# 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)

Stony Brook University – 12 Oct 2018

# Agenda

- Background: hybrid systems verification
  - What are HS? Real-world examples
  - Why verify? Safety-critical applications
  - How verify? Formal models, reachability checking, online verification.
- Contribution: Neural State Classification
  - NN-based method to approximate verification results for online analysis
  - Sampling methods
  - Statistical guarantees
  - Reducing errors via falsification
- Experimental results
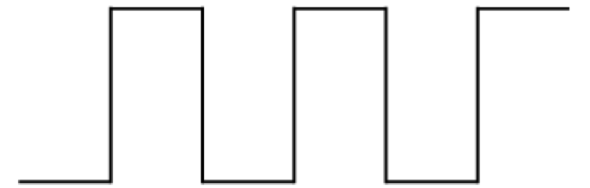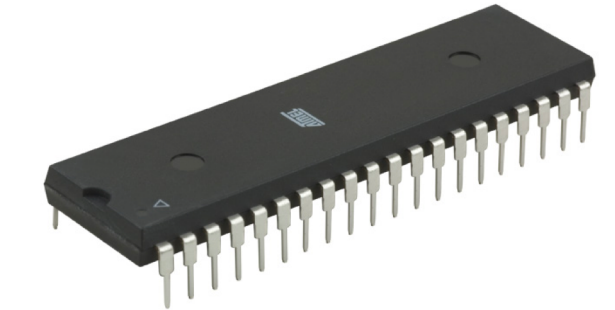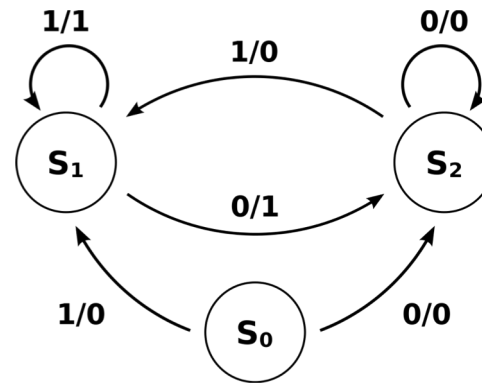
# Hybrid systems, informally

<span style="color:red">continuous / physical / analog</span> + discrete / digital components
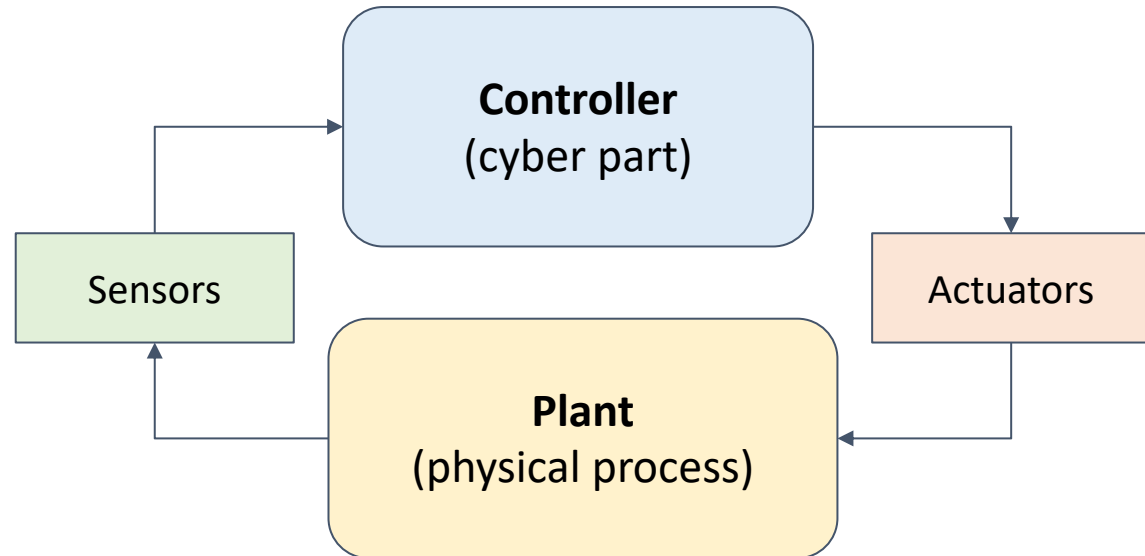


$$\frac{dn}{dt} = \alpha_n(V_m)(1-n) - \beta_n(V_m)n$$

$$\frac{dm}{dt} = \alpha_m(V_m)(1-m) - \beta_m(V_m)m$$

# Hybrid systems, examples

**Cyber-physical systems**
(aka control systems)

# Hybrid systems, examples

## Embedded systems
(building blocks of the Internet of Things)



| ADC → Microcontroller → DAC |

Sensors — Microcontroller (ASIC, FPGA) — Actuators

Physical process

# Hybrid systems, examples

**Artificial pancreas**

Glucose-insulin metabolism

Glucose monitor

Insulin pump

Sugar levels

**Cardiac devices**

**Closed-loop deep brain stimulation**

DBS leads

Cerebrum

Subthalamic nucleus

Cerebellum

Extension wire

Neurostimulator *(battery pack)*

Source: Medtronic

R.L. REB/

# Hybrid systems, examples

# Safety assurance, how?

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

## How do we ensure that they work as intended?

e.g., pacemaker always keeps its pacing rate within healthy bounds,
cruise control always maintains safety distance, collision freedom, etc

# The verification problem



System model $\mathcal{M}$ → **Verification** (aka model checking) $\mathcal{M} \models \phi?$ → ✓ (proof)

Specification $\phi$ → → ✗ counterexample

- Verification is automated and exhaustive (considers all possible system's behaviors)
- $\mathcal{M}$ is a formal, executable model
- $\phi$ is a correctness property over time
  - Liveness: "at any time, something good must eventually happen"
  - Safety: "something bad will never happen"
  - …

# Hybrid systems, formally

**Hybrid automata** [Henzinger, LICS 1996]

- Set of discrete locations: $Loc$
- Set of continuous variables: $Var$, over $X \subseteq \mathbb{R}$
- Initial set of states: $Init \subseteq Loc \times X$
- Invariant: $Inv: Loc \rightarrow 2^X$
- Flow function (continuous evolution, ODEs): $Flow: Loc \rightarrow (X \rightarrow X)$
- Transition relation (discrete jumps):
  - Jumps from source location to target location if guard condition holds
  - Updates variables before reaching target

