# Fine-Tuning Language Models Using Formal Methods Feedback: A Use Case in Autonomous Systems

Yunhao Yang*, Neel P. Bhatt*, Tyler Ingebrand*, William Ward, Steven Carr, Zhangyang Wang, and Ufuk Topcu

May 14th, 2024

**TEXAS**
The University of Texas at Austin

**autonomy.oden.utexas.edu**

**CENTER FOR**
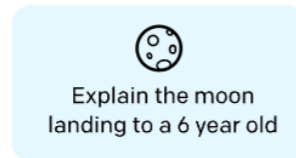**aUTonomy**

* Equal Contribution

# Reinforcement Learning via Human Feedback (RLHF)
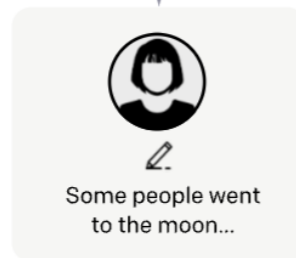
## Example: OpenAI Scheme for Instruct GPT



**Step 1**
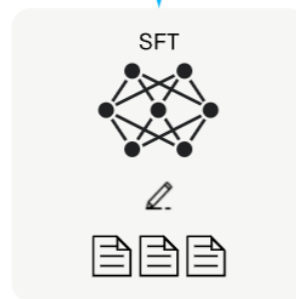
**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

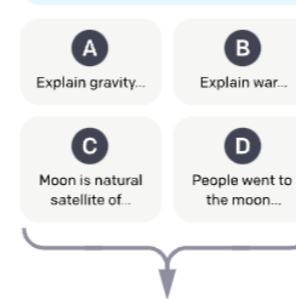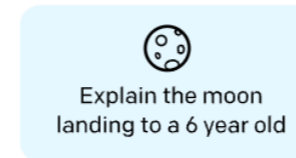Some people went to the moon...
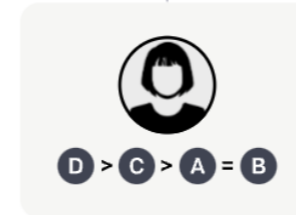
This data is used to fine-tune GPT-3 with supervised learning.

SFT

**Step 2**

**Collect comparison data, and train a reward model.**

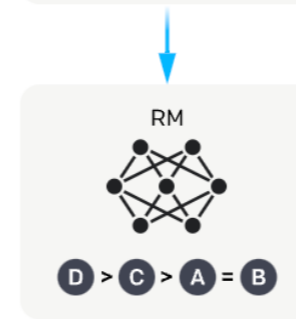A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

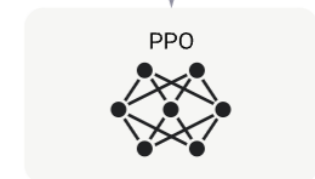This data is used to train our reward model.

RM

D > C > A = B

**Step 3**

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

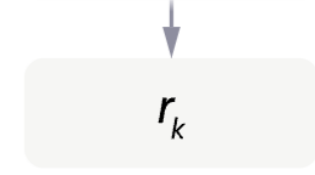Write a story about frogs

The policy generates an output.

PPO

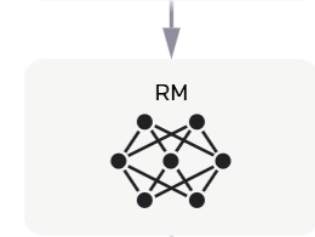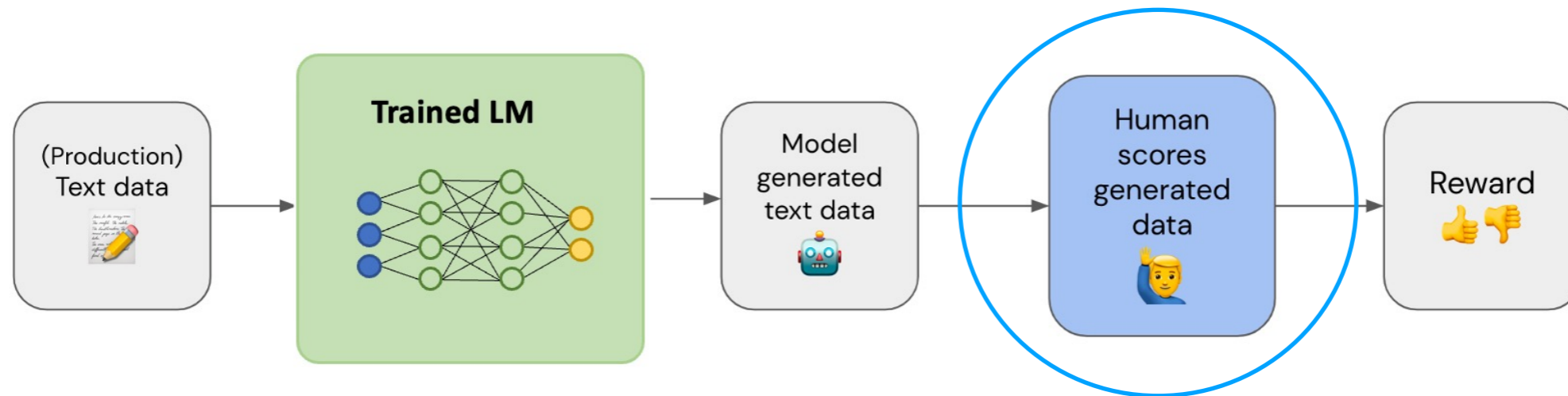Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

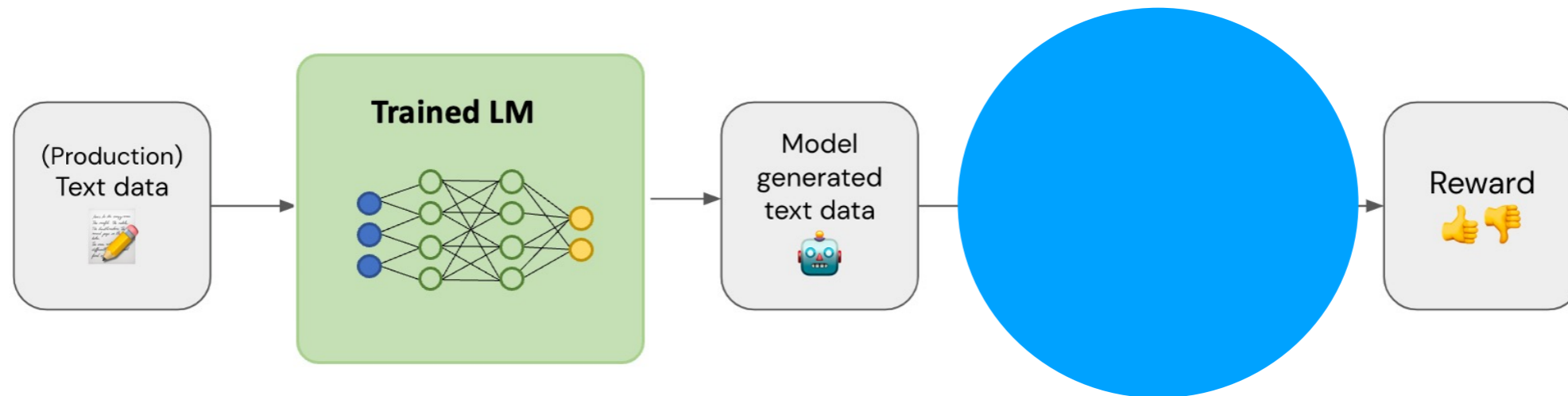# Fine-Tuning Language Models Using Human Feedback



