

Interval Reachability of Nonlinear Dynamical Systems with Neural Network Controllers

Saber Jafarpour^{*1}, Akash Harapanahalli^{*1}, and Samuel Coogan¹

(¹) Georgia Institute of Technology, {saber,aharapan,sam.coogan}@gatech.edu

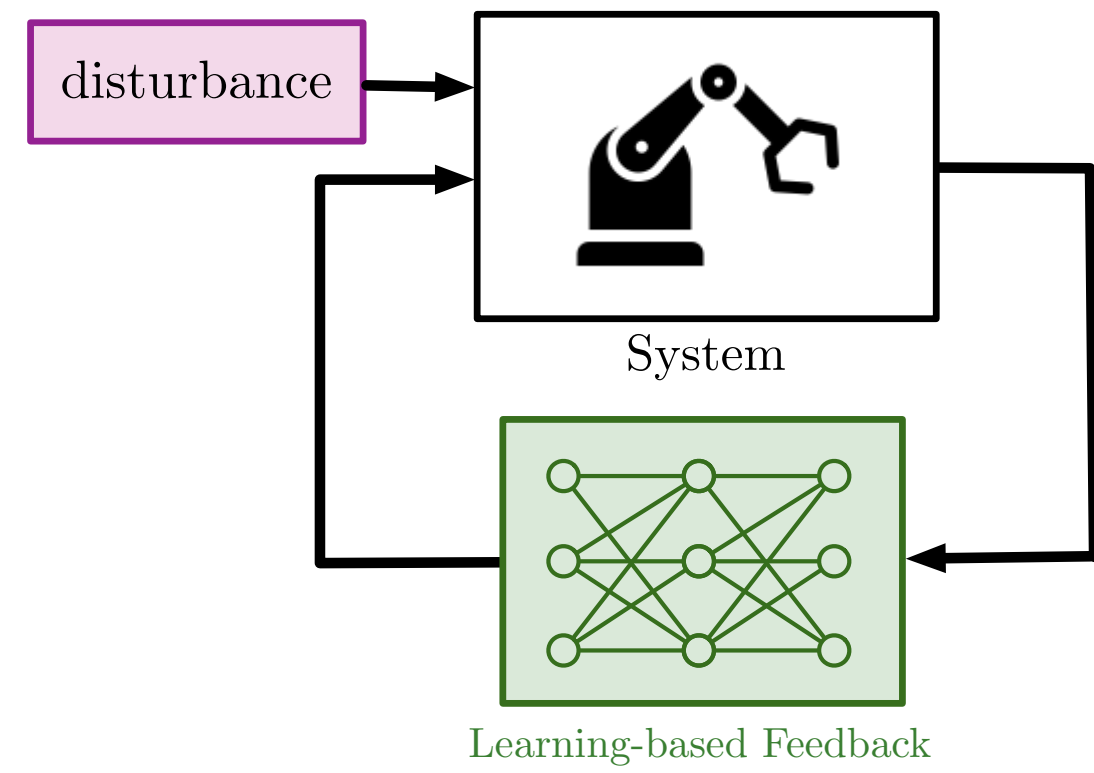


Neural Network Controllers

- Neural networks are deployed as controllers in safety-critical applications (self driving vehicle and mobile robots).

Problem

Under uncertainty, ensure safety of the closed-loop system.

