# Bisimulations for Neural Network Reduction

Pavithra Prabhakar

Kansas State University, Manhattan, KS, 66506
`pprabhakar@ksu.edu`

**Abstract.** We present a notion of bisimulation that induces a reduced network which is semantically equivalent to the given neural network. We provide a minimization algorithm to construct the smallest bisimulation equivalent network. Reductions that construct bisimulation equivalent neural networks are limited in the scale of reduction. We present an approximate notion of bisimulation that provides semantic closeness, rather than, semantic equivalence, and quantify semantic deviation between the neural networks that are approximately bisimilar. The latter provides a trade-off between the amount of reduction and deviations in the semantics.

**Keywords:** Neural Networks · Bisimulation · Verification · Reduction.

## 1 Introduction

Neural networks (NN) with small size are conducive for both automated analysis and explainability. Rigorous automated analysis using formal methods has gained momentum in recent years owing to the safety-criticality of the application domains in which NN are deployed [3,16,12,20,19,11]. For instance, NN are an integral part of control, perception and guidance of autonomous vehicles. However, the scalability of these analysis techniques, for instance, for computing output range for safety analysis [6,12], is limited by the large size of the neural networks encountered and the computational complexity due to the presence of non-linear activation functions. In this paper, we borrow ideas from formal methods to design novel network reduction techniques with formal relations between the given and the reduced networks, that can be applied to reduce the verification time. It can also potentially impact explainability by presenting to the user a smaller network with guaranteed bounds on the deviation from the original network.

Bisimulation [14] is a classical notion of equivalence between systems in process algebra that guarantees that processes that are bisimilar satisfy the same set of properties specified in certain branching time logics [2]. A bisimulation is an equivalence relation on the states of a system that requires similar behaviors with respect to one step of computation, which then inductively guarantees global behavioral equivalence. Bisimulation algorithm [2] allows one to construct the smallest systems, bisimulation quotients, that are bisimilar to a given (finite state) system.

Our first result consists of a definition of bisimulation for neural networks, namely, NN-bisimulation, that defines a notion of equivalence between neural networks. The challenge arises from the fact that neural networks semantically have multiple parallel threads of computation that are both branching and merging at each step of computation. We observe that the global equivalence can be established by imposing a *one-step backward pre-sum equivalence*, wherein we require two nodes that belong to the same class to agree on the biases, the activation functions, and the pre-sums, wherein a pre-sum corresponds to the sum of the weights on the incoming edges from a given equivalence class. Our notion resembles that of probabilistic bisimulation [13], however, our notion is based on pre-sum equivalences rather than post-sum equivalences. We define a quotienting operation on an NN with respect to a bisimulation that yields a smaller network which is input-output equivalent to the given network. We also show that there exists a coarsest bisimulation which yields the smallest neural network with respect to the quotienting operation. We provide a minimization algorithm that outputs this smallest neural network.

The notion of bisimulation can be stringent, since, it preserves the exact input-output relation. It has been observed, for instance, in the context of control systems, that a strict notion of equivalence, such as, bisimulation, does not allow for drastic reduction in state-space, thereby, motivating the notion of approximate bisimulation. Approximate bisimulations [9,8] have a notion of distance between states, and allow a bounded $\epsilon$ deviation between the executions of the systems in each step. The notion of approximate bisimulation was successfully used to construct smaller systems in the setting of dynamical systems and control synthesis [10].

We extend the notion of NN-bisimulation to an approximate notion, wherein we require nodes belonging to the same class to have bounded deviation, $\epsilon$, in the biases and the pre-sums. The quotienting operation no more results in a unique reduced network, but a set of reduced networks. Moreover, these reduced networks may not have the same input-output relation as the given neural network. However, we provide a bound on the deviation in the semantics of two approximately bisimilar NNs. It gives rise to a nice trade-off between the amount of reduction and the deviation in the semantics, that translates to a trade-off in the precision and verification time in an approximation based verification scheme.

*Related work.* Neural network reduction techniques have been explored in different contexts. There is extensive literature on compression techniques, see, for instance, surveys on network compression [5,4]. However, these techniques typically do not provide formal guarantees on the relation between the given and reduced systems. Abstraction techniques [15,17,7] computing over-approximations of the input-output relations have been explored in several works, however, they use slightly different kinds of networks such as interval neural networks and abstract neural networks, or are limited to certain kinds of activation functions such as ReLU. Notions of bisimulation for DNNs have not been explored much in the literature. Equivalence between DNNs is explored [1], however, the work is restricted to ReLU functions and does not consider approximate notions.

## 2     Preliminaries

Let $[k]$ denote the set $\{0, 1, \cdots, k\}$ and $(k]$ the set $\{1, 2, \cdots, k\}$. Let $\mathbb{R}$ denote the set of real numbers. We use $|x|$ to denote the absolute value of $x \in \mathbb{R}$. Given a set $A$, we use $|A|$ to denote the number of elements of $A$. Given a function $f : A \to \mathbb{R}$, we define the infinity norm of $f$ to be the supremum of the absolute values of elements in the range of $f$, that is, $\|f\|_\infty = \sup_{a \in A} |f(a)|$. Given functions $f : A \to B$ and $g : B \to C$, the composition of $f$ and $g$, $g \circ f : A \to C$, is given by, for all $a \in A$, $g \circ f(a) = g(f(a))$.

*Partitions.* Given a set $\mathcal{S}$, a (finite) parition of $\mathcal{S}$ is a set $\mathcal{P} = \{\mathcal{S}_1, \cdots, \mathcal{S}_n\}$, such that $\bigcup_i \mathcal{S}_i = \mathcal{P}$ and $\mathcal{S}_i \cap \mathcal{S}_j = \emptyset$ for all $i \neq j$. We refer to each element of a partition as a region or a group. A partition $\mathcal{P}$ of $\mathcal{S}$ can be seen as an equivalence relation on $\mathcal{S}$, given by the relation $s_1 \mathcal{P} s_2$ whenever $s_1$ and $s_2$ belong to the same group of the partition. Given two partitions $\mathcal{P}$ and $\mathcal{P}'$, we say that $\mathcal{P}$ is finer than $\mathcal{P}'$ (or equivalently, $\mathcal{P}'$ is coarser than $\mathcal{P}$), denoted $\mathcal{P} \preceq \mathcal{P}'$, if for every $S \in \mathcal{P}$, there exists $S' \in \mathcal{P}'$ such that $S \subseteq \mathcal{S}'$.
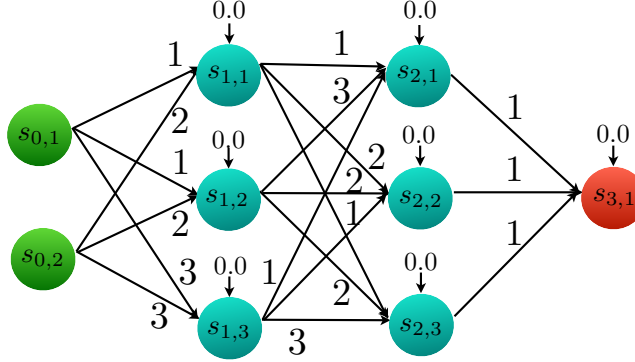
## 3     Neural Networks