**Labor-Intensive**

**Subjective/Inconsistent Feedback**

# Fine-Tuning Language Models Using Formal Feedback?
### Methods



**Formal Methods:**
 Automaton-Based Representation, Model Checking, Temporal Logic Specification, etc.



**Labor-Intensive** ✔

**Subjective/Inconsistent Feedback** ✔
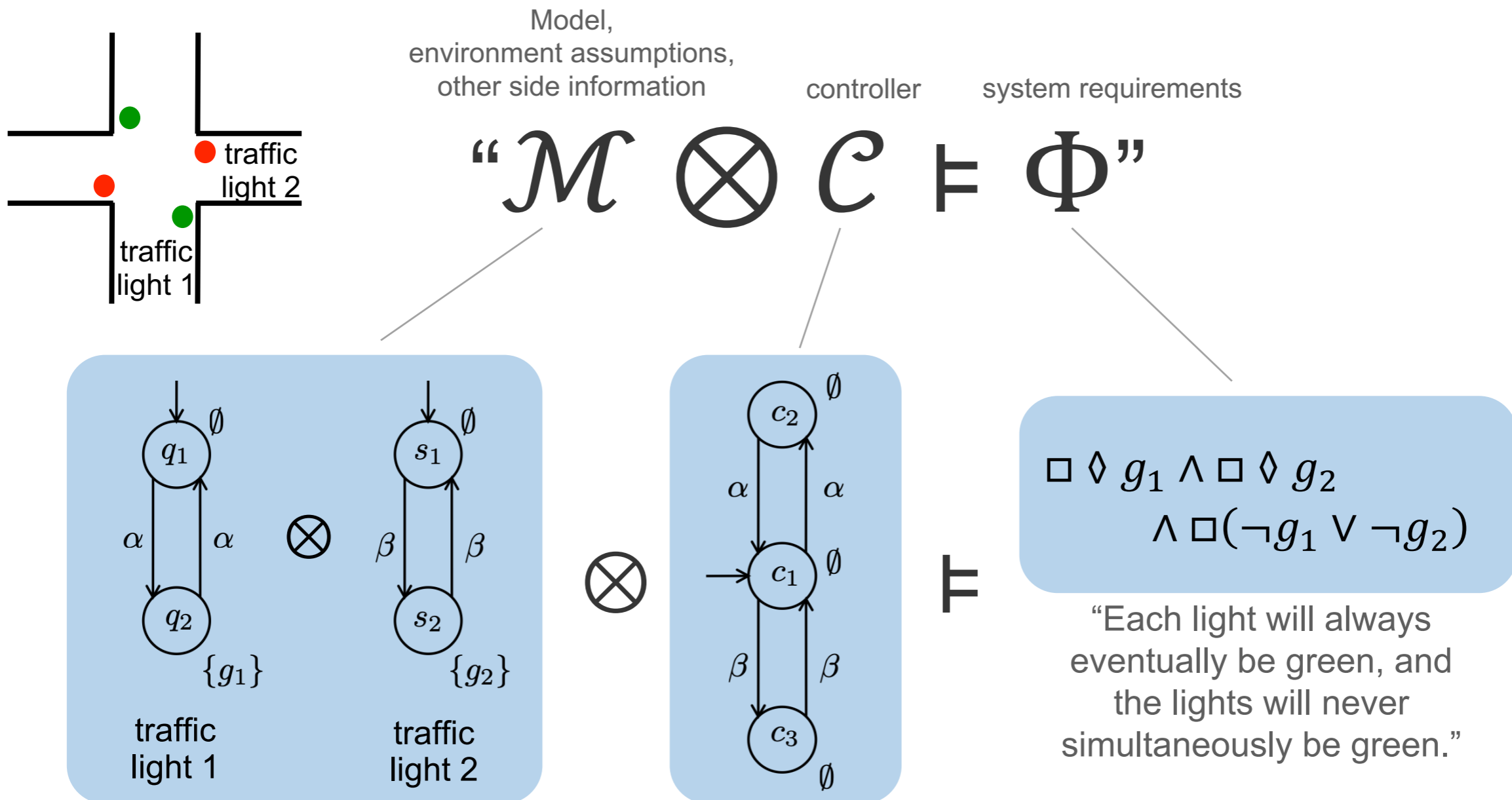
4

# Background: Automaton-Based Representations



Why automaton-based representations? They are used for

- model checking, planning,…
- reactive synthesis, games on graphs, …
- probabilistic verification and synthesis, and
- reinforcement learning.

# A (Very) Brief Introduction to Model Checking

Are the controller's outcomes **guaranteed** to satisfy user-specified requirements when implemented against a system model?



Model, environment assumptions, other side information

controller

system requirements

$$\text{``}\mathcal{M} \otimes \mathcal{C} \models \Phi\text{''}$$

$$\Box \Diamond\, g_1 \wedge \Box \Diamond\, g_2 \wedge \Box(\neg g_1 \vee \neg g_2)$$

"Each light will always eventually be green, and the lights will never simultaneously be green."
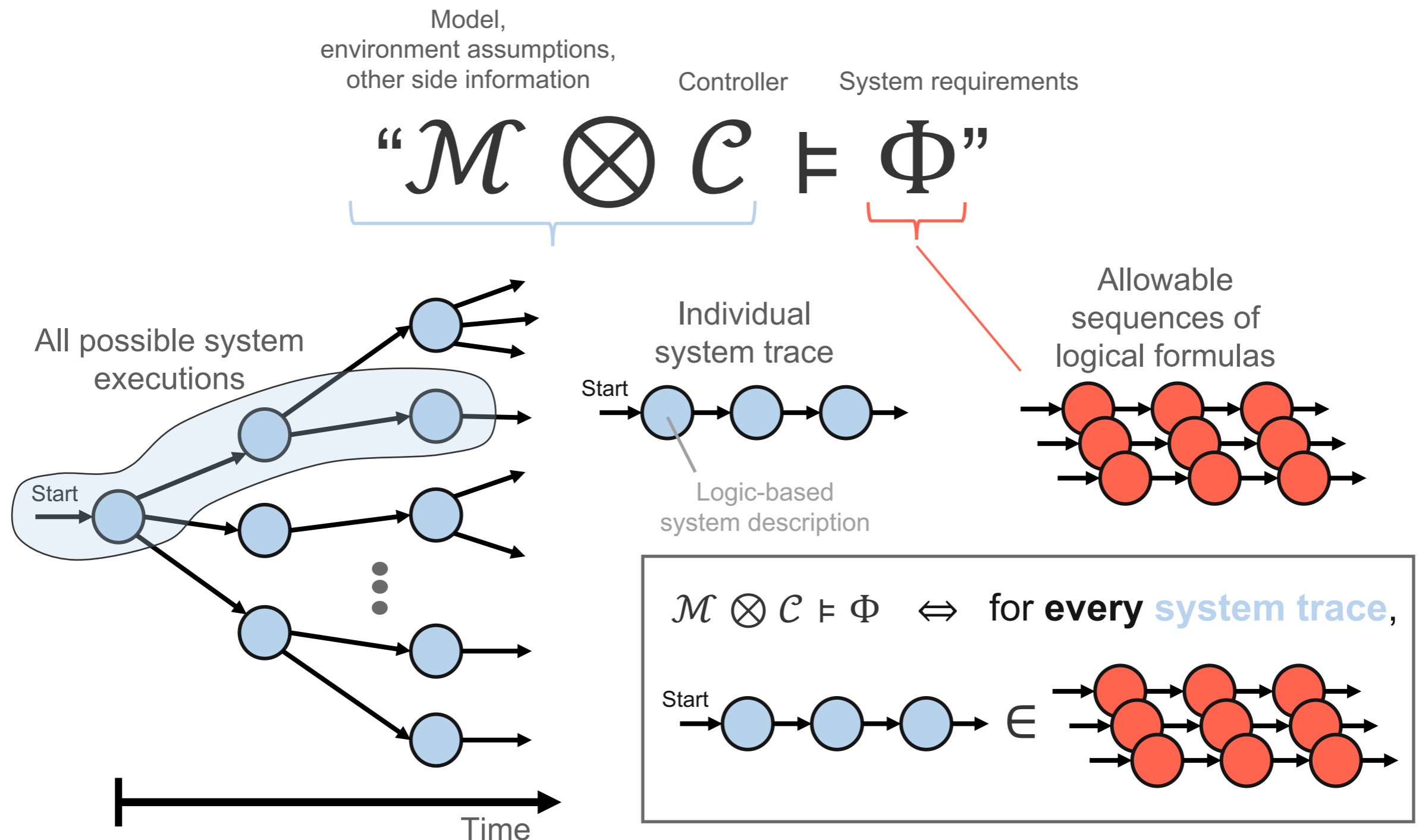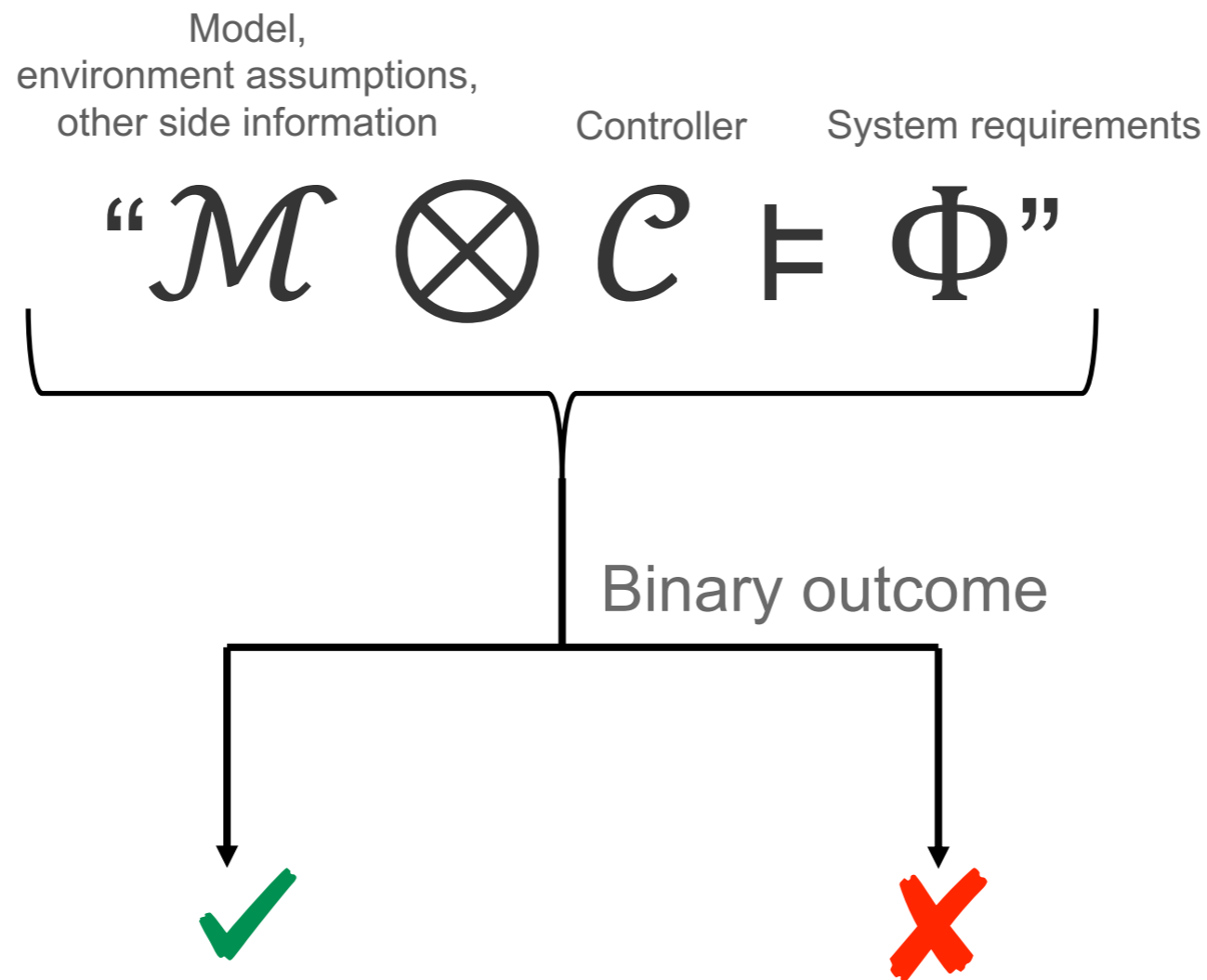
traffic light 1

