# Safe Reinforcement Learning via Formal Methods

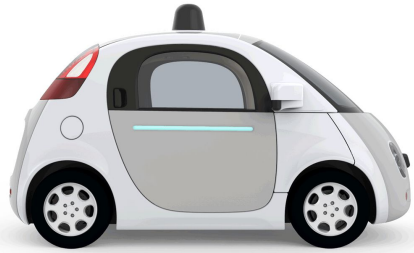**Nathan Fulton** and André Platzer

Carnegie Mellon University

# Safe Reinforcement Learning via Formal Methods

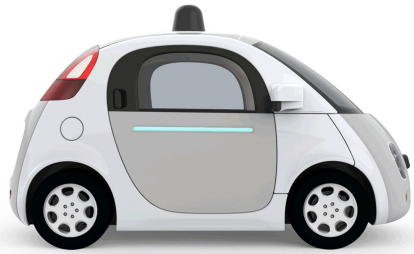**Nathan Fulton** and André Platzer

Carnegie Mellon University

# Safety-Critical Systems



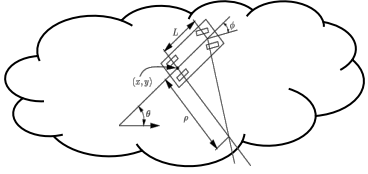"How can we provide people with cyber-physical systems they can bet their lives on?" - Jeannette Wing

# **Autonomous** Safety-Critical Systems



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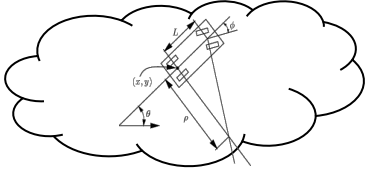
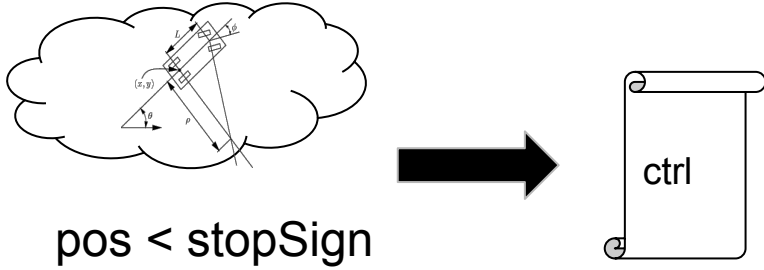# Model-Based Verification    Reinforcement Learning



φ

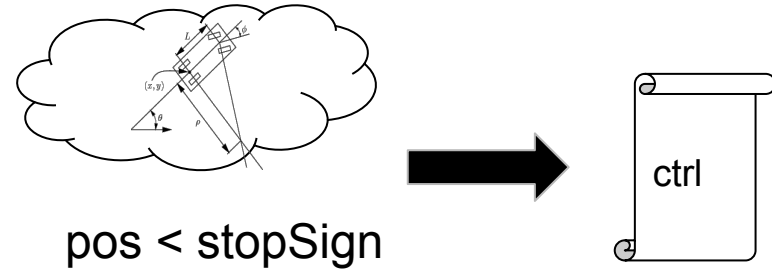# Model-Based Verification    Reinforcement Learning



pos < stopSign

# Model-Based Verification



pos < stopSign

# Reinforcement Learning

# 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    Reinforcement Learning



pos < stopSign

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

# Model-Based Verification

Reinforcement Learning



φ

**Benefits:**

- Strong safety guarantees
- Automated analysis

# 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    VERIFIED

**Benefits:**

- Strong safety guarantees
- Automated analysis

**Drawbacks:**

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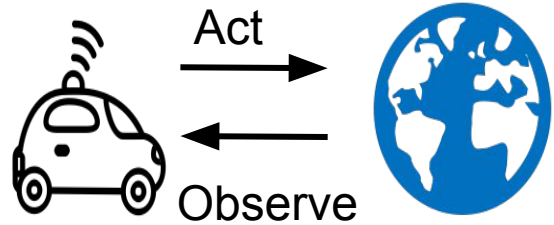- Assumes accurate model
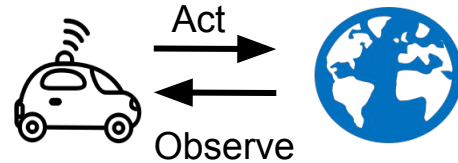
# Model-Based Verification



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

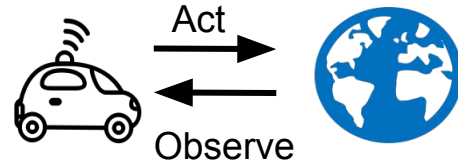Observe

# Model-Based Verification



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

# 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
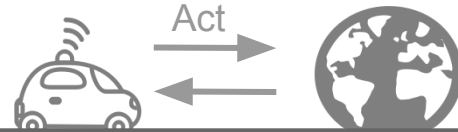- Optimal (effective) policies

**Drawbacks:**

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

# Model-Based Verification

# Reinforcement Learning

Act

**Bene**

- S                                                              del
- A                                                              s

**Draw**

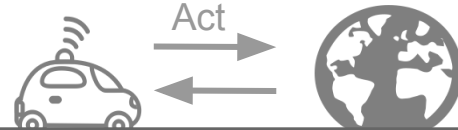## Goal: Provably correct reinforcement learning

- 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
- A

del

**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 Safety**
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}

}*]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

**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  ⊔  brake ⊔  ?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

# 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



Safety Monitor

Observe & compute reward

**(safe?)**
**Policy**

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

Safety Monitor

VERIFIED

(safe?)
**Policy**

Observe & compute reward

VERIFIED ≠ **"Trust Me"**

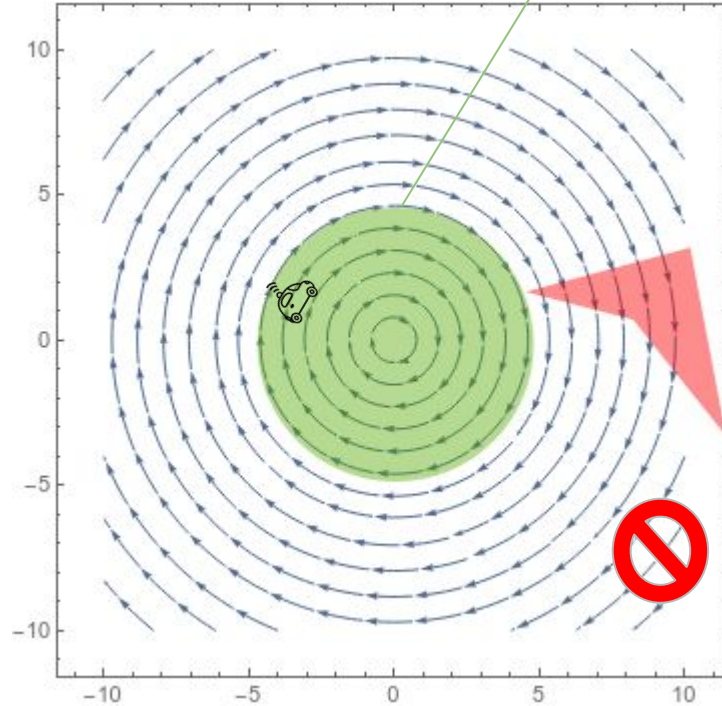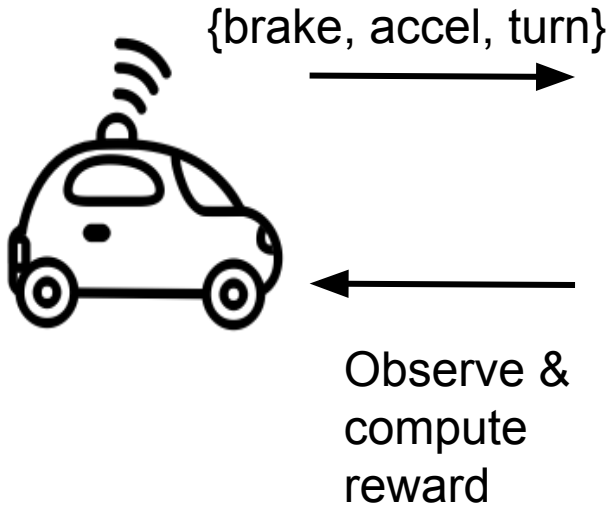# Learning to **Safely** Resolve Non-determinism



