

Counter-Example Guided Imitation Learning of Feedback Controllers From Temporal Logic Specification

Thao Dang¹, **Alexandre Donzé**², Inzemamul Haque³, Nikolaos
Kekatos⁴, Indranil Saha³

¹ Univ. Grenoble Alpes, CNRS, Grenoble INP, VERIMAG, Grenoble, France,

² Decyphir SAS, Moirans, France

³ Department of Computer Science and Engineering, IIT Kanpur, India,

⁴ Aristotle University of Thessaloniki, Thessaloniki, Greece

March 7th, 2024

Counter-Example Guided Imitation Learning of Feedback Controllers From Temporal Logic Specification

Thao Dang¹, **Alexandre Donzé**², Inzemamul Haque³, Nikolaos
Kekatos⁴, Indranil Saha³

¹ Univ. Grenoble Alpes, CNRS, Grenoble INP, VERIMAG, Grenoble, France,

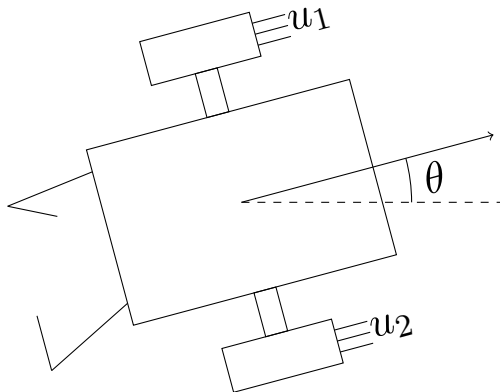
² Decyphir SAS, Moirans, France

³ Department of Computer Science and Engineering, IIT Kanpur, India,

⁴ Aristotle University of Thessaloniki, Thessaloniki, Greece

March 7th, 2024

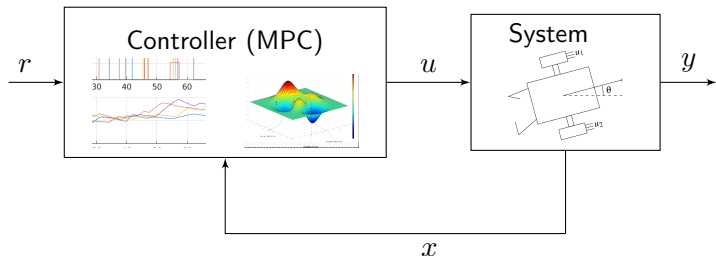
Control Problem



Flying Robot with two thrusters

Goal Go from x_0 to origin

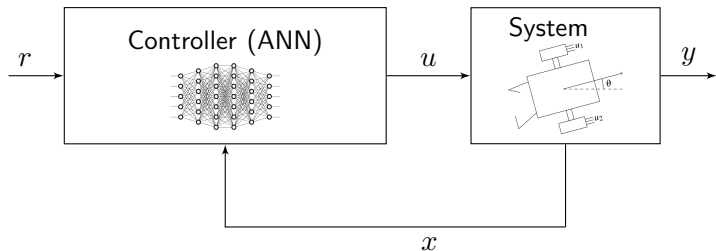
State Feedback Control



Issues

- ▶ MPC solves a NL optimization problem at each step
- ▶ Costly for online use
- ▶ Can we replace it with a NN (fast inference)?

Neural Network Control System

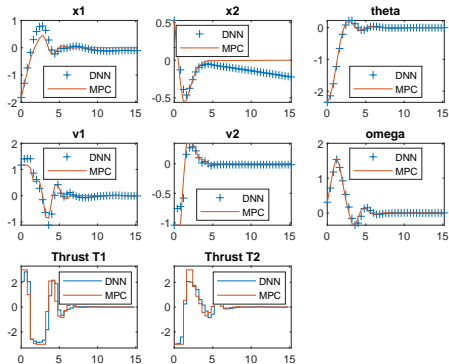


Supervised learning

- ▶ Simulate with MPC
- ▶ collect set of input/output for the NN

$$\mathcal{D} = \{(x, u), u = MPC(x)\}$$

Behavior Cloning Validation



Not So Good

Visible drift in x_2 , need more data ?

