

Interaction-aware Reachability of Neural Network Controlled Systems: A Mixed Monotone Approach

Saber Jafarpour



University of Colorado **Boulder**



**Georgia
Tech.**

September 28, 2023



Akash Harapanahalli
Georgia Tech



Samuel Coogan
Georgia Tech

SJ and A. Harapanahalli and S. Coogan. [Efficient Interaction-Aware Interval Analysis of Neural Network Feedback Loops](#). arXiv, 2023.

A. Harapanahalli and **SJ** and S. Coogan. [A Toolbox for Fast Interval Arithmetic in numpy with an Application to Formal Verification of Neural Network Controlled Systems](#). 2nd ICML workshop on Formal Verification of Machine Learning, 2023

A. Harapanahalli and **SJ** and S. Coogan. [Forward Invariance in Neural Network Controlled Systems](#). arXiv, 2023.

Neural Network Controllers

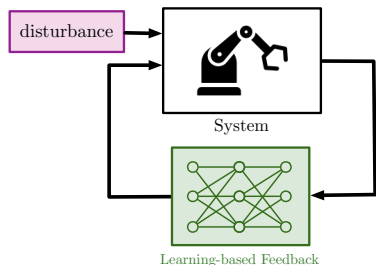
Motivations

- Neural Networks as controllers in safety-critical applications (examples: autonomous vehicles and mobile robots)

Goal: ensure and verify *safety* of the closed-loop system

Issues with neural network controllers:

- large # of parameters with nonlinearity
- sensitive wrt to input perturbations
- limited closed-loop safety guarantees



Challenges

Rigorous verification and computational efficiency vs. accuracy.

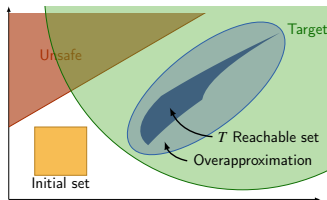
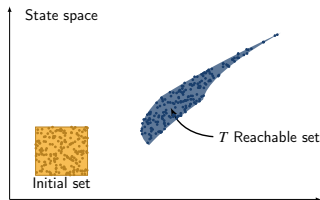
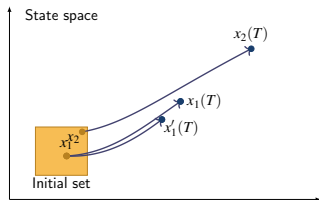
Safety Verification via Reachability Analysis

Problem Statement

System : $\dot{x} = f(x, w)$

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

Disturbance : $w \in \mathcal{W} \subseteq \mathbb{R}^m$



- reachable sets characterize evolution of the system

$$\mathcal{R}^f(t, \mathcal{X}_0) = \{x_w(t) \mid x_w(\cdot) \text{ is a traj of the system for some } w \text{ with } x_0 \in \mathcal{X}_0\}$$

- over-approximation of reachable sets for safety and verification
- reachability of dynamical system is an old problem with several classical approaches

Classical approaches are not scalable to large-scale nonlinear systems

Monotone Dynamical Systems

Definition and Characterization

A dynamical system $\dot{x} = f(x, w)$ is monotone¹(with respect to cones K, C) if

$$x_u(0) \preceq_K y_w(0) \quad \text{and} \quad u \preceq_C w \quad \implies \quad x_u(t) \preceq_K y_w(t) \quad \text{for all time}$$

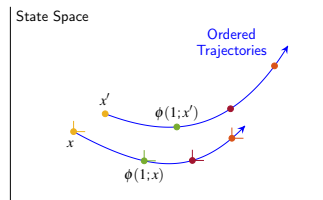
where \preceq_K (\preceq_C) is the partial order with induced by the cone K (cone C).

A **polyhedral cone** has the form

$$K = \underbrace{\{y \in \mathbb{R}^n \mid H_K y \geq 0_p\}}_{\text{halfspace rep}} = \underbrace{\{V_K y \mid y \geq 0_p\}}_{\text{vertex rep}}$$

Monotonicity test

- 1 $H_K \left(\frac{\partial f}{\partial x}(x, w) + \alpha(x, w) I_n \right) V_K \geq 0$ for some $\alpha(x, w)$
- 2 $H_K \frac{\partial f}{\partial w}(x, w) V_C \geq 0$