# Hybrid automata - Examples

## Bouncing ball



**Transition**

**Reset**

**Guard** $x_1 = 0 \land$
$x_2 < 0$

$x_2 := -c \ x_2$

**Flow function** $\longrightarrow$ $\dot{x}_1 = x_2$
$\dot{x}_2 = -g$

**Invariant** $\longrightarrow$ $x_1 \geq 0$

$q_0$

**Location**

$x_2$

$x_1$

$time$

Claire J. Tomlin AA278A lecture notes

# Hybrid automata - Examples

**Thermostat**



$$\dot{x} = -x$$
$$x > 68$$
$$q_0$$

$$x \leqslant 70$$

$$\dot{x} = -x + 100$$
$$x < 82$$
$$q_1$$

$$x \geqslant 80$$

Claire J. Tomlin AA278A lecture notes

# Hybrid automata in action

**Timed automata network of Boston Scientific dual chamber pacemaker**



(a) LRI component  (b) AVI component  (c) URI component  (d) PVARP component  (e) VRP component

Jiang et al, TACAS 2012

# Hybrid automata in action

**HA model of cardiac cell action potential (Smolka et al)**



(a) Phase 1    (b) Phase 2    (c) Phase 3    (d) Phase 4



$q_0$ : **Resting & FR**

$\dot{v}_x = \alpha_x^0 v_x$

$\dot{v}_y = \alpha_y^0 v_y$

$\dot{v}_z = \alpha_z^0 v_z$

$v = v_x - v_y + v_z$

$\{v \leq V_R\}$

$[V_{in} \neq 0]$
$v_x' = v_x$
$v_y' = v_y$
$v_z' = v_z$
$\theta' = v/V_R$

$[V_{in} \leq 0 \wedge v < V_T]$
$v_x' = v_x$
$v_y' = v_y$
$v_z' = v_z$

$q_1$ : **Stimulated**

$\dot{v}_x = \alpha_x^1 v_x + \beta_x V_{in}$

$\dot{v}_y = \alpha_y^1 v_y + \beta_y V_{in}$

$\dot{v}_z = \alpha_z^1 v_z + \beta_z V_{in}$

$v = v_x - v_y + v_z$

$\{v \leq V_T\}$

$[v \leq V_R]$
$v_x' = v_x$
$v_y' = v_y$
$v_z' = v_z$

$[v \geq V_T]$
$v_x' = v_x$
$v_y' = v_y$
$v_z' = v_z$

$q_3$ : **Plateau & ER**

$\dot{v}_x = \alpha_x^3 v_x$

$\dot{v}_y = \alpha_y^3 f(\theta) v_y$

$\dot{v}_z = \alpha_z^3 v_z$

$v = v_x - v_y + v_z$

$\{v \geq V_R\}$

$[v \geq V_O - 80.1\sqrt{\theta}]$
$v_x' = v_x$
$v_y' = v_y$
$v_z' = v_z$

$q_2$ : **Upstroke**

$\dot{v}_x = \alpha_x^2 v_x$

$\dot{v}_y = \alpha_y^2 v_y$

$\dot{v}_z = \alpha_z^2 v_z$

$v = v_x - v_y + v_z$

$\{v \leq V_O - 80.1\sqrt{\theta}\}$

# Hybrid automata in action

**HA model of prostate cancer treatment**



**Mode 1 (on-treatment)**

$\text{flow}_1:$

$$\frac{dx}{dt} = \left( \alpha_x \left( k_1 + \frac{(1-k_1)z}{z+k_2} \right) - \beta_x \left( k_3 + \frac{(1-k_3)z}{z+k_4} \right) - m_1 \left( 1 - \frac{z}{z_0} \right) \right) x$$

$$\frac{dy}{dt} = m_1 \left( 1 - \frac{z}{z_0} \right) x + \left( \alpha_y \left( 1 - d \frac{z}{z_0} \right) - \beta_y \right) y$$

$$\frac{dz}{dt} = \frac{-z}{\tau}$$

$$\frac{dv}{dt} = \left( \alpha_x \left( k_1 + \frac{(1-k_1)z}{z+k_2} \right) - \beta_x \left( k_3 + \frac{(1-k_3)z}{z+k_4} \right) - m_1 \left( 1 - \frac{z}{z_0} \right) \right) x$$
$$+ m_1 \left( 1 - \frac{z}{z_0} \right) x + \left( \alpha_y \left( 1 - d \frac{z}{z_0} \right) - \beta_y \right) y$$

$\text{jump}_{1 \rightarrow 2}:$

$$x + y \le r_0 \wedge \frac{dx}{dt} + \frac{dy}{dt} < 0$$

$\text{jump}_{2 \rightarrow 1}:$

$$x + y \ge r_1 \wedge \frac{dx}{dt} + \frac{dy}{dt} > 0$$

**Mode 2 (off-treatment)**

$\text{flow}_2:$

$$\frac{dx}{dt} = \left( \alpha_x \left( k_1 + \frac{(1-k_1)z}{z+k_2} \right) - \beta_x \left( k_3 + \frac{(1-k_3)z}{z+k_4} \right) - m_1 \left( 1 - \frac{z}{z_0} \right) \right) x$$

$$\frac{dy}{dt} = m_1 \left( 1 - \frac{z}{z_0} \right) x + \left( \alpha_y \left( 1 - d \frac{z}{z_0} \right) - \beta_y \right) y$$

$$\frac{dz}{dt} = \frac{z_0 - z}{\tau}$$

$$\frac{dv}{dt} = \left( \alpha_x \left( k_1 + \frac{(1-k_1)z}{z+k_2} \right) - \beta_x \left( k_3 + \frac{(1-k_3)z}{z+k_4} \right) - m_1 \left( 1 - \frac{z}{z_0} \right) \right) x$$
$$+ m_1 \left( 1 - \frac{z}{z_0} \right) x + \left( \alpha_y \left( 1 - d \frac{z}{z_0} \right) - \beta_y \right) y$$

