

Fine-Tuning Language Models Using Formal Methods Feedback

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

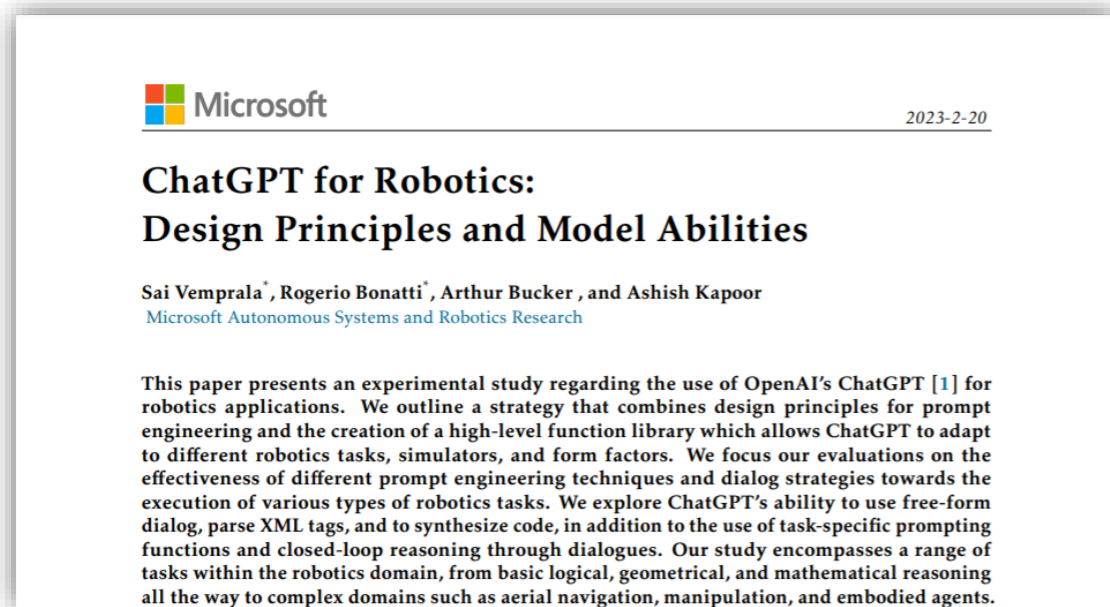
* Equal Contribution

Paper available at: <https://arxiv.org/abs/2310.18239>

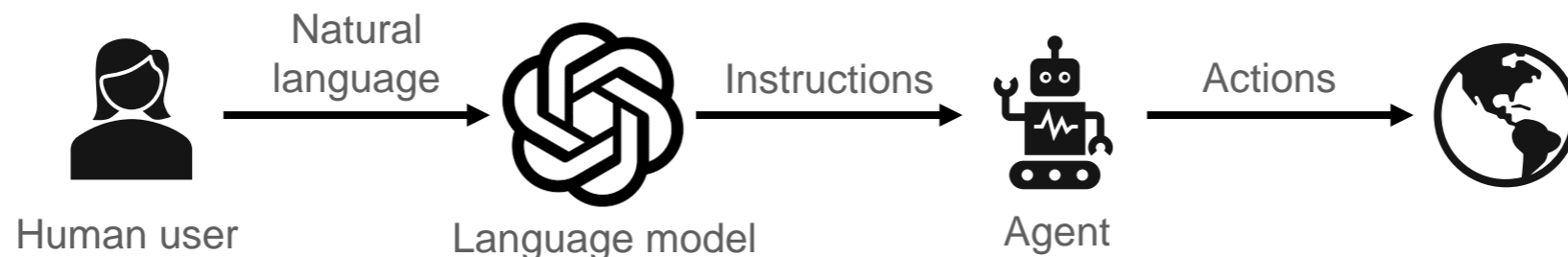


* Equal Contribution

Pre-Trained Language Models for Autonomous Systems



Pre-trained language models encode **rich world knowledge** and provide a **new interface** between humans and machines.



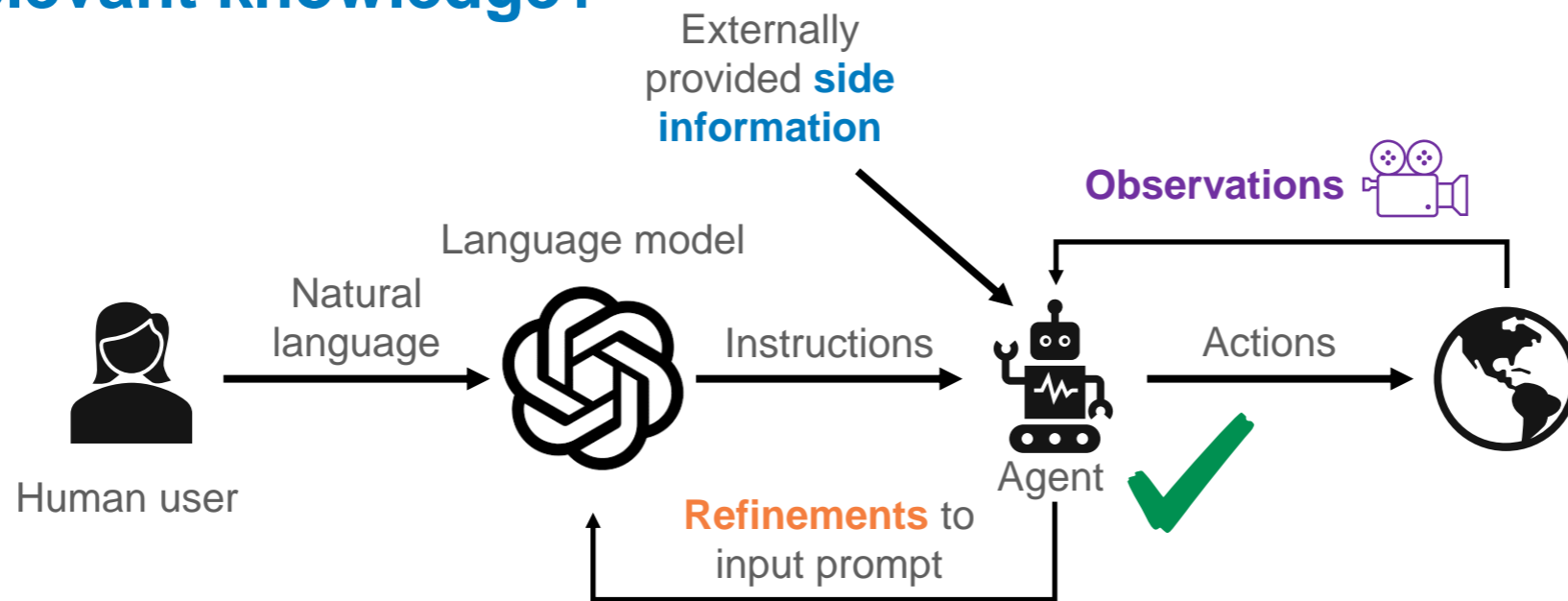
How can we incorporate language models to generate reliable high-level plans or control policies for autonomous systems?

Generative pre-trained models for autonomous systems

How can we ...

Leverage additional sources of task-relevant knowledge?

Close the decision-making loop?



Automatically refine the generated behaviors?

Verify that the generated behavior will satisfy critical requirements?

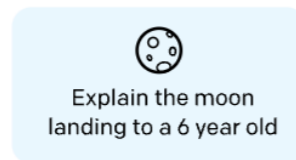
Fine-Tuning Language Models Using Human Feedback

Example: OpenAI Scheme for Instruct GPT

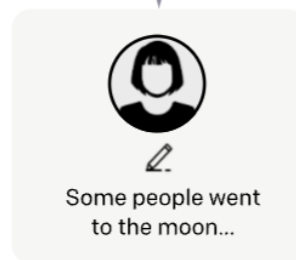
Step 1

Collect demonstration data, and train a supervised policy.

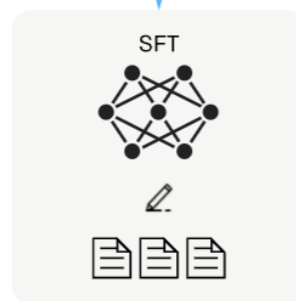
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



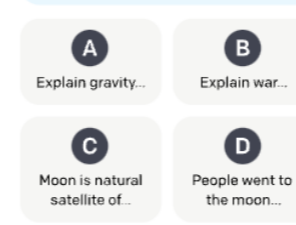
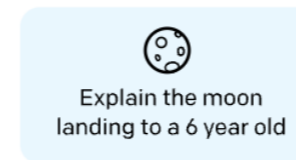
This data is used to fine-tune GPT-3 with supervised learning.



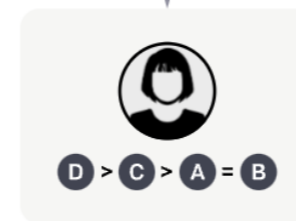
Step 2

Collect comparison data, and train a reward model.

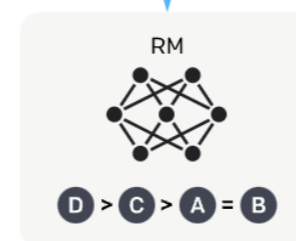
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using reinforcement learning.

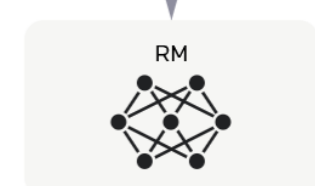
A new prompt is sampled from the dataset.



The policy generates an output.



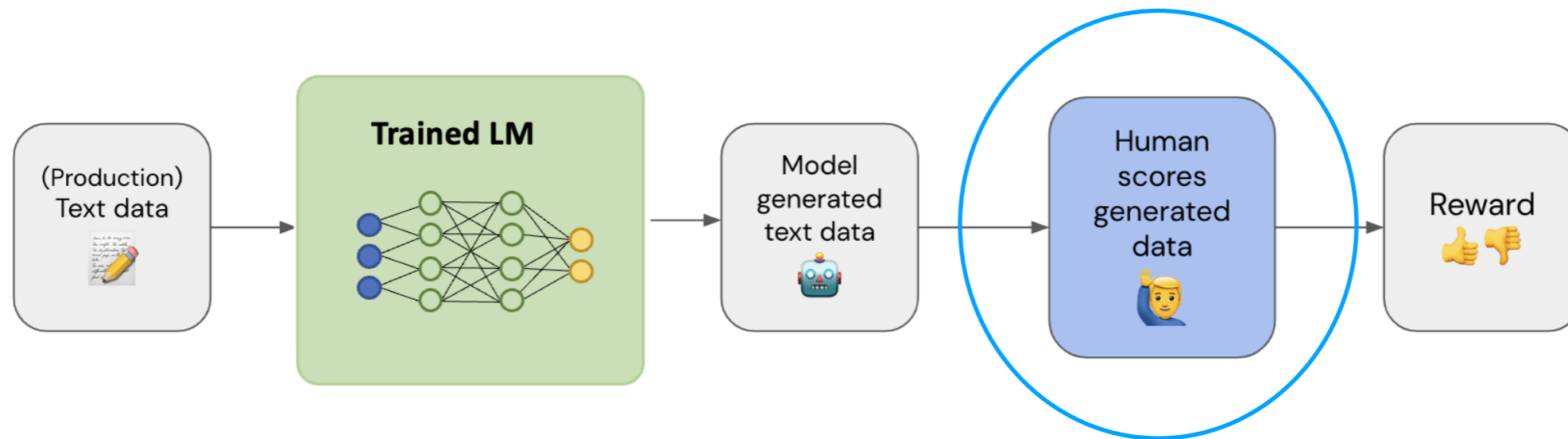
The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



Fine-Tuning Language Models Using Human Feedback

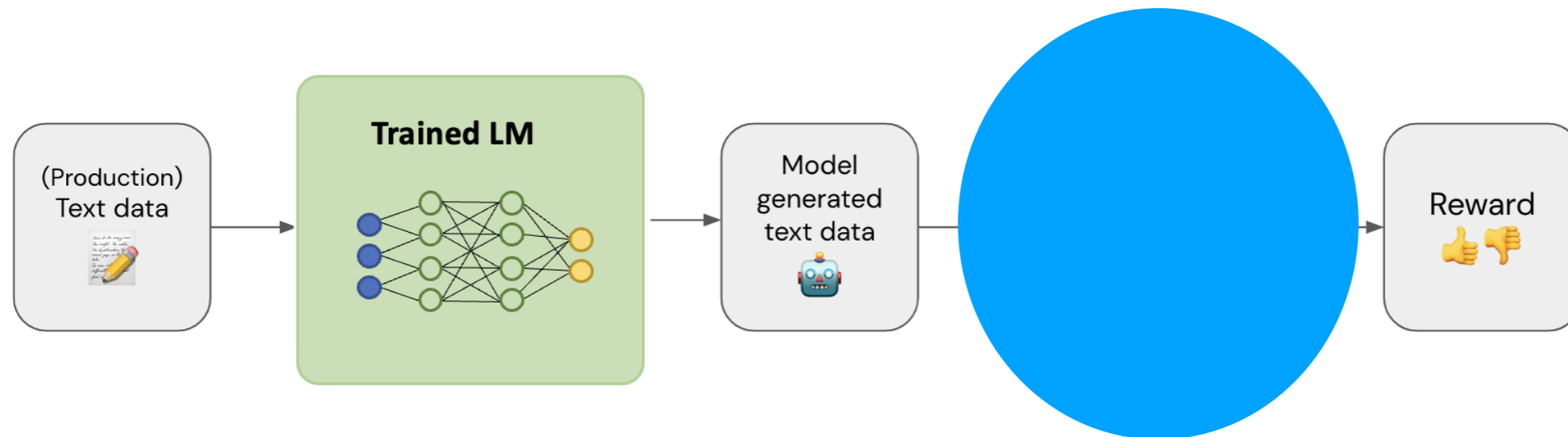


Labor-Intensive



Subjective/Inconsistent Feedback

Fine-Tuning Language Models Using Formal Feedback Methods



Formal Methods:

Automaton-based Representation, Model Checking, Temporal Logic Specification, etc.