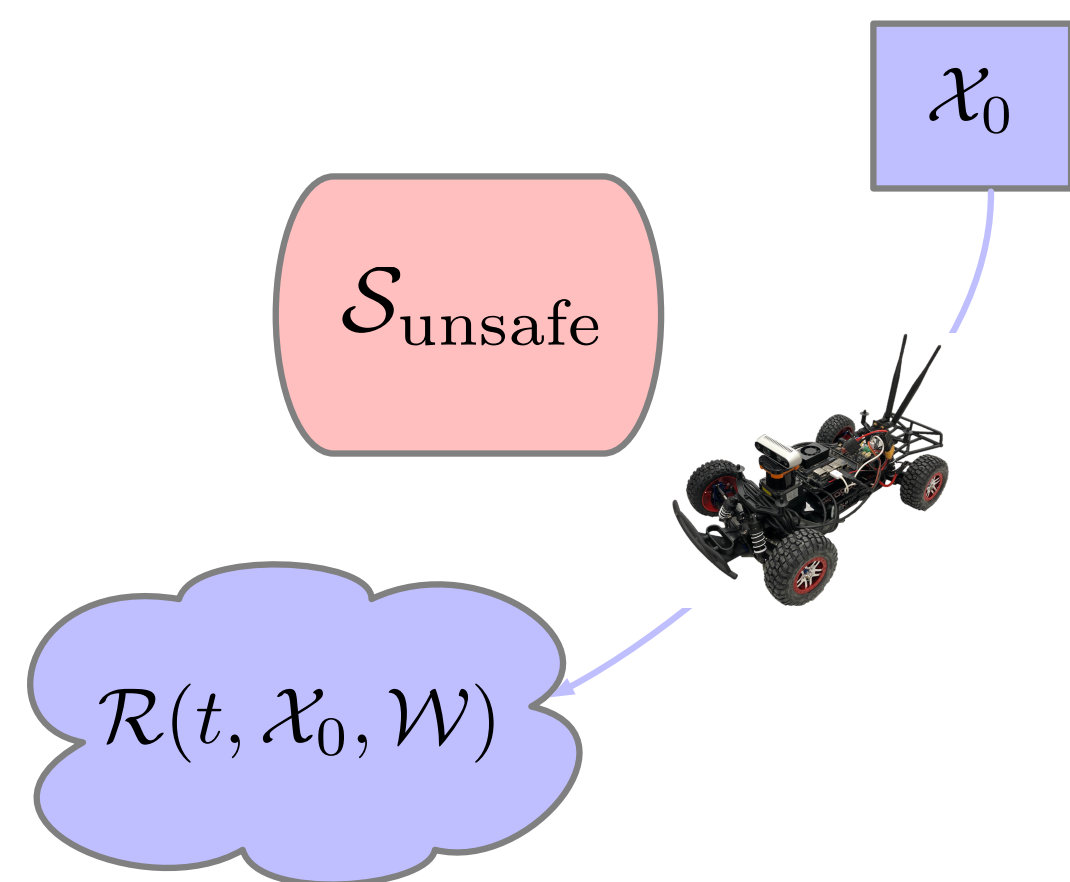


Challenges:

- Neural networks are brittle with respect to input perturbations
- The error can compound in the closed-loop interconnection.

Safety Verification via Reachability Analysis

- The disturbance \mathcal{W}
- The initial uncertainty \mathcal{X}_0
- The *reachable set*
 $\mathcal{R}(t, \mathcal{X}_0, \mathcal{W}) = \{x(t) \text{ is a trajectory}\}$
- Unsafe set $\mathcal{S}_{\text{unsafe}} \subseteq \mathbb{R}^n$