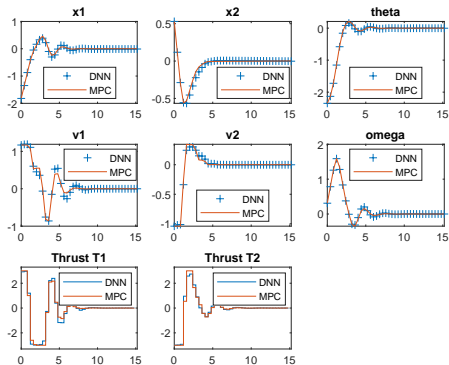
Data Aggregation Approach (Dagger)

```
 $\mathcal{N}_0 \leftarrow \emptyset, \mathcal{D}_0 \leftarrow \emptyset, k \leftarrow 1$   
repeat  
|    $\mathcal{D}_k \leftarrow \text{getNewData}(\mathcal{N}_{k-1}, \mathcal{D}_{k-1})$   
|    $\mathcal{N}_k \leftarrow \text{Train}(\mathcal{D}_k)$   
|    $k \leftarrow k + 1$   
until  $k > k_{max}$   
return “Best”  $\mathcal{N}_i$  for  $i \in 1, \dots, k_{max}$ 
```

Main Idea

New data is obtained by simulating with MPC from states visited by \mathcal{N}_{k-1}

Dagger Results



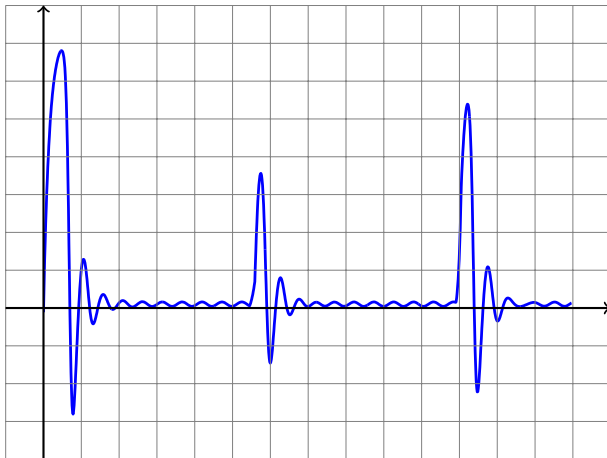
Looks better but...

- ▶ Can we trust one simulation to evaluate \mathcal{N}_i ?
- ▶ Did we need 140 iterations ?

Outline

- 1 Imitation Learning of Feedback Control
- 2 PSTL-based Evaluation of Controllers**
- 3 Counter-Example Guided Data Collection
- 4 Discussion

Temporal Logics Requirements (STL)



Temporal Logics Requirements (STL)

The signal is never above 3.5

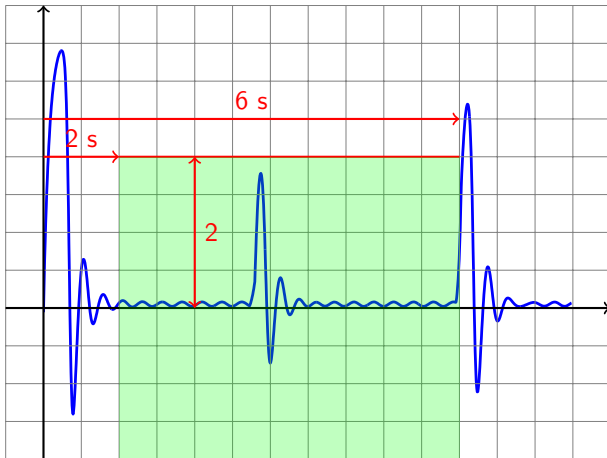
$$\varphi := \text{alw } (x[t] < 3.5)$$



Temporal Logics Requirements (STL)