traffic light 2

# A (Very) Brief Introduction to Model Checking

Are the controller's outcomes **guaranteed** to satisfy user-specified requirements when implemented against a system model?

Model, environment assumptions, other side information

Controller

System requirements

$$ ``\mathcal{M} \otimes \mathcal{C} \models \Phi" $$



All possible system executions

Start

Time

Individual system trace

Start

Logic-based system description

Allowable sequences of logical formulas

$$ \mathcal{M} \otimes \mathcal{C} \models \Phi \iff \text{for \textbf{every} system trace,} $$
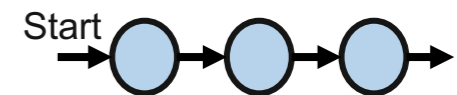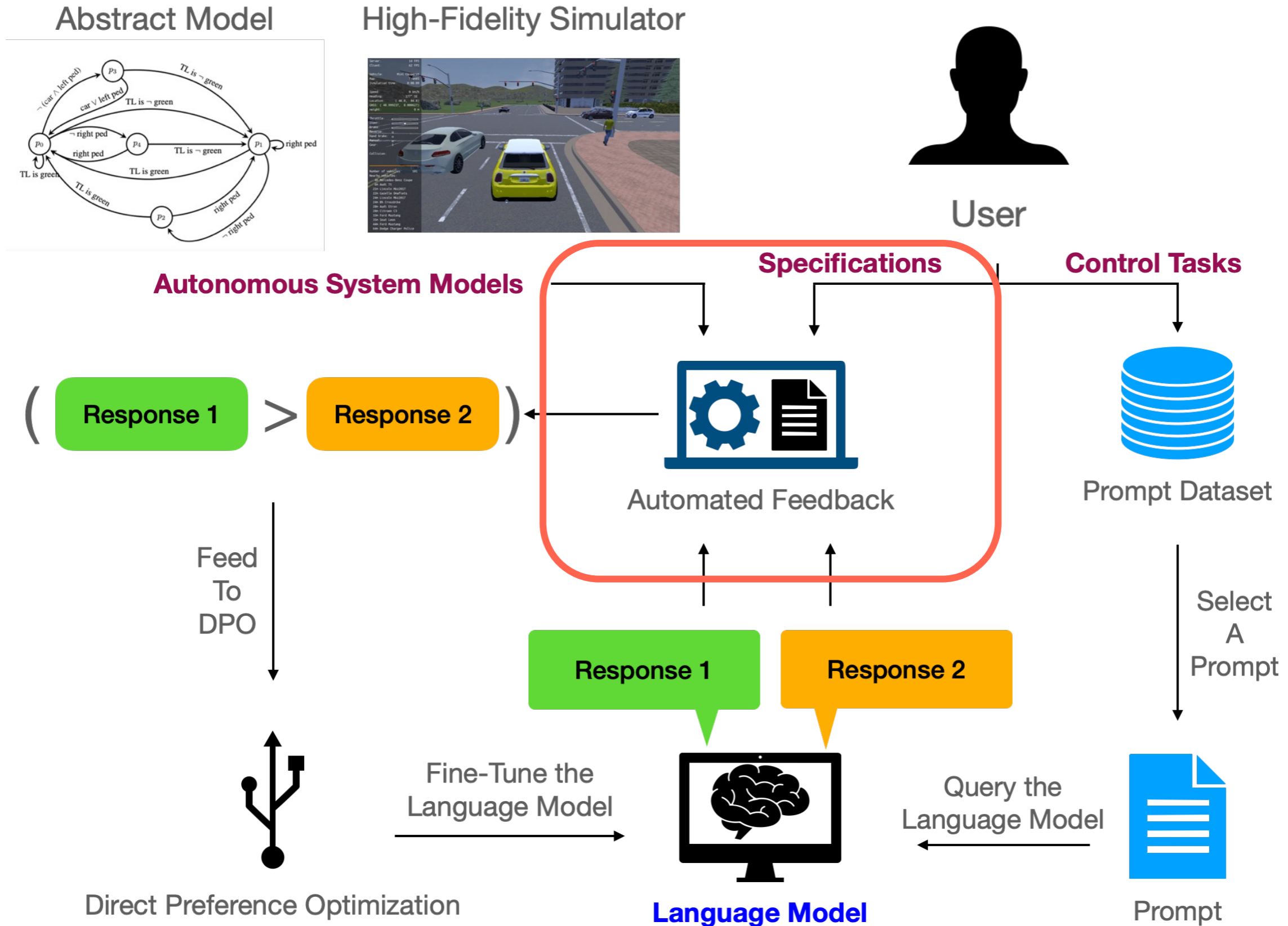
Start

$\in$

# A (Very) Brief Introduction to Model Checking

Are the controller's outcomes **guaranteed** to satisfy user-specified requirements when implemented against a system model?

Model,
environment assumptions,
other side information

Controller

System requirements

$$``\mathcal{M} \otimes \mathcal{C} \vDash \Phi"$$

Binary outcome

✔ ✘

Byproduct: Counterexample trace
that **violates** the specification.

Start

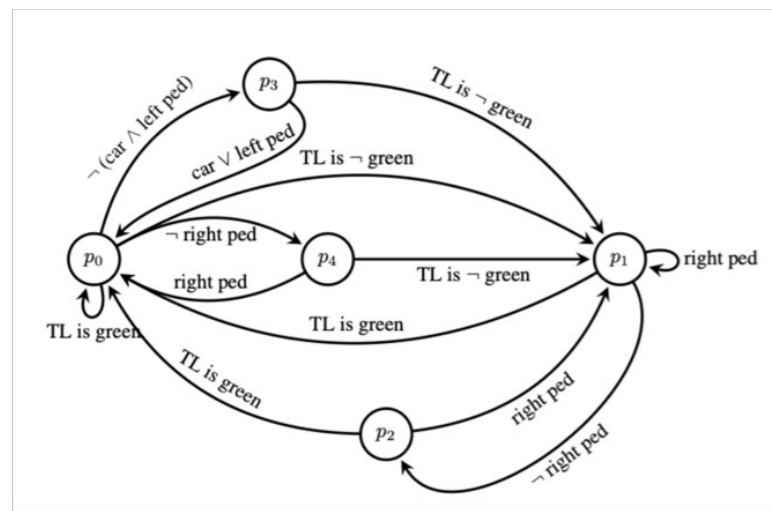# Fine-Tuning Language Models Using Formal Methods Feedback



Abstract Model

High-Fidelity Simulator

User

Autonomous System Models

Specifications

Control Tasks

( Response 1 > Response 2 )

Automated Feedback

Prompt Dataset

Feed To DPO

Response 1    Response 2

Select A Prompt

Fine-Tune the Language Model

Query the Language Model

Direct Preference Optimization

**Language Model**

Prompt

# Fine-Tuning Language Models Using Formal Methods Feedback

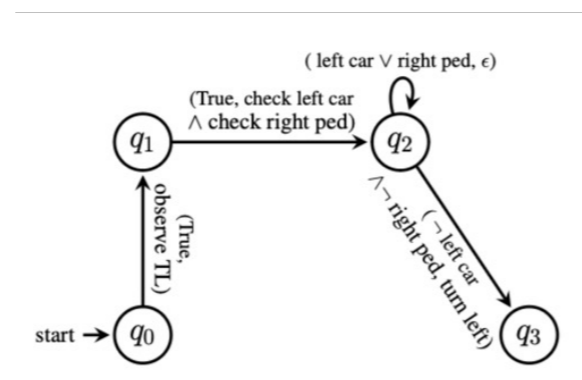**Modeling the Autonomous System**
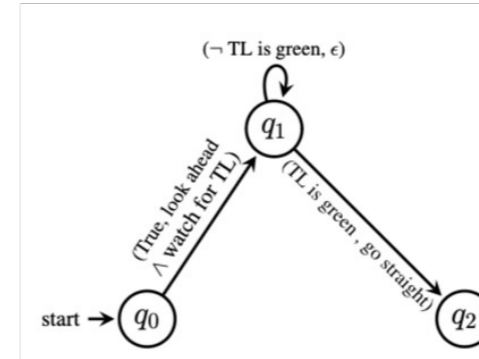


