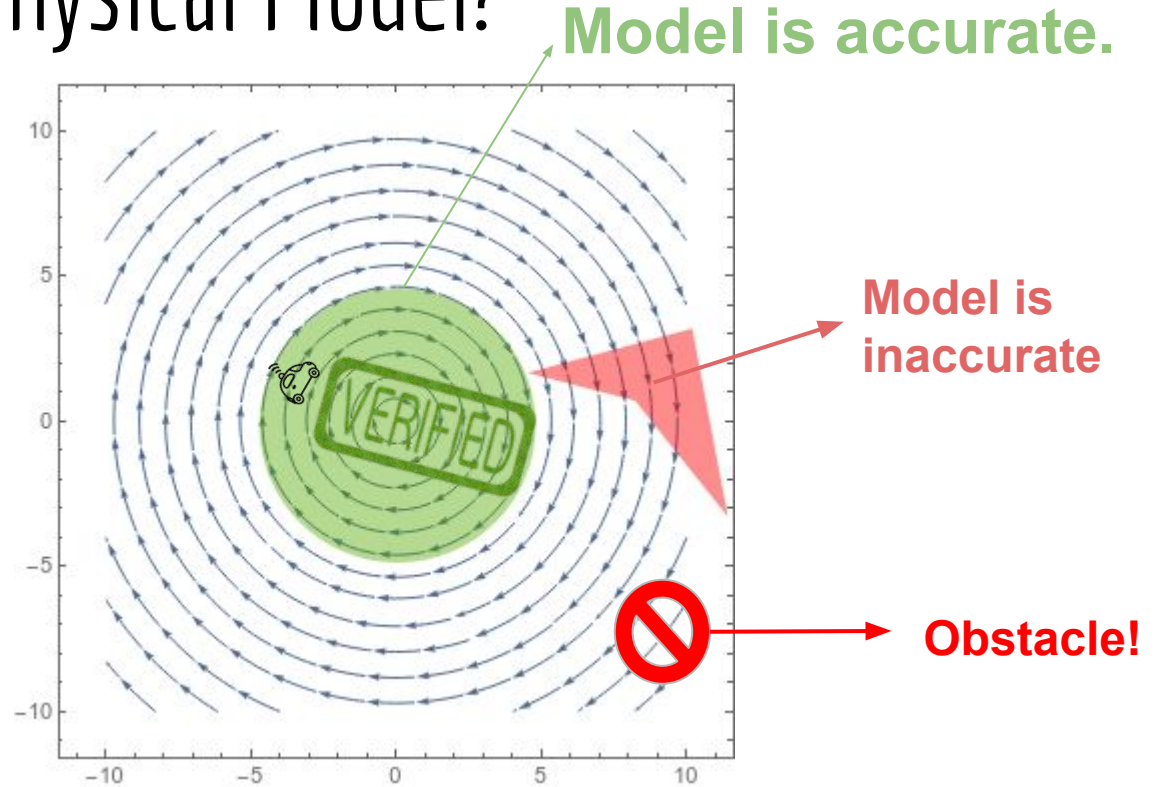
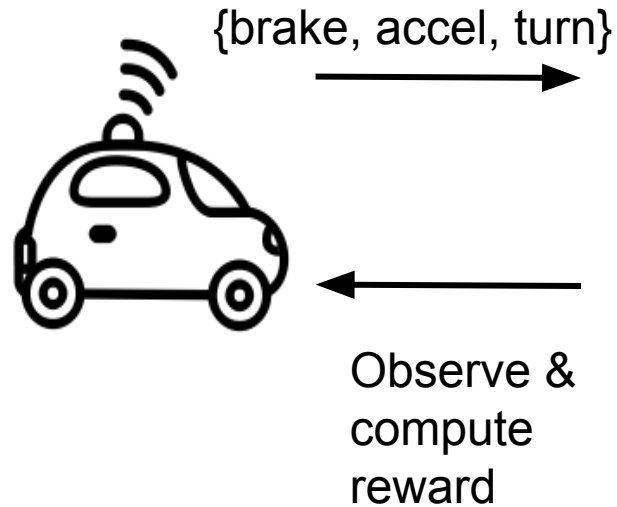**Model is accurate.**