Observe & compute reward

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



Use a theorem prover to prove:

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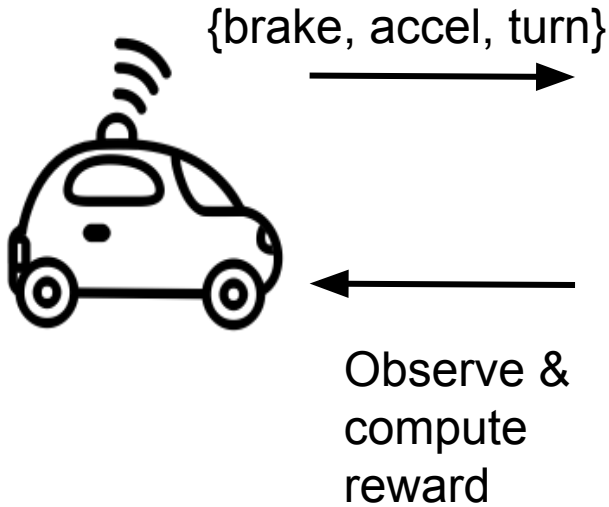
# Learning to **Safely** Resolve Non-determinism

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:

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

# 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 via the model monitor.**

Policy

Observe & compute reward

Use a theorem prover to prove:

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

# What about the physical model?

φ

VERIFIED

**{pos'=vel,vel'=acc}  ≠**

VERIFIED
**Policy**

Observe & compute
reward

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

# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward

# What About the Physical Model?

**Model is accurate.**

{brake, accel, turn}
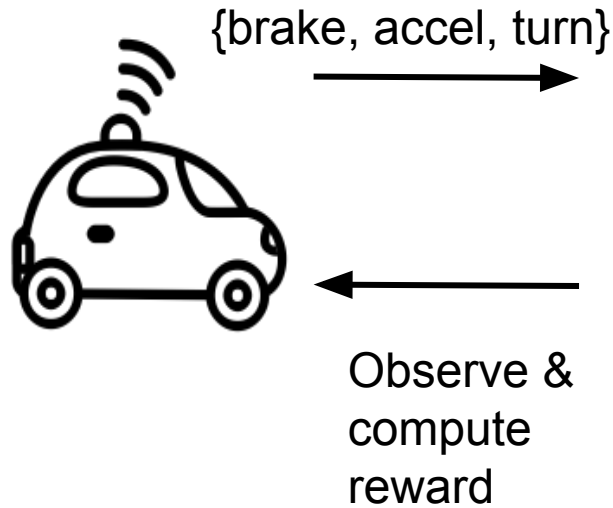
Observe & compute reward

# What About the Physical Model?

**Model is accurate.**

{brake, accel, turn}

Observe &
compute
reward

# What About the Physical Model?

**Model is accurate.**

{brake, accel, turn}

Observe & compute reward

**Model is inaccurate**

# What About the Physical Model?

**Model is accurate.**

{brake, accel, turn}

Observe & compute reward

**Model is inaccurate**

**Obstacle!**

# 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

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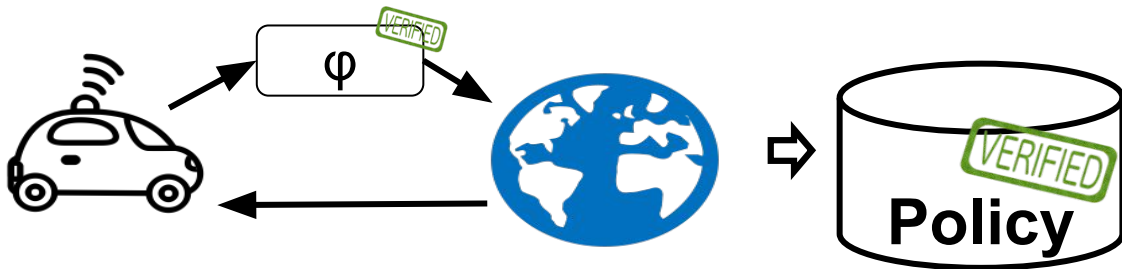
**Justified Speculative Control** provides the best of logic and learning:

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

# Conclusion

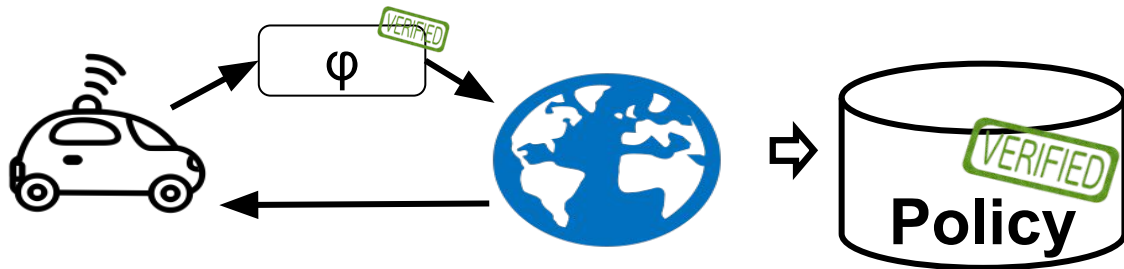**Justified Speculative Control** provides the best of logic and learning:

