# Formal Methods meet Neural Networks:
# A Selection

## Tom Henzinger
### IST Austria

# Formal Methods

1. Static: ensure at design time that a system satisfies its specification

   1a. Verification: given system $f$ and spec $\varphi$, does $f$ satisfy $\varphi$?

   - deductive reasoning (logic, decision procedures)
   - algorithmic reasoning (model checking, abstract interpretation)

# Formal Methods

1. Static: ensure at design time that a system satisfies its specification

   1a. Verification: given system $f$ and spec $\varphi$, does $f$ satisfy $\varphi$?

   - deductive reasoning (logic, decision procedures)
   - algorithmic reasoning (model checking, abstract interpretation)

   $\forall$traces

   1b. Synthesis: given spec $\varphi$, find system $f$ such that $f$ satisfies $\varphi$.

   - syntax-guided (translation, search, learning)
   - semantics-guided (game solving, control)

   $\forall$inputs $\exists$output

# Formal Methods

1. Static: ensure at design time that a system satisfies its specification

    1a. Verification: given system $f$ and spec $\varphi$, does $f$ satisfy $\varphi$?

    - deductive reasoning (logic, decision procedures)
    - algorithmic reasoning (model checking, abstract interpretation)        $\forall$traces

    1b. Synthesis: given spec $\varphi$, find system $f$ such that $f$ satisfies $\varphi$.

    - syntax-guided (translation, search, learning)
    - semantics-guided (game solving, control)        $\forall$inputs $\exists$output

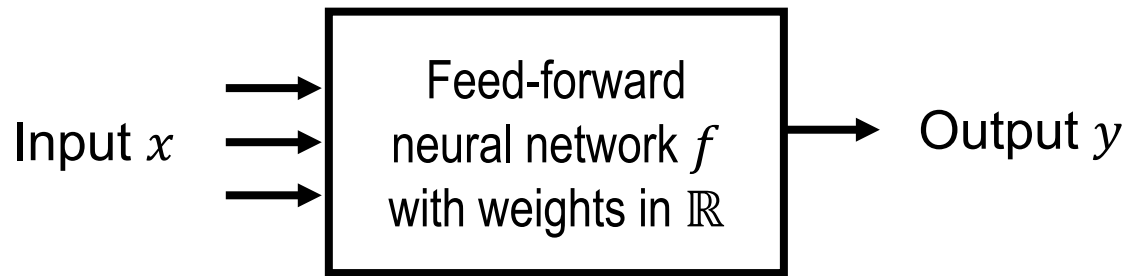2. Dynamic: watch at runtime if a system satisfies its specification

    2a. Runtime monitoring
    2b. Runtime enforcement

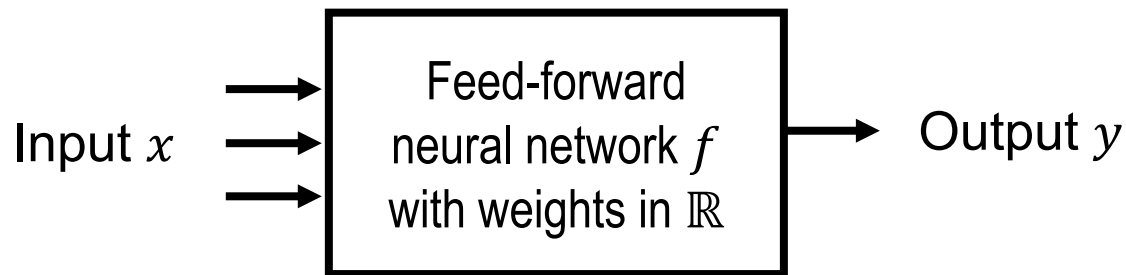# Formal Methods meet Neural Networks: A Selection

1. Open-loop verification: prove properties of neural networks [TACAS'20, AAAI'21] (joint work with Giacobbe-Lechner-Zikelic)

2. Closed-loop verification: prove properties of neural network controllers over discrete-time dynamical systems [NeurIPS'21, AAAI'22] (joint work with Lechner-Zikelic-Chatterjee)

3. Monitoring: runtime monitors as novelty detectors [ECAI'20, RV'21] (joint work with Lukina-Schilling)

4. Synthesis/enforcement: runtime monitors as specification enforcers [CAV'19] (joint work with Avni-Bloem-Chatterjee-Koenighofer-Pranger)

# Open-Loop Verification



$$\varphi(x) \;\wedge\; f(x) = y \;\Rightarrow\; \psi(y)$$

# Open-Loop Verification



Feed-forward neural network $f$ with weights in $\mathbb{R}$

Input $x$

Output $y$
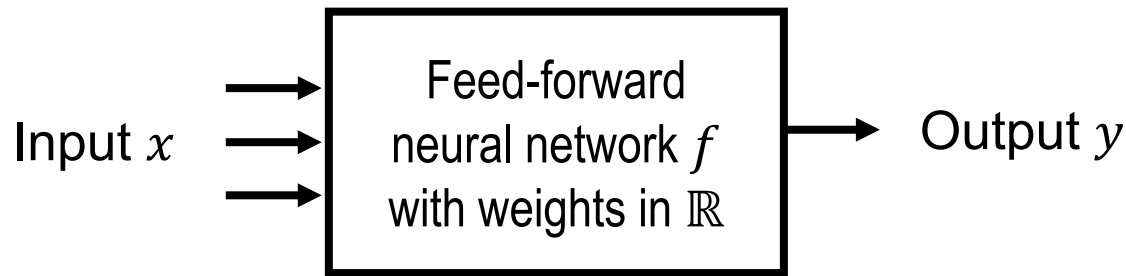
$$\varphi(x) \wedge f(x) = y \Rightarrow \psi(y)$$

NP-complete
[Katz'17]

$$\varphi(x) \wedge [f]_{\text{float32}}(x) = y \Rightarrow \psi(y)$$

$$\varphi(x) \wedge [f]_{\text{int8}}(x) = y \Rightarrow \psi(y)$$

PSPACE-hard
[AAAI'21]

# Open-Loop Verification



Input $x$ → Feed-forward neural network $f$ with weights in $\mathbb{R}$ → Output $y$

1. Abstraction domains over $\mathbb{R}^n$
(intervals, zonotopes, polyhedra)
2. Constraint solvers
(MILP, Reluplex, NRA)

[Ehlers, Katz et al.,
Bunel et al., Dutta et al.,
Mirman et al., Wang et al.,
Chang et al., Tjeng et al.,
etc.]

$$\varphi(x) \,\wedge\, f(x) = y \,\Rightarrow\, \psi(y)$$

$$\varphi(x) \,\wedge\, [f]_{\text{float32}}(x) = y \,\Rightarrow\, \psi(y)$$

Bitvector SMT

$$\varphi(x) \,\wedge\, [f]_{\text{int8}}(x) = y \,\Rightarrow\, \psi(y)$$

[Baranowski et al.,
TACAS'20, AAAI'21]

