

# Verification of Autonomous Neural Car Control with KeYmaera X

ABZ 2025 Case Study Challenge

Enguerrand Prebet, Samuel Teuber, André Platzer | 12th of June 2025

# The ABZ Case Study



Leuschel et al. 2025

Symbolic dL-model for highway car control  
→ infinite-time guarantee: absence of collision

## What does that imply for concrete controllers?

Motivation



Modelling with dL



Applications of ModelPlex



Evaluation and the Model2Sim Gap



Conclusion



References

# The ABZ Case Study



Leuschel et al. 2025

Symbolic dL-model for highway car control  
→ infinite-time guarantee: absence of collision

**What does that imply for concrete Neural Network controllers?**

Motivation



Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion



References

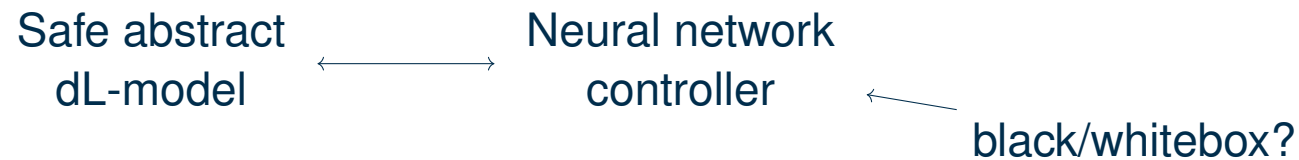
# The ABZ Case Study



Leuschel et al. 2025

Symbolic dL-model for highway car control  
→ infinite-time guarantee: absence of collision

## What does that imply for concrete Neural Network controllers?



Motivation



Modelling with dL



Applications of ModelPlex



Evaluation and the Model2Sim Gap



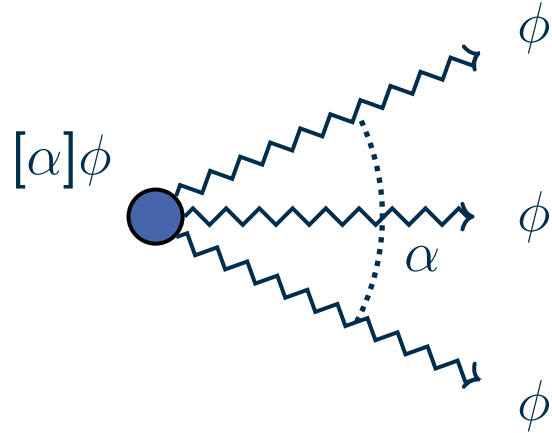
Conclusion



References

# Differential Dynamic Logic

FOL( $\mathbb{R}$ ) + program modalities



Hoare triple:  $\text{init} \rightarrow [\text{sys}]\text{post}$

Motivation

○

Modelling with dL

● ○ ○ ○ ○

Applications of ModelPlex

○ ○ ○ ○

Evaluation and the Model2Sim Gap

○ ○ ○ ○ ○ ○

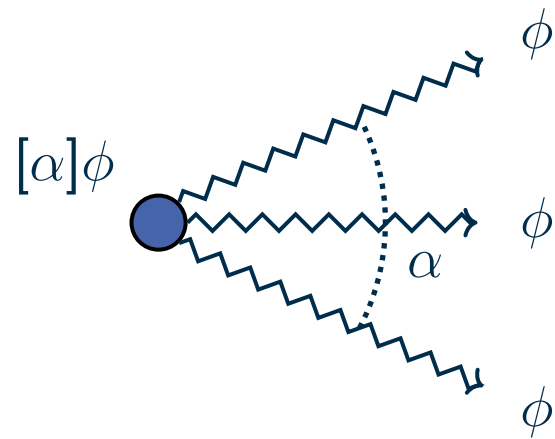
Conclusion

○

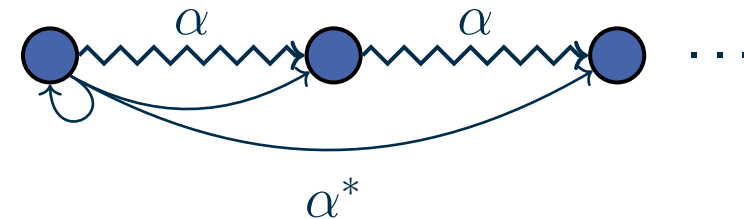
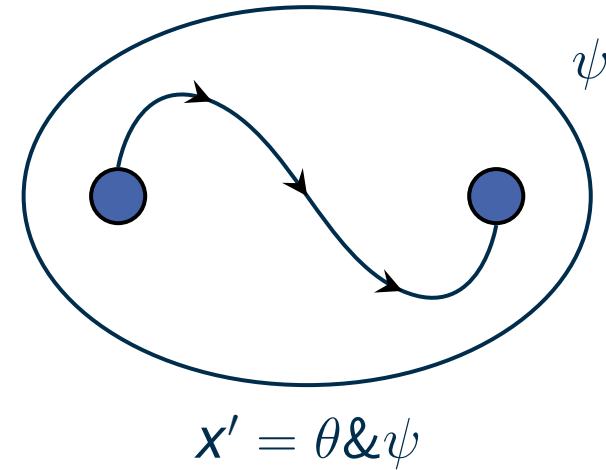
References

# Differential Dynamic Logic

FOL( $\mathbb{R}$ ) + program modalities + differential systems



Hoare triple:  $\text{init} \rightarrow [\text{sys}]\text{post}$



Motivation  
○

Modelling with dL  
● ○ ○ ○ ○

Applications of ModelPlex  
○ ○ ○ ○

Evaluation and the Model2Sim Gap  
○ ○ ○ ○ ○ ○

Conclusion  
○

References

# Proving properties

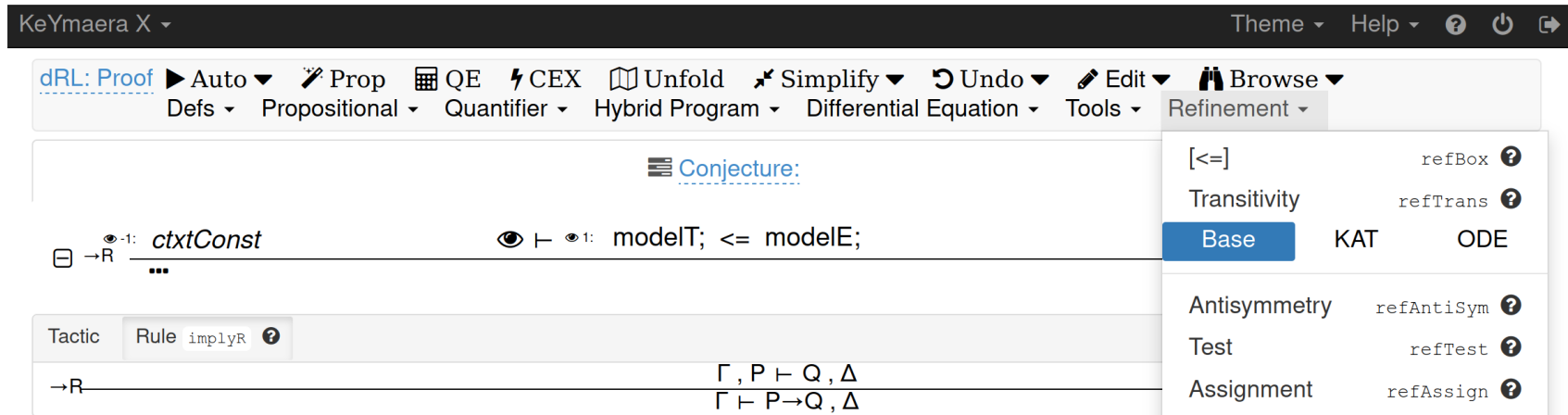
Uniform-substitution based calculus:

$$\begin{aligned} p() &\rightarrow [a]p() \\ [x := f()]p(x) &\leftrightarrow p(f()) \end{aligned}$$

$$(US) \frac{\phi}{\sigma(\phi)} \text{ if } \sigma(\phi) \text{ defined}$$

Refinements as formulas:  $\alpha \leq \beta$

All implemented in theorem prover KeYmaera X



# Our dL Model

Two unordered cars  $\rightarrow$  core question, even for multilane

$$\text{sys} ::= \underbrace{\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T)}_{\text{control}};$$

■  $\text{ctrl}_o$ : sets  $a_o$  to a value in  $[-B_{\max}, A_{\max}]$

Motivation

○

Modelling with dL

○○●○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion

○

References



# Our dL Model

Two unordered cars  $\rightarrow$  core question, even for multilane

$$\text{sys} ::= \underbrace{\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T)}_{\text{control}};$$

- $\text{ctrl}_o$ : sets  $a_o$  to a value in  $[-B_{\max}, A_{\max}]$
- $\text{ctrl}_e$ :  $a_e$ , if  $\neg \text{safe}(a_e)$ , overrides with one in  $[-B_{\max}, -B_{\min}]$  behind  
RSS-like  $[A_{\min}, A_{\max}]$  in front

Motivation

○

Modelling with dL

○○●○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion

○

References

# Our dL Model

Two unordered cars  $\rightarrow$  core question, even for multilane

$$\text{sys} ::= \underbrace{\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T)}_{\text{control}}; \underbrace{\text{accelCorr}}_{\text{plant}}; \text{dyn}$$

- $\text{ctrl}_o$ : sets  $a_o$  to a value in  $[-B_{\max}, A_{\max}]$
- $\text{ctrl}_e$ :  $a_e$ , if  $\neg \text{safe}(a_e)$ , overrides with one in  $[-B_{\max}, -B_{\min}]$  behind  
RSS-like  $[A_{\min}, A_{\max}]$  in front
- $\text{accelCorr}$ : ensures  $0 \leq v_e, v_o \leq V$

Motivation

○

Modelling with dL

○○●○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion

○

References

# Our dL Model

Two unordered cars  $\rightarrow$  core question, even for multilane

$$\text{sys} ::= \underbrace{\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T)}_{\text{control}}; \underbrace{\text{accelCorr}; \text{dyn}}_{\text{plant}}$$