Between 2s and 6s the signal is between -2 and 2

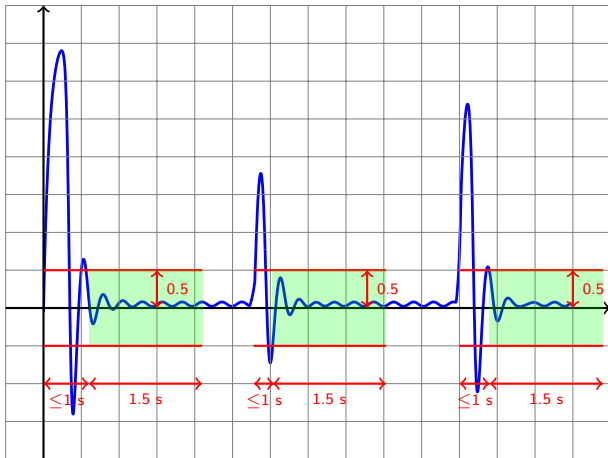
$$\varphi := \text{alw}_{[2,6]} (|x[t]| < 2)$$



Temporal Logics Requirements (STL)

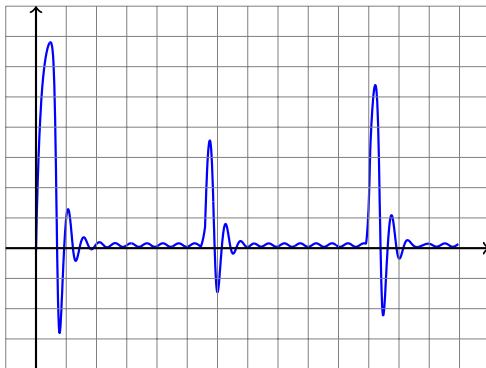
Always $|x| > 0.5 \Rightarrow$ after 1 s, $|x|$ settles under 0.5 for 1.5 s

$$\varphi := \text{alw}(x[t] > .5 \rightarrow \text{ev}_{[0,.6]} (\text{alw}_{[0,1.5]} x[t] < 0.5))$$



Parametric-STL Formulas

STL formula where numeric constants are left unspecified.

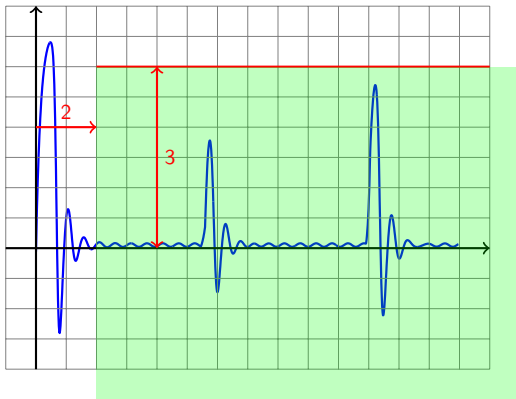


Parametric-STL Formulas

STL formula where numeric constants are left unspecified.

“After 2s, the signal is never above 3”

$$\varphi := \text{ev}_{[2,\infty]} (x[t] < 3)$$

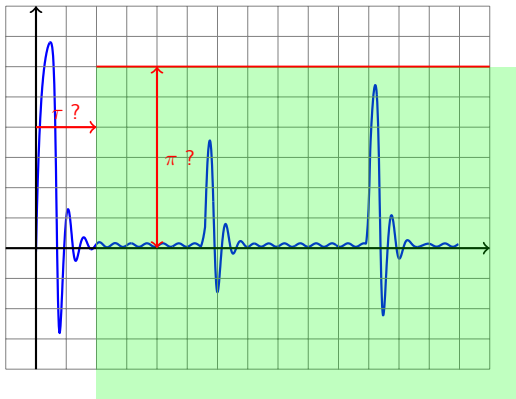


Parametric-STL Formulas

STL formula where numeric constants are left unspecified.

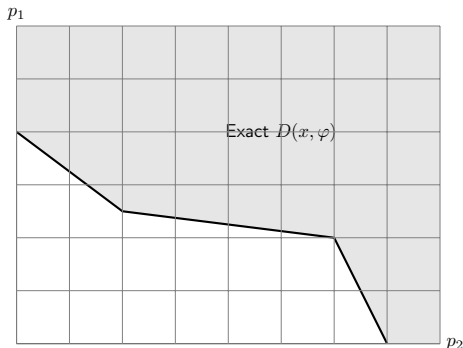
“After τ s, the signal is never above π ”