# Open-Loop Verification

Input $x$ →→→ [ Feed-forward neural network $f$ with weights in $\mathbb{R}$ ] → Output $y$
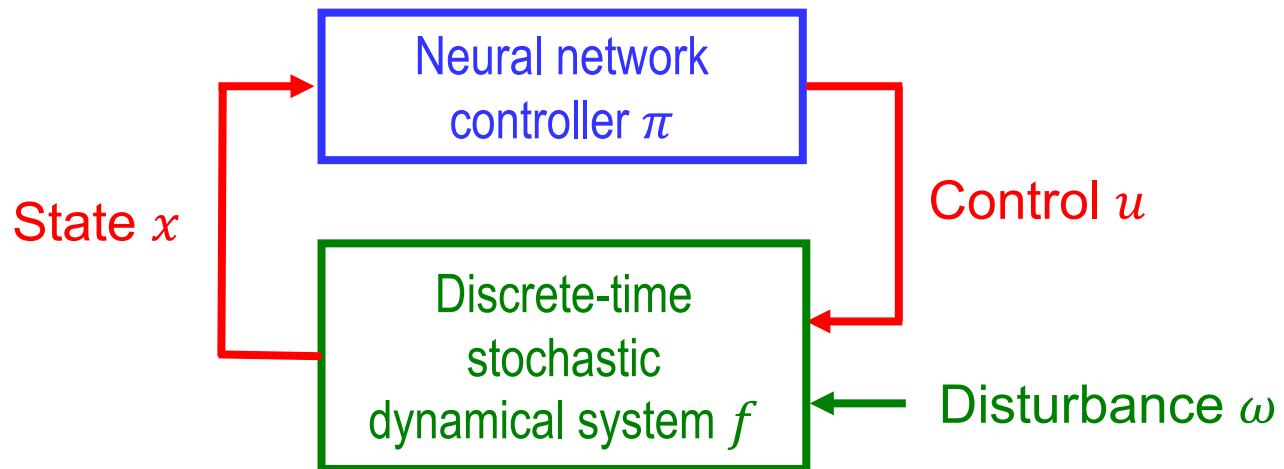
$\varepsilon$-robustness at input $x$:

$$(\forall x')\left(d(x,x') < \varepsilon \ \Rightarrow\ f(x) = f(x')\right)$$

Robustness is not preserved by $[\cdot]_{\mathrm{float}k}$ nor by $[\cdot]_{\mathrm{int}k}$ nor is robustness monotonic in $k$ [TACAS'20, Jia-Rinard].

# Closed-Loop Verification



*Safety:* unsafe state set $X_U$ is never entered (infinite time horizon!).

*Progress (stability):* stable state set $X_S$ is eventually entered with probability 1 ("almost surely") and never left.

# Discrete-Time Stochastic Dynamical System

| | |
|---|---|
| State space | $X \subseteq \mathbb{R}^n$ |
| | |
| Control policy | $\pi: X \to \mathcal{D}(U)$ |
| Disturbance | $\omega \sim \mathcal{D}(W)$ |
| | |
| Initial state | $x_0 \in X$ |
| Transition function | $f: X \times U \times W \to X$ |

Deterministic control policy: given by neural network.
Stochastic control policy: by Bayesian neural network (BNN).

# Safety Certificates: Inductive Invariants

Inductive ("positive") invariant:

continuous function $S: X \to \mathbb{R}$ such that

(1) $S(x_0) > 0$
(2) for all $x \in X, u \sim \pi(x)$, and $\omega$, if $S(x) > 0$ then $S\big(f(x, u, \omega)\big) > 0$
(3) for all unsafe states $x \in X_U$, $S(x) \leq 0$

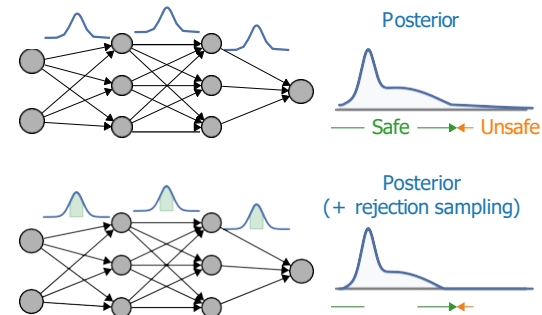If there exists an inductive invariant,
then the closed-loop system is safe.

# Safety Certificates: Inductive Invariants

Inductive ("positive") invariant:

continuous function $S: X \to \mathbb{R}$ such that

(1) $S(x_0) > 0$
(2) for all $x \in X, u \sim \pi(x)$, and $\omega$, if $S(x) > 0$ then $S\big(f(x, u, \omega)\big) > 0$
(3) for all unsafe states $x \in X_U, S(x) \leq 0$



Posterior

—— Safe ——►◄— Unsafe

Posterior
(+ rejection sampling)

If there exists an inductive invariant,
then the closed-loop system is safe.

BNNs are usually not safe (because of weight distributions with unbounded support),
but can be made safe by rejection sampling (cut off support at $\mu \pm \delta$ for mean weight $\mu$).

# Progress Certificates: Ranking Supermartingales

Ranking supermartingale [Chakarov-Sankaranarayanan]:

continuous function $P: X \to \mathbb{R}$ such that

(1) for all $x \in X$, $S(x) \geq 0$
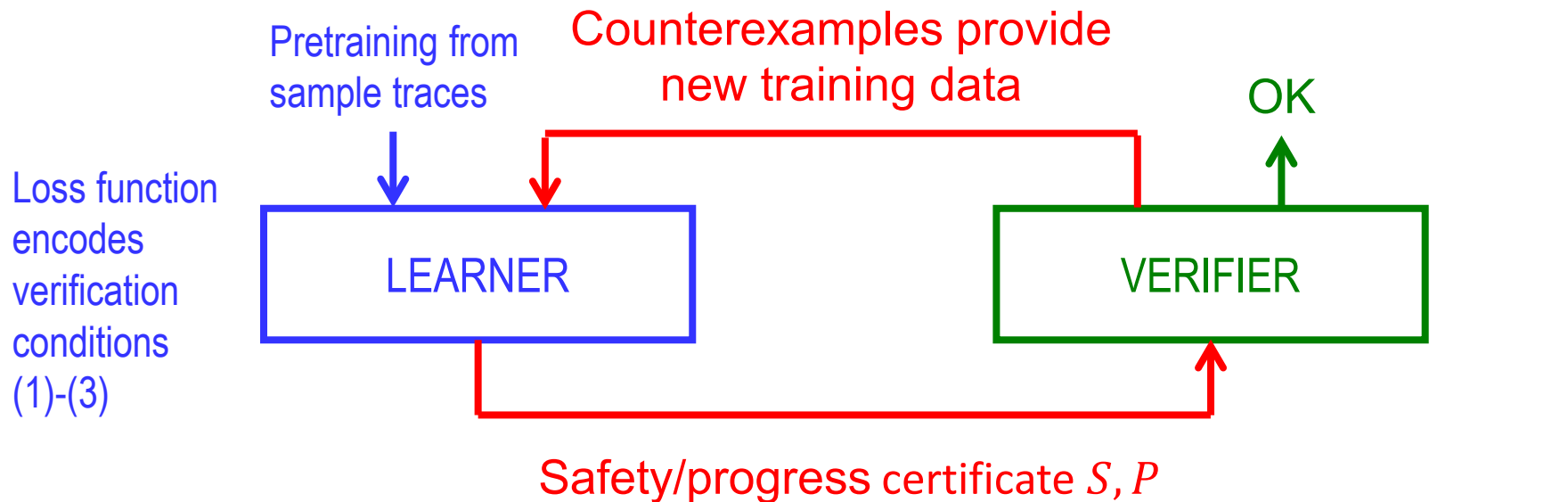(2) there exists $\varepsilon > 0$ such that for all unstable $x \in X \backslash X_S, u \sim \pi(x)$, and $\omega$,

$$\mathbb{E}_{u,\omega}\big[P\big(f(x, u, \omega)\big)\big] \leq P(x) - \varepsilon$$