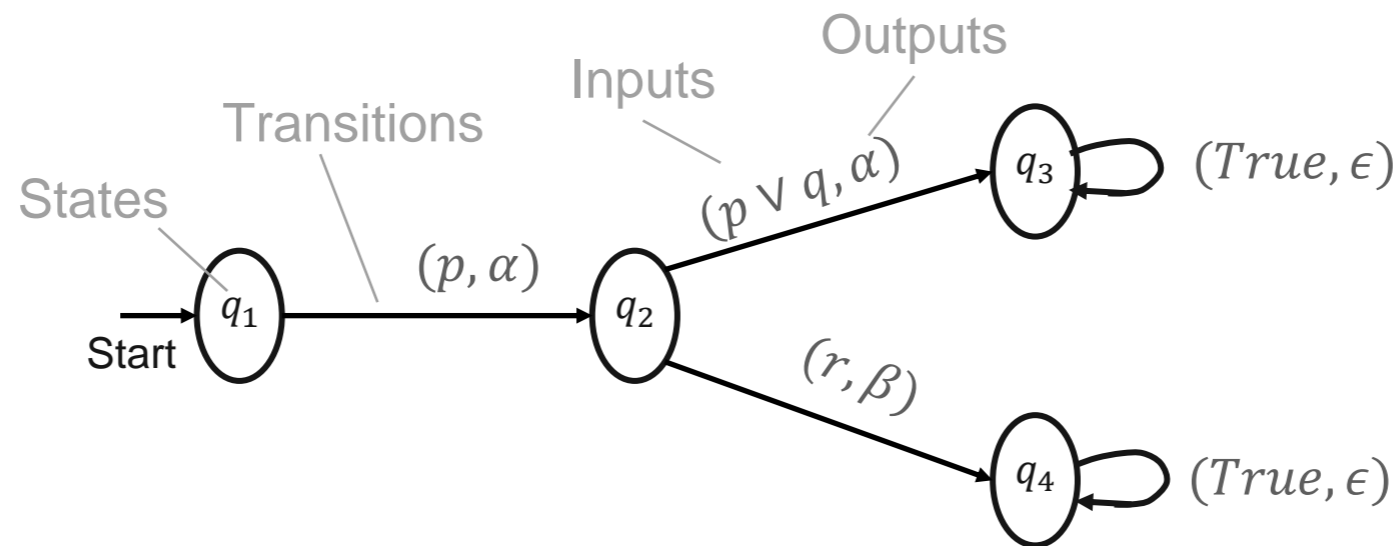
Labor-Intensive



Subjective/Inconsistent Feedback



Background: Automaton-based Representations

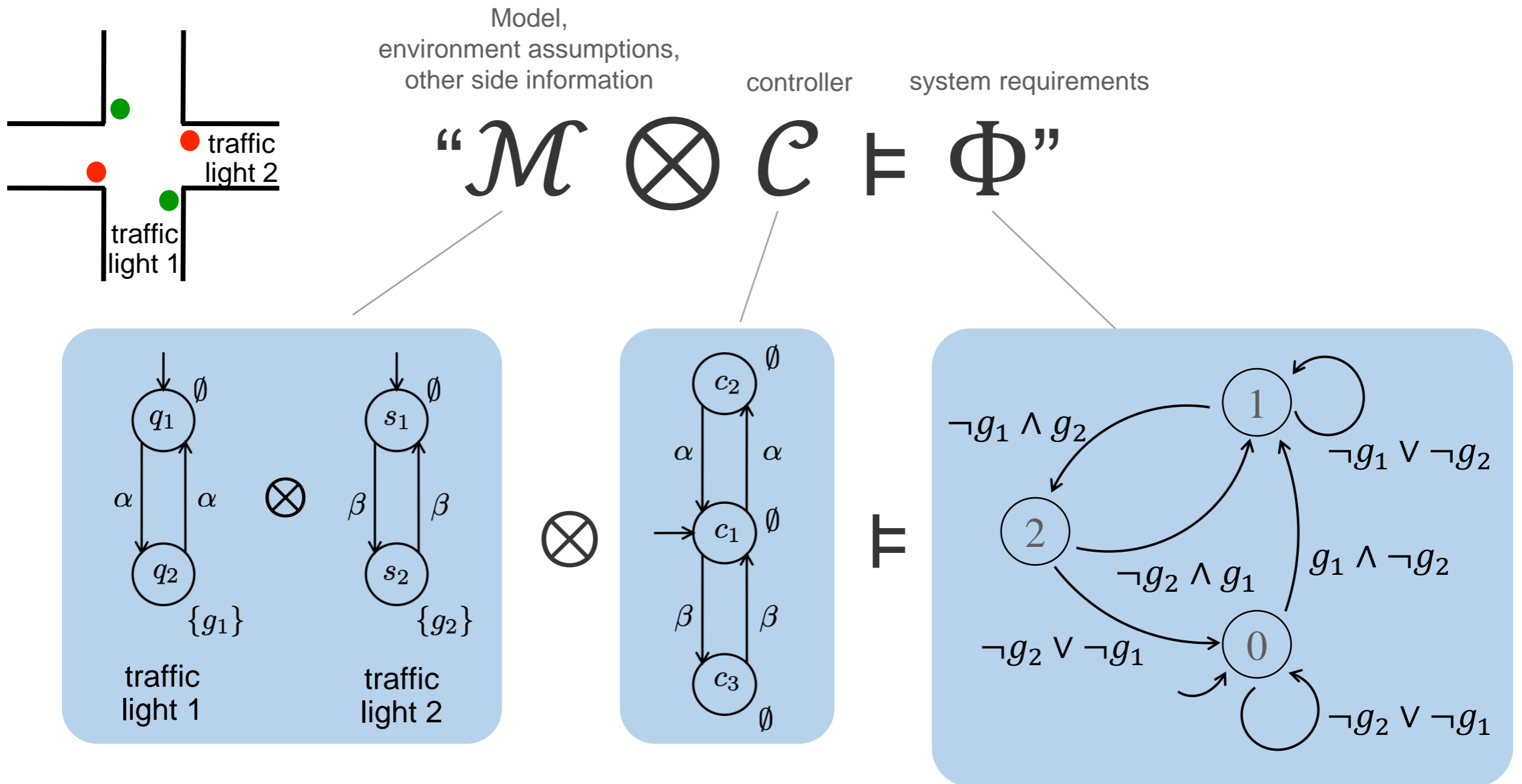


Why automaton-based representations? They are used for

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

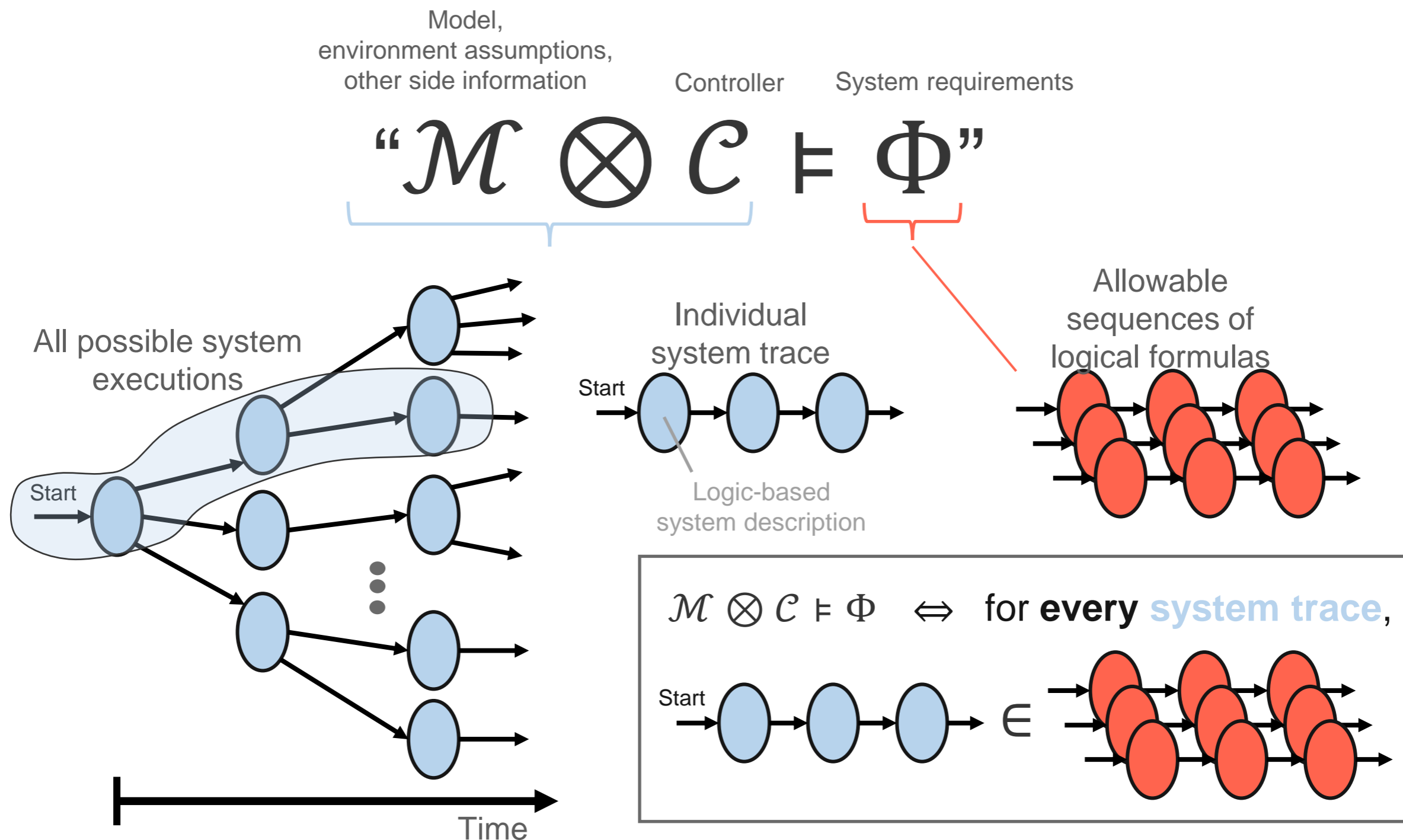
A (Very) Brief Introduction to Model Checking

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



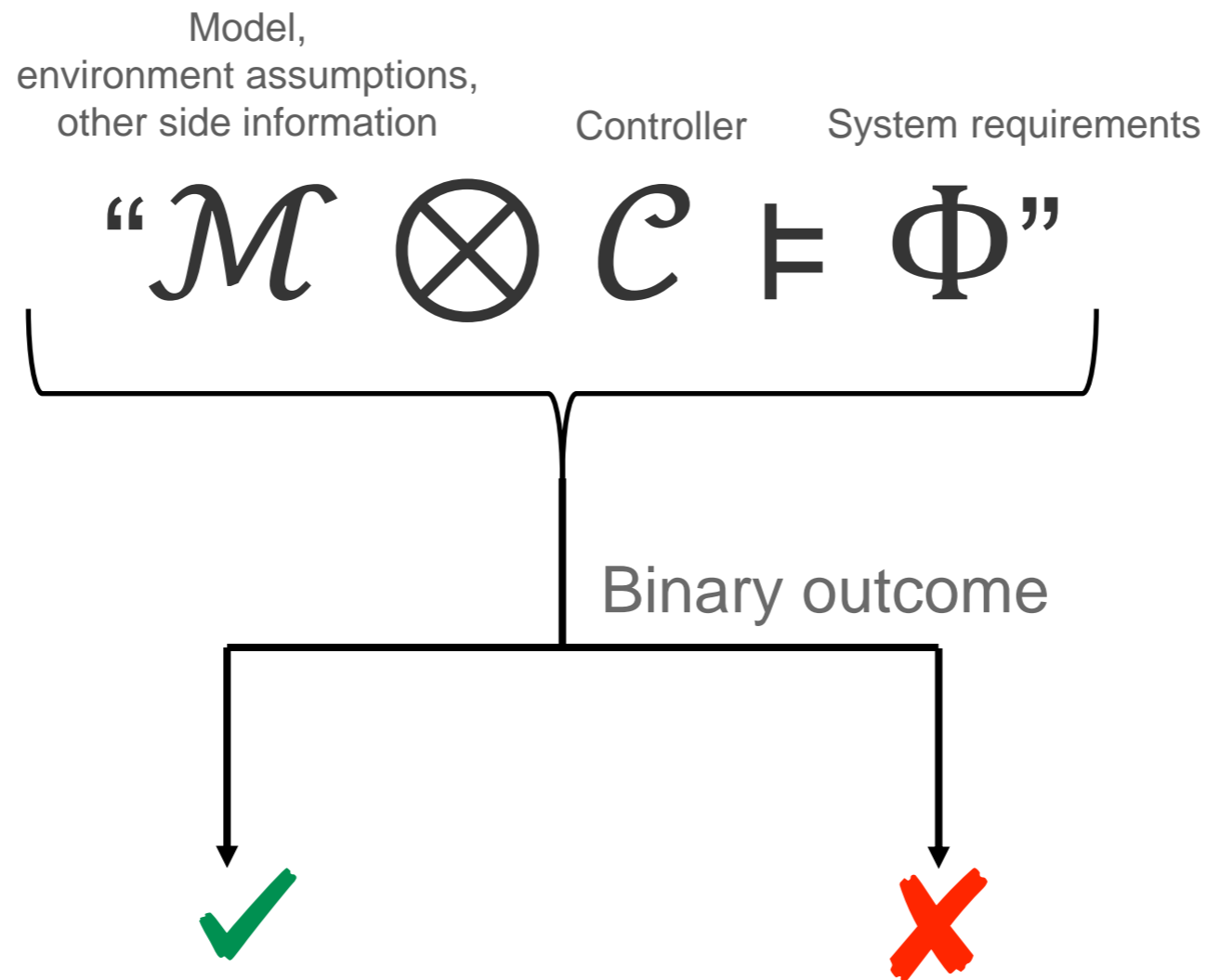
A (Very) Brief Introduction to Model Checking

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

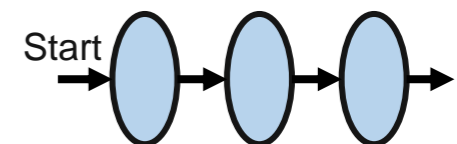


A (Very) Brief Introduction to Model Checking

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

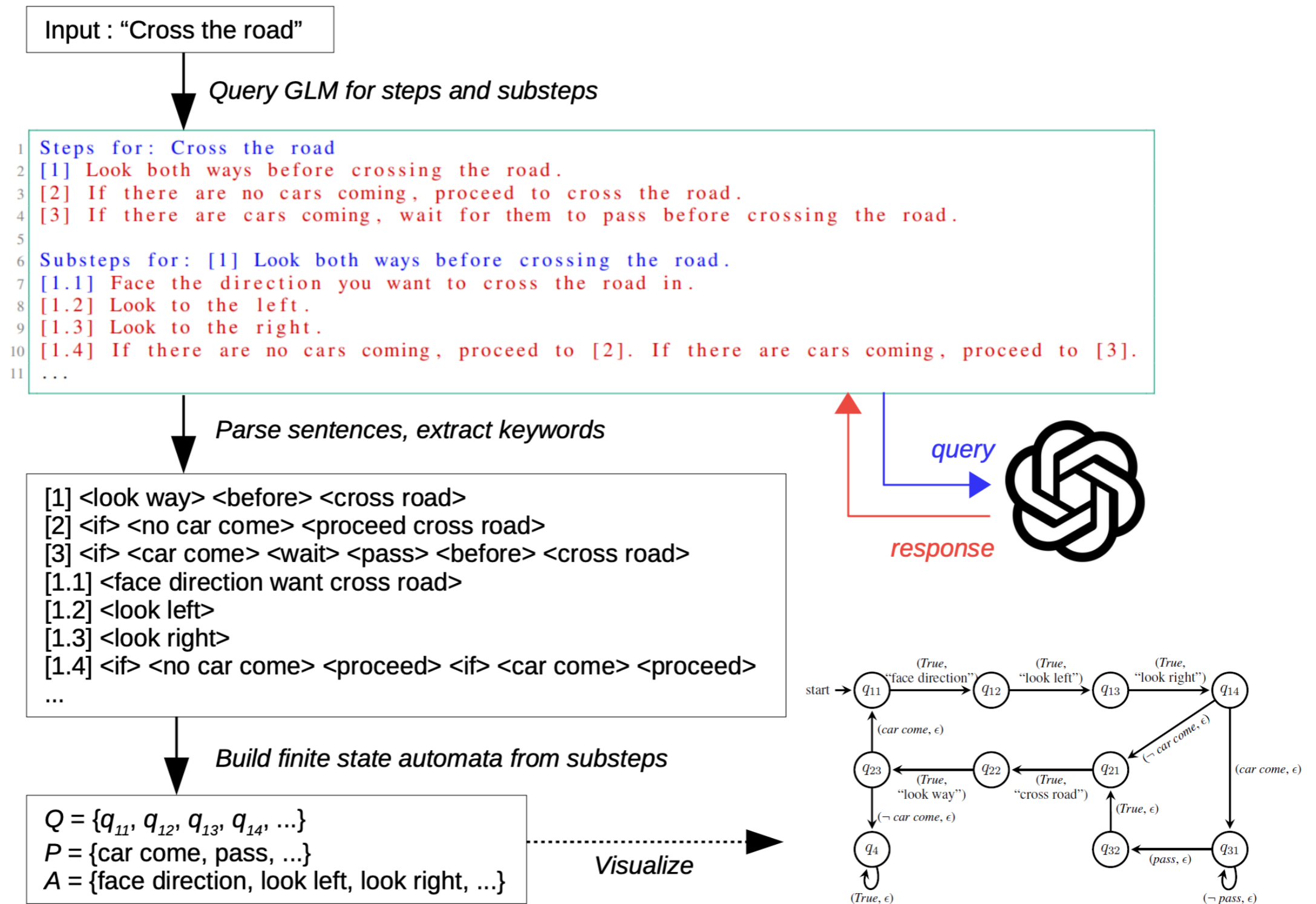


Byproduct: Counterexample trace that **violates** the specification.



How to Connect Generative Models to Automata?

(GLM2FSA: Generative Language Model to Finite-State Automaton)



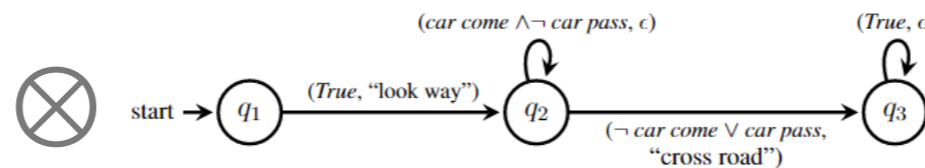
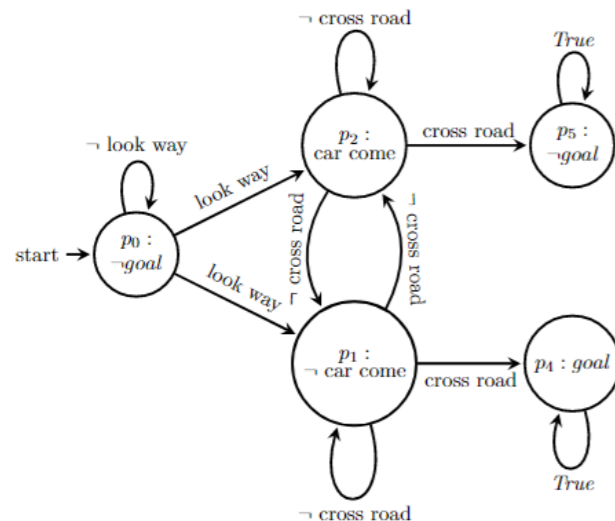
Verifying Whether the Generated Behaviors Satisfy User-Specified Requirements

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

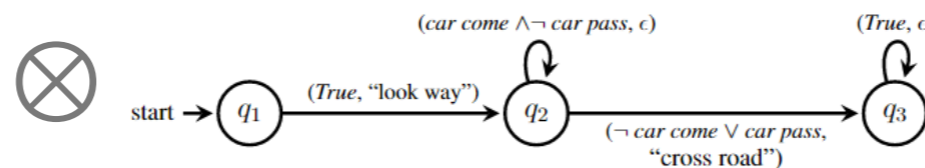
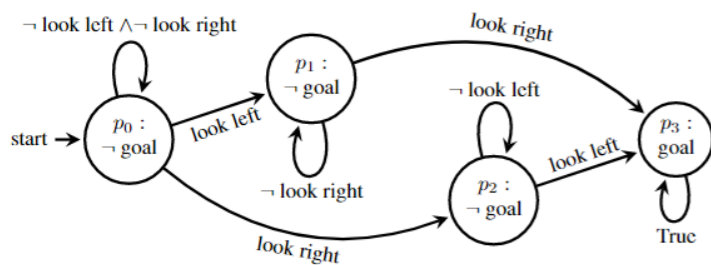
\mathcal{M} Additional available information, e.g., a model