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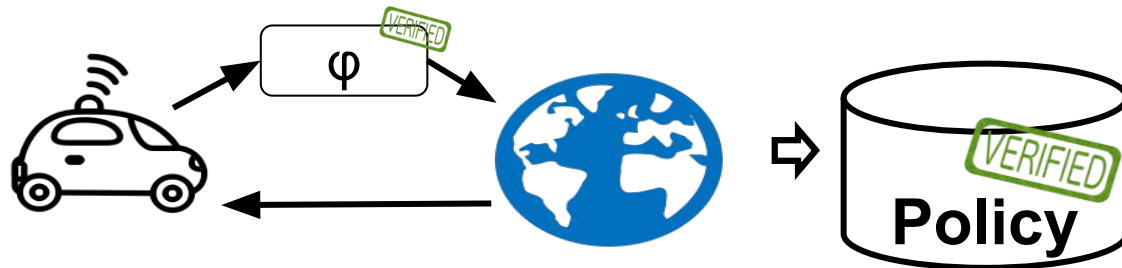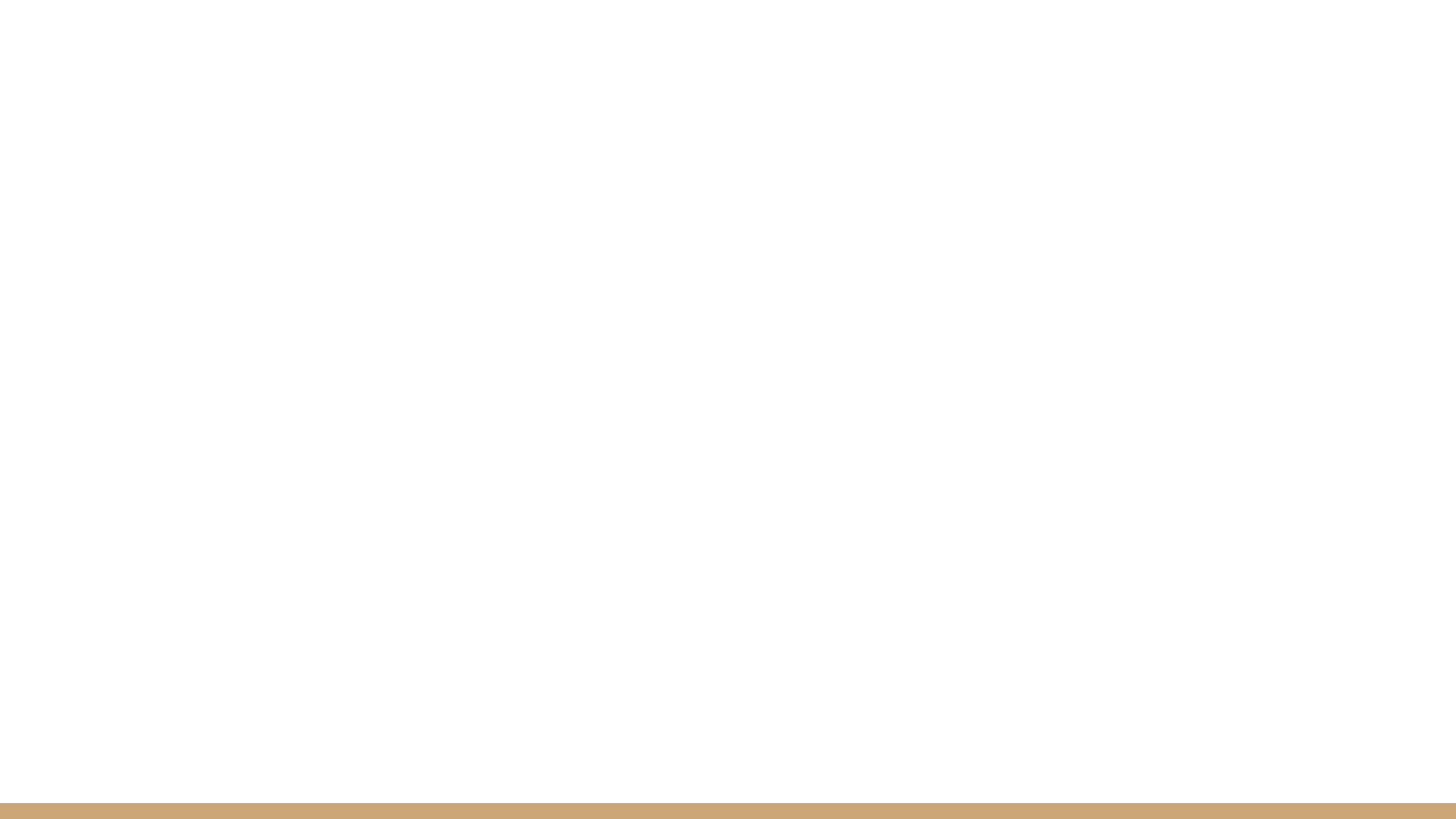- Leverage theorem proving to transfer **proofs** to learned policies.

# 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**!
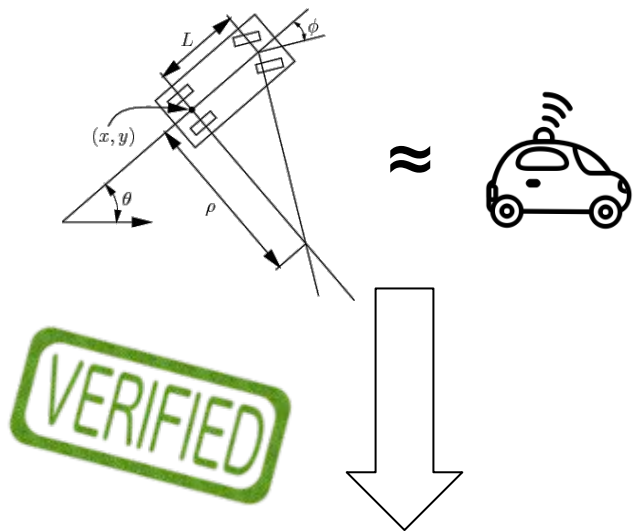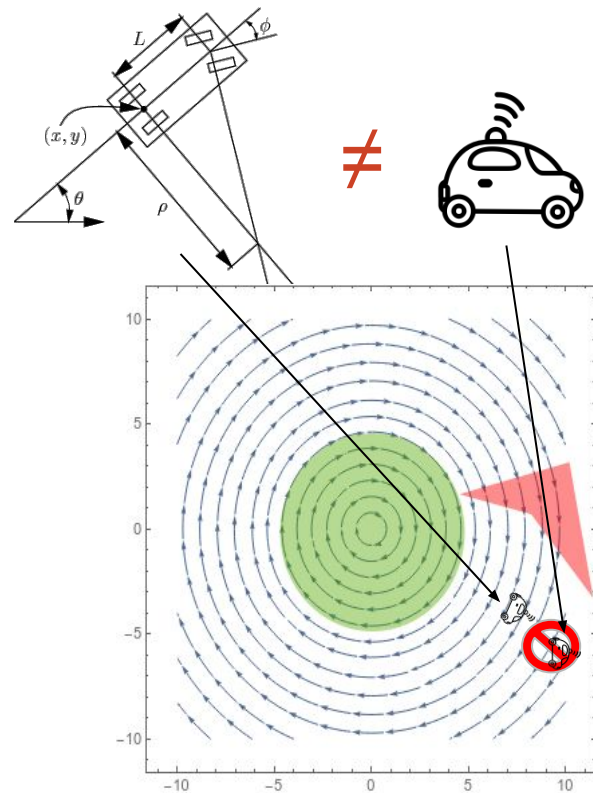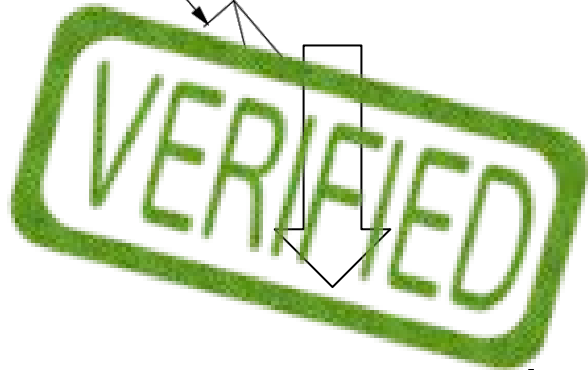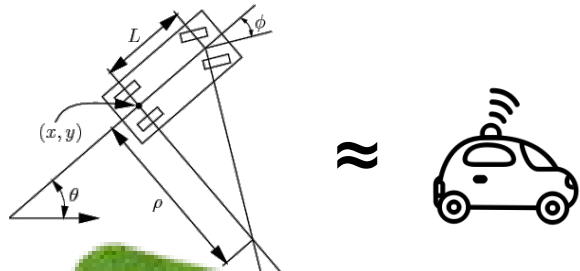
# Conclusion

