

Comparative Analysis of Interval Reachability for Robust Implicit and Feedforward Neural Networks

Alexander Davydov^{1,*}, Saber Jafarpour^{2,*}, Matthew Abate², Francesco Bullo¹, Samuel Coogan²

Abstract—Implicit neural networks (INNs) are a class of learning models that use implicit algebraic equations as layers and have been shown to exhibit several notable benefits over traditional feedforward neural networks (FFNNs). In this paper, we use interval reachability analysis to study robustness of INNs and compare them with FFNNs. We first introduce the notion of tight inclusion function and use it to provide the tightest rectangular over-approximation of the neural network’s input-output map. We also show that tight inclusion functions lead to sharper robustness guarantees than the well-studied robustness measures of Lipschitz constants. Like exact Lipschitz constants, tight inclusions functions are computationally challenging to obtain, and thus we develop a framework based upon mixed monotonicity and contraction theory to estimate the tight inclusion functions for INNs. We show that our approach performs at least as well as, and generally better than, state-of-the-art interval-bound propagation methods for INNs. Finally, we design a novel optimization problem for training robust INNs and we provide empirical evidence that suitably-trained INNs can be more robust than comparably-trained FFNNs.

I. INTRODUCTION

Implicit neural networks (INNs) are a class of implicit learning models where the hidden layers are replaced with implicit equations [1], [2], [3]. Compared to their explicit counterparts, INNs are known to have advantages including (i) being more suitable for some problems such as constrained optimization problems [1] (ii) being more memory efficient while maintaining comparable accuracy [2] (iii) allowing for new architecture possibilities [3] (vi) showing improved training due to fewer vanishing and exploding gradients [4]. Despite their benefits, INNs can suffer from well-posedness issues and convergence instabilities. Additionally, their input-output behavior may suffer from robustness issues and adversarial perturbations; indeed, such robustness vulnerabilities are a well-studied and major issue in traditional deep neural networks as well [5].

Problem statement and motivations: Obtaining provable robustness guarantees for learning algorithms has been a major goal in the machine learning literature [6], [7]. For feedforward neural networks (FFNNs), four well-established methods for producing robustness guarantees include (i) Lipschitz bound approaches, (ii) interval-bound propagation

(IBP) methods, (iii) convex-relaxation approaches, and (iv) approaches based on Satisfiability Modulo Theories (SMT). Lipschitz constants of neural networks are coarse but rigorous measures for their input-output robustness. While it has been shown that computing the exact Lipschitz constant of a neural network is NP-hard [8], several efficient methods for providing sharp estimates of Lipschitz bounds are proposed in the literature [9]. IBP methods use interval analysis to provide box over-approximations of the reachable set of neural networks. These methods have been successfully used to perform formal verification [10] and to train robust FFNNs [11], [12]. Convex-relaxation approaches are based on relaxations of nonlinear activation functions using either linear [13] or quadratic constraints [14]. Finally, the SMT-based methods generalize the existing formal verification techniques for robustness analysis of neural networks [15].

For robustness guarantees of INNs, several works provide estimates of their Lipschitz constants [16], [17], [18]. However, global Lipschitz bounds do not provide information about the local sensitivity of the networks. The paper [13] proposes an iterative IBP approach for reachability analysis and training of INNs. However, convergence of this iteration requires strong conditions which limit the expressivity of the resulting implicit models. In [19], a method based on semidefinite programming is proposed for robustness analysis of INNs. Unfortunately, this approach is computationally intensive and cannot be implemented in training.

Most works on INNs focus on comparing their representation power [3] or their memory efficiency [2] with traditional deep neural networks. While recent empirical evidence indicate that appropriately-trained INNs can be significantly more robust than deep feedforward models [16], there are very few works in the literature that rigorously compare robustness of INNs with FFNNs. In this paper, we aim to address two main challenges regarding robustness of neural networks: (i) designing algorithms to train provably robust INNs, and (ii) appropriately comparing the robustness of INNs with the robustness of their explicit counterparts.

Contributions: In this paper, we use interval reachability analysis to study robustness of INNs. We first introduce the notion of tight inclusion function associated to the INN that gives the tightest rectangular over-approximation for the input-output behavior of the INN. We show that tight inclusion functions are sharper than any robustness guarantees based on local Lipschitz bounds. Similar to Lipschitz constants, computing the tight inclusion function is computationally challenging. Instead, using mixed monotone systems theory and contraction theory, we provide computa-

* These authors contributed equally and the order is alphabetical.

¹Alexander Davydov and Francesco Bullo are with the Center for Control, Dynamical Systems, and Computation, University of California, Santa Barbara, {davydov, bullo}@ucsb.edu

²Saber Jafarpour, Matthew Abate, and Samuel Coogan are with the School of Electrical and Computer Engineering, Georgia Institute of Technology, {saber, matt.abate, sam.coogan}@gatech.edu

This work is partly funded by the Air Force Office of Scientific Research under grant FA9550-22-1-0059 and National Science Foundation under awards #1749357 and #1931980.

tionally efficient estimates of the tight inclusion functions of INNs. We use two different interpretations of implicit neural networks to compare our approach with the IBP approach for FFNNs. We show that our mixed monotone contracting approach is the natural extension of IBP methods to INNs and performs at least as well as, and generally better than, IBP methods naively applied to INNs. Lastly, we provide an algorithm to efficiently implement our mixed monotone contracting approach in the training optimization problem to design robust INNs. In numerical experiments, we compare the performance of INNs and FFNNs with a comparable number of parameters and demonstrate that suitably-trained INNs have improved certified and empirical robustness compared to their feedforward counterparts, even when trained with IBP. In the conference paper [20], we focus on verifying robustness of INNs. In contrast, in this paper, we provide an efficient means for training robust INNs, compare theoretical robustness guarantees for INNs to FFNNs, extend [20, Theorem 1] and provide the proof for it, and present an empirical study of training for robustness.