(3) for all stable $x \in X_S$, $u \sim \pi(x)$, and $\omega$, $f(x, u, \omega) \in X_S$

If there exists a ranking supermartingale,
then the closed-loop system is almost surely stable.
Moreover, we can bound the stabilization time [AAAI'22].
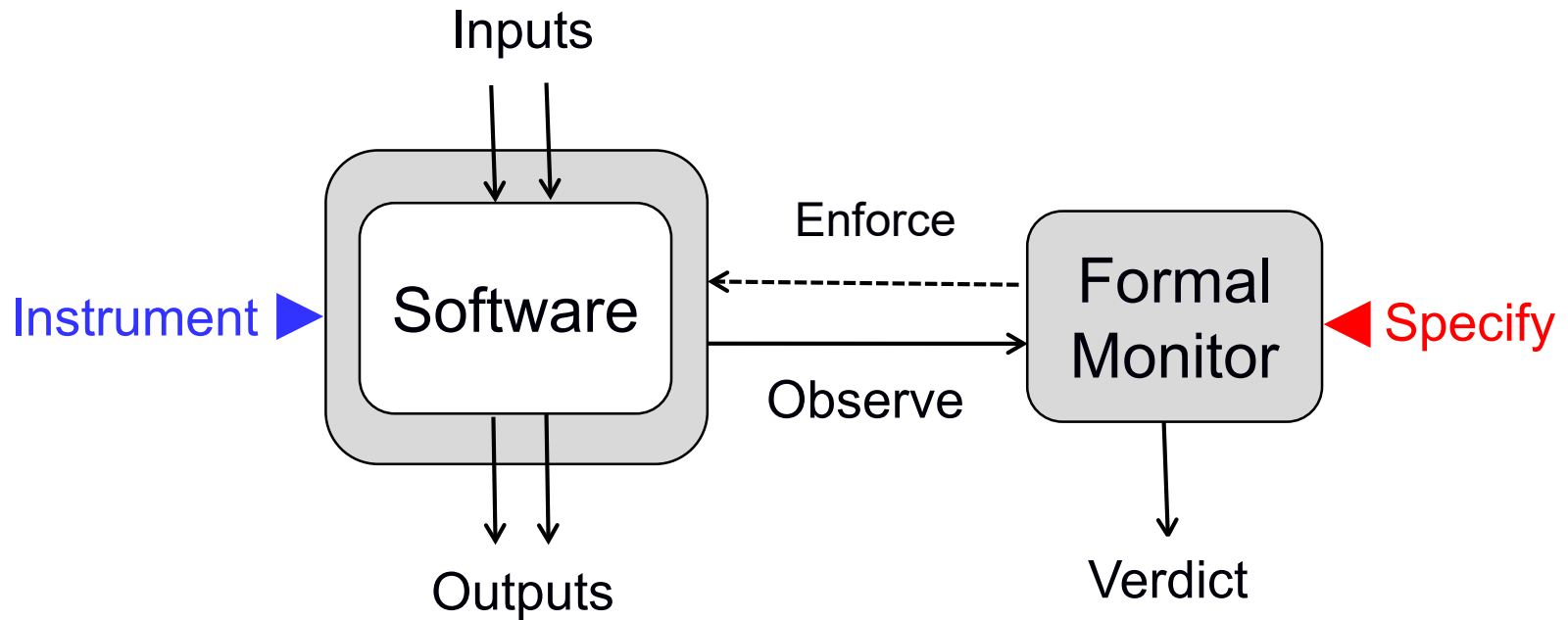
# Learning and Verifying Certificates

Pretraining from sample traces

Counterexamples provide new training data

OK

Loss function encodes verification conditions (1)-(3)

LEARNER

VERIFIER

Safety/progress certificate $S, P$

Constraint solvers:
-Mixed integer linear programming
-Reluplex [Katz et al.]
-Nonlinear real arithmetic

# Software Needs Watchdogs

## Tom Henzinger
### IST Austria

# The Vision:
## Ubiquitous Online Black-Box Monitoring of Software



Inputs

Software

Instrument ▶

Enforce

Observe

Formal Monitor

◀ Specify

Outputs

Verdict

Bounded number of monitor operations
per observation ("real time" [Rabin'63]).
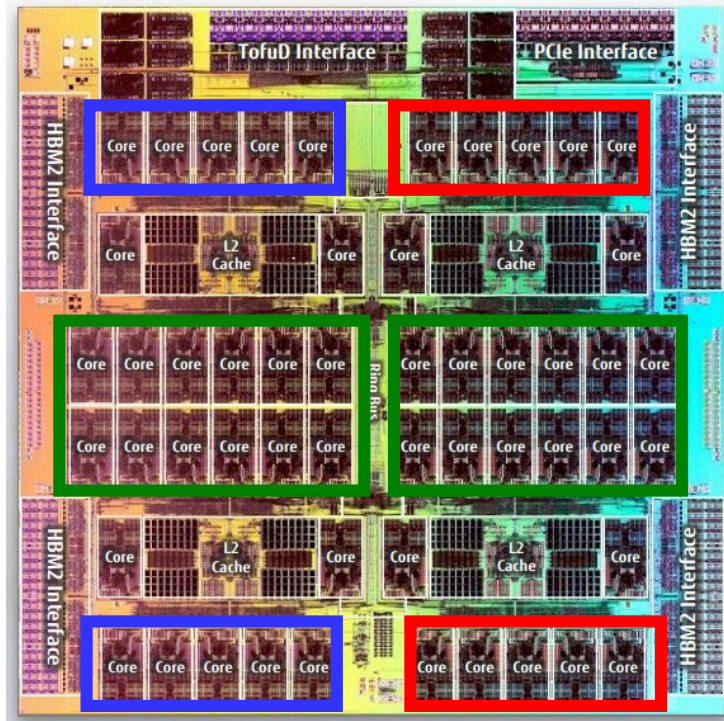
# Two Megatrends that Favor Runtime Verification

## 1. Ever increasing hardware parallelism

□ many-core processors, compute clusters, data centers
□ not all hardware resources will go into performance
  (some *should* go into assurance)

## 2. Ever increasing software complexity

□ *machine-learned components*, cloud connectivity
□ the static "verification gap" will only increase
  (software complexity grows faster than verification capability)
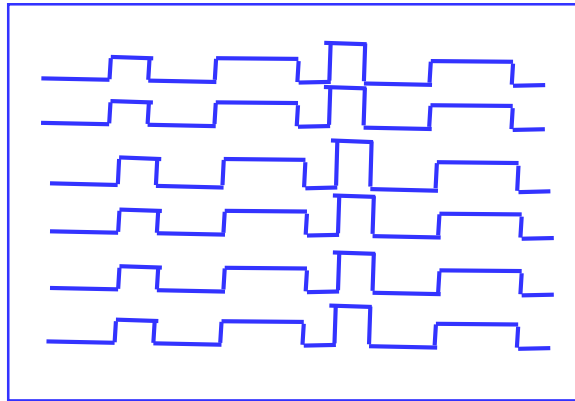
# The Future Use of CPUs and GPUs ?!