¹D. Angeli and E. Sontag, "Monotone control systems", IEEE TAC, 2003

Monotone Dynamical Systems

Definition and Characterization

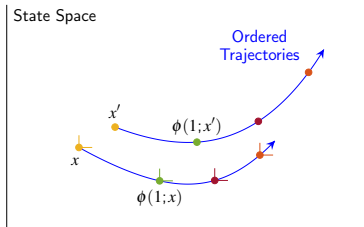
A dynamical system $\dot{x} = f(x, w)$ is monotone¹(with respect to cones K, C) if

$$x_u(0) \preceq_K y_w(0) \quad \text{and} \quad u \preceq_C w \quad \implies \quad x_u(t) \preceq_K y_w(t) \quad \text{for all time}$$

where \preceq_K (\preceq_C) is the partial order with induced by the cone K (cone C).

Monotonicity test (for the standard cone $\mathbb{R}_{\geq 0}^n$)

- 1 $\frac{\partial f}{\partial x}(x, w)$ is Metzler (off-diag ≥ 0)
- 2 $\frac{\partial f}{\partial w}(x, w) \geq 0_m$



In this talk: monotone system theory for reachability analysis

¹D. Angeli and E. Sontag, "Monotone control systems", IEEE TAC, 2003

Reachability of Monotone Dynamical Systems

Hyper-rectangular over-approximations

Theorem (classical result)

For a monotone system with $\mathcal{W} = [\underline{w}, \bar{w}]$

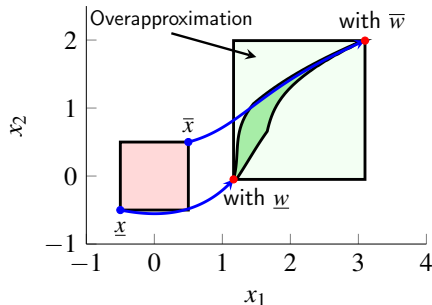
$$\mathcal{R}^f(t, [\underline{x}_0, \bar{x}_0]) \subseteq [x_{\underline{w}}(t), x_{\bar{w}}(t)]$$

where $x_{\underline{w}}(\cdot)$ (resp. $x_{\bar{w}}(\cdot)$) is the trajectory with disturbance \underline{w} (resp. \bar{w}) starting at \underline{x}_0 (resp. \bar{x}_0)

Example:

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2^3 - x_1 + w \\ x_1 \end{bmatrix}$$

$$\mathcal{W} = [2.2, 2.3] \quad \mathcal{X}_0 = \left[\begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \right]$$



Non-monotone Dynamical Systems

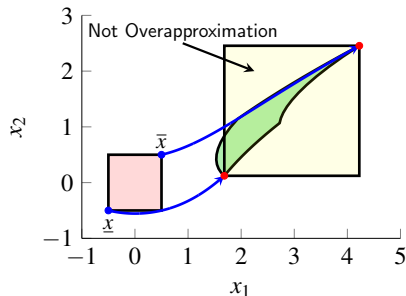
Reachability analysis

- For non-monotone dynamical systems the extreme trajectories do not provide any over-approximation of reachable sets

Example:

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2^3 - x_2 + w \\ x_1 \end{bmatrix}$$

$$\mathcal{W} = [2.2, 2.3] \quad \mathcal{X}_0 = \left[\begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \right]$$



Mixed Monotone Theory

Embedding into a larger system

- **Key idea:** embed the dynamical system on \mathbb{R}^n into a dynamical system on \mathbb{R}^{2n}
- Assume $\mathcal{W} = [\underline{w}, \bar{w}]$ and $\mathcal{X}_0 = [\underline{x}_0, \bar{x}_0]$

Original system

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

Embedding system

$$\begin{aligned}\dot{x} &= \underline{d}(x, \bar{x}, \underline{w}, \bar{w}), \\ \dot{\bar{x}} &= \bar{d}(x, \bar{x}, \underline{w}, \bar{w})\end{aligned}$$

\underline{d}, \bar{d} are **decomposition functions** s.t.

- 1 $f(x, w) = \underline{d}(x, x, w, w)$ for every x, w
- 2 **cooperative:** $(\underline{x}, \underline{w}) \mapsto \underline{d}(\underline{x}, \bar{x}, \underline{w}, \bar{w})$
- 3 **competitive:** $(\bar{x}, \bar{w}) \mapsto \bar{d}(\underline{x}, \bar{x}, \underline{w}, \bar{w})$
- 4 the same properties for \bar{d}

The embedding system is a monotone dynamical system on \mathbb{R}^{2n} with respect to the **southeast** partial order \leq_{SE} :

$$\begin{bmatrix} x \\ \hat{x} \end{bmatrix} \leq_{SE} \begin{bmatrix} y \\ \hat{y} \end{bmatrix} \iff x \leq y \quad \text{and} \quad \hat{y} \leq \hat{x}$$

In terms of cones, \leq_{SE} is induced by the cone $\mathbb{R}_{\geq 0}^n \times -\mathbb{R}_{\geq 0}^n$.

