# Sound Verification and Synthesis with Logic and Data

Alessandro Abate

Department of Computer Science

oxcav.web.ox.ac.uk

3 April 2024



$\otimes \subseteq \Delta \vee$

OXFORD CONTROL AND VERIFICATION

*[references at end of deck]*

# Outline

# Outline

# Control theory vs Formal verification

- dynamical models

$$x \in \mathbb{R}^n$$
$$\mathcal{G}_i \subset \mathbb{R}^n, i \in \{1, \ldots m\}$$
$$\forall x \in \mathcal{G}_i, \quad x^+ = f_i(x)$$

- stability,
  safety,
  reachability

- Lyapunov functions,
  barrier certificates,
  reach-set computation

- software programs

```
34: x float
35: ...
36: while Gi(x)
37:    x⁺ := fi(x)
38: endwhile
39: ...
```

- termination,
  assertion violation

- ranking functions,
  program/loop invariants,
  symbolic search

# Cyber-Physical Systems



- complex embedded systems
- interleaving of cyber/digital components with physical/analogue dynamics

- hybrid models
- dynamics, control and computation

  (and communication)

- safety-critical applications
- $\rightarrow$ correct-by-design control
- $\rightarrow$ sound and automated synthesis

# Formal verification in a nutshell

- **industrial impact** in checking the correct behaviour of

  protocols, hardware circuits, and software

# Formal verification in a nutshell

- **industrial impact** in checking the correct behaviour of

  protocols, hardware circuits, and software



- **model-based** algorithms (and SW tools)
- automated, sound, and formal proofs (e.g., via certificates)

# Formal verification in a nutshell

- industrial impact in checking the correct behaviour of

  protocols, hardware circuits, and software



- model-based algorithms (and SW tools)
- automated, sound, and formal proofs (e.g., via certificates)

# Properties: Encoding rich dynamical behaviour

- as specifications, requirements for verification, e.g., safety
- as objectives for control synthesis, e.g., reachability
- <u>without</u> manual reward engineering

# Properties: Encoding rich dynamical behaviour

$$x \in \mathbb{R}^n$$
$$\mathcal{G}_i \subset \mathbb{R}^n, i \in \{1, \dots m\}$$
$$\forall x \in \mathcal{G}_i, \quad x^+ = f_i(x)$$

# Properties: Encoding rich dynamical behaviour

$$x \in \mathbb{R}^n$$

$$x^+ = f(x)$$



$x^+ = f(x)$

# Properties: Encoding rich dynamical behaviour

$$x \in \mathbb{R}^n$$
$$\mathcal{G}_i \subset \mathbb{R}^n, i \in \{1, \dots m\}$$
$$x^+ = f(x)$$



$x^+ = f(x)$      $x^+ = f(x)$

$x:$     $x_0 \rightarrow$      $x_1 \rightarrow$      $x_2 \rightarrow$     $\dots$

$\rho:$     $\mathcal{G}_{x_0} \rightarrow$      $\mathcal{G}_{x_1} \rightarrow$      $\mathcal{G}_{x_2} \rightarrow$     $\dots$

# Properties: Encoding rich dynamical behaviour

- consider (class of) properties/requirements/specifications

$$\forall x_0 \in \mathcal{I}, \quad \exists T \in \mathbb{N}^+, \qquad \forall k \in \{0, 1, \ldots, T-1\}, \qquad \forall \tau \geq T:$$
$$x_T \in \mathcal{G}, \qquad x_k \notin \mathcal{U}, \qquad x_\tau \in \mathcal{F}$$

# Properties: Encoding rich dynamical behaviour

- consider (class of) properties/requirements/specifications

$$\forall x_0 \in \mathcal{I}, \quad \exists T \in \mathbb{N}^+, \quad \forall k \in \{0, 1, \ldots, T-1\}, \quad \forall \tau \geq T :$$
$$x_T \in \mathcal{G}, \quad x_k \notin \mathcal{U}, \quad x_\tau \in \mathcal{F}$$



a: Stability b: ROA    c: Safety      d: SWA      e: Reachability    f: RWA      g: RSWA     h: RAR

- class encompasses stability, invariance, safety, reachability, reach-avoid, . . .

# Properties: Encoding rich dynamical behaviour

- consider (class of) properties/requirements/specifications

$$\forall x_0 \in \mathcal{I}, \quad \exists T \in \mathbb{N}^+, \quad \forall k \in \{0, 1, \ldots, T-1\}, \quad \forall \tau \geq T:$$
$$x_T \in \mathcal{G}, \quad x_k \notin \mathcal{U}, \quad x_\tau \in \mathcal{F}$$



a: Stability  b: ROA    c: Safety    d: SWA    e: Reachability    f: RWA    g: RSWA    h: RAR

- connections to:
  1. automata theory
  2. temporal logics
  3. formal languages

# Outline

1 Why this Matters: Science and Technology Drivers

2 Sound Inductive Synthesis with Neural Certificates

3 Formal Verification with Neural Abstractions

4 Safe and Certified Learning

# Outline

# Decision problems: SAT and SMT

- SAT is a decision problem (yes/no question)
- find satisfying assignment of Boolean functions
- e.g., assume Boolean $x_i$, check

$$\exists x_1, x_2, x_3 : \quad (x_1 \lor \neg x_2) \land (\neg x_1 \lor x_2 \lor x_3) \land \neg x_1$$

# Decision problems: SAT and SMT

- SAT is a decision problem (yes/no question)
- find satisfying assignment of Boolean functions
- e.g., assume Boolean $x_i$, check

$$\exists x_1, x_2, x_3 : \quad (x_1 \vee \neg x_2) \wedge (\neg x_1 \vee x_2 \vee x_3) \wedge \neg x_1$$

- SMT is a decision problem for logical formulae within a theory
- instance: theory of non-linear arithmetics over real closed fields
- e.g., assume reals $x_i \in \mathbb{R}$, check

$$\exists x_1, x_2 : \quad x_1 \geq 0 \Rightarrow 3x_1 + 2x_2 + 1 > 0$$

# From decision to synthesis problems

- consider (harder) problem:
  assume integers $x_i \in \mathbb{Z}$,
  seek function $F : \mathbb{Z} \times \mathbb{Z} \to \mathbb{Z}$, s.t.

$$\exists F, \forall x_1, x_2 :$$

$$F(x_1, x_2) \geq x_1 \wedge F(x_1, x_2) \geq x_2 \wedge (F(x_1, x_2) = x_1 \vee F(x_1, x_2) = x_2)$$

# Lyapunov functions

- consider $\dot{x} = f(x)$, assume $x_e \in \mathbb{R}^n$ is an equilibrium, $f(x_e) = 0$
- ensure asymptotic stability of $x_e$ in $\mathcal{D} \subseteq \mathbb{R}^n$
- by finding Lyapunov function $V(x)$, satisfying

1. lower bound:
$$V(x_e) = 0 \tag{1}$$

2. positive definiteness:
$$V(x) > 0, \ \forall x \in \mathcal{D} \setminus \{x_e\} \tag{2}$$

3. negative Lie derivative:
$$\dot{V}(x) = \nabla V(x) \cdot f(x) < 0, \ \forall x \in \mathcal{D} \setminus \{x_e\} \tag{3}$$