# Two Observations about Runtime Verification

1. Static verification is about all possible infinite traces, runtime verification is about one *long* observed trace

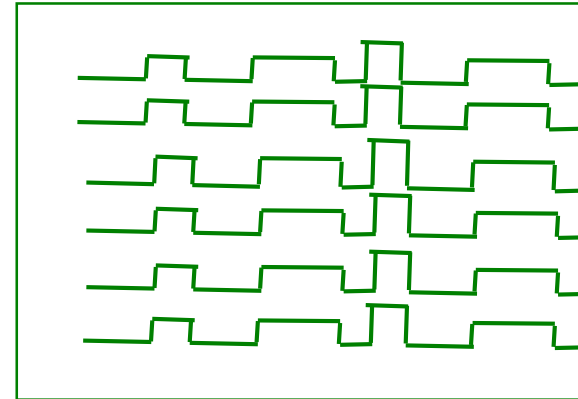&#9633; formalisms and methods should reflect this difference
(the difference between emptiness and membership checking;
over the *long* run, possibility becomes probability)
&#9633; opening for quantitative, statistical, and approximate methods
(distance measures between traces rather than between systems)

# Static versus Runtime Verification

# Two Observations about Runtime Verification

1. Static verification is about all possible infinite traces, runtime verification is about one *long* observed trace

- □ formalisms and methods should reflect this difference
  (the difference between emptiness and membership checking;
  over the *long* run, possibility becomes probability)
- □ opening for quantitative, statistical, and approximate methods
  (distance measures between traces rather than between systems)

2. To increase public acceptance and trust of monitors, they must be third-party + unintrusive + online

- □ third-party means that system and monitor are designed *independently*
- □ unintrusive means *black-box* and *low-overhead*
- □ online means *real-time* and *best-effort*

# Brief History of Boolean Monitoring

Observations                 $\Sigma = \{a, b, t, o\}$

Response                   $\Box(a \rightarrow \Diamond b)$

Trace                     $o\ t\ t\ o\ a\ t\ o\ b\ o\ o\ t\ o\ a\ o\ o\ a\ o\ t\ o\ t\ o\ o\ o\ b\ o\ t\ o \ldots$

Boolean verdict         $\bot\ \bot\ \bot \ldots$

Response is live (every finite trace can be extended to satisfy response).
Response is co-live (every finite trace can be extended to violate response).

Response is not Pnueli-Zaks monitorable
(no finite trace allows a positive or negative verdict).

PZ monitorable $\supseteq$ Boolean combinations of safety properties
[Bauer-Leucker-Schallhart, Falcone-Frnandez-Mounier, Diekert-Leucker].

# Quantitative Monitoring

Observations $\qquad\qquad \Sigma = \{a, b, t, o\}$

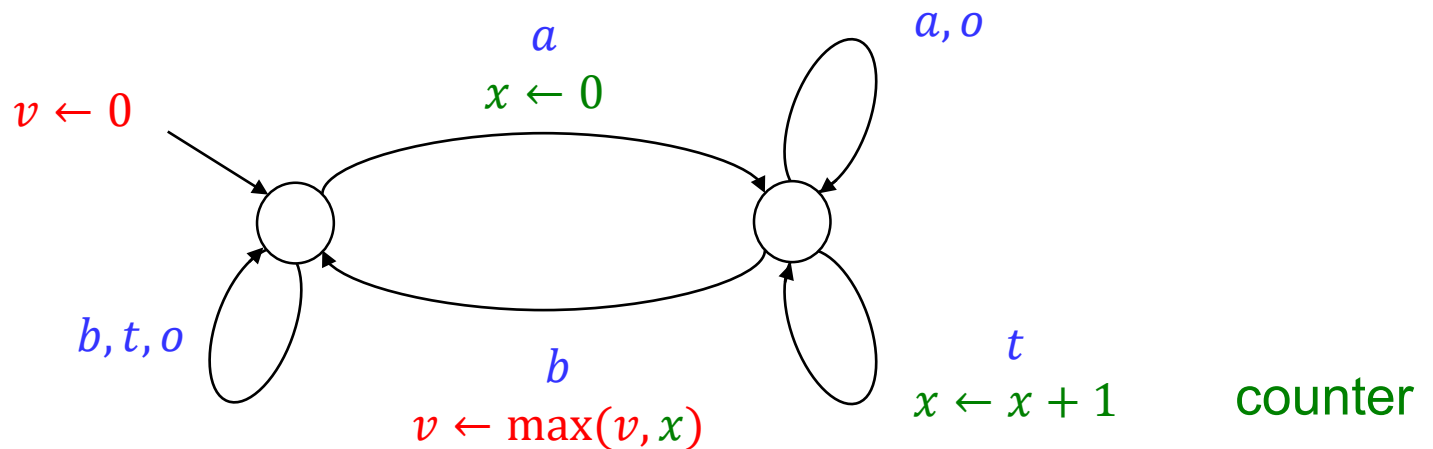Max response time $\qquad min_t \,\square(a \rightarrow \Diamond_{\leq t}\, b)$

Trace $\qquad\qquad\qquad o\,t\,t\,o\,a\,t\,o\,b\,o\,o\,t\,o\,a\,o\,o\,a\,o\,t\,o\,t\,o\,o\,o\,b\,o\,t\,o$ ...