II. MATHEMATICAL PRELIMINARY

For $x, y, z \in \mathbb{R}^n$, we write $x \leq y$ if $x_i \leq y_i$ for all $i \in \{1, \dots, n\}$ and $z \in [x, y]$ if $x \leq z \leq y$. For $\eta \in \mathbb{R}_{>0}^n$, we define the diagonal matrix $[\eta] \in \mathbb{R}^{n \times n}$ by $[\eta]_{ii} = \eta_i$, for every $i \in \{1, \dots, n\}$ and the diagonally weighted ℓ_∞ -norm by $\|x\|_{\infty, [\eta]} = \max_i |x_i|/\eta_i$, the diagonally weighted ℓ_∞ -matrix measure is defined by $\mu_{\infty, [\eta]}(A) = \max_{i \in \{1, \dots, n\}} A_{ii} + \sum_{j \neq i} \frac{\eta_j}{\eta_i} |A_{ij}|$. For any matrix $A \in \mathbb{R}^{n \times n}$, the spectral radius of A is denoted by $\rho(A)$ and the elementwise absolute value of A is denoted by $|A| \in \mathbb{R}_{\geq 0}^{n \times n}$. For two matrices A, B , let $A \otimes B$ denote their Kronecker product. Given a matrix $B \in \mathbb{R}^{n \times m}$, we denote the non-negative part of B by $[B]^+ := \max(B, 0)$ and the nonpositive part of B by $[B]^- := \min(B, 0)$. The Metzler and non-Metzler parts of a square matrix $A \in \mathbb{R}^{n \times n}$ are denoted by $[A]^{\text{Mzl}}$ and $[A]^{\text{Mzl}}$, respectively, where

$$([A]^{\text{Mzl}})_{ij} := \begin{cases} A_{ij} & \text{if } A_{ij} \geq 0 \text{ or } i = j \\ 0 & \text{otherwise,} \end{cases}$$

and $[A]^{\text{Mzl}} := A - [A]^{\text{Mzl}}$. The subset $\mathcal{T}^n \subset \mathbb{R}^{2n}$ is defined by $\mathcal{T}^n := \{(x, \hat{x}) \in \mathbb{R}^{2n} \mid x \leq \hat{x}\}$. Given a map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, a set $\mathcal{U} \subset \mathbb{R}^n$, and $p \in [1, \infty]$, the ℓ_p -Lipschitz constant of f on \mathcal{U} is the smallest $\text{Lip}_p^{\mathcal{U}}(f) \in \mathbb{R}_{\geq 0}$ such that $\|f(x) - f(y)\|_p \leq \text{Lip}_p^{\mathcal{U}}(f) \|x - y\|_p$, for all $x, y \in \mathcal{U}$. Given a map $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\alpha \in [0, 1]$, the α -average map is defined by $f_\alpha(x) = (1 - \alpha)x + \alpha f(x)$, for every $x \in \mathbb{R}^n$.

III. INCLUSION FUNCTIONS

Given a mapping f , an ℓ_∞ -box over-approximation of the image of f is attainable via an inclusion function.

Definition 3.1 (Inclusion function): Let $f : \mathbb{R}^r \rightarrow \mathbb{R}^q$ be a mapping. Then $F = \begin{bmatrix} \underline{F} \\ \overline{F} \end{bmatrix} : \mathcal{T}^r \rightarrow \mathbb{R}^{2q}$ is an inclusion function for f , if, for every $x \leq \hat{x}$,

- (i) $\underline{F}(y, y) \geq \underline{F}(x, \hat{x})$ and $\overline{F}(y, y) \leq \overline{F}(x, \hat{x})$;
- (ii) $\underline{F}(x, x) = \overline{F}(x, x) = f(x)$.

Moreover, the inclusion function F for f is called *tight*, if

- (iii) for every inclusion function $G = \begin{bmatrix} \underline{G} \\ \overline{G} \end{bmatrix} : \mathcal{T}^r \rightarrow \mathbb{R}^{2q}$ of f , we have

$\underline{G}(x, \hat{x}) \leq \underline{F}(x, \hat{x})$, $\overline{F}(x, \hat{x}) \geq \overline{G}(x, \hat{x})$, for all $x \leq \hat{x}$
If F is an inclusion function for f , then it is easy to see that,

$$f([x, \hat{x}]) \subseteq [\underline{F}(x, \hat{x}), \overline{F}(x, \hat{x})], \quad \text{for all } x \leq \hat{x}. \quad (1)$$

Given a map $f : \mathbb{R}^r \rightarrow \mathbb{R}^q$, one can use [21, Theorem 1] to compute the tight inclusion function, component-wise. Indeed, for every $i \in \{1, \dots, n\}$, one can show that:

$$\underline{F}_i(x, \hat{x}) = \min_{z \in [x, \hat{x}]} f_i(z), \quad \overline{F}_i(x, \hat{x}) = \max_{z \in [x, \hat{x}]} f_i(z) \quad (2)$$

The next Theorem studies the connection between the local Lipschitz constants and the tight inclusion functions.

Theorem 3.2 (Inclusion function vs. Lipschitz constant):

Let $f : \mathbb{R}^r \rightarrow \mathbb{R}^q$ be a continuous mapping and $F = \begin{bmatrix} \underline{F} \\ \overline{F} \end{bmatrix} : \mathcal{T}^r \rightarrow \mathbb{R}^{2q}$ be the tight inclusion function for f . Then, for every $\underline{x} \leq \overline{x}$, we have

$$\|\overline{F}(\underline{x}, \overline{x}) - \underline{F}(\underline{x}, \overline{x})\|_\infty \leq \text{Lip}_{\infty}^{[\underline{x}, \overline{x}]}(f) \|\underline{x} - \overline{x}\|_\infty.$$

Proof: Let $i \in \{1, \dots, k\}$ be such that $\|\overline{F}(\underline{x}, \overline{x}) - \underline{F}(\underline{x}, \overline{x})\|_\infty = |\overline{F}_i(\underline{x}, \overline{x}) - \underline{F}_i(\underline{x}, \overline{x})|$. Note that since f is continuous and the box $[\underline{x}, \overline{x}]$ is compact, there exist $\eta^*, \xi^* \in [\underline{x}, \overline{x}]$ such that $\max_{y \in [\underline{x}, \overline{x}]} f_i(y) = f_i(\eta^*)$, $\min_{y \in [\underline{x}, \overline{x}]} f_i(y) = f_i(\xi^*)$. This implies that $\|\overline{F}(\underline{x}, \overline{x}) - \underline{F}(\underline{x}, \overline{x})\|_\infty = |f_i(\eta^*) - f_i(\xi^*)| \leq \|f(\xi^*) - f(\eta^*)\|_\infty \leq \text{Lip}_{\infty}^{[\underline{x}, \overline{x}]} \|\xi^* - \eta^*\|_\infty \leq \text{Lip}_{\infty}^{[\underline{x}, \overline{x}]}(f) \|\underline{x} - \overline{x}\|_\infty. \quad \blacksquare$

Remark 3.3 (Tight inclusion functions): Theorem 3.2 shows that a tight inclusion function for $y = f(x)$ will provide a tighter over-approximation of the image of f than is attainable from local Lipschitz constants of f . In general, finding tight inclusion functions using (2) is not computationally tractable. This motivates developing efficient methods for estimating the tight inclusion function.