## Lyapunov functions

- consider $\dot{x} = f(x)$, assume $x_e \in \mathbb{R}^n$ is an equilibrium, $f(x_e) = 0$
- ensure asymptotic stability of $x_e$ in $\mathcal{D} \subseteq \mathbb{R}^n$
- by finding Lyapunov function $V(x)$, satisfying

  **1** lower bound:
  $$V(x_e) = 0 \tag{1}$$

  **2** positive definiteness:
  $$V(x) > 0, \ \forall x \in \mathcal{D} \setminus \{x_e\} \tag{2}$$

  **3** negative Lie derivative:
  $$\dot{V}(x) = \nabla V(x) \cdot f(x) < 0, \ \forall x \in \mathcal{D} \setminus \{x_e\} \tag{3}$$

- that is, solve following synthesis problem:

  $$\exists V \colon \mathcal{D} \to \mathbb{R} \quad s.t. \ \ \forall x \in \mathcal{D}, \quad \text{conditions } (1) \wedge (2) \wedge (3) \text{ hold}$$

# Counterexample-guided inductive synthesis (CEGIS)



1. **Learner**

   generates candidates $V$ over finite set

2. **Verifier**

   certifies validity on $\mathcal{D}$, or provides counterexample(s) $c$

# Counterexample-guided inductive synthesis (CEGIS)



$f(x), \mathcal{D}$

1. Learner

   generates candidates $V$ over finite set

2. Verifier

   certifies validity on $\mathcal{D}$, or provides counterexample(s) $c$

- inductive synthesis loop

  1. sample (finite) set $S \subset \mathcal{D}$

  2. Learner generates $V(\theta)$ via query SMT solver on formula:
     $\exists \theta : (1) \wedge (2) \wedge (3)$ on points $s \in S$

  3. Verifier checks either $V(x)$ valid over dense $\mathcal{D}$, or counterexample $c$ :
     query SMT solver on formula $\exists c \in \mathcal{D} : \neg(1) \vee \neg(2) \vee \neg(3)$

  4. $S \leftarrow S \cup c$, loop back to 2

# Counterexample-guided inductive synthesis (CEGIS)



1. Learner

   generates candidates $V$ over finite set

2. Verifier

   certifies validity on $\mathcal{D}$, or provides counterexample(s) $c$

- inductive synthesis loop

  1. sample (finite) set $S \subset \mathcal{D}$

  2. Learner generates $V(\theta)$ via query SMT solver on formula:
     $\exists \theta : (1) \wedge (2) \wedge (3)$ on points $s \in S$

  3. Verifier checks either $V(x)$ valid over dense $\mathcal{D}$, or counterexample $c$ :
     query SMT solver on formula $\exists c \in \mathcal{D} : \neg(1) \vee \neg(2) \vee \neg(3)$

  4. $S \leftarrow S \cup c$, loop back to 2

- sound, but not complete: infinite search space ($\theta$ in $V$) and domain $\mathcal{D}$

# Lyapunov functions as neural networks

- neural nets are general and flexible (universal function approximators)
- Learner trains shallow neural network

$$V(x) = W_2 \cdot \sigma_1(W_1 x + b_1)$$

  ($W_i$ weights, $(\sigma_1)$ activation fcns)



- loss function enforces Lyapunov conditions in (2) and (3) on points in $S$:

$$L(S) = \sum_{s \in S} \max\{0, -V(s)\} + \sum_{s \in S} \max\{0, \dot{V}(s)\}$$

- loss function $L$ is *"pretty good"* proxy of synthesis formula

# Lypunov functions as neural networks



- surprisingly effective! Communication *Learner ↔ Verifier* is crucial

- loss function enforces Lyapunov conditions in (2) and (3) on points in $S$:

$$L(S) = \sum_{s \in S} \max\{0, -V(s)\} + \sum_{s \in S} \max\{0, \dot{V}(s)\}$$

- loss function $L$ is *"pretty good"* proxy of synthesis formula

# Synthesis of Lyapunov functions - example

# Barrier certificates

- consider sets $\mathcal{I}$ (initial) and $\mathcal{U}$ (unsafe)
- ensure there exists no trajectory starting in $\mathcal{I}$ ever entering $\mathcal{U}$

1. negativity within initial set $\mathcal{I}$:

$$B(x) \leq 0 \ \forall x \in \mathcal{I}$$

2. positivity within unsafe set $\mathcal{U}$:

$$B(x) > 0 \ \forall x \in \mathcal{U}$$

3. set invariance property via Lie derivative:

$$\dot{B}(x) < 0 \ \forall x \text{ s.t. } B(x) = 0$$

# Barrier certificates

1. negativity within initial set $\mathcal{I}$:

$$B(x) \leq 0 \ \forall x \in \mathcal{I}$$

2. positivity within unsafe set $\mathcal{U}$:

$$B(x) > 0 \ \forall x \in \mathcal{U}$$

3. set invariance property via Lie derivative:

$$\dot{B}(x) < 0 \ \forall x \text{ s.t. } B(x) = 0$$

# Synthesis of barrier certificates - examples



Barrier Certificate



Barrier Border

$$\begin{cases} \dot{x} & = y + 2xy, \\ \dot{y} & = -x + 2x^2 - y^2 \end{cases}$$

[10] · Linear

# Synthesis of barrier certificates - examples



Barrier Certificate



Barrier Border

$$\begin{cases} \dot{x} & = \exp(-x) + y - 1, \\ \dot{y} & = -\sin(x)^2 \end{cases}$$

[20] · Softplus

# Synthesis of barrier certificates - examples



Barrier Certificate



Barrier Border

$$\begin{cases} \dot{x} & = y, \\ \dot{y} & = -x - y + \frac{1}{3}x^3 \end{cases}$$

$[20, 20] \cdot$ Sigmoid, Sigmoid

# Synthesis of barrier certificates - benchmarks

| Benchmark | CEGIS (this work) | | | | BC[1] | | | SOS[2] | |
|---|---|---|---|---|---|---|---|---|---|
| | Learn | Verify | Samples | Iters | Learn | Verify | Samples | Synth | Verify |
| Darboux | 31.6 | 0.01 | 0.5 k | 2 | 54.9 | 20.8 | 65 k | $\times$ | – |
| Exponential | 15.9 | 0.07 | 1.5 k | 2 | 234.0 | 11.3 | 65 k | $\times$ | – |
| Obstacle | 55.5 | 1.83 | 2.0 k | 9 | 3165.3 | 1003.3 | 2097 k | $\times$ | – |
| Polynomial | 64.5 | 4.20 | 2.3 k | 2 | 1731.0 | 635.3 | 65 k | 8.10 | $\times$ |
| Hybrid mod | 0.58 | 2.01 | 0.5 k | 1 | – | – | – | 12.30 | 0.11 |
| 4-d ODE | 29.31 | 0.07 | 1 k | 1 | – | – | – | 12.90 | OOT |
| 6-d ODE | 89.52 | 1.61 | 1 k | 3 | – | – | – | 16.60 | OOT |
| 8-d ODE | 104.5 | 82.51 | 1 k | 3 | – | – | – | 26.10 | OOT |

- time for Learning and Verification steps in [sec]
- 'Samples' = size of input data for Learner (in thousands)
- 'Iters' = number of iterations of CEGIS loop
- $\times$ = synthesis or verification failure, OOT = verification timeout

---

[1] H. Zhao, X. Zeng, T. Chen, and Z. Liu. Synthesizing Barrier Certificates Using Neural Networks. In Proceedings of the 23rd International Conference on Hybrid Systems: Computation and Control, HSCC, 2020.