\mathcal{C} Controller constructed using GLM2FSA

Φ Specification

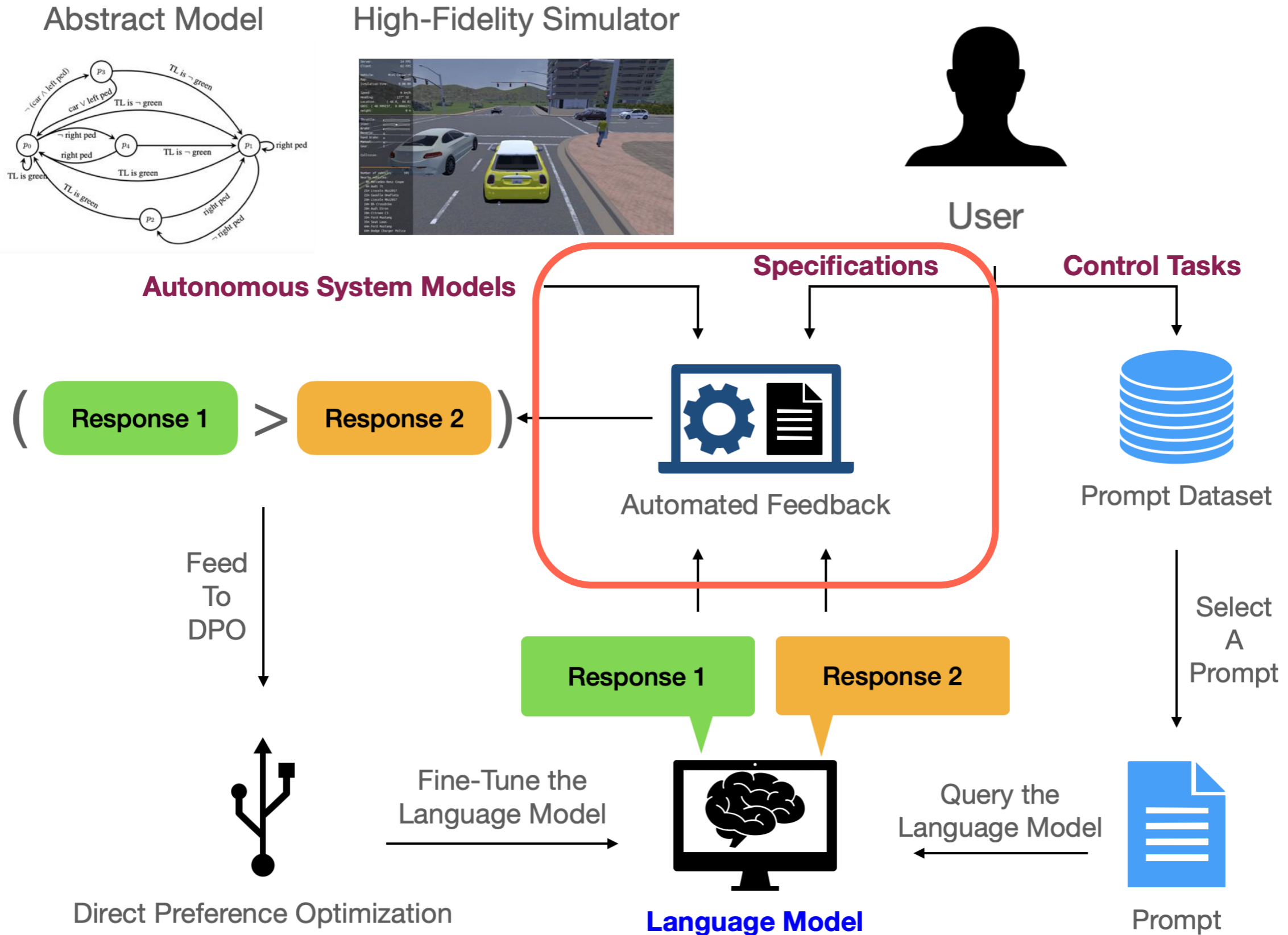


$\models \neg \text{traffic light} \rightarrow \diamond \text{goal}$
 "If not at a traffic light, eventually complete the task as specified by the model"



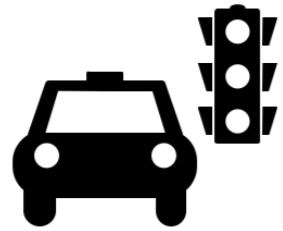
$\not\models \neg \text{traffic light} \rightarrow \diamond \text{goal}$

Fine-Tuning Language Models Using Formal Methods Feedback

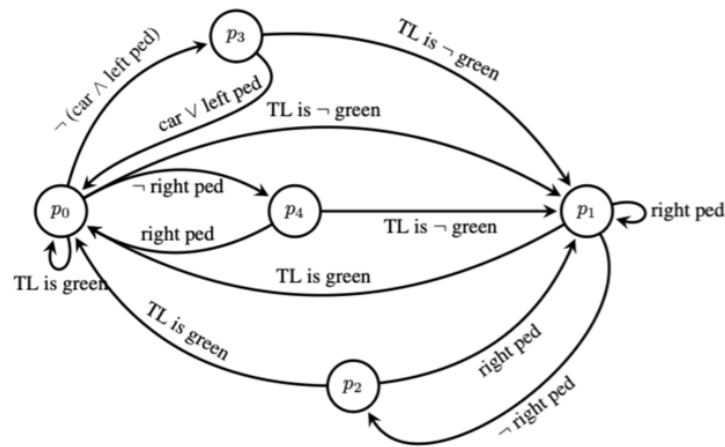


Fine-Tuning Language Models Using Formal Methods Feedback

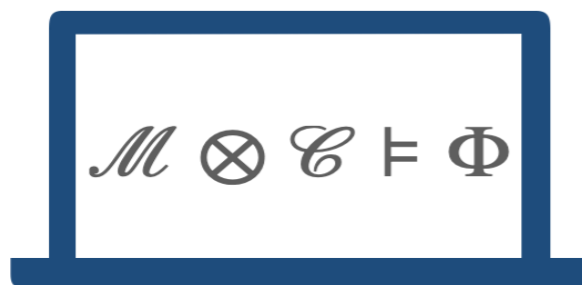
Modeling the Autonomous System



Autonomous System

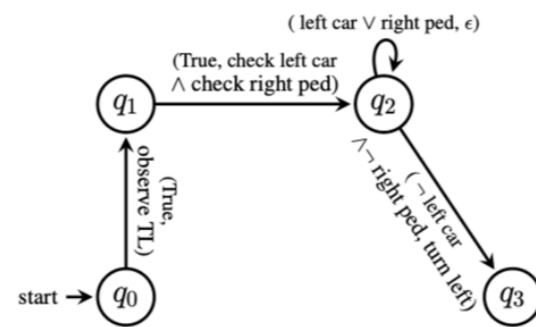


Autonomous System Model \mathcal{M}

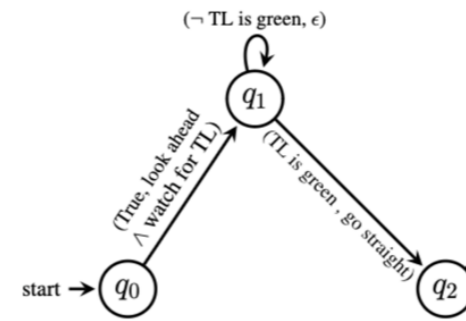


Model Checker

Controller Construction



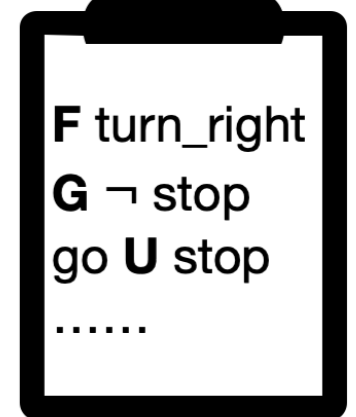
Controller \mathcal{C}_1



Controller \mathcal{C}_2



User



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

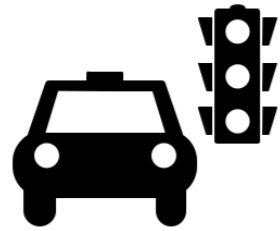
Compare the # of Satisfied Specifications



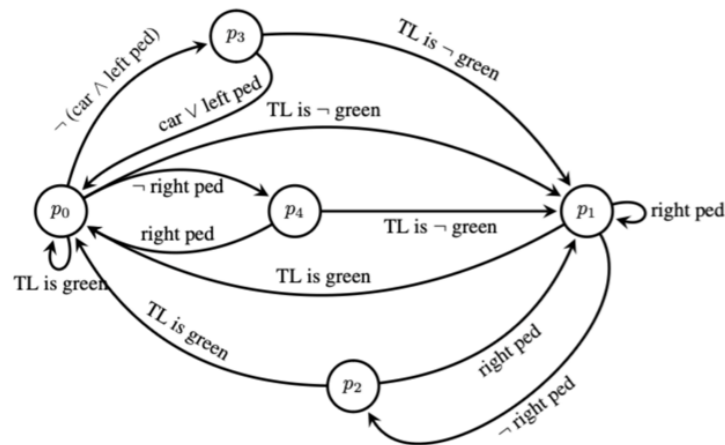
Formal Verification

Fine-Tuning Language Models Using Formal Methods Feedback

Modeling the Autonomous System



Autonomous System



Autonomous System Model \mathcal{M}

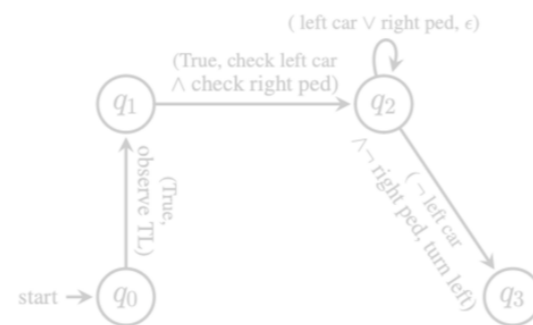
Controller Construction

Response 1

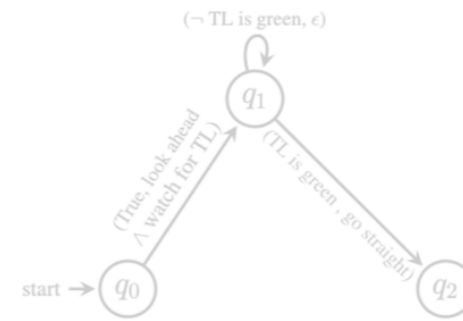
Response 2



User



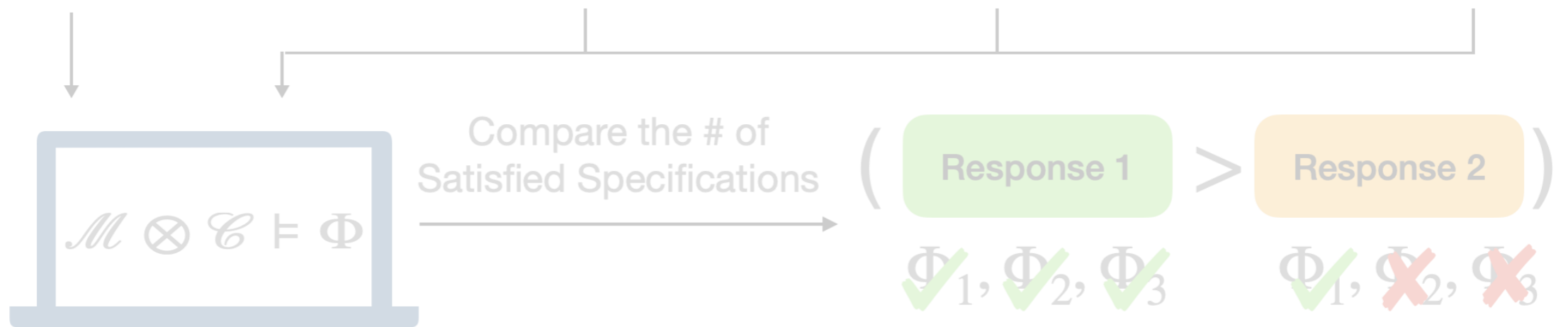
Controller \mathcal{C}_1



Controller \mathcal{C}_2



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

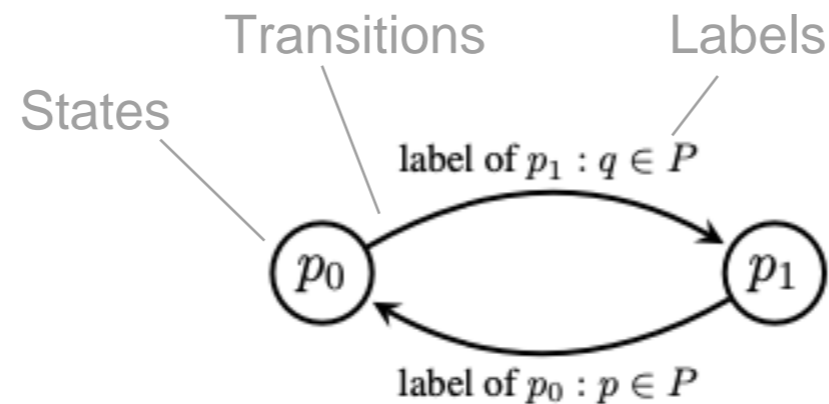
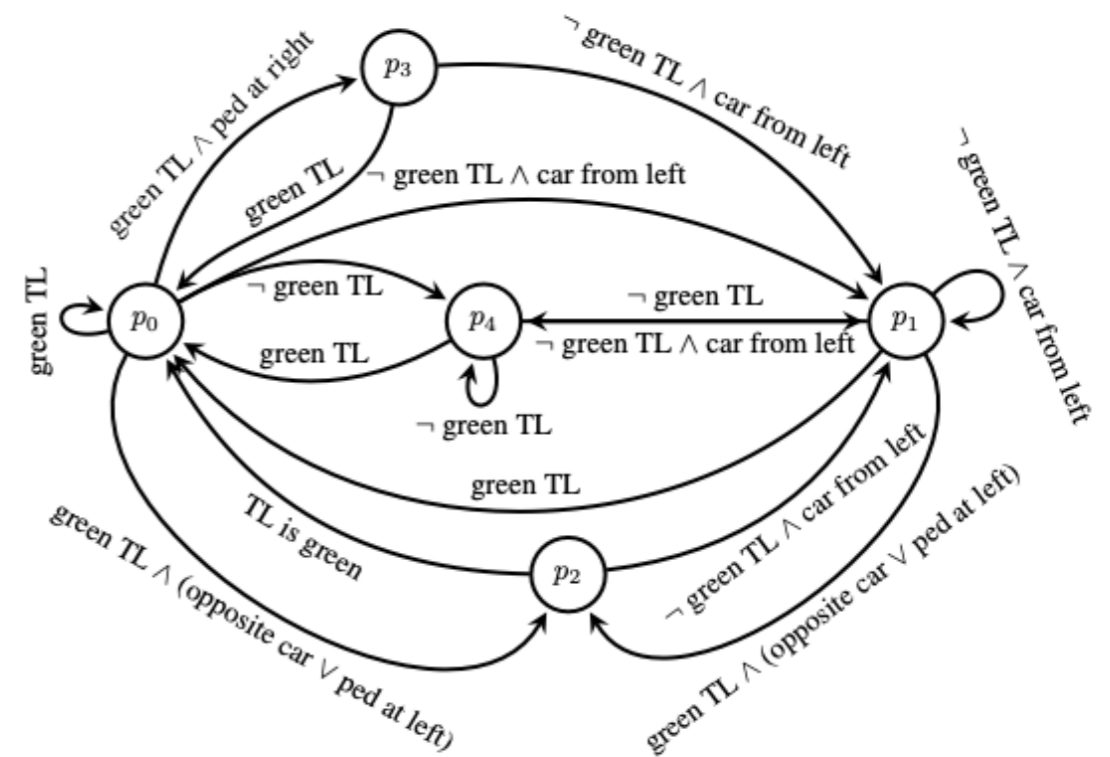
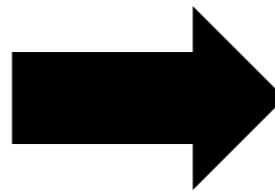
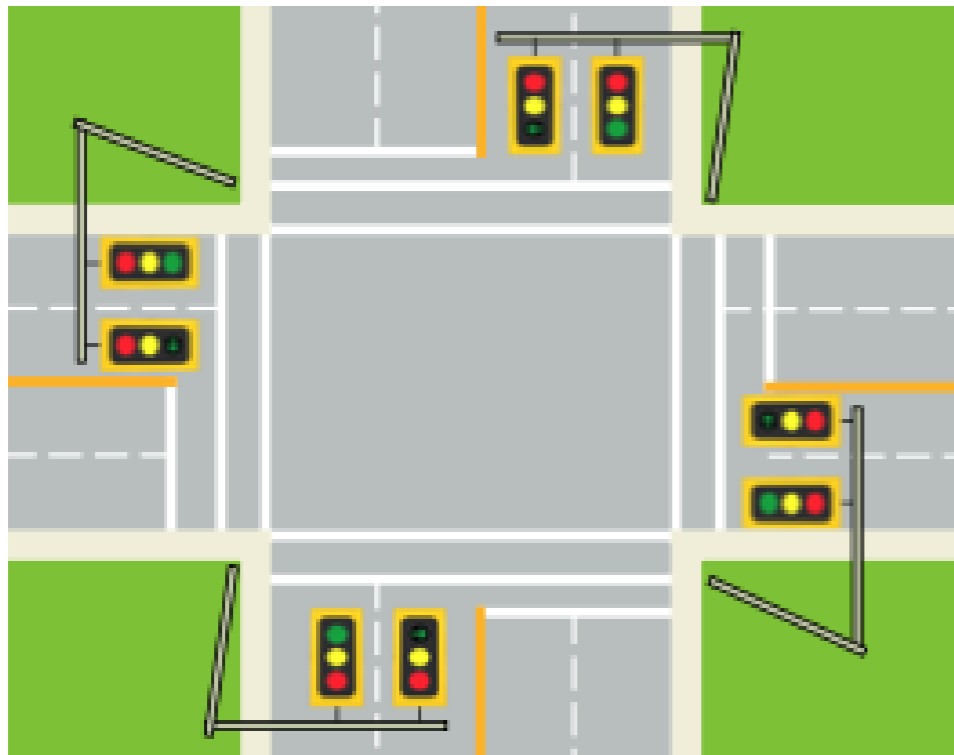