In this section, we present the preliminaries regarding the neural network. Recall that a neural network (NN) consists of a layered set of nodes or neurons, including an input layer, an output layer and one or more hidden layers. Each node except those in the input layer are annotated with a bias and an activation function, and there are weighted edges between nodes of adjacent layers. We capture these elements of a neural network using a tuple in the following definition.

**Definition 1.** *A neural network (NN) is a tuple* $\mathcal{N} = \big(k, \mathcal{A}ct, \{\mathcal{S}_i\}_{i \in [k]}, \{W_i\}_{i \in (k]},$ $\{b_i\}_{i \in (k]}, \{A_i\}_{i \in (k]}\big)$*, where:*

- *$k$ represents the number of layers (except the input layer);*
- *$\mathcal{A}ct$ is a set of activation functions;*
- *for each $i \in [k]$, $\mathcal{S}_i$ is a set of nodes of layer $i$, we assume $\mathcal{S}_i \cap \mathcal{S}_j = \emptyset$ for $i \neq j$;*
- *for each $i \in (k]$, $W_i : \mathcal{S}_{i-1} \times \mathcal{S}_i \to \mathbb{R}$ is the weight function that captures the weight on the edges between nodes at layer $i-1$ and $i$;*
- *for each $i \in (k]$, $b_i : \mathcal{S}_i \to \mathbb{R}$ is the bias function that associates a bias with nodes of layer $i$;*
- *for each $i \in (k]$, $A_i : \mathcal{S}_i \to \mathcal{A}ct$ is an activation association function that associates an activation function with each node of layer $i$.*

$\mathcal{S}_0$ and $\mathcal{S}_k$ are the set of nodes corresponding to the input and output layers, respectively. We will fix the NN $\mathcal{N} = \big(k, \mathcal{A}ct, \{\mathcal{S}_i\}_{i \in [k]}, \{W_i\}_{i \in (k]}, \{b_i\}_{i \in (k]}, \{A_i\}_{i \in (k]}\big)$ for the rest of the paper.

**Fig. 1.** Neural Network $\mathcal{N}$

*Example 1.* The neural network $\mathcal{N}$ shown in Figure 1 consists of an input layer with 2 nodes, 2 hidden layers with 3 nodes each, and an output layer. The weights on the edges are shown, for instance, $W_2(s_{1,2}, s_{2,2}) = 2$. The biases are all 0s and the activation functions are all ReLUs (not shown).

In the sequel, the central notion to the definition of bisimulation will be the total weight on the incoming edges for a node $s'$ of the $i$-th layer from a set of nodes $\mathcal{S}$ of the $i-1$-st layer. We will capture this using the notion of a pre-sum, denoted $PreSum_i^{\mathcal{N}}(\mathcal{S}, s')$. For instance, $PreSum_2^{\mathcal{N}}(\{s_{1,1}, s_{1,2}\}, s_{2,2}) = 2 + 2 = 4$.

**Definition 2.** *Given a set $\mathcal{S} \subseteq \mathcal{S}_{i-1}$ and $s' \in \mathcal{S}_i$, we define $PreSum_i^{\mathcal{N}}(\mathcal{S}, s') = \sum_{s \in \mathcal{S}} W_i(s, s')$.*

Next, we capture the operational behavior of a neural network. A valuation $v$ for the $i$-th layer of $\mathcal{N}$ refers to an assignment of real-values to all the nodes in $\mathcal{S}_i$, that is, $v : \mathcal{S}_i \to \mathbb{R}$. Let $Val(\mathcal{S}_i)$ denote the set of all valuations for the $i$-th layer of $\mathcal{N}$. By the operational semantics of $\mathcal{N}$, we mean the assignments for all the layers of $\mathcal{N}$, that are obtained from an assignment for the input layer. We define $[\![\mathcal{N}]\!]_i(v)$, which given a valuation $v$ for layer $i - 1$, returns the corresponding valuation for layer $i$ according to the semantics of $\mathcal{N}$. The valuation for the output layer of $\mathcal{N}$ is then obtained by the composition of the functions $[\![\mathcal{N}]\!]_i$.

**Definition 3.** *The semantics of the $i$-the layer is the function $[\![\mathcal{N}]\!]_i : Val(\mathcal{S}_{i-1}) \to Val(\mathcal{S}_i)$, where for any $v \in Val(\mathcal{S}_{i-1})$, $[\![\mathcal{N}]\!]_i(v) = v'$, given by*

$$\forall s' \in \mathcal{S}_i, v'(s') = A_i(s')\Big( \sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')v(s) + b_i(s')\Big).$$

To capture the input-output semantics, we compose these one layer semantics. More precisely, we define $[\![\mathcal{N}]\!]^i$ to be the composition of the first $i$ layers, that

is, $[\![\mathcal{N}]\!]^i(v)$ provides the valuation of the $i$-th layer given $v$ as input. It is defined inductively as:

$$[\![\mathcal{N}]\!]^1 = [\![\mathcal{N}]\!]_1$$

$$\forall i \in (k) \backslash \{1\}, [\![\mathcal{N}]\!]^i = [\![\mathcal{N}]\!]_i \circ [\![\mathcal{N}]\!]^{i-1}$$

**Definition 4.** *The input-output semantic function, represented by* $[\![\mathcal{N}]\!] : Val(\mathcal{S}_0) \rightarrow Val(\mathcal{S}_k)$, *is defined as:*

$$[\![\mathcal{N}]\!] = [\![\mathcal{N}]\!]^k$$

The notion of bisimulation requires the notion of a partition of the nodes of $\mathcal{N}$. We define a partition on $\mathcal{N}$ as an indexed set of partitions each corresponding to a layer.