**Model is inaccurate**

# What About the Physical Model?



**Model is accurate.**

**Model is inaccurate**
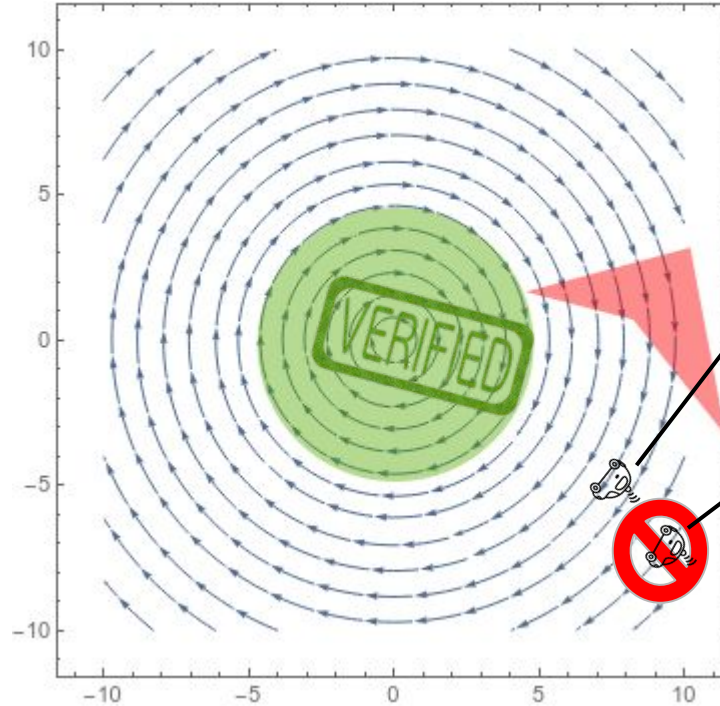
**Obstacle!**

{brake, accel, turn}

Observe & compute reward

# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward



Expected

Reality

# Speculation is Justified



{brake, accel, turn}

Observe & compute reward

Expected (safe)

Reality (crash!)

# Leveraging Verification Results to Learn Better

{brake, accel, turn}

Observe & compute reward

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

# An Example

# An Example: The System

init → [{

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

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

}*]safe

# An Example: The Monitor

init → [{

    {?safeAccel;accel ∪ brake ∪ ?safeMaintain; maintainVel};

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

}*]safe

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

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

# An Example: The Monitor

init → [{

    `{?safeAccel;accel` ∪ `brake` ∪ `?safeMaintain; maintainVel};`

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

}*]safe

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

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

# An Example: The Monitor

init → [{

  {?safeAccel;accel ∪ brake ∪ ?safeMaintain; maintainVel};

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

}*]safe

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

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

# An Example: The Monitor

init → [{

    {?safeAccel;accel ∪ brake ∪ ?safeMaintain; maintainVel};

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

}*]safe

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

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

# An Example: The Monitor  VERIFIED

init → [{

   {?safeAccel;accel ∪ brake ∪ ?safeMaintain; maintainVel};

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

}*]safe

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

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

# An Example: The Monitor

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

init → [{

   {?safeAccel;accel ∪ brake ∪ ?safeMaintain; maintainVel};

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

}*]safe

$(t_{post} >= 0 \land a_{post} = acc \land v_{post} = accel\ t_{post} + v \land \mathbf{p_{post} = acc\ t_{post}^2/2 + v\ t_{post} + p})\ \lor$

$(t_{post} >= 0 \land a_{post} = 0 \land v_{post} = v \land \boldsymbol{p_{post} = v t_{post} + p})\ \lor$ Etc.