- $\text{ctrl}_o$ : sets  $a_o$  to a value in  $[-B_{\max}, A_{\max}]$
- $\text{ctrl}_e$ :  $a_e$ , if  $\neg \text{safe}(a_e)$ , overrides with one in  $[-B_{\max}, -B_{\min}]$  behind  
RSS-like  $[A_{\min}, A_{\max}]$  in front
- $\text{accelCorr}$ : ensures  $0 \leq v_e, v_o \leq V$
- $\text{dyn}$ : cars move

$$\begin{aligned}x'_e &= v_e, v'_e = a_e, \\x'_o &= v_o, v'_o = a_o, \ \& \ t \leq t_e + T \\t' &= 1\end{aligned}$$

# Our dL Model

Two unordered cars  $\rightarrow$  core question, even for multilane

$$\text{sys} ::= \underbrace{(\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T))}_{\text{control}}; \underbrace{\text{accelCorr}; \text{dyn}}_{\text{plant}})^*$$

- $\text{ctrl}_o$ : sets  $a_o$  to a value in  $[-B_{\max}, A_{\max}]$
- $\text{ctrl}_e$ :  $a_e$ , if  $\neg \text{safe}(a_e)$ , overrides with one in  $[-B_{\max}, -B_{\min}]$  behind  
RSS-like  $[A_{\min}, A_{\max}]$  in front
- $\text{accelCorr}$ : ensures  $0 \leq v_e, v_o \leq V$
- $\text{dyn}$ : cars move

$$\begin{aligned}x'_e &= v_e, v'_e = a_e, \\x'_o &= v_o, v'_o = a_o, \ \& \ t \leq t_e + T \\t' &= 1\end{aligned}$$

# Our dL Model

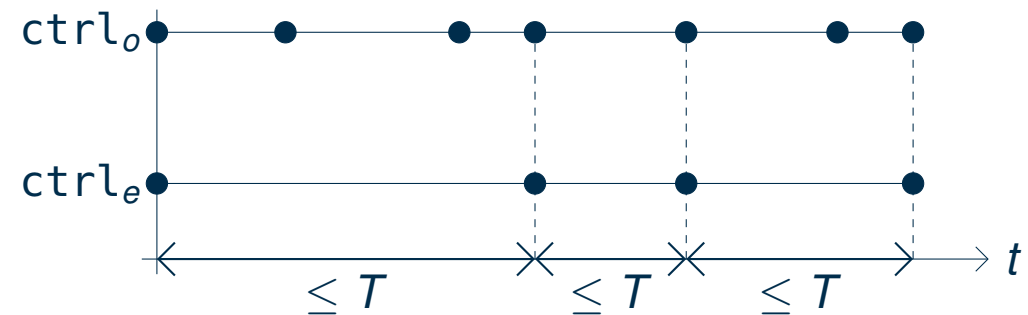
Two unordered cars  $\rightarrow$  core question, even for multilane

$$\text{sys} ::= \underbrace{(\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T))}_{\text{control}}; \underbrace{\text{accelCorr; dyn}}_{\text{plant}})^*$$

- $\text{ctrl}_o$ : sets  $a_o$  to a value in  $[-B_{\max}, A_{\max}]$
- $\text{ctrl}_e$ :  $a_e$ , if  $\neg \text{safe}(a_e)$ , overrides with one in  $[-B_{\max}, -B_{\min}]$  behind  $[A_{\min}, A_{\max}]$  in front
- $\text{accelCorr}$ : ensures  $0 \leq v_e, v_o \leq V$  RSS-like
- $\text{dyn}$ : cars move

$$\begin{aligned} x'_e &= v_e, v'_e = a_e, \\ x'_o &= v_o, v'_o = a_o, \text{ \& } t \leq t_e + T \\ t' &= 1 \end{aligned}$$

Desynchronised controllers:



# Safety proofs

## Theorem

*These formulas are proved in dL:*

$$\text{ctx} \wedge x_e + L \leq x_o \wedge \text{init} \rightarrow [\text{sys}]x_e + L \leq x_o$$

$$\text{ctx} \wedge x_o + L \leq x_e \wedge \widetilde{\text{init}} \rightarrow [\text{sys}]x_o + L \leq x_e$$

$$\text{init} ::= x_e + \frac{v_e^2}{2B_{\min}} + L \leq x_o + \frac{v_o^2}{2B_{\max}}$$

$$\widetilde{\text{init}} ::= x_e + \frac{(v_e - V)^2}{2(-A_{\min})} + L \leq x_o + \frac{(v_o - V)^2}{2(-A_{\max})}$$

KeYmaera X proofs and experiments online: <https://doi.org/10.5281/zenodo.14959858>

Motivation

○

Modelling with dL

○○○●○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion

○

References

# ModelPlex for Neural Network Controllers



Motivation



Modelling with dL



Applications of ModelPlex



Evaluation and the Model2Sim Gap

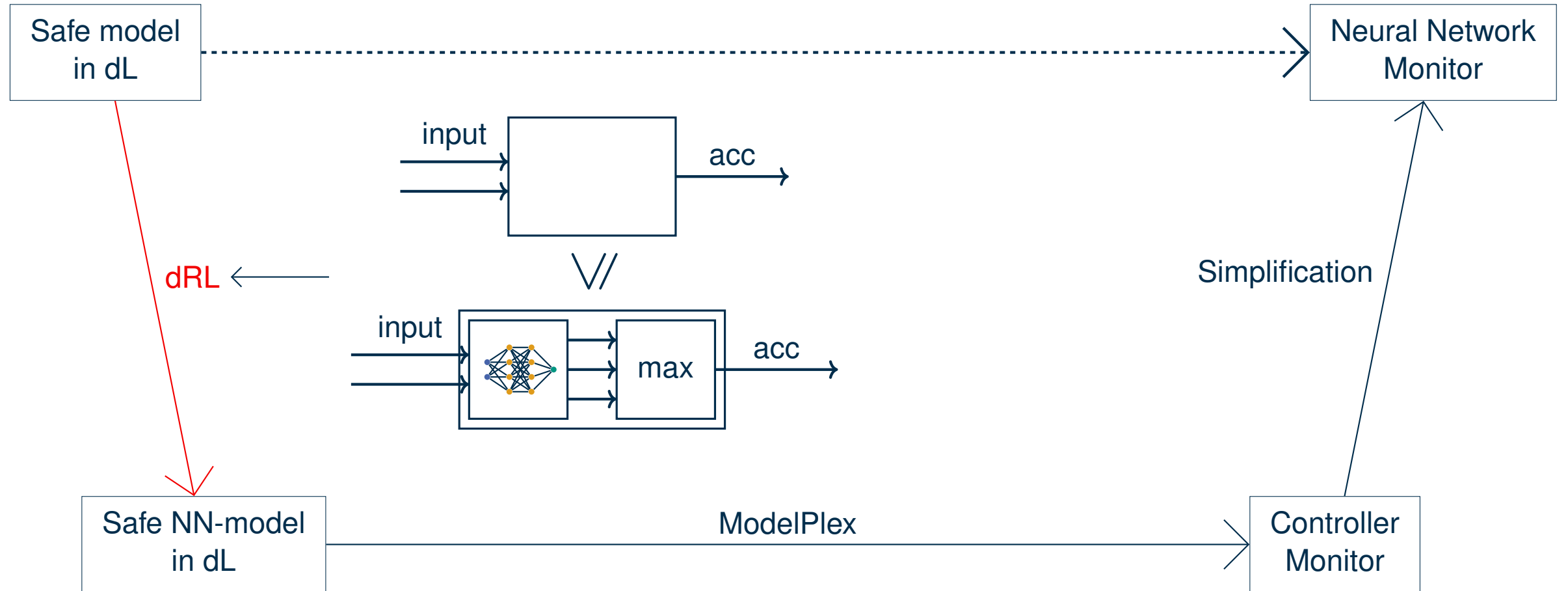


Conclusion



References

# ModelPlex for Neural Network Controllers



Motivation  
○

Modelling with dL  
○○○○●

Applications of ModelPlex  
○○○○

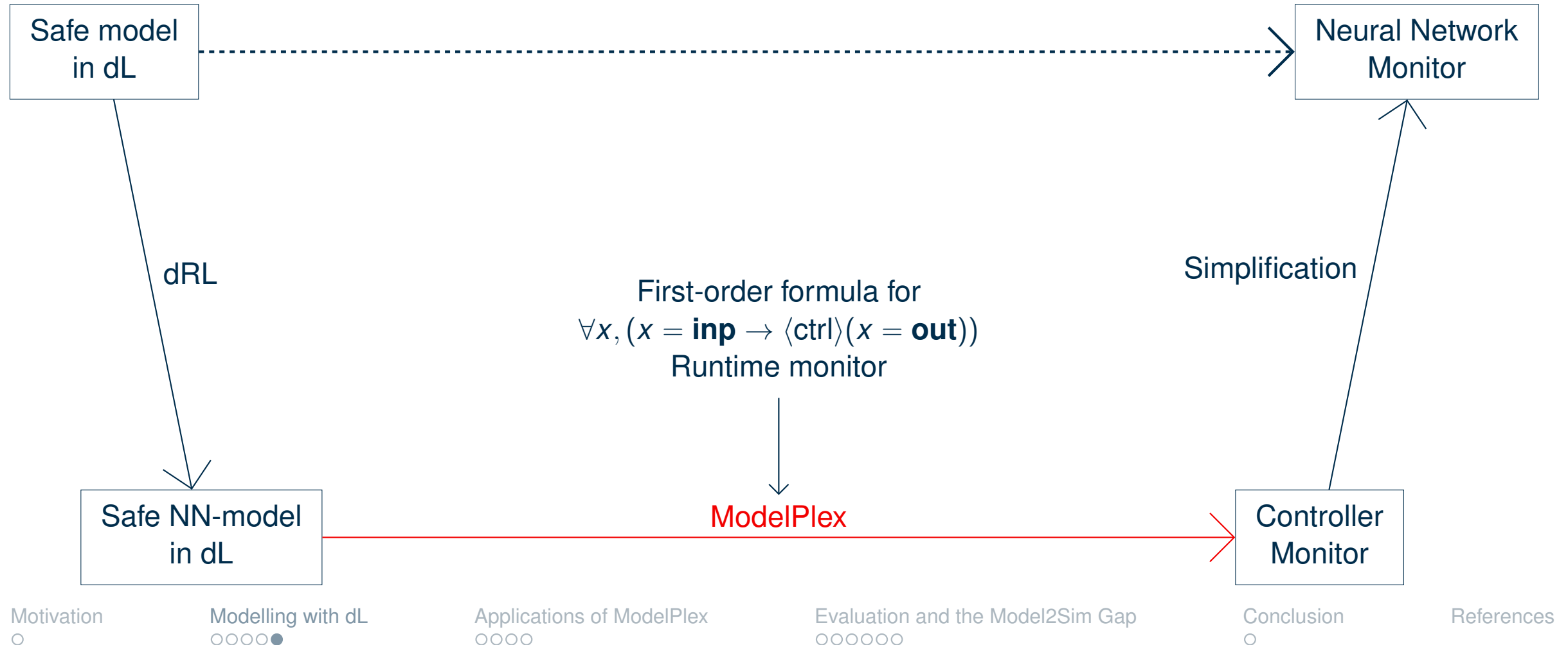
Evaluation and the Model2Sim Gap  
○○○○○