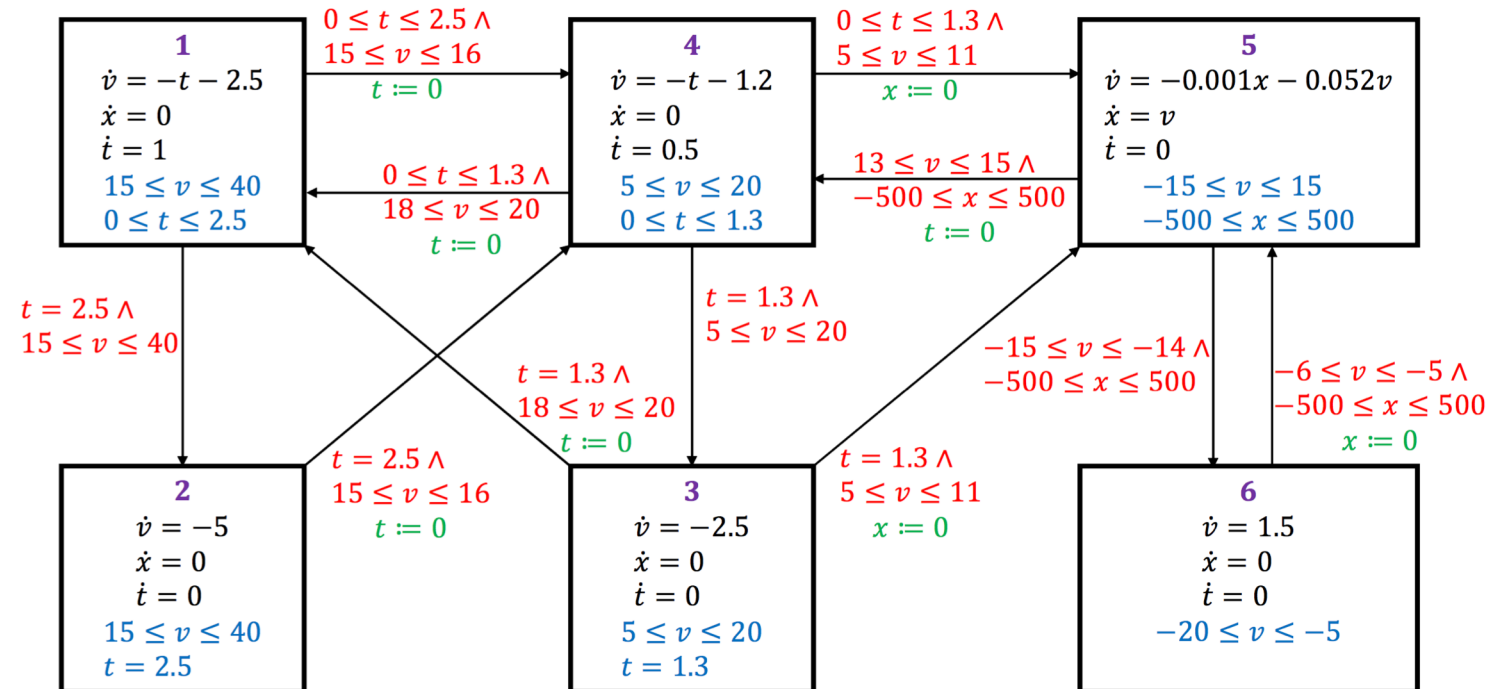Ideta et al, J. Nonlinear Sci. 18 (2008)

# Hybrid automata in action

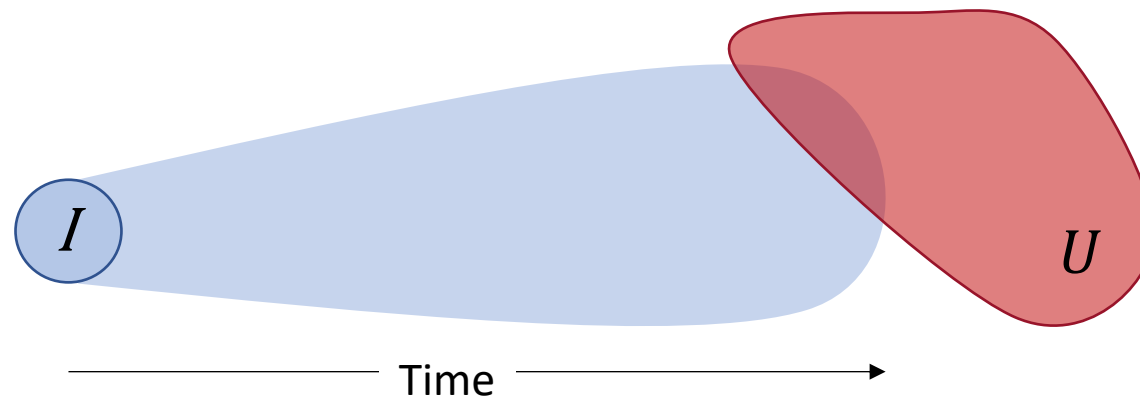**Powertain system by Toyota**



**Cruise control HA model**

# Hybrid automata verification

HA verification problem usually formulated as reachability

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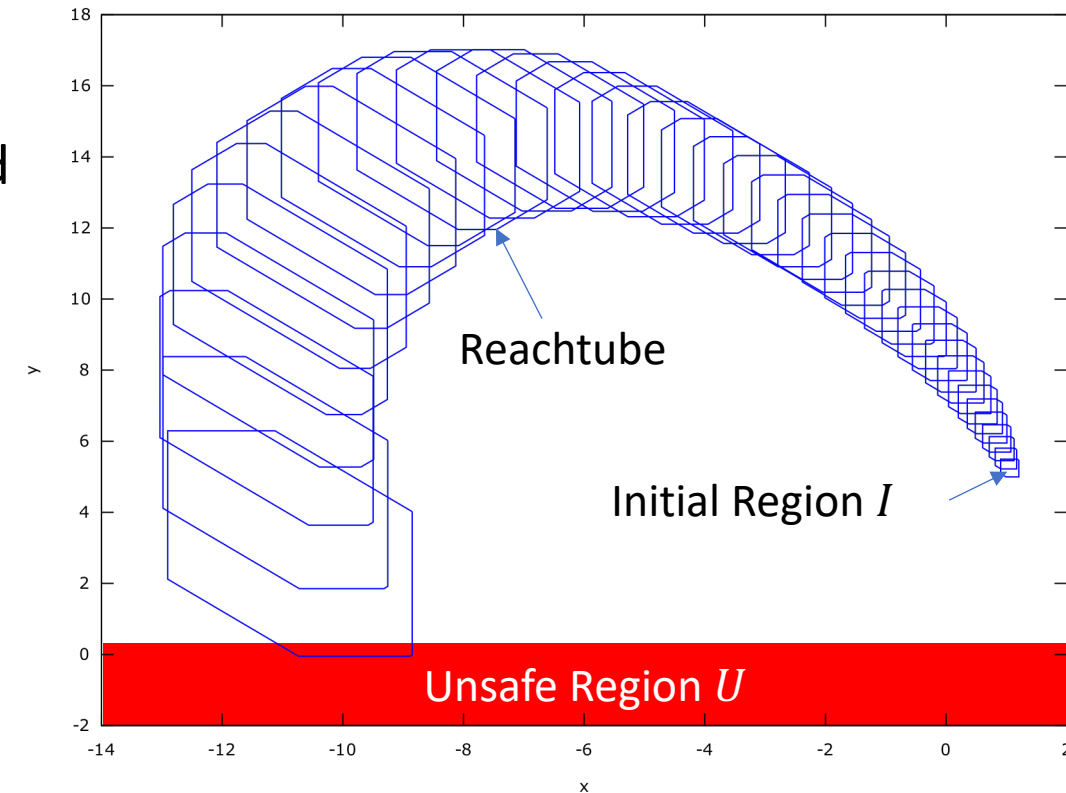can an HA $\mathcal{M}$, starting in an initial region $I$, reach a state $u \in U$ (within time $T$)?