[2] A. Papachristodoulou, J. Anderson, G. Valmorbida, S. Prajna, P. Seiler, and P. A. Parrilo. SOSTOOLS: Sum of squares optimization toolbox for MATLAB, 2013.

# Synthesis of control certificates for complex tasks

- dynamical models with inputs (a.k.a., external non-determinism)

$$\dot{x} = f(x, u)$$

$\rightarrow$ synthesis of "control certificates"

- modify known synthesis problem:

$$\exists V \colon \mathcal{D} \to \mathbb{R} \quad s.t. \quad \forall x \in \mathcal{D} \quad \text{conditions } (1) \wedge (2) \wedge (3) \text{ hold}$$

# Synthesis of control certificates for complex tasks

- dynamical models with inputs (a.k.a., external non-determinism)

$$\dot{x} = f(x, u)$$

$\rightarrow$ synthesis of "control certificates"

- approach:
    1. control policies are NN-templated
    2. concurrent synthesis controls & certificates

# Synthesis of control certificates for complex tasks

- dynamical models with inputs (a.k.a., external non-determinism)

$$\dot{x} = f(x, u)$$

$\rightarrow$ synthesis of "control certificates"

- (back to) broad class of properties/requirements

$$\forall x_0 \in \mathcal{I}, \quad \exists T \in \mathbb{N}^+, \quad \forall t \in \{0, \ldots, T-1\}, \quad \forall \tau \geq T :$$
$$x_T \in \mathcal{G}, \quad x_t \notin \mathcal{U}, \quad x_\tau \in \mathcal{F}$$



a: Stability b: ROA    c: Safety      d: SWA      e: Reachability      f: RWA      g: RSWA      h: RAR

# Synthesis of control certificates for complex tasks

- dynamical models with inputs (a.k.a., external non-determinism)

$$\dot{x} = f(x, u)$$

$\rightarrow$ synthesis of "control certificates"

| | $N_s$ | $N_u$ | Property | Neurons | Activations | $T$ (s) | | | Success (%) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | min | $\mu$ | max | S |
| 1 | 2 | 0 | Stability | [6] | $[\varphi_2]$ | 0.01 ($\approx 0.00$) | 0.16 (0.15) | 1.50 (1.48) | 100 |
| 2 | 3 | 0 | Stability | [8] | $[\varphi_2]$ | 0.28 ($\approx 0.00$) | 2.22 (0.45) | 12.57 (3.31) | 100 |
| 3 | 2 | 2 | Stability | [4] | $[\varphi_2]$ | 0.07 (0.01) | 0.19 (0.02) | 0.47 (0.04) | 100 |
| 4 | 2 | 2 | Stability | [5] | $[\varphi_2]$ | 0.09 (0.01) | 0.26 (0.02) | 0.54 (0.03) | 100 |
| 5 | 2 | 0 | ROA | [5] | $[\sigma_{\text{soft}}]$ | 0.21 (0.12) | 14.09 (12.59) | 25.32 (22.13) | 40 |
| 6 | 3 | 3 | ROA | [8] | $[\varphi_2]$ | 1.24 (0.02) | 39.08 (0.03) | 287.89 (0.04) | 100 |
| 7 | 2 | 0 | Safety | [15] | $[\sigma_{\text{t}}]$ | 0.44 (0.35) | 3.36 (2.90) | 7.61 (7.11) | 100 |
| 9 | 8 | 0 | Safety | [10] | $[\varphi_1]$ | 12.63 (7.71) | 51.97 (32.75) | 70.59 (44.66) | 70 |
| 10 | 3 | 1 | Safety | [15] | $[\sigma_{\text{t}}]$ | 1.57 (0.19) | 11.87 (2.50) | 51.08 (7.52) | 90 |
| 11 | 3 | 0 | SWA | [6], [5] | $[\varphi_2], [\sigma_{\text{t}}]$ | 0.19 (0.05) | 2.46 (0.100) | 12.10 (0.20) | 90 |
| 12 | 3 | 0 | SWA | [5], [5, 5] | $[\varphi_2], [\sigma_{\text{sig}}, \varphi_2]$ | 0.13 (0.06) | 0.27 (0.14) | 0.39 (0.20) | 100 |
| 13 | 2 | 1 | SWA | [8], [5] | $[\varphi_2], [\varphi_2]$ | 0.06 (0.03) | 0.20 (0.10) | 0.58 (0.24) | 90 |
| 14 | 3 | 1 | SWA | [10], [8] | $[\varphi_2], [\sigma_{\text{t}}]$ | 4.06 (0.87) | 19.81 (2.73) | 103.49 (7.23) | 90 |
| 15 | 2 | 0 | RWA | [4] | $[\varphi_2]$ | 0.14 (0.09) | 1.81 (1.75) | 4.70 (4.63) | 100 |
| 16 | 3 | 0 | RWA | [16] | $[\varphi_2]$ | 1.36 (0.09) | 14.10 (0.14) | 72.97 (0.20) | 90 |
| 17 | 2 | 1 | RWA | [4, 4] | $[\sigma_{\text{sig}}, \varphi_2]$ | 0.59 (0.27) | 6.82 (3.32) | 20.07 (11.46) | 100 |
| 18 | 3 | 1 | RWA | [5] | $[\varphi_2]$ | 0.46 (0.11) | 16.06 (5.81) | 72.47 (44.64) | 80 |
| 19 | 2 | 2 | RWA | [5] | $[\sigma_{\text{sig}}]$ | 0.69 (0.40) | 1.38 (0.94) | 2.14 (1.90) | 100 |
| 20 | 2 | 0 | RSWA | [4] | $[\varphi_2]$ | 0.19 (0.03) | 1.29 (1.04) | 3.79 (3.37) | 100 |
| 21 | 3 | 0 | RSWA | [16] | $[\varphi_2]$ | 4.81 (0.13) | 27.14 (0.19) | 80.95 (0.25) | 100 |
| 22 | 2 | 0 | RSWA | [5, 5] | $[\sigma_{\text{sig}}, \varphi_2]$ | 1.52 (0.06) | 4.45 (0.19) | 10.97 (0.35) | 100 |
| 23 | 2 | 1 | RSWA | [8] | $[\varphi_2]$ | 0.21 (0.05) | 0.67 (0.25) | 1.19 (0.91) | 100 |
| 24 | 2 | 2 | RSWA | [5, 5] | $[\sigma_{\text{sig}}, \varphi_2]$ | 0.98 (0.16) | 1.23 (0.28) | 1.61 (0.46) | 100 |
| 25 | 2 | 0 | RAR | [6], [6] | $[\sigma_{\text{soft}}], [\varphi_2]$ | 6.65 (1.08) | 24.74 (6.46) | 77.80 (15.06) | 100 |
| 26 | 2 | 2 | RAR | [6, 6], [6, 6] | $[\sigma_{\text{sig}}, \varphi_2], [\sigma_{\text{sig}}, \varphi_2]$ | 5.13 (1.34) | 26.99 (9.90) | 101.23 (60.14) | 100 |