- f locally Lipschitz \implies a decomposition function exists

The **best (tightest)** decomposition function is given by

$$\underline{d}_i(\underline{x}, \bar{x}, \underline{w}, \bar{w}) = \min_{\substack{z \in [\underline{x}, \bar{x}], z_i = \underline{x}_i \\ u \in [\underline{w}, \bar{w}]}} f_i(z, u), \quad \bar{d}_i(\underline{x}, \bar{x}, \underline{w}, \bar{w}) = \max_{\substack{z \in [\underline{x}, \bar{x}], z_i = \bar{x}_i \\ u \in [\underline{w}, \bar{w}]}} f_i(z, u)$$

A short (and incomplete) history:

J-L. Gouze and L. P. Hadeler. [Monotone flows and order intervals](#). Nonlinear World, 1994

G. Enciso, H. Smith, and E. Sontag. [Nonmonotone systems decomposable into monotone systems with negative feedback](#). Journal of Differential Equations, 2006.

H. Smith. [Global stability for mixed monotone systems](#). Journal of Difference Equations and Applications, 2008

Reachability using Embedding Systems

Hyper-rectangular over-approximations

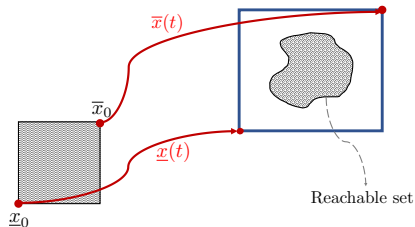
Theorem²

Assume $\mathcal{W} = [\underline{w}, \bar{w}]$ and $\mathcal{X}_0 = [\underline{x}_0, \bar{x}_0]$ and

$$\dot{\underline{x}} = \underline{d}(\underline{x}, \bar{x}, \underline{w}, \bar{w}), \quad \underline{x}(0) = \underline{x}_0$$

$$\dot{\bar{x}} = \bar{d}(\bar{x}, \underline{x}, \bar{w}, \underline{w}), \quad \bar{x}(0) = \bar{x}_0$$

Then $\mathcal{R}^f(t, \mathcal{X}_0) \subseteq [\underline{x}(t), \bar{x}(t)]$



(Scalable) a single trajectory of embedding system provides **lower bound** (\underline{x}) and **upper bound** (\bar{x}) for the trajectories of the original system.

²Coogan and Arcak, “Efficient finite abstraction of mixed monotone systems”, HSCC, 2015.

Reachability using Embedding Systems

Example

Original System:

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2^3 - x_2 + w \\ x_1 \end{bmatrix}$$

$$\mathcal{W} = [2.2, 2.3] \quad \mathcal{X}_0 = \left[\begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \right]$$

red = cooperative, blue = competitive

Decomposition function

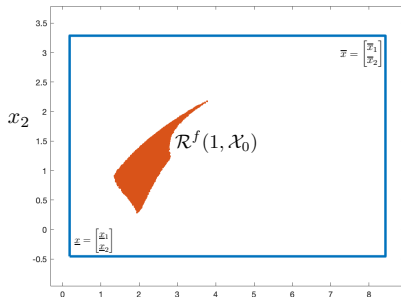
$$\underline{d}(x, \bar{x}, \underline{w}, \bar{w}) = \begin{bmatrix} x_2^3 + \underline{w} \\ x_1 \end{bmatrix} + \begin{bmatrix} -\bar{x}_2 \\ 0 \end{bmatrix}$$

$$\bar{d}(x, \bar{x}, \underline{w}, \bar{w}) = \begin{bmatrix} \bar{x}_2^3 + \bar{w} \\ \bar{x}_1 \end{bmatrix} + \begin{bmatrix} -x_2 \\ 0 \end{bmatrix}$$