IV. IMPLICIT NEURAL NETWORKS

We consider the implicit neural network

$$\begin{aligned} z &= \sigma(Wz + Ux + b) := \text{N}(z, x), \\ y &= Cz + c, \end{aligned} \quad (3)$$

where $z \in \mathbb{R}^n$ is the hidden variable, $x \in \mathbb{R}^r$ is the input, $y \in \mathbb{R}^q$ is the output, $W \in \mathbb{R}^{n \times n}$, $U \in \mathbb{R}^{n \times r}$, and $C \in \mathbb{R}^{q \times n}$ are the weight matrices and $b \in \mathbb{R}^n$ and $c \in \mathbb{R}^q$ are bias vectors. Moreover, σ is a diagonal activation function (e.g., ReLU) defined by $\sigma(z_1, \dots, z_n) = (\sigma_1(z_1), \dots, \sigma_n(z_n))^T$, where for every $i \in \{1, \dots, n\}$, the activation function $\sigma_i : \mathbb{R} \rightarrow \mathbb{R}$ satisfies $0 \leq \frac{\sigma_i(x) - \sigma_i(y)}{x - y} \leq 1$ for all $x \neq y \in \mathbb{R}$. Compared to traditional neural networks, INNs replace the layers with a fixed-point equation. This change in the structure is known to allow for new architecture possibilities and provide alternative approaches to deep modeling.

Generalized architecture: Notably, FFNNs can be considered as special cases of INNs [3]. Consider the FFNN

$$\begin{aligned} z^i &= \sigma(W_i z^{i-1} + b_i) =: \text{FN}_i(z^{i-1}), \quad i \in \{1, \dots, k\}, \\ y &= Cz^k + c \end{aligned} \quad (4)$$

where $z_0 = x \in \mathbb{R}^r$ is the input. For every $i \in \{1, \dots, k\}$, $z^i \in \mathbb{R}^{n_i}$, $W_i \in \mathbb{R}^{n_{i-1} \times n_i}$, and $b_i \in \mathbb{R}^{n_{i-1}}$ are the weights, the biases, and the hidden variables in the i -th layer of the network, respectively. Finally, $y \in \mathbb{R}^q$ is the output, C and c are the output layer's weight matrix and bias vector. The FFNN (4) is equivalent to the following INN

$$\begin{aligned} z &= \sigma(W^{\text{FN}}z + U^{\text{FN}}x + b) =: \text{IFN}(z, x), \\ y &= C^{\text{N}}z + c \end{aligned} \quad (5)$$

where $z = [z_k, \dots, z_1]^T$, $b = [b_k, \dots, b_1]^T$, and $W^{\text{FN}}, U^{\text{FN}}$, and C^{FN} are defined as follows:

$$\begin{aligned} W^{\text{FN}} &= \begin{bmatrix} 0 & W_k & 0 & \cdots & 0 \\ 0 & 0 & W_{k-1} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & W_1 \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}, \quad U^{\text{FN}} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ W_0 \end{bmatrix} \\ C^{\text{FN}} &= [C \quad 0 \quad 0 \quad \cdots \quad 0]. \end{aligned}$$

Using this perspective, implicit neural networks generalize FFNNs by allowing arbitrary interconnections between layers leading to full weight matrices W , U , and C .

Alternative deep modeling: By replacing the notion of layer with an algebraic equation, implicit neural networks provide a novel perspective toward deep modeling. Consider the class of FFNN where weights and biases are equal for each layer (i.e., the network is weight-tied) and the input is injected to each layer given by

$$\begin{aligned} z^i &= \sigma(Wz^{i-1} + Ux + b), \quad i \in \{1, \dots, k\}, \\ y &= Cz_k + c \end{aligned} \quad (6)$$

where $z^0 = x$. While weight-tying may appear restrictive, it is usually considered as a form of regularization that stabilizes training and significantly reduces the model size [2]. If the depth of the network increases, i.e., $k \rightarrow \infty$ and the iteration (6) converges, then the weight-tied input-injected neural network (6) is equivalent to the implicit neural network (3). Using this perspective, INNs provide a depth-independent alternative to deep FFNNs.

Suppose that, for every input $x \in \mathbb{R}^r$, the implicit neural network (3) has a unique fixed point $z^*(x) \in \mathbb{R}^n$. Then, the *input-output map* $f: \mathbb{R}^r \rightarrow \mathbb{R}^q$ is given by

$$f(x) := y = Cz^*(x) + c. \quad (7)$$

In the next section, our goal is to provide estimates for the tight inclusion function of the input-output map f .

V. REACHABILITY ANALYSIS OF IMPLICIT NEURAL NETWORKS

In this section, we use mixed monotone system theory to present a framework for estimating the tight inclusion function for the input-output map of INNs. Given an implicit neural network (3) and input bounds $\underline{x} \leq \bar{x} \in \mathbb{R}^r$, we first introduce the embedding map $\mathbf{N}^{\text{E}}: \mathbb{R}^{2n} \times \mathbb{R}^{2r} \rightarrow \mathbb{R}^n$ by

$$\begin{aligned} \mathbf{N}^{\text{E}}(z, \bar{z}, \underline{x}, \bar{x}) &= \sigma([W]^{\text{Mzl}}z + [W]^{\text{Mzl}}\bar{z} \\ &\quad + [U]^+ \underline{x} + [U]^- \bar{x} + b). \end{aligned}$$

Using the embedding map \mathbf{N}^{E} , we define the *embedded implicit neural network* associated to (3) by

$$\begin{bmatrix} z \\ \bar{z} \end{bmatrix} = \begin{bmatrix} \mathbf{N}^{\text{E}}(z, \bar{z}, \underline{x}, \bar{x}) \\ \mathbf{N}^{\text{E}}(\bar{z}, z, \bar{x}, \underline{x}) \end{bmatrix}, \quad \begin{bmatrix} y \\ \bar{y} \end{bmatrix} = \begin{bmatrix} [C]^+ & [C]^- \\ [C]^- & [C]^+ \end{bmatrix} \begin{bmatrix} z \\ \bar{z} \end{bmatrix} + \begin{bmatrix} c \\ c \end{bmatrix}. \quad (8)$$

The embedded INN (8) can be considered as an INN with the box input $[\underline{x}, \bar{x}]$ and the box output $[\underline{y}, \bar{y}]$. In the next theorem, we use the embedded system (8) to obtain an inclusion function for the input-output map (7).