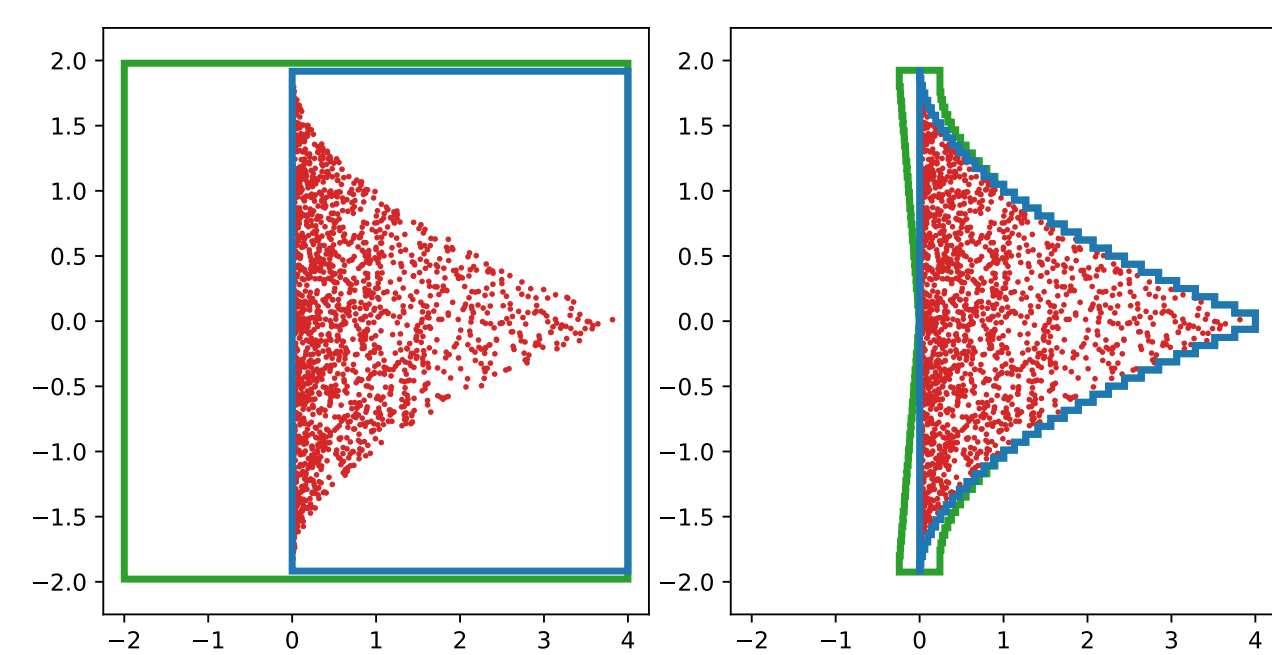
Find an over-approximation $\bar{\mathcal{R}}(t, \mathcal{X}_0, \mathcal{W})$, and check if

$$\bar{\mathcal{R}}(t, \mathcal{X}_0, \mathcal{W}) \cap \mathcal{S}_{\text{unsafe}} = \emptyset$$

Interval Analysis

Goal: over-approximate the output of a mapping using intervals.

- $G = \left[\frac{G}{G} \right]$ is an *inclusion function* for g if for every $x \in [\underline{x}, \bar{x}]$,
 $g(x) \in [G(\underline{x}, \bar{x}), \bar{G}(\underline{x}, \bar{x})]$.
- Inclusion functions can capture localized behaviors of functions—they preserve the structure when the intervals are small.
- Different approaches exist for constructing inclusion functions.