Quantitative verdict $\qquad$ 0 $\qquad$ 1 $\qquad\qquad\qquad\qquad\qquad$ 2
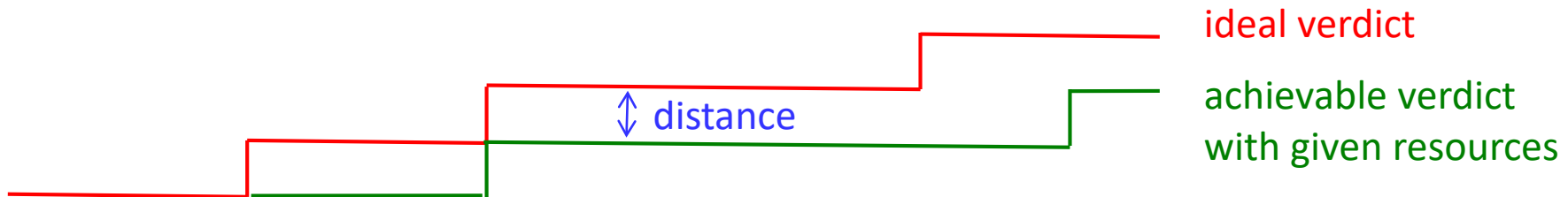
# Quantitative Monitoring

Observations $\qquad \Sigma = \{a, b, t, o\}$

Max response time $\qquad min_t \, \square(a \rightarrow \diamondsuit_{\leq t} \, b)$

Trace $\qquad\qquad\qquad o\ t\ t\ o\ a\ t\ o\ b\ o\ o\ t\ o\ a\ o\ o\ a\ o\ t\ o\ t\ o\ o\ o\ b\ o\ t\ o$ ...

<span style="color:red">Quantitative verdict</span> $\qquad$ <span style="color:red">0 $\qquad\qquad$ 1 $\qquad\qquad\qquad\qquad\qquad\qquad$ 2</span>

<span style="color:green">Approximate verdict</span> $\qquad$ <span style="color:green">0 $\qquad\qquad$ 1 $\qquad\qquad\qquad\qquad\qquad\qquad$ 1</span>

<span style="color:green">Delayed verdict</span> $\qquad\qquad$ <span style="color:green">0 $\qquad\qquad\qquad\qquad\qquad$ 1 $\qquad\qquad\qquad\qquad\qquad$ 2</span>



distance

ideal verdict

achievable verdict
with given resources

# Quantitative Monitoring [LICS'21]

| | |
|---|---|
| Observation alphabet | $\Sigma$ |
| Value cpo | $\Lambda$ |
| Quantitative property | $p: \Sigma^\omega \to \Lambda$ |
| Verdict function | $v: \Sigma^* \to \Lambda$ |
| Infinite trace | $w \in \Sigma^\omega$ |

smallest value such that every larger value is seen only finitely often ("least eventual upper bound")

## Quantitative monitoring can be about limits:

verdict $v$ monitors property $p$ on trace $w$ from below if $limsup_i \{v(w_{0..i})\} = p(w)$.

## Quantitative monitoring can be universal:

property $p$ is monitorable from below if $\exists v : \forall w : v$ monitors $w$ from below.

## Quantitative monitoring can be approximate:

verdict $u$ approximates verdict $v$ with error $\varepsilon$ if $\forall w : \forall i : |v(w_{0..i}) - u(w_{0..i})| \leq \varepsilon$.

# Quantitative Monitoring

Compare monitors with regard to

<span style="color:red">1. precision</span>
<span style="color:blue">2. resource use</span>

◇

# Expressiveness of Counter Monitors

Observations   $\Sigma_k = \{0, \dots, k\}$
Traces         $\Sigma_k^\omega$

*3 3 2 3 3 2 2 1 3 2 1 0 1 2 3 3 2 1 0 0 3 3 …*

Property       $P_k = \{w \in \Sigma_k^\omega : \forall \pi \prec w : \forall i < k : \#(i+1, \pi) \geq \#(i, \pi)\}$
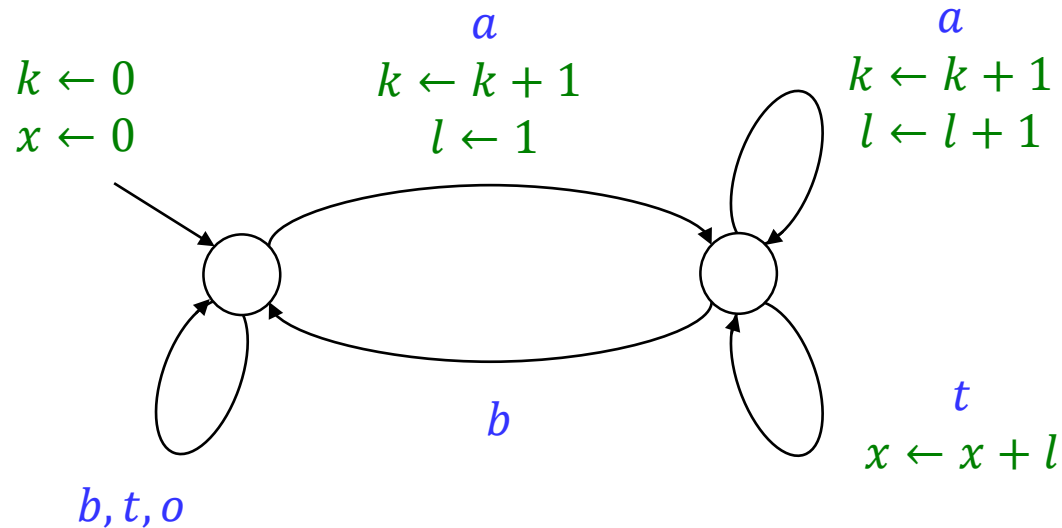
*"Every prefix contains as many 1s as 0s, as many 2s as 1s, as many 3s as 2s, etc."*

Theorem [LICS'18]:

For all $k$, the property $P_k$ can be real-time monitored with $k+1$ counters but not with $k$ counters.

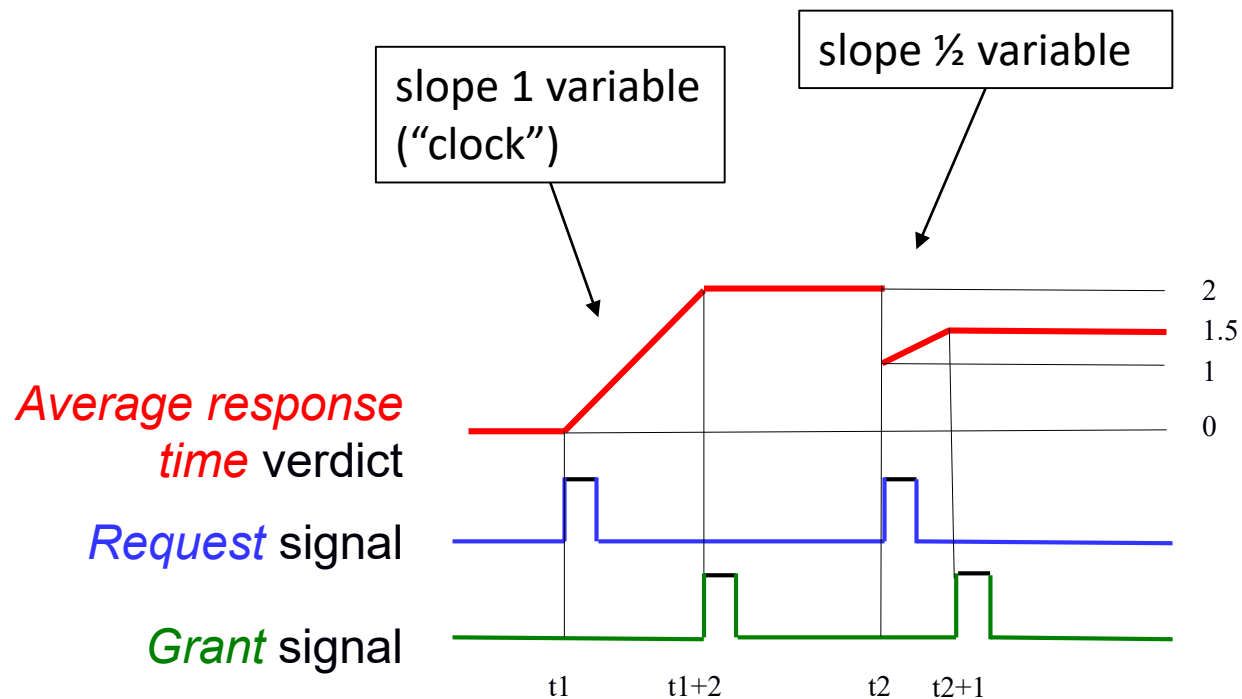(In contrast, in computability every Turing machine can be simulated by two counters.)