Theorem 5.1 (Inclusion function via embedded network): Consider the implicit neural network (3) and its associated embedded implicit neural network (8). Suppose that there exists $\eta \in \mathbb{R}_{>0}^n$ such that $\mu_{\infty, [\eta]^{-1}}(W) < 1$. For every $\underline{x} \leq x \leq \bar{x}$, and every $\alpha \in (0, \alpha^* := [1 - \min_{i \in \{1, \dots, n\}} (W_{ii})_-]^{-1}]$, the following statements hold:

- (i) the iterations $\begin{bmatrix} z^{k+1} \\ \bar{z}^{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{N}^{\text{E}}_\alpha(z^k, \bar{z}^k, \underline{x}, \bar{x}) \\ \mathbf{N}^{\text{E}}_\alpha(\bar{z}^k, z^k, \bar{x}, \underline{x}) \end{bmatrix}$ are contracting with respect to the norm $\|\cdot\|_{\infty, I_2 \otimes [\eta]^{-1}}$ and converge to the unique fixed-point $\begin{bmatrix} z^* \\ \bar{z}^* \end{bmatrix}$ of the embedded INN (8);
- (ii) the iterations $z^{k+1} = \mathbf{N}_\alpha(z^k, x)$ are contracting with respect to the norm $\|\cdot\|_{\infty, [\eta]^{-1}}$ and converges to the unique fixed point $z^* \in [z^*, \bar{z}^*]$ of the INN (3);
- (iii) the map $\mathbf{F}^{\text{N}}: \mathcal{T}^r \rightarrow \mathbb{R}^{2q}$ defined by

$$\begin{aligned} \underline{\mathbf{F}}^{\text{N}}(\underline{x}, \bar{x}) &= [C]^+ z^* + [C]^- \bar{z}^* + c \\ \bar{\mathbf{F}}^{\text{N}}(\underline{x}, \bar{x}) &= [C]^+ \bar{z}^* + [C]^- z^* + c \end{aligned} \quad (9)$$

is an inclusion function for f defined in (7).

Proof: Regarding part (i), we define $\tilde{\sigma} = I_2 \otimes \sigma$, the map $\mathbf{G}: \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$ by $\mathbf{G}(z, \bar{z}) = \begin{bmatrix} [W]^{\text{Mzl}}z + [W]^{\text{Mzl}}\bar{z} \\ [W]^{\text{Mzl}}\bar{z} + [W]^{\text{Mzl}}z \end{bmatrix}$ and the matrices $D = \begin{bmatrix} [U]^+ & [U]^- \\ [U]^- & [U]^+ \end{bmatrix}$ and $w = [\underline{x}, \bar{x}]^T$. Then define $\tilde{\sigma}^{\mathbf{G}}: \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$ as follows

$$\tilde{\sigma}^{\mathbf{G}}(z, \bar{z}, w) := \begin{bmatrix} \mathbf{N}^{\text{E}}(z, \bar{z}, \underline{x}, \bar{x}) \\ \mathbf{N}^{\text{E}}(\bar{z}, z, \bar{x}, \underline{x}) \end{bmatrix} = \tilde{\sigma}(\mathbf{G}(z, \bar{z}) + Dw + I_2 \otimes b).$$

The assumptions on each scalar activation function imply that (i) $\tilde{\sigma}: \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$ is non-expansive with respect to $\|\cdot\| := \|\cdot\|_{\infty, I_2 \otimes [\eta]^{-1}}$ and (ii) for every $p, q \in \mathbb{R}$, there exists $\theta_i \in [0, 1]$ such that $\sigma_i(p) - \sigma_i(q) = \theta_i(p - q)$ or in the matrix form $\tilde{\sigma}(p) - \tilde{\sigma}(q) = \Theta(p - q)$ where $\Theta \in \mathbb{R}^{2n \times 2n}$

is a diagonal matrix with diagonal elements $\theta_i \in [0, 1]$ and $\mathbf{p}, \mathbf{q} \in \mathbb{R}^{2n}$. As a result, for every $y_1, y_2 \in \mathbb{R}^{2n}$, we have

$$\begin{aligned} & \|\tilde{\sigma}_\alpha^G(y_1, w) - \tilde{\sigma}_\alpha^G(y_2, w)\| \\ &= \|(1 - \alpha)(y_1 - y_2) + \alpha\Theta(G(y_1) - G(y_2))\| \\ &\leq \sup_{y \in \mathbb{R}^{2n}} \|(1 - \alpha)I_{2n} + \alpha\Theta DG(y)\| \|y_1 - y_2\| \end{aligned}$$

where the inequality holds by the mean value theorem. Then, for every $\alpha \in (0, [1 - \min_i \inf_{y \in \mathbb{R}^{2n}} (\Theta DG(y))_{ii}]^{-1}]$,

$$\begin{aligned} & \|I_{2n} + \alpha(-I_{2n} + \Theta DG(y))\| \\ &= 1 + \alpha\mu_{\infty, I_2 \otimes [\eta]^{-1}}(-I_{2n} + \Theta DG(y)) \\ &= 1 + \alpha(-1 + \mu_{\infty, I_2 \otimes [\eta]^{-1}}(\Theta DG(y))) \\ &\leq 1 + \alpha(-1 + \mu_{\infty, I_2 \otimes [\eta]^{-1}}(DG(y))^+) \\ &\leq 1 - \alpha(1 - \mu_{\infty, [\eta]^{-1}}(W)^+) < 1, \end{aligned}$$

where the first equality holds by [17, Lemma 7(i)], the second equality holds by translation property of matrix measures, the third inequality holds by [17, Lemma 8(i)], and the fourth inequality holds by the definition of matrix measure. Moreover, since $\theta_i \in [0, 1]$, we have $\theta_i(DG)_{ii} \geq (DG)_{ii}^-$, for every $i \in \{1, \dots, 2n\}$. This means that

$$\inf_{y \in \mathbb{R}^{2n}} (\Theta DG(y))_{ii} \geq \inf_{y \in \mathbb{R}^{2n}} (DG_{ii}(y))^- = (W_{ii})^-.$$