Model Checker

Formal Verification

Formal Methods Feedback in an Autonomous Driving System

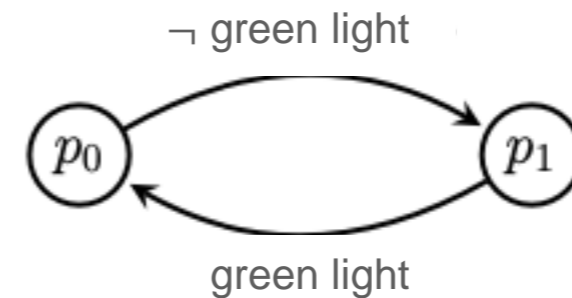
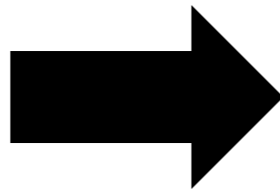
Modeling the Autonomous System



Formal Methods Feedback in an Autonomous Driving System

Modeling the Autonomous System

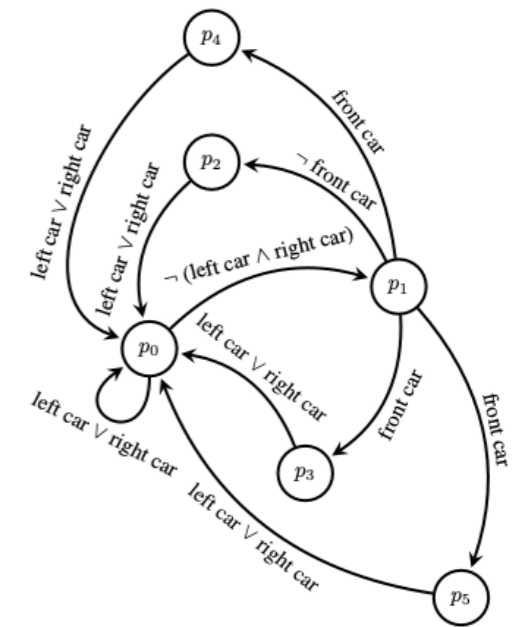
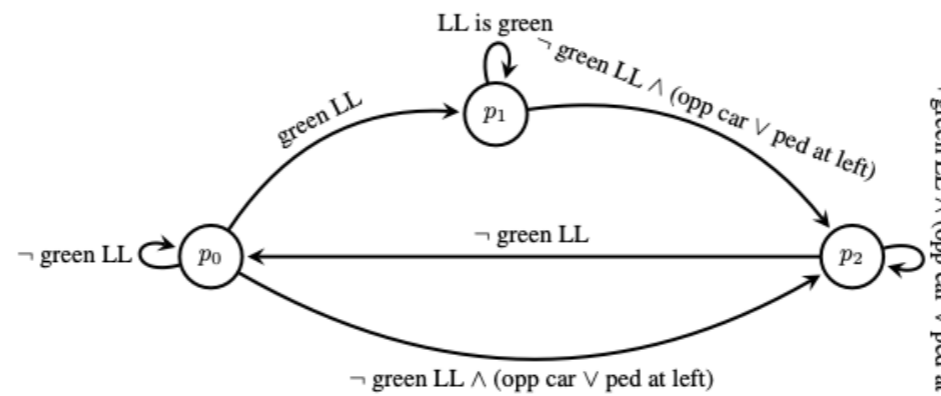
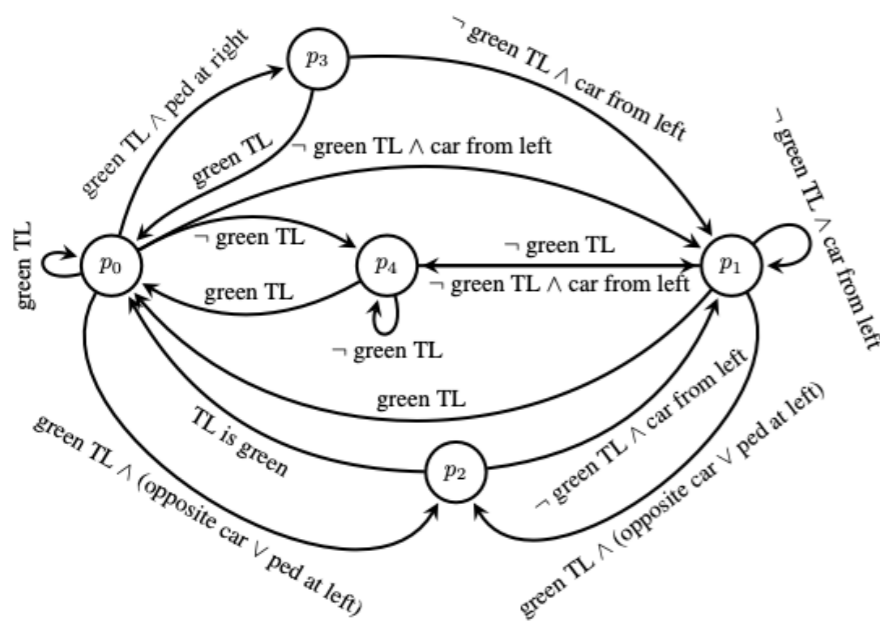
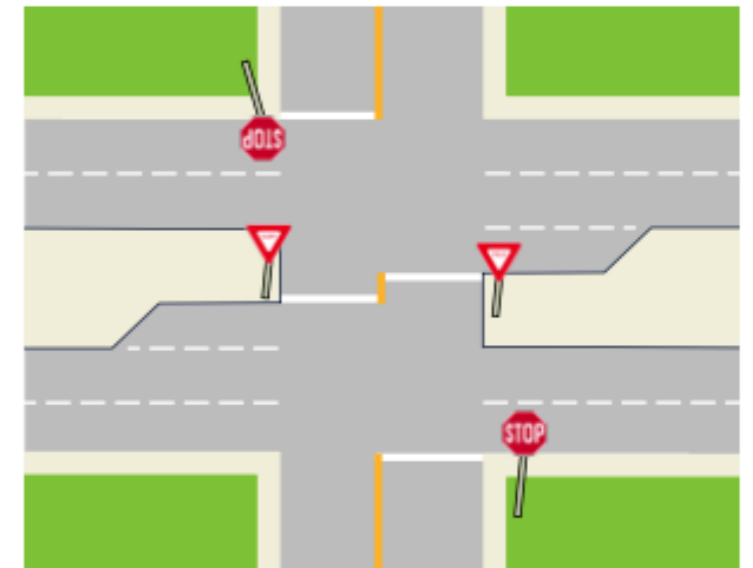
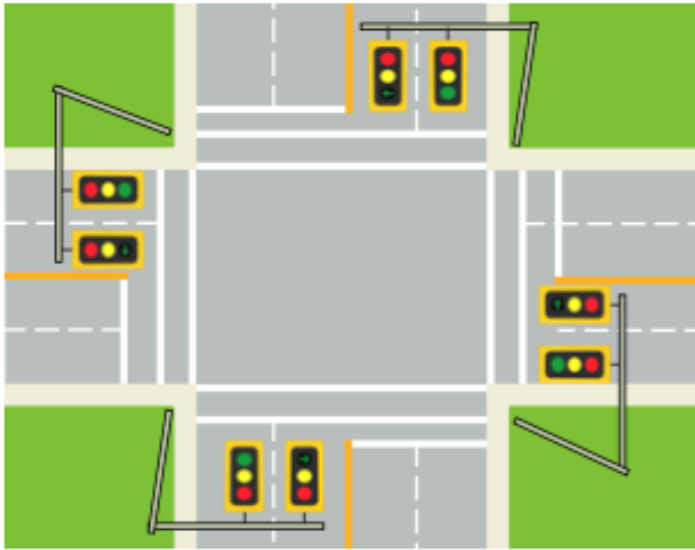
$P = \{\text{green light}\}$
 $S \sim \text{Traffic Light}$



$Q_M = \{p_0, p_1\}$
 $\lambda_M(p_0) = \text{green light}$
 $\lambda_M(p_1) = \neg \text{green light}$
 $\delta_M(p_0, p_1) = 1$
 $\delta_M(p_1, p_0) = 1$

Formal Methods Feedback in an Autonomous Driving System

Modeling the Autonomous System

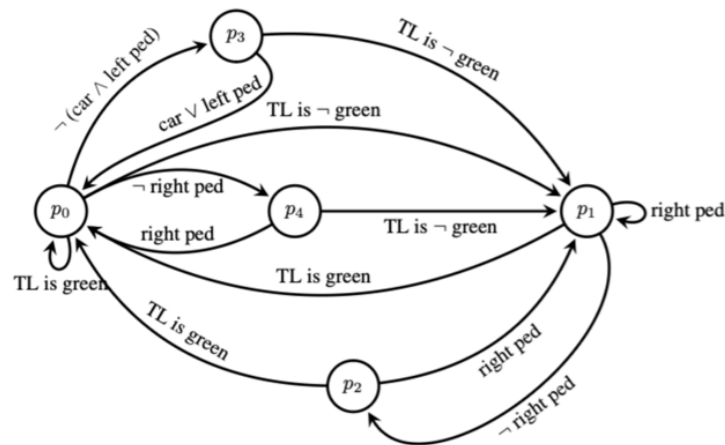


Fine-Tuning Language Models Using Formal Methods Feedback

Modeling the Autonomous System



Autonomous System

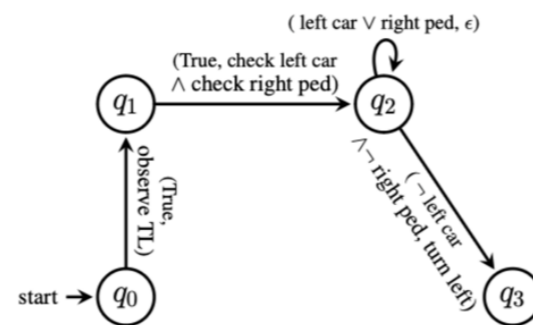


Autonomous System Model \mathcal{M}

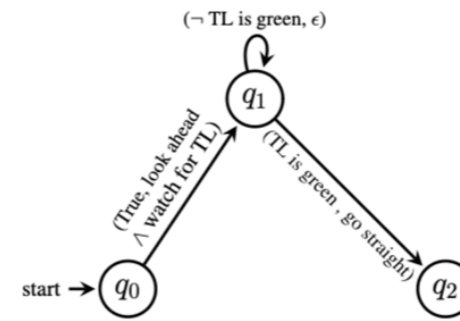
Controller Construction

Response 1

Response 2



Controller \mathcal{C}_1



Controller \mathcal{C}_2



User



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



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.

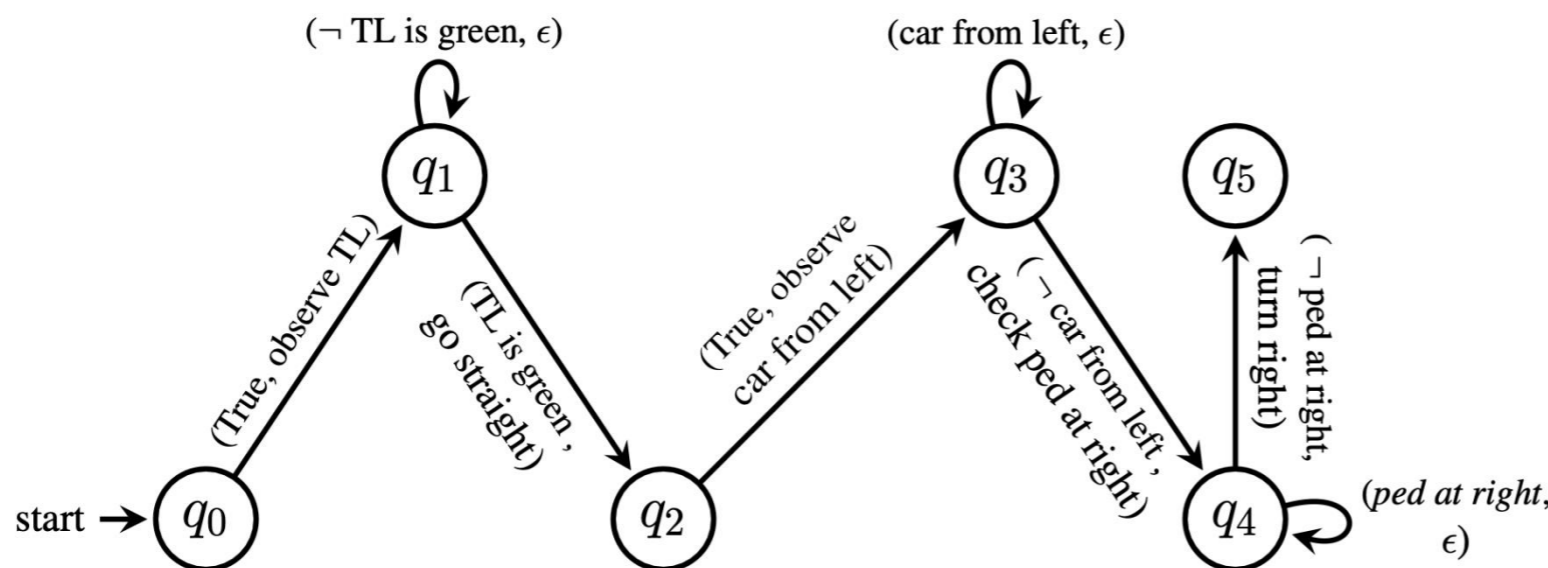
Formal Methods Feedback in an Autonomous Driving System

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.

Parse the sentence
Align the vocabulary to *P* and *A*

1. `<observe traffic light>`.
2. `<if> <green traffic light>, <go straight>`.
3. `<observe car from left>`.
4. `<if> <no car from left>, <check pedestrian at right>`.
5. `<if> <no pedestrian at right>, <turn right>`.

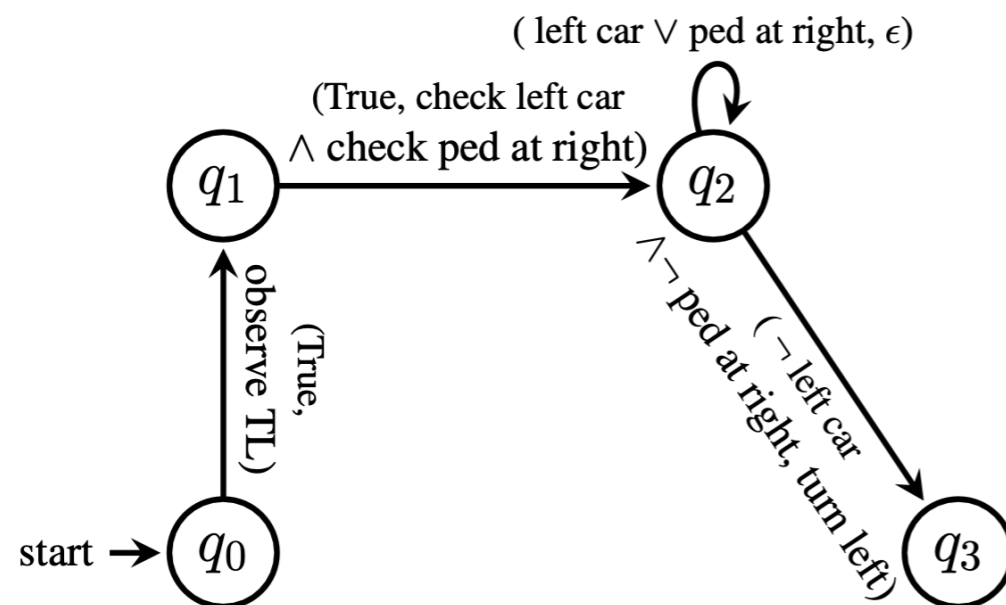
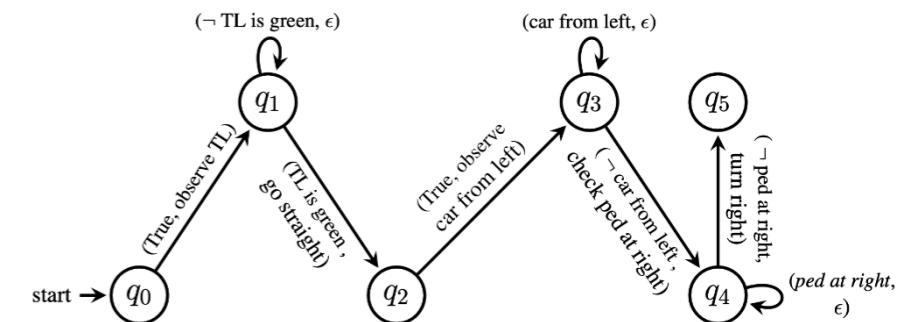


