

# Exploiting Structure in Feedback Systems with Learning-based Components

Saber Jafarpour



Decision and Control Laboratory  
Georgia Institute of Technology

January 18, 2023

# Modern societal autonomous systems

## Introduction



Power grids



Transportation networks



Learning-based systems

- large penetration of distributed renewable units in power grids
- unprecedented demand is pushing transportation networks to their capacity
- increasing deployment of learning algorithms in safety-critical systems

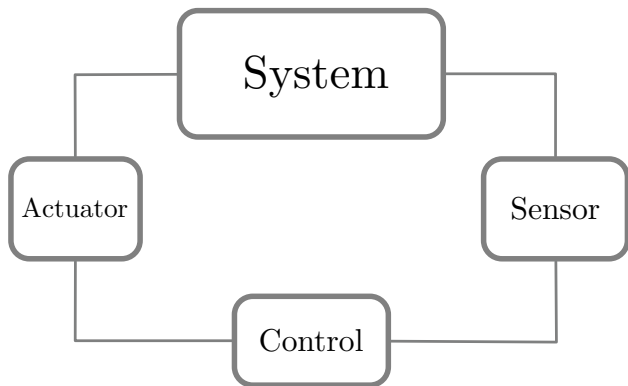
societal autonomous systems are becoming **large-scale** with **interconnected** and **complex** components

reconsider the traditional approaches for **monitoring** and **control** of autonomous systems

# Feedback control of autonomous systems

## Opportunities and challenges

Feedback is a central paradigm in control theory

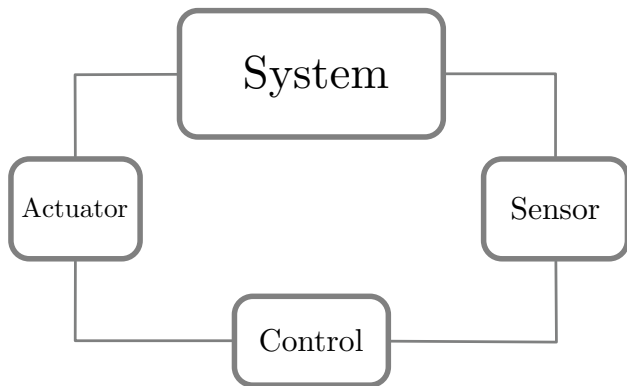


Magic of feedback<sup>a</sup>: robustness, shape behavior, command tracking, etc.

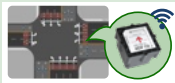
<sup>a</sup>Karl J. Astrom, Automatic Control - A Perspective, 2019

# Feedback control of autonomous systems

## Opportunities and challenges



Agents have wide range of **communication** capabilities



Enhanced processing units allow new **computational** approaches



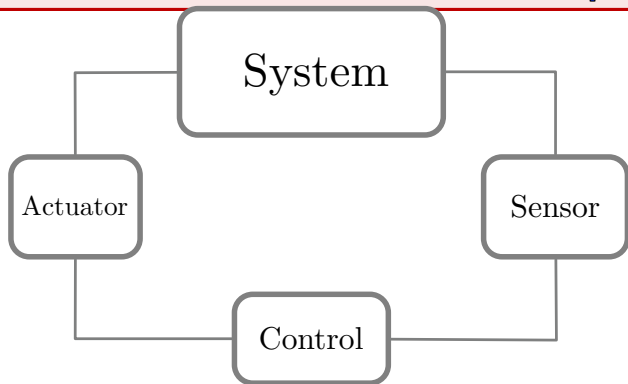
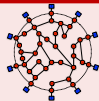
**Large number** of measurement devices for sensing



# Feedback control of autonomous systems

## Opportunities and challenges

Systems are becoming **large-scale** with **heterogeneous** and **interconnected** components



Controllers contain **high-dimensional**, **learning-based**, and **complex** parts



# My research

Safety and robustness in control of autonomous systems

## A critical task

Desired performance while ensuring their **safety** and **robustness**.



2011 US Southwest blackout



Traffic congestion in Beijing



Self-driving car accident

## My Contribution

Exploit **structure** to ensure safety and resilience in control of large-scale autonomous systems

**Tools:** control theory, dynamical systems, optimization

### Stability of large-scale power grids

- threshold of frequency synchronization (TAC 2018, SICON 2019)
- multi-stability via partitioning the state-space (SIAM Review 2021, Nature Com 2022)
- dynamic stability of low-inertia power grids (TCNS 2019)

### Geometric control

- small time local controllability (SICON 2020)
- locally convex topologies and control theory (MCSS 2016)

### Contraction theory

- weak and semi-contraction theory (TAC 2021)
- non-Euclidean contraction theory (TAC 2022)
- time-varying optimization (TAC 2021)
- non-Euclidean monotone operator theory (CDC 2022)

### Robustness of learning algorithms

- implicit neural networks (NeurIPS 2021, L4DC 2022)
- interval reachability of neural networks (L4DC 2022)
- safety verification of feedback loops

### Learning-based feedback

Feedback controller or some elements of it are learned from data



Aerial vehicles



Self-driving cars



GaTech A1 robot

### Why data-driven feedback?

- models are complicated or not available
- environment is unknown or varying
- traditional methods are cumbersome



### Learning-based feedback

Feedback controller or some elements of it are learned from data



Aerial vehicles



Self-driving cars



GaTech A1 robot

### Why data-driven feedback?

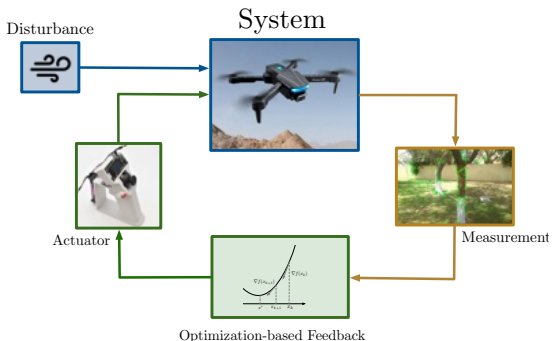
- models are complicated or not available
- environment is unknown or varying
- traditional methods are cumbersome

# Learning-based feedback

A data-driven approach to controller design

**Assumption:** An (approximate) model of the system is available

## Optimization-based control



- **Example method:** Model Predictive Control (MPC)

$$\min_{u(0), \dots, u(N-1)} \sum_{i=0}^{N-1} \ell(x(t), u(t)) + \phi(x(N)),$$
$$x(t+1) = x(t) + \alpha f(x(t), u(t)),$$
$$x(t) \in \mathcal{X}, \quad t \in \{1, \dots, N\}$$
$$u(t) \in \mathcal{U}, \quad t \in \{0, \dots, N-1\}$$
$$x(0) = x$$

- $\mathcal{X}$  and  $\mathcal{U}$  are the safety constraints

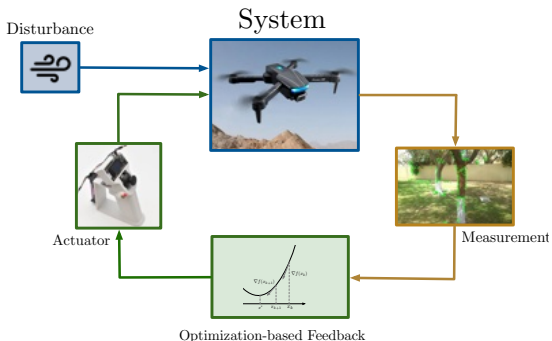
Feedback law:  $u(0) = K(x)$

# Learning-based feedback

A data-driven approach to controller design

**Assumption:** An (approximate) model of the system is available

## Optimization-based control



- Example issues: set  $\mathcal{X}$  is learned online

$$\min_{u(0), \dots, u(N-1)} \sum_{i=0}^{N-1} \ell(x(t), u(t)) + \phi(x(N)),$$
$$x(t+1) = x(t) + \alpha f(x(t), u(t)),$$
$$\mathbf{x}(t) \in \mathcal{X}, \quad t \in \{1, \dots, N\}$$
$$u(t) \in \mathcal{U}, \quad t \in \{0, \dots, N-1\}$$
$$x(0) = x$$