Both bounded and unbounded versions are undecidable
[Henzinger et al, *JCSS 57 1* (1998); Brihaye et al, *ICALP* (2011)]
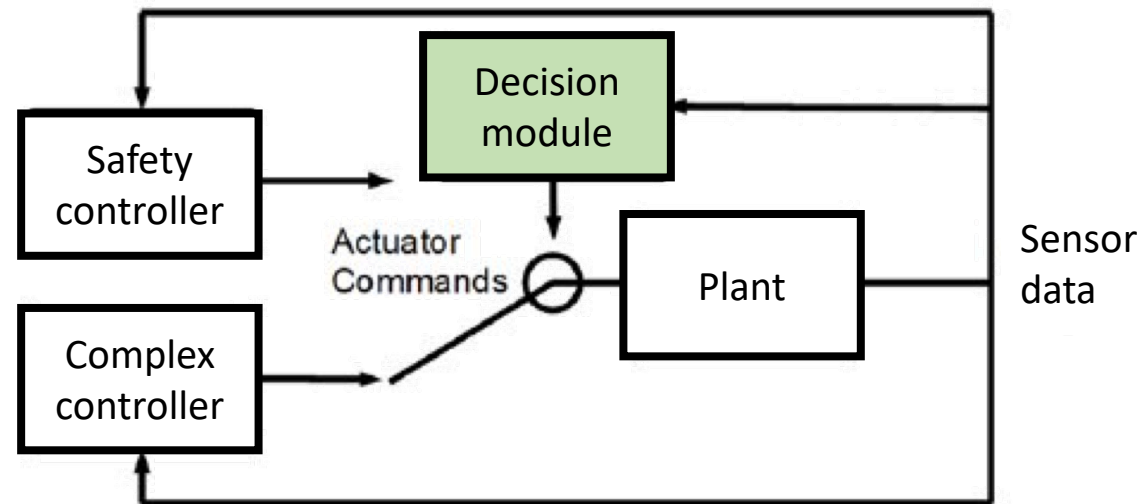
# Reachability checkers for HAs

- 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)])

# Motivation - Online model checking (OMC)

## Offline

- Reachability from a (large) region

- One-off analysis, potentially long time horizons (blow-up of over-approximation)

- No hard time constraints

- Controlled settings
  - Model is ground truth

## Online

- Reachability from a single state

- Analysis run periodically → short time horizons

- Strict time constraints

- Less predictable settings
  - Real system might differ from model
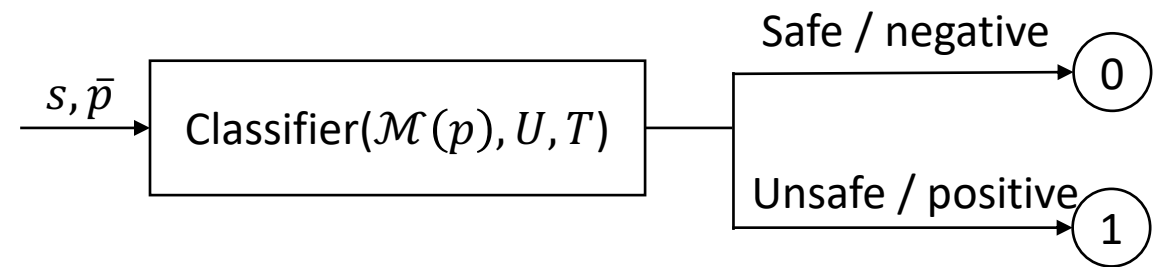  - Noisy observations

# 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
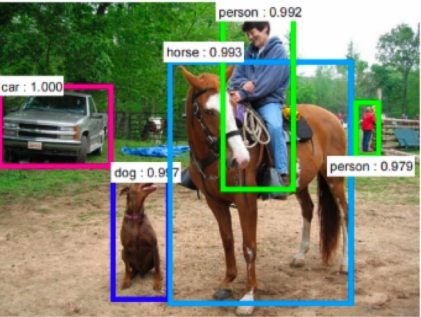
# State Classification Problem (SCP)

- 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 $T$;
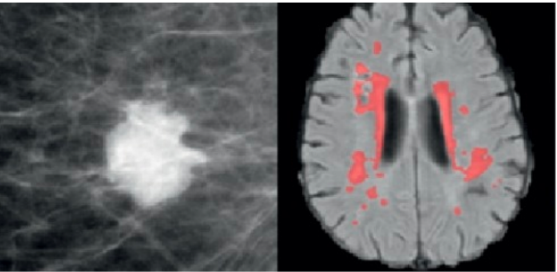
  - *negative* o/w

$$s, \bar{p} \longrightarrow \boxed{\text{Classifier}(\mathcal{M}(p), U, T)}$$

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 $p$

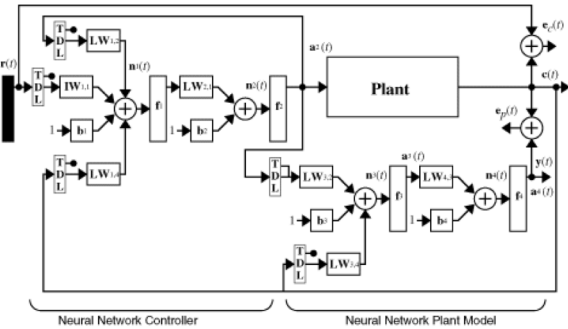# Neural networks (NNs) as state classifiers