One state per step
Transition input: condition
Transition output: action

Formal Methods Feedback in an Autonomous Driving System

Controller Construction

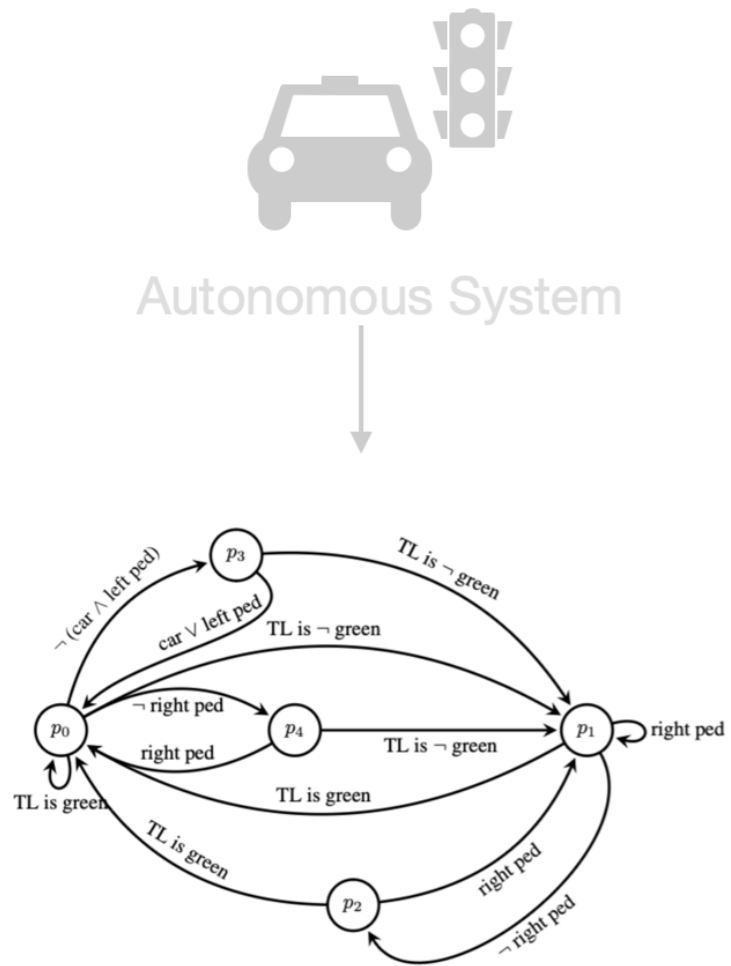
1. Look straight ahead and watch for traffic light.
2. If the traffic light turns green, start moving forward.
3. As you approach the intersection, look to your left for oncoming traffic.
4. If there is no traffic coming from your left, check pedestrians on your right.
5. If it is safe, turn your vehicle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

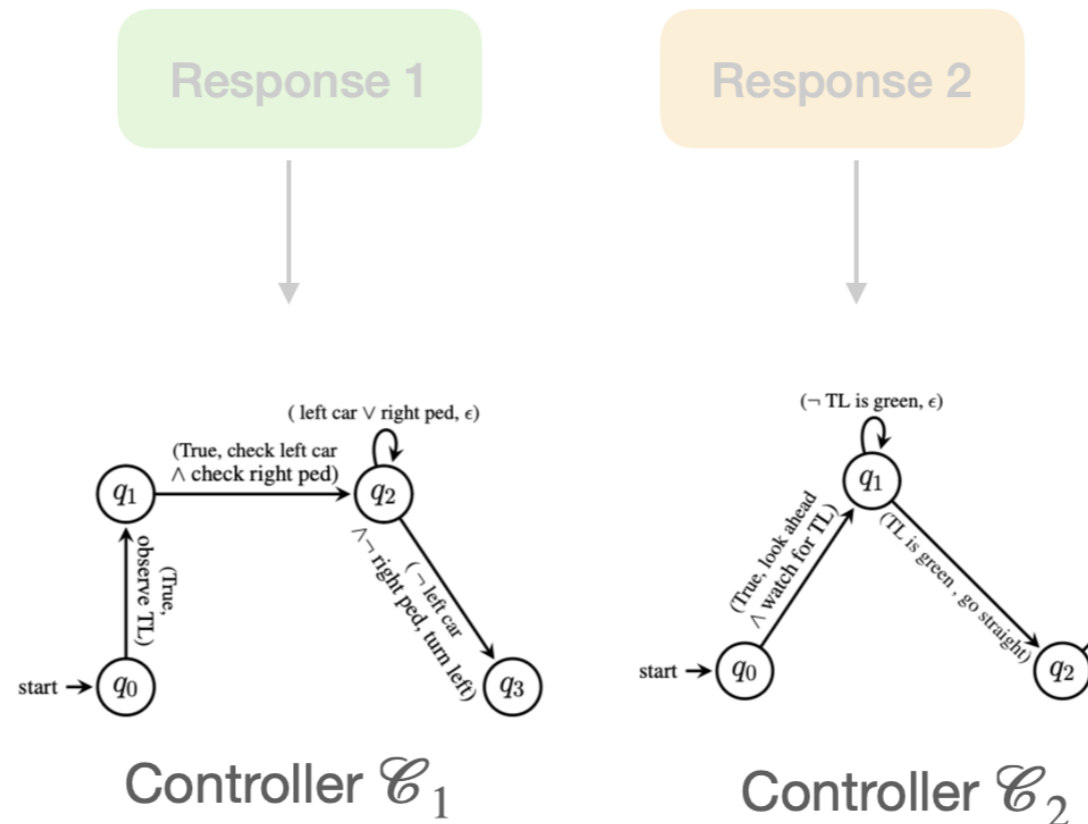
Fine-Tuning Language Models Using Formal Methods Feedback

Modeling the Autonomous System



Autonomous System Model \mathcal{M}

Controller Construction

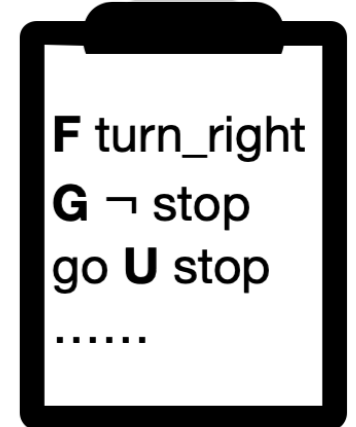


Controller \mathcal{C}_1

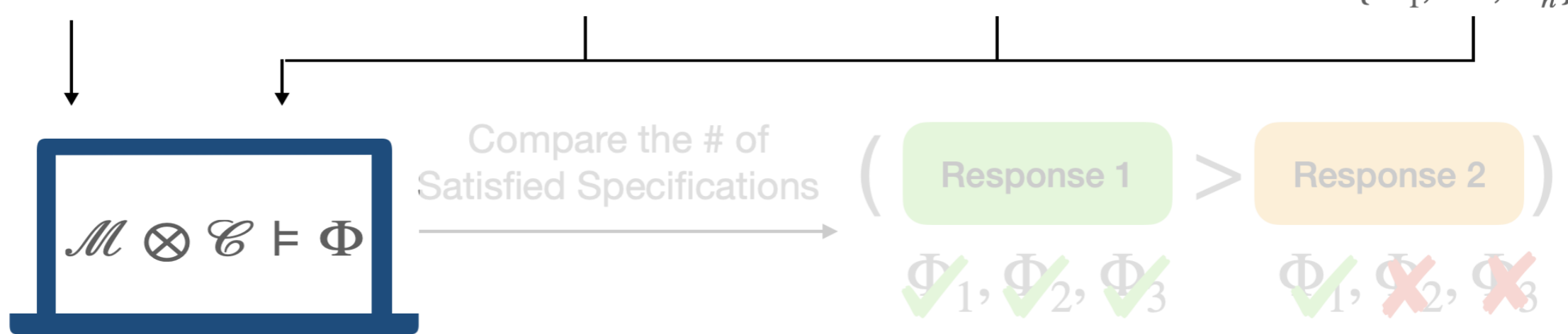
Controller \mathcal{C}_2



User



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

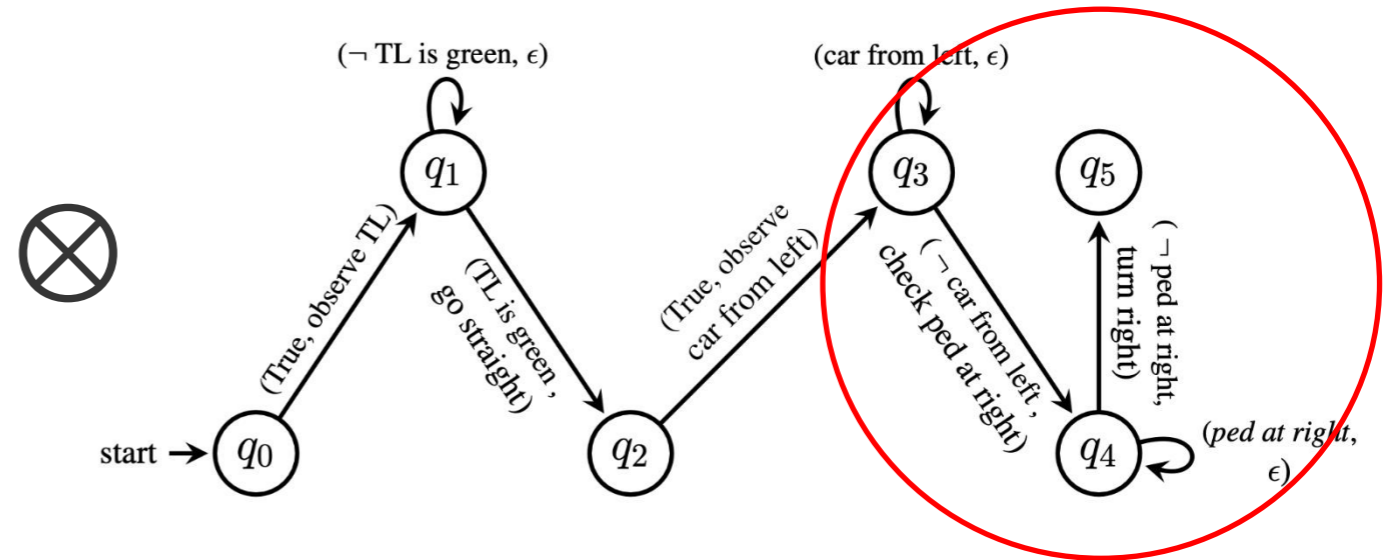
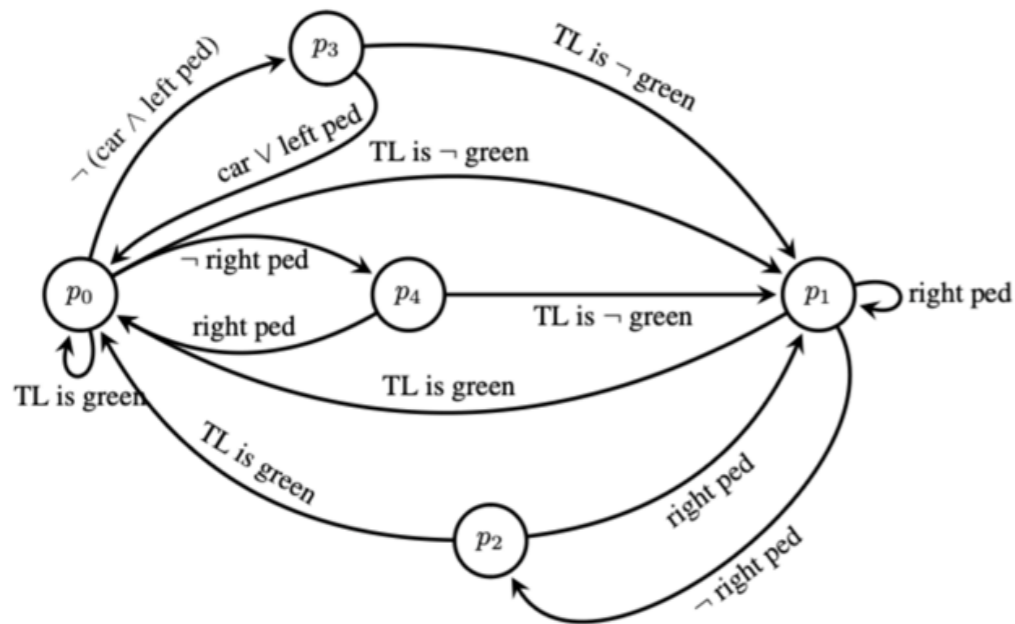


Model Checker

Formal Verification

Formal Methods Feedback in an Autonomous Driving System

Formal Verification



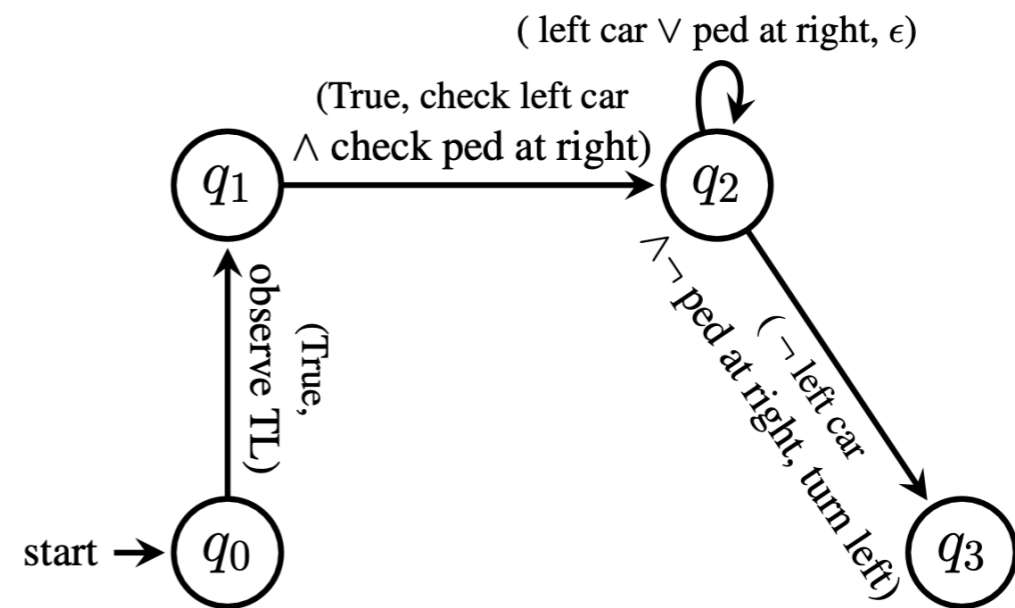
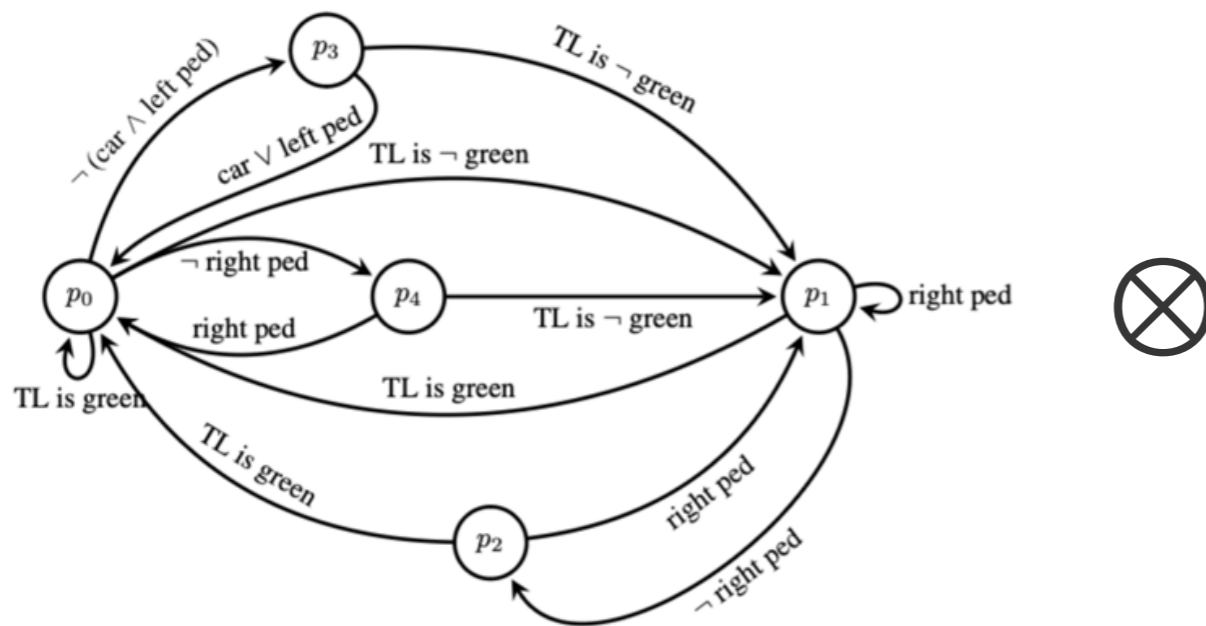
(\neg green traffic light \rightarrow \neg go straight), ✓