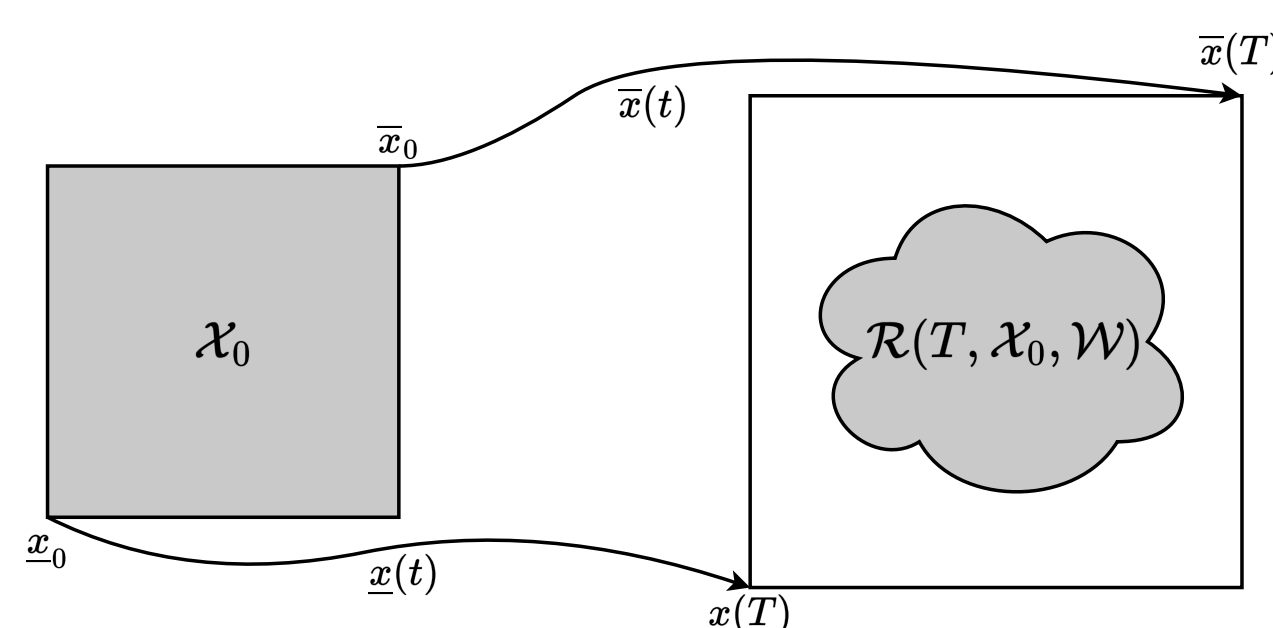


Interval Reachability of Dynamical Systems

- Consider $\dot{x} = f(x, w)$ with an inclusion function $F = \left[\frac{F}{F} \right]$ for f , the *embedding system* is

$$\begin{aligned} \dot{\underline{x}}_i &= \underline{F}_i(\underline{x}, \bar{x}_{i:\bar{x}}, \underline{w}, \bar{w}), \\ \dot{\bar{x}}_i &= \bar{F}_i(\underline{x}_{i:\underline{x}}, \bar{x}, \underline{w}, \bar{w}) \end{aligned}$$

A single trajectory of the embedding system provides lower bound \underline{x} and upper bound \bar{x} on reachable set of original system at time t .



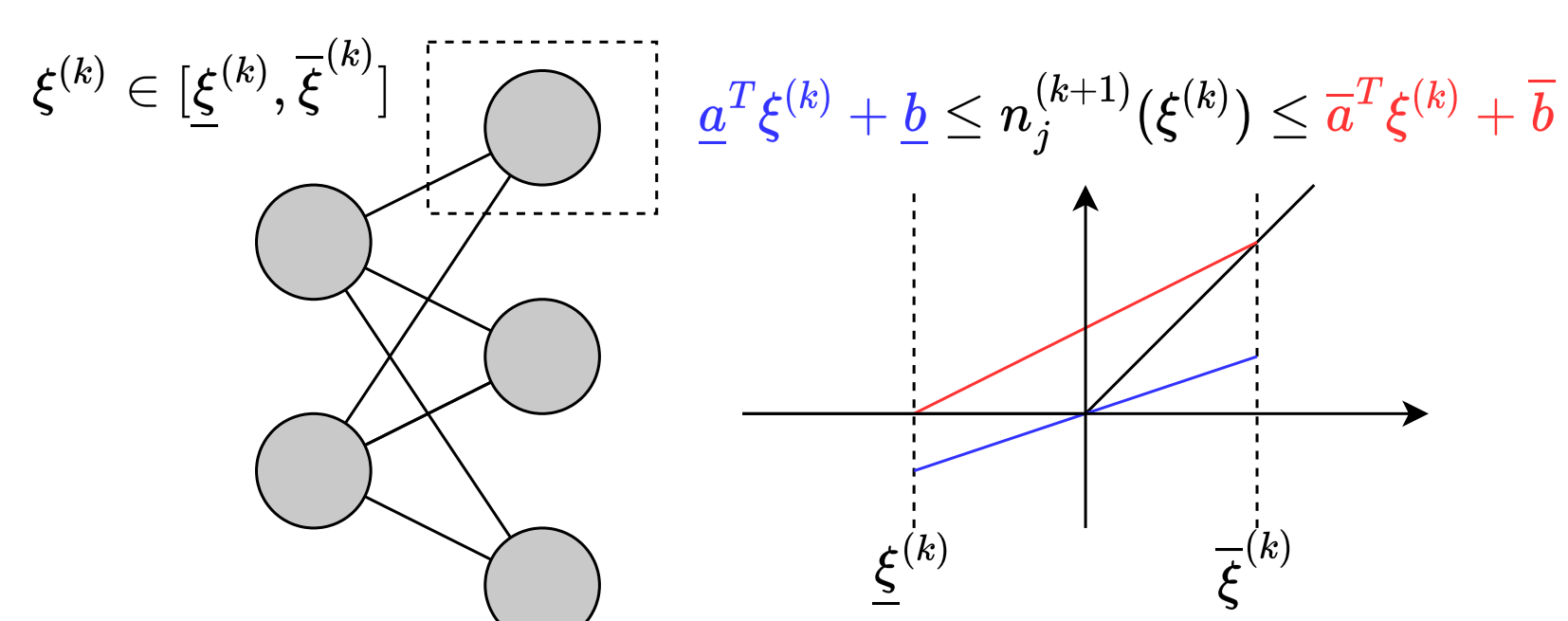
Main Question: How to construct an embedding system for the neural network controlled system?