This implies that, for every $\alpha \in (0, \alpha^*]$,

$$\begin{aligned} & \|\tilde{\sigma}_\alpha^G(x_1, u) - \tilde{\sigma}_\alpha^G(x_2, u)\| \\ &\leq (1 - \alpha(1 - \mu_{\infty, [\eta]^{-1}}(W)^+)) \|x_1 - x_2\| \end{aligned}$$

Since $1 - \alpha(1 - \mu_{\infty, [\eta]^{-1}}(W)^+) < 1$, $\tilde{\sigma}_\alpha^G(\cdot, w)$ is a contraction mapping with respect to $\|\cdot\|_{\infty, I_2 \otimes [\eta]^{-1}}$ for every $\alpha \in (0, \alpha^*]$. It is easy to see that $\tilde{\sigma}_\alpha^G$ and $\tilde{\sigma}^G$ have the identical fixed-points, for every $\alpha \in [0, 1]$. Therefore the iterations in part (i) converge to the unique fixed point of the embedded INN (8). Regarding part (ii), the proof follows by applying the same argument as in the proof of part (i) and using $\sigma(Wz + Ux + b)$ instead of $\tilde{\sigma}^G(z, \bar{z}, \underline{x}, \bar{x})$. Now, we show that $\underline{z}^* \leq z^* \leq \bar{z}^*$. We choose the initial condition $\begin{bmatrix} \underline{z}^0 \\ \bar{z}^0 \end{bmatrix}$ for the iterations in part (i) and choose an initial condition $z^0 \in \mathbb{R}^n$ satisfying $\underline{z}^0 \leq z^0 \leq \bar{z}^0$ for the iterations in part (ii). We prove by induction that, for every $k \in \mathbb{Z}_{\geq 0}$, we have $\underline{z}^k \leq z^k \leq \bar{z}^k$. Suppose that this claim is true for $k \in \{1, \dots, m\}$ and we show that this claim is true for $k = m + 1$. We first define $p = \lceil W \rceil^{\text{Mzl}} \underline{z}^m + \lfloor W \rfloor^{\text{Mzl}} \bar{z}^m + [U]^+ \underline{x} + [U]^- \bar{x} + b$ and $q = Wz^m + Ux + b$. Then we have

$$\begin{aligned} \underline{z}^{m+1} - z^{m+1} &= (1 - \alpha^*)(\underline{z}^m - z^m) + \alpha^*(\sigma(p) - \sigma(q)) \\ &= ((1 - \alpha^*)I_n + \alpha^*\Theta \lceil W \rceil^{\text{Mzl}})(\underline{z}^m - z^m) \\ &\quad + \alpha^*\Theta \lfloor W \rfloor^{\text{Mzl}}(\bar{z}^m - z^m) \\ &\quad + \alpha^*\Theta [U]^+(\underline{x} - x) + \alpha^*\Theta [U]^-(\bar{x} - x), \end{aligned}$$

where the non-negative diagonal matrix $\Theta = \text{diag}(\theta_i) \in \mathbb{R}^n$ is defined as follows: for every $i \in \{1, \dots, n\}$, $\theta_i \in [0, 1]$ is such that $\sigma_i(p_i) - \sigma_i(q_i) = \theta_i(p_i - q_i)$. Moreover, we know that $\Theta \lceil W \rceil^{\text{Mzl}} \leq \mathbb{0}_{n \times n}$ and, for every $i \in \{1, \dots, n\}$,

$$(1 - \alpha^*) + \alpha^*\theta_i W_{ii} \geq (1 - \alpha^*) + \alpha^*W_{ii}^- \geq 0.$$

This implies that $(1 - \alpha^*)I_n + \alpha^*\Theta \lceil W \rceil^{\text{Mzl}} \geq \mathbb{0}_{n \times n}$. Additionally, we have $\Theta [U]^+ \geq \mathbb{0}_{n \times r}$ and $\Theta [U]^- \leq \mathbb{0}_{n \times r}$. Therefore, using the induction assumption, we get $\underline{z}^{m+1} - z^{m+1} \leq \mathbb{0}_n$. Similarly, one can show that $z^{m+1} - \bar{z}^{m+1} \leq \mathbb{0}_n$. As a consequence, $\underline{z}^* = \lim_{k \rightarrow \infty} \underline{z}^k \leq \lim_{k \rightarrow \infty} z^k = z^* \leq \lim_{k \rightarrow \infty} \bar{z}^k = \bar{z}^*$. This proves part (ii). The proof of part (iii) follows easily from parts (i) and (ii) and by checking the properties of inclusion functions from Definition 3.1. ■

Remark 5.2 (Mixed monotone contracting approach):

Theorem 5.1 can be interpreted as a dynamical system approach to study robustness of INNs. Indeed, the α -average iteration in part (ii) is the forward Euler discretization of the dynamical system $\frac{dz}{dt} = -z + N(z, x)$. The convergence of the iterations is due to the contraction property of the dynamical system and the estimate for the inclusion function is due to the mixed monotonicity of the dynamical system associated with the embedded INN [20].