**Definition 5.** *A partition of an NN $\mathcal{N}$ is an indexed set $\mathcal{P} = \{\mathcal{P}_i\}_{i \in [k]}$, where for every $i$, $\mathcal{P}_i$ is a partition of $\mathcal{S}_i$.*

*A note on Lipschitz Continuity.* A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is said to be Lipschitz continuous if there exists a constant $L(f)$, referred to as Lipschitz constant for $f$, such that for all $x, y \in \mathbb{R}^m$,

$$\|f(x) - f(y)\|_\infty \leq L(f)\|x - y\|_\infty.$$

Several activation functions including ReLU, Leaky ReLU, SoftPlus, Tanh, Sigmoid, ArcTan and Softsign are known to be 1-Lipschitz continuous [18], that is, satisfy the above constraint with $L(f) = 1$. In fact, the function $[\![\mathcal{N}]\!]$ is itself Lipschitz continuous, when the activation functions are Lipschitz continuous [18]. We will use $L(\mathcal{N})$ to denote an upper bound on $L([\![\mathcal{N}]\!]^i)$ over all $i$. Hence, given an input $v$, we know that $\|[\![\mathcal{N}]\!]^i(v)\|_\infty \leq L(\mathcal{N})\|v\|_\infty$.

## 4 NN-bisimulation and Semantic Equivalence

In this section, we define a notion of bisimulation on neural networks, which induces a reduced system that is equivalent to the given network. A partition of $\mathcal{N}$ is an NN-bisimulation if the biases and activation functions associated with the nodes in any region are the same, and the pre-sums of nodes in any region with respect to any region of the previous layer are the same.

**Definition 6.** *An NN-bisimulation for $\mathcal{N}$ is a partition $\mathcal{P} = \{\mathcal{P}_i\}_{i \in [k]}$ such that for all $i \in (k]$, $\mathcal{S} \in \mathcal{P}_{i-1}$ and $s'_1, s'_2 \in \mathcal{S}_i$ with $s'_1 \mathcal{P}_i s'_2$, the following hold:*

1. *$A_i(s'_1) = A_i(s'_2)$,*
2. *$b_i(s'_1) = b_i(s'_2)$, and*
3. *$PreSum_i^{\mathcal{N}}(\mathcal{S}, s'_1) = PreSum_i^{\mathcal{N}}(\mathcal{S}, s'_2)$.*

Our notion is inspired by the well-known notion of probabilistic bisimulation [13], where post-sums are used instead of pre-sums to characterize which nodes have the same branching structure. Though neural networks consist of branching in

both forward and backward directions, surprisingly, just pre-sum equivalence suffices to guarantee input-output relation equivalence.

Bisimulation naturally induces a reduced system, which corresponds to merging the nodes in a group of the partition, and choosing a representative node from the group to assign the activation functions, biases and pre-sums. We represent the reduced system obtained by taking the quotient of $\mathcal{N}$ with respect to a bisimulation $\mathcal{P}$ as $\mathcal{N}/\mathcal{P}$.

**Definition 7.** *Given an NN-bisimulation $\mathcal{P}$ for $\mathcal{N}$, the reduced system $\mathcal{N}/\mathcal{P} = \left(k, \mathcal{A}ct, \{\widehat{\mathcal{S}}_i\}_{i\in[k]}, \{\widehat{W}_i\}_{i\in(k]}, \{\widehat{b}_i\}_{i\in(k]}, \{\widehat{A}_i\}_{i\in(k]}\right)$, where:*

1. $\forall i \in [k], \widehat{\mathcal{S}}_i = \mathcal{P}_i$;
2. $\forall i \in (k], \widehat{s} \in \widehat{\mathcal{S}}_{i-1}, \widehat{s}' \in \widehat{\mathcal{S}}_i, \widehat{W}_i(\widehat{s}, \widehat{s}') = PreSum_i^{\mathcal{N}}(\widehat{s}, s')$ *for some $s' \in \widehat{s}'$.*
3. $\forall i \in (k], \widehat{s}' \in \widehat{\mathcal{S}}_i, \widehat{b}_i(\widehat{s}') = b_i(s')$ *for some $s' \in \widehat{s}'$.*
4. $\forall i \in (k], \widehat{s}' \in \widehat{\mathcal{S}}_i, \widehat{A}_i(\widehat{s}') = A_i(s')$ *for some $s' \in \widehat{s}'$.*

Note that though the definition depends on the choice of $s'$, the reduced system is unique, since, from the definition of NN-bisimulation, the values of biases, activation functions and pre-sums, corresponding to different choices of $s'$ within a group are the same. We also use just bisimulation to refer to NN-bisimulation.

In order to formally establish the connection between the NN $\mathcal{N}$ and its reduction $\mathcal{N}/\mathcal{P}$, we define a mapping from the valuations of $\mathcal{N}$ to those of $\mathcal{N}/\mathcal{P}$, but only for certain valuations that are consistent in that they map all the related nodes in $\mathcal{P}$ to the same value.

**Definition 8.** *A valuation $v \in Val(\mathcal{S}_i)$ is $\mathcal{P}$-consistent, if for all $s_1, s_2 \in \mathcal{S}_i$, if $s_1 \mathcal{P}_i s_2$, then $v(s_1) = v(s_2)$.*

Our first result is that a consistent input valuation leads to a consistent output valuation, when $\mathcal{P}$ is a bisimulation. We show this for a particular layer; the extension to the whole network follows from a simple inductive reasoning.