# An Example: The Reward Signal (simplified)

$$x \geq 0 \land v \geq 0 \land A \geq 0 \rightarrow [\{x' = v, v' = A\}]x \geq 0$$

# An Example: The Reward Signal (simplified)

$$x\text{>=}0 \land v\text{>=}0 \land A\text{>=}0 \rightarrow [\{x' = v, v' = A\}]x\text{>=}0$$

Minimize **max(vError, xError)** where

$\text{vError} = \max(v_{post} - (A*t_{post} + v), A*t_{post} + v - v_{post})$

$\text{xError} = \max($
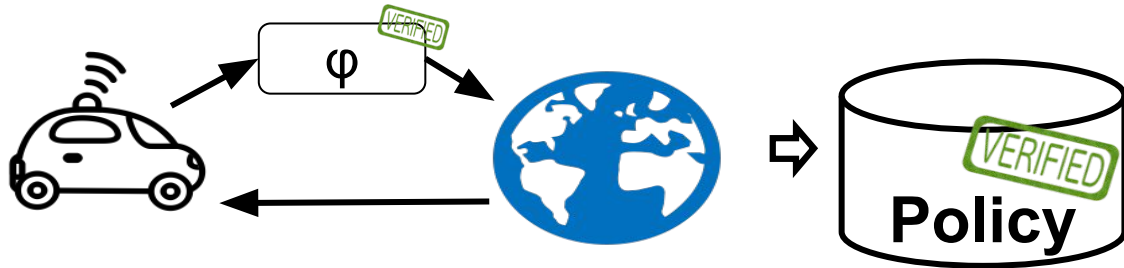
$x_{post} - (A*t_{post}^2/2 + v * t_{post} + x)$

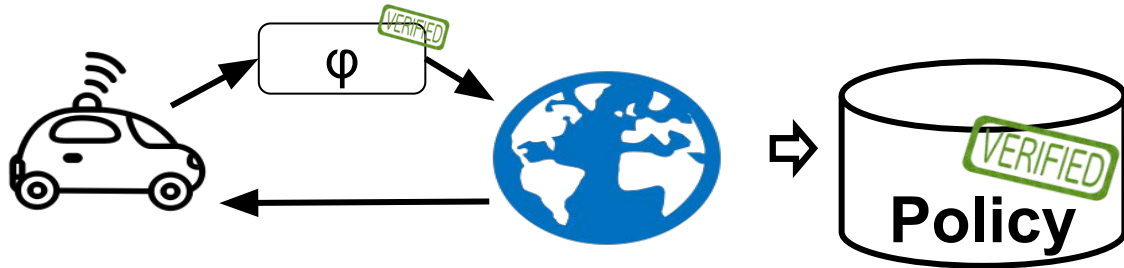$(A*t_{post}^2/2 + v*t_{post} + x) - x_{post}$

$)$

# Conclusion

**Justified Speculative Control** provides the best of logic and learning:

# Conclusion

**Justified Speculative Control** provides the best of logic and learning:
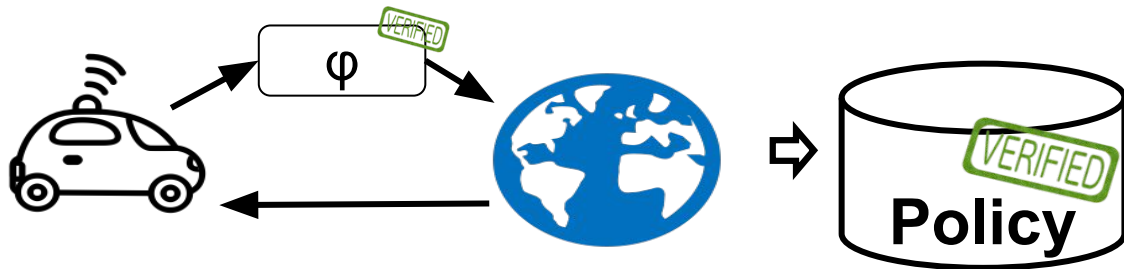
- Formally model the control system (**control + physics**)