- $\mathcal{X}$  and  $\mathcal{U}$  are the safety constraints

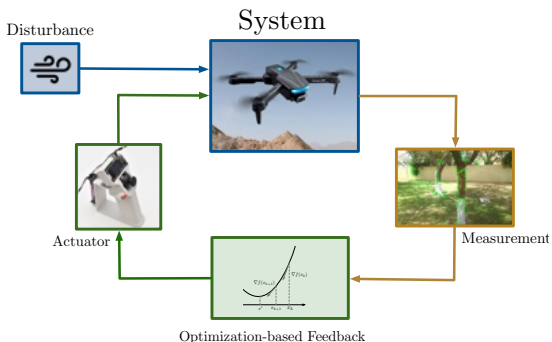
Feedback law:  $u(0) = K(x)$

# Learning-based feedback

A data-driven approach to controller design

**Assumption:** An (approximate) model of the system is available

## Optimization-based control



- **Example issues:** the optimization problem is **computationally complicated**

$$\min_{u(0), \dots, u(N-1)} \sum_{i=0}^{N-1} \ell(x(t), u(t)) + \phi(x(N)),$$
$$x(t+1) = x(t) + \alpha f(x(t), u(t)),$$
$$x(t) \in \mathcal{X}, \quad t \in \{1, \dots, N\}$$
$$u(t) \in \mathcal{U}, \quad t \in \{0, \dots, N-1\}$$
$$x(0) = x$$

- $\mathcal{X}$  and  $\mathcal{U}$  are the safety constraints

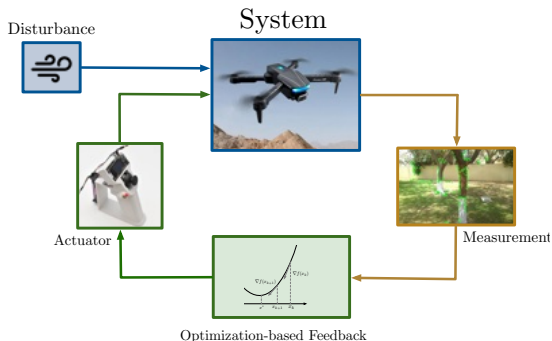
Feedback law:  $u(0) = K(x)$

# Learning-based feedback

A data-driven approach to controller design

**Assumption:** An (approximate) model of the system is available

Optimization-based control



- full knowledge of environment
- computationally complexity

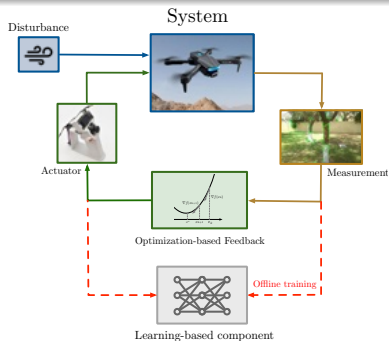
replace (some part of) the controller with a learning-based component

# Learning-based feedback

A data-driven approach to controller design

**Assumption:** An (approximate) model of the system is available

## Offline training



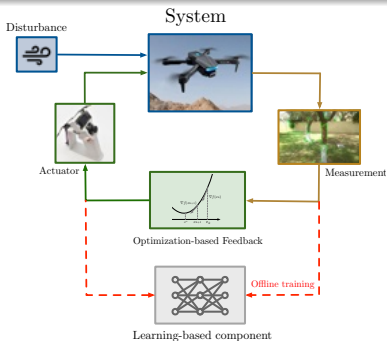
- overly conservative constraints
- solve the optimization offline
- data to train the learning algorithm

# Learning-based feedback

A data-driven approach to controller design

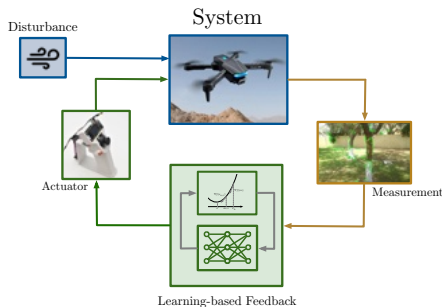
**Assumption:** An (approximate) model of the system is available

## Offline training



- overly conservative safety guarantees
- solve the optimization offline
- data to train the learning algorithm

## Online implementation



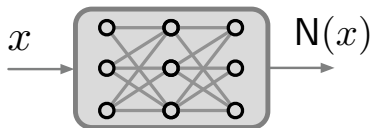
- efficient implementation
- partial knowledge of environment
- **limited** safety guarantees

# Reachability analysis

A paradigm for safety assurance

## Isolated learning component

- Robustness of learning algorithms



- An input perturbation set  $\mathcal{X}$
- Safe output domain  $\mathcal{Y}$

## Output perturbations

$$N(\mathcal{X}) = \{N(x) \mid x \in \mathcal{X}\}$$

**Goal:** ensure that  $N(\mathcal{X}) \subset \mathcal{Y}$ .

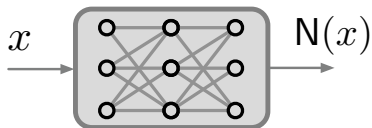


# Reachability analysis

A paradigm for safety assurance

## Isolated learning component

- Robustness of learning algorithms



- An input perturbation set  $\mathcal{X}$
- Safe output domain  $\mathcal{Y}$

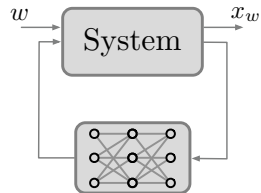
## Output perturbations

$$N(\mathcal{X}) = \{N(x) \mid x \in \mathcal{X}\}$$

**Goal:** ensure that  $N(\mathcal{X}) \subset \mathcal{Y}$ .

## Interconnected learning-based system

- Safety of closed-loop system



- An input perturbation set  $\mathcal{W}$
- Safe output domain  $\mathcal{S}$

## Reachable sets

$$\mathcal{R}(\mathcal{W}, t) = \{x_w(t) \mid w \in \mathcal{W}\}$$

**Goal:** ensure that  $\mathcal{R}(\mathcal{W}, t) \subset \mathcal{S}$

# Robustness of learning algorithms

## Verification and training

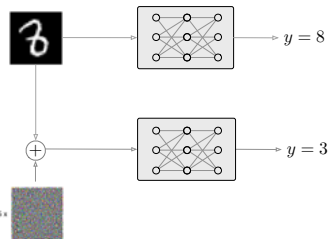
- 1 learning algorithms are fragile wrt input perturbations

### Adversarial Perturbations

Small changes in the input



Large changes in the output

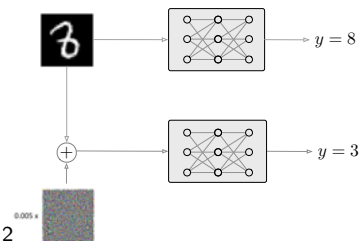
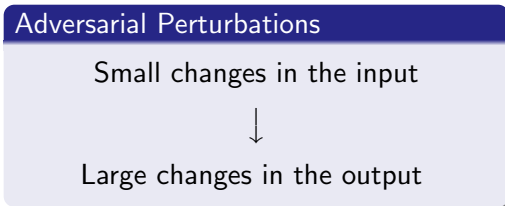


C. Szegedy and et. al. Intriguing properties of neural networks. In *ICLR*, 2

# Robustness of learning algorithms

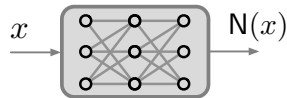
## Verification and training

- 1 learning algorithms are fragile wrt input perturbations



C. Szegedy and et. al. Intriguing properties of neural networks. In *ICLR*, 2

- 2 learning algorithms have large number of parameters and are highly nonlinear



# Reachability of learning-based systems

## The role of the structure

- Reachability of dynamical system is an old problem:  $\sim 1980$ 
  - ▶ Example approaches: [Hamilton-Jacobi](#), [Ellipsoidal methods](#)

# Reachability of learning-based systems

## The role of the structure

- Reachability of dynamical system is an old problem:  $\sim$  1980
  - ▶ Example approaches: [Hamilton-Jacobi](#), [Ellipsoidal methods](#)

not scalable to large-scale systems

# Reachability of learning-based systems

## The role of the structure

- Reachability of dynamical system is an old problem:  $\sim$  1980
  - ▶ Example approaches: [Hamilton-Jacobi](#), [Ellipsoidal methods](#)

not scalable to large-scale systems

- Reachability of learning algorithms is more recent:  $\sim$  2010
  - ▶ Example approaches: [Interval arithmetic](#), [Semi-definite programming](#)

# Reachability of learning-based systems

## The role of the structure

- Reachability of dynamical system is an old problem:  $\sim$  1980
  - ▶ Example approaches: [Hamilton-Jacobi](#), [Ellipsoidal methods](#)

not scalable to large-scale systems

- Reachability of learning algorithms is more recent:  $\sim$  2010
  - ▶ Example approaches: [Interval arithmetic](#), [Semi-definite programming](#)

- 1 structure of the learning algorithm
- 2 Interconnection structure of the system