$$\varphi := \text{ev}_{[\tau, \infty]} (x[t] < \pi)$$



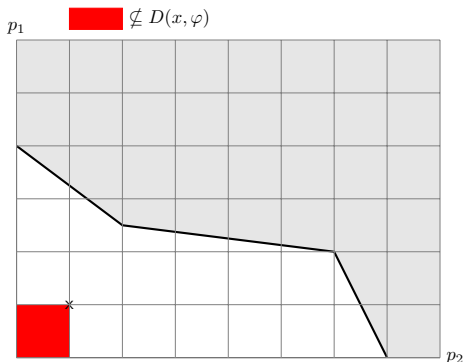
PSTL Validity Domains

- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**



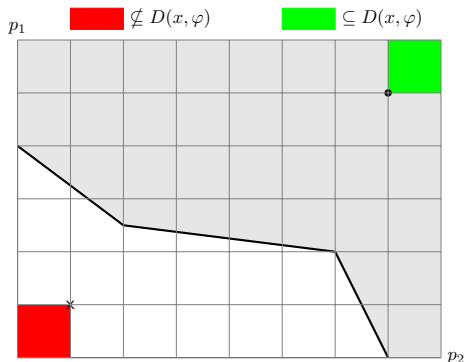
PSTL Validity Domains

- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**



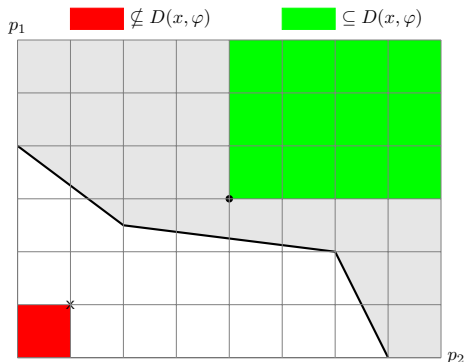
PSTL Validity Domains

- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**



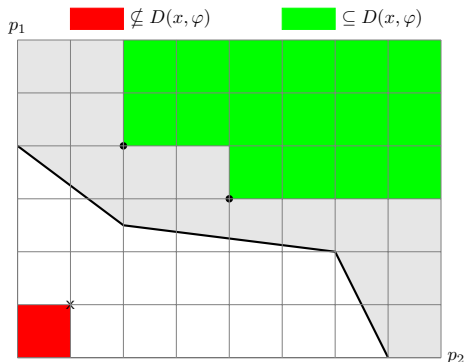
PSTL Validity Domains

- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**



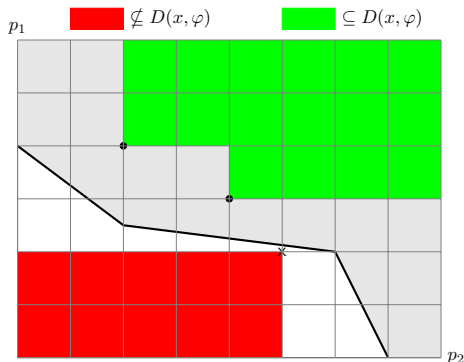