(stop sign \rightarrow \diamond stop), ✓

\neg turn right \vee \neg (car from left \vee pedestrian at right), ✗

Formal Methods Feedback in an Autonomous Driving System

Formal Verification



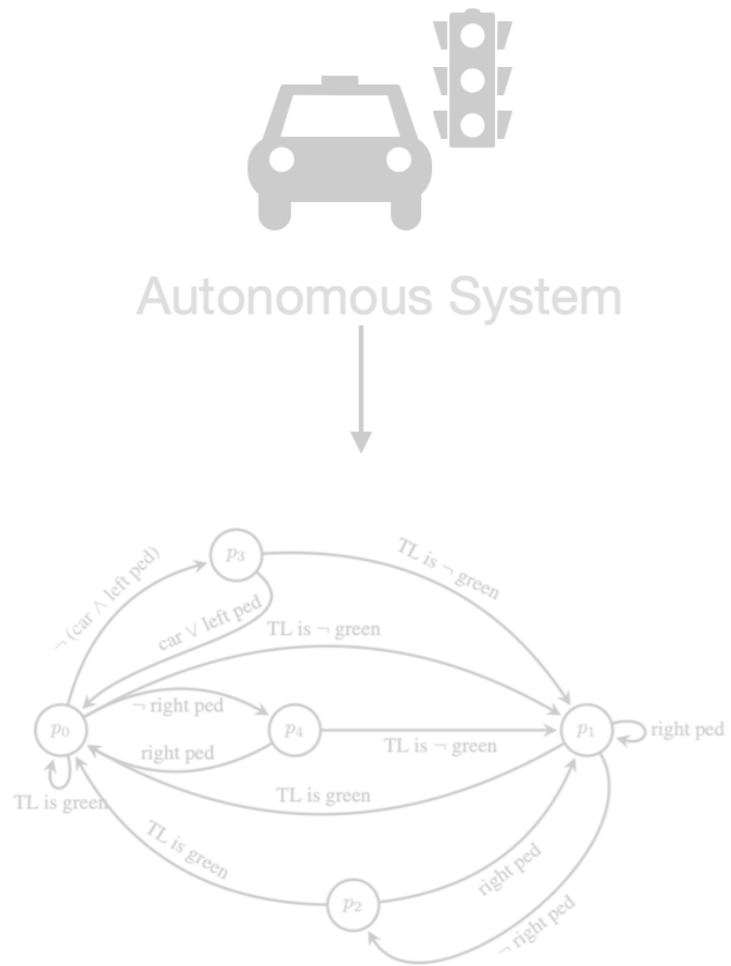
$(\neg \text{green traffic light} \rightarrow \neg \text{go straight})$, ✓

$(\text{stop sign} \rightarrow \diamond \text{stop})$, ✓

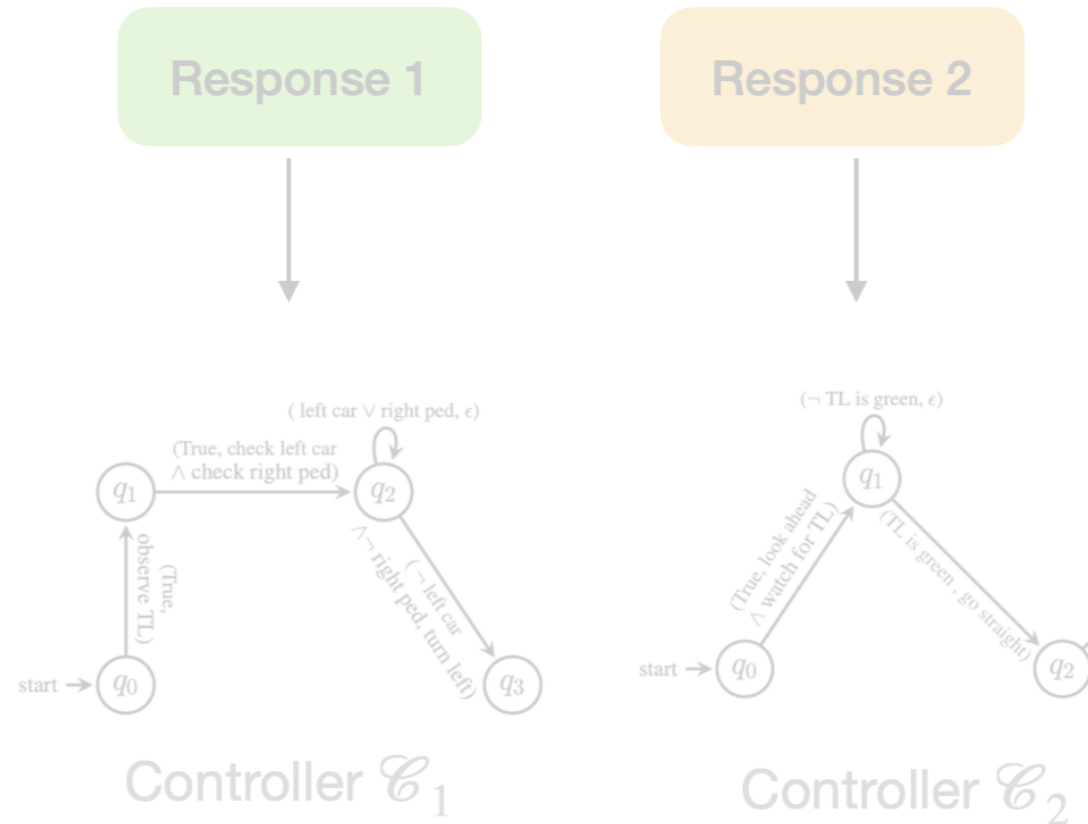
$\neg \text{turn right} \vee \neg(\text{car from left} \vee \text{pedestrian at right})$, ✓

Fine-Tuning Language Models Using Formal Methods Feedback

Modeling the Autonomous System



Controller Construction



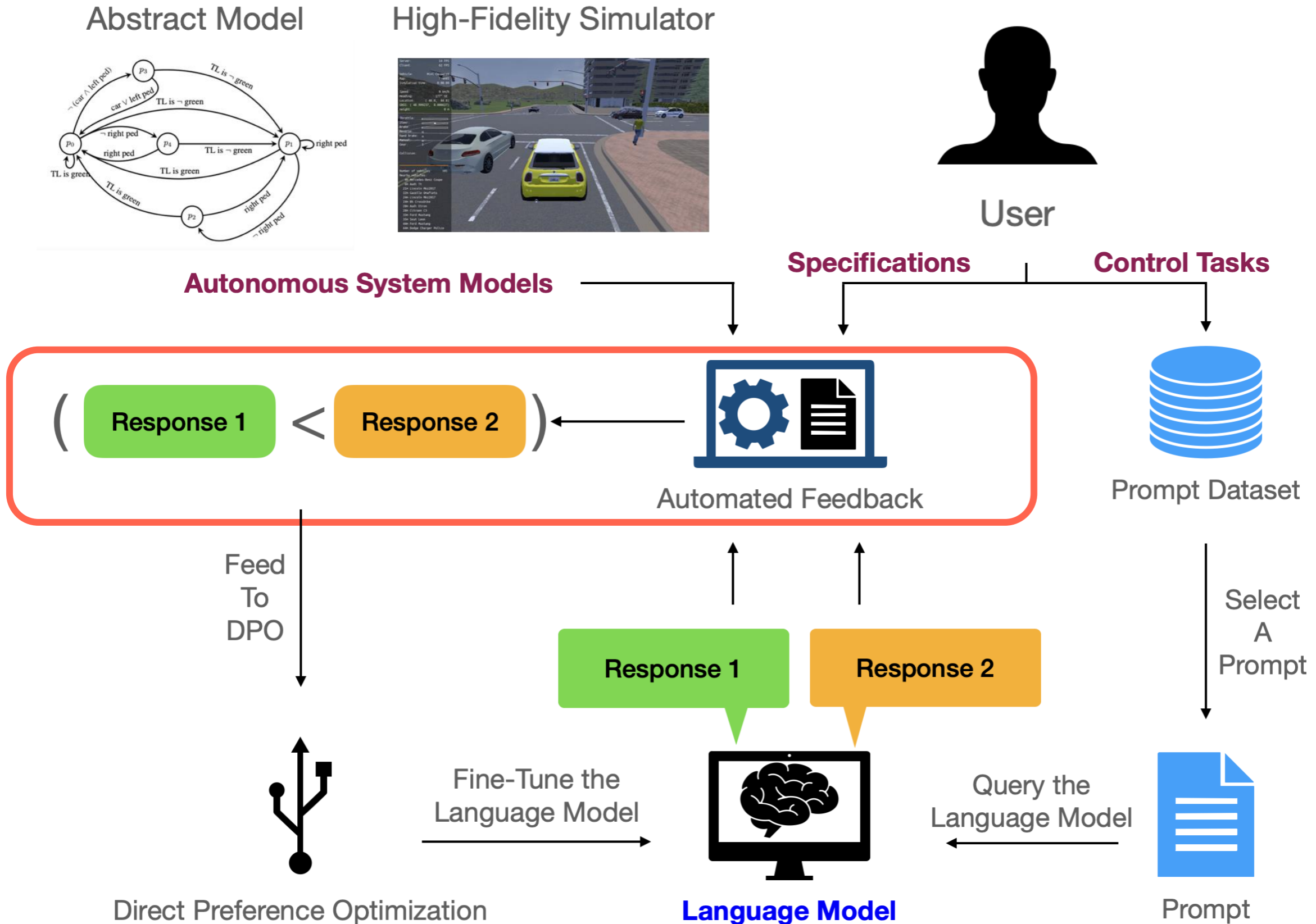
User



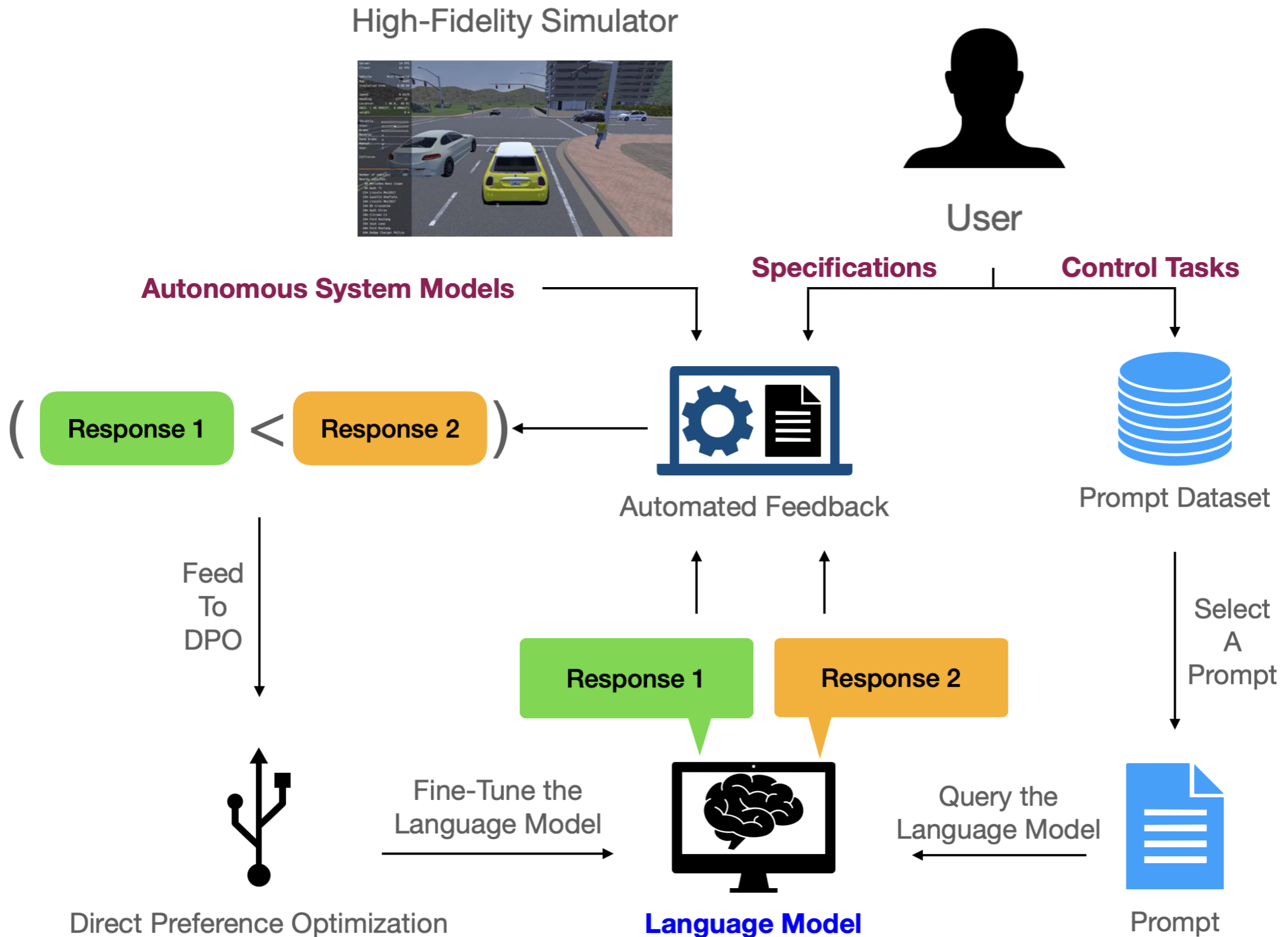
Autonomous System Model \mathcal{M}



Fine-Tuning Language Models Using Formal Methods Feedback

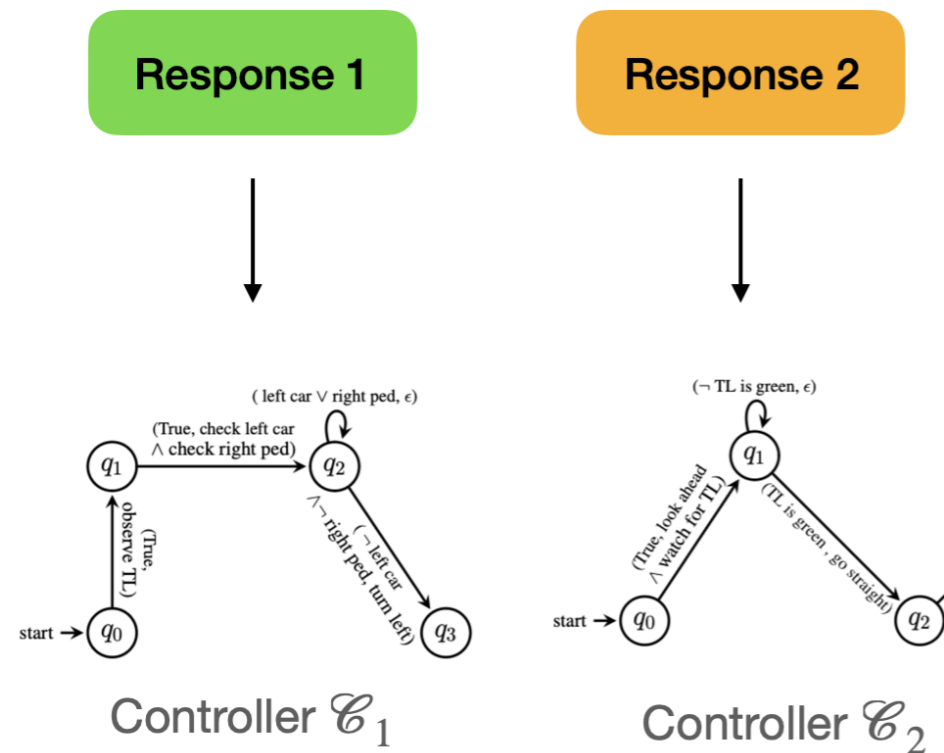


Fine-Tuning Language Models Using Formal Methods Feedback



Fine-Tuning Language Models Using Formal Methods Feedback

Empirical Evaluation via Simulation



High-Fidelity Simulator

Information from the Simulator

```

Vector3D(x=10.495529, y=5.662697, z=2.781272) Class: TL or TS
Vector3D(x=2.850960, y=18.144859, z=2.781269) Class: TL or TS
Vector3D(x=-9.272949, y=29.014732, z=2.781277) Class: TL or TS
-----
Vector3D(x=3.002785, y=17.807884, z=2.781323) Class: TL or TS
Vector3D(x=-9.121124, y=28.677757, z=2.781330) Class: TL or TS
-----
Vector3D(x=3.161827, y=17.457993, z=2.782099) Class: TL or TS
Vector3D(x=-8.962082, y=28.327866, z=2.782107) Class: TL or TS
    
```

Execution Trace

$(left_car \wedge go_straight, pedestrian \wedge stop, \neg pedestrian \wedge go_straight)$