# Synthesis of control certificates for complex tasks

- dynamical models with inputs (a.k.a., external non-determinism)

$$\dot{x} = f(x, u)$$

$\rightarrow$ synthesis of "control certificates"



ROA for NL model,
non-poly Lyapunov,
2 disjoint initial sets

RWA: reach-while-avoid

RAR certificate for
closed-loop NL model

dashed lines: level sets; dark blue: $\mathcal{I}$; light blue: $\mathcal{S}$; green: $\mathcal{G}$; orange: $\mathcal{F}$

UNIVERSITY OF OXFORD



github.com/oxford-oxcav/fossil

# Extension: discrete-time, prob. programs/models

- discrete-time models (e.g. SW programs)

$$\texttt{while g(x), } x^+ \texttt{ := f(x)}$$

$\rightarrow$ similar Lyapunov-like conditions, except concerning "next step":

$$V(f(x)) < V(x), \quad \forall x \in \mathcal{D} \setminus \{x_e\}$$

- stochastic models:

$$x^+ = f(x) + \sigma(x), \quad \sigma \sim \mathcal{N}(0, \Sigma(x))$$

$\rightarrow$ same story, "next step"-condition in expectation (super-martingale):

$$\mathbb{E}[V(f(x)) \mid x] < V(x), \quad \forall x \in \mathcal{D} \setminus \{x_e\}$$

# Outline

# Outline

# Formal abstractions

| complex model | specification |
|---|---|

# Formal abstractions

$\zeta$-quantitative
abstraction

complex
model    specification

# Formal abstractions

# Formal abstractions

# Formal abstractions

# Formal abstractions

SAT,
model
checking

abstract
model | $\tilde{\zeta}$-specification

↓
automated
verification

⟶

↑ $\tilde{\zeta}$-quantitative
abstraction

complex
model | specification

# Formal abstractions

SAT,
model
checking

↓

automated
verification

| abstract model | $\tilde\zeta$-specification |

→

$\tilde\zeta$-spec holds,
policy $\mu_{\tilde\zeta} \models \tilde\zeta$-spec

↑

$\tilde\zeta$-quantitative
abstraction

| complex model | specification |

# Formal abstractions



SAT,
model
checking

↓

automated
verification

abstract
model | $\tilde{\zeta}$-specification

$\tilde{\zeta}$-spec holds,
policy $\mu_{\tilde{\zeta}} \models \tilde{\zeta}$-spec

$\tilde{\zeta}$-quantitative
abstraction

refine back

complex
model | specification

# Formal abstractions



SAT,
model
checking

↓

abstract model | $\tilde{\zeta}$-specification

automated verification

$\tilde{\zeta}$-spec holds,
policy $\mu_{\tilde{\zeta}} \models \tilde{\zeta}$-spec

$\tilde{\zeta}$-quantitative abstraction

refine back

complex model | specification

spec holds,
policy $\mu \models$ spec

# Formal abstractions

# From uncountable to finite stochastic models

∞-space Markov process

$s \in \mathbb{R}^n$

$s^+ = f(s) + \sigma(s), \quad \sigma \sim \mathcal{N}(0, \Sigma(s))$

finite-space Markov chain

$\{z_1, z_2, z_3, \ldots, z_p\}$

$$\mathbb{T} = \begin{bmatrix} p_{11} & \cdots & p_{1p} \\ \cdots & \cdots & \cdots \\ p_{p1} & \cdots & \cdots \end{bmatrix}$$

# From uncountable to finite stochastic models

# From uncountable to finite stochastic models

# From uncountable to finite stochastic models



- error $\xi \sim h_s \delta T$, where
  - $\delta$ - max diameter of partitions
  - $T$ - time horizon
  - $h_s$ - local kernel stiffness (Lipschitz constant)

# From uncountable to finite stochastic models

- error $\xi \sim h_s \delta T$, where
    - $\delta$ - max diameter of partitions
    - $T$ - time horizon
    - $h_s$ - local kernel stiffness (Lipschitz constant)

- probabilistic safety:
    - prob. $p_s$ that execution, started at $s \in \mathcal{I}$, stays in set $A = \mathcal{U}^c$ within $[0, T]$,

# From uncountable to finite stochastic models



- error $\xi \sim h_s \delta T$, where

    $\delta$ - max diameter of partitions

    $T$ - time horizon

    $h_s$ - local kernel stiffness (Lipschitz constant)

- probabilistic safety:

    *prob. $p_s$ that execution, started at $s \in \mathcal{I}$, stays in set $A = \mathcal{U}^c$ within $[0, T]$,*
    *can be computed on abstract model as $\tilde{p}_z$, so that $p_s = \tilde{p}_z \pm \xi$*

# Software for formal abstractions - FAUST$^2$ & StocHy



- sequential, adaptive, anytime

a: Stability  b: ROA  c: Safety  d: SWA  e: Reachability  f: RWA  g: RSWA  h: RAR

# Software for formal abstractions - FAUST$^2$ & StocHy

- sequential, adaptive, anytime



a: Stability  b: ROA    c: Safety      d: SWA      e: Reachability    f: RWA      g: RSWA      h: RAR

# Software for formal abstractions - FAUST$^2$ & StocHy

- sequential, adaptive, anytime



a: Stability  b: ROA   c: Safety   d: SWA   e: Reachability   f: RWA   g: RSWA   h: RAR

# Software for formal abstractions - FAUST$^2$ & StocHy

- sequential, adaptive, anytime



a: Stability   b: ROA   c: Safety   d: SWA   e: Reachability   f: RWA   g: RSWA   h: RAR

# Software for formal abstractions - FAUST$^2$ & StocHy

- sequential, adaptive, anytime



a: Stability  b: ROA    c: Safety    d: SWA    e: Reachability    f: RWA    g: RSWA    h: RAR

# Software for formal abstractions - FAUST[2] & StocHy

- sequential, adaptive, anytime



`sourceforge.net/projects/faust2`

`gitlab.com/natchi92/StocHy`

- numerous extensions and applications
- wide ecosystem of SHS abstractions
- annual ARCH competition
  `cps-vo.org/group/ARCH`

Alessandro Abate[1], Henk Blom[2], Marc Bouissou[3], Nathalie Cauchi[1], Hassane Chraibi[3], Joanna Delicaris[4], Sofie Haesaert[5], Arnd Hartmanns[6], Mahmoud Khaled[7], Abolfazl Lavaei[8], Hao Ma[4], Kaushik Mallik[8], Mathis Niehage[4], Anne Remke[4], Stefan Schupp[10], Fedor Shmarov[11], Sadegh Soudjani[12], Adam Thorpe[13], Vlad Turcuman[1], and Paolo Zuliani[12]

[1] University of Oxford, Oxford, UK
[2] Delft University of Technology, Delft, The Netherlands
[3] R&D Division of Electricité de France (EDF), France
[4] University of Münster, Germany
[5] TU Eindhoven, Eindhoven, The Netherlands
[6] University of Twente, Enschede, The Netherlands
[7] Technical University of Munich, Germany
[8] ETH Zürich, Germany
[9] Max Planck Institute for Software Systems, Germany
[10] TU Wien, Vienna, Austria
[11] University of Manchester, Manchester, UK
[12] Newcastle University, Newcastle upon Tyne, UK
[13] University of New Mexico, USA

**Abstract**

This report presents the results of a friendly competition for formal verification and policy synthesis of stochastic models. It also introduces new benchmarks within this category, and recommends next steps for this category towards next year's edition of the competition. The friendly competition took place as part of the workshop Applied Verification for Continuous and Hybrid Systems (ARCH) in Spring/Summer 2021.