**Lemma 1.** *Let $\mathcal{P}$ be a bisimulation on $\mathcal{N}$. If $v_1 \in Val(\mathcal{S}_{i-1})$ is $\mathcal{P}$-consistent, then $v_2 = [\![\mathcal{N}]\!]_i(v_1)$ is $\mathcal{P}$-consistent.*

*Proof.* Let $s', s'' \in \mathcal{S}_i$ such that $s'\mathcal{P}_i s''$. We need to show that $v_2(s') = v_2(s'')$. $v_2(s') = A_i(s')(\sum_{s\in\mathcal{S}_{i-1}} W_i(s, s')v_1(s) + b_i(s')) = A_i(s')(\sum_{\mathcal{S}\in\mathcal{P}_{i-1}}\sum_{s\in\mathcal{S}} W_i(s, s') v_1(s) + b_i(s'))$. Since, $v_1$ is $\mathcal{P}$-consistent, for each $\mathcal{S}$, we have a value $v_1^{\mathcal{S}}$ that all elements of $\mathcal{S}$ are mapped to by $v_1$, that is, $v_1^{\mathcal{S}} = v_1(s)$ for all $s \in \mathcal{S}$. Replacing $v_1(s)$ for each $s$, by $v_1^{\mathcal{S}}$, we obtain $v_2(s') = A_i(s')(\sum_{\mathcal{S}\in\mathcal{P}_{i-1}}\sum_{s\in\mathcal{S}} W_i(s, s')v_1^{\mathcal{S}} + b_i(s'))$. From the definition of pre-sum, we can replace $\sum_{s\in\mathcal{S}} W_i(s, s')$ by $PreSum_i^{\mathcal{N}}(\mathcal{S}, s')$, which is also equal to $PreSum_i^{\mathcal{N}}(\mathcal{S}, s'')$ from the definition of bisimulation, since $\mathcal{P}$ is a bisimulation and $s'\mathcal{P}_i s''$. Also, $A_i(s') = A_i(s'')$ and $b_i(s') = b_i(s'')$. So, we obtain, $v_2(s') = A_i(s'')(\sum_{\mathcal{S}\in P_{i-1}} PreSum_i^{\mathcal{N}}(\mathcal{S}, s'')v_1^{\mathcal{S}} + b_i(s''))$. Expanding back $PreSum_i^{\mathcal{N}}(\mathcal{S}, s'')$, and $v_1^{\mathcal{S}} = v_1(s)$ for all $s \in \mathcal{S}$, we obtain $v_2(s') = A_i(s'')(\sum_{\mathcal{S}\in P_{i-1}}\sum_{s\in\mathcal{S}} W_i(s, s'')v_1(s) + b_i(s'')) = v_2(s'')$.

Note that if we do not group together the nodes in the input and output layers, there is a bijection between $\mathcal{S}_0$ and $\widehat{\mathcal{S}}_0$ and $\mathcal{S}_k$ and $\widehat{\mathcal{S}}_k$, and hence, a bijection between their valuations. We will show that both $\mathcal{N}$ and $\mathcal{N}/\mathcal{P}$ have the "same" input-output relation modulo the bijection between their nodes. First, we define a formal relation between $\mathcal{P}$-consistent valuations of $\mathcal{N}$ and valuations of $\mathcal{N}/\mathcal{P}$.

**Definition 9.** *Let $\mathcal{P}$ be a bisimulation on $\mathcal{N}$, and $v \in Val(\mathcal{S}_i)$ be a $\mathcal{P}$-consistent valuation. The abstraction of $v$, denoted, $\alpha(v)_{\mathcal{N},\mathcal{P}} \in Val(\widehat{\mathcal{S}}_i)$, is defined as, for every $\hat{s} \in \widehat{\mathcal{S}}_i$, $\alpha(v)_{\mathcal{N},\mathcal{P}}(\hat{s}) = v(s)$ for some $s \in \hat{s}$.*

Note that $\alpha(v)_{\mathcal{N},\mathcal{P}}$ is well defined, since, from the $\mathcal{P}$-consistency of $v$, $v(s)$ is the same for any choice of $s \in \hat{s}$. When $\mathcal{N}$ and $\mathcal{P}$ are clear from the context, we will drop the subscript and write $\alpha(v)_{\mathcal{N},\mathcal{P}}$ as just $\alpha(v)$. The next result states that the output of the $i$-th layer of $\mathcal{N}/\mathcal{P}$ with the abstraction of a $\mathcal{P}$-consistent valuation $v$ of the $i - 1$-st layer of $\mathcal{N}$ as input, results in a valuation that is the abstraction of the output of the $i$-th layer of $\mathcal{N}$ on input $v$. In other words, it says that propagating a valuation for one-step in $\mathcal{N}$ is the same as propagating its abstraction in $\mathcal{N}/\mathcal{P}$.

**Lemma 2.** *Let $\mathcal{P}$ be a bisimulation on $\mathcal{N}$, and $v \in Val(\mathcal{S}_i)$ be $\mathcal{P}$-consistent. Then, $\alpha([\![\mathcal{N}]\!]_i(v)) = [\![\mathcal{N}/\mathcal{P}]\!]_i(\alpha(v))$.*

*Proof.* From Lemma 1, we know that $v' = [\![\mathcal{N}]\!]_i(v)$ is $\mathcal{P}$-consistent. Hence, for any $\hat{s}' \in \widehat{\mathcal{S}}_i$, $\alpha(v')(\hat{s}') = v'(s')$ for some (any) $s' \in \hat{s}'$. Let us fix $s' \in \hat{s}'$.

$$\alpha(v')(\hat{s}') = v'(s') = A_i(s')\Big( \sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')v(s) + b_i(s') \Big)$$

from the semantics of $\mathcal{N}$. Further,

$$\sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')v(s) = \sum_{\mathcal{S} \in \mathcal{P}_{i-1}} \sum_{s \in \mathcal{S}} W_i(s, s')v(s)$$

From $\mathcal{P}$-consistency of $v$, $v(s) = \alpha(v)(\mathcal{S})$ for any $s \in \mathcal{S}$. Hence,

$$\sum_{s \in \mathcal{S}} W_i(s, s')v(s) = \sum_{s \in \mathcal{S}} W_i(s, s')\alpha(v)(\mathcal{S}) = PreSum_i^{\mathcal{N}}(\mathcal{S}, s')\alpha(v)(\mathcal{S}).$$
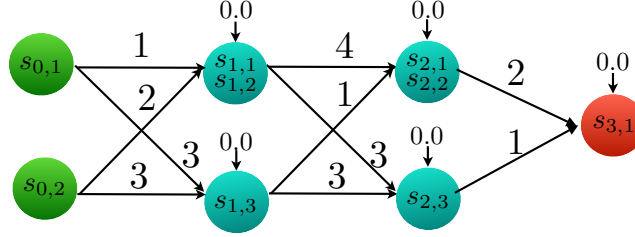
From the definition of $\mathcal{N}/\mathcal{P}$, $\widehat{A}_i(\hat{s}') = A_i(s'), \widehat{W}_i(\mathcal{S}, \hat{s}') = PreSum_i^{\mathcal{N}}(\mathcal{S}, s'), \widehat{b}_i(\hat{s}') = b_i(s')$, and $\mathcal{P}_{i-1} = \widehat{\mathcal{S}}_{i-1}$. From $\mathcal{P}$-consistency of $v$, $\alpha(v)(\mathcal{S}) = v(s)$ for any $s \in \mathcal{S}$. Therefore, for any $\hat{s}' \in \widehat{\mathcal{S}}_i$,

