# Neural State Classification for Hybrid Systems

## Nicola Paoletti
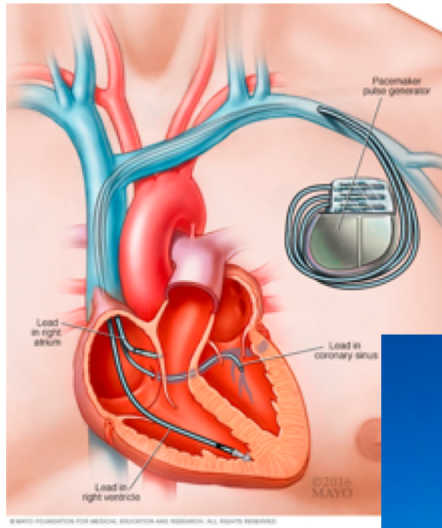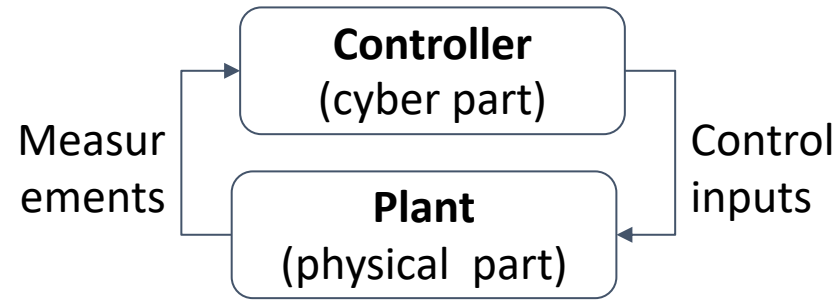
Royal Holloway, University of London, UK

JWW: D Phan, T Zhang, SA Smolka, SD Stoller (Stony Brook University) and R Grosu (TU Wien)

ATVA 2018 – University of Southern California, LA, 8 Oct 2018
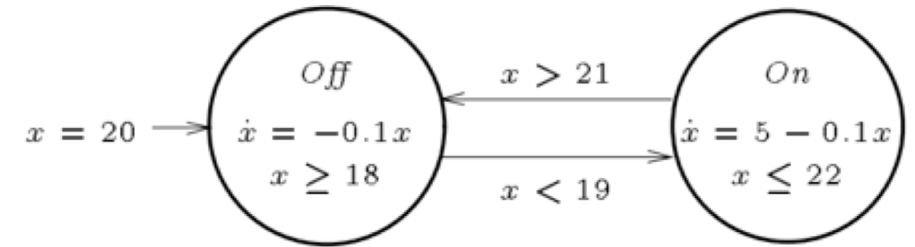
# Hybrid system verification

Hybrid systems are ubiquitous and found in many safety-critical applications



Cyber-physical system

Controller (cyber part)

Plant (physical part)

Measurements

Control inputs

# Hybrid system verification



Thermostat from *Henzinger, The Theory of Hybrid Automata*

- Hybrid automata (HA) are a common formal model for hybrid systems

- HA verification problem usually formulated as reachability
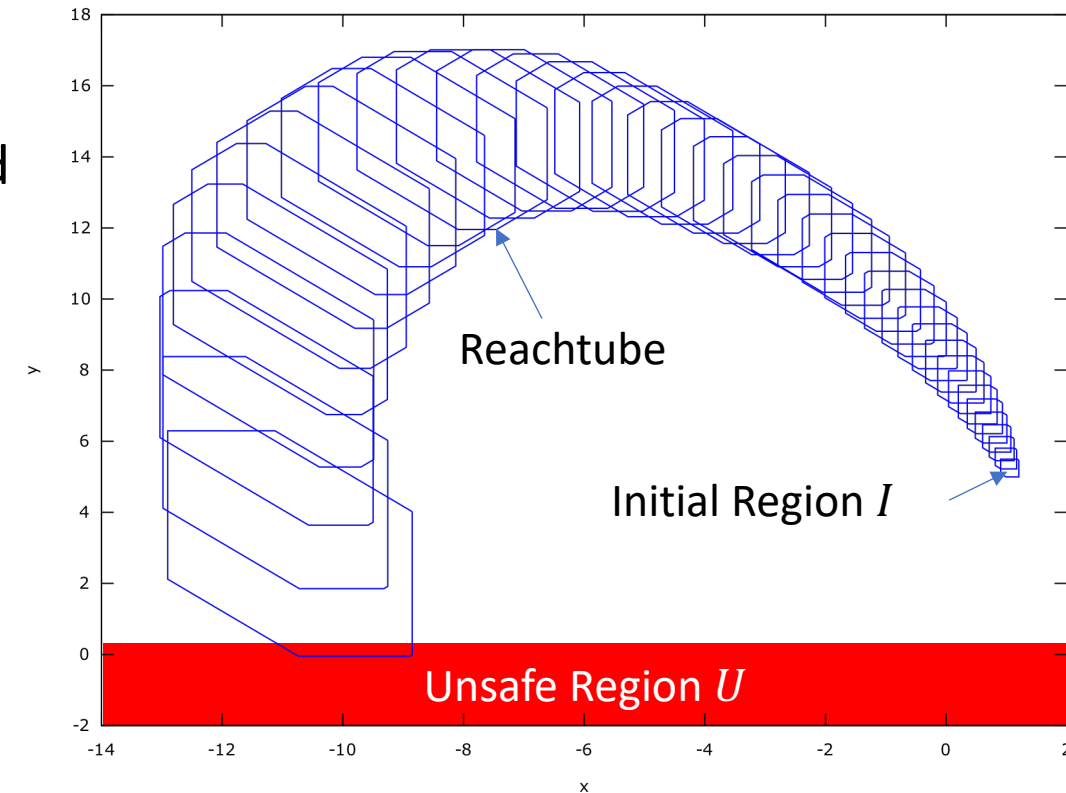
*(Time-bounded)* reachability:
can an HA $\mathcal{M}$, starting in an initial region $I$, reach a state $u \in U$ (within time $T$)?