Conclusion  
○

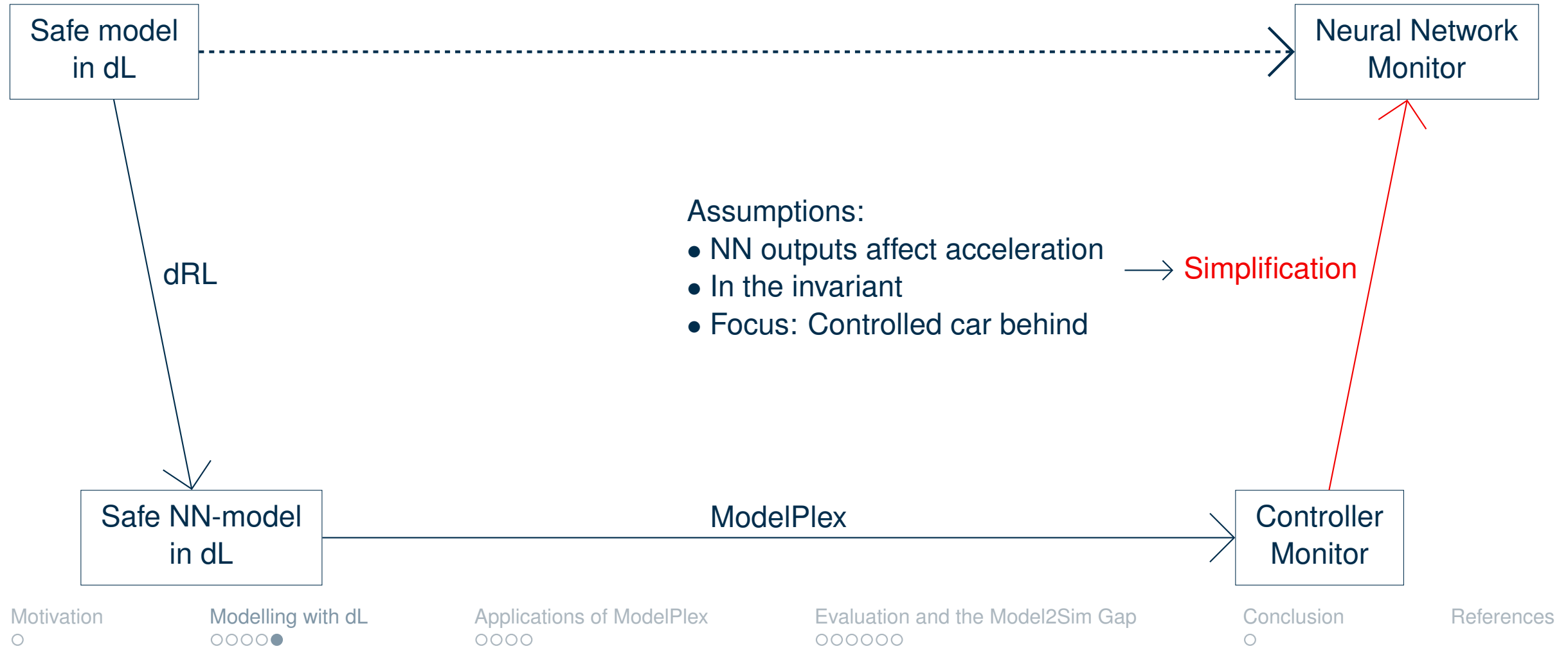
References



# ModelPlex for Neural Network Controllers



# ModelPlex for Neural Network Controllers



# Applications of ModelPlex

## Via dL: Correct by Construction Monitoring Condition

$$\begin{aligned} & y_1^+ \geq y_2^+ \wedge y_1^+ \geq y_3^+ \\ & \vee y_2^+ > y_1^+ \wedge y_2^+ \geq y_3^+ \wedge \\ & \quad (B_{\min} \leq 0 \leq A_{\max} \wedge v_e \geq 0 \wedge \text{pos}_e(B_{\min}) + (\frac{0}{B_{\min}} + 1)Tv_e + L < \text{pos}_o) \\ & \vee y_3^+ > y_1^+ \wedge y_3^+ > y_2^+ \wedge (B_{\min} \leq A_{\max} \wedge v_e + A_{\max}T < 0 \wedge \text{pos}_e(A_{\max}) + L < \text{pos}_o) \\ & \vee B_{\min} \leq A_{\max} \wedge v_e + A_{\max}T \geq 0 \wedge \text{pos}_e(B_{\min}) + (\frac{-A_{\max}}{B_{\min}} + 1)(\frac{A_{\max}}{2}T^2 + Tv_e) + L < \text{pos}_o) \end{aligned}$$

Given concrete inputs and outputs, this form tells us what actions are provably safe.

**But how do we put this knowledge into practice?**

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

●○○○

Evaluation and the Model2Sim Gap

○○○○○

Conclusion

○

References

# Applications of ModelPlex

## Via dL: Correct by Construction Monitoring Condition

$$\begin{aligned} & y_1^+ \geq y_2^+ \wedge y_1^+ \geq y_3^+ \\ & \vee y_2^+ > y_1^+ \wedge y_2^+ \geq y_3^+ \wedge \\ & \quad (B_{\min} \leq 0 \leq A_{\max} \wedge v_e \geq 0 \wedge \text{pos}_e(B_{\min}) + (\frac{0}{B_{\min}} + 1)Tv_e + L < \text{pos}_o) \\ & \vee y_3^+ > y_1^+ \wedge y_3^+ > y_2^+ \wedge (B_{\min} \leq A_{\max} \wedge v_e + A_{\max}T < 0 \wedge \text{pos}_e(A_{\max}) + L < \text{pos}_o) \\ & \vee B_{\min} \leq A_{\max} \wedge v_e + A_{\max}T \geq 0 \wedge \text{pos}_e(B_{\min}) + (\frac{-A_{\max}}{B_{\min}} + 1)(\frac{A_{\max}}{2}T^2 + Tv_e) + L < \text{pos}_o) \end{aligned}$$

Given concrete inputs and outputs, this form tells us what actions are provably safe.

**But how do we put this knowledge into practice?**

Monitoring

Shielding

Verification

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

●○○○

Evaluation and the Model2Sim Gap

○○○○○

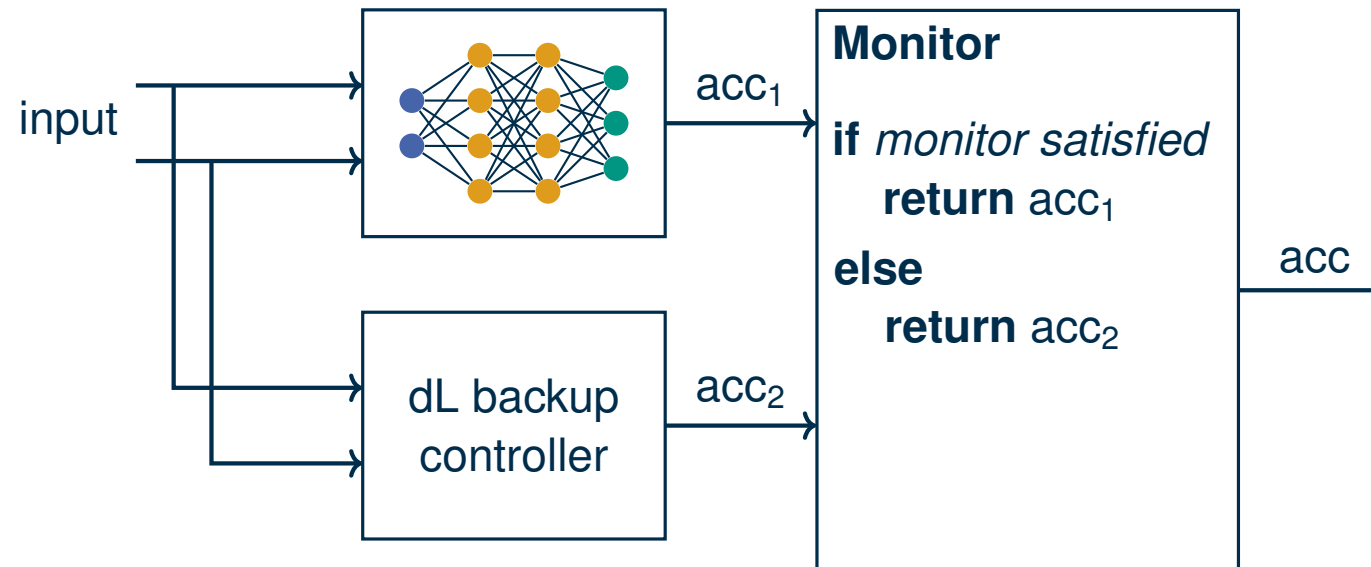
Conclusion

○

References

# Monitoring / Sandboxing (VeriPhy)

Check NN actions during runtime at each step



Can be combined with correct-by-construction sandbox synthesis (Bohrer et al. 2018)