$$\alpha([\![\mathcal{N}]\!]_i(v))(\hat{s}') = \alpha(v')(\hat{s}')$$

$$= \widehat{A}_i(\hat{s}')\Big( \sum_{\mathcal{S} \in \widehat{\mathcal{S}}_{i-1}} \widehat{W}_i(\mathcal{S}, \hat{s}')\alpha(v)(\mathcal{S}) + \widehat{b}_i(\hat{s}') \Big) = [\![\mathcal{N}/\mathcal{P}]\!]_i(\alpha(v))(\hat{s}').$$

The following theorem follows by composing the results from Lemma 2 for the different layers.

**Theorem 1.** *Given $\mathcal{P}$ a bisimulation on $\mathcal{N}$, and $v \in Val(\mathcal{S}_0)$ that is $\mathcal{P}$-consistent, we have $\alpha(\llbracket \mathcal{N} \rrbracket(v)) = \llbracket \mathcal{N}/\mathcal{P} \rrbracket(\alpha(v))$.*

*Proof.* We can show by induction on $i$ that $\alpha(\llbracket \mathcal{N} \rrbracket^i(v)) = \llbracket \mathcal{N}/\mathcal{P} \rrbracket^i(\alpha(v))$.



**Fig. 2.** Reduced System $\mathcal{N}/\mathcal{P}$

*Example 2.* Consider a partition $\mathcal{P}$ for the NN $\mathcal{N}$ in Figure 1 where each node appears as a region by itself except for the regions $\mathcal{S}_1 = \{s_{1,1}, s_{1,2}\}$, and $\mathcal{S}_2 = \{s_{2,1}, s_{2,2}\}$. We can verify that this is a bisimulation. For instance, $PreSum_2^{\mathcal{N}}(\mathcal{S}_1, s_{2,1}) = 1 + 3$ and $PreSum_2^{\mathcal{N}}(\mathcal{S}_1, s_{2,2}) = 2 + 2$, which are the same. The reduced system is given by the NN $\mathcal{N}/\mathcal{P}$ in Figure 2. Here, $\widehat{W}_2(\mathcal{S}_1, \mathcal{S}_2) = 4$.

## 5    $\delta$-NN-bisimulation and Semantic Closeness

NN-bisimulation provides a foundation for reducing a neural network while preserving the input-output relation. However, existence of such bisimulations leading to equivalent reduced networks with much fewer neurons is limited in that for many networks no bisimulation quotient may lump together lot of nodes. Hence, we relax the notion of bisimulation to an approximate notion wherein we allow potentially large reductions, however, the reduced systems may not be semantically equivalent, but only be semantically close to the given neural network. We quantify the deviation of the reduced system in terms of the "deviation" of the approximate notion from the exact bisimulation.
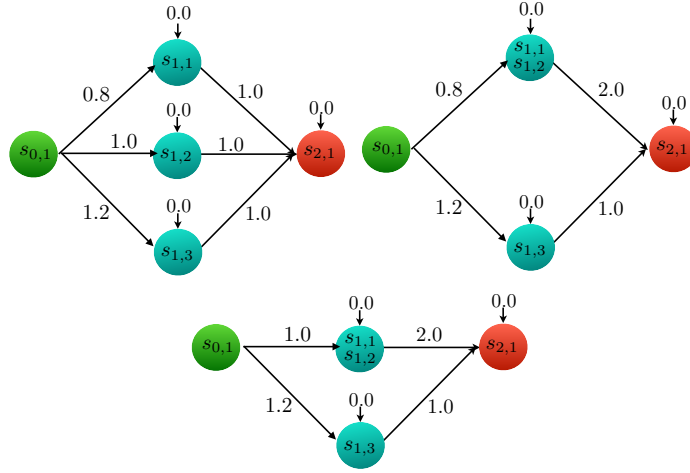
The approximation notion of bisimulation we consider is inspired by the notion of approximate bisimulation in the context of dynamical systems [9,8]. We essentially relax the requirement of the NN-bisimulation that the biases and pre-sums match by allowing them to be within a $\delta$. This is formalized in the following definition.

**Definition 10.** *A $\delta$-NN-bisimulation for an NN $\mathcal{N}$ and $\delta \geq 0$ is a partition $\mathcal{P} = \{\mathcal{P}_i\}_{i \in [k]}$ such that for all $i \in (k]$, $\mathcal{S} \in \mathcal{P}_{i-1}$ and $s'_1, s'_2 \in \mathcal{S}_i$ with $s'_1 \mathcal{P}_i s'_2$, the following hold:*

1. $A_i(s'_1) = A_i(s'_2)$,
2. $|b_i(s'_1) - b_i(s'_2)| \leq \delta$, and
3. $|PreSum_i^{\mathcal{N}}(\mathcal{S}, s'_1) - PreSum_i^{\mathcal{N}}(\mathcal{S}, s'_2)| \leq \delta$.

We will also use $\delta$-bisimulation to refer to $\delta$-NN-bisimulation. The reduced system can be constructed similar to that for NN-bisimulation. However, the choice of the nodes $s' \in \hat{s}'$ used to construct the weights and biases of the reduced system could lead to different neural networks. Hence, we obtain a finite set of possibilities for the reduced system that we denote by $\mathcal{N}/_\delta \mathcal{P}$.