### Structure lead to tractable algorithms

- **Contractivity**, to ensure computational efficiency
- **Mixed monotonicity**, a key property of neural network loops

- Contraction theory and mixed monotonicity
- Isolated learning algorithms
- Learning-based feedback loops



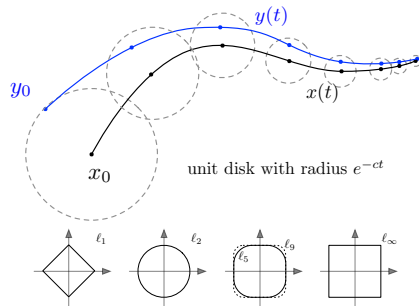
# Tool #1: Contraction theory

A framework for stability analysis

## Definition (Contraction)

$\dot{x} = f(x, u)$  is contracting wrt  $\| \cdot \|$  if

the distance between every two trajectory is decreasing exponentially with rate  $c$  wrt  $\| \cdot \|$



# Tool #1: Contraction theory

A framework for stability analysis

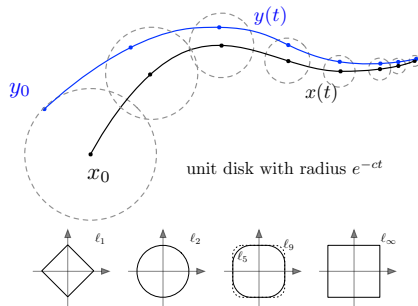
## Definition (Contraction)

$\dot{x} = f(x, u)$  is contracting wrt  $\| \cdot \|$  if

the distance between every two trajectory is decreasing exponentially with rate  $c$  wrt  $\| \cdot \|$

### Transient and asymptotic behavior:

- unique globally exponential stable equilibrium
- efficient equilibrium point computation
- input-output robustness
- modularity and interconnection properties
- ...



# Tool #1: Contraction theory

Characterization for Euclidean norms

## Main Result

$\dot{x} = f(x, u)$  is contracting wrt  $\|\cdot\|_2$  with rate  $c$  iff

**Differential:**  $D_x f(x, u)^\top + D_x f(x, u) + 2cI \preceq 0,$  for all  $x, u$

**Integral:**  $\langle f(x, u) - f(y, u), x - y \rangle \leq -c\|x - y\|_2^2,$  for all  $x, y, u$

# Tool #1: Contraction theory

Characterization for Euclidean norms

## Main Result

$\dot{x} = f(x, u)$  is contracting wrt  $\|\cdot\|_2$  with rate  $c$  iff

**Differential:**  $D_x f(x, u)^\top + D_x f(x, u) + 2cI \preceq 0$ , for all  $x, u$

**Integral:**  $\langle f(x, u) - f(y, u), x - y \rangle \leq -c\|x - y\|_2^2$ , for all  $x, y, u$

- Connection between **Euclidean contraction theory** and **monotone operator theory**

$f$  is contracting with respect to  $\|\cdot\|_2$   
iff

–  $f$  is a monotone operator with respect to the inner product  $\langle \cdot, \cdot \rangle$ .

- How about general norms?



# Tool #1: Contraction theory

## Logarithmic norm and weak pairings

### Differential condition

#### Logarithmic norm

Given a matrix  $A \in \mathbb{R}^{n \times n}$  and a norm  $\|\cdot\|$ :

$$\mu_{\|\cdot\|}(A) := \lim_{h \rightarrow 0^+} \frac{\|I_n + hA\| - 1}{h}$$

- Directional derivative of norm  $\|\cdot\|$  in direction of  $A$ ,

$$\mu_2(A) = \frac{1}{2} \lambda_{\max}(A + A^T)$$

$$\mu_1(A) = \max_j (a_{jj} + \sum_{i \neq j} |a_{ij}|)$$

$$\mu_\infty(A) = \max_i (a_{ii} + \sum_{j \neq i} |a_{ij}|)$$

<sup>1</sup>A. Davydov, S. Jafarpour, F. Bullo, TAC 2022

# Tool #1: Contraction theory

## Logarithmic norm and weak pairings

### Differential condition

#### Logarithmic norm

Given a matrix  $A \in \mathbb{R}^{n \times n}$  and a norm  $\|\cdot\|$ :

$$\mu_{\|\cdot\|}(A) := \lim_{h \rightarrow 0^+} \frac{\|I_n + hA\| - 1}{h}$$

- Directional derivative of norm  $\|\cdot\|$  in direction of  $A$ ,

$$\mu_2(A) = \frac{1}{2} \lambda_{\max}(A + A^\top)$$

$$\mu_1(A) = \max_j (a_{jj} + \sum_{i \neq j} |a_{ij}|)$$

$$\mu_\infty(A) = \max_i (a_{ii} + \sum_{j \neq i} |a_{ij}|)$$

<sup>1</sup>A. Davydov, S. Jafarpour, F. Bullo, TAC 2022

### Integral condition

#### Weak pairing<sup>1</sup>

Given a norm  $\|\cdot\|$ , the associated weak pairing is  $\llbracket \cdot, \cdot \rrbracket : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ :

- Subadditive and weakly homogeneity
- Positive definite
- Cauchy-Schwarz inequality
- $\llbracket x, x \rrbracket = \|x\|^2$

$$\llbracket x, y \rrbracket_2 = y^\top x$$

$$\llbracket x, y \rrbracket_1 = \text{sign}(y)^\top x$$

$$\llbracket x, y \rrbracket_\infty = \max_{i \in I_\infty(x)} x_i y_i$$

$$I_\infty(x) = \{i \mid |x_i| = \|x\|_\infty\}$$

# Tool #1: Contraction theory

Characterization for Non-Euclidean norms

## Main Result

$\dot{x} = f(x, u)$  is contracting wrt  $\|\cdot\|_2$  with rate  $c$  iff

**Differential:**  $D_x f(x, u)^\top + D_x f(x, u) + 2cI \preceq 0,$  for all  $x, u$

**Integral:**  $\langle f(x, u) - f(y, u), x - y \rangle \leq -c\|x - y\|_2^2,$  for all  $x, y, u$

<sup>2</sup>A. Davydov, S. Jafarpour, F. Bullo, TAC 2022

# Tool #1: Contraction theory

Characterization for Non-Euclidean norms

## Theorem<sup>2</sup>

$\dot{x} = f(x, u)$  is contracting wrt  $\|\cdot\|$  with rate  $c$  iff

**Differential:**  $\mu_{\|\cdot\|}(D_x f(x, u)) \leq -c$ , for all  $x, u$

**Integral:**  $\llbracket f(x, u) - f(y, u), x - y \rrbracket \leq -c\|x - y\|^2$ , for all  $x, y, u$

## Why non-Euclidean?

- well suited for conservative systems
- computational advantages
- structural robustness

## Non-Euclidean monotone operators

$f$  is contracting with respect to  $\|\cdot\|$   
iff

$-f$  is a monotone operator with respect to  $\llbracket \cdot, \cdot \rrbracket$ .

<sup>2</sup>A. Davydov, S. Jafarpour, F. Bullo, TAC 2022



# Tool #2: Mixed monotonicity

Cooperative and competitive dynamics

Original system

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

Embedding system

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

$g$  is a **decomposition function** s.t.

- 1  $f(x, u) = g(x, x, u, u)$  for every  $x, u$
- 2 **cooperative:**  $(\underline{x}, \underline{u}) \mapsto g(\underline{x}, \bar{x}, \underline{u}, \bar{u})$
- 3 **competitive:**  $(\bar{x}, \bar{u}) \mapsto g(\underline{x}, \bar{x}, \underline{u}, \bar{u})$

# Tool #2: Mixed monotonicity

Cooperative and competitive dynamics

Original system

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

Embedding system

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

$g$  is a **decomposition function** s.t.

- 1  $f(x, u) = g(x, x, u, u)$  for every  $x, u$
- 2 **cooperative**:  $(\underline{x}, \underline{u}) \mapsto g(\underline{x}, \bar{x}, \underline{u}, \bar{u})$
- 3 **competitive**:  $(\bar{x}, \bar{u}) \mapsto g(\underline{x}, \bar{x}, \underline{u}, \bar{u})$

Embedding system is monotone with respect to  $\leq_{SE}$

$$\begin{aligned}\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq_{SE} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \\ \text{iff} \\ x_1 \leq y_1 \text{ and } y_2 \leq x_2\end{aligned}$$

- $f$  locally Lipschitz  $\implies$  mixed monotonicity

Monotone systems are being studied by Hirsch, Smith, Sontag, Angeli, ...