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○●○○○

Evaluation and the Model2Sim Gap

○○○○○○

Mitsch and Platzer 2016

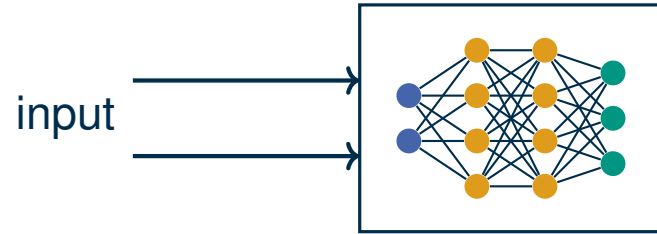
Conclusion

○

References

# Shielding (Justified Speculative Control)

**Insight:** RL Agents often learn a **distribution** of actions  
⇒ Constrain action space



Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○●○

Evaluation and the Model2Sim Gap

○○○○○

Fulton and Platzer 2018

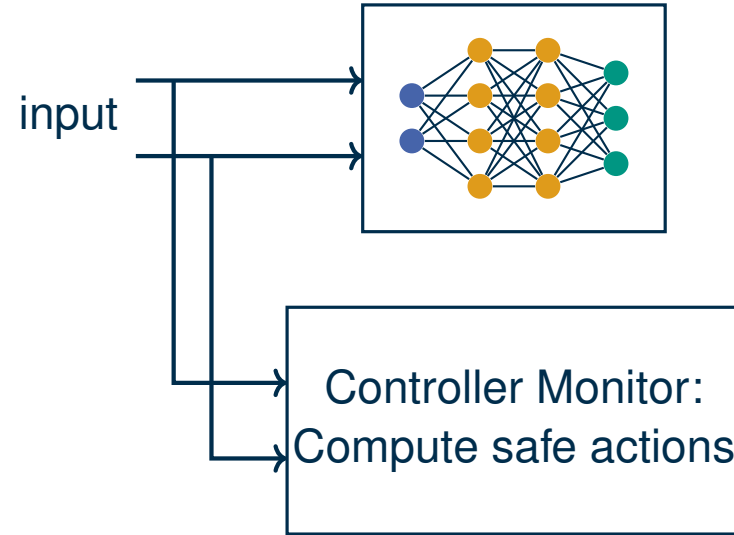
Conclusion

○

References

# Shielding (Justified Speculative Control)

**Insight:** RL Agents often learn a **distribution** of actions  
⇒ Constrain action space



Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○●○

Evaluation and the Model2Sim Gap  
○○○○○

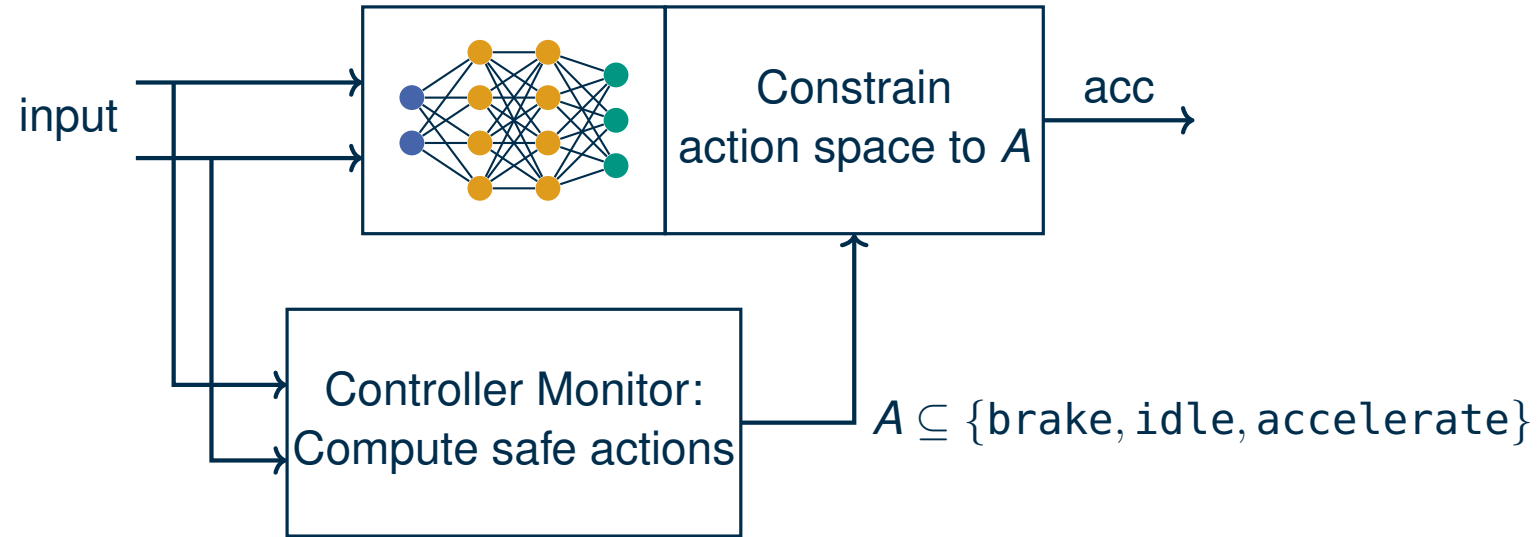
Fulton and Platzer 2018

Conclusion  
○

References

# Shielding (Justified Speculative Control)

**Insight:** RL Agents often learn a **distribution** of actions  
 $\Rightarrow$  Constrain action space



Provably safe actions during **training & deployment!**  
Can also take into account **model monitoring**

Fulton and Platzer 2018

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○●○

Evaluation and the Model2Sim Gap  
○○○○○

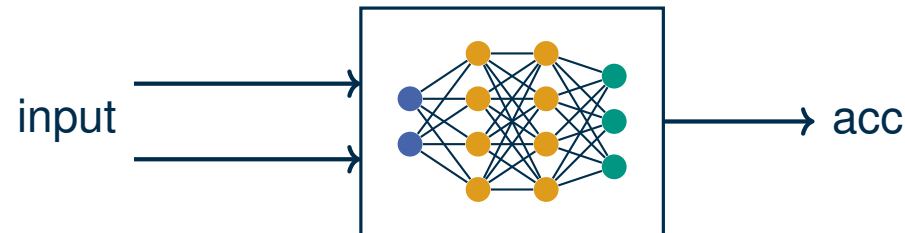
Conclusion  
○

References



# Verification (VerSAILLE)

**Objective:** A priori guarantees on safety of NN controller



Teuber, Mitsch, and Platzer 2024

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○●

Evaluation and the Model2Sim Gap

○○○○○

Conclusion

○

References

# Verification (VerSAILLE)

**Objective:** A priori guarantees on safety of NN controller

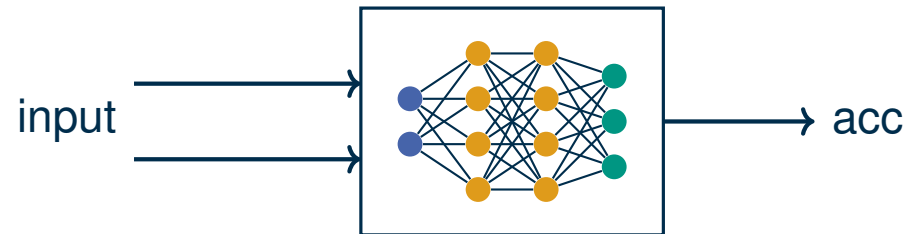
Neural Network Verification:

$$\forall \bar{x} \quad \phi(\bar{x}, g(\bar{x}))$$

**Before Deployment**

---

**At Runtime**



Teuber, Mitsch, and Platzer 2024

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○●

Evaluation and the Model2Sim Gap

○○○○○

Conclusion

○

References

# Verification (VerSAILLE)

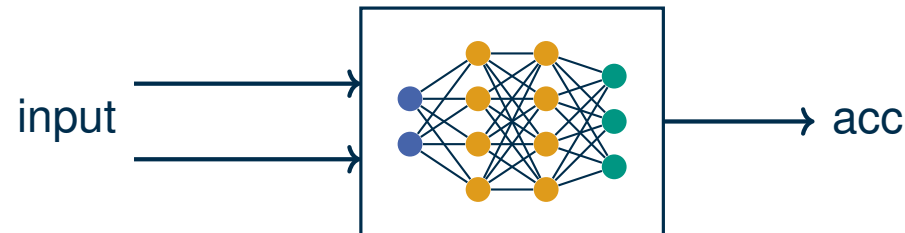
**Objective:** A priori guarantees on safety of NN controller

$$(\alpha_{\text{ctrl}} ; \alpha_{\text{plant}})^* \quad \text{Safe}$$

**Before Deployment**

---

**At Runtime**



Teuber, Mitsch, and Platzer 2024

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○●

Evaluation and the Model2Sim Gap

○○○○○

Conclusion

○

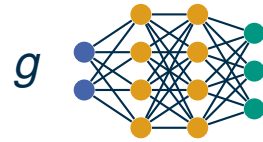
References

# Verification (VerSAILLE)

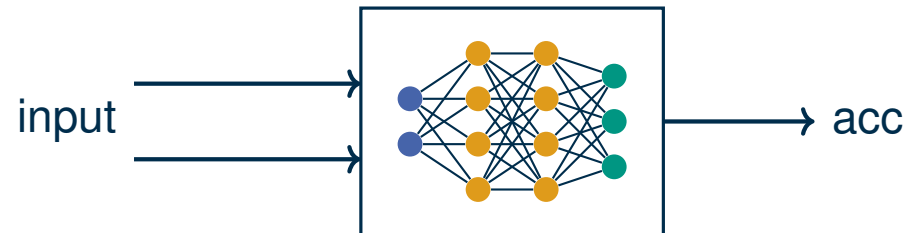
**Objective:** A priori guarantees on safety of NN controller

$$(\alpha_{\text{ctrl}} ; \alpha_{\text{plant}})^* \text{ Safe}$$

**Before Deployment**



**At Runtime**



Teuber, Mitsch, and Platzer 2024

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○●

Evaluation and the Model2Sim Gap

○○○○○

Conclusion

○

References

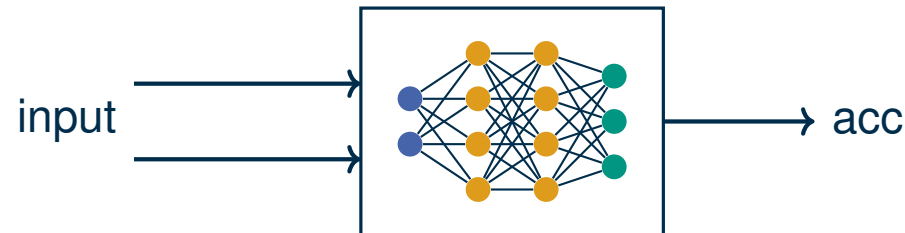
# Verification (VerSAILLE)