VI. FEEDFORWARD VS. IMPLICIT NEURAL NETWORKS

In this section, we compare the robust training framework developed in Sections (V) and (VII) with the IBP approach developed in [11]. Consider a k -layer FFNN with input-output map $f(x) =: y$ defined by (4). For every $i \in \{1, \dots, k\}$, following [11], [12], we can obtain layer-wise upper and lower bounds for the hidden variables as follows:

$$\begin{aligned} \underline{z}^{i+1} &:= \text{FN}_i^E(\underline{z}^i, \bar{z}^i) = \sigma([W_i]^+ \underline{z}^i + [W_i]^- \bar{z}^i + b_i), \\ \bar{z}^{i+1} &:= \text{FN}_i^E(\bar{z}^i, \underline{z}^i) = \sigma([W_i]^+ \bar{z}^i + [W_i]^- \underline{z}^i + b_i). \end{aligned}$$

By applying this bounding technique recursively, one can obtain the upper bound \bar{y} and lower bound \underline{y} for the output of the FFNN. The IBP inclusion function $\text{F}^{\text{FN}} = \begin{bmatrix} \text{F}^{\text{FN}} \\ \text{F}^{\text{FN}} \end{bmatrix} : \mathcal{T}^r \rightarrow \mathbb{R}^{2q}$ for the input-output map f is then defined by:

$$\underline{\text{F}}^{\text{FN}}(\underline{x}, \bar{x}) = \underline{y}, \quad \bar{\text{F}}^{\text{FN}}(\underline{x}, \bar{x}) = \bar{y}, \quad (10)$$

In the next two subsections, we use the two perspective *generalized architecture* and *alternative deep modeling* toward INNs, to establish connections between the IBP approach in [11] and our mixed monotone contracting approach.

A. Generalized architecture

By considering finite-depth FFNNs as a special case of implicit neural networks, one can show that our mixed monotone contracting approach is a generalization of the IBP approach in [11] to INNs.

Theorem 6.1 (Embedded feedforward neural networks):

Consider the FFNN (4) with the associated implicit neural network (5). The following statements hold:

- (i) for every $i \in \{1, \dots, k\}$, the function $(z, \bar{z}) \mapsto \begin{bmatrix} \text{FN}_i^E(z, \bar{z}) \\ \text{FN}_i^E(\bar{z}, z) \end{bmatrix}$ is a tight inclusion function for the i th layer evaluation map $\text{FN}_i(z) := \sigma(W_i z + b_i)$;

- (ii) there exists $\eta \in \mathbb{R}_{>0}^n$ such that $\mu_{\infty, [\eta]^{-1}}(W^{\text{FN}}) < 1$.

If F^{FN} is the inclusion function obtained from (5) using mixed monotone contracting approach and F^{FN} is the inclusion function (10) obtained using IBP approach, then

- (iii) $\underline{\text{F}}^{\text{FN}}(\underline{x}, \bar{x}) = \underline{\text{F}}^{\text{FN}}(\underline{x}, \bar{x})$, for every $\underline{x} \leq \bar{x} \in \mathbb{R}^r$.

The proof of Theorem 6.1(ii) and (iii) is based on choosing $[\eta]^{-1} = \text{diag}(\delta^{-1}, \dots, \delta^{-n})$ for $\delta > 0$ sufficiently small and then applying Theorem 5.1.

B. Alternative deep modeling

By separating the notion of depth from the layer-wise evaluation, our mixed monotone contracting approach can be used to estimate the reachable sets of deep weight-tied FFNNs. For the weight-tied FFNN (6), we define

$$\text{WFN}_i^E(z, \bar{z}, \underline{x}, \bar{x}) = \sigma([W]^+ \underline{z} + [W]^- \bar{z} + [U]^+ \underline{x} + [U]^- \bar{x} + b).$$

By replacing FN_i^E with WFN_i^E in (10), we the IBP inclusion function $\mathbb{F}^{\text{WFN}} : \mathcal{T}^r \rightarrow \mathbb{R}^{2q}$.

Theorem 6.2 (Weight-tied infinite-layer neural networks): Suppose that $\rho(|W|) < 1$ and let $\eta \in \mathbb{R}_{>0}^n$ be the right Perron eigenvector of $|W|$. Then for the weight-tied FFNN (6),

- (i) $\lim_{i \rightarrow \infty} \begin{bmatrix} z^i \\ \bar{z}^i \end{bmatrix} = \begin{bmatrix} \underline{w}^* \\ \bar{w}^* \end{bmatrix}$ for some $\underline{w}^* \leq \bar{w}^* \in \mathbb{R}^n$;
- (ii) $\lim_{i \rightarrow \infty} z^i = w^*$ for some $w^* \in [\underline{w}^*, \bar{w}^*]$;

Moreover, for the implicit neural network (3),

- (iii) $\mu_{\infty, [\eta]^{-1}}(W) < 1$;

If \mathbb{F}^{WFN} is the inclusion function by the IBP approach as $k \rightarrow \infty$ and \mathbb{F}^N is the inclusion function from (3), then

- (iv) for every $\underline{x} \leq \bar{x} \in \mathbb{R}^r$, $\mathbb{F}^N(\underline{x}, \bar{x}) \geq \mathbb{F}^{\text{WFN}}(\underline{x}, \bar{x})$, and $\bar{\mathbb{F}}^N(\underline{x}, \bar{x}) \leq \bar{\mathbb{F}}^{\text{WFN}}(\underline{x}, \bar{x})$

The proof of Theorem 6.2(iii) is based on the inequality $\mu_{\infty, [\eta]^{-1}}(W) \leq \rho(|W|) < 1$ and the proof of Theorem 6.2(iv) follows from $[W]^+ \leq [W]^{\text{Mzl}}$.

VII. TRAINING ROBUST IMPLICIT NEURAL NETWORKS

A. Certified adversarial robustness for classification tasks

We say an INN is certifiably robust for input x if its prediction at x is verifiably constant within a given ℓ_∞ ball around x . We refer to [20] for a rigorous definition of certified adversarial robustness. We use the embedded INN (8) to obtain a sufficient condition for certified robustness. Given a robustness radius $\epsilon > 0$, for every input $x \in \mathbb{R}^r$, we define $\underline{x} = x - \epsilon \mathbb{1}_r, \bar{x} = x + \epsilon \mathbb{1}_r$. Following [12, Eq. 3] and [20], for each input $x' \in [\underline{x}, \bar{x}]$, we define the *relative classifier variable*, $m^x(x') \in \mathbb{R}^q$ by

$$m^x(x') = f(x')_i \mathbb{1}_q - f(x'), \quad (11)$$

where i is the correct label of x . Note that $m^x(x')_j > 0$ for all $j \neq i$ if and only if x' is labeled the same as x by the neural network. Therefore, we write $m^x(x') = T^x f(x') = T^x C z^*(x') + T^x c$, for suitable specification matrix $T^x \in \{-1, 0, 1\}^{q \times q}$ defined via the linear transformation (11). Moreover, if there exists $\eta \in \mathbb{R}_{>0}^n$ so that $\mu_{\infty, [\eta]^{-1}}(W) < 1$, then we can use Theorem 5.1 to define

$$\underline{m}^x(\underline{x}, \bar{x}) = [T^x C]^+ z^*(\underline{x}, \bar{x}) + [T^x C]^- \bar{z}^*(\underline{x}, \bar{x}) + T^x c.$$

Moreover, $\min_{j \neq i} (\underline{m}^x(\underline{x}, \bar{x}))_j > 0$ is a sufficient condition for certified adversarial robustness of the INN [20].