Autonomous System

Autonomous System Model $\mathcal{M}$
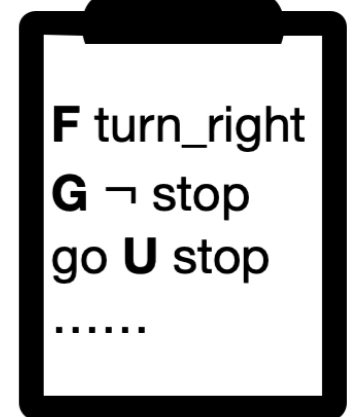
**Controller Construction**

Response 1

Response 2

Controller $\mathcal{C}_1$

Controller $\mathcal{C}_2$

User

**F** turn_right
**G** ¬ stop
go **U** stop
......

A Set of Specifications
$\{\Phi_1, \ldots, \Phi_n\}$

$$\mathcal{M} \otimes \mathcal{C} \models \Phi$$

Model Checker

Compare the # of
Satisfied Specifications

$\left( \begin{array}{ccc} \text{Response 1} & > & \text{Response 2} \end{array} \right)$
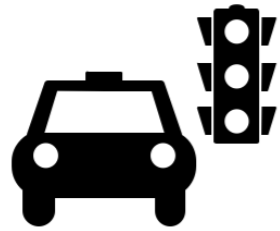
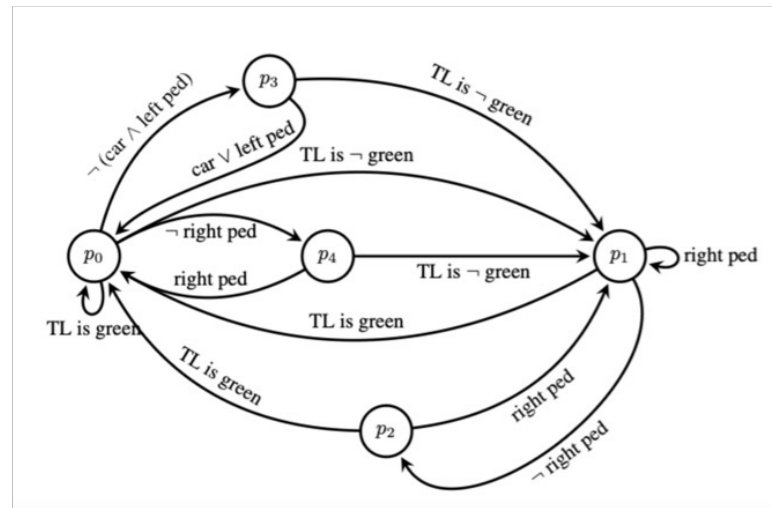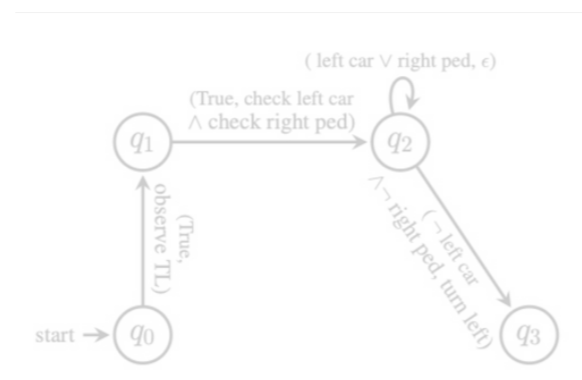$\Phi_1, \Phi_2, \Phi_3$    $\Phi_1, \Phi_2, \Phi_3$

**Formal Verification**

# Fine-Tuning Language Models Using Formal Methods Feedback

**Modeling the Autonomous System**
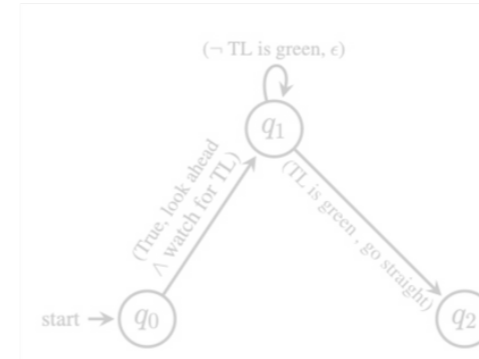


Autonomous System

Autonomous System Model $\mathcal{M}$

Controller Construction

Response 1

Response 2
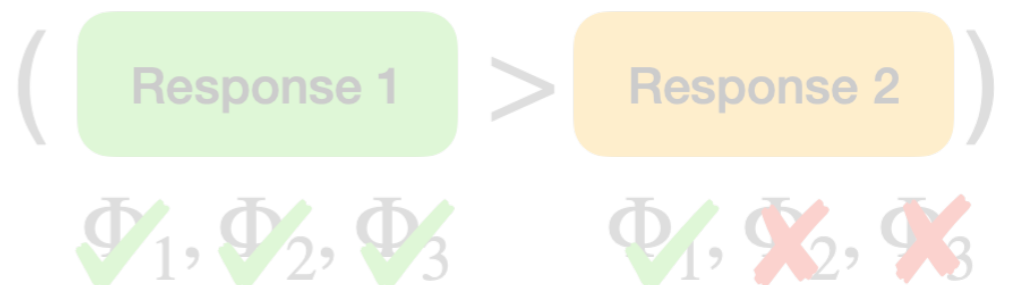
User

Controller $\mathscr{C}_1$

Controller $\mathscr{C}_2$

**F** turn_right
**G** ¬ stop
go **U** stop
......

A Set of Specifications
$\{\Phi_1, \ldots, \Phi_n\}$

$\mathcal{M} \otimes \mathscr{C} \vDash \Phi$

Model Checker

Compare the # of
Satisfied Specifications

Formal Verification

$\left( \text{Response 1} > \text{Response 2} \right)$

$\checkmark\Phi_1, \checkmark\Phi_2, \checkmark\Phi_3$    $\checkmark\Phi_1, \times\Phi_2, \times\Phi_3$

# Formal Methods Feedback in an Autonomous Driving System

## Modeling the Autonomous System
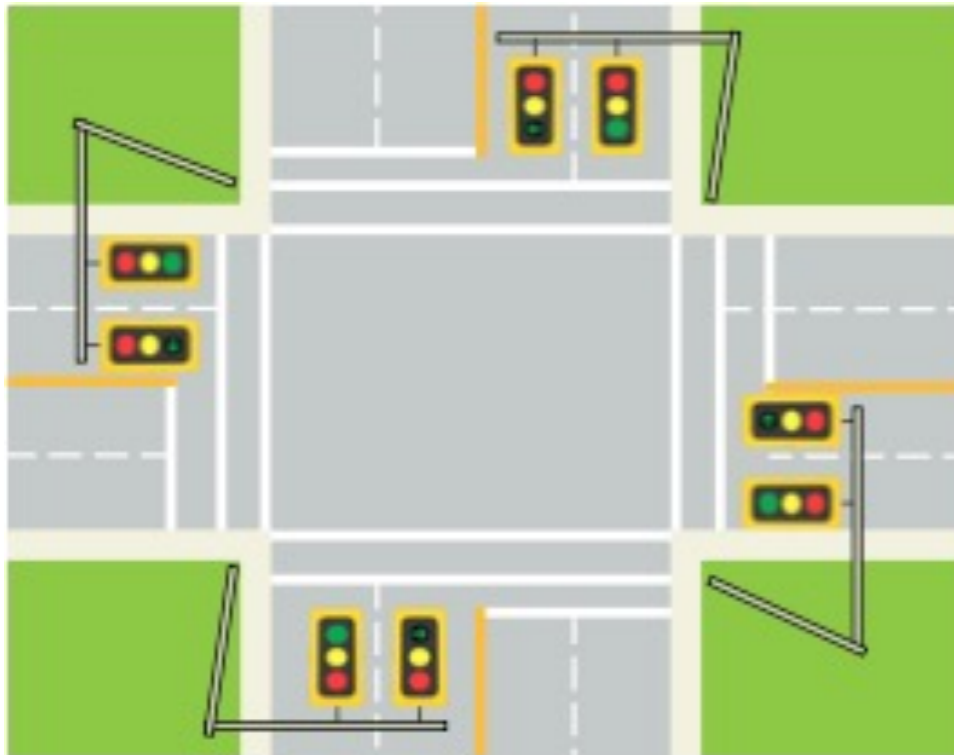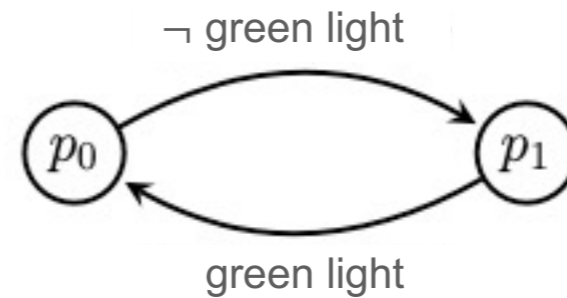
# Formal Methods Feedback in an Autonomous Driving System