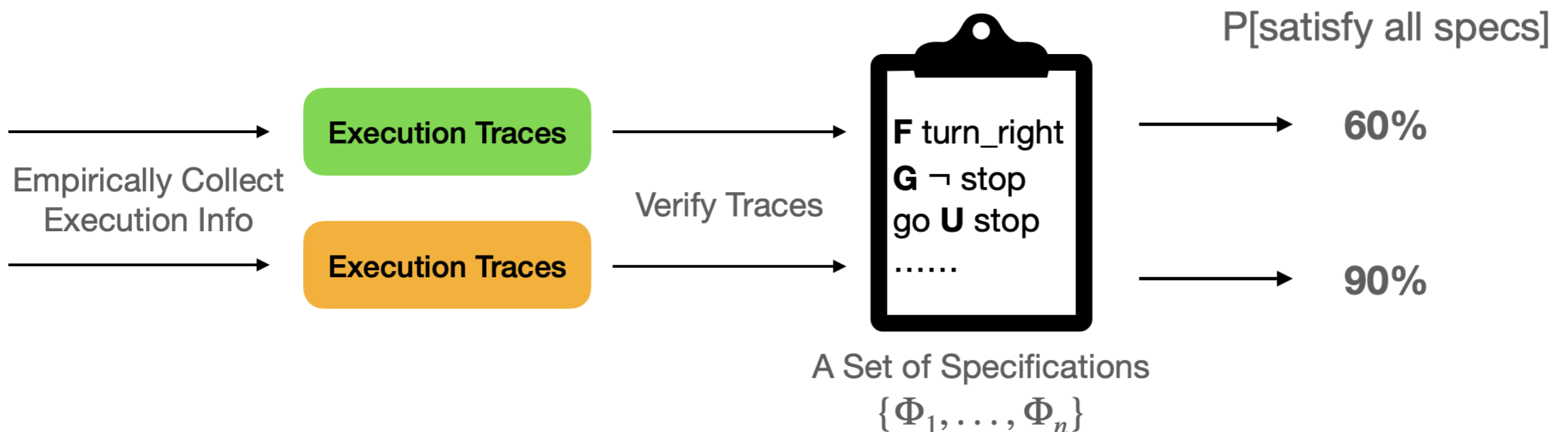
Empirically Collect Execution Info

Execution Traces

Execution Traces

Fine-Tuning Language Models Using Formal Methods Feedback

Empirical Evaluation via Simulation



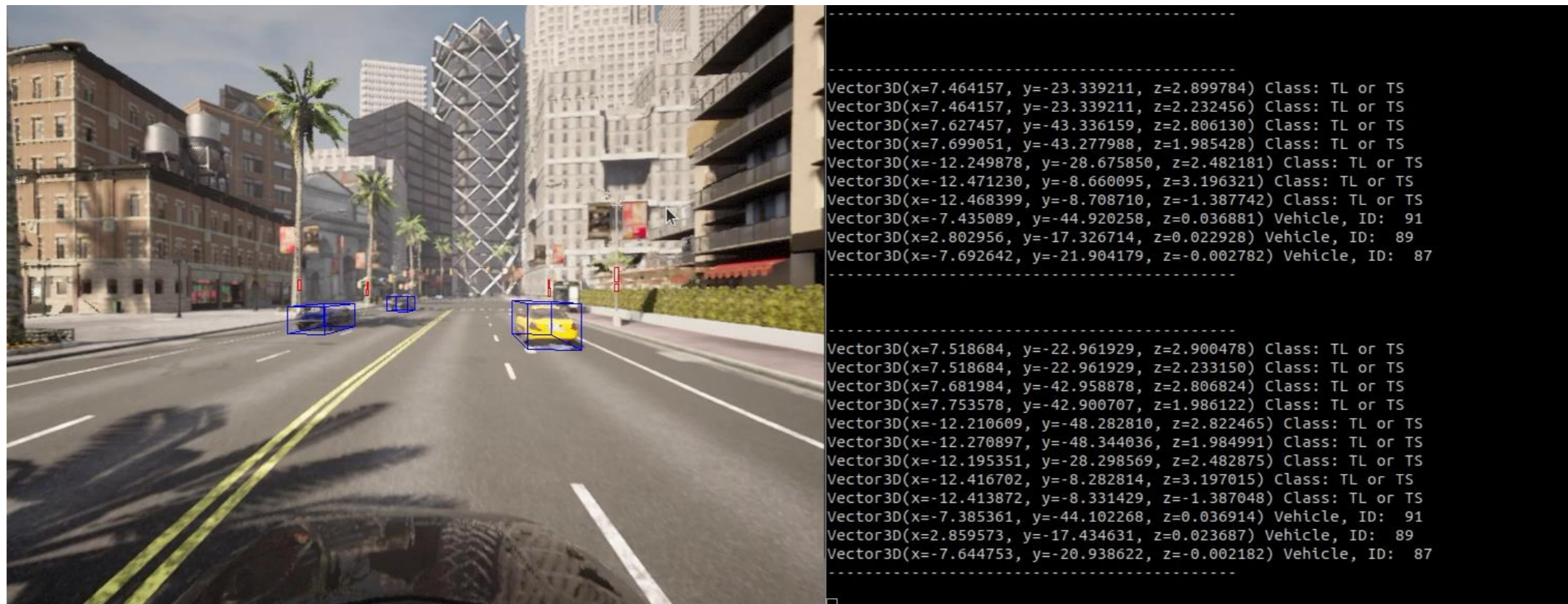
Quantitative Analysis

Empirical Evaluation via Simulation

Carla Simulator: Extract execution traces.

Carla Simulation Video

Object and Position Information

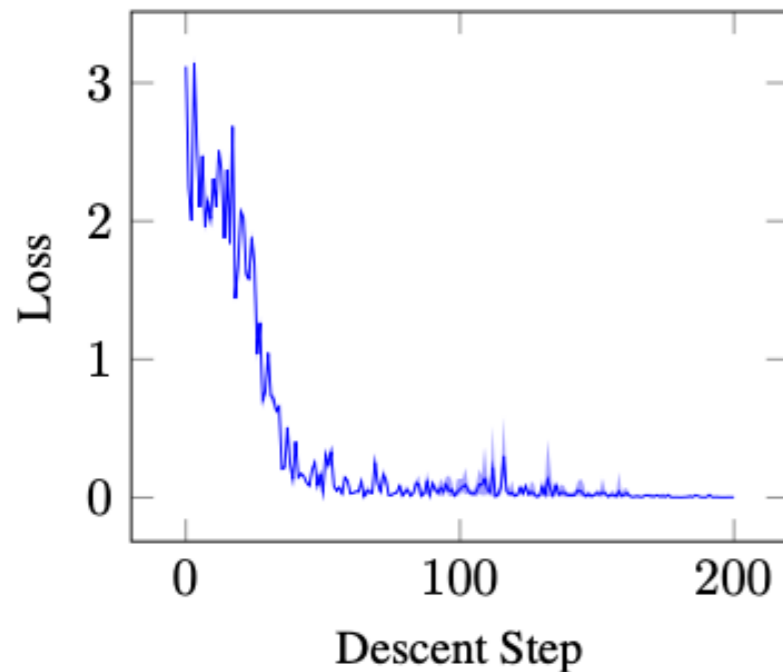


Execution Trace: (desired objects with positions, action),.....

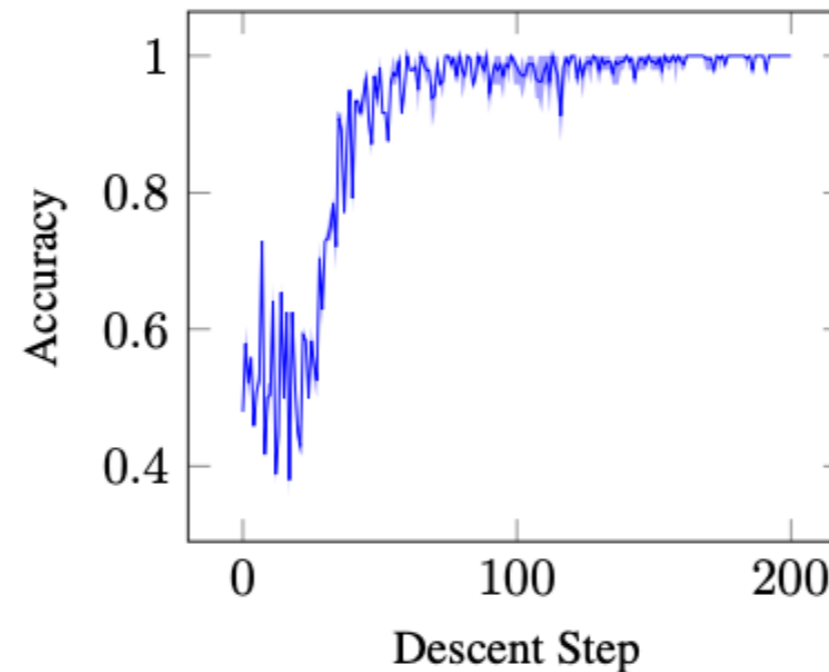
Quantitative Analysis

Formal Verification

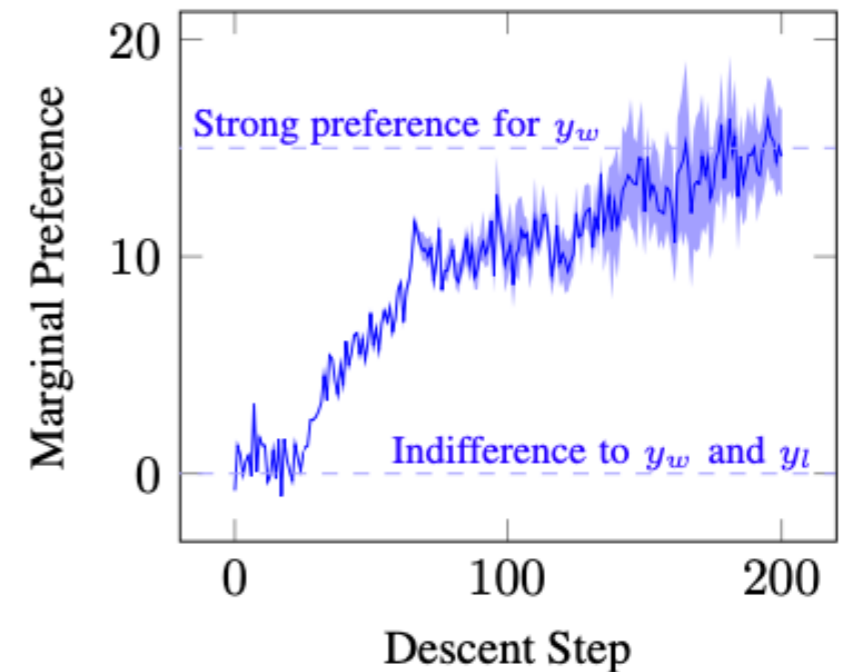
Training Loss



The percentage of the LLM outputs with preferred response rather than the unpreferred response.



How strong the LLM intends to output the preferred response (y_w) rather than the unpreferred response (y_l).

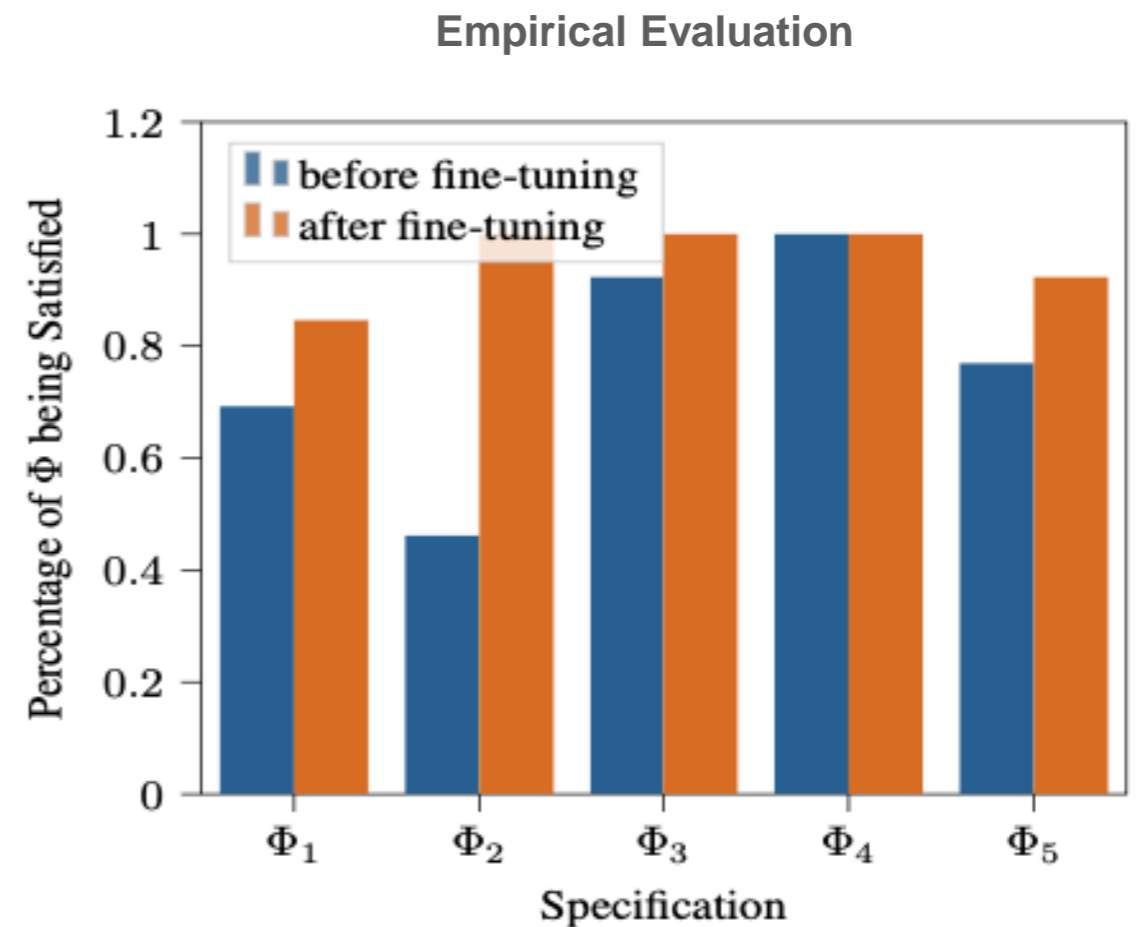
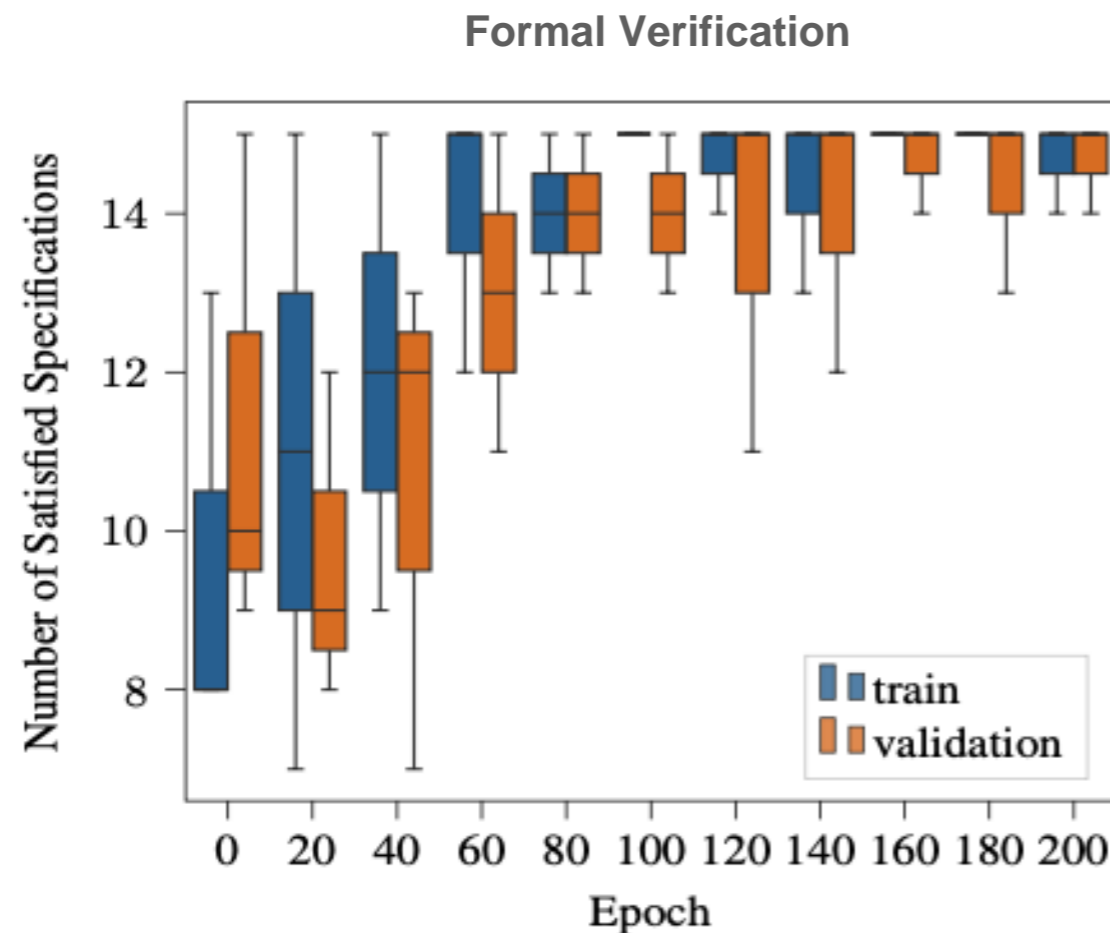


Training Performance

1. Training loss converges after 100 epochs.
2. Nearly 100% **preference accuracy**.
3. The fine-tuned language model strongly prefers the “correct” responses.