**Fig. 3.** Illustration of $\mathcal{N}^*/_\delta \mathcal{P}$ on NN $\mathcal{N}^*$

*Example 3.* Consider the NN $\mathcal{N}^*$ in Figure 3 (top left) and a partition $\mathcal{P} = \{\mathcal{P}_i\}_i$, where $\mathcal{P}_0 = \{\{s_{0,1}\}\}$, $\mathcal{P}_2 = \{\{s_{2,1}\}\}$ and $\mathcal{P}_1 = \{\{s_{1,1}, s_{1,2}\}, \{s_{1,3}\}\}$, that is, $\mathcal{P}$ merges nodes $s_{1,1}$ and $s_{1,2}$. Note that $\mathcal{P}$ is a $\delta$-bisimulation on $\mathcal{N}^*$ for $\delta = 0.2$. For instance, $PreSum_1^{\mathcal{N}^*}(\{s_{0,1}\}, s_{1,1}) = 0.8$ and $PreSum_1^{\mathcal{N}^*}(\{s_{0,1}\}, s_{1,2}) = 1.0$ whose difference is $\leq 0.2 = \delta$. $\mathcal{N}/_\delta \mathcal{P}$ consists of $\mathcal{N}_1^*$ and $\mathcal{N}_2^*$ in Figure 3 (top right and bottom), one which is obtained by choosing the pre-sum corresponding to $s_{1,1}$ and other by choosing the pre-sum corresponding to $s_{1,2}$.

Our objective is to give a bound on the deviation of the semantics of any $\mathcal{N}' \in \mathcal{N}/_\delta \mathcal{P}$ from that of $\mathcal{N}$. We start by quantifying this deviation in one step of computation. For that, we extend the notion of consistent valuations to an

approximate notion, wherein we require the valuations of related states to be within a bound rather than match exactly.

**Definition 11.** *A valuation $v \in Val(\mathcal{S}_i)$ is $\epsilon, \mathcal{P}$-consistent, if for all $s_1, s_2 \in \mathcal{S}_i$ with $s_1 \mathcal{P}_i s_2$, $|v(s_1) - v(s_2)| \leq \epsilon$.*

Our next step is to establish a relation between the valuation propagation in $\mathcal{N}$ and any $\mathcal{N}' \in \mathcal{N}/_\delta \mathcal{P}$ analogous to Lemma 2. First, we will need to relax the notion of the abstraction of a valuation, however, unlike in the previous case, we obtain a set of abstractions $\alpha^\epsilon(v)$.

**Definition 12.** *Let $\mathcal{P}$ be a partition of $\mathcal{N}$, and $v \in Val(\mathcal{S}_i)$. The $\epsilon$-abstraction of $v$, denoted, $\alpha^\epsilon(v)_{\mathcal{N},\mathcal{P}}$, consists of $\hat{v} \in Val(\mathcal{P}_i)$ such that for all $\hat{s} \in \mathcal{P}_i, s \in \hat{s}$, $|\hat{v}(\hat{s}) - v(s)| \leq \epsilon$.*

When $\mathcal{N}$ and $\mathcal{P}$ are clear from the context, we will drop the subscript and write $\alpha^\epsilon(v)_{\mathcal{N},\mathcal{P}}$ as just $\alpha^\epsilon(v)$. The next result states that the $\epsilon$-abstraction for any $\epsilon, \mathcal{P}$-consistent valuation is non-empty.

**Proposition 1.** *Let $v \in Val(\mathcal{S}_i)$ be an $\epsilon, \mathcal{P}$-consistent valuation. Then $\alpha^\epsilon(v)$ is non-empty.*

*Proof.* Note that the valuation $\hat{v}$, given by $\hat{v}(\hat{s}) = v(s)$ for some $s \in \hat{s}$ gives a valuation in $\alpha^\epsilon(v)_{\mathcal{N},\mathcal{P}}$.

The converse of the above theorem also holds with a slight modification of the error.

**Proposition 2.** *Let $v \in Val(\mathcal{S}_i)$, such that $\alpha^\epsilon(v)_{\mathcal{N},\mathcal{P}}$ is non-empty. Then, $v$ is a $2\epsilon, \mathcal{P}$-consistent valuation.*

*Proof.* Note that there is some $\hat{v}$, such that $\forall \hat{s} \in \mathcal{P}_i, s \in \hat{s}, |\hat{v}(\hat{s}) - v(s)| \leq \epsilon$. Then for all $s, s' \in \hat{s}, |v(s) - v(s')| \leq 2\epsilon$.

Now, we give a bound on the deviation of the output of the $i$-th layer of $\mathcal{N}/_\delta \mathcal{P}$ from that of $\mathcal{N}$ in terms of the deviation in their inputs. Let $L(A_i) = \max_{s' \in \mathcal{S}_i} L(A_i(s'))$.

**Lemma 3.** *Let $\mathcal{P}$ be a $\delta$-bisimulation on $\mathcal{N}$, and $v \in Val(\mathcal{S}_{i-1})$ be $\epsilon, \mathcal{P}$-consistent. Then, for every $\hat{v} \in \alpha^\epsilon(v)$, and $\mathcal{N}' \in \mathcal{N}/_\delta \mathcal{P}$,*

$$[\![\mathcal{N}']\!]_i(\hat{v}) \in \alpha^{\epsilon'}([\![\mathcal{N}]\!]_i(v)),$$

*where $\epsilon' = a_i \epsilon + b_i$, $a_i = L(A_i)|\mathcal{S}_{i-1}|\|W_i\|_\infty$, and $b_i = L(A_i)(|\mathcal{P}_{i-1}|\|v\|_\infty + 1)\delta$.*

*Proof.* Let $v' = [\![\mathcal{N}]\!]_i(v)$ and $\hat{v}' = [\![\mathcal{N}']\!]_i(\hat{v})$. We need to show that $\hat{v}' \in \alpha^{\epsilon'}(v')$. Consider any $\hat{s}' \in \widehat{\mathcal{S}}_i$ and $s' \in \hat{s}'$. We need to show that $|\hat{v}'(\hat{s}') - v'(s')| \leq \epsilon'$.