## Modeling the Autonomous System

$P = \{green\ light\}$
$S \sim Traffic\ Light$
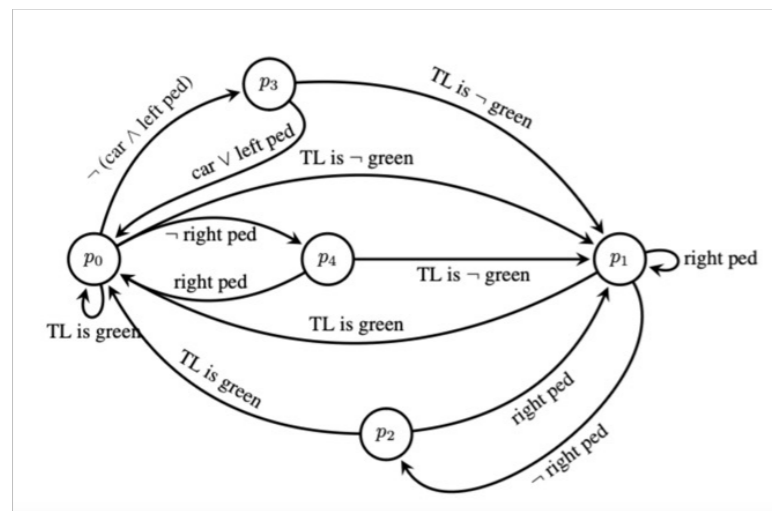
$Q_M = \{p_0\ p_1\}$
$\lambda_M(p_0) = green\ light$
$\lambda_M(p_1) = \neg\ green\ light$
$\delta_M(p_0, p_1) = 1$
$\delta_M(p_1, p_0) = 1$



¬ green light

$p_0$    $p_1$

green light

# Fine-Tuning Language Models Using Formal Methods Feedback

**Modeling the Autonomous System**

**Controller Construction**

Autonomous System

**Response 1**

**Response 2**

User



Autonomous System Model $\mathcal{M}$

Controller $\mathcal{C}_1$

Controller $\mathcal{C}_2$

**F** turn_right
**G** ¬ stop
go **U** stop
......

A Set of Specifications
$\{\Phi_1, \ldots, \Phi_n\}$

Compare the # of
Satisfied Specifications

$\mathcal{M} \otimes \mathcal{C} \models \Phi$

$\left( \quad \text{Response 1} \quad > \quad \text{Response 2} \quad \right)$

$\Phi_1, \Phi_2, \Phi_3$

$\Phi_1, \Phi_2, \Phi_3$

Model Checker

Formal Verification

# Formal Methods Feedback in an Autonomous Driving System

## Controller Construction

**User**

> Steps for turning right at the traffic light

**Language Model**

> 1. Look straight ahead and watch for traffic light.
> 2. If the traffic light turns green, start moving forward.
> 3. As you approach the intersection, look to your left for oncoming traffic.
> 4. If there is no traffic coming from your left, check pedestrians on your right.
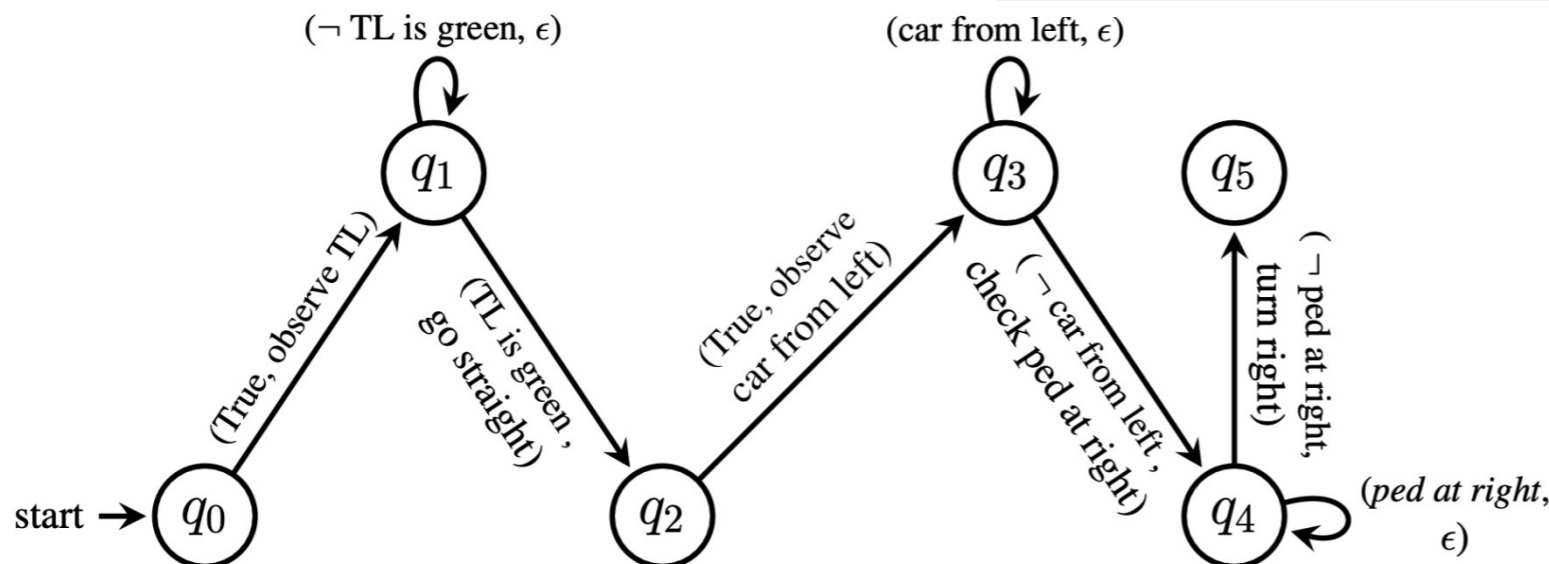> 5. If it is safe, turn your vehicle right.

> 1. Observe the traffic light in front of you.
> 2. Check for the left approaching car and right side pedestrian.
> 3. If no car from the left is approaching and no pedestrian on the right, proceed to turn right.

15

## Controller Construction

1. Look straight ahead and watch for traffic light.
2. If the traffic light turns green, start moving forward.
3. As you approach the intersection, look to your left for oncoming traffic.
4. If there is no traffic coming from your left, check pedestrians on your right.
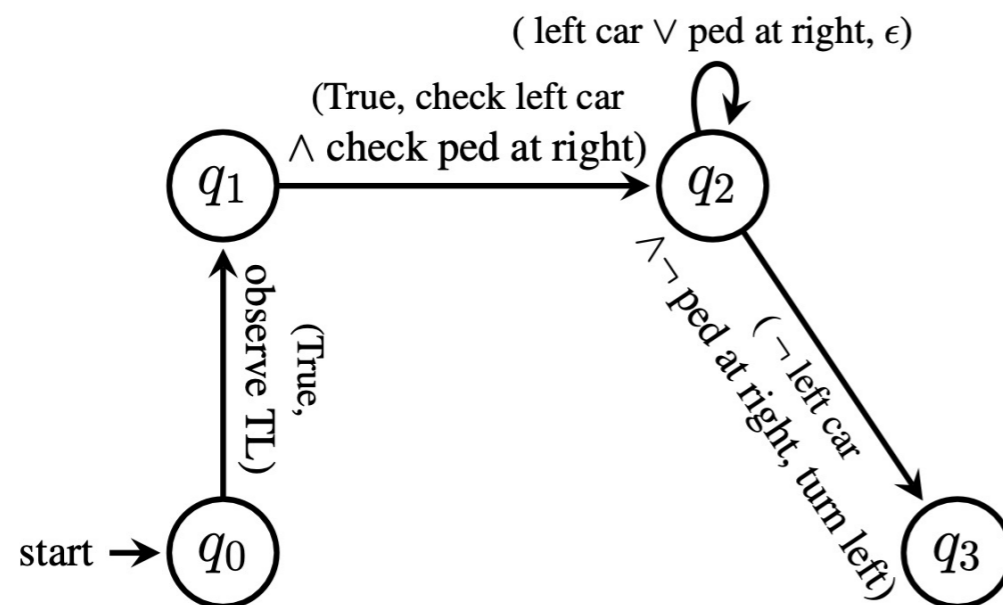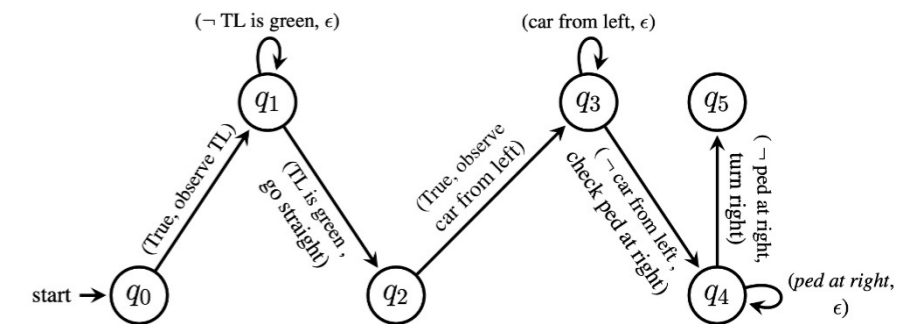5. If it is safe, turn your vehicle right.



| Grammar | Transition Rule |
|---|---|
| $VP^A$ | $q_i \xrightarrow{(True, VP^A)} q_{i+1}$ |
| $VP^A \leftarrow VP^A [j]$ | $q_i \xrightarrow{(True, \epsilon)} q_j$ |
| **if** $VP^C$, $VP^A$ <br> $VP^A$ **if** $VP^C$ | $(\neg VP^C, \epsilon) \circlearrowleft q_i \xrightarrow{(VP^C, VP^A)} q_j$ |
| **if** $VP^C$, $VP^A_1$. **if** $\neg VP^C$, $VP^A_2$ <br> **if** $VP^C$ $VP^A_1$ **else** $VP^A_2$ <br> $VP^A_1$ **if** $VP^C$, **else** $VP^A_2$ | $q_k \xleftarrow{(\neg VP^C, VP^A_2)} q_i \xrightarrow{(VP^C, VP^A_1)} q_j$ |
| $VP^A \leftarrow$ **wait** $VP^C$ $VP^A$ <br> $VP^A \leftarrow VP^A$ **after** $VP^C$ | $(\neg VP^C, \epsilon) \circlearrowleft q_i \xrightarrow{(VP^C, VP^A)} q_{i+1}$ |
| $VP^A \leftarrow VP^A$ **until** $VP^C$ | $(\neg VP^C, VP^A) \circlearrowleft q_i \xrightarrow{(VP^C, \epsilon)} q_{i+1}$ |

1. `<obse`...VP^A...`traight>`.
2. `<if>`
3. `<obse`...
4. `<if>`...`edestrian`
   `at`
5. `<if>`...`rn`
   `righ`...