Both bounded and unbounded versions are undecidable

[Henzinger et al, *JCSS 57 1* (1998); Brihaye et al, *ICALP* (2011)]

# Reachability checkers for HAs

- Over-approximate the set of states reachable from the initial region

  - Given initial region $I$ of an HA $\mathcal{M}$ and a time bound $T$, compute $ReachTube(\mathcal{M}, I, T)$

  - Check if $ReachTube(\mathcal{M}, I, T)$ intersects the unsafe region $U$

    - No: 100% safe

    - Yes: *maybe* unsafe, s.t. false positives

- Tools: HyCreate, Flow*, SpaceEx, iSAT, dReal, etc.

- HA reachability is computationally expensive



Reachtube

Initial Region $I$

Unsafe Region $U$

# Motivation - Online model checking (OMC)

- OMC – *predicting at runtime future violations from current state* – is as important as offline model verification for HSs and CPSs
  - switch to fail-safe operation mode when failure is imminent
    (e.g. Simplex architecture of [Sha, *IEEE Software* (2001)])

- OMC focus is on reachability from a single state, and not from a (large) region

- OMC runs the the analysis periodically → short time horizons
  - Avoids blow-up of reach-set over-approximation

- Runtime settings are less predictable
  - system might differ from model, noisy observations
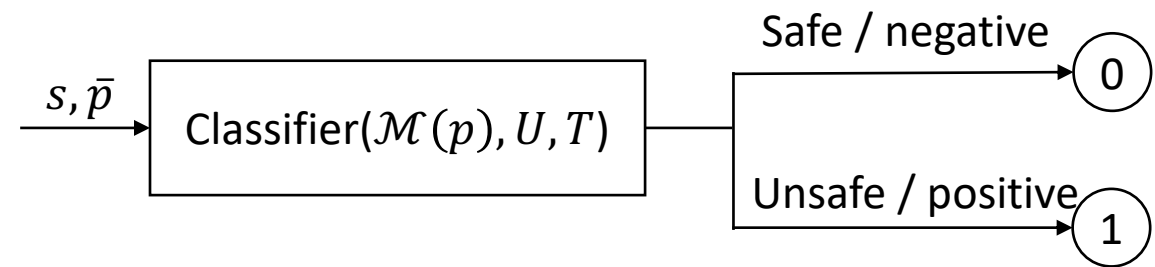
# Motivation - Online model checking (OMC)

- OMC focus is on reachability from a single state, and not from a (large) region
- OMC runs the the analysis periodically → short time horizons
- Runtime settings are less predictable

*Does OMC need fully-fledged reachability checking?*

- We rather need methods that can work under real-time constraints
  - Reachability checking is too expensive for online analysis
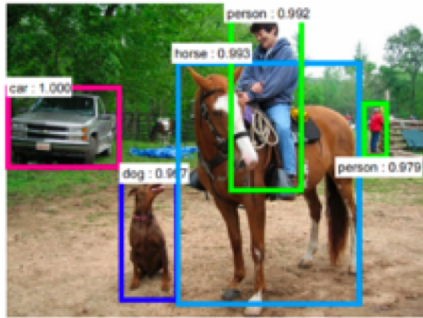
# State Classification Problem (SCP)

- We want a function that, given HA $\mathcal{M}$ with state space $S$, set of unsafe states $U$, and time bound $T$, classifies every state $s \in S$ as either *positive* or *negative*

    - $s$ is *positive* if $\mathcal{M}$, starting in $s$, can reach a state in $U$ within time $T$;

    - *negative* o/w

$s, \bar{p}$ → $\text{Classifier}(\mathcal{M}(p), U, T)$ → Safe / negative → ⓪
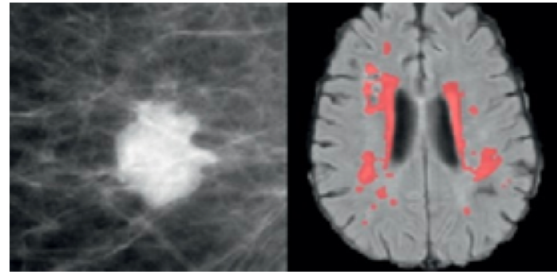
Unsafe / positive → ①