Inclusion Functions for Neural Networks

Find \underline{N}, \bar{N} such that for every $x \in [\underline{x}, \bar{x}] \subseteq [y, \bar{y}]$,

$$\underline{N}_{[y, \bar{y}]}(\underline{x}, \bar{x}) \leq N(x) \leq \bar{N}_{[y, \bar{y}]}(\underline{x}, \bar{x}).$$

- \underline{N}, \bar{N} : neural network verification algorithms such as CROWN, IBP, LipSDP.
- CROWN [1] provides linear bounds $\underline{N}_{[y, \bar{y}]}$ and $\bar{N}_{[y, \bar{y}]}$.

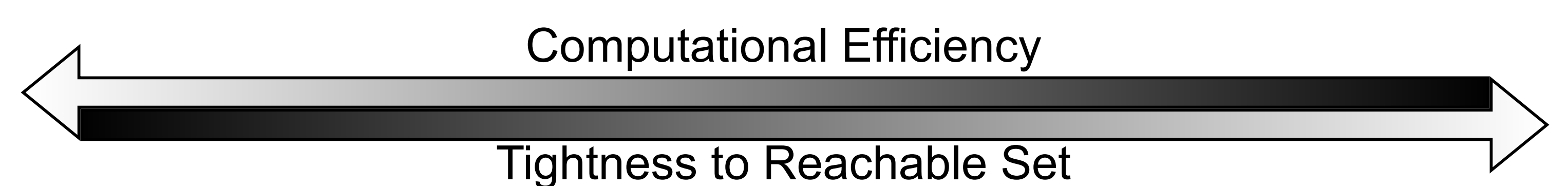


Compositional Reachability of the Closed-loop System

- open-loop system $\dot{x} = f(x, u, w)$ with inclusion function $F = \left[\frac{F}{F} \right]$ for f ,
- neural network controller $u = N(x)$ with inclusion function $\left[\frac{N}{N} \right]$
- the embedding system for the closed-loop system is

$$\begin{aligned} \dot{\underline{x}}_i &= \underline{F}_i(\underline{x}, \bar{x}_{i:\underline{x}}, \underline{\eta}_i, \bar{\eta}_i, \underline{w}, \bar{w}) := \underline{F}_i^S(\underline{x}, \bar{x}, \underline{w}, \bar{w}), \\ \dot{\bar{x}}_i &= \bar{F}_i(\underline{x}_{i:\underline{x}}, \bar{x}, \underline{\nu}_i, \bar{\nu}_i, \underline{w}, \bar{w}) := \bar{F}_i^S(\underline{x}, \bar{x}, \underline{w}, \bar{w}). \end{aligned}$$