# Formal Methods Feedback in an Autonomous Driving System

## Controller Construction

1. Look straight ahead and watch for traffic light.
2. If the traffic light turns green, start moving forward.
3. As you approach the intersection, look to your left for oncoming traffic.
4. If there is no traffic coming from your left, check pedestrians on your right.
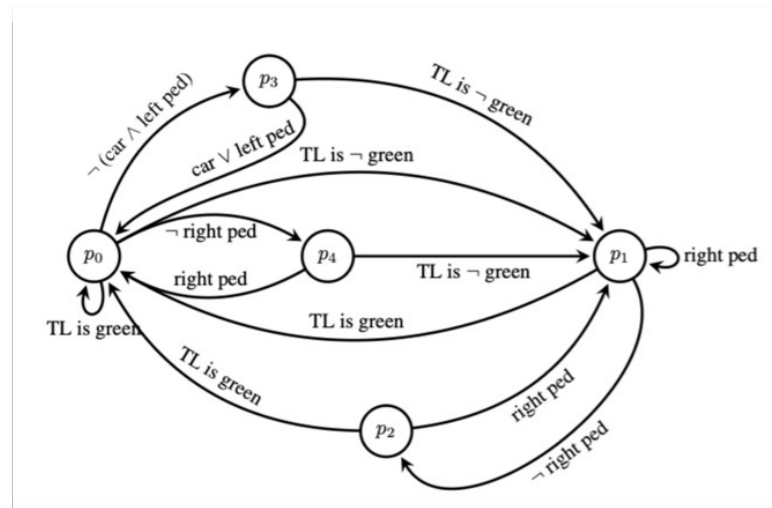5. If it is safe, turn your vehicle right.





1. Observe the traffic light in front of you.
2. Check for the left approaching car and right side pedestrian.
3. If no car from the left is approaching and no pedestrian on the right, proceed to turn right.

# Fine-Tuning Language Models Using Formal Methods Feedback
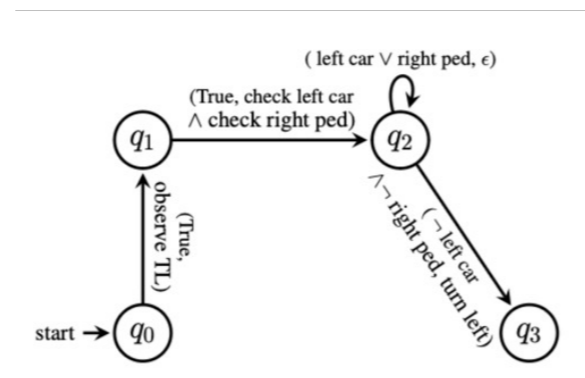
Modeling the Autonomous System