- We call such a function a *state classifier,* a solution to the SCP

- $\mathcal{M}$ can be parameterized by a set of parameters $p$

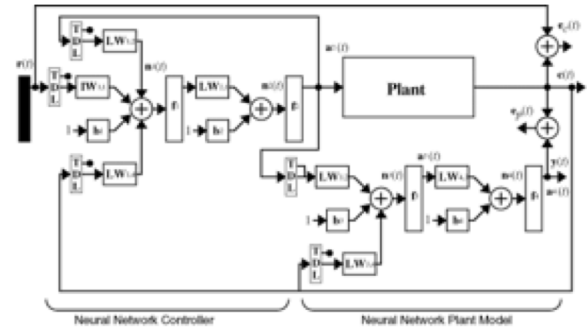# Neural networks (NNs) as state classifiers

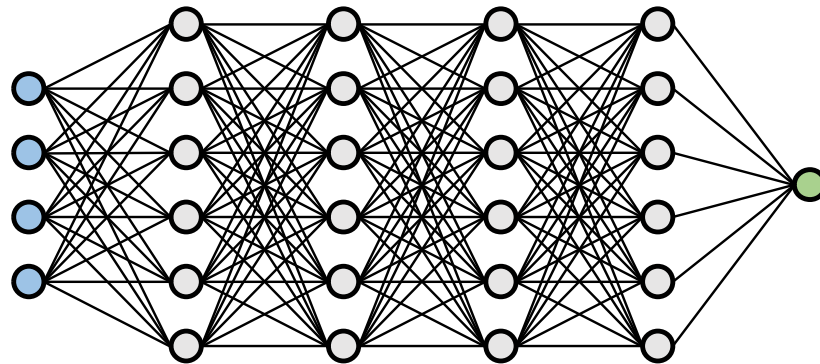(Deep) NNs are extremely successful at complex classification and regression tasks



Object detection



Classification of tumor and diseases from medical images



System identification and control



**Verification**

# Neural networks (NNs) as state classifiers

- *Can we train a NN to learn a HA reachability function, i.e., solve the SCP?*

- In principle, **YES:** NNs are universal approximators [Hornik et al, *Neural networks 2(5)* (1989)]

- In practice, good accuracy but prediction errors can't be avoided

- Trained NN state classifier runs in **constant time** ->  suitable for online model checking

Two kinds of errors in **neural state classification**:

- **False positives:** a negative state is predicted to be positive (conservative decision)

- **False negatives:**  a positive state is predicted to be negative (can compromise system's safety!)

# Neural State Classification (NSC)

# Oracles

Positive          Negative



Oracle

- Simulator (deterministic)
- Reachability checker (dReal [Gao et al, *CADE* (2013)])
- Backwards simulator

Unsafe

Unsafe

# Sampling methods



**Uniform Sampling**
- all states equally important

**Balanced Sampling**
- balanced number of pos. and neg. samples
- suitable when unsafe set U is small
- based on **backwards HA simulation**

**Dynamics-Aware Sampling**
- reflects the likelihood of visiting a state from the initial region
- based on estimating state distribution from random HA runs

# Backwards simulator

- For generating arbitrarily many positive samples for a balanced dataset

- Given an unsafe state $u \in U$, simulate $\overleftarrow{\mathcal{M}}$, the *reverse HA* of $\mathcal{M}$, for up to time $T$

- Every state in the reverse trajectory is positive

- We provide a constructive definition of reverse HA and prove its correctness (more general than [Henzinger et al, *STOC* (1995)] for rectangular automata)



Initial state of the reverse trajectory

# Statistical guarantees via hypothesis testing

- We provide guarantees on classifier's performance on unseen (test) states using the *sequential probability ratio test* (SPRT):

  - Accuracy (probability of correct prediction): $P_A \geq \theta_A$

  - FN rate (probability that prediction is an FN): $P_{FN} \leq \theta_{FN}$

  - Subject to user-defined strength of test (prob. of type-I and type-II errors)

- *Sequential* means that we only need the number of test samples necessary for SPRT to make a decision

- Idea borrowed from statistical model checking [Younes et al, *STTT 8.3* (2006)]

  - Where SPRT is for verifying $P(M \vDash \phi) \sim \theta$ for a probabilistic system

# Reducing FN rate via falsification

- Make the classifier more conservative (reduce FN) through re-training with new FN samples
  - **Dual of CEGAR** [Clarke et al, CAV (2000)]: CEGAR refines an overapproximation using counterexamples (FPs)

- FNs found via a falsifier / adversarial sampling, an algorithm that finds states maximizing the discrepancy between predictions and true labels

- Under assumptions on falsifier and classifier, the *algorithm converges to an empty set of FNs with high probability* (proof based on bounds on generalization error of ML models [Vapnik, *The nature of statistical learning theory* (2013)])

**Input:**   classifier (NN) $F$,
      training samples $D$
**Output:**  "conservative" classifier $F$
**do**
- $\widehat{FN} \leftarrow$ subset of the true FN set of $F$
  /*found via falsifier (genetic alg)*/
- $D \leftarrow D \cup \widehat{FN}$
- $F \leftarrow$ **train**$(D)$
**while** $\widehat{FN} \neq \emptyset$ **or** *max_iter*