Global (F^G)	Hybrid (F^H)	Local (F^L)
$\underline{\eta}_i = \underline{N}_{[\underline{x}, \bar{x}]}(\underline{x}, \bar{x})$ $\bar{\eta}_i = \bar{N}_{[\underline{x}, \bar{x}]}(\underline{x}, \bar{x})$	$\underline{\eta}_i = \underline{N}_{[\underline{x}, \bar{x}]}(\underline{x}, \bar{x}_{i:\underline{x}})$ $\bar{\eta}_i = \bar{N}_{[\underline{x}, \bar{x}]}(\underline{x}, \bar{x}_{i:\underline{x}})$	$\underline{\eta}_i = \underline{N}_{[\underline{x}, \bar{x}_{i:\underline{x}}]}(\underline{x}, \bar{x}_{i:\underline{x}})$ $\bar{\eta}_i = \bar{N}_{[\underline{x}, \bar{x}_{i:\underline{x}}]}(\underline{x}, \bar{x}_{i:\underline{x}})$
$\underline{\nu}_i = \underline{N}_{[\underline{x}, \bar{x}]}(\underline{x}, \bar{x})$ $\bar{\nu}_i = \bar{N}_{[\underline{x}, \bar{x}]}(\underline{x}, \bar{x})$	$\underline{\nu}_i = \underline{N}_{[\underline{x}, \bar{x}]}(\underline{x}_{i:\underline{x}}, \bar{x})$ $\bar{\nu}_i = \bar{N}_{[\underline{x}, \bar{x}]}(\underline{x}_{i:\underline{x}}, \bar{x})$	$\underline{\nu}_i = \underline{N}_{[\underline{x}_{i:\underline{x}}, \bar{x}]}(\underline{x}_{i:\underline{x}}, \bar{x})$ $\bar{\nu}_i = \bar{N}_{[\underline{x}_{i:\underline{x}}, \bar{x}]}(\underline{x}_{i:\underline{x}}, \bar{x})$