Embedding System:

$$\frac{d}{dt} \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \bar{x}_1 \\ \bar{x}_2 \end{bmatrix} = \begin{bmatrix} \underline{x}_2^3 - \bar{x}_2 + \underline{w} \\ x_1 \\ \bar{x}_2^3 - \underline{x}_2 + \bar{w} \\ \bar{x}_1 \end{bmatrix} \quad \begin{bmatrix} \underline{w} \\ \bar{w} \end{bmatrix} = \begin{bmatrix} 2.2 \\ 2.3 \end{bmatrix}$$

$$\begin{bmatrix} \underline{x}_1(0) \\ \underline{x}_2(0) \end{bmatrix} = \begin{bmatrix} -0.5 \\ -0.5 \end{bmatrix} \quad \begin{bmatrix} \bar{x}_1(0) \\ \bar{x}_2(0) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$



Systems with NN Controllers

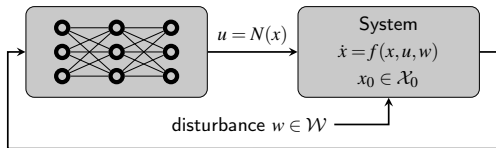
A Mixed Monotone Approach

Given the open-loop nonlinear system with a neural network controller

$$\begin{aligned}\dot{x} &= f(x, u, w), \\ u &= N(x),\end{aligned}$$

study reachability of the closed-loop system

$$\dot{x} = f(x, N(x), w) := f^c(x, w)$$



Challenge: find a decomposition function for closed-loop system

Key observation: Interval bounds for neural networks combines nicely with **mixed monotone theory** for the open-loop system!

- Interval bounds for NN using verification algorithms (CROWN, LipSDP, IBP, etc)

Question: How to capture the interaction between NN and the system

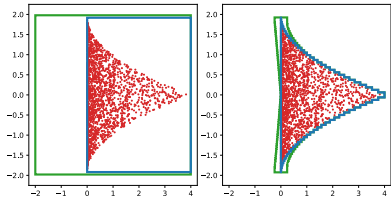
Decomposition Functions for Systems

A Jacobian-based Approach

Jacobian-based: $\dot{x} = f(x, w)$ such that $\frac{\partial f}{\partial x} \in [J_{[x, \bar{x}]}, \bar{J}_{[x, \bar{x}]}]$ and $\frac{\partial f}{\partial w} \in [J_{[w, \bar{w}]}, \bar{J}_{[w, \bar{w}]}]$, then

$$\begin{bmatrix} \underline{d}(\underline{x}, \bar{x}, \underline{w}, \bar{w}) \\ \bar{d}(\underline{x}, \bar{x}, \underline{w}, \bar{w}) \end{bmatrix} = \begin{bmatrix} -[J_{[x, \bar{x}]}]^- & [J_{[x, \bar{x}]}]^- \\ -[J_{[x, \bar{x}]}]^+ & [J_{[x, \bar{x}]}]^+ \end{bmatrix} \begin{bmatrix} \underline{x} \\ \bar{x} \end{bmatrix} + \begin{bmatrix} -[J_{[w, \bar{w}]}]^- & [J_{[w, \bar{w}]}]^- \\ -[J_{[w, \bar{w}]}]^+ & [J_{[w, \bar{w}]}]^+ \end{bmatrix} \begin{bmatrix} \underline{w} \\ \bar{w} \end{bmatrix} + \begin{bmatrix} f(\underline{x}, \underline{w}) \\ f(\underline{x}, \underline{w}) \end{bmatrix}$$

- Interval arithmetic allows computing Jacobian bounds efficiently using *inclusion functions*.
- `npinterval`³: Toolbox that implements intervals as native data-type in `numpy`.



$$g(x_1, x_2) = [(x_1 + x_2)^2, 4 \sin((x_1 - x_2)/4)]^T$$

vs.

$$g(x_1, x_2) = [x_2^2 + 2x_1x_2 + x_1^2, 4 \sin(x_1/4) \cos(x_2/4) - 4 \cos(x_1/4) \sin(x_2/4)]^T$$

³Harapanahalli, Jafarpour, Coogan. "A Toolbox for Fast Interval Arithmetic in `numpy` with an Application to Formal Verification of Neural Network Controlled Systems", 2nd WFVML, ICML, 2023

Interval Bounds for Neural Networks

Input-output Bounds vs. Functional Bounds

Input-output bounds: Given a neural network controller $u = N(x)$

$$\underline{u}_{[\underline{x}, \bar{x}]} \leq N(x) \leq \bar{u}_{[\underline{x}, \bar{x}]}, \quad \text{for all } x \in [\underline{x}, \bar{x}]$$

Functional bounds: Given a neural network controller $u = N(x)$

$$\underline{N}_{[\underline{x}, \bar{x}]}(x) \leq N(x) \leq \bar{N}_{[\underline{x}, \bar{x}]}(x), \quad \text{for all } x \in [\underline{x}, \bar{x}]$$

- Example: CROWN⁴ can provide both input-output and functional bounds.