### Theorem<sup>3</sup>

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

---

<sup>3</sup>S. Coogan, M. Arcak, HSCC 2015

# Tool #2: Mixed monotonicity

## Reachability analysis

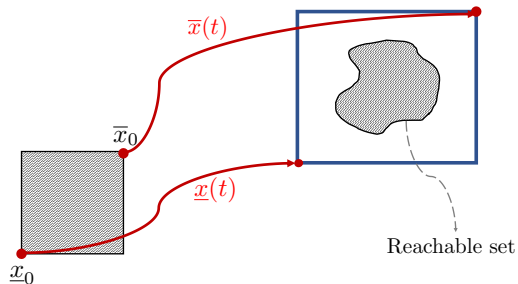
### Theorem<sup>3</sup>

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

### Embedding system

$$\dot{\underline{x}} = g(\underline{x}, \bar{x}, \underline{u}, \bar{u}), \quad \underline{x}(0) = \underline{x}_0$$

$$\dot{\bar{x}} = g(\bar{x}, \underline{x}, \bar{u}, \underline{u}), \quad \bar{x}(0) = \bar{x}_0$$



<sup>3</sup>S. Coogan, M. Arcak, HSCC 2015

# Tool #2: Mixed monotonicity

## Tight decomposition functions

How to find a decomposition function?

- decomposition function might not be unique
- different approaches exist for certain class of systems
  - ▶ Jacobian of  $f$
  - ▶ polynomial structure of  $f$

# Tool #2: Mixed monotonicity

## Tight decomposition functions

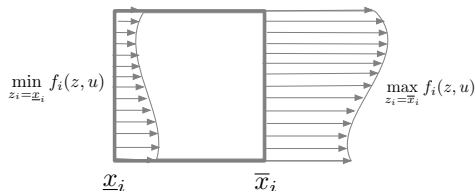
How to find a decomposition function?

- decomposition function might not be unique
- different approaches exist for certain class of systems
  - ▶ Jacobian of  $f$
  - ▶ polynomial structure of  $f$

The sharpest bounds = **tight decomposition function**

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

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



# Tool #2: Mixed monotonicity

Connection with contraction theory

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

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

## Theorem<sup>4</sup>

The original system is contracting wrt to  $\|\cdot\|_\infty$  with rate  $c$   
iff  
the **tight** embedding system is contracting wrt to  $\|\cdot\|_\infty$  with rate  $c$

<sup>4</sup>S. Jafarpour, S. Coogan, arXiv 2022

# Tool #2: Mixed monotonicity

Connection with contraction theory

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

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

## Theorem<sup>4</sup>

The original system is contracting wrt to  $\|\cdot\|_\infty$  with rate  $c$   
iff  
the **tight** embedding system is contracting wrt to  $\|\cdot\|_\infty$  with rate  $c$

The unique role of non-Euclidean  $\ell_\infty$ -norm

<sup>4</sup>S. Jafarpour, S. Coogan, arXiv 2022



# Tool #2: Mixed monotonicity

Connection with contraction theory

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

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

## Theorem<sup>4</sup>

The original system is contracting wrt to  $\|\cdot\|_\infty$  with rate  $c$   
iff  
the **tight** embedding system is contracting wrt to  $\|\cdot\|_\infty$  with rate  $c$

The unique role of non-Euclidean  $\ell_\infty$ -norm

## Corollary (for contracting systems)

Mixed-monotone reachability is **sharper** than global input-to-state bounds:

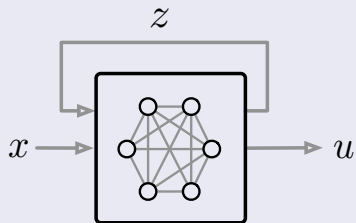
$$\|x(t) - y(t)\|_\infty \leq e^{-ct} \|x_0 - y_0\|_\infty + \frac{\ell(1-e^{-ct})}{c} \|u - w\|_\infty.$$

<sup>4</sup>S. Jafarpour, S. Coogan, arXiv 2022

- Contraction theory and mixed monotonicity
- Isolated learning algorithms
- Learning-based feedback loops

# Generalized neural networks

A general learning model via fixed-point equations





- Generalized neural networks:

$$z = \Phi(Az + Bx + b)$$

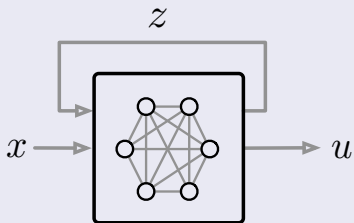
$$u = Cz + c$$

- $\Phi(y_1, \dots, y_n) = (\phi_1(y_1), \dots, \phi_n(y_n))^T$  with  $\phi_i$  satisfies  $0 \leq \frac{\phi_i(x) - \phi_i(y)}{x - y} \leq 1$ .

-  S. Bai, J. Z. Kolter, and V. Koltun. Deep equilibrium models. In *NeurIPS*, 2019
-  L. El Ghaoui, F. Gu, B. Travacca, A. Askari, and A. Y. Tsai. Implicit deep learning. *SIMODS*, 2019

# Generalized neural networks

A general learning model via fixed-point equations



- Generalized neural networks:

$$z = \Phi(Az + Bx + b)$$

$$u = Cz + c$$



- $\Phi(y_1, \dots, y_n) = (\phi_1(y_1), \dots, \phi_n(y_n))^T$  with  $\phi_i$  satisfies  $0 \leq \frac{\phi_i(x) - \phi_i(y)}{x - y} \leq 1$ .

## Notion of layer

Output is an **implicit** function of input  
(e.g., fixed-point equation, differential equations, optimization problem)

## Why implicit models?

- Representation
- Performance
- Memory

-  S. Bai, J. Z. Kolter, and V. Koltun. Deep equilibrium models. In *NeurIPS*, 2019
-  L. El Ghaoui, F. Gu, B. Travacca, A. Askari, and A. Y. Tsai. Implicit deep learning. *SIMODS*, 2019

## Main Questions

$$z = \Phi(Az + Bx + b)$$

$$u = Cz + c$$

- 1 Existence and computation of solutions?
- 2 How to estimate the input-output  $x \mapsto u$  robustness?

## Main Questions

$$z = \Phi(Az + Bx + b)$$

$$u = Cz + c$$

- 1 Existence and computation of solutions?
- 2 How to estimate the input-output  $x \mapsto u$  robustness?

## Key insight

Fixed-point equation

$$z = \Phi(Az + Bx + b)$$



Dynamical system

$$\dot{z} = -z + \Phi(Az + Bx + b)$$

**fixed-points**



**equilibrium points**

**robustness**



**forward reachability** ( $t = \infty$ )

- We can use tools from dynamical systems to study generalized neural networks

# Fixed-points of neural network

A non-Euclidean contracting approach

$$\begin{array}{ccc} \text{Fixed-point of} & \iff & \text{Equilibrium point of} \\ z = \Phi(Az + Bx + b) & & \dot{z} = -z + \Phi(Az + Bx + b) \end{array}$$

- **Contraction theory:** Sufficient condition for existence a globally stable equilibrium point.

<sup>5</sup>S. Jafarpour, A. Davydov, A. Proskurnikov, F. Bullo, NeurIPS 2022

# Fixed-points of neural network

A non-Euclidean contracting approach

$$\begin{array}{ccc} \text{Fixed-point of} & \iff & \text{Equilibrium point of} \\ z = \Phi(Az + Bx + b) & & \dot{z} = -z + \Phi(Az + Bx + b) \end{array}$$

- **Contraction theory:** Sufficient condition for existence a globally stable equilibrium point.

$$\llbracket \Phi(Az_1 + Bx + b) - \Phi(Az_2 + Bx + b), z_1 - z_2 \rrbracket_{\infty} < \|z_1 - z_2\|_{\infty}^2$$

<sup>5</sup>S. Jafarpour, A. Davydov, A. Proskurnikov, F. Bullo, NeurIPS 2022



# Fixed-points of neural network

A non-Euclidean contracting approach

$$\begin{array}{ccc} \text{Fixed-point of} & \iff & \text{Equilibrium point of} \\ z = \Phi(Az + Bx + b) & & \dot{z} = -z + \Phi(Az + Bx + b) \end{array}$$

- **Contraction theory:** Sufficient condition for existence a globally stable equilibrium point.

$$a_{ii} + \sum_{j \neq i} |a_{ij}| < 1 \implies \|\Phi(Az_1 + Bx + b) - \Phi(Az_2 + Bx + b), z_1 - z_2\|_{\infty} < \|z_1 - z_2\|_{\infty}^2$$