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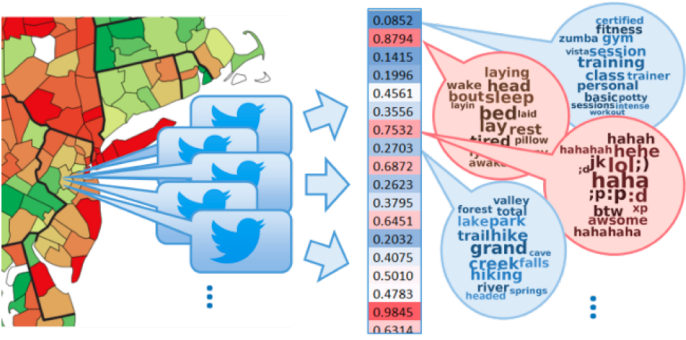


Object detection



Classification of tumor and
diseases from medical images
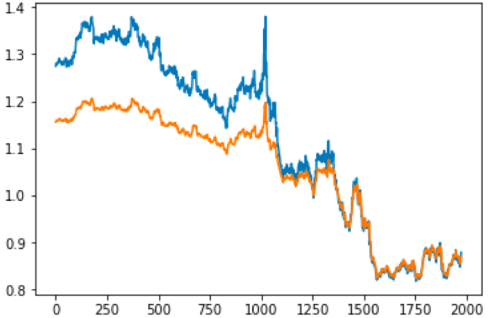


System identification and control



Natural language processing, sentiment analysis

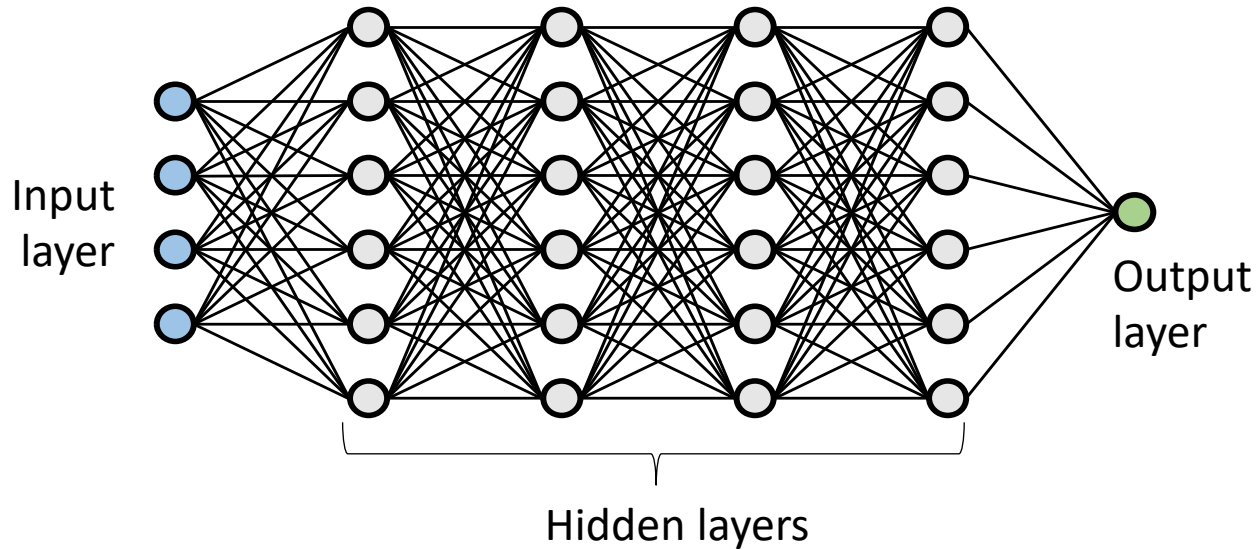Image credits: H. Andrew Schwartz

**?**

**Verification**



Time-series analysis and prediction

# Feedforward neural networks



Input layer

Output layer

Hidden layers

$$F = f_l \circ f_{l-1} \circ \ldots \circ f_1 \circ f_0$$

$$f_i\left(p_{i-1}\right) = g_i\left(W_{i,i-1} \cdot p_{i-1} + b_i\right), \qquad i = 1, \ldots, l$$

weights

Output of layer $i$

Activation function (sigmoid, ReLU, …)

Output of layer $i\text{-}1$

biases

Supervised learning of NN =
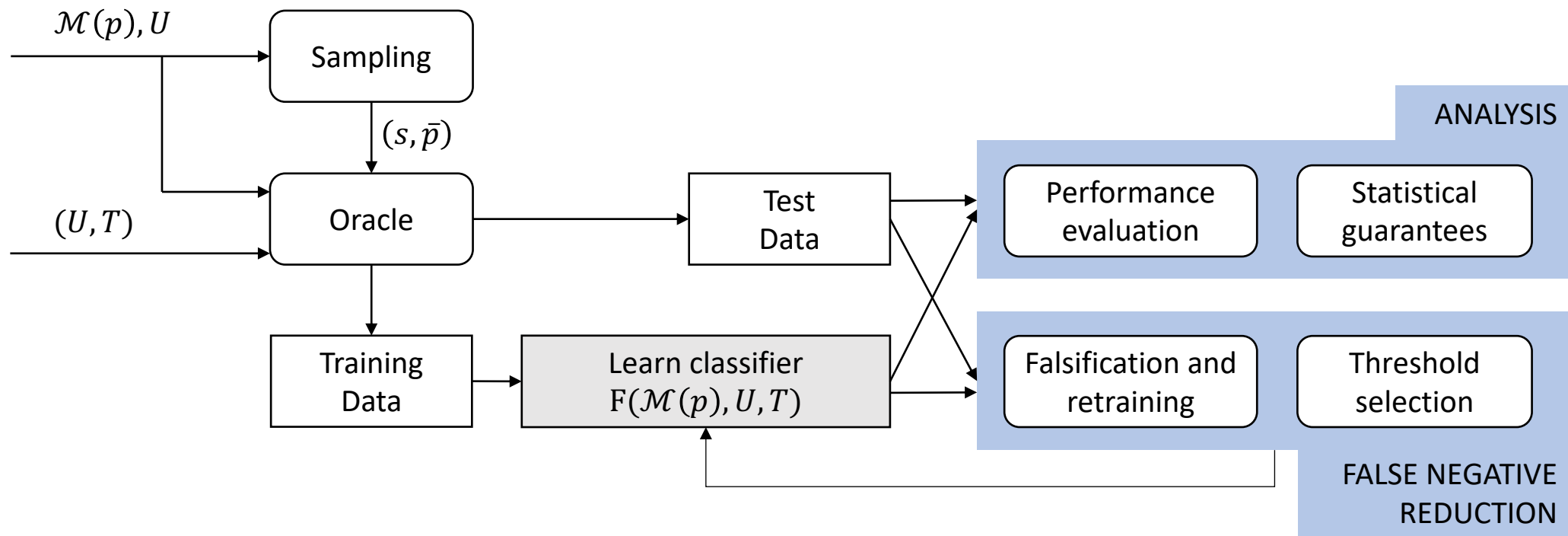finding weights and biases that maximize the fit between predictions and training data