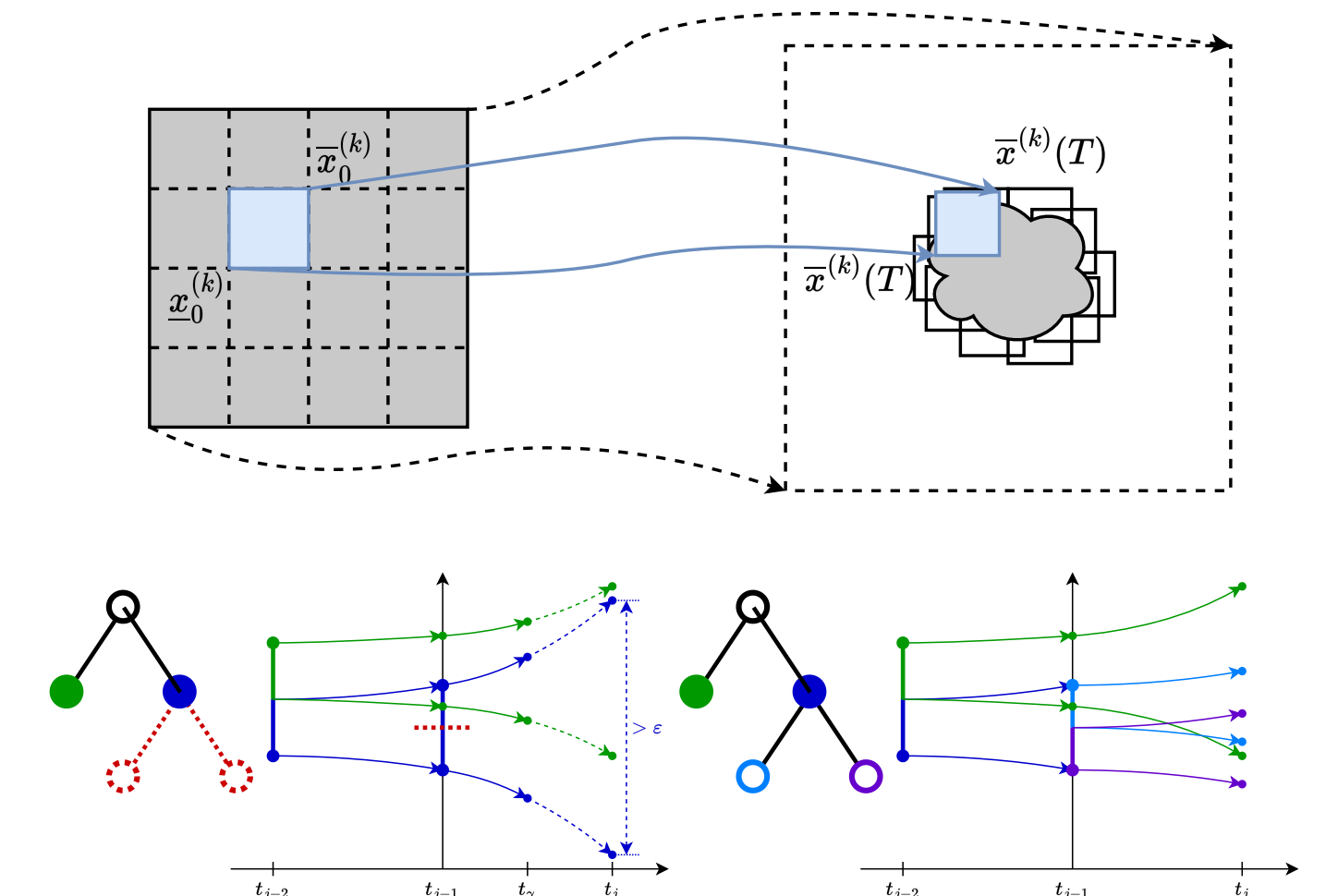


Theorem

For disturbance $\mathcal{W} = [\underline{w}, \bar{w}]$ and initial set $\mathcal{X}_0 = [\underline{x}_0, \bar{x}_0]$,
 $\mathcal{R}(t, \mathcal{X}_0, \mathcal{W}) \subseteq [\underline{x}^S(t), \bar{x}^S(t)]$,
 where $(\underline{x}^S(t), \bar{x}^S(t))$ is the trajectory of F^S for $S \in \{G, H, L\}$.