B. Training optimization problem

We aim to design optimization problems to train a neural network which is robust to input perturbations with ℓ_∞ -norm smaller than some ϵ . Let \mathcal{L} be the cross-entropy loss function and assume that $\{(\hat{x}^l, \hat{y}^l)\}_{l=1}^N$ is a set of N labeled data points used for training. For every $l \in \{1, \dots, N\}$, we define the following upper and the lower bounds on the input \hat{x}^l by $\underline{x}^l = \hat{x}^l - \epsilon \mathbb{1}_r$ and $\bar{x}^l = \hat{x}^l + \epsilon \mathbb{1}_r$. We use the robust optimization framework [7] for designing robust neural networks. Our objective is to minimize the robust loss function $\sum_{l=1}^N \max_{x \in [\underline{x}^l, \bar{x}^l]} \mathcal{L}(f(\hat{x}^l), \hat{y}^l)$ on the training data. Using [6, Theorem 2], for the cross-entropy loss, and for $\underline{m}^l := \underline{m}^{\hat{x}^l}(\underline{x}^l, \bar{x}^l)$ and every $l \in \{1, \dots, N\}$,

$$\mathcal{L}(f(\hat{x}^l), \hat{y}^l) \leq \mathcal{L}(-\underline{m}^l, \hat{y}^l), \quad \text{for all } x \in [\underline{x}^l, \bar{x}^l].$$

As pointed out in [11] for FFNNs, using the loss function $\mathcal{L}(-\underline{m}^l, \hat{y}^l)$ in the training can lead to convergence instability and difficulty in training. To improve the stability of the training, following [11], we instead use a convex combination of the empirical risk loss and the robust loss. Therefore, for $T^l := T^{\hat{x}^l}$ we get the following training problem:

$$\begin{aligned} \min_{W, U, C, b, c, \eta} \quad & \sum_{l=1}^N (1 - \kappa) \mathcal{L}(y^l, \hat{y}^l) + \kappa \mathcal{L}(-\underline{m}^l, \hat{y}^l), \\ \begin{bmatrix} z^l \\ \bar{z}^l \end{bmatrix} = \quad & \begin{bmatrix} \mathbb{N}^E(z^l, \bar{z}^l, \underline{x}^l, \bar{x}^l) \\ \mathbb{N}^E(\bar{z}^l, z^l, \bar{x}^l, \underline{x}^l) \end{bmatrix}, \\ m^l = [T^l C]^+ z^l + [T^l C]^- \bar{z}^l + T^l c, \quad & z^l = \mathbb{N}(z^l, \hat{x}^l), \\ y^l = C z^l + c, \quad & \mu_{\infty, [\eta]^{-1}}(W) \leq \gamma. \end{aligned} \quad (12)$$

where $\kappa \in [0, 1]$ and $\gamma \in (-\infty, 1)$ are hyperparameters.

VIII. NUMERICAL EXPERIMENTS

In this section we provide an experimental comparison between the robustness of FFNNs and INNs trained with and without IBP and mixed monotonicity, respectively¹.

Experimental setup: We consider the MNIST dataset, which contains 70000 28×28 pixel images of handwritten digits. For training, pixels are normalized into the range $[0, 1]$. All INNs have $n = 100$ neurons with ReLU activation, while we consider five-layer FFNNs ($784 \rightarrow 100 \rightarrow 75 \rightarrow 50 \rightarrow 40 \rightarrow 25 \rightarrow 10$) with ReLU activation.

Each model was trained for 40 epochs using the Adam optimizer. INNs that were trained using mixed monotonicity and FFNNs that were trained using IBP have $\epsilon_{\text{test}} = 0.1$ and $\kappa_{\text{nom}} = 0.75$. From epochs 1 to 10, κ and ϵ are set to 0 so the models undergo regular (nonrobust) training. From epochs 11 to 20, ϵ and κ are linearly increased such that at epoch 20, $\epsilon = \epsilon_{\text{test}}$ and $\kappa = \kappa_{\text{nom}}$. Regarding training INNs, we follow the non-Euclidean monotone operator framework described in [17]. We impose $\mu_{\infty, [\eta]^{-1}}(W) \leq 0$ for some $\eta \in \mathbb{R}_{>0}^n$.

10 FFNNs and 10 INNs were trained; 5 of each were trained using IBP or mixed monotonicity (for feedforward and implicit, respectively) and 5 of each were trained without

¹Code to reproduce the experiments is available at <https://github.com/davydovalexander/robust-inn-mn>.

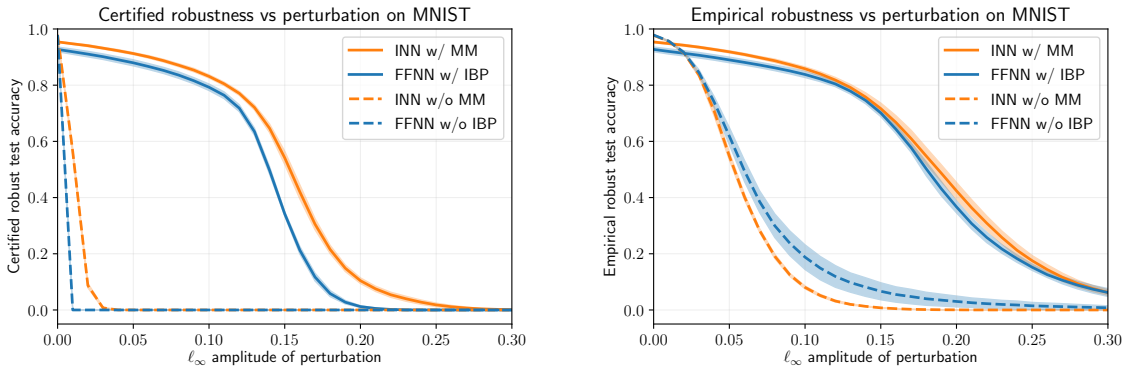


Fig. 1: Performance comparison on the MNIST test data between INNs trained with and without mixed monotonicity (MM) and 5-layer FFNNs trained with and without IBP. The INNs have 89710 trainable parameters and FFNNs have 93200 trainable parameters. The left plot shows the certified robust accuracy of the models computed using either MM or IBP while the right plot shows the empirical robustness of the models against a PGD attack. In each plot, dark lines correspond to the mean accuracy across 5 neural networks while light envelopes around the dark lines correspond to one standard deviation.