# 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

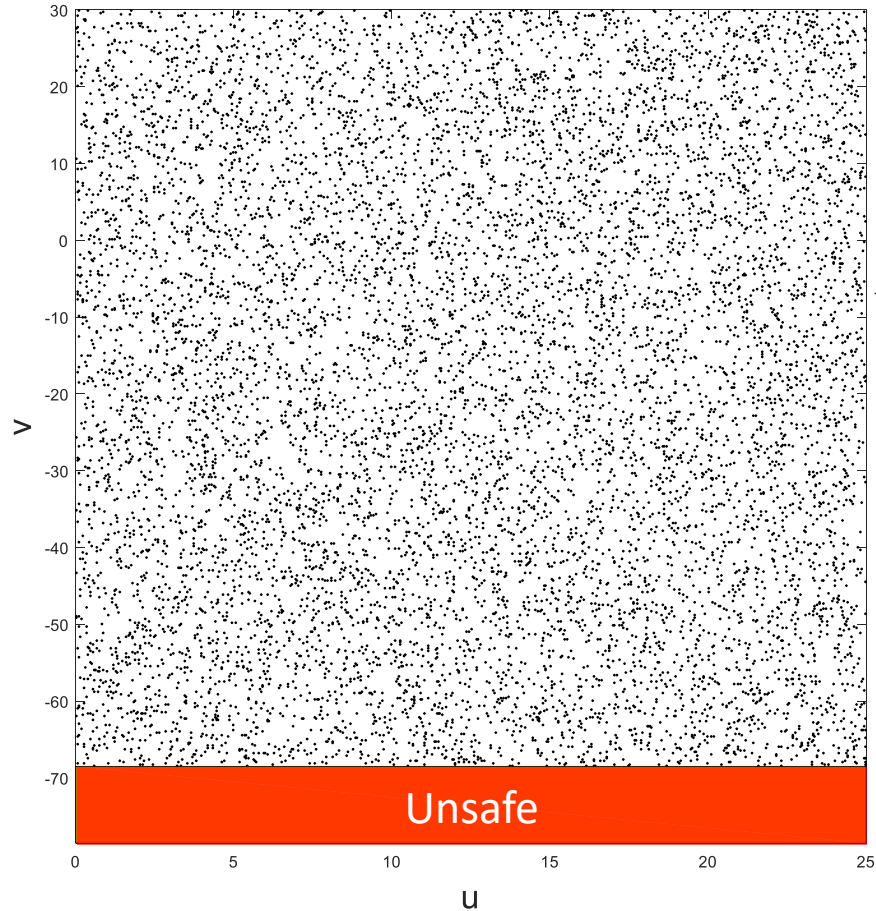Two kinds of errors in **neural state classification**:

- **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



Oracle

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

# 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 don't just want empirical performance, but also to establish guaranteed performance requirements

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

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

- Deriving absolute guarantees is infeasible

- statistical guarantees (precise up to a small error probability) via the *sequential probability ratio test* (SPRT) [Wald and Wolfowitz (1948)]

# Sequential probability ratio test

- *Sequential* means that we only need the number of test samples necessary to decide the threshold

- Precise up to arbitrary error bounds $\alpha$ (prob of type-I errors) and $\beta$ (prob of type-II errors)

- To ensure both bounds, the test $P \geq \Theta$ vs $P < \Theta$ is relaxed to

  - $H_0 : P \geq p_0$ vs $H_1 : P \leq p_1$ where $p_1 < \Theta < p_0$ (but both close to $\Theta$)

- $H_0$ accepted if $\frac{p_{1m}}{p_{0m}} \leq \frac{1-\beta}{\alpha}$; $H_1$ accepted if $\frac{p_{1m}}{p_{0m}} \geq \frac{\beta}{1-\alpha}$

- $\frac{p_{1m}}{p_{0m}} = \frac{p_1^{t_m}(1-p_1)^{f_m}}{p_0^{t_m}(1-p_0)^{f_m}}$ , $t_m$ : # pos. samples; $f_m$ : # neg. samples

# 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

$$\max_{s \in S} |b(s) - F(s)|$$

True label of $s$    Network prediction for $s$

**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

# Reducing FN rate via falsification

- 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)])

$$for\ all\ \eta \in (0,1),\ \Pr(\lim_{k \to \infty} FN_k = \emptyset) \geq 1 - \eta$$