PSTL Validity Domains

- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**



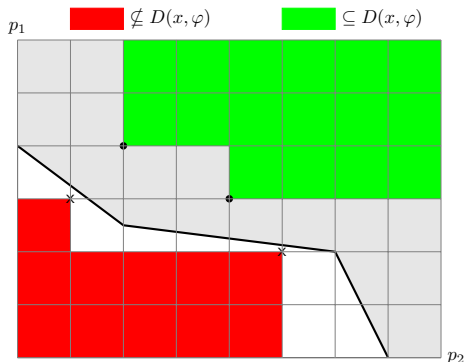
PSTL Validity Domains

- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**



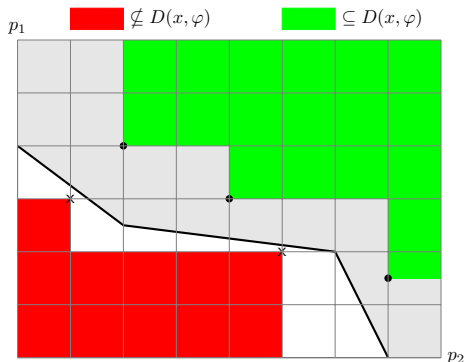
PSTL Validity Domains

- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**

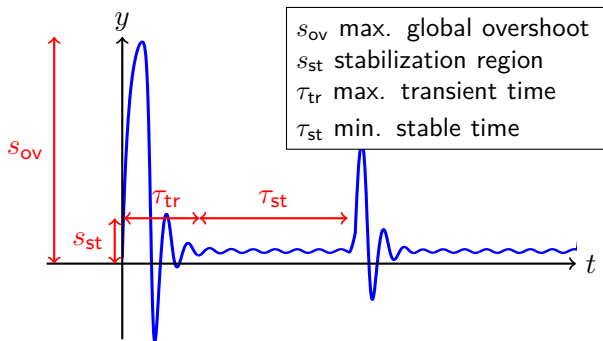


PSTL Validity Domains

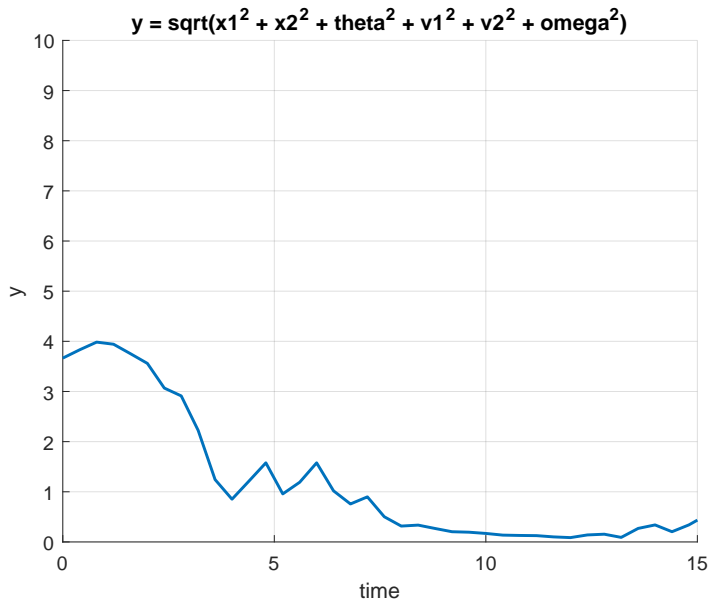
- ▶ The validity domain D of φ and x is the set of valuations v s.t. $x \models \varphi(v)$
- ▶ In case of monotonicity, ∂D has the structure of a **Pareto front**