Since, $\hat{v}' = [\![\mathcal{N}']\!]_i(\hat{v})$, from the semantics of $\mathcal{N}'$, we have

$$\hat{v}'(\hat{s}') = \widehat{A}_i(\hat{s}')\Big(\sum_{\hat{s} \in \widehat{\mathcal{S}}_{i-1}} \widehat{W}_i(\hat{s}, \hat{s}')\hat{v}(\hat{s}) + \widehat{b}_i(\hat{s}')\Big),$$

and from the fact that $\mathcal{N}' \in \mathcal{N}/_\delta\mathcal{P}$, we have $\widehat{W}_i(\hat{s}, \hat{s}') = PreSum_i^{\mathcal{N}}(\hat{s}, s'_1)$ for some $s'_1 \in \hat{s}'$, $\hat{b}_i(\hat{s}') = b_i(s'_2)$ for some $s'_2 \in \hat{s}'$. Since $\mathcal{P}$ is a $\delta$-bisimulation, $|PreSum_i^{\mathcal{N}}(\hat{s}, s'_1) - PreSum_i^{\mathcal{N}}(\hat{s}, s')| \le \delta$, $|b_i(s'_2) - b_i(s')| \le \delta$, and $\widehat{A}_i(\hat{s}') = A_i(s')$. Therefore,

$$\hat{v}'(\hat{s}') = A_i(s')\Big(\sum_{\hat{s} \in \mathcal{P}_{i-1}} (PreSum_i^{\mathcal{N}}(\hat{s}, s') + \delta_{\hat{s}})\hat{v}(\hat{s}) + b_i(s') + \delta_{s'}\Big),$$

$$= A_i(s')\Big(\sum_{\hat{s} \in \mathcal{P}_{i-1}} PreSum_i^{\mathcal{N}}(\hat{s}, s')\hat{v}(\hat{s}) + b_i(s') + \epsilon_1 + \delta_{s'}\Big),$$

where $\epsilon_1 = \sum_{\hat{s} \in \mathcal{P}_{i-1}} \delta_{\hat{s}}\hat{v}(\hat{s})$ and $\delta_{\hat{s}}, \delta_{s'} \in [-\delta, \delta]$. We will examine the terms in the above expression in more detail.

$$\sum_{\hat{s} \in \mathcal{P}_{i-1}} PreSum_i^{\mathcal{N}}(\hat{s}, s')\hat{v}(\hat{s}) = \sum_{\hat{s} \in \mathcal{P}_{i-1}} [\sum_{s \in \hat{s}} W_i(s, s')\hat{v}(\hat{s})]$$

(Further, since, $\hat{v} \in \alpha^\epsilon(v)$, we have for any $s \in \hat{s}$, $|\hat{v}(\hat{s}) - v(s)| \le \epsilon$.)

$$= \sum_{\hat{s} \in \mathcal{P}_{i-1}} [\sum_{s \in \hat{s}} W_i(s, s')(v(s) + \epsilon_s)] = \sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')(v(s) + \epsilon_s)]$$

$$= \sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')v(s) + \sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')\epsilon_s$$

Plugging the above into the expression for $\hat{v}'(\hat{s}')$, we obtain

$$\hat{v}'(\hat{s}') = A_i(s')\Big(\sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')v(s) + b_i(s') + \epsilon_1 + \epsilon_2 + \delta_{s'}\Big)$$

where $\epsilon_2 = \sum_{\hat{s} \in \mathcal{P}_{i-1}} \delta_{\hat{s}}\hat{v}(\hat{s})$. Note that the expression for $\hat{v}'(\hat{s}')$ looks similar to $v'(s') = A_i(s')\big(\sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')v(s) + b_i(s')\big)$ except for the additional error terms $\epsilon_1 + \epsilon_2 + \delta_{s'}$. From the Lipschitz continuity of $A_i(s')$, we obtain

$$|\hat{v}'(\hat{s}') - v'(s')| \le L(A_i)(s')(|\epsilon_1 + \epsilon_2 + \delta_{s'}|)$$

Note that $L(A_i)(s') \le L(A_i)$, $|\epsilon_1| = |\sum_{s \in \mathcal{S}_{i-1}} W_i(s, s')\epsilon_s| \le |\mathcal{S}_{i-1}|\|W_i\|_\infty\epsilon$, $|\epsilon_1| = |\sum_{\hat{s} \in \mathcal{P}_{i-1}} \delta_{\hat{s}}\hat{v}(\hat{s})| \le |\mathcal{P}_{i-1}|\delta\|v\|_\infty$, and $|\delta_{s'}| \le \delta$. Hence,

$$|\hat{v}'(\hat{s}') - v'(s')| \le L(A_i)(s')(|\epsilon_1 + \epsilon_2 + \delta_{s'}|) \le L(A_i)(|\mathcal{S}_{i-1}|\|W_i\|_\infty\epsilon + |\mathcal{P}_{i-1}|\delta\|v\|_\infty + \delta)$$

$$= L(A_i)|\mathcal{S}_{i-1}|\|W_i\|_\infty\epsilon + L(A_i)(|\mathcal{P}_{i-1}|\|v\|_\infty + 1)\delta = \epsilon'$$

Lemma 3 provides a bound on the error propagation in one step. The next theorem provides a global bound on the deviation of the output of the reduced system with respect to that of the given neural network. Let $L(A) = \max_i L(A_i)$, $|\mathcal{P}| = \max_i |\mathcal{P}_i|$, $\|W\|_\infty = \max_i \|W_i\|_\infty$ and $|\mathcal{S}| = |\max_i \mathcal{S}_i|$.

**Theorem 2.** *Let $\mathcal{P}$ be a $\delta$-bisimulation on $\mathcal{N}$, and $v \in Val(\mathcal{S}_0)$ be $\epsilon, \mathcal{P}$-consistent. Then, for every $\hat{v} \in \alpha^\epsilon(v)$, and $\mathcal{N}' \in \mathcal{N}/_\delta\mathcal{P}$,*

$$[\![\mathcal{N}']\!](\hat{v}) \in \alpha^{\epsilon''}([\![\mathcal{N}]\!](v)),$$

*where $\epsilon'' = [(2/a)^k - 1]b/(2a-1)$, $a = L(A)|\mathcal{S}|\|W\|_\infty$, and $b = L(A)(|\mathcal{P}|L(\mathcal{N})\|v\|_\infty + 1)\delta$.*

*Proof.* Let us define:
$$v_0 = v, \hat{v}_0 = \hat{v}, \epsilon_0 = \epsilon'_0 = 0$$

and for all $i \in (k]$,
$$v_i = [\![\mathcal{N}]\!]_i(v), \hat{v}_i = [\![\mathcal{N}']\!]_i(\hat{v}).$$
$$\epsilon'_i = a\epsilon_{i-1} + b, \epsilon_i = 2\epsilon'_i.$$

We will show by induction on $i$ that for all $i \in [k]$, $v_i$ is $\epsilon_i, \mathcal{P}$-consistent and $\hat{v}_i \in \alpha^{\epsilon_i}(v_i)$.