# Monitoring Average Response Time



$k \leftarrow 0$
$x \leftarrow 0$

$a$
$k \leftarrow k + 1$
$l \leftarrow 1$

$a$
$k \leftarrow k + 1$
$l \leftarrow l + 1$

$b$

$b, t, o$

$t$
$x \leftarrow x + l$

$v(x, k) = x/k$

Requires addition and division.
Led to original results in static verification [SAS'16].

# Quantitative Monitoring in Continuous Time

# Quantitative Monitoring

Compare monitors with regard to

        1. precision
        2. resource use
        3. strength of assumptions

An assumption $A \subseteq \Sigma^\omega$ restricts the universe of possible traces [RV'20].
How do assumptions arise?

        1. Knowledge about the system (predictive monitoring)
        2. Knowledge about the environment (e.g. $\square \diamond t$ )
        3. Information from other monitors (decentralized monitoring)

# Limit Monitoring [CSL'20]

☐ mode $m(\pi)$ is the most common letter in $\pi \in \Sigma^*$

☐ for $w \in \Sigma^\omega$ generated by finite connected Markov chain (the "assumption"),
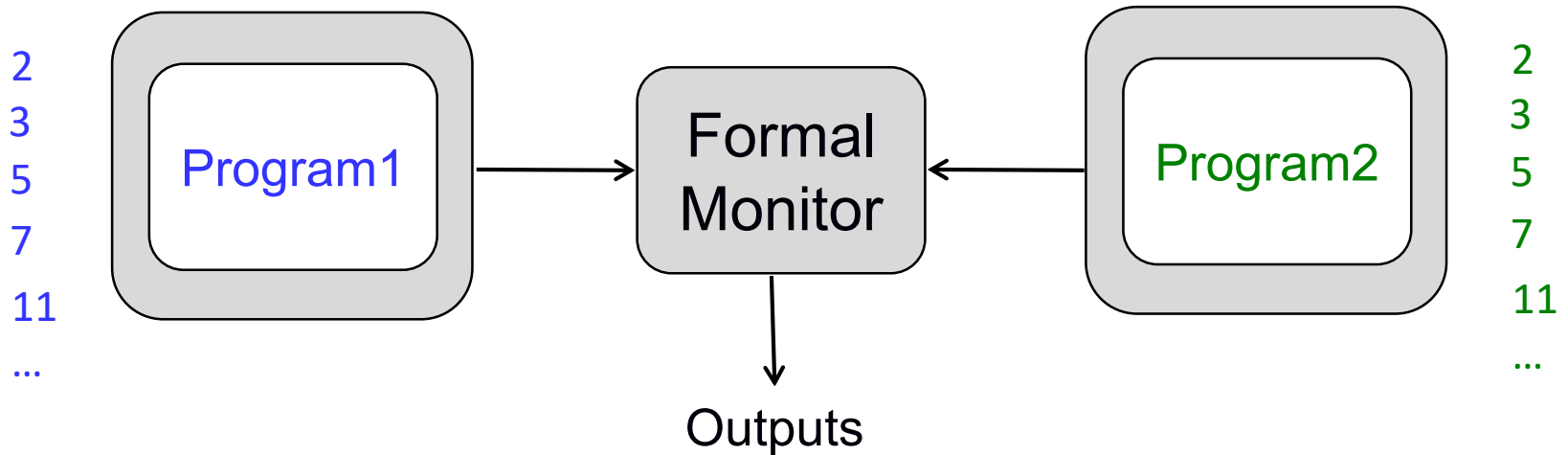$m(w) = \lim_i m(w_{0..i})$ converges with probability 1

**Real-time monitoring the mode: requires |Σ| counters**
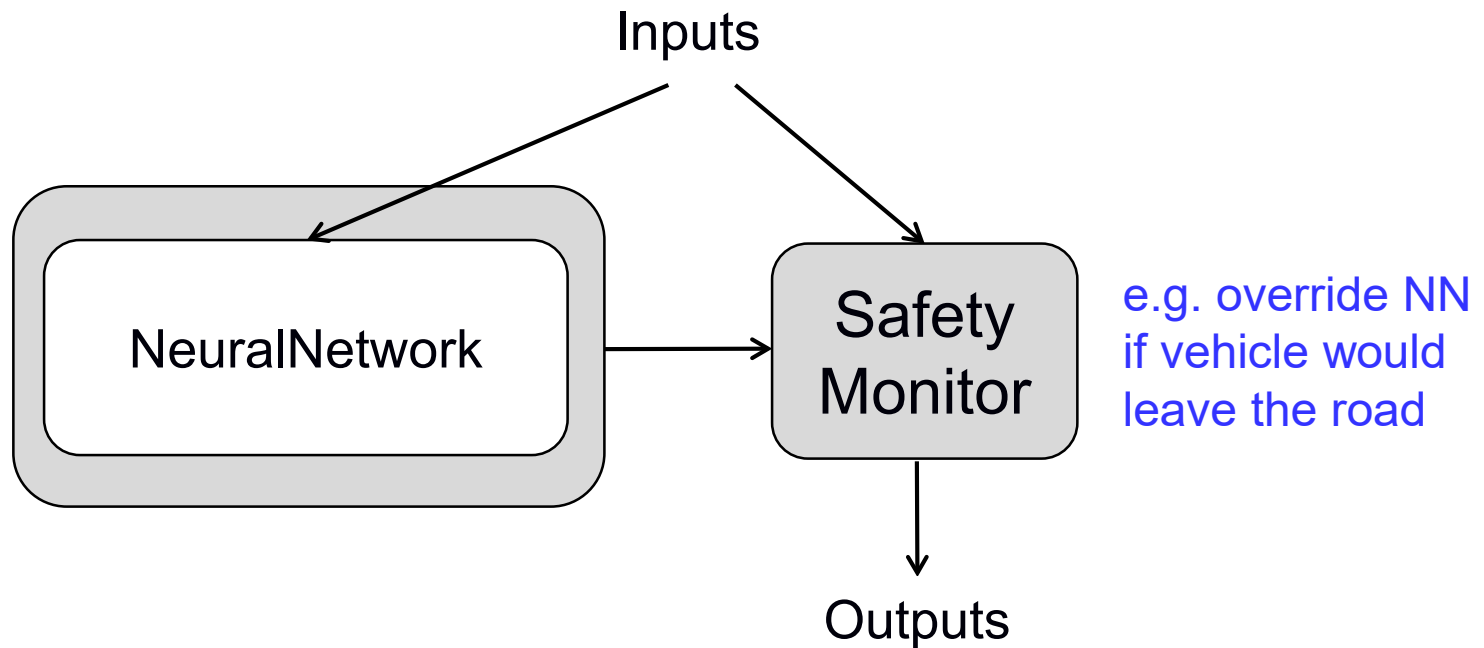
**Limit monitoring the mode: can be done by 4 counters**