# Model hybridisations

- safety verification of non-linear models $\dot{x} = f(x)$ over $x \in \mathcal{X} \subset \mathbb{R}^n$,
- it is in general <u>hard</u> - not automated, not scalable



$$\begin{cases} \dot{x} & = -y - 1.5x^2 - 0.5x^3 - 0.5 \\ \dot{y} & = 3x - y \end{cases}$$



$$\begin{cases} \dot{x} & = x^2 + y \\ \dot{y} & = \sqrt[3]{x^2} - x \end{cases}$$

$$\mathcal{X} = [-1, 1]^2$$

# Model hybridisations

- safety verification of non-linear models $\dot{x} = f(x)$ over $x \in \mathcal{X} \subset \mathbb{R}^n$,
- it is in general <u>hard</u> - not automated, not scalable

- leverage <u>formal abstractions</u> (simulations) for verification

- abstraction as hybridisation:
  partition $\mathcal{X}$, locally approximate $f(x)$ as $\tilde{f}(x)$
  each partition has own flow $\tilde{f}(x)$ & transitions to other partitions

# Model hybridisations

- safety verification of non-linear models $\dot{x} = f(x)$ over $x \in \mathcal{X} \subset \mathbb{R}^n$,
- it is in general <u>hard</u> - not automated, not scalable

- leverage <u>formal abstractions</u> (simulations) for verification

- abstraction as hybridisation:

  partition $\mathcal{X}$, locally approximate $f(x)$ as $\tilde{f}(x)$

  each partition has own flow $\tilde{f}(x)$ & transitions to other partitions
- compute upper-bound $\xi$ to error; obtain simulation as

$$\dot{x} = \tilde{f}(x) + d, \quad \|d\| \le \xi, \quad x \in \mathcal{X}$$

# Model hybridisations

- safety verification of non-linear models $\dot{x} = f(x)$ over $x \in \mathcal{X} \subset \mathbb{R}^n$,
- it is in general <u>hard</u> - not automated, not scalable

- leverage <u>formal abstractions</u> (simulations) for verification

- abstraction as hybridisation:
  partition $\mathcal{X}$, locally approximate $f(x)$ as $\tilde{f}(x)$
  each partition has own flow $\tilde{f}(x)$ & transitions to other partitions
- compute upper-bound $\xi$ to error; obtain simulation as

$$\dot{x} = \tilde{f}(x) + d, \quad \|d\| \leq \xi, \quad x \in \mathcal{X}$$

- more partitions $\rightarrow$ larger abstraction
- ! mesh size & shape important for small error bound $\xi$

# Model hybridisations as neural abstractions



- neural network $\mathcal{N}$ as abstraction $\tilde{f}$ of nonlinear vector field $f$
- $\mathcal{N}(x) : \mathbb{R}^n \to \mathbb{R}^n$ approximates $f(x)$
- $H$ neurons $\to$ at most $2^H$ total partitions

# Model hybridisations as neural abstractions

- synthesis of neural abstractions via CEGIS
  1. learn parameters of NN $\mathcal{N}$ w/ MSE loss $\mathcal{L} = \|f(S) - \mathcal{N}(S)\|$, $S$ finite

  2. SMT solver formally checks upper bound $\xi$ on approximation error:

  $$\exists c \in \mathcal{X} \;\; s.t. \; \|f(c) - \mathcal{N}(c)\| > \xi$$

$$\begin{cases} \dot{x} & = y - 1.5x^2 - 0.5x^3 \\ \dot{y} & = 3x - y \end{cases}$$



$$\begin{cases} \dot{x} & = \exp(-x) + y - 1 \\ \dot{y} & = -\sin(x)^2 \end{cases}$$



$$\begin{cases} \dot{x} & = x^2 + y \\ \dot{y} & = \sqrt[3]{x^2} - x \end{cases}$$

$$\begin{cases} \dot{x} &= y - 1.5x^2 - 0.5x^3 \\ \dot{y} &= 3x - y \end{cases}$$

$$\begin{cases} \dot{x} &= \exp(-x) + y - 1 \\ \dot{y} &= -\sin(x)^2 \end{cases}$$

$$\begin{cases} \dot{x} &= x^2 + y \\ \dot{y} &= \sqrt[3]{x^2} - x \end{cases}$$

# Model hybridisations as neural abstractions - examples

$$\begin{cases} \dot{x} & = y - 1.5x^2 - 0.5x^3 \\ \dot{y} & = 3x - y \end{cases}$$

$$\begin{cases} \dot{x} & = \exp(-x) + y - 1 \\ \dot{y} & = -\sin(x)^2 \end{cases}$$

$$\begin{cases} \dot{x} & = x^2 + y \\ \dot{y} & = \sqrt[3]{x^2} - x \end{cases}$$

$$\begin{cases} \dot{x} & = y - 1.5x^2 - 0.5x^3 \\ \dot{y} & = 3x - y \end{cases}$$



$$\begin{cases} \dot{x} & = \exp(-x) + y - 1 \\ \dot{y} & = -\sin(x)^2 \end{cases}$$



$$\begin{cases} \dot{x} & = x^2 + y \\ \dot{y} & = \sqrt[3]{x^2} - x \end{cases}$$

# Safety verification via neural abstractions

# Neural abstractions: alternative templates

| | Piecewise constant | Piecewise affine | Nonlinear |
|---|---|---|---|
| **Concrete model** | Nonlinear, non-Lipschitz | Nonlinear, non-Lipschitz | Nonlinear, non-Lipschitz |
| **Activation functions** |  |  |  |
| **Training procedure** | Particle swarm | Gradient descent | Gradient descent |
| **Loss function** | $\max_{s \in S} l^{\infty}(s)$ | $\frac{1}{|S|} \sum_{s \in S} l^2(s)$ | $\frac{1}{|S|} \sum_{s \in S} l^2(s)$ |
| **Abstract model** | PWA with disturbance | PWC with disturbance | NL-ODE with disturbance |
| **Safety verification tech** | Symbolic model checking | Reach algorithm | Flowpipe propagation (Taylor models) |
| **Safety verification tool** | PHAVer | SpaceEx | Flow* |

# Neural abstractions: alternative templates



(a) Neural PWC abstraction

(b) Neural PWA abstraction

(c) Sigmoidal abstraction.