Iterative falsification / re-training algorithm

# Experimental design

## Hybrid system benchmark:

- Spiking neuron
- Inverted pendulum
- Quadcopter dynamics
- Cruise control
- Powertrain
- Helicopter

## State classifier models:

- Feed-forward deep NNs (3 hidden layers, 10 neurons each, sigmoid and ReLU)
- Feed-forward shallow NNs (1 hidden layer, 20 neurons, sigmoid)
- Support Vector Machines (SVMs)
- Binary Decision Trees (BDTs)
- Nearest neighbor (returns label of closest training sample)

# Accuracy and FNs

| | Neuron Acc | Neuron FN | Pendulum Acc | Pendulum FN | Quadcopter Acc | Quadcopter FN | Cruise Acc | Cruise FN | Powertrain Acc | Powertrain FN | Helicopter Acc | Helicopter FN | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DNN-S | **99.81** | **0.1** | **99.98** | **0** | 99.83 | 0.1 | 99.95 | 0.01 | **96.68** | 1.28 | **98.49** | **0.84** | Uniform |
| DNN-R | 99.52 | 0.29 | 99.93 | 0.04 | **99.89** | **0.06** | **99.98** | **0** | 96.21 | **1.08** | 98 | 0.96 | |
| SNN | 99.17 | 0.43 | 99.81 | **0** | 99.85 | 0.08 | 99.84 | 0.15 | 96.02 | 1.37 | 97.69 | 1.25 | |
| SVM | 98.73 | 0.75 | 99.84 | **0** | 97.33 | 0.69 | 99.88 | 0.1 | 92.26 | 3.48 | 95.58 | 2.42 | |
| BDT | 99.3 | 0.37 | 99.6 | 0.17 | 99.52 | 0.2 | 99.84 | 0.08 | 95.59 | 2.19 | 80.07 | 9.8 | |
| NBOR | 97.03 | 1.22 | 99.69 | 0.14 | 99.53 | 0.25 | 99.49 | 0.33 | 71.44 | 14.51 | 67.39 | 16.98 | |

| | Neuron Acc | Neuron FN | Pendulum Acc | Pendulum FN | Quadcopter Acc | Quadcopter FN | Cruise Acc | Cruise FN | Powertrain Acc | Powertrain FN | Helicopter Acc | Helicopter FN | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DNN-S | **99.83** | **0.12** | **99.89** | **0** | **99.82** | 0.04 | 99.94 | **0** | **97.2** | **0.86** | **98.24** | **0.79** | Balanced |
| DNN-R | 99.48 | 0.24 | 99.63 | 0.01 | 99.67 | 0.09 | **99.95** | **0** | 96.07 | 1.24 | 97.91 | 1.2 | |
| SNN | 98.89 | 0.69 | 99.2 | **0** | 99.49 | **0.01** | 99.6 | **0** | 95.21 | 1.79 | 97.58 | 1.16 | |
| SVM | 98.63 | 0.78 | 99.37 | **0** | 96.93 | 0.2 | 99.61 | **0** | 91.84 | 3.3 | 95.36 | 1.85 | |
| BDT | 99.07 | 0.45 | 99.46 | 0.05 | 99.36 | 0.22 | 99.9 | 0.03 | 95.86 | 2.4 | 79.03 | 10.26 | |
| NBOR | 96.95 | 1.62 | 99.51 | 0.04 | 99.11 | 0.56 | 99.47 | 0.11 | 71.33 | 13.99 | 65.18 | 17.48 | |

20K training samples, 10K test samples

**DNN-S**: Sigmoid DNN

**SVM**: Support Vector Machine

**SNN**: Shallow NN

**DNN-R**: ReLU DNN

**BDT**: Binary Decision Tree

**SNN**: Shallow NN

# Accuracy and FNs

If we increase training samples from 20K to 1M: **99.25 0.33 99.92 0.04**

20K training samples, 10K test samples

| | Neuron | | Pendulum | | Quadcopter | | Cruise | | Powertrain | | Helicopter | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | FN | |
| DNN-R | 99.52 | 0.29 | 99.93 | 0.04 | **99.89** | **0.06** | **99.98** | **0** | 96.21 | **1.08** | 98 | 0.96 | Uniform |
| SNN | 99.17 | 0.43 | 99.81 | **0** | 99.85 | 0.08 | 99.84 | 0.15 | 96.02 | 1.37 | 97.69 | 1.25 | |
| SVM | 98.73 | 0.75 | 99.84 | **0** | 97.33 | 0.69 | 99.88 | 0.1 | 92.26 | 3.48 | 95.58 | 2.42 | |
| BDT | 99.3 | 0.37 | 99.6 | 0.17 | 99.52 | 0.2 | 99.84 | 0.08 | 95.59 | 2.19 | 80.07 | 9.8 | |
| NBOR | 97.03 | 1.22 | 99.69 | 0.14 | 99.53 | 0.25 | 99.49 | 0.33 | 71.44 | 14.51 | 67.39 | 16.98 | |
| | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | FN | |
| DNN-S | **99.83** | **0.12** | **99.89** | **0** | **99.82** | 0.04 | 99.94 | **0** | **97.2** | **0.86** | **98.24** | **0.79** | Balanced |
| DNN-R | 99.48 | 0.24 | 99.63 | 0.01 | 99.67 | 0.09 | **99.95** | **0** | 96.07 | 1.24 | 97.91 | 1.2 | |
| SNN | 98.89 | 0.69 | 99.2 | **0** | 99.49 | **0.01** | 99.6 | **0** | 95.21 | 1.79 | 97.58 | 1.16 | |
| SVM | 98.63 | 0.78 | 99.37 | **0** | 96.93 | 0.2 | 99.61 | **0** | 91.84 | 3.3 | 95.36 | 1.85 | |
| BDT | 99.07 | 0.45 | 99.46 | 0.05 | 99.36 | 0.22 | 99.9 | 0.03 | 95.86 | 2.4 | 79.03 | 10.26 | |
| NBOR | 96.95 | 1.62 | 99.51 | 0.04 | 99.11 | 0.56 | 99.47 | 0.11 | 71.33 | 13.99 | 65.18 | 17.48 | |