$$\alpha_0 \quad \alpha_1 \quad \alpha_2 \quad \alpha_3 \quad \alpha_4 \quad \alpha_5 \quad \alpha_6$$

☐ partition $w \in \Sigma^\omega$ into segments $\alpha_0 \alpha_1 \alpha_2 \alpha_3$ ... of increasing length

☐ while reading $w$, with each $\alpha_i$ maintain a current "hypothesis" $m_i \in \Sigma$ for the mode

☐ cycle round-robin through the letters in Σ, looking at one letter $\sigma_i$ per segment $\alpha_i$

☐ if $\sigma_i$ occurs more often than $m_i$ in $\alpha_i$, then $m_{i+1} = \sigma_i$ else $m_{i+1} = m_i$

# Differential Monitoring [RV'21]

2
3
5
7
11
...

Program1

Formal Monitor

Program2

2
3
5
7
11
...

Outputs
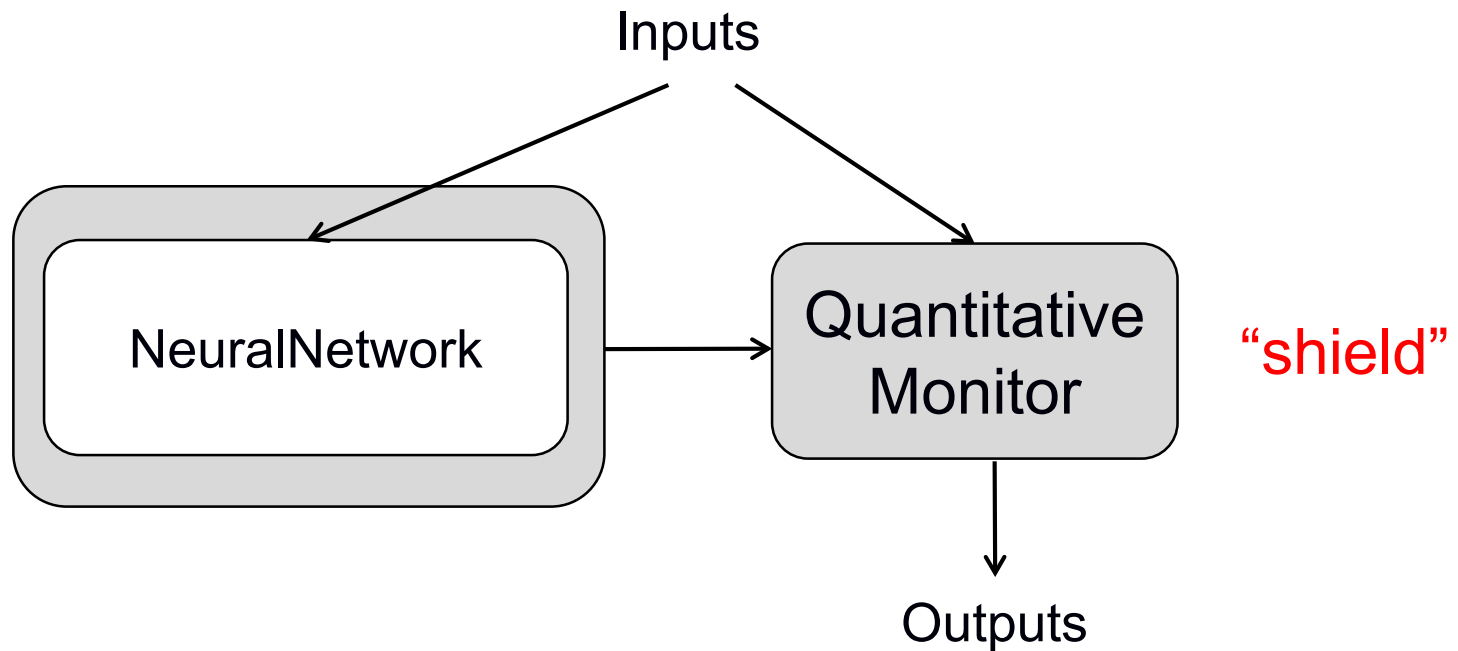
□ assurance through redundancy ("overengineering")
□ programs as specifications

# Autonomous Vehicles: Monitors as Shields



Inputs

NeuralNetwork

Safety Monitor

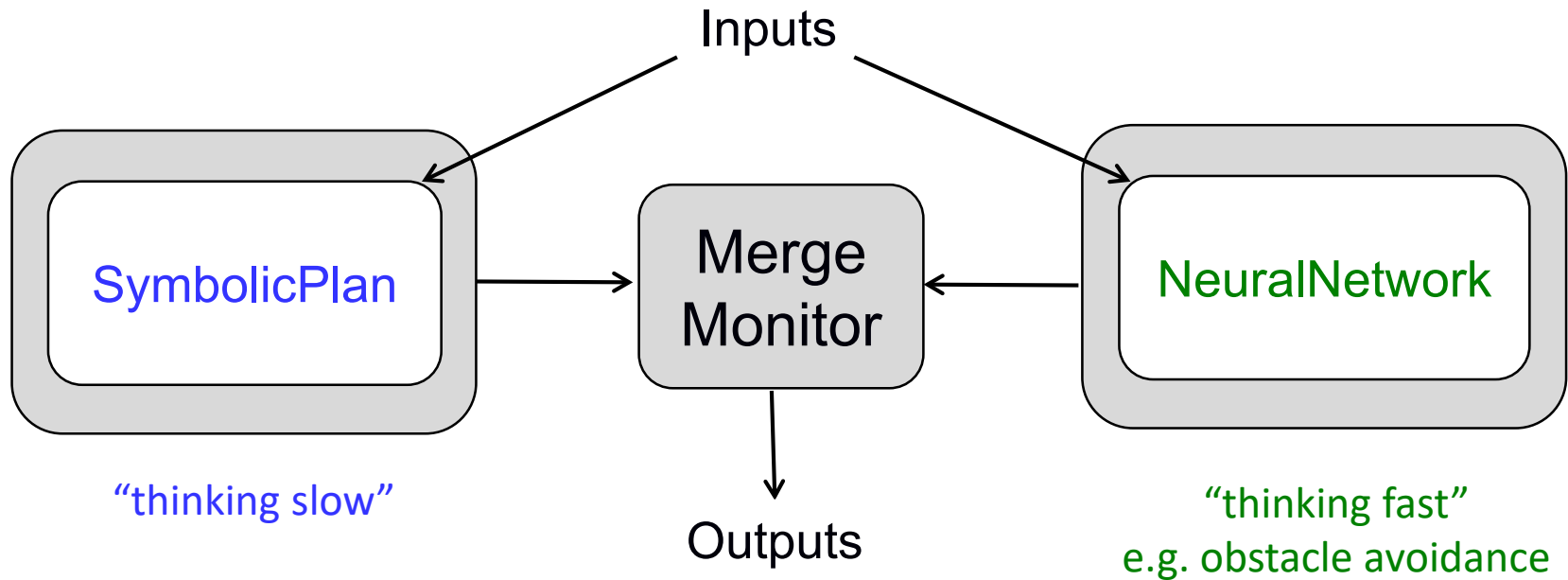e.g. override NN if vehicle would leave the road