CROWN functional bounds:

$$\begin{aligned} \underline{N}_{[\underline{x}, \bar{x}]}(x) &= \underline{A}_{[\underline{x}, \bar{x}]}x + \underline{b}_{[\underline{x}, \bar{x}]}, \\ \bar{N}_{[\underline{x}, \bar{x}]}(x) &= \bar{A}_{[\underline{x}, \bar{x}]}x + \bar{b}_{[\underline{x}, \bar{x}]} \end{aligned}$$

CROWN input-output bounds:

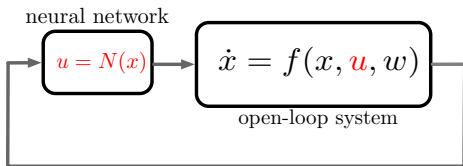
$$\begin{aligned} \underline{u}_{[\underline{x}, \bar{x}]} &= \underline{A}_{[\underline{x}, \bar{x}]}^+ \bar{x} + \bar{A}_{[\underline{x}, \bar{x}]}^- \underline{x} + \underline{b}_{[\underline{x}, \bar{x}]}, \\ \bar{u}_{[\underline{x}, \bar{x}]} &= \bar{A}_{[\underline{x}, \bar{x}]}^+ \bar{x} + \underline{A}_{[\underline{x}, \bar{x}]}^- \underline{x} + \bar{b}_{[\underline{x}, \bar{x}]} \end{aligned}$$

⁴Zhang, Weng, Chen, Hsieh, Daniel. "Efficient neural network robustness certification with general activation functions." NeurIPS, 2018.

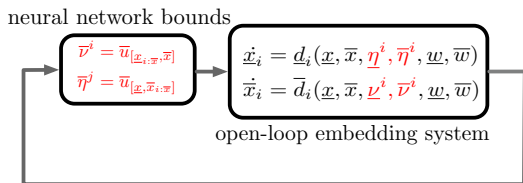
Approach #1: Interconnection-based Approach

A pictorial explanation

Original system:

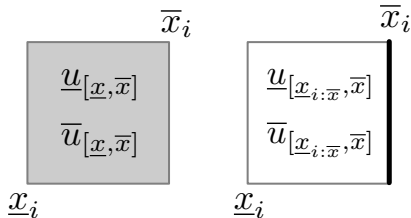


Embedding system:



How does the **interconnection-based approach** work?

- closed-loop embedding system = interconnection of NN interval bounds + open-loop embedding system
- NN bounds are evaluated on each edge instead of the whole box, i.e., we use $\underline{u}_{[x_i, \bar{x}]}$ and $\bar{u}_{[x_i, \bar{x}]}$ instead of $\underline{u}_{[x, \bar{x}]}$ and $\bar{u}_{[x, \bar{x}]}$.



Theorem⁵

- 1 **Decomposition function** \underline{d}, \bar{d} for the open-loop system $\dot{x} = f(x, u, w)$
- 2 **Interval input-output bounds** $\underline{u}_{[\underline{x}, \bar{x}]}, \bar{u}_{[\underline{x}, \bar{x}]}$ for the neural network controller $u = N(x)$,

Then

$$\underline{d}_i^c(\underline{x}, \bar{x}, \underline{w}, \bar{w}) = \underline{d}_i(\underline{x}, \bar{x}, \underline{\eta}^i, \bar{\eta}^i, \underline{w}, \bar{w})$$

$$\bar{d}_i^c(\underline{x}, \bar{x}, \underline{w}, \bar{w}) = \bar{d}_i(\underline{x}, \bar{x}, \underline{\nu}^i, \bar{\nu}^i, \underline{w}, \bar{w})$$

where

$$\underline{\eta}^i = \underline{u}_{[\underline{x}, \bar{x}_{i:\underline{x}}]} \quad \bar{\eta}^i = \bar{u}_{[\underline{x}, \bar{x}_{i:\underline{x}}]}, \quad \underline{\nu}^i = \underline{u}_{[\underline{x}_{i:\bar{x}}, \bar{x}]} \quad \bar{\nu}^i = \bar{u}_{[\underline{x}_{i:\bar{x}}, \bar{x}]},$$

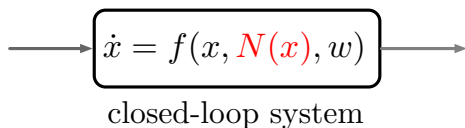
is a **decomposition function for the closed-loop system** where $v_{i:w}$ is the vector v with i th component replaced with i th component of w .