(a) Flowpipes for neural PWC model. 11.6s

(b) Flowpipe for neural PWA model. 76.5s

(c) Flowpipe for sigmoidal model. 1084.3s

# Outline

1 Why this Matters: Science and Technology Drivers

2 Sound Inductive Synthesis with Neural Certificates

3 Formal Verification with Neural Abstractions

4 Safe and Certified Learning

# Outline

# Reinforcement learning

- learning algorithm, relies on reward signal from environment
- synthesises policies (actions) maximising cumulative reward

# Reinforcement learning

- **learning** algorithm, relies on **reward signal** from environment
- synthesises **policies** (actions) maximising **cumulative reward**



- rewards are **not enough**!
- verification goal: **certified** synthesis of **policies** satisfying **requirement, task**

# Certified reinforcement learning: LCRL

**IN** requirement, task
- encode task, e.g. temporal requirement in LTL formula, as automaton
- synchronise automaton with environment                    % via labels
- synthesise policies via RL    % automaton guides/rewards exploration

**OUT** certified policies: max probability of task satisfaction

# Certified reinforcement learning: LCRL



- model-free → extracts information efficiently
- guided learning → faster convergence, high-dimensional environments
- flexible → numerous extensions and applications

# Ambiguity and Misspecification in Inverse RL



- inverse RL: from expert behaviour to rewards
- preference elicitation and alignment
- formalising reward learning with
  1. invariances

# Ambiguity and Misspecification in Inverse RL



- inverse RL: from expert behaviour to rewards
- preference elicitation and alignment
- formalising reward learning with
  1. invariances
  2. metrics

*AAAI23 BPA*

Thank you for your attention

`oxcav.web.ox.ac.uk`

All images used are under Wikimedia CCAS license, or by author

**Selected References on Sound Neural Synthesis**

A. Abate, M. Giacobbe, and D. Roy, "Stochastic Omega-Regular Verification and Control with Supermartingales," CAV24, In Press, 2024.

A. Edwards, A. Peruffo and A. Abate, "A General Verification Framework for Dynamical and Control Models via Certificate Synthesis," arXiv:2309.06090, 2023.

A. Abate, A. Edwards, M. Giacobbe, H. Punchihewa, and D. Roy, "Quantitative Neural Verification of Probabilistic Programs," CONCUR23, arXiv:2301.06136, 2023.

D. Roy, M. Giacobbe, and A. Abate, "Learning Probabilistic Termination Proofs," CAV21, LNCS 12760, pp. 3–26, 2021.

A. Abate, D. Ahmed, A. Edwards, M. Giacobbe and A. Peruffo, "FOSSIL: A Software Tool for the Formal Synthesis of Lyapunov Functions and Barrier Certificates using Neural Networks," HSCC, pp. 1-11, 2021.

A. Abate, D. Ahmed and A. Peruffo, "Automated Formal Synthesis of Neural Barrier Certificates for Dynamical Models," TACAS21, LNCS 12651, pp. 370–388, 2021.

D. Ahmed, A. Peruffo and A. Abate, "Automated and Sound Synthesis of Lyapunov Functions with SMT Solvers," TACAS20, LNCS 12078, pp. 97-114, 2020.

A. Abate, D. Ahmed, M. Giacobbe and A. Peruffo "Automated Formal Synthesis of Lyapunov Neural Networks," IEEE Control Systems Letters, 5 (3), 773-778, 2020.


A. Edwards, M. Giacobbe, and A. Abate, "On the Trade-off Between Efficiency and Precision of Neural Abstraction," QEST23, LNCS 14287, pp. 152-171, 2023.

A. Abate, A. Edwards, and M. Giacobbe, "Neural Abstractions," NeurIPS22, Advances in Neural Information Processing Systems 35, 26432-26447, 2022.


A. Abate, I. Bessa, D. Cattaruzza, L. Cordeiro, C. David, P. Kesseli, D. Kroening and E. Polgreen, "Automated Formal Synthesis of Provably Safe Digital Controllers for Continuous Plants," Acta Informatica, 57(3), 2020.

**Selected Journal References on (Model- and Sample-Based) Formal Abstractions**

T. Badings, L Romao, A. Abate, D. Parker, H. Poonawala, M. Stoelinga and N. Jansen, "Robust Control for Dynamical Systems with Non-Gaussian Noise via Formal Abstractions," JAIR, vol 76, pp.341-391, 2023.

T.S. Badings, A. Abate, N. Jansen, D. Parker, H.A. Poonawala, and M. Stoelinga, "Sampling-Based Robust Control of Autonomous Systems with Non-Gaussian Noise," AAAI22, 36 (9), pp. 9669-9678, 2022.

A. Lavaei, S. Soudjani, A. Abate, and M. Zamani, "Automated Verification and Synthesis of Stochastic Hybrid Systems: A Survey," Automatica, vol. 146, Dec. 2022.

L. Laurenti, M. Lahijanian, A. Abate, L. Cardelli and M. Kwiatkowska, "Formal and Efficient Control Synthesis for Continuous-Time Stochastic Processes," IEEE Transactions on Automatic Control, vol. 66, no. 1, pp. 17-32, Jan 2021.

S. Haesaert, S.E.Z. Soudjani, and A. Abate, "Verification of general Markov decision processes by approximate similarity relations and policy refinement," SIAM Journal on Control and Optimisation, vol. 55, nr. 4, pp. 2333-2367, 2017.

I. Tkachev, A. Mereacre, J.-P. Katoen, and A. Abate, "Quantitative Model Checking of Controlled Discrete-Time Markov Processes," Information and Computation, vol. 253, nr. 1, pp. 1–35, 2017.

S. Haesaert, N. Cauchi and A. Abate, "Certified policy synthesis for general Markov decision processes: An application in building automation systems," Performance Evaluation, vol. 117, pp. 75-103, 2017.

S.E.Z. Soudjani and A. Abate, "Aggregation and Control of Populations of Thermostatically Controlled Loads by Formal Abstractions," IEEE Transactions on Control Systems Technology. vol. 23, nr. 3, pp. 975–990, 2015.

S.E.Z. Soudjani and A. Abate, "Quantitative Approximation of the Probability Distribution of a Markov Process by Formal Abstractions," Logical Methods in Computer Science, Vol. 11, nr. 3, Oct. 2015.

M. Zamani, P. Mohajerin Esfahani, R. Majumdar, A. Abate, and J. Lygeros, "Symbolic control of stochastic systems via approximately bisimilar finite abstractions," IEEE Transactions on Automatic Control, vol. 59 nr. 12, pp. 3135-3150, Dec. 2014.

I. Tkachev and A. Abate, "Characterization and computation of infinite horizon specifications over Markov processes," Theoretical Computer Science, vol. 515, pp. 1-18, 2014.

S. Soudjani and A. Abate, "Adaptive and Sequential Gridding for Abstraction and Verification of Stochastic Processes," SIAM Journal on Applied Dynamical Systems, vol. 12, nr. 2, pp. 921-956, 2013.

A. Abate, J.P Katoen, J. Lygeros and M. Prandini, "Approximate Model Checking of Stochastic Hybrid Systems," European Journal of Control, 16(6), 624-641, 2010.

A. Abate, M. Prandini, J. Lygeros and S. Sastry, "Probabilistic Reachability and Safety Analysis of Controlled Discrete-Time Stochastic Hybrid Systems," Automatica, 44(11), 2724-2734, Nov. 2008.

**Selected references - Certified and Cautious Reinforcement Learning**

M. Hasanbeig, A. Abate, D. Kroening, "Certified Reinforcement Learning with Logic Guidance," AIJ, In Press, 2023. arXiv:1801.08099.

R. Mitta, H. Hasanbeig, Jun W, D. Kroening, Y. Kantaros, and A. Abate, "Safeguarded Progress in Reinforcement Learning: Safe Bayesian Exploration for Control Policy Synthesis," AAAI24, 2024.