Outputs

□ neural network computes control actions
□ monitor overrides actions that violate safety
[Bloem-Koenighofer-Koenighofer-Wang]
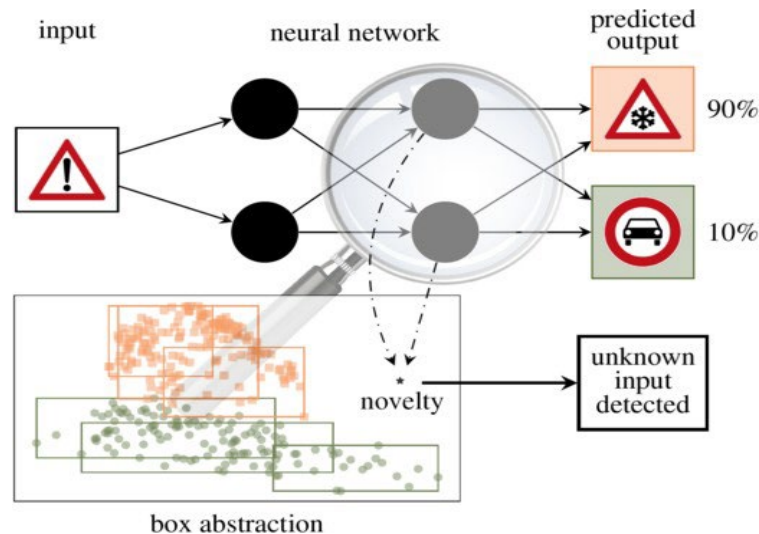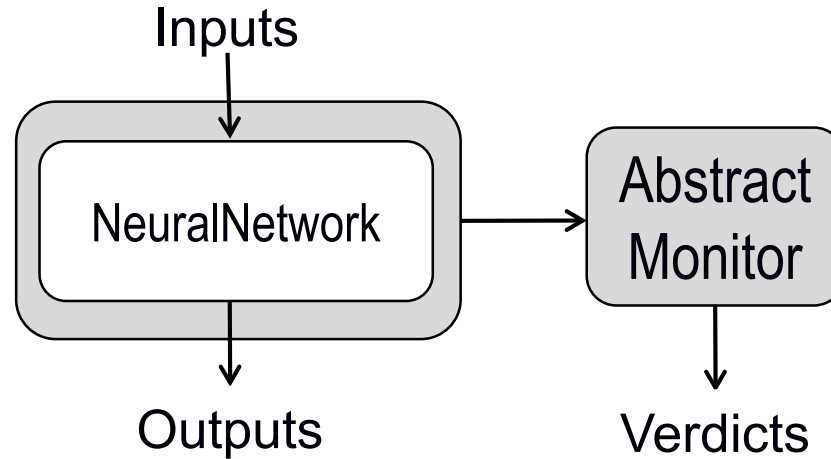
# Beyond Safety: Quantitative Shields [CAV'19]



□ quantitative rewards and interference penalties
□ shield specification using weighted automaton
□ shield synthesis by solving stochastic games (controller+plant against shield)
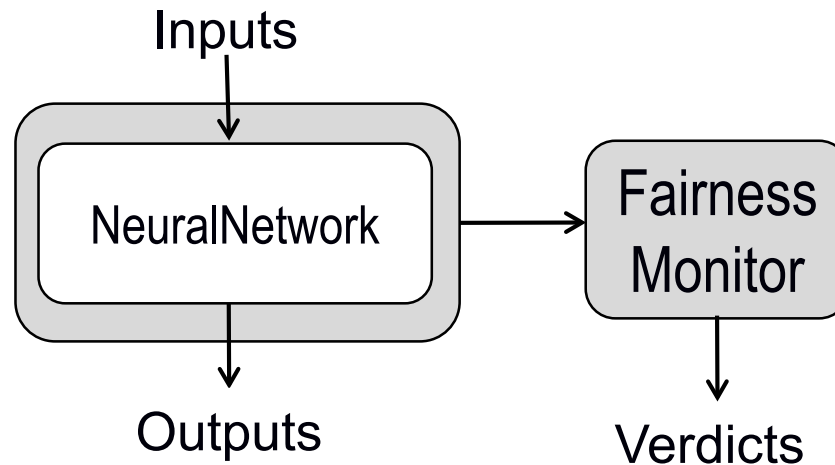
# Planning: Monitors as Arbiters



- □ neural network computes short-term actions
- □ symbolic controller follows long-term plan
- □ monitor arbitrates both decisions

# Classification: Monitors as Novelty Detectors [RV'21]





box abstraction

□ monitor watches if predictions are similar to previously observed patterns

□ when novelties are detected, network may need retraining

# Algorithmic Fairness: Monitors as Watchdogs

Inputs

NeuralNetwork → Fairness Monitor

Outputs

Verdicts

Demographic parity

$$\frac{\mathrm{Pr}(\, output = 1 \mid protectedInputAttribute = 0\,)}{\mathrm{Pr}(\, output = 1 \mid protectedInputAttribute = 1\,)} \cong 1$$

can be monitored using frequency counters.

# Summary

Static formal methods have scalability issues for neural networks.

<span style="color:red">A promising approach is the LEARNER + VERIFIER architecture, which uses machine learning for hypothesis generation of correctness certificates.</span>

<span style="color:blue">Runtime methods can monitor not only safety violations, but also *quality, limits*, and *fairness.*</span>
<span style="color:green">They can enforce not only safety (shielding), but also *progress* (planning).</span>

# References

[NeuIPS'21] Infinite time horizon safety of Bayesian neural networks
[AAAI'22] Stability verification in stochastic control systems via neural network
supermartingales

[TACAS'20] How many bits does it take to quantize your neural network?
[AAAI'21] Scalable verification of quantized neural networks

[CAV'19]   Runtime optimization for learned controllers
[ECAI'20]  Outside the box: Abstraction-based monitoring of neural networks
[RV'21]    Into the unknown: Active monitoring of neural networks