under assumptions that:
- Falsifier always finds a FN if it exists
- Classifier doesn't make mistakes on positive training samples
- FP rate for test data is not below that for training data

```
Input:     classifier (NN) F,
           training samples D
Output:    "conservative" classifier F
do
•   FN ← subset of the true FN set of F
    /*found via falsifier (genetic alg)*/
•   D ← D ∪ FN
•   F ← train(D)
while FN ≠ ∅ 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 | | Pendulum | | Quadcopter | | Cruise | | Powertrain | | Helicopter | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | 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** |
| 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 |

*Uniform*

| | 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** |
| 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 |

*Balanced*

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:

|        | Neuron | | Pendulum | | Quadcopter | | Cruise | | Powertrain | | Helicopter | |
|--------|--------|------|----------|------|------------|------|--------|------|------------|------|------------|------|
|        | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | FN | Acc | FN |
|        |  |  |  |  |  |  |  |  | **99.25** | **0.33** | **99.92** | **0.04** |
| 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 |

*Uniform*

|        | 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** |
| 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 |

*Balanced*

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

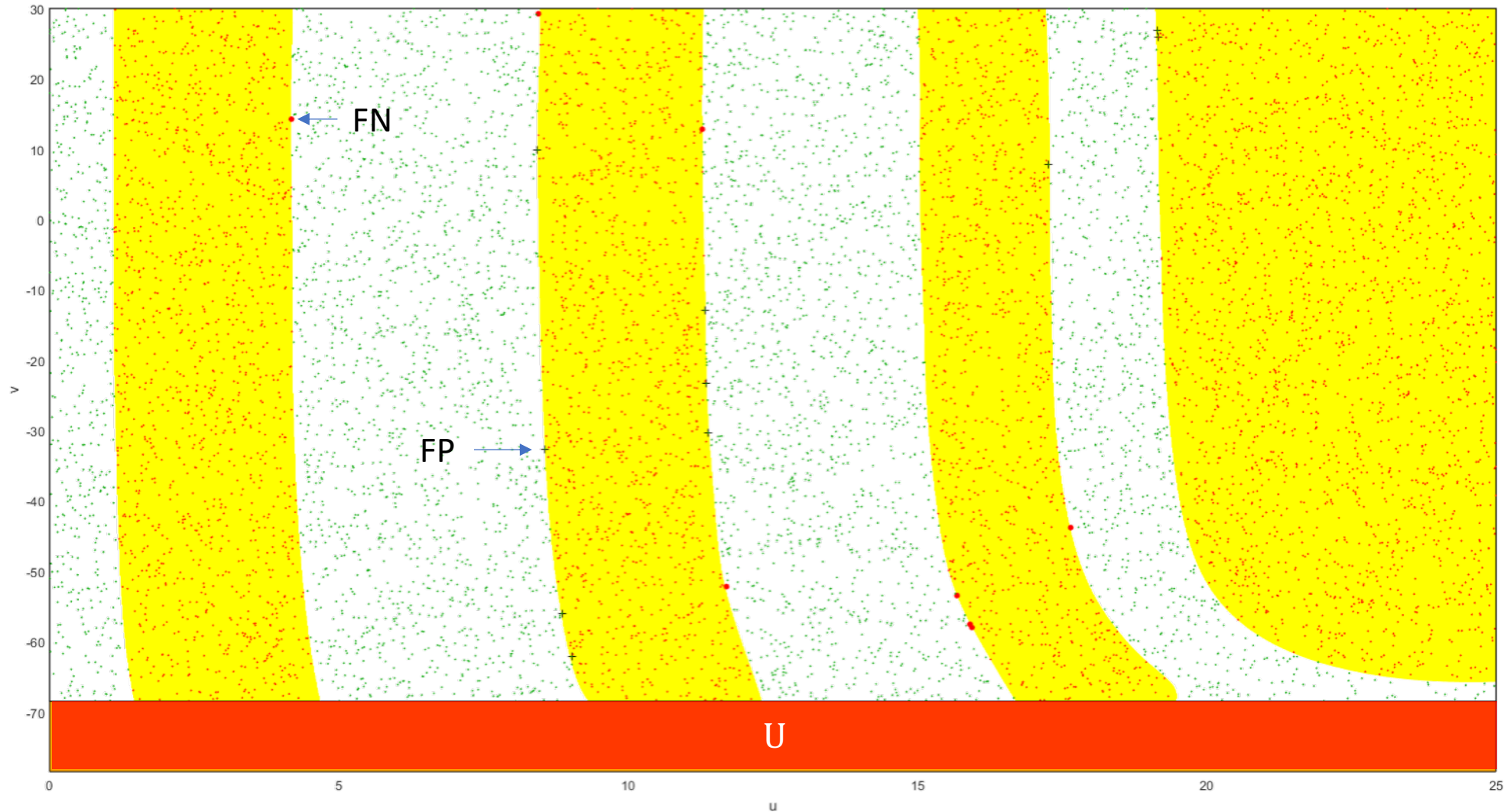**SNN**: Shallow NN

# 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…

# …with falsification and re-training

FNs and FPs