**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**!
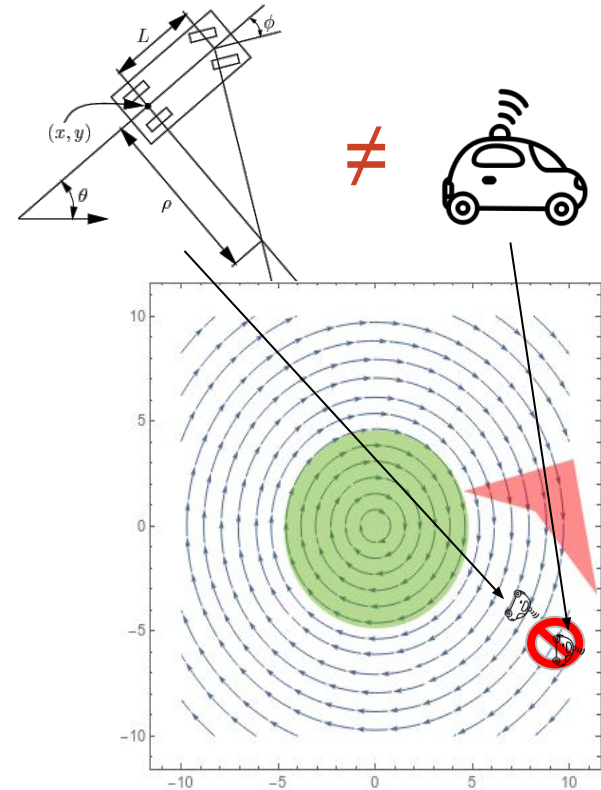
# Justified Speculative Control



≈

≠

**VERIFIED**

Learn over a constrained action space

# Justified Speculative Control



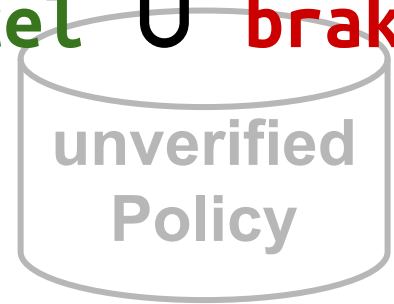Learn over a constrained action space

# Safe Reinforcement Learning?

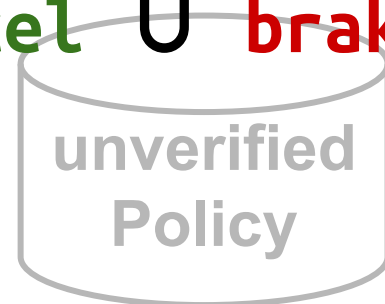{accel,brake,turn}

⇨ **unverified Policy**

Observe & compute reward

Policy deviates from model:

1. Policy is deterministic, verification result is set-valued.

# Some Actions Aren't Always Safe

**{accel,brake,turn} ≠ ?safeAccel; accel ∪ brake**



unverified Policy

Observe & compute reward

Policy deviates from model:
1. Policy is deterministic, verification result is set-valued.

# Some Actions Aren't Always Safe

**{accel,brake,turn} ≠ ?safeAccel; accel ∪ brake**



Observe & compute reward

unverified Policy

Policy deviates from model:
1. Policy is deterministic, verification result is set-valued.

# Safe Reinforcement Learning?



**?safeAccel; accel** U **brake** ≠ **unverified Policy**

Observe & compute reward

Policy deviates from model:

1. Policy is deterministic, verification result is set-valued.

# Physical Models are Approximations

{accel,brake,turn}

**≠ pos'=vel, vel'=acc**

unverified Policy

Observe & compute reward

Policy deviates from model:
1. Policy is deterministic, verification result is set-valued.
2. Environment may not be accurately modeled.

# Safety resolving non-determinism

**?safeAccel;** **accel** ∪ **brake** ≠ [unverified Policy]

# Sandboxing Reinforcement Learning



"Accurate modulo determinism"

```
init → [{ {accel ∪ brake}; t:=0; continuousMotion }*](safe)
```
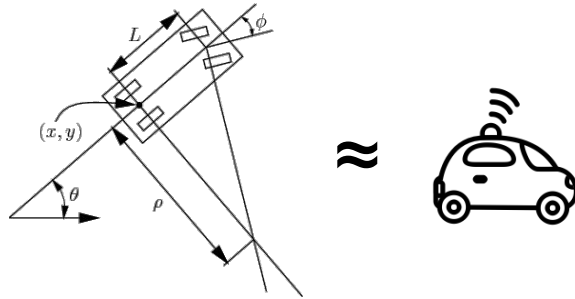
# Sandboxing Reinforcement Learning
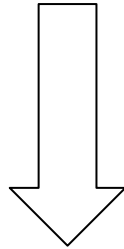


"Accurate modulo determinism"

Learn over a constrained action space

# Sandboxing Reinforcement Learning



≈

"Accurate modulo determinism"

VERIFIED

Learn over a constrained action space

# Sandboxing Reinforcement Learning

**Constrained**

Actions

Observe & compute reward

**Policy** VERIFIED

Theorem: **If the physical model is accurate** then verification results are preserved during learning and by learned policies.

# Sandboxing Reinforcement Learning

**<u>Constrained</u>**
Actions

Observe & compute reward

**Policy**

```
init → [{ {accel ∪ brake}; t:=0; continuousMotion }*](safe)
```

Theorem: **If the physical model is accurate** then verification results are preserved during learning and by learned policies.

# Sandboxing Reinforcement Learning

**Constrained**

Actions

VERIFIED

**Policy**

VERIFIED

Observe & compute reward

$$\text{init} \rightarrow [\{ \{accel \cup brake\}; t:=0; continuousMotion \}*](safe)$$

Theorem: If the physical model is accurate then **verification results are preserved during learning and by learned policies**.