**Objective:** A priori guarantees on safety of NN controller

$$(\alpha_{\text{ctrl}} ; \alpha_{\text{plant}})^* \quad \text{Safe}$$



**At Runtime**



Teuber, Mitsch, and Platzer 2024

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○●

Evaluation and the Model2Sim Gap

○○○○○

Conclusion

○

References

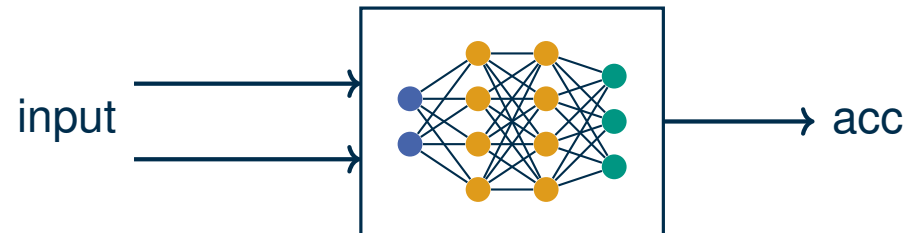
# Verification (VerSAILLE)

**Objective:** A priori guarantees on safety of NN controller

$$\text{invariant} \wedge \neg \text{monitor} \xleftarrow[\text{based on ModelPlex \& Loop Invariants}]{\text{VerSAILLE}} (\alpha_{\text{ctrl}} ; \alpha_{\text{plant}})^* \quad \text{Safe}$$



**At Runtime**



Teuber, Mitsch, and Platzer 2024

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○●

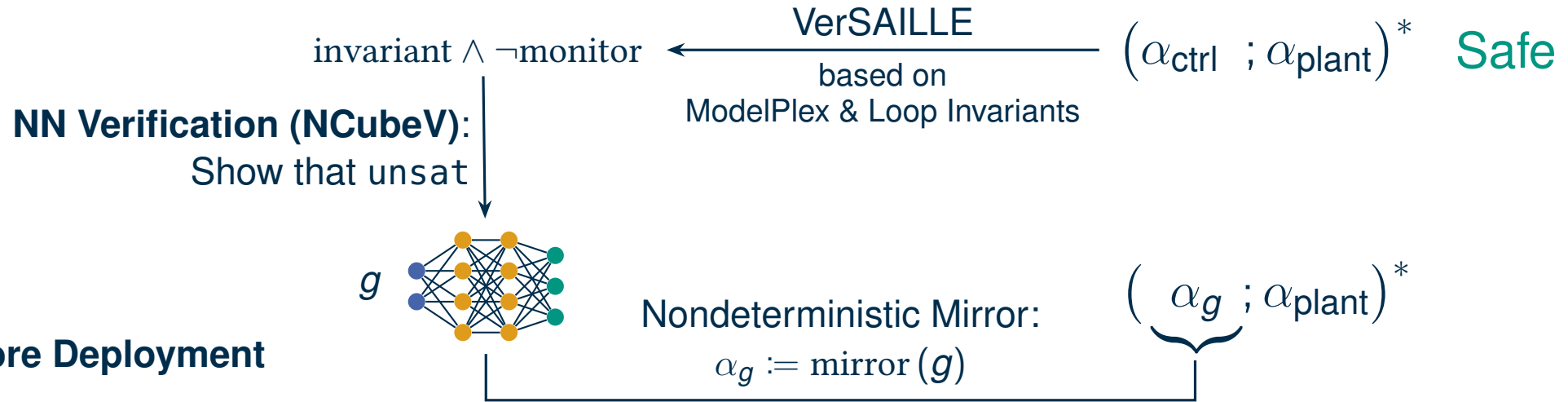
Evaluation and the Model2Sim Gap  
○○○○○

Conclusion  
○

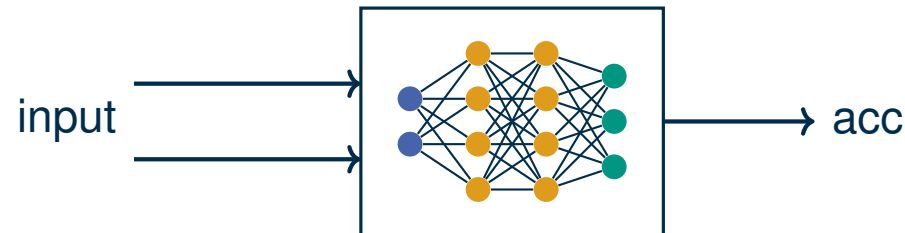
References

# Verification (VerSAILLE)

**Objective:** A priori guarantees on safety of NN controller



**At Runtime**



Teuber, Mitsch, and Platzer 2024

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○●

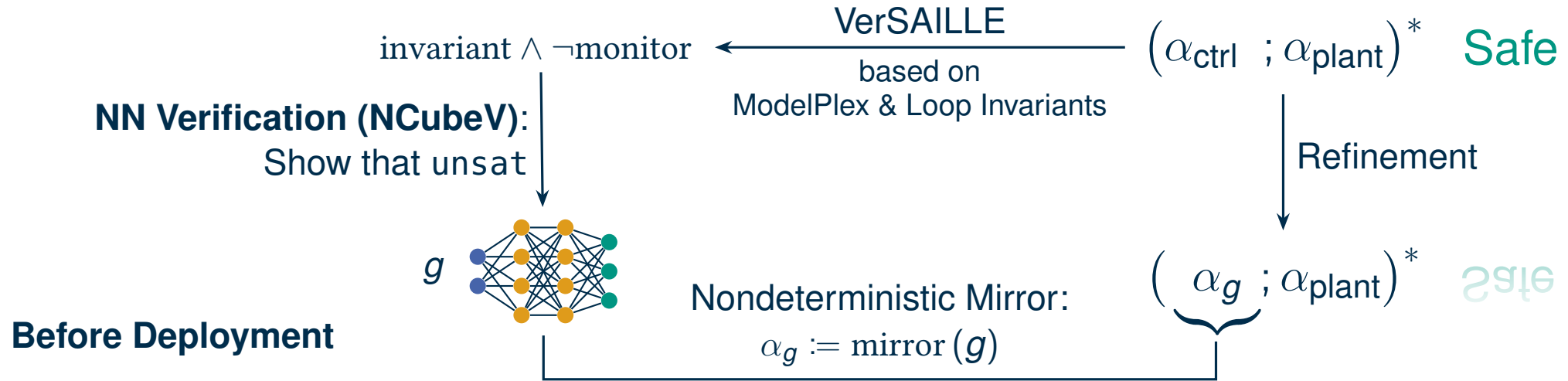
Evaluation and the Model2Sim Gap  
○○○○○

Conclusion  
○

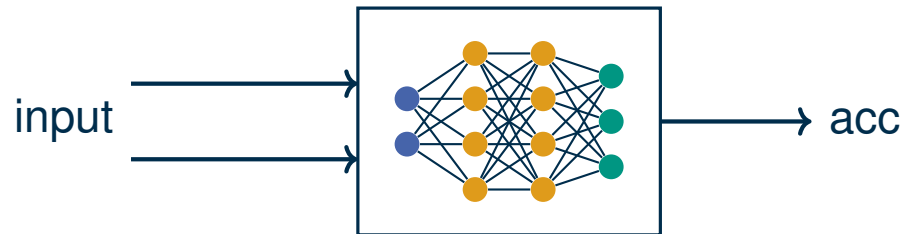
References

# Verification (VerSAILLE)

**Objective:** A priori guarantees on safety of NN controller



**At Runtime**



**A priori and infinite-time horizon safety guarantees**

Teuber, Mitsch, and Platzer 2024

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○●

Evaluation and the Model2Sim Gap  
○○○○○

Conclusion  
○

References



# Model-To-Simulation Gap (1)

## Behaviour of FASTER

### Spec:

*“This action **increases the speed (up to  $v_{max}$ ) with an acceleration up to  $a_{max}$  m/s<sup>2</sup>.** Once the car reaches  $v_{max}$ , the acceleration is 0 m/s<sup>2</sup>.”*

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○

Evaluation and the Model2Sim Gap  
●○○○○○

Conclusion  
○

References

# Model-To-Simulation Gap (1)

## Behaviour of FASTER

### Spec:

*“This action **increases the speed (up to  $v_{max}$ ) with an acceleration** up to  $a_{max} \text{ m/s}^2$ .  
Once the car reaches  $v_{max}$ , the acceleration is  $0 \text{ m/s}^2$ .”*

### Simulator:

Uses the configuration DiscreteMetaAction: FASTER increases the **reference velocity**  $v_r$ .  
Subsequently, a **low-level** proportional controller adjusts the acceleration.

# Model-To-Simulation Gap (1)

## Behaviour of FASTER

### Spec:

*“This action **increases the speed (up to  $v_{max}$ ) with an acceleration up to  $a_{max} \text{ m/s}^2$ .**  
Once the car reaches  $v_{max}$ , the acceleration is  $0 \text{ m/s}^2$ .”*

### Simulator:

Uses the configuration DiscreteMetaAction: FASTER increases the **reference velocity  $v_r$** .  
Subsequently, a **low-level** proportional controller adjusts the acceleration.

⇒ **FASTER can lead to braking if  $v_r < v$ !**

# Model-To-Simulation Gap (1)

## Behaviour of FASTER

### Spec:

*“This action **increases the speed (up to  $v_{max}$ ) with an acceleration** up to  $a_{max} \text{ m/s}^2$ .  
Once the car reaches  $v_{max}$ , the acceleration is  $0 \text{ m/s}^2$ .”*

### Simulator:

Uses the configuration DiscreteMetaAction: FASTER increases the **reference velocity**  $v_r$ .  
Subsequently, a **low-level** proportional controller adjusts the acceleration.