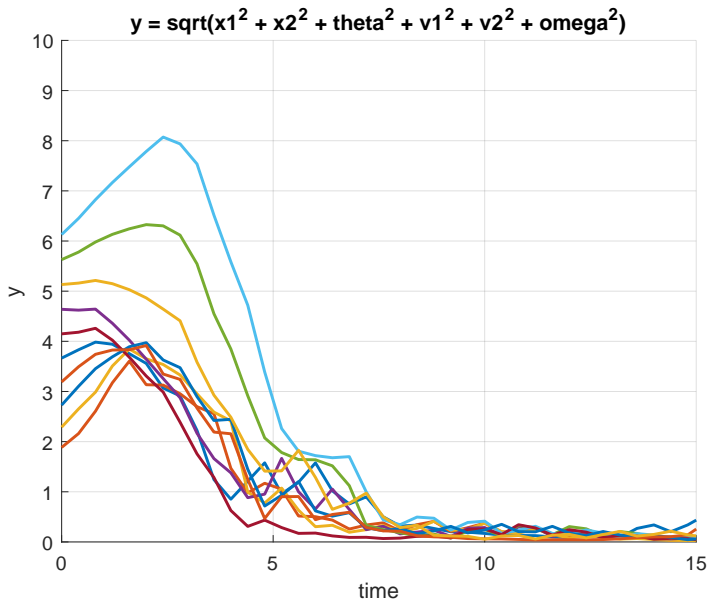
PSTL for Stabilization



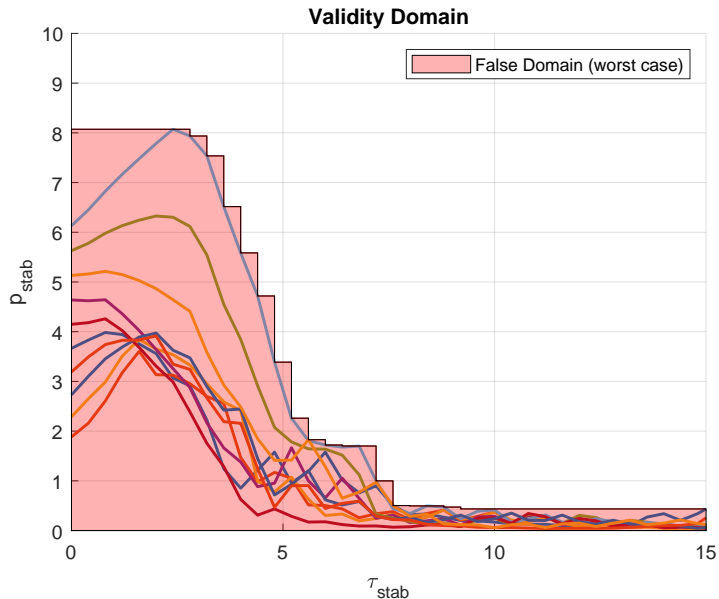
Flying Paretos



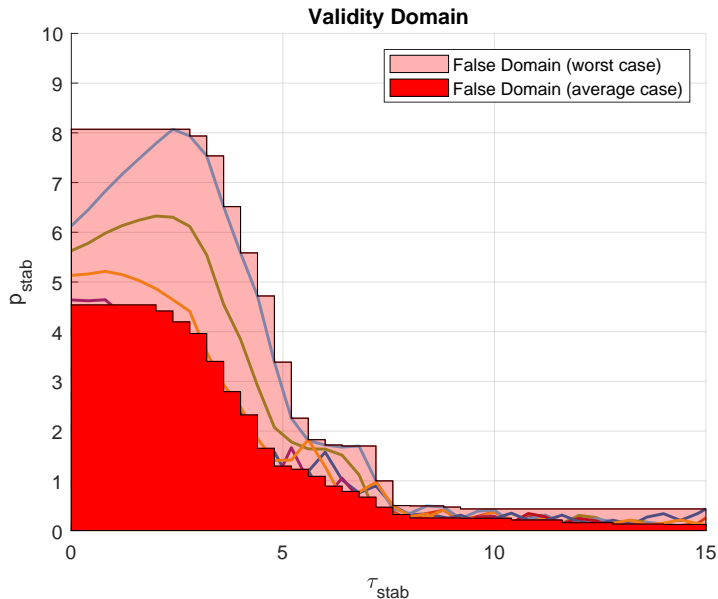
Flying Paretos



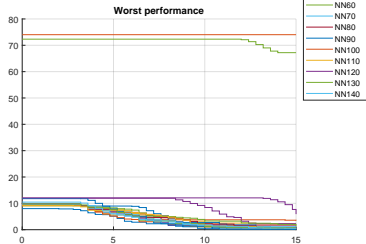
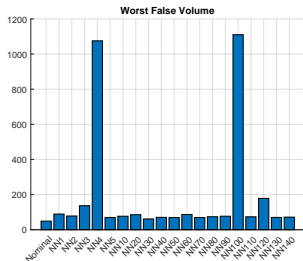
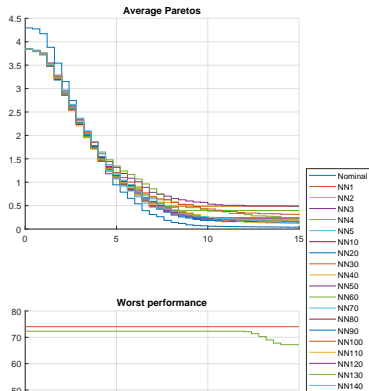
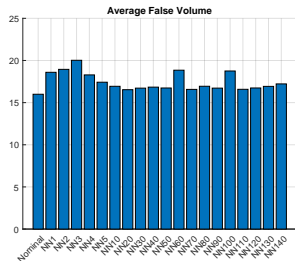
Flying Paretos



Flying Paretos



Dagger "Progress"



Outline

- 1 Imitation Learning of Feedback Control
- 2 PSTL-based Evaluation of Controllers
- 3 Counter-Example Guided Data Collection**
- 4 Discussion

How can we improve the data collection process?