# ~~Sandboxing~~ Safe Reinforcement Learning

**Theorem 1** (JSCGeneric Explores Safely in Modeled Environments). *Assume a valid safety specification*

$$\models \; init \rightarrow [\{ctrl; plant\}^*]safe \qquad (3)$$

*i.e., any repetition of* $\{ctrl; plant\}$ *starting from a state in init will end in a state described by safe. Then* $u_i(s_i) \models safe$ *for all* $u_i, s_i$ *satisfying the learning process for*

$$(init, (S, A, R, E), choose, update, done, CM, MM)$$

*where CM and MM are the controller and model monitor for* $init \rightarrow [\{ctrl; plant\}^*]safe.$

init →                                                                    (safe)

Theorem                                                              ation
results are preserved by learned policies.

# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward
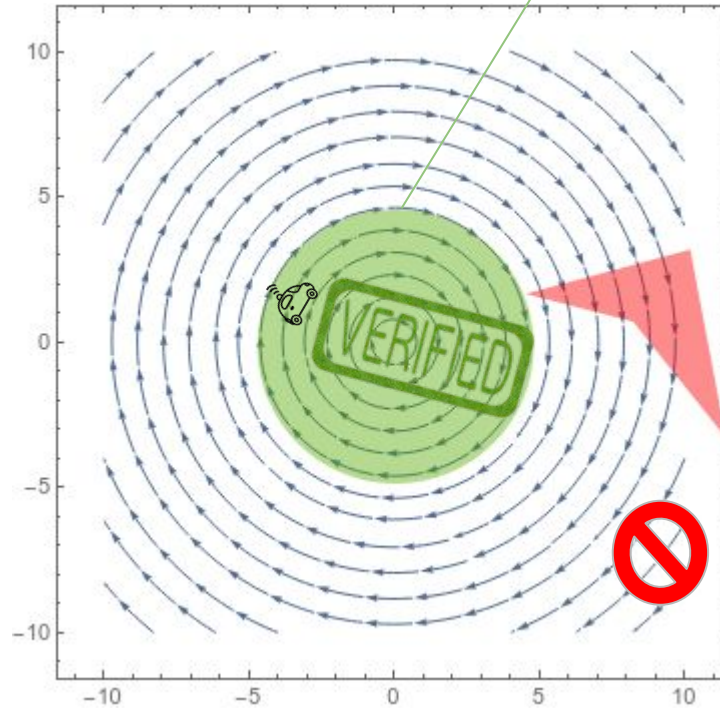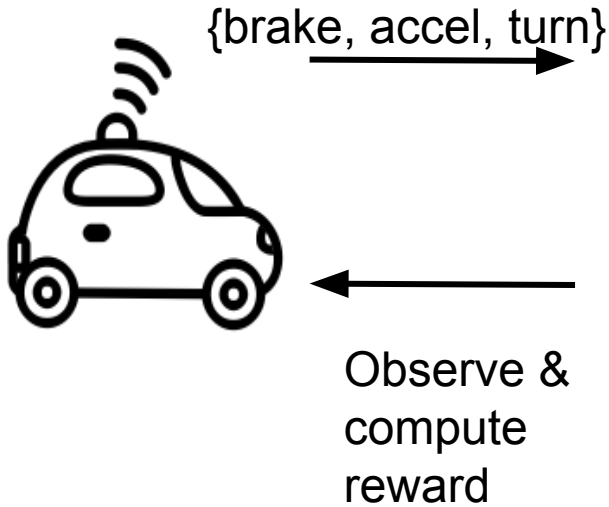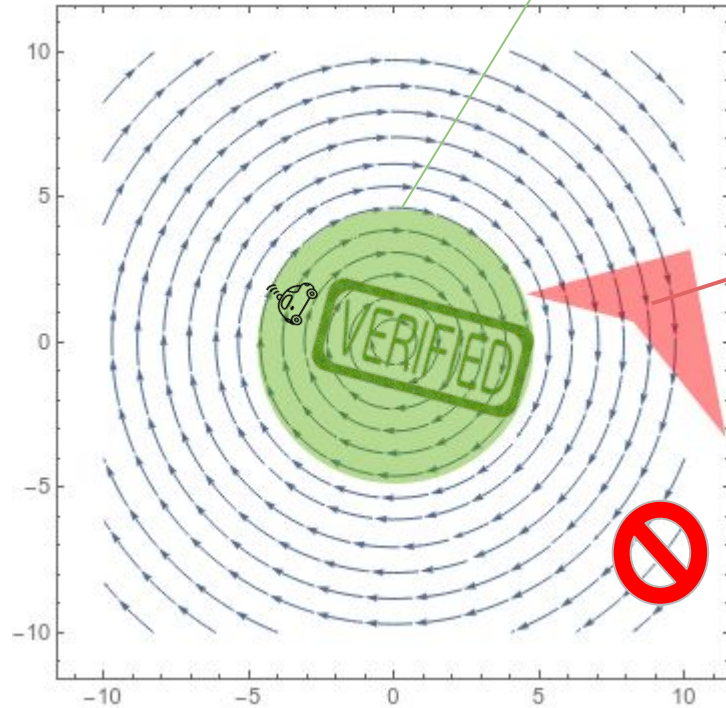
# What About the Physical Model?

**Model is accurate.**

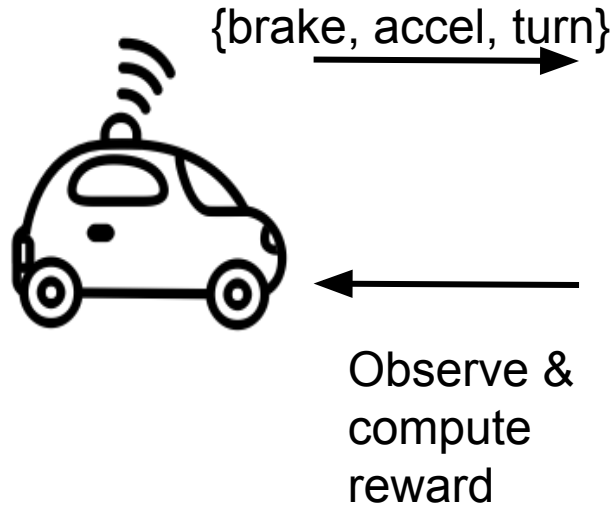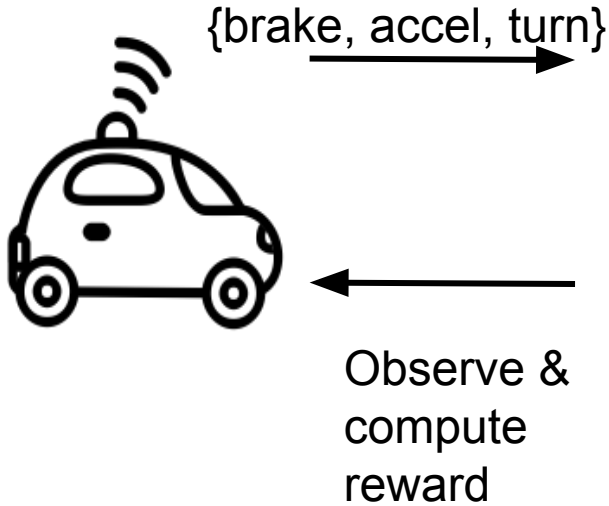{brake, accel, turn}