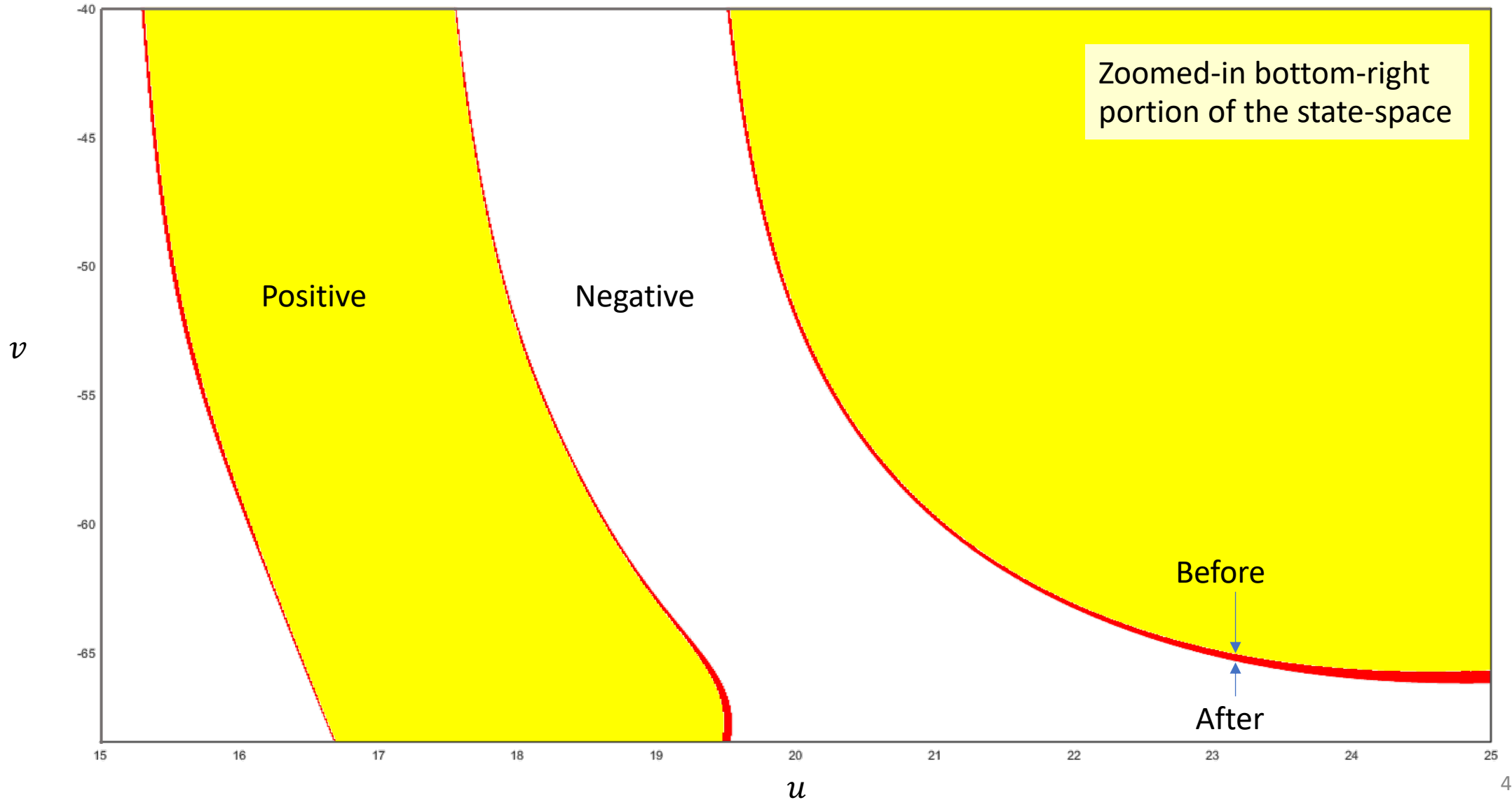
Accuracy

# Reducing FNs



Before

Test FNs are eliminated and
the state classifier becomes more conservative



After

# Pushing the DNN decision boundary

# 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
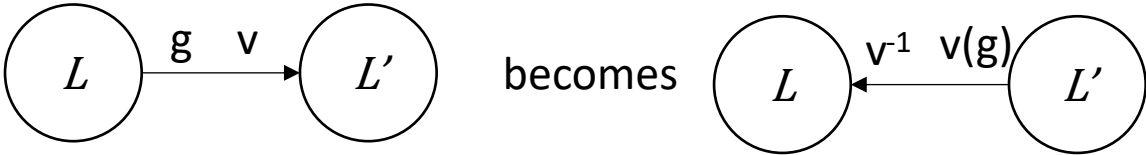- Falsification-based techniques to reduce FNs and make classifier more conservative

**Future work:**

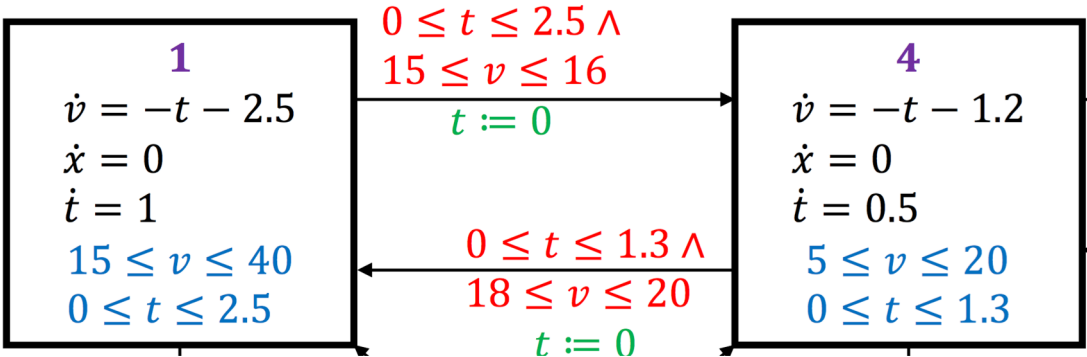- More expressive properties, quantitative semantics, confidence intervals of point predictions

# Backup slides

# Reverse HA automaton
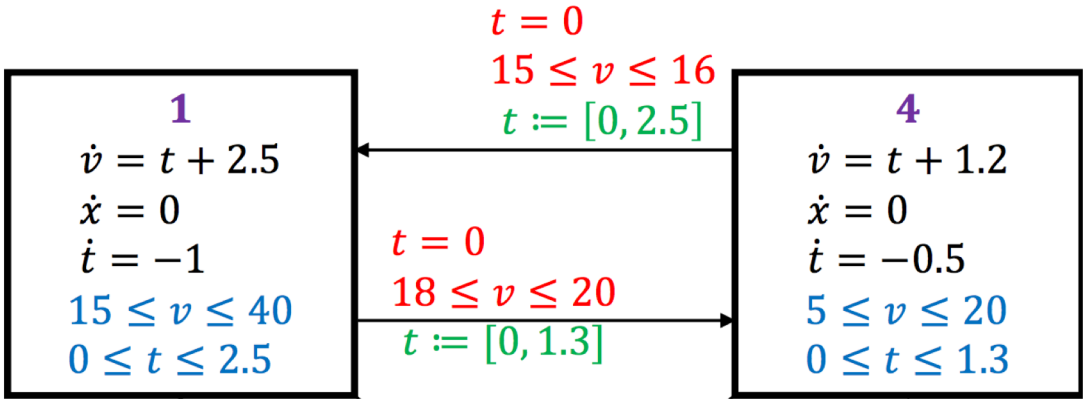
- Locations and invariants stay the same
- Flows are reversed (sign changes)



becomes



**Forward**

**Reverse**



$$\begin{array}{c}\textbf{1}\\ \dot{v}=-t-2.5\\ \dot{x}=0\\ \dot{t}=1\\ 15\leq v\leq 40\\ 0\leq t\leq 2.5\end{array}$$

$$\begin{array}{c}0\leq t\leq 2.5\land\\ 15\leq v\leq 16\\ t:=0\end{array}$$

$$\begin{array}{c}\textbf{4}\\ \dot{v}=-t-1.2\\ \dot{x}=0\\ \dot{t}=0.5\\ 5\leq v\leq 20\\ 0\leq t\leq 1.3\end{array}$$

$$\begin{array}{c}0\leq t\leq 1.3\land\\ 18\leq v\leq 20\\ t:=0\end{array}$$

$$\begin{array}{c}\textbf{1}\\ \dot{v}=t+2.5\\ \dot{x}=0\\ \dot{t}=-1\\ 15\leq v\leq 40\\ 0\leq t\leq 2.5\end{array}$$

$$\begin{array}{c}t=0\\ 15\leq v\leq 16\\ t:=[0,2.5]\end{array}$$

$$\begin{array}{c}\textbf{4}\\ \dot{v}=t+1.2\\ \dot{x}=0\\ \dot{t}=-0.5\\ 5\leq v\leq 20\\ 0\leq t\leq 1.3\end{array}$$

$$\begin{array}{c}t=0\\ 18\leq v\leq 20\\ t:=[0,1.3]\end{array}$$

# Hybrid automata in action

**Timed automata network of task scheduling in Boeing Bold Stroke platform**



Madl et al, RTSS 2006