⁵Jafarpour, Harapanahalli, Coogan. "Efficient Interaction-aware Interval Reachability of Neural Network Feedback Loops", arXiv, 2003

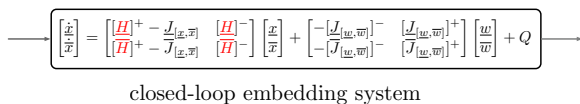
Approach #2: Interaction-based Approach

A pictorial explanation

Original system:



Embedding system:



How does the **interaction-based approach** work?

- Closed-loop decomposition function = Jacobian based for $f(x, N(x), w)$.
- Neural Network affine functional bounds

$$\underline{N}_{[x,\bar{x}]} = \underline{A}_{[x,\bar{x}]}x + \underline{b}_{[x,\bar{x}]},$$

$$\overline{N}_{[x,\bar{x}]} = \overline{A}_{[x,\bar{x}]}x + \overline{b}_{[x,\bar{x}]}$$

are used to compute the interactions.

$$\underline{H} = \underline{J}_{[x,\bar{x}]} + [\underline{J}_{[u,\bar{u}}]^+ \underline{A}_{[x,\bar{x}]} + [\underline{J}_{[u,\bar{u}}]^- \overline{A}_{[x,\bar{x}]}]$$

$$\overline{H} = \overline{J}_{[x,\bar{x}]} + [\underline{J}_{[u,\bar{u}}]^+ \overline{A}_{[x,\bar{x}]} + [\underline{J}_{[u,\bar{u}}]^- \underline{A}_{[x,\bar{x}]}]$$

Theorem⁶

Let $\frac{\partial f}{\partial x} \in [J_{[\underline{x}, \bar{x}]}, \bar{J}_{[\underline{x}, \bar{x}]}]$, $\frac{\partial f}{\partial u} \in [J_{[\underline{u}, \bar{u}]}, \bar{J}_{[\underline{u}, \bar{u}]}]$, and $\frac{\partial f}{\partial w} \in [J_{[\underline{w}, \bar{w}]}, \bar{J}_{[\underline{w}, \bar{w}]}]$. Then

$$\begin{bmatrix} \underline{d}_i^c(\underline{x}, \bar{x}, \underline{w}, \bar{w}) \\ \bar{d}_i^c(\underline{x}, \bar{x}, \underline{w}, \bar{w}) \end{bmatrix} = \begin{bmatrix} [\underline{H}]^+ - J_{[\underline{x}, \bar{x}]} & [\underline{H}]^- \\ [\bar{H}]^+ - \bar{J}_{[\underline{x}, \bar{x}]} & [\bar{H}]^- \end{bmatrix} \begin{bmatrix} \underline{x} \\ \bar{x} \end{bmatrix} + \begin{bmatrix} -[J_{[\underline{w}, \bar{w}]}]^- & [J_{[\underline{w}, \bar{w}]}]^+ \\ -[\bar{J}_{[\underline{w}, \bar{w}]}]^- & [\bar{J}_{[\underline{w}, \bar{w}]}]^+ \end{bmatrix} \begin{bmatrix} \underline{w} \\ \bar{w} \end{bmatrix} + Q$$

where

$$\begin{aligned} \underline{H} &= J_{[\underline{x}, \bar{x}]} + [J_{[\underline{u}, \bar{u}]}]^+ \underline{A}_{[\underline{x}, \bar{x}]} + [J_{[\underline{u}, \bar{u}]}]^- \bar{A}_{[\underline{x}, \bar{x}]} \\ \bar{H} &= \bar{J}_{[\underline{x}, \bar{x}]} + [J_{[\underline{u}, \bar{u}]}]^+ \bar{A}_{[\underline{x}, \bar{x}]} + [J_{[\underline{u}, \bar{u}]}]^- \underline{A}_{[\underline{x}, \bar{x}]} \end{aligned}$$

is a **decomposition function for the closed-loop system**.