⇒ **FASTER can lead to braking if  $v_r < v$ !**

*We adjusted the simulator's configuration and retrained a new set of NNs using the provided scripts.*

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

●○○○○○

Conclusion

○

References

# Experimental Results (1)

## A first concrete NN

Performance: Standalone NN / Monitoring / Shielding					
Original NN		Monitoring (VeriPhy)		Shielding (JSC)	
Reward	Crash	Reward	Crash	Reward	Crash
<b>17.63</b> $\pm$ 0.21	<b>0</b> %	16.72 $\pm$ 0.32	<b>0</b> %	<b>17.63</b> $\pm$ 0.21	<b>0</b> %

This looks good – let's verify it!

(1000 simulations)

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○●○○○○

Conclusion

○

References

# Experimental Results (1)

## A first concrete NN

Performance: Standalone NN / Monitoring / Shielding					
Original NN		Monitoring (VeriPhy)		Shielding (JSC)	
Reward	Crash	Reward	Crash	Reward	Crash
<b>17.63</b> $\pm$ 0.21	<b>0</b> %	16.72 $\pm$ 0.32	<b>0</b> %	<b>17.63</b> $\pm$ 0.21	<b>0</b> %

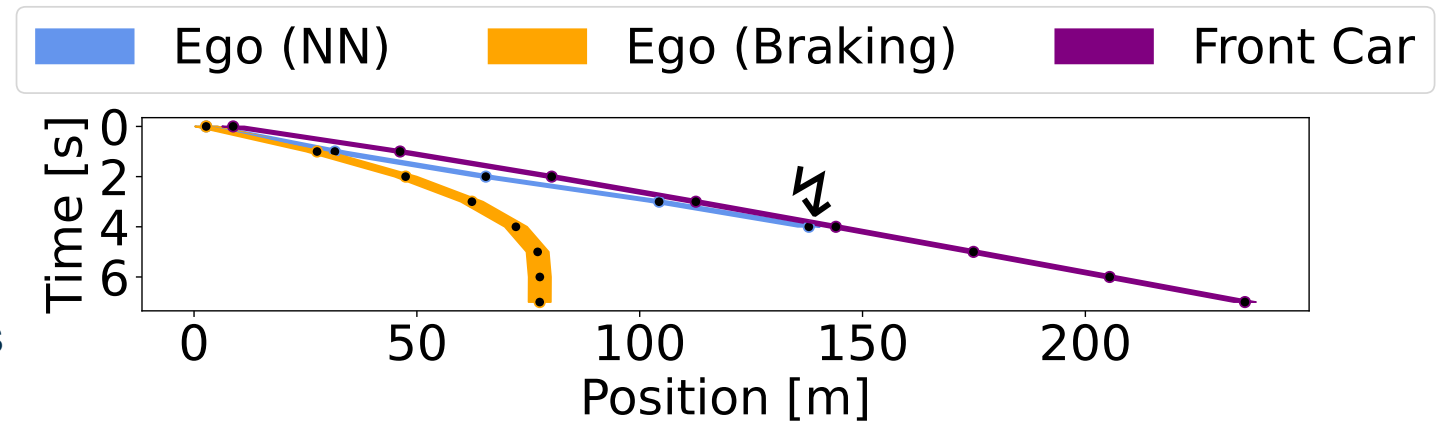
This looks good – let's verify it!

(1000 simulations)

For starters: 2 cars

- Verifier (NCubeV): 3.6 hours  
NN size: 2x256 ReLU nodes
- **14,917 counterexample regions**  
(exhaustive!)
- Sampling trajectories: 538 concrete crashes

**What went wrong?**



Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○

Evaluation and the Model2Sim Gap  
○●○○○○

Conclusion  
○

References

# Model-To-Simulation Gap (2)

## Behaviour of other cars

### Spec:

“Maximum **braking** acceleration of **front** vehicle:  $\beta_{max}$ ”

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○

Evaluation and the Model2Sim Gap  
○○●○○○

Conclusion  
○

References

# Model-To-Simulation Gap (2)

## Behaviour of other cars

### Spec:

“Maximum **braking** acceleration of **front** vehicle:  $\beta_{max}$ ”

### Simulator (highway-env):

Other cars are controlled by the *Intelligent Driver Model*

Originally used for **congestion modelling**; Cars **rarely/never brake!**

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○●○○○

Conclusion

○

References



# Model-To-Simulation Gap (2)

## Behaviour of other cars

### Spec:

“Maximum **braking** acceleration of **front** vehicle:  $\beta_{max}$ ”

### Simulator (highway-env):

Other cars are controlled by the *Intelligent Driver Model*

Originally used for **congestion modelling**; Cars **rarely/never brake!**

*We adjusted the implementation of the other cars to increase likelihood of braking.*

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○●○○○

Conclusion

○

References

# Experimental Results (2)

## A first concrete NN

Performance: Standalone NN / Monitoring / Shielding						
Env	Original NN		Monitoring (VeriPhy)		Shielding (JSC)	
	Reward	Crash	Reward	Crash	Reward	Crash
default (IDM)	<b>17.63</b> $\pm$ 0.21	<b>0</b> %	16.72 $\pm$ 0.32	<b>0</b> %	<b>17.63</b> $\pm$ 0.21	<b>0</b> %
braking	5.44 $\pm$ 1.27	99.6%	<b>16.47</b> $\pm$ 0.05	<b>0</b> %	<b>16.47</b> $\pm$ 0.05	<b>0</b> %

(1000 simulations)

# Experimental Results (2)

## A first concrete NN

Env	Performance: Standalone NN / Monitoring / Shielding					
	Original NN		Monitoring (VeriPhy)		Shielding (JSC)	
	Reward	Crash	Reward	Crash	Reward	Crash
default (IDM)	<b>17.63</b> $\pm$ 0.21	<b>0</b> %	16.72 $\pm$ 0.32	<b>0</b> %	<b>17.63</b> $\pm$ 0.21	<b>0</b> %
braking	5.44 $\pm$ 1.27	99.6%	<b>16.47</b> $\pm$ 0.05	<b>0</b> %	<b>16.47</b> $\pm$ 0.05	<b>0</b> %

Can we train a better NN?

(1000 simulations)

# Experimental Results (2)

## A first concrete NN

Env	Performance: Standalone NN / Monitoring / Shielding					
	Original NN		Monitoring (VeriPhy)		Shielding (JSC)	
	Reward	Crash	Reward	Crash	Reward	Crash
default (IDM)	<b>17.63</b> $\pm$ 0.21	<b>0</b> %	16.72 $\pm$ 0.32	<b>0</b> %	<b>17.63</b> $\pm$ 0.21	<b>0</b> %
braking	5.44 $\pm$ 1.27	99.6%	<b>16.47</b> $\pm$ 0.05	<b>0</b> %	<b>16.47</b> $\pm$ 0.05	<b>0</b> %

Can we train a better NN?

### Modifications:

- 80% of initial states: **within controllable region**
- Front Car: Initiates emergency brake with 15% likelihood
- Smaller NN for better verifiability (2 layers with 16 neurons)

**Performance for braking:** 16.08  $\pm$  0.07 reward / 0 crashes

(1000 simulations)

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○

Evaluation and the Model2Sim Gap  
○○●○○

Conclusion  
○

References

# Experimental Results (3)

## A better NN?

Verification w.r.t. **full specification** for front scenario:

- 2-5 cars in the front
- Assume  $B_{\min} = B_{\max}$

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○●○

Conclusion

○

References

# Experimental Results (3)

## A better NN?

Verification w.r.t. **full specification** for front scenario:

- 2-5 cars in the front
- Assume  $B_{\min} = B_{\max}$

### Verification:

- 1.9 hours
- 11,059 counterexample regions
- default: 4852 crashes
- braking: 8713 crashes

**Would braking have saved the car?**

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○●○

Conclusion

○

References

# Experimental Results (3)

## A better NN?

Verification w.r.t. **full specification** for front scenario:

- 2-5 cars in the front
- Assume  $B_{\min} = B_{\max}$

### Verification:

- 1.9 hours
- 11,059 counterexample regions
- default: 4852 crashes
- braking: 8713 crashes

### Would braking have saved the car?

- default: still 181 crashes
- braking: still 40 crashes

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○●○

Conclusion

○

References

# Experimental Results (3)

## A better NN?

Verification w.r.t. **full specification** for front scenario:

- 2-5 cars in the front
- Assume  $B_{\min} = B_{\max}$

### Verification:

- 1.9 hours
- 11,059 counterexample regions
- default: 4852 crashes
- braking: 8713 crashes

### Would braking have saved the car?

- default: still 181 crashes
- braking: still 40 crashes

?!?!