Observe & compute reward

# What About the Physical Model?

**Model is accurate.**

{brake, accel, turn}

Observe &
compute
reward

# What About the Physical Model?

**Model is correct.**

{brake, accel, turn}

Observe &
compute
reward

**Model is
inaccurate**

# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward

**Model is correct.**

**Model is inaccurate**

**Obstacle!**

# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward

Expected

Reality

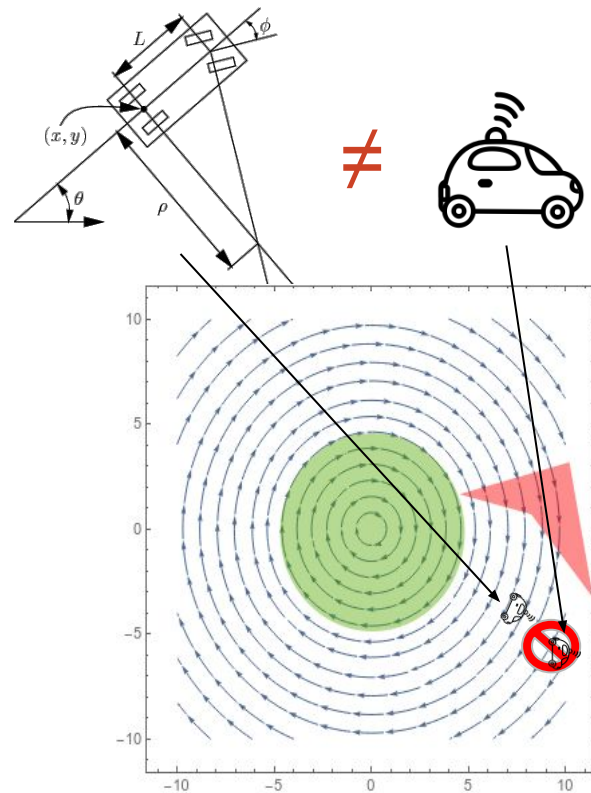# What About the Physical Model?

{brake, accel, turn}

Observe & compute reward

Expected (safe)

Reality (crash!)

# Justified Speculative Control



≈



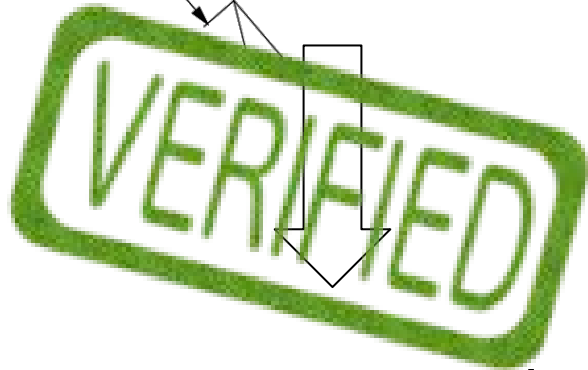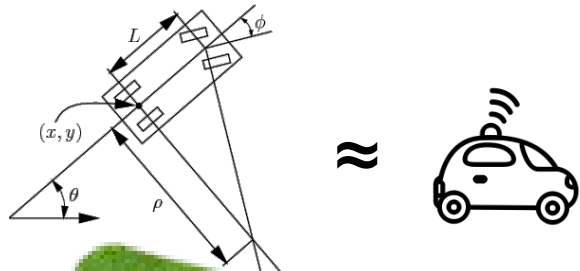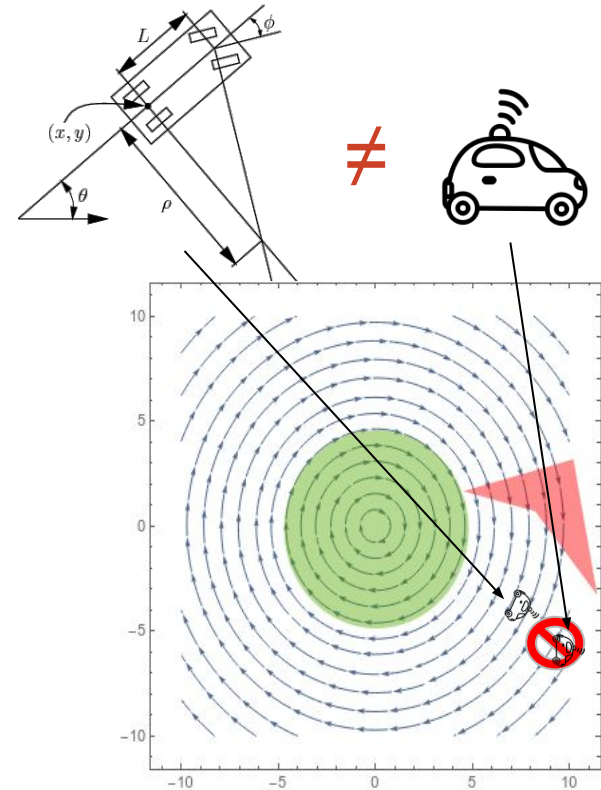VERIFIED

Learn over a constrained action space

≠

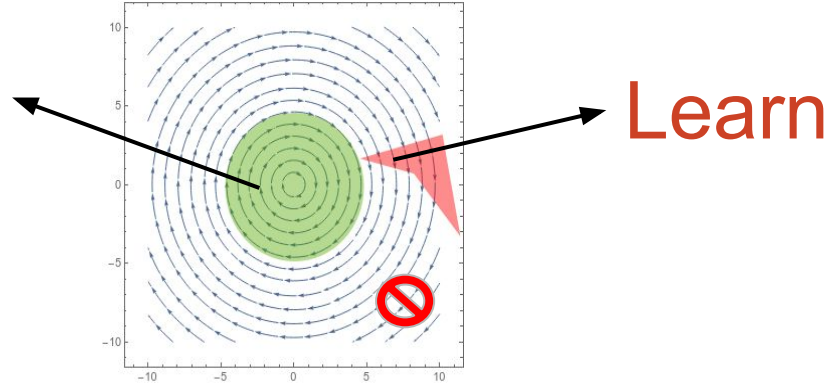# Justified Speculative Control



≈

≠

Learn over a constrained action space

# Justified Speculative Control

Learn over a constrained action space

**VERIFIED**



→ Learn

Some Questions:
1. How do we **know** when we're in unmodeled state space?
2. What do we **do** when we *are* in modeled state space?