<sup>5</sup>S. Jafarpour, A. Davydov, A. Proskurnikov, F. Bullo, NeurIPS 2022

# Fixed-points of neural network

A non-Euclidean contracting approach

$$\begin{array}{ccc} \text{Fixed-point of} & \iff & \text{Equilibrium point of} \\ z = \Phi(Az + Bx + b) & & \dot{z} = -z + \Phi(Az + Bx + b) \end{array}$$

- **Contraction theory:** Sufficient condition for existence a globally stable equilibrium point.

$$a_{ii} + \sum_{j \neq i} |a_{ij}| < 1 \implies \|\Phi(Az_1 + Bx + b) - \Phi(Az_2 + Bx + b), z_1 - z_2\|_\infty < \|z_1 - z_2\|_\infty^2$$

## Theorem<sup>5</sup>

If  $a_{ii} + \sum_{j \neq i} |a_{ij}| < 1$  then

- 1  $z = \Phi(Az + Bx + b)$  has a unique solution  $z_x^*$
- 2  $z_x^*$  can be computed using average iterations for  $z = \Phi(Az + Bx + b)$

<sup>5</sup>S. Jafarpour, A. Davydov, A. Proskurnikov, F. Bullo, NeurIPS 2022

# Robustness of neural network

A mixed monotone contracting approach

$$\begin{array}{ccc} \text{robustness of} & \iff & \text{forward reachability of} \\ z = \Phi(Az + Bx + b) & & \dot{z} = -z + \Phi(Az + Bx + b) \end{array}$$

<sup>6</sup>S.Jafarpour, M. Abate, A. Davydov, F. Bullo, S. Coogan, L4DC 2022

# Robustness of neural network

A mixed monotone contracting approach

$$\begin{array}{ccc} \text{robustness of} & \iff & \text{forward reachability of} \\ z = \Phi(Az + Bx + b) & & \dot{z} = -z + \Phi(Az + Bx + b) \end{array}$$

- **Metzler/non-Metzler** decomposition:  $A = \lceil A \rceil^{\text{Mzl}} + \lfloor A \rfloor^{\text{Mzl}}$

- Example:  $A = \begin{bmatrix} 2 & 0 & -1 \\ 1 & -3 & 0 \\ 0 & 0 & 1 \end{bmatrix} \implies \lceil A \rceil^{\text{Mzl}} = \begin{bmatrix} 2 & 0 & 0 \\ 1 & -3 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \lfloor A \rfloor^{\text{Mzl}} = \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$

<sup>6</sup>S.Jafarpour, M. Abate, A. Davydov, F. Bullo, S. Coogan, L4DC 2022

# Robustness of neural network

A mixed monotone contracting approach

$$\begin{array}{ccc} \text{robustness of} & \iff & \text{forward reachability of} \\ z = \Phi(Az + Bx + b) & & \dot{z} = -z + \Phi(Az + Bx + b) \end{array}$$

- **Metzler/non-Metzler** decomposition:  $A = \lceil A \rceil^{\text{Mzl}} + \lfloor A \rfloor^{\text{Mzl}}$

- Example:  $A = \begin{bmatrix} 2 & 0 & -1 \\ 1 & -3 & 0 \\ 0 & 0 & 1 \end{bmatrix} \implies \lceil A \rceil^{\text{Mzl}} = \begin{bmatrix} 2 & 0 & 0 \\ 1 & -3 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \lfloor A \rfloor^{\text{Mzl}} = \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$

## Theorem<sup>6</sup>

The neural network is mixed monotone with the tight decomposition function:

$$G(\underline{z}, \bar{z}, \underline{x}, \bar{x}) = -\underline{z} + \Phi(\lceil A \rceil^{\text{Mzl}} \underline{z} + \lfloor A \rfloor^{\text{Mzl}} \bar{z} + [B]^+ \underline{x} + [B]^- \bar{x} + b)$$

<sup>6</sup>S.Jafarpour, M. Abate, A. Davydov, F. Bullo, S. Coogan, L4DC 2022

# Robustness of neural network

A mixed monotone contracting approach

## Theorem<sup>7</sup>

If  $a_{ii} + \sum_{j \neq i} |a_{ij}| < 1$  and  $x \in [\underline{x}, \bar{x}]$

- 1  $z = \Phi(Az + Bx + b)$  has a unique solution  $z_u^*$
- 2  $\begin{bmatrix} \underline{z} \\ \bar{z} \end{bmatrix} = \begin{bmatrix} G(\underline{z}, \bar{z}, \underline{x}, \bar{x}) \\ G(\bar{z}, \underline{z}, \bar{x}, \underline{x}) \end{bmatrix}$  has a unique solution  $\begin{bmatrix} \underline{z}^* \\ \bar{z}^* \end{bmatrix}$
- 3  $\underbrace{([C]^+ [C]^-) \begin{bmatrix} \underline{z}^* \\ \bar{z}^* \end{bmatrix}}_u + c \leq u \leq \underbrace{([C]^- [C]^+) \begin{bmatrix} \underline{z}^* \\ \bar{z}^* \end{bmatrix}}_{\bar{u}} + c$

<sup>7</sup>S.Jafarpour, M. Abate, A. Davydov, F. Bullo, S. Coogan, L4DC 2022

# Robustness of neural network

A mixed monotone contracting approach

## Theorem<sup>7</sup>

If  $a_{ii} + \sum_{j \neq i} |a_{ij}| < 1$  and  $x \in [\underline{x}, \bar{x}]$

①  $z = \Phi(Az + Bx + b)$  has a unique solution  $z_u^*$

②  $\begin{bmatrix} \underline{z} \\ \bar{z} \end{bmatrix} = \begin{bmatrix} G(\underline{z}, \bar{z}, \underline{x}, \bar{x}) \\ G(\bar{z}, \underline{z}, \bar{x}, \underline{x}) \end{bmatrix}$  has a unique solution  $\begin{bmatrix} \underline{z}^* \\ \bar{z}^* \end{bmatrix}$

③  $\underbrace{([C]^+ [C]^-) \begin{bmatrix} \underline{z}^* \\ \bar{z}^* \end{bmatrix}}_u + c \leq u \leq \underbrace{([C]^- [C]^+) \begin{bmatrix} \underline{z}^* \\ \bar{z}^* \end{bmatrix}}_{\bar{u}} + c$

- **Verification:** find robustness margin of generalized neural networks
- **Training:** design robust generalized neural networks

- $a_{ii} + \sum_{j \neq i} |a_{ij}| < 1$  as a constraint to the training problem
- a regularization term  $\mathcal{R}(\underline{u}, \bar{u})$  to the training cost

<sup>7</sup>S.Jafarpour, M. Abate, A. Davydvov, F. Bullo, S. Coogan, L4DC 2022

# Numerical experiments

## MNIST dataset classification

- MNIST dataset:  $28 \times 28$  pixel handwritten digits between 0 – 9.
- hidden layer of neural network  $n = 100$
- $\epsilon =$  size of perturbation,  $\mathcal{X} = [x - \epsilon \mathbb{1}_{784}, x + \epsilon \mathbb{1}_{784}]$ .

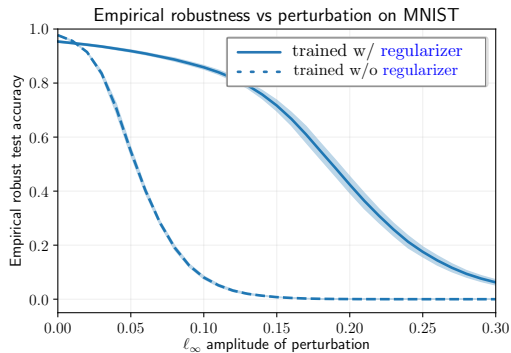
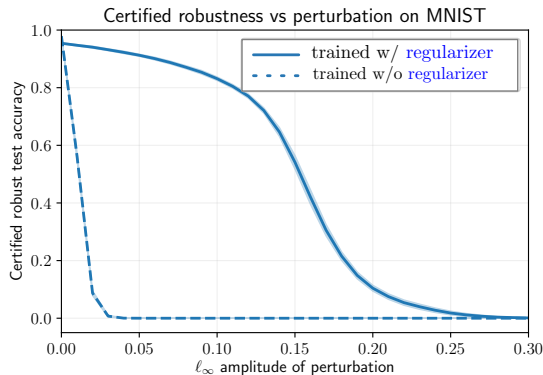




# Numerical experiments

## MNIST dataset classification

- MNIST dataset:  $28 \times 28$  pixel handwritten digits between 0 – 9.
- hidden layer of neural network  $n = 100$
- $\epsilon =$  size of perturbation,  $\mathcal{X} = [x - \epsilon \mathbb{1}_{784}, x + \epsilon \mathbb{1}_{784}]$ .