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○●○

Conclusion

○

References



# Model-To-Simulation Gap (3)

## Environment Model

**Spec:** Continuous evolution of environment

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○●

Conclusion

○

References

# Model-To-Simulation Gap (3)

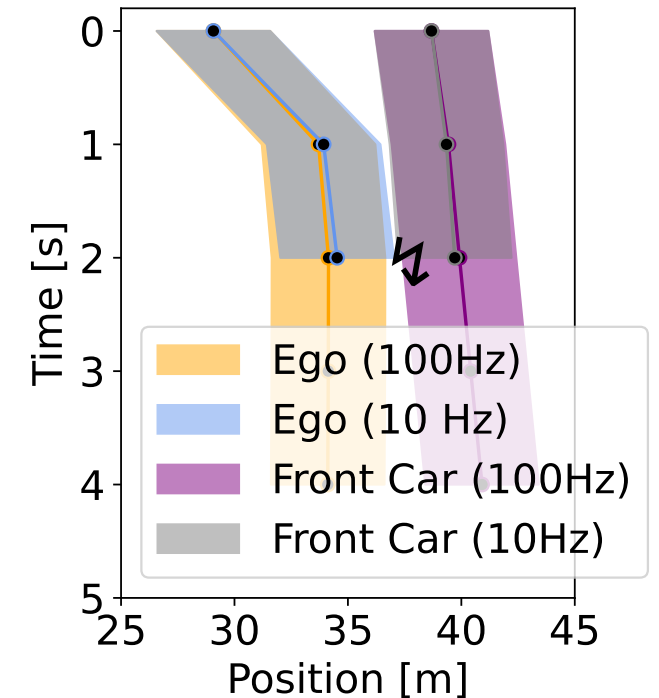
## Environment Model

**Spec:** Continuous evolution of environment

**Simulator** (highway-env): **Euler Approximations**

⇒ **Euler Crashes:**

Occurrence of crash dependent on precision of approximation



Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○○

Evaluation and the Model2Sim Gap  
○○○○○●

Conclusion  
○

References

# Model-To-Simulation Gap (3)

## Environment Model

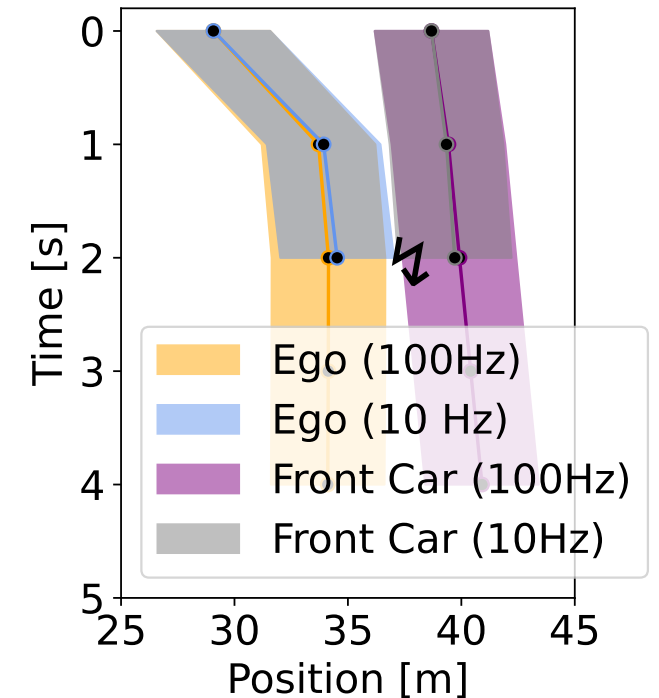
**Spec:** Continuous evolution of environment

**Simulator** (highway-env): **Euler Approximations**

⇒ **Euler Crashes:**

Occurrence of crash dependent on precision of approximation

**Additionally:** Simulator seems to initialize environment on **small** subset of admissible states.



# Model-To-Simulation Gap (3)

## Environment Model

**Spec:** Continuous evolution of environment

**Simulator** (highway-env): **Euler Approximations**

⇒ **Euler Crashes:**

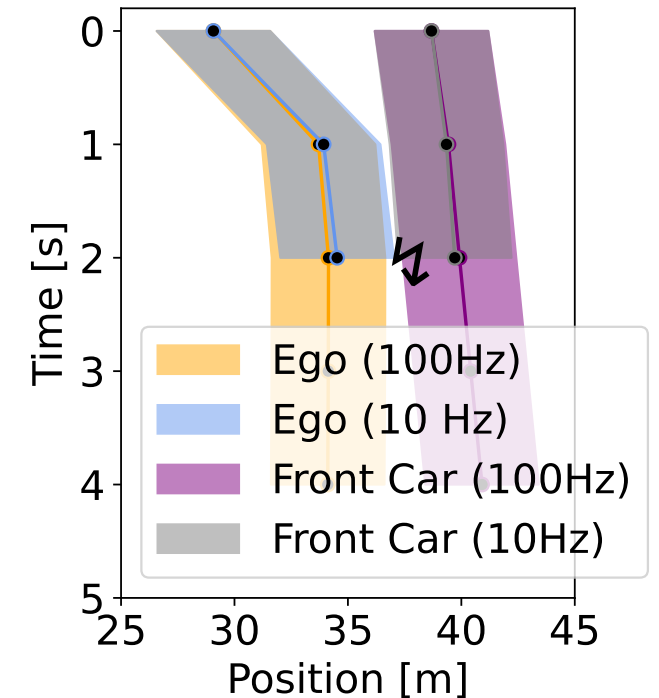
Occurrence of crash dependent on precision of approximation

**Additionally:** Simulator seems to initialize environment on **small** subset of admissible states.

## The Model-to-Simulation Gap

- Unifying assumptions across formal models & simulations is **challenging**
- Safe control requires simulators showing full breadth of possible behaviour
- As is, highway-env is no reliable basis for training safe car control NNs.

**This is a problem beyond this concrete case study!**



Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○○

Evaluation and the Model2Sim Gap  
○○○○○●

Conclusion  
○

References

# Conclusion

## Contributions

- General **dL model** for highway car control

$$\text{sys} ::= \underbrace{(\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T))}_{\text{control}} \underbrace{\text{accelCorr}; \text{dyn}}_{\text{plant}}^*$$

Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion

●

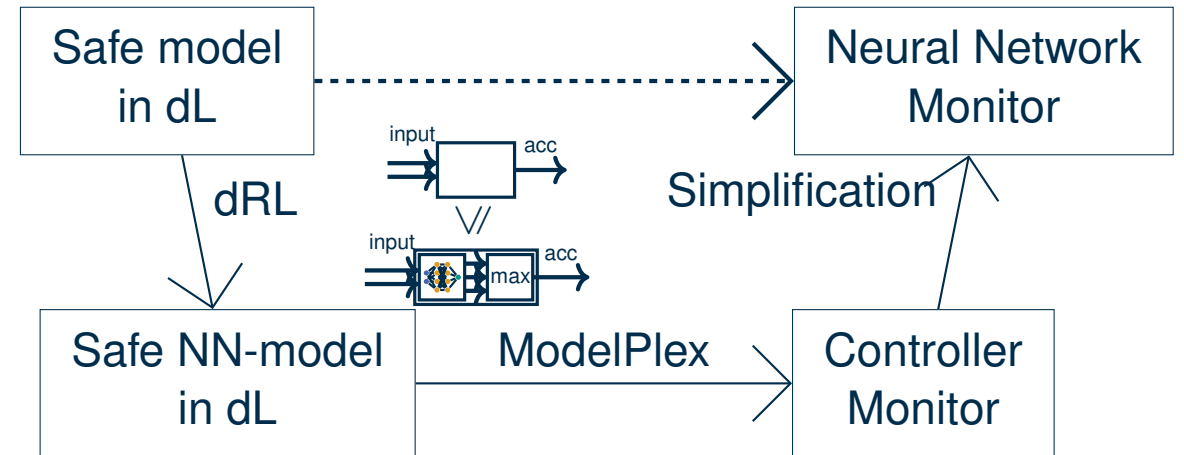
References

# Conclusion

## Contributions

- General **dL model** for highway car control
- Derivation of **real arithmetic constraints** for monitoring/shielding/verification

$$\text{sys} ::= \underbrace{(\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T))}_{\text{control}} \underbrace{\text{accelCorr; dyn}}_{\text{plant}}^*$$

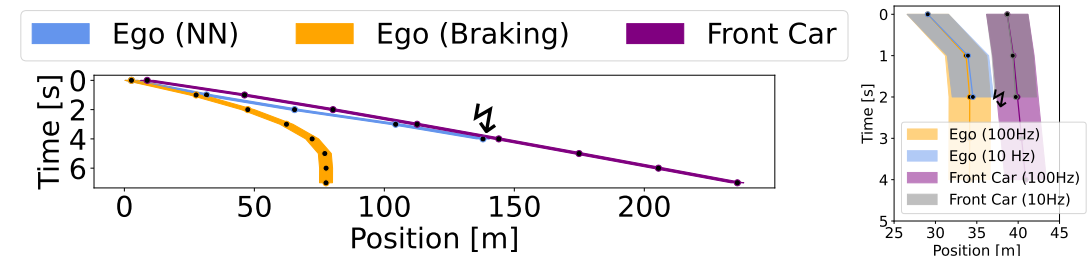
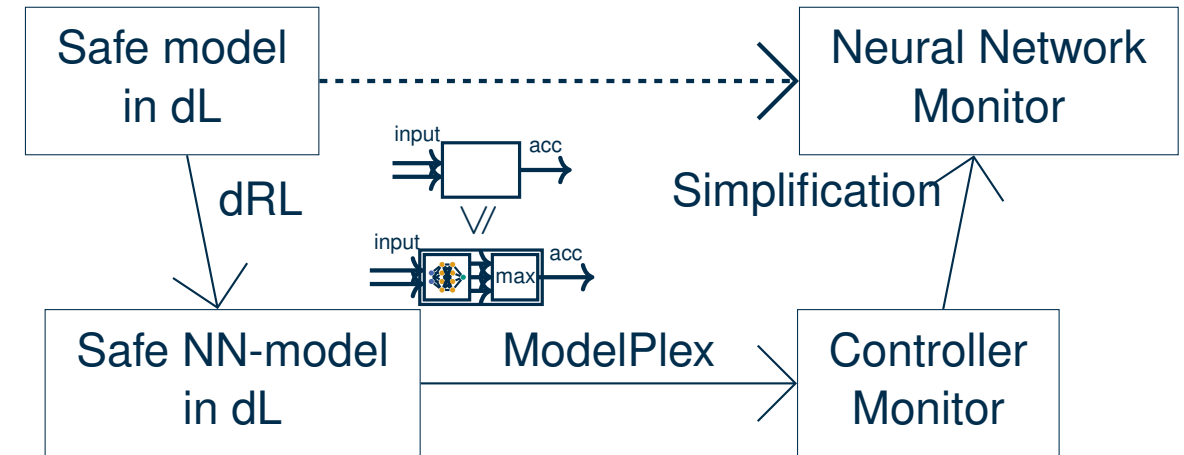


# Conclusion

## Contributions

- General **dL model** for highway car control
- Derivation of **real arithmetic constraints** for monitoring/shielding/verification
- An **empirical validation** of all three dL-based safeguarding techniques

$$\text{sys} ::= \underbrace{(\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T))}_{\text{control}} \underbrace{\text{accelCorr; dyn}}_{\text{plant}}^*$$



Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion

●

References

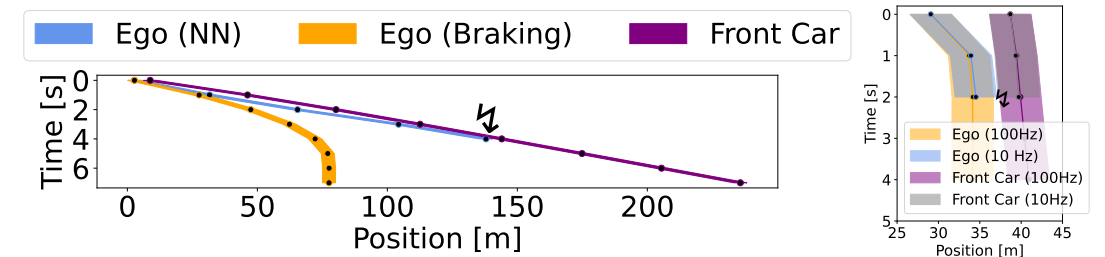
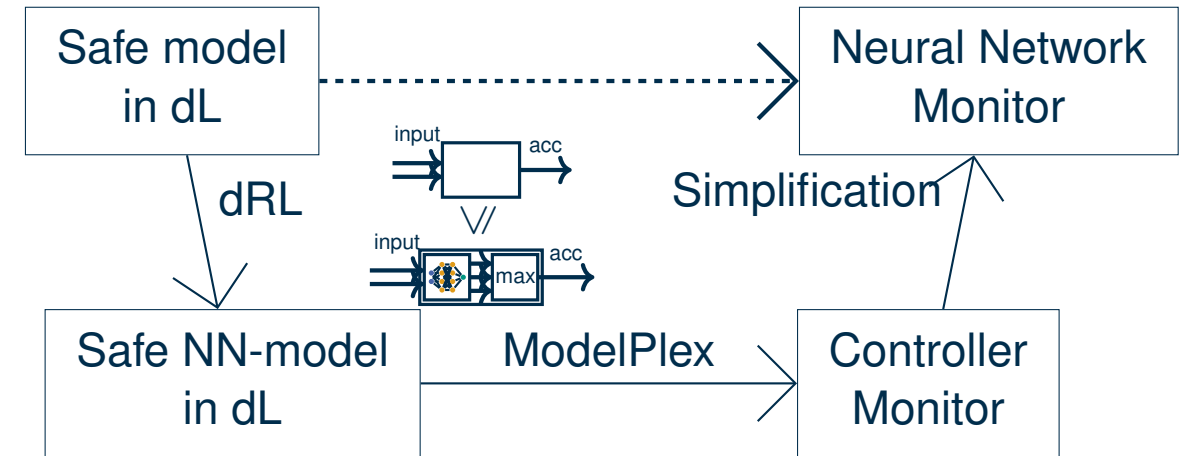
# Conclusion

## Contributions

- General **dL model** for highway car control
- Derivation of **real arithmetic constraints** for monitoring/shielding/verification
- An **empirical validation** of all three dL-based safeguarding techniques

All presented techniques are **general**!

$$\text{sys} ::= \underbrace{(\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T))}_{\text{control}} \underbrace{\text{accelCorr; dyn}}_{\text{plant}}^*$$



Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○

Evaluation and the Model2Sim Gap  
○○○○○

Conclusion  
●

References



# Conclusion

## Contributions

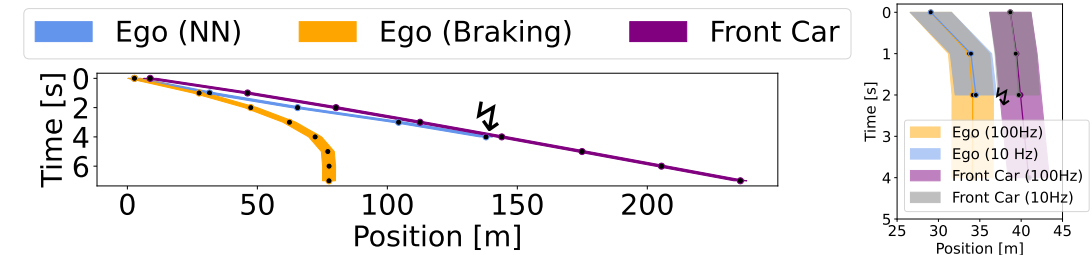
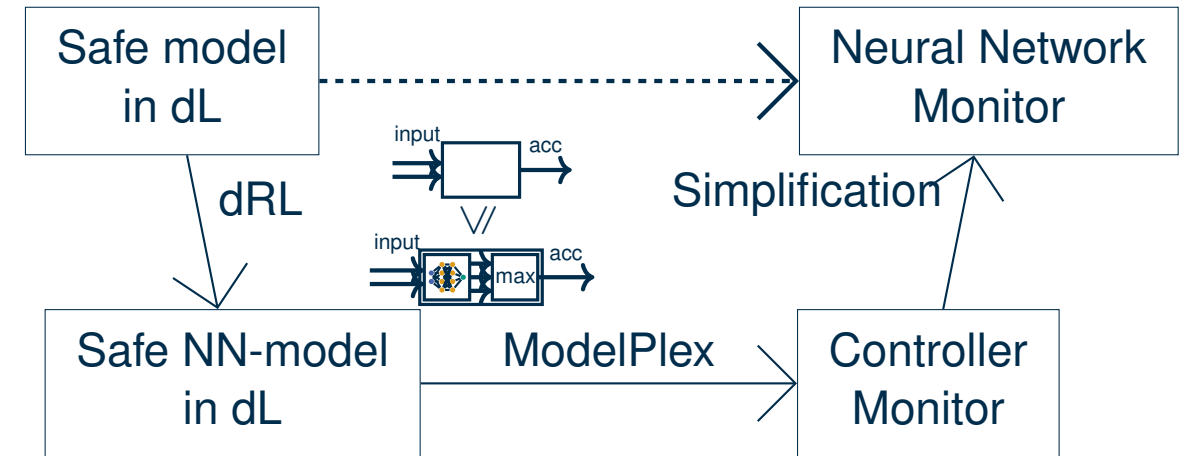
- General **dL model** for highway car control
- Derivation of **real arithmetic constraints** for monitoring/shielding/verification
- An **empirical validation** of all three dL-based safeguarding techniques

All presented techniques are **general**!

## Observations

- Consistency between different views of the system (model, simulation,...) is challenging
- **BUT:** Consistency is paramount to train provably safe ML systems

$$\text{sys} ::= \underbrace{(\text{ctrl}_o; (\text{ctrl}_e \cup ?t < t_e + T))}_{\text{control}} \underbrace{\text{accelCorr; dyn}}_{\text{plant}}^*$$



Motivation

○

Modelling with dL

○○○○○

Applications of ModelPlex

○○○○

Evaluation and the Model2Sim Gap

○○○○○○

Conclusion

●

References

# Literature I

- [1] Rose Bohrer et al. “VeriPhy: verified controller executables from verified cyber-physical system models”. In: *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018, Philadelphia, PA, USA, June 18-22, 2018*. Ed. by Jeffrey S. Foster and Dan Grossman. ACM, 2018, pp. 617–630. DOI: 10.1145/3192366.3192406.
- [2] Nathan Fulton and André Platzer. “Safe Reinforcement Learning via Formal Methods: Toward Safe Control Through Proof and Learning”. In: *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*. Ed. by Sheila A. McIlraith and Kilian Q. Weinberger. AAAI Press, 2018, pp. 6485–6492. DOI: 10.1609/aaai.v32i1.12107.
- [3] Stefan Mitsch and André Platzer. “ModelPlex: verified runtime validation of verified cyber-physical system models”. In: *Formal Methods Syst. Des.* 49.1-2 (2016), pp. 33–74. DOI: 10.1007/s10703-016-0241-z.
- [4] Samuel Teuber, Stefan Mitsch, and André Platzer. “Provably Safe Neural Network Controllers via Differential Dynamic Logic”. In: *Advances in Neural Information Processing Systems*. Ed. by A. Globerson et al. Curran Associates, Inc., 2024. URL: <https://doi.org/10.48550/arXiv.2402.10998>.

Motivation  
○

Modelling with dL  
○○○○○

Applications of ModelPlex  
○○○○

Evaluation and the Model2Sim Gap  
○○○○○○

Conclusion  
○

References

# Safety formula when behind

$$\begin{aligned} & x_e + L \leq x_o \wedge (a_e \leq B_{\min} \wedge \text{pos}_e(B_{\min}) + L < \text{pos}_o \\ & \vee B_{\min} \leq a_e \wedge v_e + a_e T < 0 \wedge \text{pos}_e(a_e) + L < \text{pos}_o \\ & \vee B_{\min} \leq a_e \wedge v_e + a_e T \geq 0 \wedge \text{pos}_e(B_{\min}) + \text{corrDist} + L < \text{pos}_o) \end{aligned}$$

$$\begin{aligned} \text{pos}_e(a_e) &= x_e - \frac{v_e^2}{2a_e} \\ \text{pos}_o &= x_o - \frac{v_o^2}{2B_{\max}} \\ \text{corrDist} &= \left(\frac{-a_e}{B_{\min}} + 1\right) \left(\frac{a_e}{2} T^2 + T v_e\right) \end{aligned}$$

# Full dL Model

$\text{ctrl}_o$	$a_o := *; ?(B_{\max} \leq a_o \leq A_{\max});$
$\text{ctrl}_e$	$a_e := *; ?(B_{\max} \leq a_e \leq A_{\max}); t_e := t;$ $\text{if}(\neg(\text{safeBack} \vee \text{safeFront}))$ $\text{if}(x_e \leq x_o)$ $a_e := *; ?(B_{\max} \leq a_e \leq B_{\min});$ $\text{else}$ $a_e := *; ?(A_{\min} \leq a_e \leq A_{\max});$
$\text{accelCorr}$	$\text{if} (v_o = 0 \wedge a_o < 0) \vee (v_o = V \wedge a_o > 0) a_o := 0$ $\text{if} (v_e = 0 \wedge a_e < 0) \vee (v_e = V \wedge a_e > 0) a_e := 0$
$\text{dyn}$	$x'_e = v_e, v'_e = a_e, x'_o = v_o, v'_o = a_o, t' = 1$ $\& t \leq t_e + T \wedge 0 \leq v_e \leq V \wedge 0 \leq v_o \leq V$