Controller Construction



Autonomous System

Response 1

Response 2

User

Autonomous System Model $\mathcal{M}$

Controller $\mathscr{C}_1$

Controller $\mathscr{C}_2$

**F** turn_right
**G** ¬ stop
go **U** stop
......

A Set of Specifications
$\{\Phi_1, \ldots, \Phi_n\}$

$$\mathcal{M} \otimes \mathscr{C} \vDash \Phi$$

Compare the # of
Satisfied Specifications

$\left(\ \text{Response 1} > \text{Response 2}\ \right)$

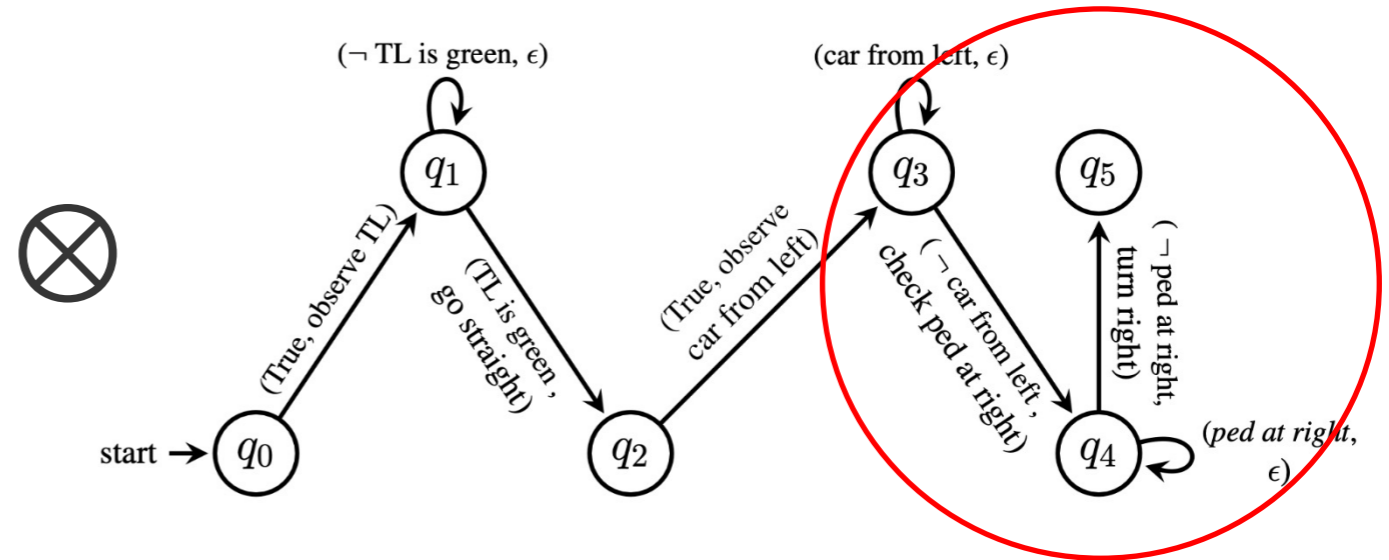$\Phi_1, \Phi_2, \Phi_3$

$\Phi_1, \Phi_2, \Phi_3$
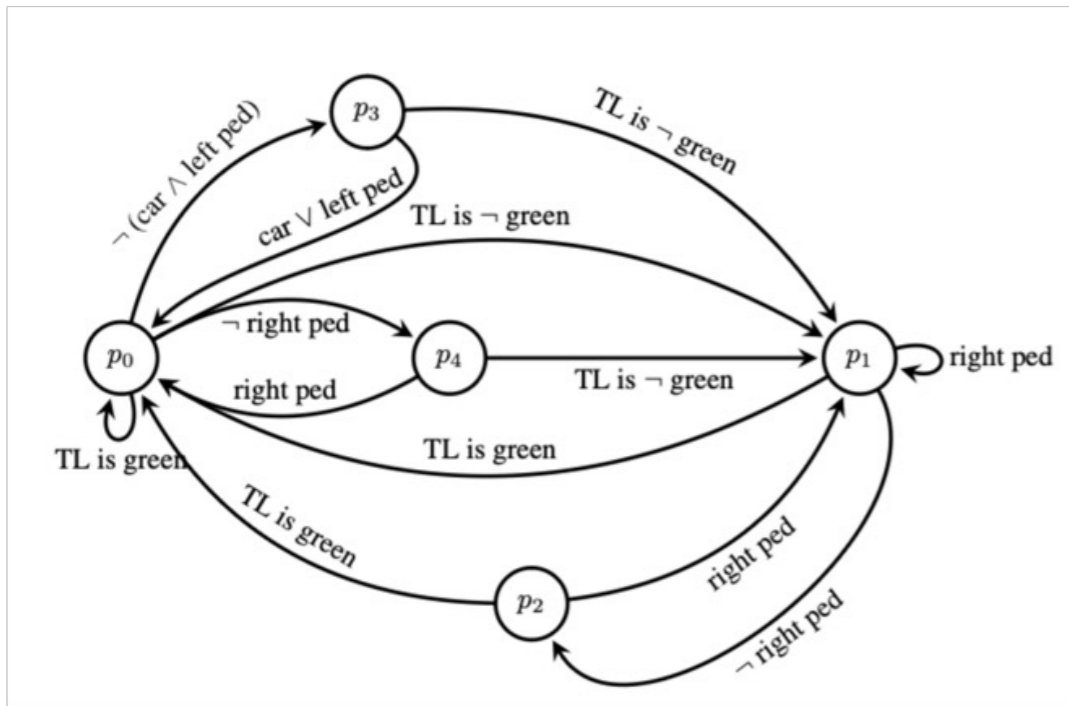
Model Checker

**Formal Verification**

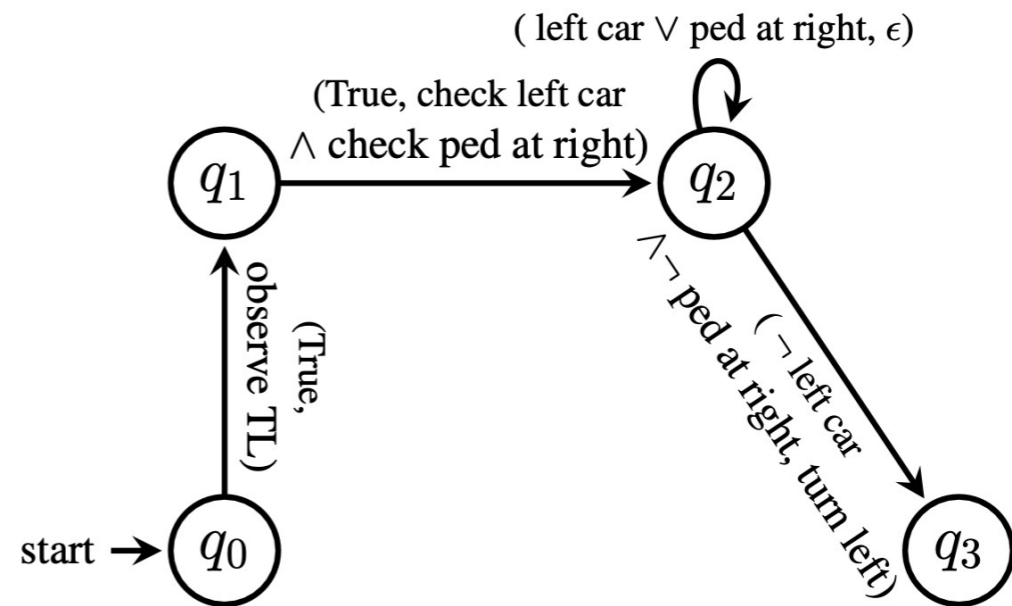# Formal Methods Feedback in an Autonomous Driving System

## Formal Verification



$$\Box(\neg\text{green traffic light} \rightarrow \neg\text{go straight}), \; \checkmark$$

$$\Box(\text{stop sign} \rightarrow \Diamond\text{stop}), \; \checkmark$$

$$\Box\neg\text{turn right} \lor \neg(\text{car from left} \lor \text{pedestrian at right}), \; \times$$

## Formal Verification



$$\Box(\neg\text{green traffic light} \rightarrow \neg\text{go straight}),$$

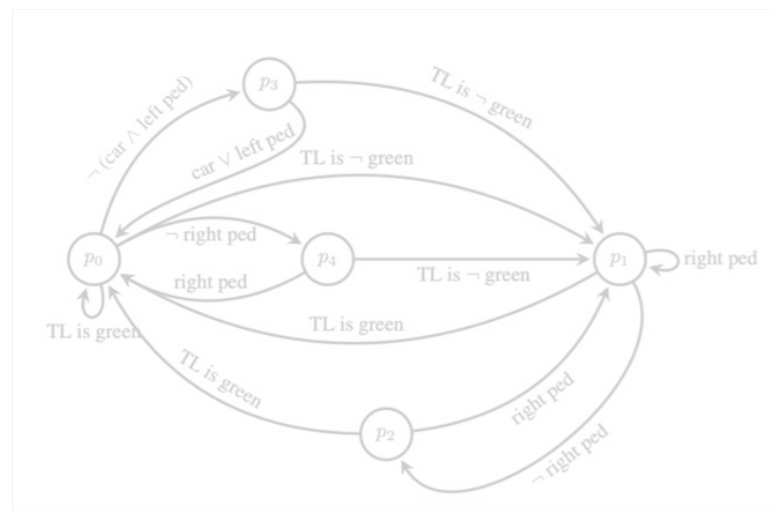$$\Box(\text{stop sign} \rightarrow \Diamond\text{stop}),$$

$$\Box\neg\text{turn right} \vee \neg(\text{car from left} \vee \text{pedestrian at right}),$$

# Fine-Tuning Language Models Using Formal Methods Feedback
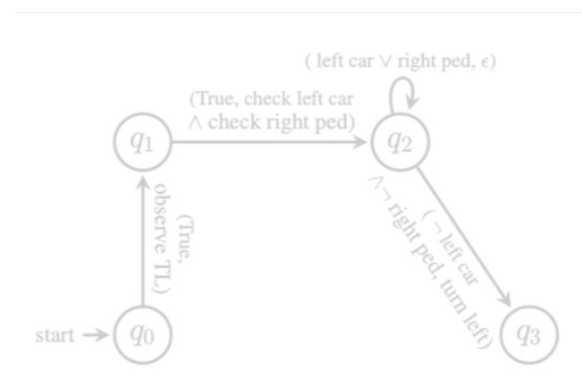


**Modeling the Autonomous System**

Autonomous System

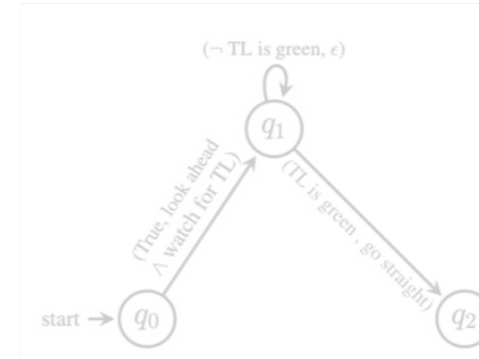Autonomous System Model $\mathcal{M}$

**Controller Construction**
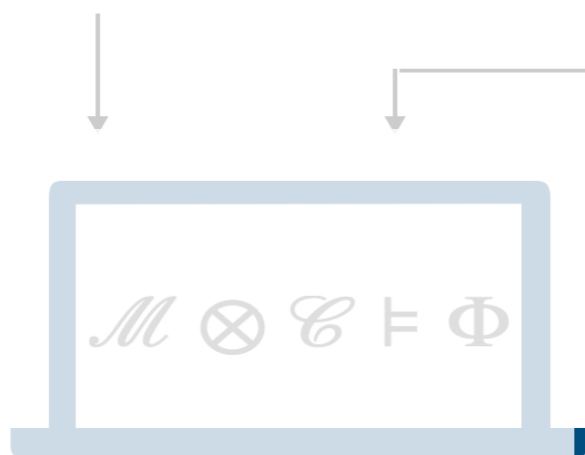
Response 1

Response 2

Controller $\mathcal{C}_1$

Controller $\mathcal{C}_2$

User

$\mathbf{F}$ turn_right
$\mathbf{G} \neg$ stop
go $\mathbf{U}$ stop
......

A Set of Specifications
$\{\Phi_1, \ldots, \Phi_n\}$

$$\mathcal{M} \otimes \mathcal{C} \models \Phi$$

Model Checker

Compare the # of
Satisfied Specifications

$\left( \quad \text{Response 1} \quad < \quad \text{Response 2} \quad \right)$

$\Phi_1, \Phi_2, \Phi_3$    $\Phi_1, \Phi_2, \Phi_3$

# Fine-Tuning Language Models Using Formal Methods Feedback

Abstract Model

High-Fidelity Simulator

User

**Autonomous System Models**

**Specifications**

**Control Tasks**

( Response 1 < Response 2 ) ← Automated Feedback

Prompt Dataset

Feed To DPO

Response 1    Response 2

Select A Prompt

Direct Preference Optimization

Fine-Tune the Language Model → **Language Model**

Query the Language Model

Prompt

# Fine-Tuning Language Models Using Formal Methods Feedback

High-Fidelity Simulator

User

**Autonomous System Models**

**Specifications**　　**Control Tasks**

( **Response 1** < **Response 2** )

Automated Feedback

Prompt Dataset

Feed To DPO

Select A Prompt

**Response 1**　　**Response 2**

Fine-Tune the Language Model

Query the Language Model

Direct Preference Optimization

**Language Model**

Prompt

# Fine-Tuning Language Models Using Formal Methods Feedback

## Empirical Evaluation via Simulation



Response 1

Response 2

Controller $\mathcal{C}_1$

Controller $\mathcal{C}_2$

High-Fidelity Simulator

Empirically Collect Execution Info

Information from the Simulator

Execution Trace
$(left\_car \wedge go\_straight, pedestrian \wedge stop, \neg pedestrian \wedge go\_straight)$

Execution Traces

Execution Traces

# Fine-Tuning Language Models Using Formal Methods Feedback

**Empirical Evaluation via Simulation**

$$\left( \boxed{\text{Response 1}} < \boxed{\text{Response 2}} \right)$$

Empirically Collect Execution Info

Execution Traces

Execution Traces

Verify Traces

**F** turn_right
**G** ¬ stop
go **U** stop
......

A Set of Specifications
$\{\Phi_1, \ldots, \Phi_n\}$

P[satisfy all specs]