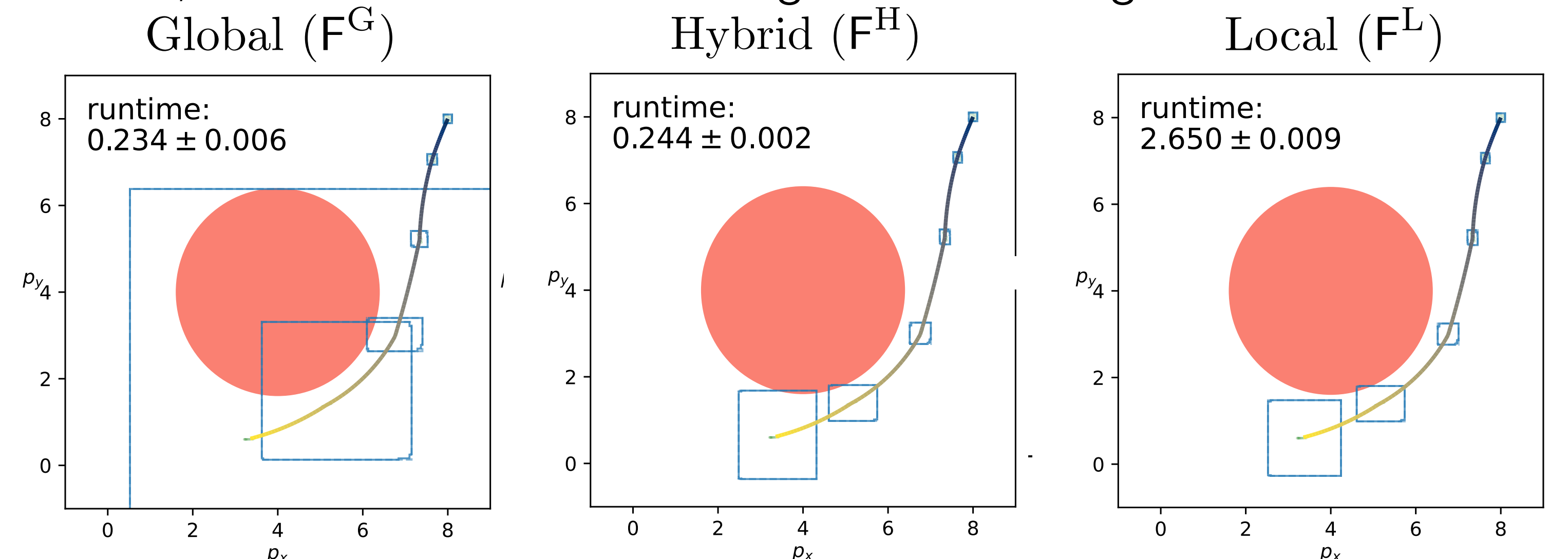
Numerical Experiments

- Partitioning improves the accuracy of interval analysis.
- separation between i) partitions that query neural network verification algorithm, and ii) partitions that only do integration.



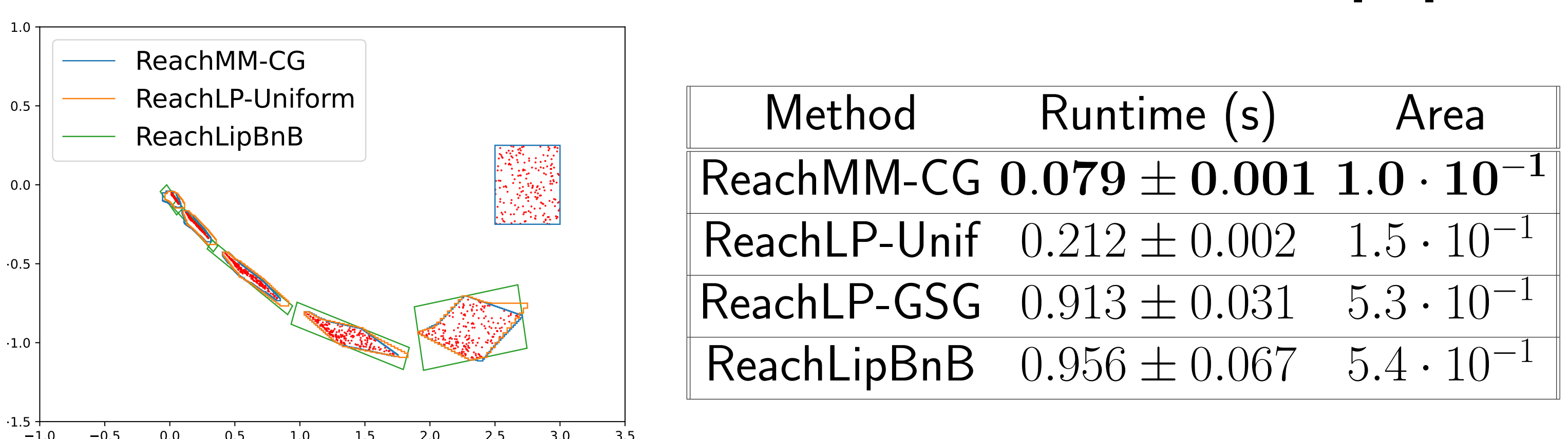
Vehicle Model:

Kinematic bicycle model, controlled by a $4 \times 100 \times 100 \times 2$ ReLU neural network, trained to stabilize to the origin while avoiding an obstacle.



Double Integrator Model:

Controlled by a $2 \times 10 \times 5 \times 1$ ReLU neural network, compare to [2,3].



References

- H. Zhang et al., *Efficient neural network robustness certification with general activation function*, NeurIPS, 2018
- M. Everett et al., *Reachability analysis of neural feedback loops*, IEEE Access, 2021
- T. Entesari et al., *ReachLipBnB: A branch-and-bound method for reachability analysis of neural autonomous systems using lipschitz bounds*, arXiv, 2022