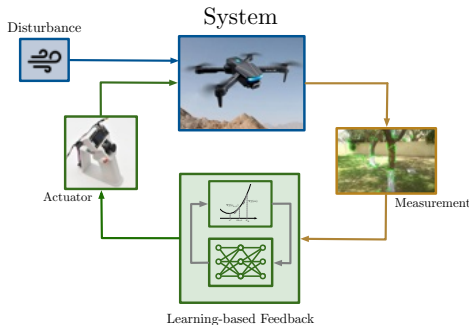
- Certified robustness = all the elements of  $[\underline{u}(\epsilon), \bar{u}(\epsilon)]$  classify as the correct digit
- Empirical robustness = Projected Gradient Descent (PGD) attack

- Contraction theory and mixed monotonicity
- Isolated learning algorithms
- Learning-based feedback loops

# Learning-based feedback

## Safety guarantees for the feedback loops

**Run-Time Assurance mechanism (RTA):** monitor + predict



Closed-loop safety  
for  $t \mapsto t + T$

$$x \leq \underline{x}_b$$

OR

$$x \geq \bar{x}_b$$



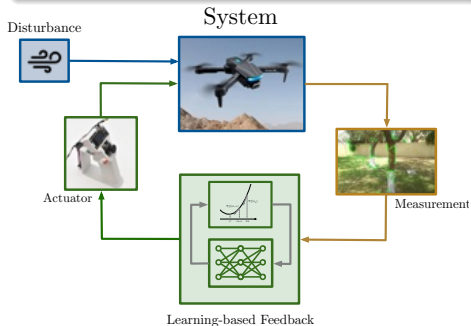
Online Safety Mechanism

# Learning-based feedback

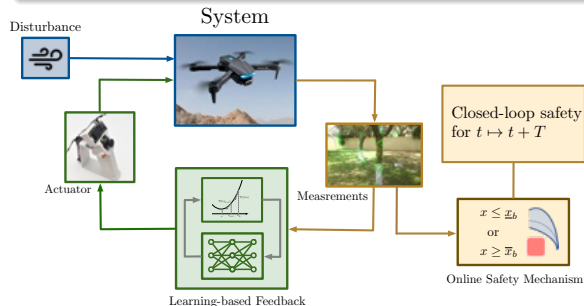
Safety guarantees for the feedback loops

**Run-Time Assurance mechanism (RTA):** monitor + predict

Closed-loop system without RTA



Closed-loop system with RTA

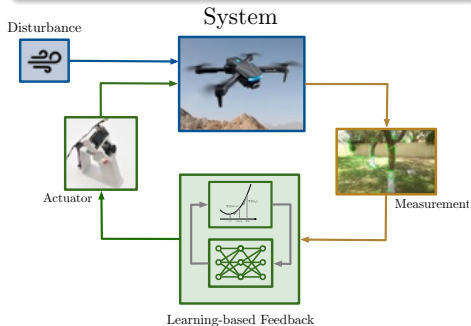


# Learning-based feedback

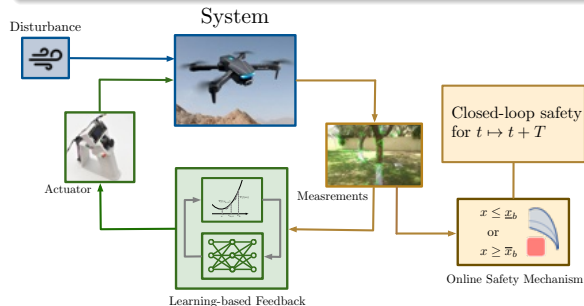
Safety guarantees for the feedback loops

**Run-Time Assurance mechanism (RTA):** monitor + predict

Closed-loop system without RTA



Closed-loop system with RTA



**Mixed monotonicity** offers a computationally efficient framework

# Design of RTA mechanism

A mixed monotone compositional approach

**Idea:** find a decomposition function for closed-loop system

# Design of RTA mechanism

A mixed monotone compositional approach

**Idea:** find a decomposition function for closed-loop system

System dynamics is mixed monotone with a decomposition function  $g$

# Design of RTA mechanism

A mixed monotone compositional approach

**Idea:** find a decomposition function for closed-loop system

System dynamics is mixed monotone with a decomposition function  $g$

A neural network verification algorithm, for all

$$x \in [\underline{x}, \bar{x}],$$

$$\underline{L}(\underline{x}, \bar{x}) \leq \mathbf{N}(x) \leq \bar{L}(\underline{x}, \bar{x})$$



# Design of RTA mechanism

A mixed monotone compositional approach

**Idea:** find a decomposition function for closed-loop system

System dynamics is mixed monotone with a decomposition function  $g$

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

Embedding system

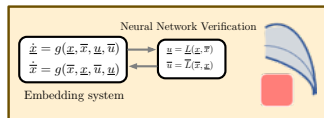
A neural network verification algorithm, for all  $x \in [\underline{x}, \bar{x}]$ ,

$$\underline{L}(\underline{x}, \bar{x}) \leq N(x) \leq \bar{L}(\underline{x}, \bar{x})$$

$$\begin{aligned}\underline{u} &= \underline{L}(\underline{x}, \bar{x}) \\ \bar{u} &= \bar{L}(\bar{x}, \underline{x})\end{aligned}$$

Neural Network Verification

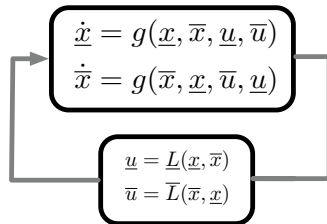
RTA mechanism = Embedding system + Neural network verification



# Design of RTA mechanism

A mixed monotone compositional approach

For  $x \in [\underline{x}, \bar{x}]$  feed the output of neural network verification algorithm into the embedding system

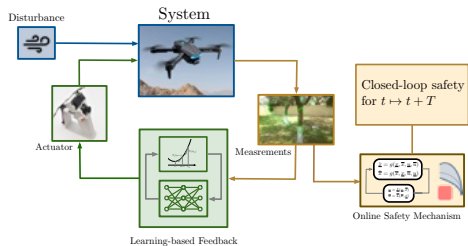


## Theorem<sup>8</sup>

The mapping

$$h(x, \bar{x}) = g(x, \bar{x}, \underline{L}(x, \bar{x}), \bar{L}(x, \bar{x}))$$

is a decomposition function for closed-loop system

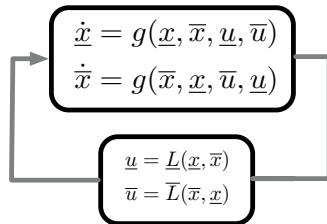


<sup>8</sup>S.Jafarpour, A. Harapanahalli, S. Coogan, arXiv 2022

# Design of RTA mechanism

A mixed monotone compositional approach

For  $x \in [\underline{x}, \bar{x}]$  feed the output of neural network verification algorithm into the embedding system

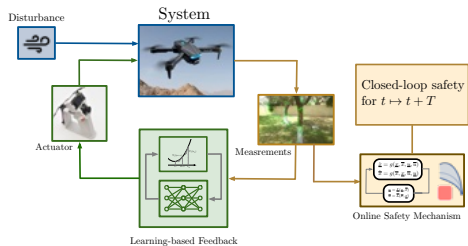


For the dynamical system

$$\dot{x} = g(x, \bar{x}, \underline{L}(x, \bar{x}), \bar{L}(x, \bar{x})) \quad x(0) = x_0$$

$$\dot{\bar{x}} = g(\bar{x}, x, \bar{L}(x, \bar{x}), \underline{L}(x, \bar{x})) \quad \bar{x}(0) = \bar{x}_0$$

we have  $\mathcal{R}([\underline{x}_0, \bar{x}_0], t) \subseteq [x(t), \bar{x}(t)]$  for all  $t \geq 0$ .