**60%**

**90%**

# Quantitative Analysis

## Empirical Evaluation via Simulation

**Carla Simulator: Extract execution traces.**

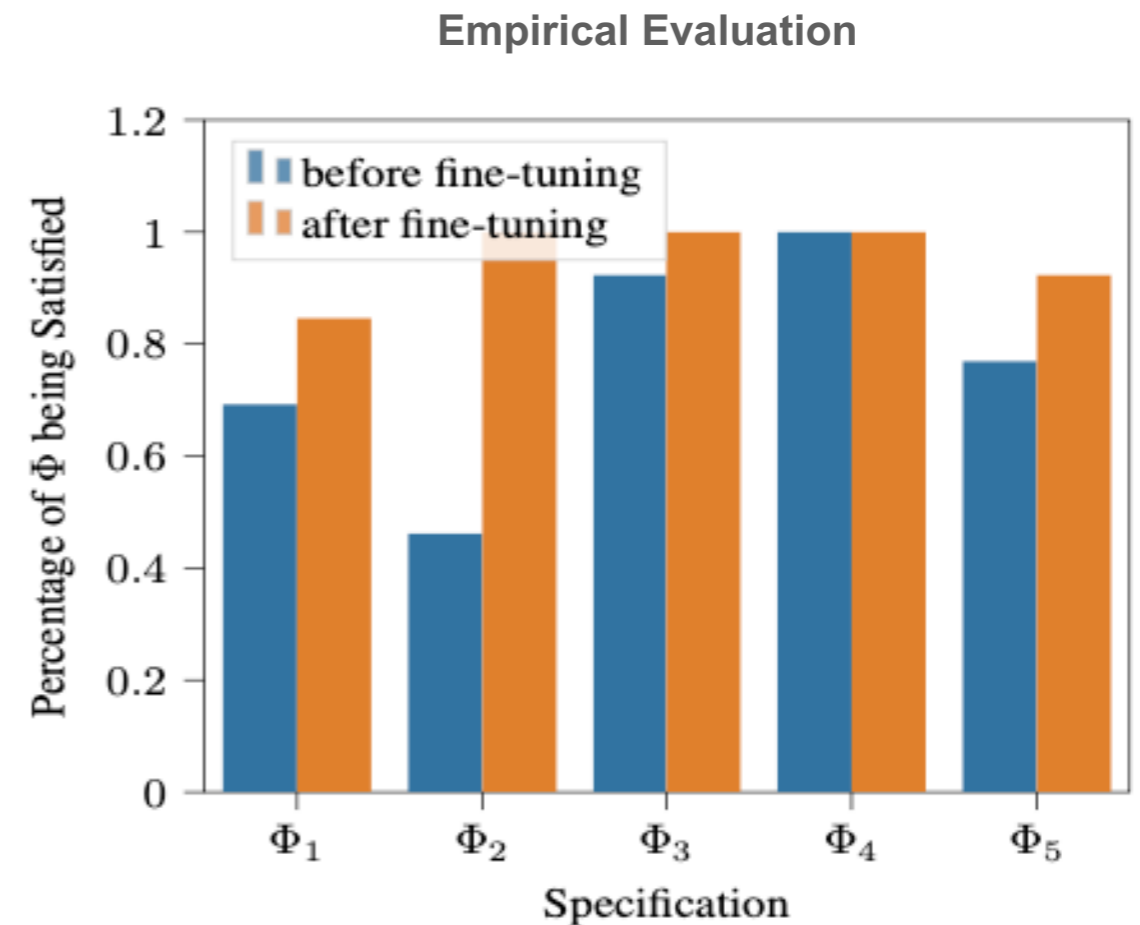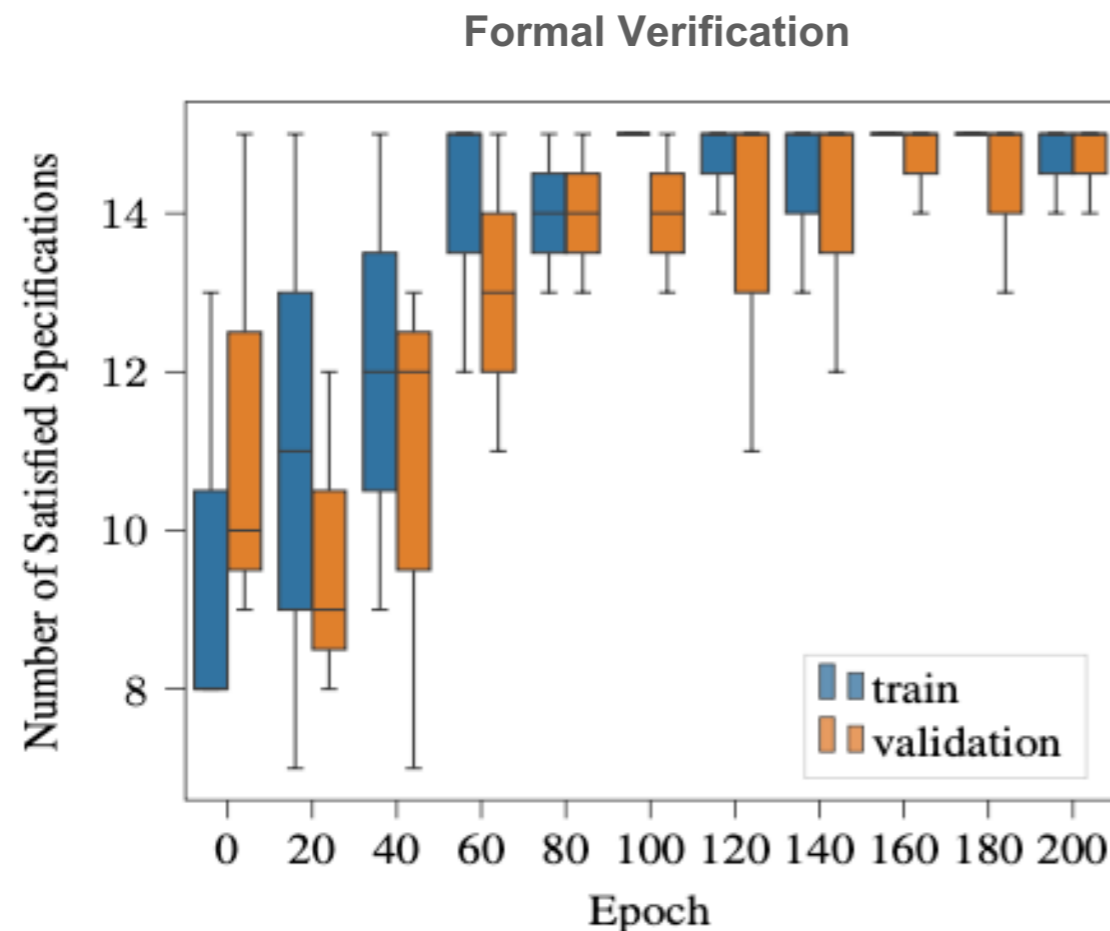Carla Simulation Video

Object and Position Information



**Execution Trace: (desired objects with positions, action),……**

**Desired objects**: pedestrian, car, red/green traffic light, stop sign,……
**Actions**: go straight, turn left, stop, turn right,……

# Quantitative Analysis

## Empirical Evaluation via Simulation



The results indicate that our approach can improve the language model's ability to satisfy critical requirements.

Our approach can act as a starting point to guide the design process for real-world implementations of autonomous driving systems.

# Next Step: Verifiable Grounding

**Ground the controller to real autonomous driving robots**

```python
class AutonomousVehicle():
    def __init__(self):
        self.pedestrian = False
        self.car = False

    @abstractmethod
    def move_forward(self):
        # vehicle starts moving forward
        pass

    @abstractmethod
    def turn_left(self):
        # vehicle turns left
        pass

    @abstractmethod
    def turn_right(self):
        # vehicle turns right
        pass

    @abstractmethod
    def stop(self):
        # vehicle slows down and stops
        pass
```
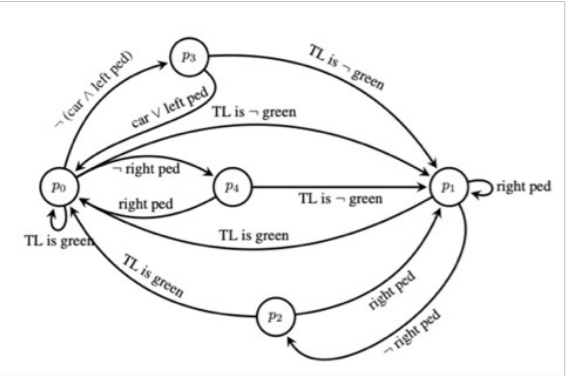
( left car $\lor$ ped at right, $\epsilon$)

(True, check left car

```python
def CrossRoad(vehicle):
    vehicle.observe()  # observe the environment first
    if vehicle.pedestrian_observed():  # Step 1
        vehicle.stop()  # stop if pedestrian is observed
        return "Stop"
    elif vehicle.car_observed():  # Step 2
        vehicle.stop()  # stop if car is observed
        return "Stop"
    else:  # Step 3
        vehicle.move_forward()
        return "Move_forward"
```

$\Box(\neg\text{green traffic light} \rightarrow \neg\text{go straight})$, ✅

$\Box(\text{stop sign} \rightarrow \Diamond\text{stop})$, ✅

$\Box\,\neg\text{turn right} \lor \neg(\text{car from left} \lor \text{pedestrian at right})$, ✅

# Q&A