**DNN-S**: Sigmoid DNN

**SVM**: Support Vector Machine

**SNN**: Shallow NN

**DNN-R**: ReLU DNN

**BDT**: Binary Decision Tree

**SNN**: Shallow NN

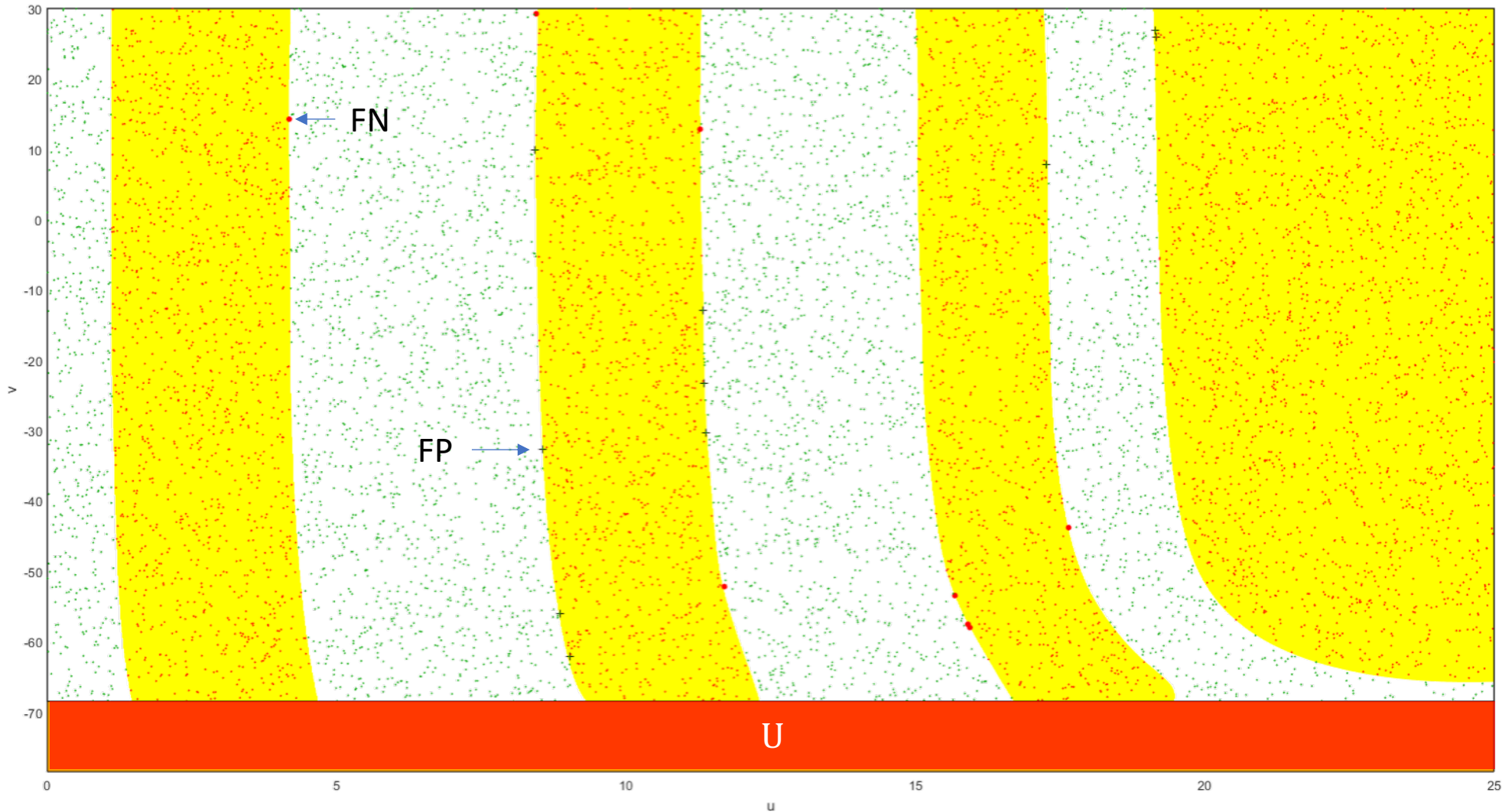# Statistical guarantees based on SPRT

$\theta_A = 99.7\%, \theta_{FN} = 0.2\%$

In parenthesis: number of samples needed to reach the decision

| | Neuron | | Pendulum | | Quadcopter | | Cruise | |
|---|---|---|---|---|---|---|---|---|
| | $P_A \geq \theta_A$ | $P_{FN} \leq \theta_{FN}$ | $P_A \geq \theta_A$ | $P_{FN} \leq \theta_{FN}$ | $P_A \geq \theta_A$ | $P_{FN} \leq \theta_{FN}$ | $P_A \geq \theta_A$ | $P_{FN} \leq \theta_{FN}$ |
| DNN-S | ✓ (5800) | ✓ (2900) | ✓ (2300) | ✓ (2300) | ✓ (4400) | ✓ (2300) | ✓ (3000) | ✓ (2300) |
| DNN-R | ✗ (3600) | ✗ (8600) | ✓ (15500) | ✓ (4000) | ✗ (1400) | ✓ (7300) | ✓ (3000) | ✓ (2300) |
| SNN | ✗ (700) | ✗ (1000) | ✗ (2900) | ✓ (2300) | ✗ (1500) | ✓ (3400) | ✗ (3600) | ✓ (2300) |
| SVM | ✗ (400) | ✗ (600) | ✗ (6600) | ✓ (2300) | ✗ (200) | ✗ (5300) | ✗ (3400) | ✓ (2300) |
| BDT | ✗ (1700) | ✗ (3300) | ✗ (6300) | ✓ (15000) | ✗ (800) | ✗ (1100) | ✓ (2700) | ✓ (2900) |
| NBOR | ✗ (300) | ✗ (300) | ✗ (28500) | ✓ (2900) | ✗ (1000) | ✗ (1300) | ✗ (3400) | ✗ (2300) |

Strength of test: $\alpha = \beta = 0.01$.