<sup>8</sup>S.Jafarpour, A. Harapanahalli, S. Coogan, arXiv 2022

# Vehicle experiment

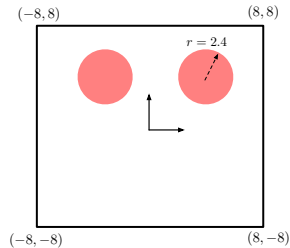
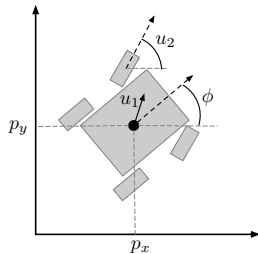
## Neural network controller

### Dynamics of vehicle

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

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

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



**Goal:** steer the vehicle to the origin avoiding the obstacles

# Vehicle experiment

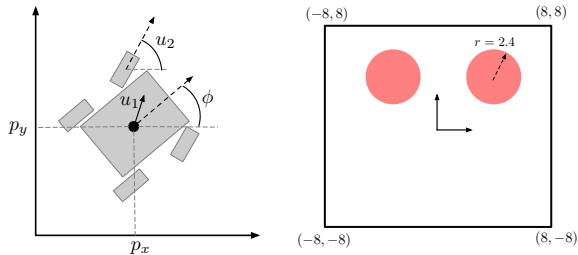
## Neural network controller

### Dynamics of vehicle

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

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

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



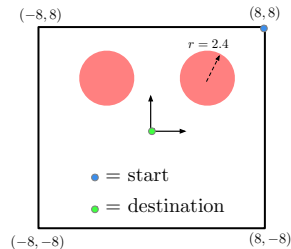
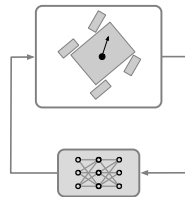
**Goal:** steer the vehicle to the origin avoiding the obstacles

- **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

# Vehicle experiment

## Design of RTA mechanism

- start from  $(8, 8)$  toward  $(0, 0)$
- $\mathcal{X}_0 = [\underline{x}_0, \bar{x}_0]$  with
$$\underline{x}_0 = (7.9 \quad 7.9 \quad -\frac{2\pi}{3} - 0.01 \quad 1.99)^\top$$
$$\bar{x}_0 = (8.1 \quad 8.1 \quad -\frac{2\pi}{3} + 0.01 \quad 2.01)^\top$$
- CROWN<sup>9</sup> for verification of neural network
- partition the states to improve accuracy

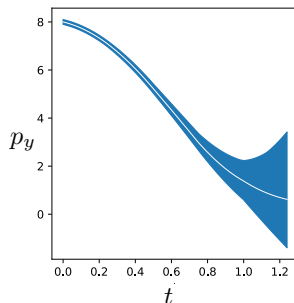
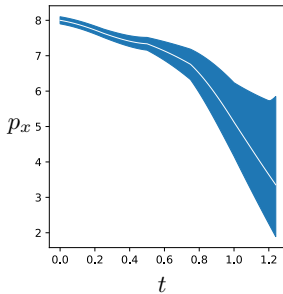
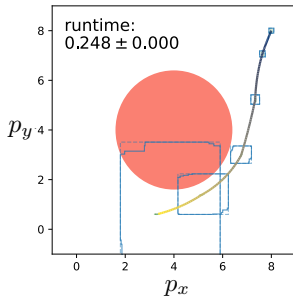
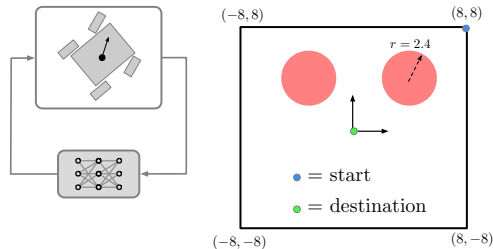


<sup>9</sup>H. Zhang, T-W. Weng, P-Y. Chen, C-J. Hsieh, L. Daniel, NeurIPS 2018

# Vehicle experiment

## Design of RTA mechanism

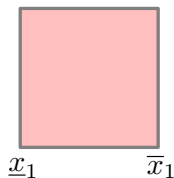
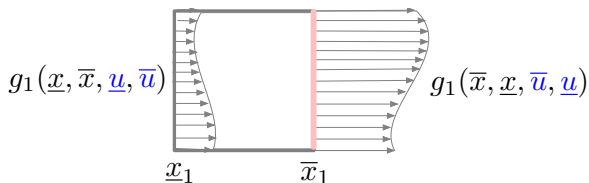
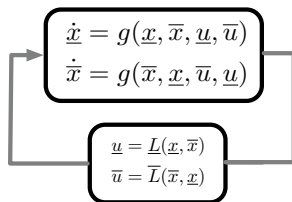
- start from  $(8, 8)$  toward  $(0, 0)$
- $\mathcal{X}_0 = [\underline{x}_0, \bar{x}_0]$  with
$$\underline{x}_0 = (7.9 \quad 7.9 \quad -\frac{2\pi}{3} - 0.01 \quad 1.99)^\top$$
$$\bar{x}_0 = (8.1 \quad 8.1 \quad -\frac{2\pi}{3} + 0.01 \quad 2.01)^\top$$
- CROWN<sup>9</sup> for verification of neural network
- partition the states to improve accuracy



<sup>9</sup>H. Zhang, T-W. Weng, P-Y. Chen, C-J. Hsieh, L. Daniel, NeurIPS 2018

# Design of RTA mechanism

The source of conservativeness



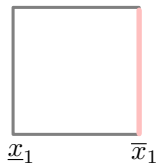
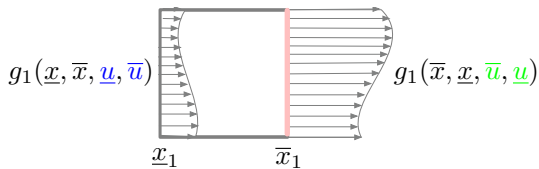
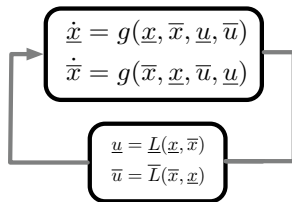
$$\underline{u} = \underline{L}(\underline{x}, \bar{x}) \leq N(x), \quad \text{for all } x \in [\underline{x}, \bar{x}]$$

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



# Design of RTA mechanism

The source of conservativeness



$$\underline{u} = \underline{L}(\underline{x}, \bar{x}) \leq N(x) \quad \text{for all } x \in \left[ \begin{pmatrix} \bar{x}_1 \\ \underline{x}_2 \end{pmatrix}, \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \end{pmatrix} \right]$$

$$N(x) \leq \bar{L}(\underline{x}, \bar{x}) = \bar{u} \quad \text{for all } x \in \left[ \begin{pmatrix} \bar{x}_1 \\ \underline{x}_2 \end{pmatrix}, \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \end{pmatrix} \right]$$

# Design of RTA mechanism

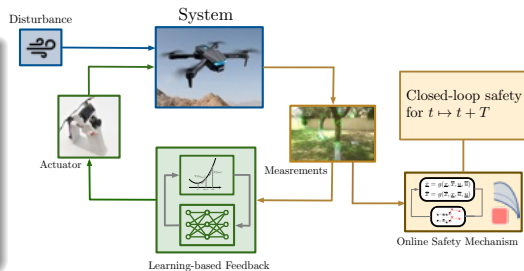
## Revised compositional approach

### Theorem<sup>10</sup>

The mapping

$$h_i(\underline{x}, \bar{x}) = g_i(\underline{x}, \bar{x}, \underline{L}(\underline{x}, \bar{x}), \bar{L}(\underline{x}, \bar{x}))$$

is a decomposition function for closed-loop system



<sup>10</sup>S.Jafarpour, A. Harapanahalli, S. Coogan, arXiv 2022

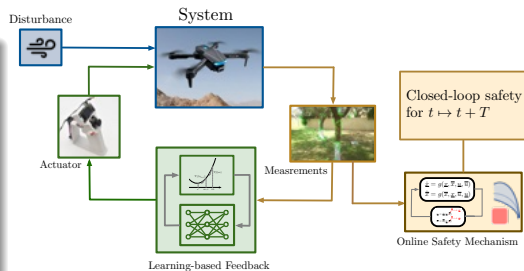
# Design of RTA mechanism

## Revised compositional approach

For the dynamical system

$$\begin{aligned} \dot{\underline{x}}_i &= g_i(\underline{x}, \bar{x}, \underline{L}(\underline{x}, \bar{x}), \bar{L}(\underline{x}, \bar{x})) & \underline{x}_i(0) &= (\underline{x}_0)_i \\ \dot{\bar{x}}_i &= g_i(\bar{x}, \underline{x}, \bar{L}(\bar{x}, \underline{x}), \underline{L}(\bar{x}, \underline{x})) & \bar{x}_i(0) &= (\bar{x}_0)_i \end{aligned}$$

we have  $\mathcal{R}(\mathcal{X}_0, t) \subseteq [\underline{x}(t), \bar{x}(t)]$  for all  $t \geq 0$ .