A. Abate, Y. Almulla, J. Fox, D. Hyland, and M. Wooldridge, "Learning Task Automata for RL Using Hidden Markov Models," ECAI23, 2023.

M. Hasanbeig, A. Abate, and D. Kroening, "Logically-Constrained Neural Fitted Q-Iteration," AAMAS19, pp. 2012-2014, 2019.

M. Hasanbeig, A. Abate and D. Kroening, "Cautious Reinforcement Learning with Logical Constraints," AAMAS20, pp. 483-491, 2020.

M. Hasanbeig, D. Kroening and A. Abate, "Deep Reinforcement Learning with Temporal Logics," FORMATS20, LNCS 12288, pp. 1-22, 2020.

M. Hasanbeig, N. Jeppu, A. Abate, T. Melham and D. Kroening, "DeepSynth: Program Synthesis for Automatic Task Segmentation in Deep Reinforcement Learning," AAAI 2021.

L. Hammond, A. Abate, J. Gutierrez, and M. Wooldridge, "Multi-Agent Reinforcement Learning with Temporal Logic Specifications," AAAMAS 2021.

J. Skalse, L. Hammond, and A. Abate, "Lexicographic Multi-Objective Reinforcement Learning," IJCAI-ECAI 2022.

J. Skalse, L. Farnik, S. Motwani, E, Jenner, A. Gleave, and A. Abate, "STARC: A General Framework For Quantifying Differences Between Reward Functions," ICLR24, In Print, 2023.

J. Skalse, M. Farrugia-Roberts, S. Russell, A. Abate, and A. Gleave, "Invariance in Policy Optimisation and Partial Identifiability in Reward Learning," ICML23, PMLR, pp. 32033-32058, 2023.

J. Skalse and A. Abate, "Misspecification in Inverse Reinforcement Learning," AAAI23, vol. 37, nr. 12, pp. 15136-15143, 2023.

Backup slides

# Formal abstractions: algorithm

- approximate stochastic process $(\mathcal{S}, \mathcal{T})$ as Markov chain $(S, \mathbb{T})$, where
  - $S = \{z_1, z_2, \ldots, z_p\}$ – finite set of abstract states
  - $\mathbb{T} : \mathcal{S} \times \mathcal{S} \to [0, 1]$ – transition probability matrix

# Formal abstractions: algorithm

- approximate stochastic process $(\mathcal{S}, \mathcal{T})$ as Markov chain $(S, \mathbb{T})$, where
  - $S = \{z_1, z_2, \ldots, z_p\}$ – finite set of abstract states
  - $\mathbb{T} : \mathcal{S} \times \mathcal{S} \to [0,1]$ – transition probability matrix

- algorithm:

  **input:** stochastic process $(\mathcal{S}, \mathcal{T})$

  **output:** Markov chain $(S, \mathbb{T})$

# Formal abstractions: algorithm

- approximate stochastic process $(\mathcal{S}, \mathcal{T})$ as Markov chain $(S, \mathbb{T})$, where
    - $S = \{z_1, z_2, \ldots, z_p\}$ – finite set of abstract states
    - $\mathbb{T} : \mathcal{S} \times \mathcal{S} \rightarrow [0,1]$ – transition probability matrix
- algorithm:

---

**input:** stochastic process $(\mathcal{S}, \mathcal{T})$

1 select finite partition $\mathcal{S} = \cup_{i=1}^{p} S_i$          [aligned with $\mathcal{G}_i$]

**output:** Markov chain $(S, \mathbb{T})$

---

# Formal abstractions: algorithm

- approximate stochastic process $(\mathcal{S}, \mathcal{T})$ as Markov chain $(S, \mathbb{T})$, where
  - $S = \{z_1, z_2, \ldots, z_p\}$ – finite set of abstract states
  - $\mathbb{T} : \mathcal{S} \times \mathcal{S} \to [0,1]$ – transition probability matrix

- algorithm:

  > **input:** stochastic process $(\mathcal{S}, \mathcal{T})$
  >
  > 1 select finite partition $\mathcal{S} = \cup_{i=1}^{p} S_i$     [aligned with $\mathcal{G}_i$]
  > 2 select representative points $z_i \in S_i$
  > 3 define finite state space $S := \{z_i, i = 1, ..., p\}$
  >
  > **output:** Markov chain $(S, \mathbb{T})$

# Formal abstractions: algorithm

- approximate stochastic process $(\mathcal{S}, \mathcal{T})$ as Markov chain $(S, \mathbb{T})$, where
    - $S = \{z_1, z_2, \ldots, z_p\}$ – finite set of abstract states
    - $\mathbb{T} : \mathcal{S} \times \mathcal{S} \to [0, 1]$ – transition probability matrix

- algorithm:

---
**input:** stochastic process $(\mathcal{S}, \mathcal{T})$

1 select finite partition $\mathcal{S} = \cup_{i=1}^{p} S_i$        [aligned with $\mathcal{G}_i$]

2 select representative points $z_i \in S_i$

3 define finite state space $S := \{z_i, i = 1, \ldots, p\}$

4 compute transition probability matrix: $\mathbb{T}(z_i, z_j) = \mathcal{T}(S_j \mid z_i)$

**output:** Markov chain $(S, \mathbb{T})$
---

# Formal abstractions: algorithm

- approximate stochastic process $(\mathcal{S}, \mathcal{T})$ as Markov chain $(S, \mathbb{T})$, where
    - $S = \{z_1, z_2, \ldots, z_p\}$ – finite set of abstract states
    - $\mathbb{T} : \mathcal{S} \times \mathcal{S} \rightarrow [0, 1]$ – transition probability matrix

- algorithm:

    **input:** stochastic process $(\mathcal{S}, \mathcal{T})$
    1. select finite partition $\mathcal{S} = \cup_{i=1}^{p} S_i$          [aligned with $\mathcal{G}_i$]
    2. select representative points $z_i \in S_i$
    3. define finite state space $S := \{z_i, i = 1, \ldots, p\}$
    4. compute transition probability matrix: $\mathbb{T}(z_i, z_j) = \mathcal{T}(S_j \mid z_i)$
    **output:** Markov chain $(S, \mathbb{T})$

# Formal abstractions: error $\tilde{\zeta}$

- consider $\mathcal{T}(d\bar{s}|s) = \mathfrak{t}(\bar{s}|s)d\bar{s}$; assume $\mathfrak{t}$ is Lipschitz continuous, namely

$$\exists\, 0 \leq h_s < \infty: \quad \left| \mathfrak{t}(\bar{s}|s) - \mathfrak{t}(\bar{s}|s') \right| \leq h_s \left\| s - s' \right\|, \quad \forall s, s', \bar{s} \in \mathcal{S}$$

# Formal abstractions: error $\xi$

- consider $\mathcal{T}(d\bar{s}|s) = \mathfrak{t}(\bar{s}|s)d\bar{s}$; assume $\mathfrak{t}$ is Lipschitz continuous, namely

$$\exists\, 0 \leq h_s < \infty : \quad \left| \mathfrak{t}(\bar{s}|s) - \mathfrak{t}(\bar{s}|s') \right| \leq h_s \left\| s - s' \right\|, \quad \forall s, s', \bar{s} \in \mathcal{S}$$