⁶Jafarpour, Harapanahalli, Coogan. "Efficient Interaction-aware Interval Reachability of Neural Network Feedback Loops", arXiv, 2003

Case Study: Bicycle Model

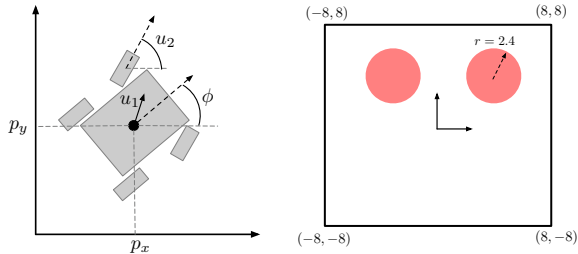
Design of the neural network

Dynamics of bicycle

$$\dot{p}_x = v \cos(\phi + \beta(u_2)) \quad \dot{\phi} = \frac{v}{l_r} \sin(\beta(u_2))$$

$$\dot{p}_y = v \sin(\phi + \beta(u_2)) \quad \dot{v} = u_1$$

$$\beta(u_2) = \arctan\left(\frac{l_r}{l_f + l_r} \tan(u_2)\right)$$



Goal: steer the bicycle to the origin avoiding the obstacles

- **offline controller:** MPC with hard constraint to avoid the obstacles
- run MPC for 65000 randomly chosen initial condition (20 sample per trajectory)
- train a feedforward neural network $4 \mapsto 100 \mapsto 100 \mapsto 2$ with this data

Case Study: Bicycle Model

Numerical Experiments

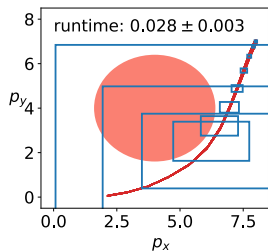
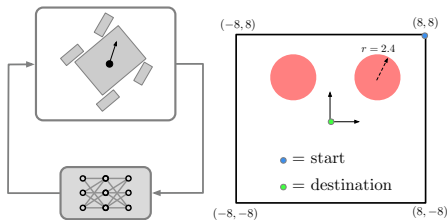
- start from $(8, 8)$ toward $(0, 0)$

- $\mathcal{X}_0 = [\underline{x}_0, \bar{x}_0]$ with

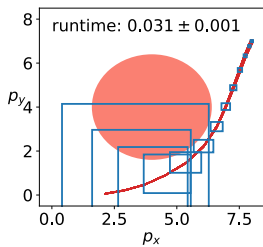
$$\underline{x}_0 = \left(7.95 \quad 6.95 \quad -\frac{2\pi}{3} - 0.01 \quad 1.99 \right)^\top$$

$$\bar{x}_0 = \left(8.05 \quad 7.05 \quad -\frac{2\pi}{3} + 0.01 \quad 2.01 \right)^\top$$

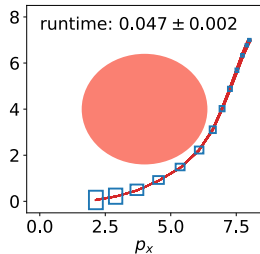
- CROWN for verification of neural network



naive combination



interconnection



interaction

- reachability using mixed monotone theory
- mixed monotone theory for reachability of NN controlled systems
- two methods for capturing the interaction between system and NN controller
 - 1 interconnection-based approach
 - 2 interaction-based approach

Future directions:

- forward invariance of systems with NN controllers
- design of suitable correction actions
- ensuring safety in the training of NN

Embedding System for Linear Dynamical System

A structure preserving decomposition

- Metzler/non-Metzler decomposition: $A = [A]^{Mzl} + [A]^{Mzl}$

- Example: $A = \begin{bmatrix} 2 & 0 & -1 \\ 1 & -3 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow [A]^{Mzl} = \begin{bmatrix} 2 & 0 & 0 \\ 1 & -3 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ $[A]^{Mzl} = \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$

Linear systems

Original system

$$\dot{x} = Ax + Bw$$

Embedding system

$$\dot{\underline{x}} = [A]^{Mzl} \underline{x} + [A]^{Mzl} \bar{x} + B^+ \underline{w} + B^- \bar{w}$$

$$\dot{\bar{x}} = [A]^{Mzl} \bar{x} + [A]^{Mzl} \underline{x} + B^+ \bar{w} + B^- \underline{w}$$