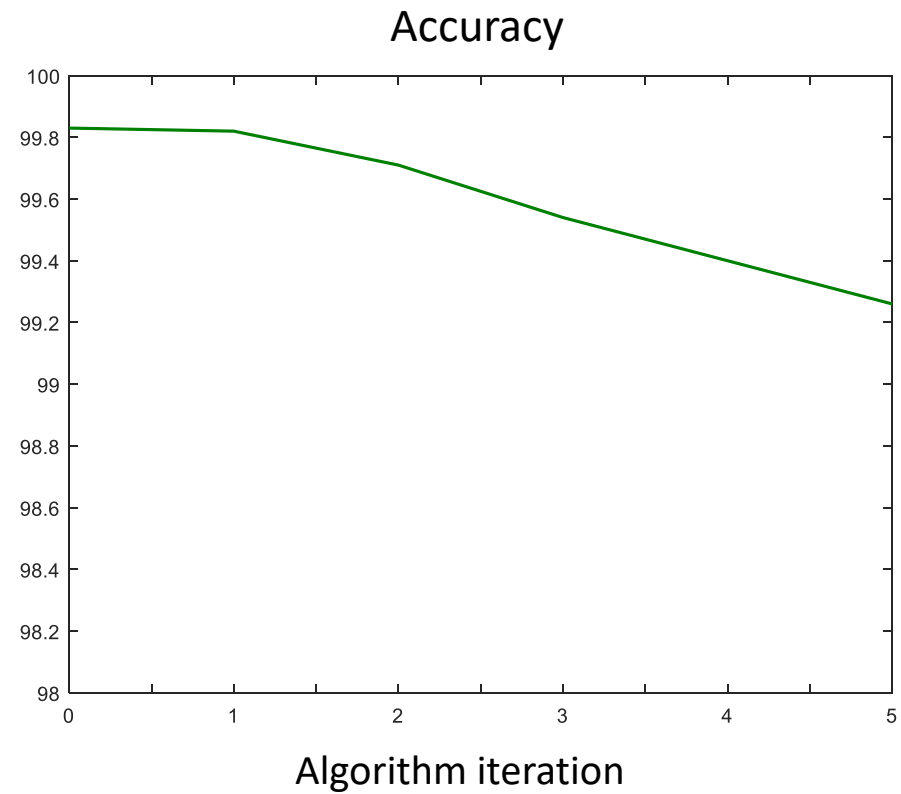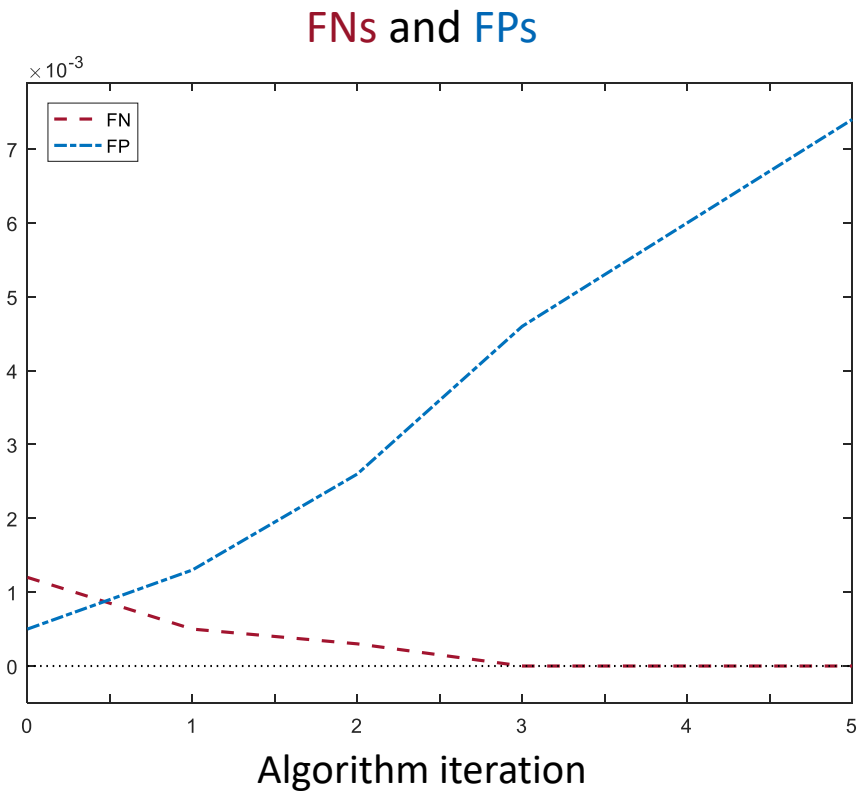