Quantitative Analysis

Empirical Evaluation via Simulation



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

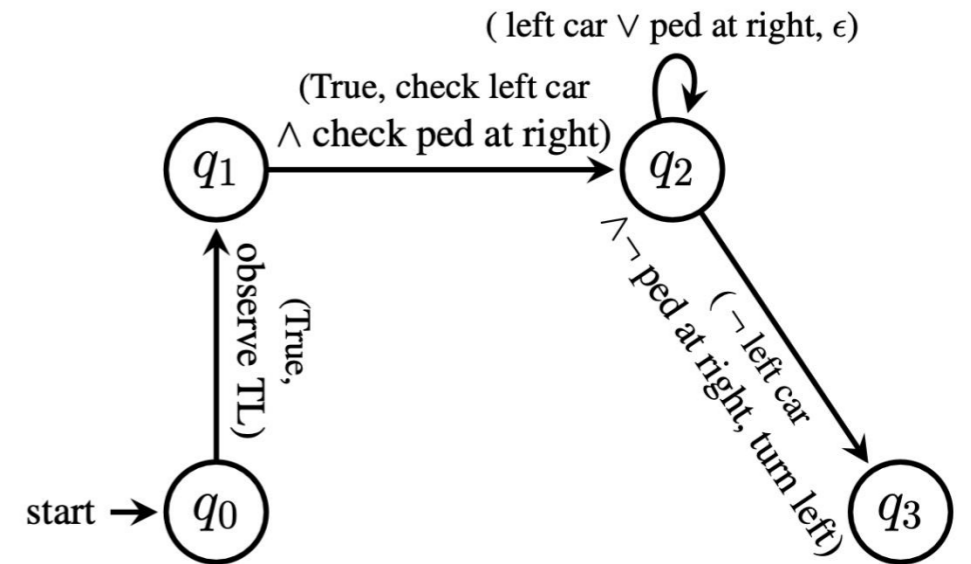
Our approach can act as a starting point to guide the design process for real-world implementations of autonomous driving systems (60% \rightarrow 90%).

A Step Toward Real-World Grounding

Empirical Evaluation via Simulation \rightarrow Real-World Performance

Statements

1. The controller's decisions are solely based on visual observations collected from the environment.
2. The vision model performs consistently in simulation and reality.

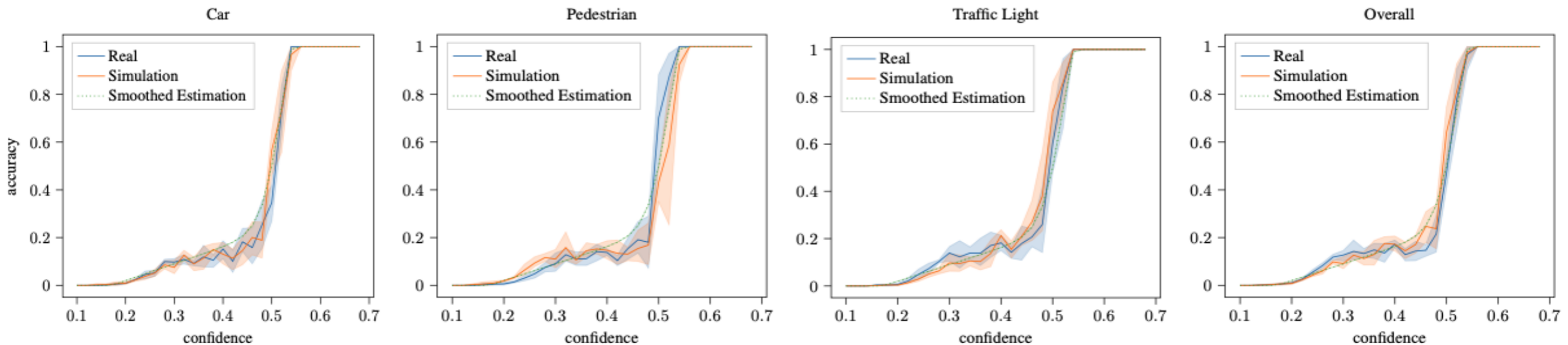
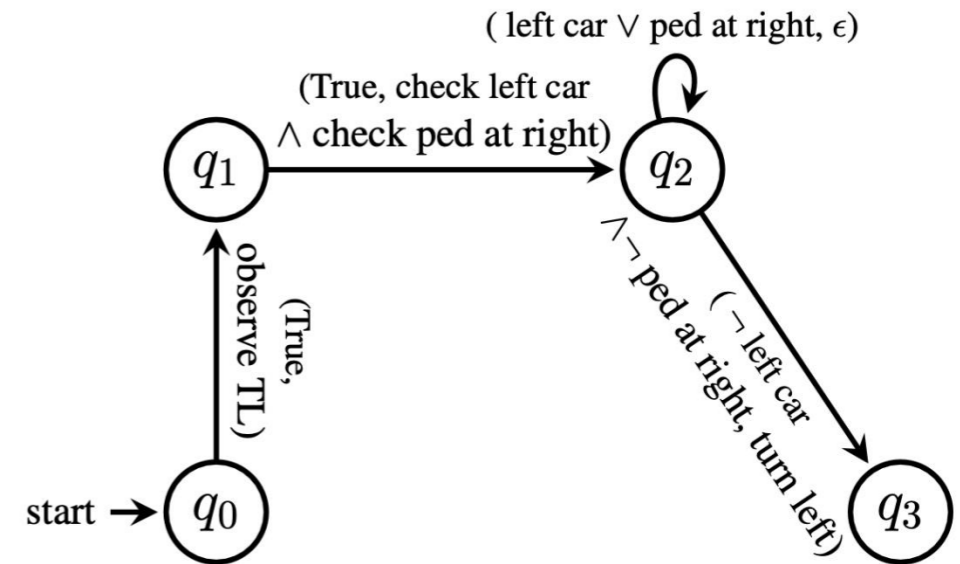


A Step Toward Real-World Grounding

Empirical Evaluation via Simulation \rightarrow Real-World Performance

Statements

1. The controller's decisions are solely based on visual observations collected from the environment.
2. The vision model performs consistently in simulation and reality.

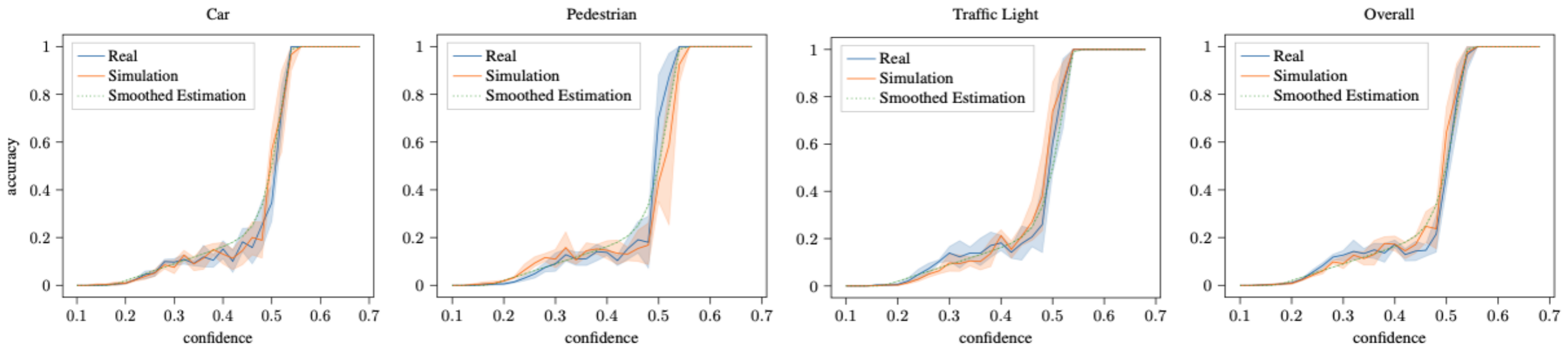
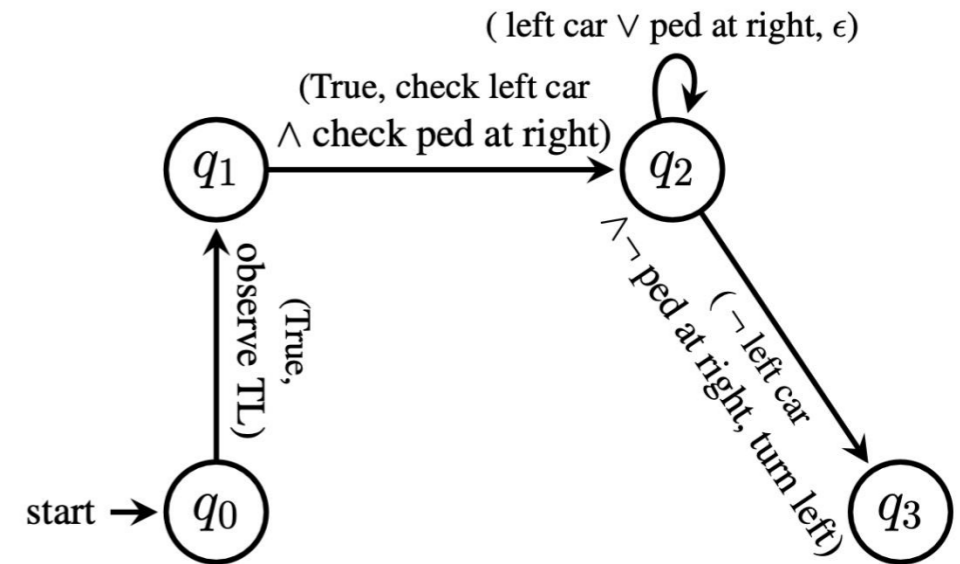


Extension: A Step Toward Real-World Grounding

Empirical Evaluation via Simulation \rightarrow Real-World Performance

Statements

1. The controller's decisions are solely based on visual observations collected from the environment.
2. The vision model performs consistently in simulation and reality.
3. If Statement 2 holds and if the controllers satisfy the critical specifications in simulation, then the controllers also satisfy the specifications in reality.



Extension: A Step Toward Real-World Grounding

```
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
```



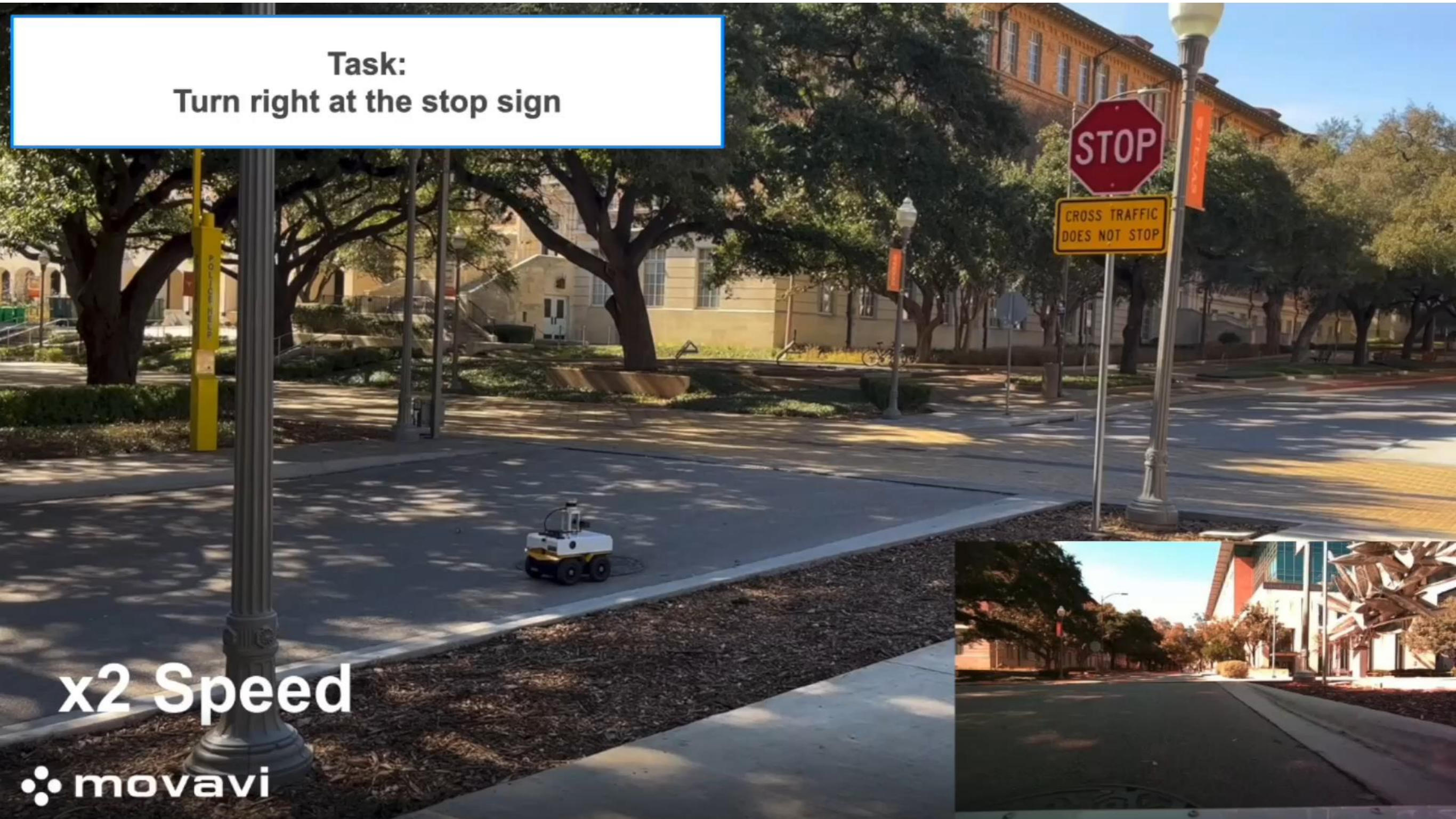
(left car \vee ped at right, ϵ)
 (True, check left car)

```
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"
```

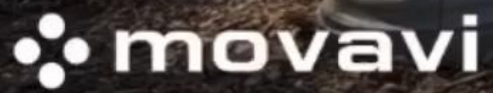
- (\neg green traffic light \rightarrow \neg go straight), ✓
- (stop sign \rightarrow \diamond stop), ✓
- \neg turn right \vee \neg (car from left \vee pedestrian at right), ✓

Experimental Evaluation

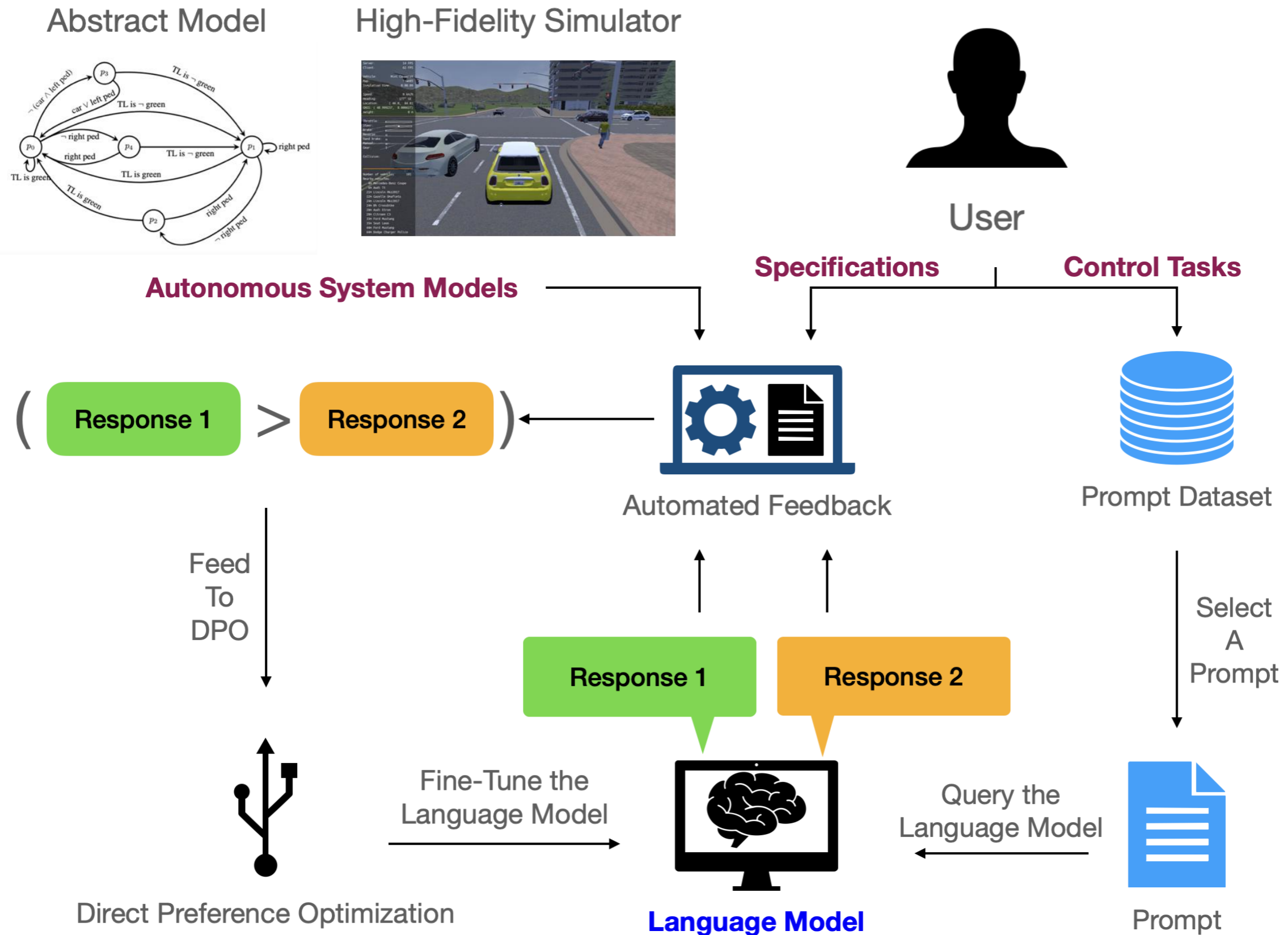
Task:
Turn right at the stop sign



x2 Speed



Fine-Tuning Language Models Using Formal Methods Feedback



Thank you!



Paper available at: <https://arxiv.org/abs/2310.18239>