# Justified Speculative Control

Learn over a constrained action space **VERIFIED**
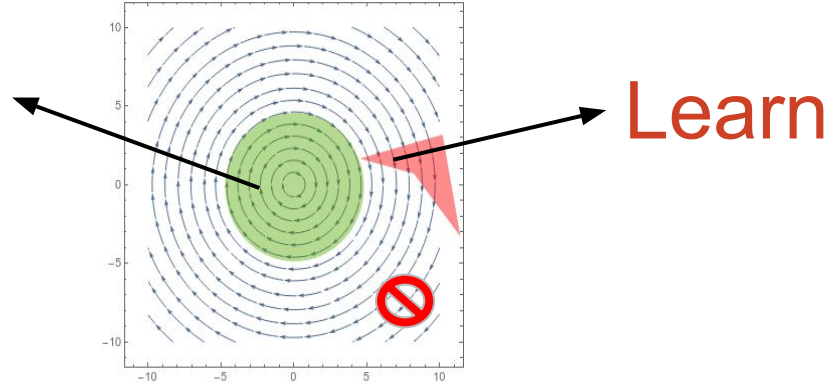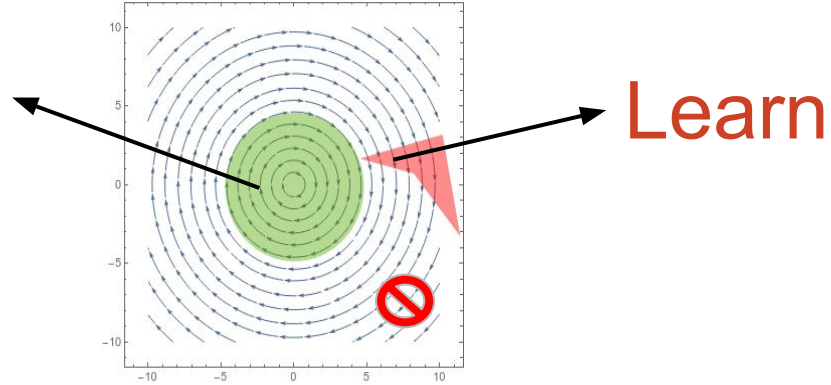


Learn

Some Questions:
1. How do we **know** when we're in unmodeled state space?
2. What do we **do** when we *are* in modeled state space?

# Justified Speculative Control
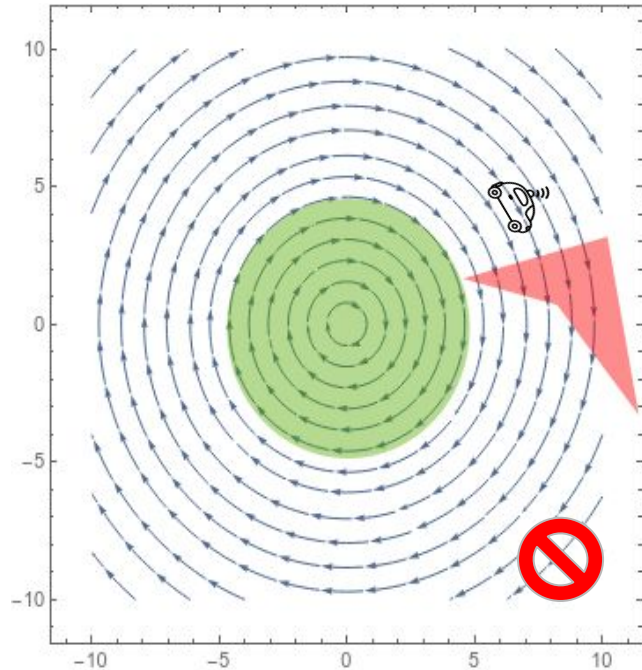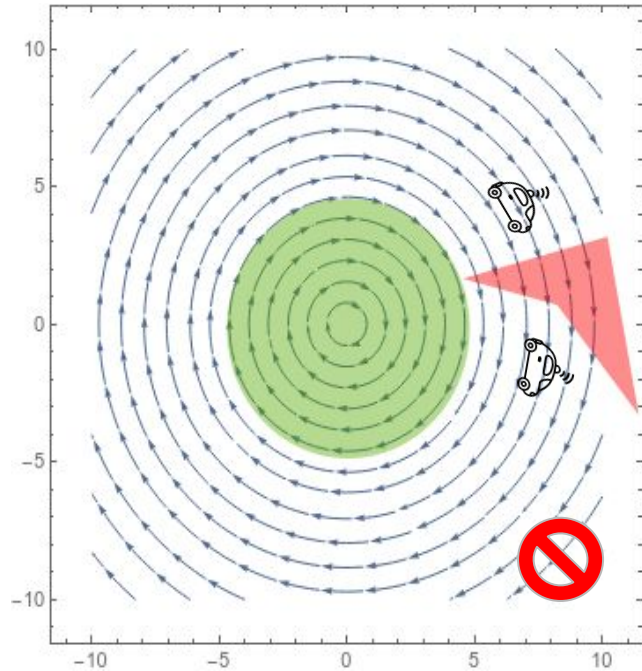
Learn over a constrained action space

VERIFIED

Learn

Theorem: Verification results are preserved outside of red region. But:

☒ How do we know when we're in unmodeled state space?

☐ **What do we do when we *are* in modeled state space?**
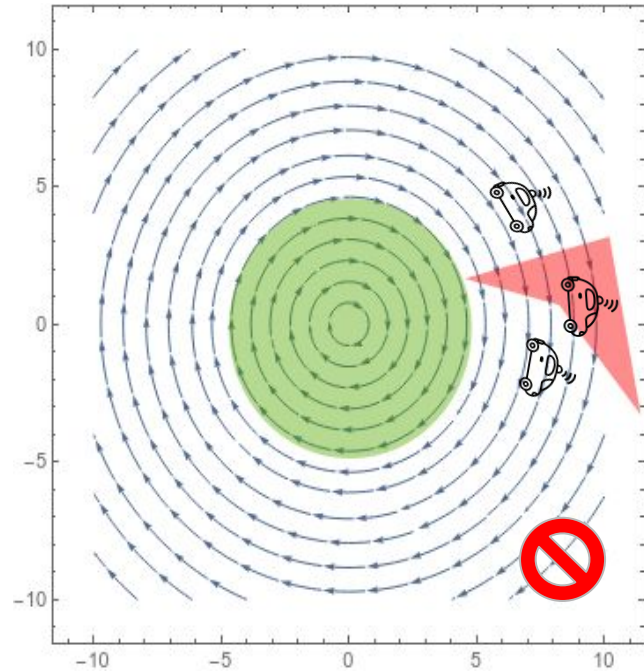
# What do we do in unmodeled state-space?