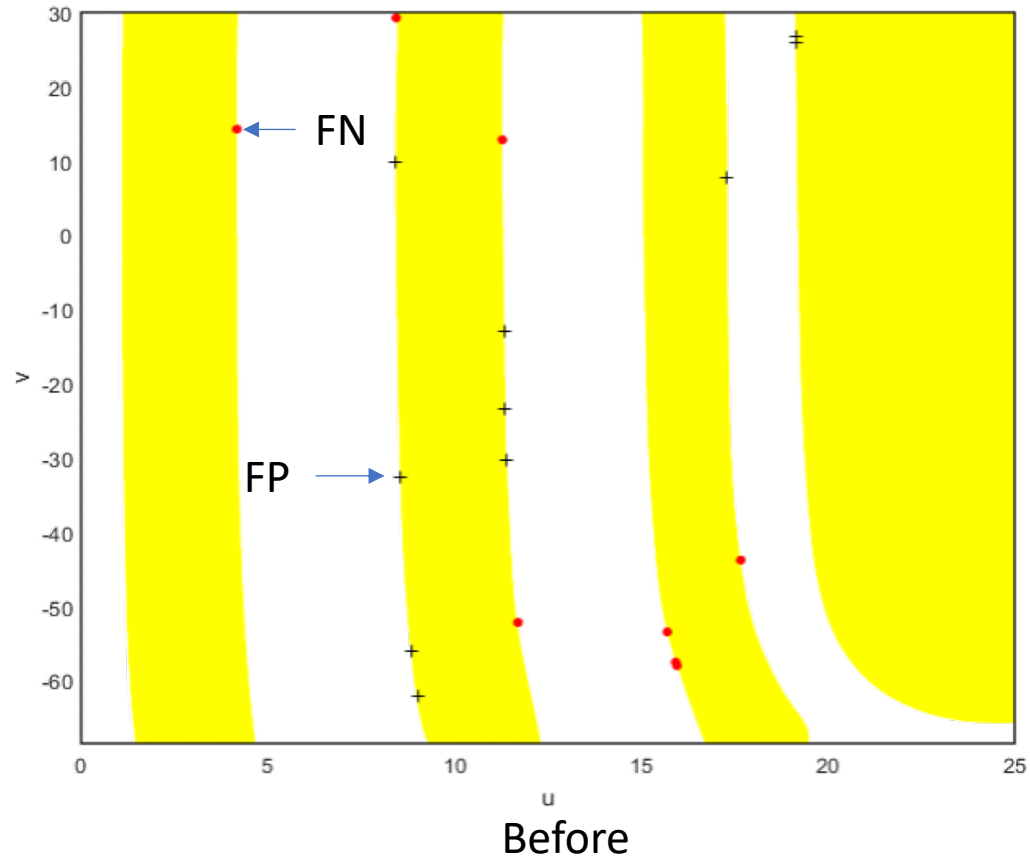# Reducing FNs...