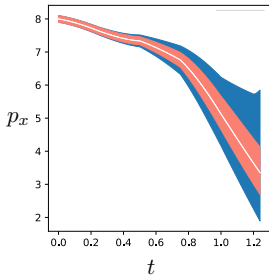
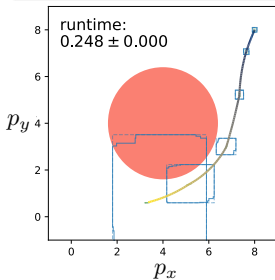
<sup>10</sup>S.Jafarpour, A. Harapanahalli, S. Coogan, arXiv 2022

# Vehicle experiment revisited

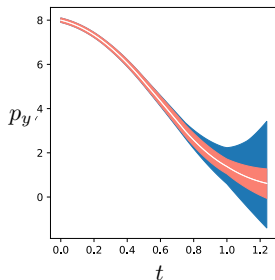
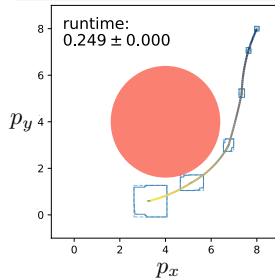
## Design of RTA mechanism

- CROWN to verify the neural network
- partition the states to improve accuracy
- blue = reachable set of **old**
- red = reachable set of **new**
- very small increase in computational time
- significant improvement in accuracy
- New decomposition function **certify** that closed-loop system is avoiding the obstacle

### Old decomposition function



### New decomposition function



- A computationally efficient framework for reachability of dynamical systems
- Exploit the structure in **isolated neural networks** and **Neural network feedback loops**
- Suitable for generalized neural networks
  - ▶ Sufficient conditions for their well-posedness
  - ▶ Hyper-rectangular over-approximation of reachable sets
- Safety assurance mechanism for monitoring neural network feedback loops
- Mixed monotonicity to design a computationally efficient run-time assurance algorithm

# Acknowledgment

## Collaborators



Alexander Davydov

UCSB



Matthew Abate

Georgia Tech



Pedro Cisneros-Velarde

UIUC



Akash Harapanahalli

Georgia Tech



Anton Proskurnikov

Politecnico di Torino



Francesco Bullo

UCSB



Samuel Coogan

Georgia Tech

Thank you for your attention!

# Generalized Structure

## Comparison with feedforward neural networks

- Feedforward neural networks:

$$x^{(\ell+1)} = \Phi(A_\ell x^{(\ell)} + b_\ell), \quad x^{(0)} = u$$
$$y = A_k x^{(k)} + b_k$$

$$\begin{bmatrix} x^{(k)} \\ x^{(k-1)} \\ \vdots \\ x^{(2)} \\ x^{(1)} \end{bmatrix} = \Phi \left( \begin{bmatrix} 0 & A_{k-1} & 0 & \dots & 0 \\ 0 & 0 & A_{k-2} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & A_1 \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} x^{(k)} \\ x^{(k-1)} \\ \vdots \\ x^{(2)} \\ x^{(1)} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ A_0 \end{bmatrix} u + \begin{bmatrix} b_{k-1} \\ b_{k-2} \\ \vdots \\ b_1 \\ b_0 \end{bmatrix} \right)$$
$$y = [A_k \quad 0 \quad 0 \quad \dots \quad 0] \begin{bmatrix} x^{(k)} \\ x^{(k-1)} \\ \vdots \\ x^{(2)} \\ x^{(1)} \end{bmatrix} + b_k$$

- Generalized neural networks:

$$x = \Phi(Ax + Bu + b)$$
$$y = Cx + c$$

$$\begin{bmatrix} x^{(k)} \\ x^{(k-1)} \\ \vdots \\ x^{(2)} \\ x^{(1)} \end{bmatrix} = \Phi \left( \begin{bmatrix} A_{11} & A_{12} & A_{13} & \dots & A_{1k} \\ A_{21} & A_{22} & A_{23} & \dots & A_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ A_{(k-1)1} & A_{(k-1)2} & A_{(k-1)3} & \dots & A_{(k-1)k} \\ A_{k1} & A_{k2} & A_{k3} & \dots & A_{kk} \end{bmatrix} \begin{bmatrix} x^{(k)} \\ x^{(k-1)} \\ \vdots \\ x^{(2)} \\ x^{(1)} \end{bmatrix} + \begin{bmatrix} B_1 \\ B_2 \\ \vdots \\ B_{k-1} \\ B_k \end{bmatrix} u + \begin{bmatrix} b_{k-1} \\ b_{k-2} \\ \vdots \\ b_1 \\ b_0 \end{bmatrix} \right)$$
$$y = [C_1 \quad C_2 \quad C_3 \quad \dots \quad C_k] \begin{bmatrix} x^{(k)} \\ x^{(k-1)} \\ \vdots \\ x^{(2)} \\ x^{(1)} \end{bmatrix} + c$$



# Generalized Structure

## Comparison with feedforward neural networks

- Feedforward neural networks:

$$x^{(\ell+1)} = \Phi(A_\ell x^{(\ell)} + b_\ell), \quad x^{(0)} = u$$

$$y = A_k x^{(k)} + b_k$$

$$x = \Phi \left( \begin{array}{|c|} \hline \square \\ \hline \end{array} x + \begin{array}{|c|} \hline \square \\ \hline \end{array} u + b \right)$$

$$y = \begin{array}{|c|} \hline \square \\ \hline \end{array} x + b_k$$

- Generalized neural networks:

$$x = \Phi(Ax + Bu + b)$$

$$y = Cx + c$$

$$x = \Phi \left( \begin{array}{|c|} \hline \square \\ \hline \end{array} x + \begin{array}{|c|} \hline \square \\ \hline \end{array} u + b \right)$$

$$y = \begin{array}{|c|} \hline \square \\ \hline \end{array} x + c$$

# Training of INNs

Promoting robustness via regularization

- 1 loss function  $\mathcal{L}$  and training data  $(\hat{u}_i, \hat{y}_i)_{i=1}^N$
- 2  $\epsilon =$  size of  $\ell_\infty$ -perturbation in input:  $\mathcal{U} = \underbrace{[u - \epsilon \mathbb{1}_r, u + \epsilon \mathbb{1}_r]}_{\underline{u}} \quad \underbrace{\quad}_{\bar{u}}$

## Training INNs

$$\min_{A,B,b,c} \sum_{i=1}^N \mathcal{L}(\hat{y}_i, Cx_i + c)$$

$$x_i = \Phi(Ax_i + B\hat{u}_i + b),$$

$$a_{ii} + \sum_{j=1} |a_{ij}| \leq \gamma \quad \text{well-posedness}$$

## Training FFNNs

$$\min_{A,B,b,c} \sum_{i=1}^N \mathcal{L}(\hat{y}_i, Cx_i^{(k)} + c)$$

$$x_i^{(\ell+1)} = \Phi(A_\ell x_i^{(\ell)} + b_\ell), \quad \ell \in \{1, \dots, k-1\}$$

# Training of INNs

Promoting robustness via regularization

- 1 loss function  $\mathcal{L}$  and training data  $(\hat{u}_i, \hat{y}_i)_{i=1}^N$
- 2  $\epsilon =$  size of  $\ell_\infty$ -perturbation in input:  $\mathcal{U} = \underbrace{[u - \epsilon \mathbb{1}_r, u + \epsilon \mathbb{1}_r]}_{\underline{u}} \underbrace{\hspace{1.5cm}}_{\bar{u}}$

Our main result

output  $y \in [\underline{y}(\epsilon), \bar{y}(\epsilon)]$

## Training INNs

$$\min_{A,B,b,c} \sum_{i=1}^N \mathcal{L}(\hat{y}_i, Cx_i + c) + \underbrace{\kappa \mathcal{R}(\underline{y}_i(\epsilon), \bar{y}_i(\epsilon))}_{\text{robustness}}$$
$$x_i = \Phi(Ax_i + B\hat{u}_i + b),$$
$$a_{ii} + \sum_{j=1}^r |a_{ij}| \leq \gamma < 1 \quad \text{well-posedness}$$

## Training FFNNs (S. Gowal, et. al., 2018)

$$\min_{A,B,b,c} \sum_{i=1}^N \mathcal{L}(\hat{y}_i, Cx_i^{(k)} + c) + \underbrace{\kappa \mathcal{R}(\underline{y}_i(\epsilon), \bar{y}_i(\epsilon))}_{\text{robustness}}$$
$$x_i^{(\ell+1)} = \Phi(A_\ell x_i^{(\ell)} + b_\ell), \quad \ell \in \{1, \dots, k-1\}$$

- $\mathcal{R}(\underline{y}(\epsilon), \bar{y}(\epsilon))$  uses  $\underline{y}(\epsilon)$  and  $\bar{y}(\epsilon)$  to estimate robustness margin
- $\kappa, \epsilon, \gamma$  are hyperparameters