# Conclusion

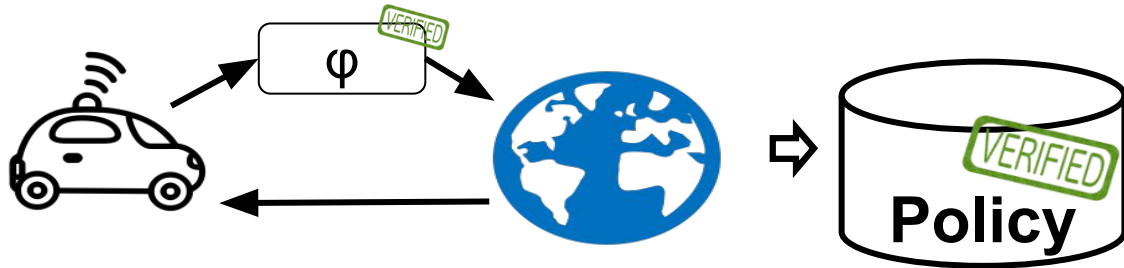**Justified Speculative Control** provides the best of logic and learning:

- Formally model the control system (**control + physics**)
- Learn how to resolve non-determinism in models.

# Conclusion

**Justified Speculative Control** provides the best of logic and learning:

- Formally model the control system (**control + physics**)
- Learn how to resolve non-determinism in models.
- Leverage theorem proving to transfer **proofs** to learned policies.

# Conclusion

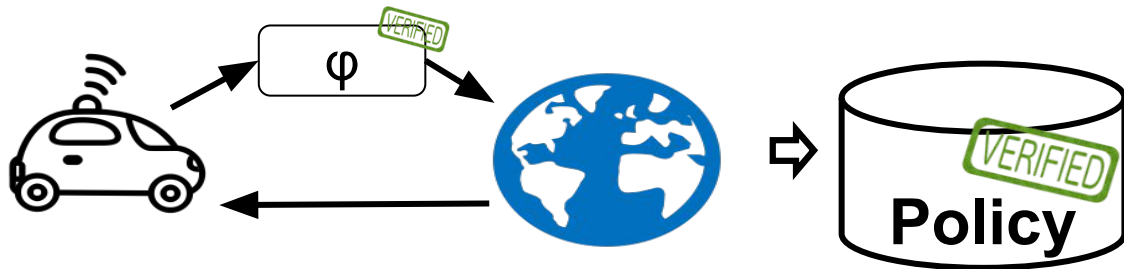**Justified Speculative Control** provides the best of logic and learning:

- Formally model the control system (**control + physics**)
- Learn how to resolve non-determinism in models.
- Leverage theorem proving to transfer **proofs** to learned policies.
- Unsafe **speculation is justified** when model deviates from reality

# Conclusion

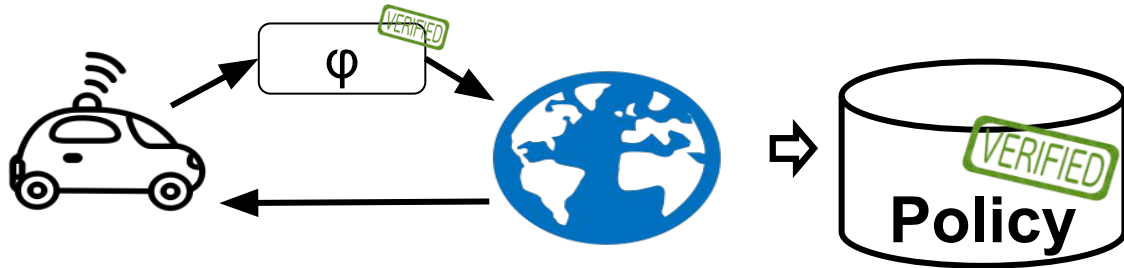**Justified Speculative Control** provides the best of logic and learning:

- Formally model the control system (**control + physics**)
- Learn how to resolve non-determinism in models
- Leverage theorem proving to transfer **proofs** to learned policies
- Unsafe **speculation is justified** when model deviates from reality, but **verification results can still be helpful**!

# Conclusion

**Justified Speculative Control** provides the best of logic and learning:

- Formally model the control system (**control + physics**)
- Learn how to resolve non-determinism in models
- Leverage theorem proving to transfer **proofs** to learned policies
- Unsafe **speculation is justified** when model deviates from reality, but **verification results can still be helpful**!