NN prediction:
- positive (yellow)
- negative (white)

Unseen (test) state:
- positive (red)
- negative (green)

FN

FP

U

# …with falsification and re-training

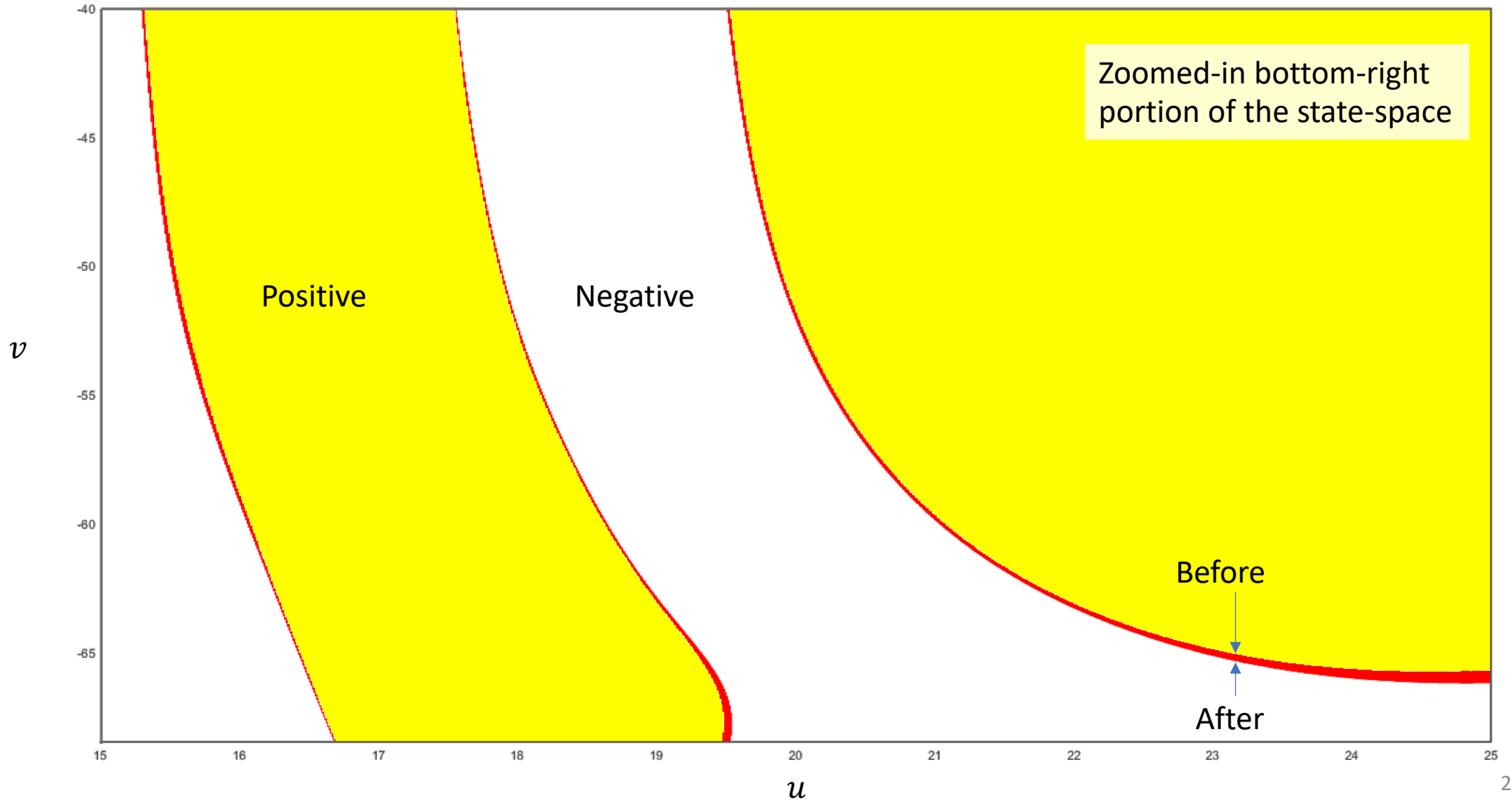FNs and FPs

Accuracy

# Reducing FNs



Before

Test FNs are eliminated and
the state classifier becomes more conservative



After

# Pushing the DNN decision boundary



Positive

Negative

Zoomed-in bottom-right portion of the state-space

Before

After

$v$

$u$

# Related work

## Machine-learning-aided verification

- Gaussian processes to approximate the satisfaction function of continuous-time Markov chains
  [Bortolussi et al, *Information and Computation 247* (2016)]
- NeuroSAT, learning to solve SAT problems from examples
  [Selsam et al, *arXiv:1802.03685* (2018)]
- Reinforcement learning of DNN policies for heuristics in QBF solvers [Lederman et al, *arXiv:1807.08058* (2018)]
- NN-based program synthesis from I/O examples
  [Parisotto et al, *arXiv:1611.01855* (2016)]

## Verification of NNs

- Robustness (absence of adversarial inputs)
  [Huang et al, *CAV* (2017); Gopinath et al, *ATVA* (2018)]
- Convex specifications
  [Katz et al, *CAV* (2017); Ehlers, *ATVA* (2017)]
- Analysis of NN components in-the-loop with CPS models
  [Dreossi et al, *NFM* (2017)]
- Range estimation for NNs (compute "reach set" of NN function)
  [Dutta et al, *NFM* (2018); Xiang et al, *IEEE Trans on Neural Networks and Learning Systems* (2018)]

# Conclusion

- State classification problem for hybrid systems
- NSC, a solution based on neural networks, efficient and with high accuracy
- Reverse HA construction for balanced sampling
- Statistical guarantees on classifier accuracy and FN rate
- Falsification-based techniques to reduce FNs and make classifier more conservative

**Future work:**
- More expressive properties, quantitative semantics, confidence intervals of point predictions