# What do we do in unmodeled state-space?
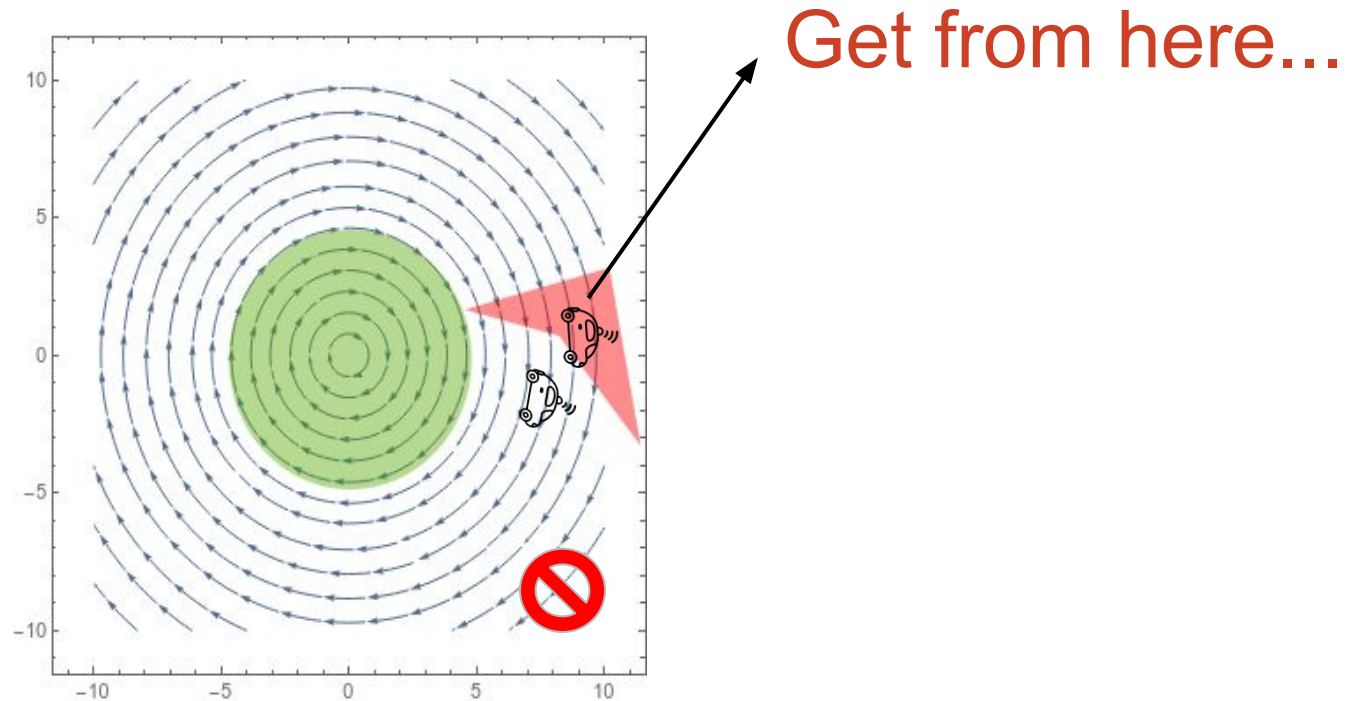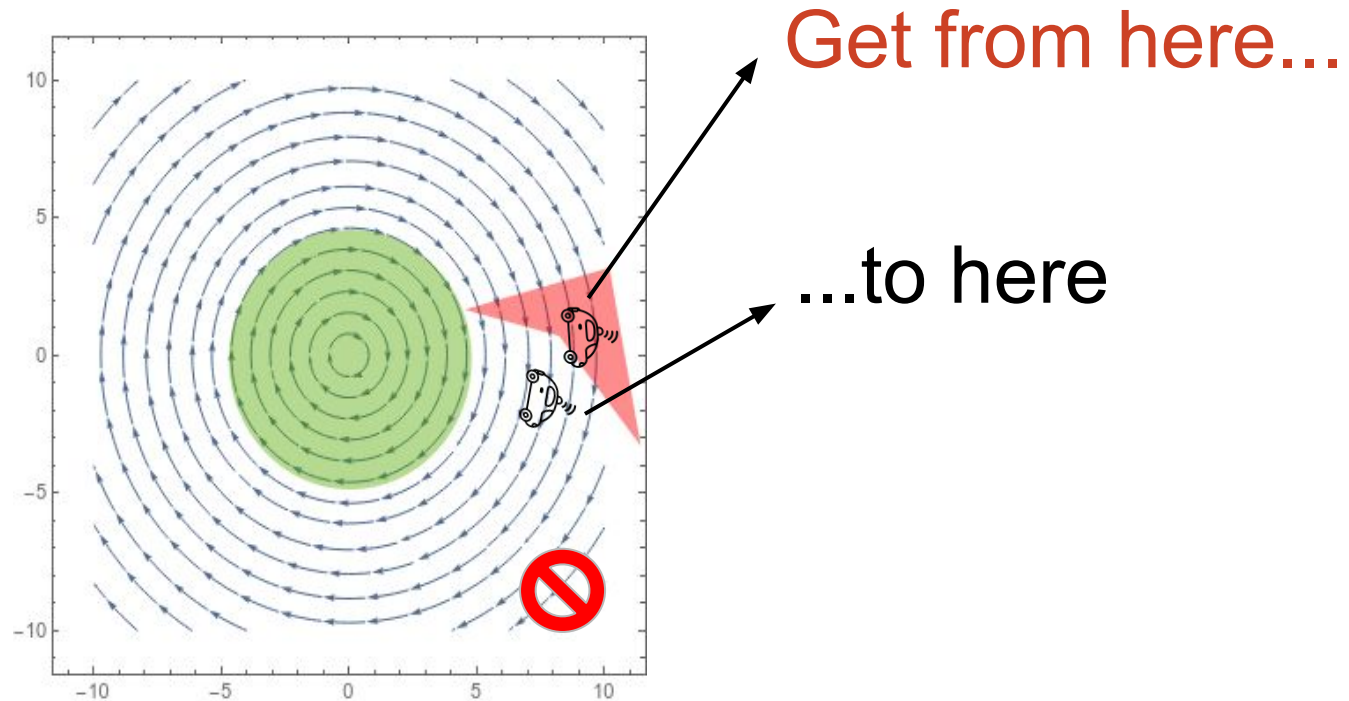
# What do we do in unmodeled state-space?

# What do we do in unmodeled state-space?



Get from here...

# What do we do in unmodeled state-space?



Get from here...

...to here

# Leveraging Formal Methods during Learning

Own Car

Leader

# Leveraging Formal Methods during Learning

Own Car

Leader

| Perturbation | "Don't hit the leader" | "Get back to modeled state space" |
|---|---|---|
| 5% | 3 | 2 |
| 25% | 18 | 16 |
| 50% | 41 | 24 |

# Conclusion

KeYmaera X  + Justified Speculative Control:

1.  Transfer **formal** verification results for **non-deterministic** control policies to policies obtained via a generic reinforcement learning algorithm.
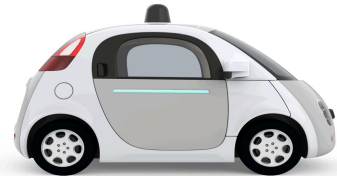
# Conclusion

KeYmaera X + Justified Speculative Control:

1. Transfer **formal** verification results for **non-deterministic** control policies to policies obtained via a generic reinforcement learning algorithm.
2. Leverages insights obtained during verification to direct future learning.

# Model-Based Verification



Reinforcement Learning

pos < stopSign

```
init → [{

    {?safeAccel; accel

      ∪ brake};

    t:=0; {pos'=vel,vel'=acc}

}*](pos < stopSign)
```