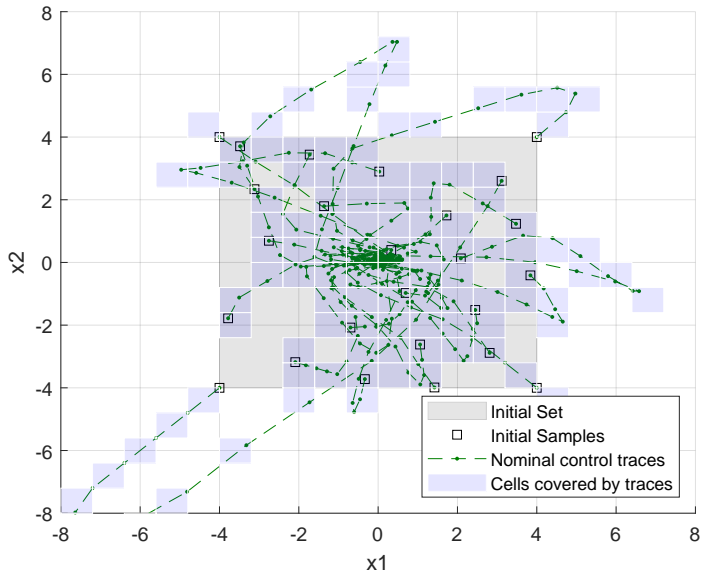
It's the data stupid

- ▶ Idea 1: Better control over coverage (sample efficiency)
- ▶ Idea 2: Generate data from counter-examples

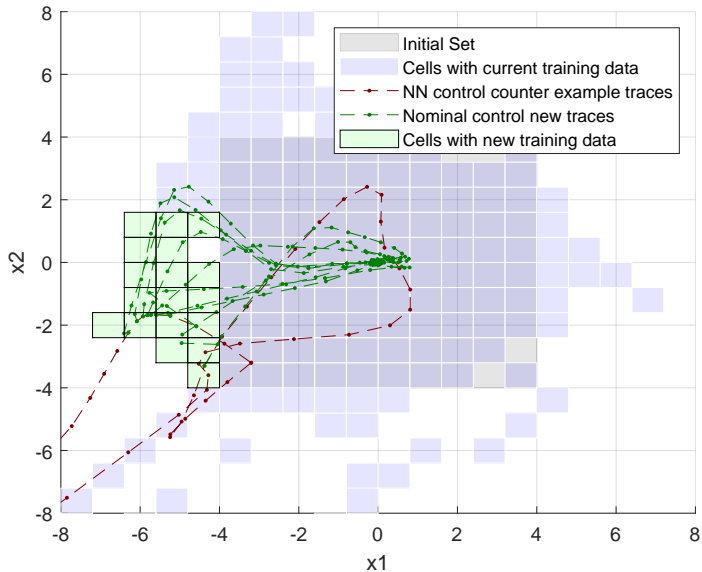
New Data Acquisition Procedure

```
1: procedure GETNEWDATA( $\mathcal{N}, \mathcal{D}$ )
2:   if  $\mathcal{D}$  is  $\emptyset$  then                                ▷ Sample traces from initial set
3:      $InitSamples \leftarrow \text{getInitSamples}()$ 
4:      $InitTraces \leftarrow \text{simNominal}(InitSamples)$ 
5:      $\mathcal{D}_{\text{new}} \leftarrow \text{gridFilter}(InitTraces)$ 
6:      $Status \leftarrow$  "New data available for training."
7:   else                                                  ▷ Search for counter-examples
8:      $CexTraces \leftarrow \text{falsify}(\mathcal{N})$ 
9:     if  $CexTraces \neq \emptyset$  then
10:      |  $(\mathcal{D}_{\text{new}}, Status) \leftarrow \text{fixAndMerge}(\mathcal{D}, CexTraces)$ 
11:    else                                                ▷ Success.
12:      |  $\mathcal{D}_{\text{new}} \leftarrow \mathcal{D}$ 
13:      |  $Status \leftarrow$  "No counter-example found."
```