any robust optimization (i.e., $\epsilon_{\text{test}} = 0$). Figure 1 provides plots of certified adversarial robustness via the corresponding interval reachability technique and the empirically-observed robustness against a projected gradient descent (PGD) attack.

Evaluation summary: Regarding certified robustness, at an ℓ_∞ perturbation radius of 0.1, we observe that INNs trained using mixed monotonicity had, on average, an accuracy of 83.13%, while FFNNs trained using IBP had, on average an accuracy of 79.26%. We additionally observe that at the cost of a few percentage points in clean accuracy, both INNs and FFNNs trained robustly vastly outperform non-robustly trained models in both certified and empirical robustness. For example, at an ℓ_∞ perturbation radius of 0.1, INNs trained without mixed monotonicity have an empirical accuracy of 8.04%, while INNs trained with mixed monotonicity have an accuracy of 85.84%, indicating an order of magnitude improvement in empirical robustness.

IX. CONCLUSION

We develop a computationally efficient algorithm for training robust INNs. Moreover, we provide theoretical and empirical evidence in support of the following claims: (i) robustly-trained INNs are more robust than comparably-trained FFNNs, (ii) inclusion functions provide tighter estimates than Lipschitz constants, (iii) robustly-trained networks enjoy much stronger robustness properties than their non-robustly trained counterparts.

REFERENCES

- [1] B. Amos and J. Z. Kolter, “OptNet: Differentiable optimization as a layer in neural networks,” in *International Conference on Machine Learning*, 2017.
- [2] S. Bai, J. Z. Kolter, and V. Koltun, “Deep equilibrium models,” in *Advances in Neural Information Processing Systems*, 2019.
- [3] L. El Ghaoui, F. Gu, B. Travacca, A. Askari, and A. Tsai, “Implicit deep learning,” *SIAM Journal on Mathematics of Data Science*, vol. 3, no. 3, pp. 930–958, 2021.
- [4] A. Kag, Z. Zhang, and V. Saligrama, “RNNs incrementally evolving on an equilibrium manifold: A panacea for vanishing and exploding gradients?” in *International Conference on Learning Representations*, 2020.
- [5] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” in *International Conference on Learning Representations*, 2014.
- [6] E. Wong and J. Z. Kolter, “Provable defenses against adversarial examples via the convex outer adversarial polytope,” in *International Conference on Machine Learning*, 2018, pp. 5286–5295.
- [7] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” in *International Conference on Machine Learning*, 2018.
- [8] A. Virmaux and K. Scaman, “Lipschitz regularity of deep neural networks: analysis and efficient estimation,” in *Advances in Neural Information Processing Systems*, vol. 31, 2018, p. 3839–3848.
- [9] P. Pauli, A. Koch, J. Berberich, P. Kohler, and F. Allgower, “Training robust neural networks using Lipschitz bounds,” *IEEE Control Systems Letters*, vol. 6, pp. 121–126, 2022.
- [10] S. Wang, K. Pei, J. Whitehouse, J. Yang, and S. Jana, “Formal security analysis of neural networks using symbolic intervals,” in *USENIX Conference on Security Symposium*, 2018, p. 1599–1614.
- [11] S. Goyal, K. Dvijotham, R. Stanforth, R. Bunel, C. Qin, J. Uesato, R. Arandjelovic, T. A. Mann, and P. Kohli, “Scalable verified training for provably robust image classification,” in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 4841–4850.
- [12] H. Zhang, H. Chen, C. Xiao, S. Goyal, R. Stanforth, B. Li, D. Boning, and C.-J. Hsieh, “Towards stable and efficient training of verifiably robust neural networks,” in *International Conference on Learning Representations*, 2020.
- [13] C. Wei and J. Z. Kolter, “Certified robustness for deep equilibrium models via interval bound propagation,” in *International Conference on Learning Representations*, 2022.
- [14] M. Fazlyab, M. Morari, and G. J. Pappas, “Safety verification and robustness analysis of neural networks via quadratic constraints and semidefinite programming,” *IEEE Transactions on Automatic Control*, vol. 67, no. 1, pp. 1–15, 2022.
- [15] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer, “Reluplex: An efficient smt solver for verifying deep neural networks,” in *Computer Aided Verification*. Springer International Publishing, 2017, pp. 97–117.
- [16] C. Pabbaraju, E. Winston, and J. Z. Kolter, “Estimating Lipschitz constants of monotone deep equilibrium models,” in *International Conference on Learning Representations*, 2021.
- [17] S. Jafarpour, A. Davydov, A. V. Proskurnikov, and F. Bullo, “Robust implicit networks via non-Euclidean contractions,” in *Advances in Neural Information Processing Systems*, 2021.
- [18] M. Revay, R. Wang, and I. R. Manchester, “Lipschitz bounded equilibrium networks,” *arXiv preprint arXiv:2010.01732*, 2020.
- [19] T. Chen, J. B. Lasserre, V. Magron, and E. Pauwels, “Semialgebraic representation of monotone deep equilibrium models and applications to certification,” in *Advances in Neural Information Processing Systems*, 2021.
- [20] S. Jafarpour, M. Abate, A. Davydov, F. Bullo, and S. Coogan, “Robustness certificates for implicit neural networks: A mixed monotone contractive approach,” in *Learning for Dynamics & Control Conference*, June 2022.
- [21] L. Yang and N. Ozay, “Tight decomposition functions for mixed monotonicity,” in *IEEE Conference on Decision and Control (CDC)*, 2019, pp. 5318–5322.