- **one-step error**

  $\epsilon = h_s \delta, \quad \delta$ max diameter of partition sets

- **$T$-step error**      *(tuneable via $\delta$)*

  $\xi(\delta, T) = \epsilon T$

# Formal abstractions: error $\xi$

- consider $\mathcal{T}(d\bar{s}|s) = \mathfrak{t}(\bar{s}|s)d\bar{s}$; assume $\mathfrak{t}$ is Lipschitz continuous, namely

$$\exists\, 0 \leq h_s < \infty: \quad \left|\mathfrak{t}(\bar{s}|s) - \mathfrak{t}(\bar{s}|s')\right| \leq h_s \left\|s - s'\right\|, \quad \forall s, s', \bar{s} \in \mathcal{S}$$

- **one-step error**　　　　　　*(related to approximate probabilistic bisimulation)*

  $\epsilon = h_s \delta,$　　$\delta$ max diameter of partition sets

- **$T$-step error**　　　*(tuneable via $\delta$)*

  $\xi(\delta, T) = \epsilon T$

$\rightarrow$ improved and generalised error $\tilde{\xi}$

# Formal abstractions: error $\tilde{\zeta}$

- consider $\mathcal{T}(d\bar{s}|s) = \mathfrak{t}(\bar{s}|s)d\bar{s}$; assume $\mathfrak{t}$ is Lipschitz continuous, namely

$$\exists\, 0 \le h_s < \infty: \quad \left|\mathfrak{t}(\bar{s}|s) - \mathfrak{t}(\bar{s}|s')\right| \le h_s \left\|s - s'\right\|, \quad \forall s, s', \bar{s} \in \mathcal{S}$$

- **one-step error** *(related to approximate probabilistic bisimulation)*

  $\epsilon = h_s\delta$,    $\delta$ max diameter of partition sets

- **$T$-step error**    *(tuneable via $\delta$)*

  $\tilde{\zeta}(\delta, T) = \epsilon T$

$\rightarrow$ improved and generalised error $\tilde{\zeta}$

# Formal abstractions: probabilistic safety



- recall temporal logic properties, e.g. probabilistic safety:

*probability that execution, started at $s \in \mathcal{I}$,*
*stays in safe set $A = \mathcal{U}^c$ within $[0, T]$*

$$\mathcal{P}_s(A) = \mathbb{P}_s(s_k \in A, \forall k \in [0, T])$$

- probabilistic safe set with safety level $\theta \in [0, 1]$ is

$$S(\theta) = \{s \in \mathcal{S} : \mathcal{P}_s(A) \geq \theta\}$$

# Formal abstractions: probabilistic safety



- recall temporal logic properties, e.g. probabilistic safety:

*probability that execution, started at $s \in \mathcal{I}$,*
*stays in safe set $A = \mathcal{U}^c$ within $[0, T]$*

$$\mathcal{P}_s(A) = \mathbb{P}_s(s_k \in A, \forall k \in [0, T])$$

- probabilistic safe set with safety level $\theta \in [0, 1]$ is

$$S(\theta) = \{s \in \mathcal{S} : \mathcal{P}_s(A) \geq \theta\}$$

- whenever stochastic process $(\mathcal{S}, \mathcal{T})$ is controlled, $\sup_\pi \mathcal{P}_s(A)$

# Formal abstractions: probabilistic safety



- $\delta$-abstract $(\mathcal{S}, \mathcal{T})$ as MC $(S, \mathbb{T})$, so that $A \to A_\delta$, quantify error $\xi(\delta, T)$ as above

$\Rightarrow$ probabilistic safe set on $(\mathcal{S}, \mathcal{T})$

$$S(\theta) = \{s \in \mathcal{S} : \mathcal{P}_s(A) \geq \theta\}$$

# Formal abstractions: probabilistic safety



- $\delta$-abstract $(\mathcal{S}, \mathcal{T})$ as MC $(S, \mathbb{T})$, so that $A \to A_\delta$, quantify error $\xi(\delta, T)$ as above

$\Rightarrow$ probabilistic safe set on $(\mathcal{S}, \mathcal{T})$

$$S(\theta) = \{s \in \mathcal{S} : \mathcal{P}_s(A) \geq \theta\}$$

is automatically computed with model checker (e.g. PRISM) on $(S, \mathbb{T})$ as

$$Z_\delta(\theta + \xi) \doteq \mathsf{Sat}\left(\mathbb{P}_{\geq \theta + \xi}\left(\square^{\leq T} A_\delta\right)\right)$$
$$= \left\{z \in S : z \models \mathbb{P}_{\geq \theta + \xi}\left(\square^{\leq T} A_\delta\right)\right\}$$

- whenever stochastic process $(\mathcal{S}, \mathcal{T})$ is controlled, obtain $\arg\sup_\pi \mathcal{P}_s(A)$

# Formal abstractions with data

- alluring idea: can we abstract models by sampling their dynamics?

# Formal abstractions with data

*"Timeō dāta, et dōna ferentēs"* [Laocoon, Aeneid]

- alluring idea: can we abstract models by sampling their dynamics?
- Beware many subtle issues: zero-measure sets, memory dependencies, . . .



$y_0 y_1 \ldots y_{19} = 01110111011101110111$

$$x^+ = x + \theta \bmod 2\pi$$

# Formal abstractions with data

$$x^+ = A(\alpha)x + B(\alpha)u + \sigma$$

- $\sigma \sim \mathcal{P}$ unknown - aleatoric uncertainty
- $\alpha \in \Theta$ - epistemic uncertainty



*($\rho$ is trace of closed-loop trajectory)*

*(probabilistic reach-avoid specification)*

Given $T \in \mathbb{N}$, and sets $\mathcal{G}$ (goal) and $\mathcal{U}^C$ (safe), find controller s.t., $\forall x_0 \in \mathcal{I}$,

$$\mathbb{P}_{\mathcal{I}}\{\rho \models \mathcal{U}^C \mathsf{U}^{\leq T} \mathcal{G}\} \geq \theta, \quad \text{with confidence} \geq 1 - \beta$$

*AAAI22 BPA*

$$x^+ = A(\bar{\alpha})x + B(\bar{\alpha})u + \sigma$$

# Formal abstractions with data

$$x^+ = A(\alpha)x + B(\alpha)u + \sigma$$

- $\sigma \sim \mathcal{P}$ unknown - aleatoric uncertainty



- scenario approach for convex optimisation: $\mathbb{P}\{\underline{p} \leq P(s' \mid s_i, a) \leq \bar{p}\} \geq 1 - \beta$
- abstraction as iMDP

# Formal abstractions with data

$$x^+ = A(\alpha)x + B(\alpha)u + \sigma$$

- $\alpha \in \Theta$ - epistemic uncertainty



- abstraction as iMDP

# Formal abstractions with data

$$x^+ = A(\alpha)x + B(\alpha)u + \sigma$$

- $\alpha \in \Theta$ - epistemic uncertainty



- abstraction as iMDP

# Formal abstractions with data

$$x^+ = A(\alpha)x + B(\alpha)u + \sigma$$

- $\alpha \in \Theta$ - epistemic uncertainty



- abstraction as iMDP

# Formal abstractions with data

*AAAI22 BPA*