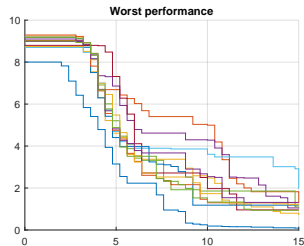
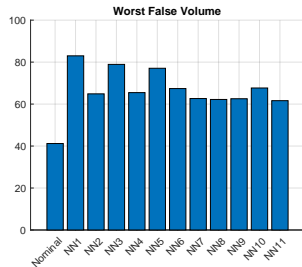
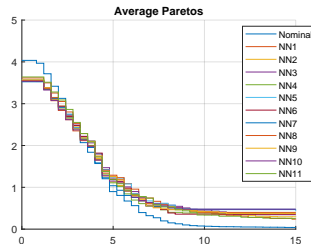
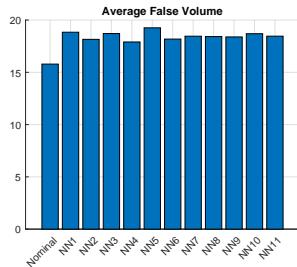
Grid Coverage - Initial Sampling



Counter-Example Data Extraction



Results



Outline

- 1 Imitation Learning of Feedback Control
- 2 PSTL-based Evaluation of Controllers
- 3 Counter-Example Guided Data Collection
- 4 Discussion

Discussion

More sensible stopping criterion and performance evaluation

- ▶ PSTL provides some control over what criterion to train for in priority
- ▶ Falsification stopping criterion gives some confidence over result

Discussion

More sensible stopping criterion and performance evaluation

- ▶ PSTL provides some control over what criterion to train for in priority
- ▶ Falsification stopping criterion gives some confidence over result

Approach “Probably” more sample efficient than vanilla DAgger

- ▶ “Converges” with less iteration and samples
- ▶ Evidence that performance is not affected by filtering (like quantization)

Discussion

More sensible stopping criterion and performance evaluation

- ▶ PSTL provides some control over what criterion to train for in priority
- ▶ Falsification stopping criterion gives some confidence over result

Approach “Probably” more sample efficient than vanilla DAgger

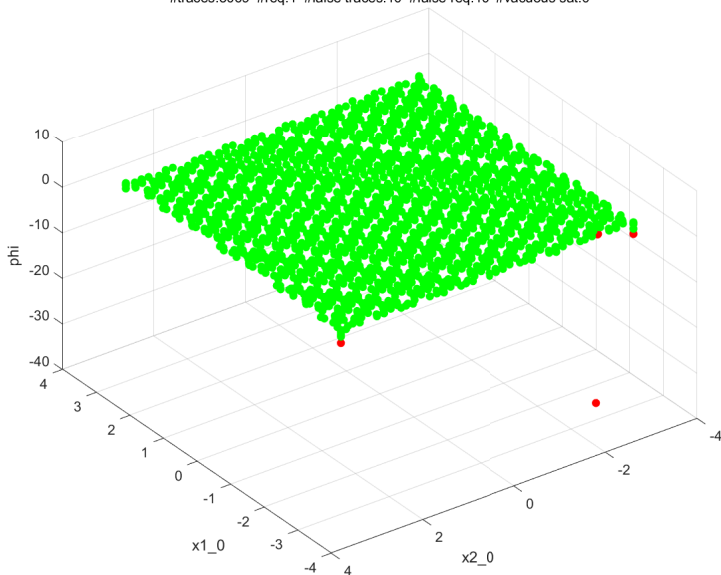
- ▶ “Converges” with less iteration and samples
- ▶ Evidence that performance is not affected by filtering (like quantization)

But...

- ▶ Reproducibility is an issue. High variability even with same training data
- ▶ Performance still comparable to original `final_policy` provided in Mathworks examples, though...

Mathworks' Policy Falsification

#traces:3069 #req:1 #false traces:10 #false req:10 #vacuous sat:0



Personal Conclusion

- ▶ Flying robot actually very unstable,
- ▶ I now doubt it can be robustly solved with pure NN state feedback controller
- ▶ Our method seem to work well with more well-behaved problems,
- ▶ More experiments in the making