**Base case:** Base case trivially holds from the assumptions of the theorem statement.

**Induction Step:** For $i \in (k]$, we know from Lemma 3, that if $v_{i-1} \in Val(\mathcal{S}_{i-1})$ is $\epsilon_{i-1}, \mathcal{P}$-consistent and $\hat{v}_{i-1} \in \alpha^{\epsilon_{i-1}}(v_{i-1})$, then $\hat{v}_i = [\![\mathcal{N}']\!]_i(\hat{v}_{i-1}) \in \alpha^{\epsilon'}([\![\mathcal{N}]\!]_i(v_{i-1})) = \alpha^{\epsilon'}(v_i)$, where $\epsilon' = a_i\epsilon_{i-1} + b_i$.

$$a_i = L(A_i)|\mathcal{S}_{i-1}|\|W_i\|_\infty \leq L(A)|\mathcal{S}|\|W\|_\infty = a,$$

$$b_i = L(A_i)(|\mathcal{P}_{i-1}|\|v_i\|_\infty + 1)\delta \leq L(A)(|\mathcal{P}|L(\mathcal{N})\|v\|_\infty + 1)\delta = b$$

Hence, $\epsilon \leq \epsilon'_i$ and $\hat{v}_i \in \alpha^{\epsilon_i}(v_i)$. Further from Proposition 2, we obtain $v_i$ is $\epsilon_i, \mathcal{P}$-consistent.

We will show that $\epsilon'_k = \epsilon''$. Unrolling the recursive equation, we obtain $\epsilon'_i = 2a\epsilon'_{i-1} + b = (2a)^i\epsilon'_0 + [(2a)^{i-1} + \cdots + 1]b = [(2/a)^i - 1]b/(2a - 1)$. Hence,

$$\epsilon'_k = [(2/a)^k - 1]b/(2a - 1) = \epsilon''$$

We finish the proof by noting that $[\![\mathcal{N}']\!](\hat{v}) = \hat{v}_k \in \alpha^{\epsilon'_k}(v_k) = \alpha^{\epsilon''}([\![\mathcal{N}]\!](v))$.

*Remark 1.* Note that for $\delta = 0$, all the notions and results reduce to that of NN-bisimulation.

## 6   Minimization algorithm

In this section, we show that there is a coarsest NN-bisimulation for a given NN, that encompasses all other bisimulations. This implies that the induced reduced network with respect to this coarsest bisimulation is the smallest NN-bisimulation equivalent network. We will provide an algorithm that outputs the coarsest NN-bisimulation.

We note that a coarsest $\delta$-NN-bisimulation may not exist in general. For instance, consider the NN $\mathcal{N}^*$ from Figure 3, along with the 0.2-bisimulation $\mathcal{P}$

that induces the reduced systems $\mathcal{N}_1^*$ and $\mathcal{N}_2^*$. There is another 0.2-bisimulation $\mathcal{P}'$ which is obtained by merging $s_{1,2}$ and $s_{1,3}$ instead of $s_{1,1}$ and $s_{1,2}$ as in $\mathcal{P}$. Note that the reduced networks in $\mathcal{N}^*/_{0.2}\mathcal{P}$ and $\mathcal{N}^*/_{0.2}\mathcal{P}'$ have the same size. However, there is no 0.2-bisimulation that is coarser than both $\mathcal{P}$ and $\mathcal{P}'$, since, that would require merging $s_{1,1}, s_{1,2}$ and $s_{1,3}$, which would violate the 0.2 bound on the difference between the pre-sums of $s_{1,1}$ and $s_{1,3}$.

The broad algorithm for minimization consists of starting with a partition in which all the nodes in a layer are merged together and then splitting them such that the regions in the partition respect the activation functions, biases and the pre-sums. We use the function $SplitActBias(\mathcal{S})$ in the algorithm that splits a set of nodes $\mathcal{S}$ into maximal groups such that the elements in each group agree on the activation functions and the biases. More precisely, $SplitActBias(\mathcal{S})$ takes $\mathcal{S} \subseteq \mathcal{P}_i$ as input and returns a partition $\mathcal{P}_\mathcal{S}$ such that for all $s_1, s_2 \in \mathcal{S}$, $s_1 \mathcal{P}_\mathcal{S} s_2$ if and only if $A_i(s_1) = A_i(s_2)$ and $b_i(s_1) = b_i(s_2)$. Further, we split those regions that have nodes with inconsistent pre-sums. Next, we define inconsistent pairs of regions with respect to pre-sums and the corresponding splitting operations.

**Definition 13.** *Given a partition $\mathcal{P} = \{\mathcal{P}_i\}_i$ of NN $\mathcal{N}$, a region $\mathcal{S}' \in \mathcal{P}_i$ is inconsistent in $\mathcal{N}$ with respect to $\mathcal{S} \in \mathcal{P}_{i-1}$, written $(\mathcal{S}', \mathcal{S})$ inconsistent, if there exist $s_1', s_2' \in \mathcal{S}'$, such that $PreSum_i^\mathcal{N}(\mathcal{S}, s_1') \neq PreSum_i^\mathcal{N}(\mathcal{S}, s_2')$.*

The algorithm searches for inconsistent pairs $(\mathcal{S}', \mathcal{S})$ and splits $\mathcal{S}'$ into maximal groups such that all nodes in a group have the same pre-sum with respect to $\mathcal{S}$. More precisely, $SplitPre(\mathcal{S}', \mathcal{S})$ takes $\mathcal{S}'$ and $\mathcal{S}$ as input and returns a partition $\mathcal{P}'$ of $\mathcal{S}'$ such that $PreSum_i^\mathcal{N}(\mathcal{S}, s_1') = PreSum_i^\mathcal{N}(\mathcal{S}, s_1')$ if and only if $s_1' \mathcal{P}' s_2'$.

---

**Algorithm 1**: **MinNN:** Minimization Algorithm

**Input**: A NN $\mathcal{N}$
**Output**: Coarsest Bisimulation $\mathcal{P}$, and Minimized NN $\mathcal{N}/\mathcal{P}$

1 **begin**
2     $\mathcal{P} = \{\mathcal{S}_0\}$
3     **for** $i \in (k]$ **do**
4         $\lfloor$  $\mathcal{P} = \mathcal{P} \cup SplitActBias(\mathcal{S}_i)$
5     **while** *Exists $\mathcal{S}, \mathcal{S}' \in \mathcal{P}$, such that $(\mathcal{S}, \mathcal{S}')$ inconsistent* **do**
6         $\lfloor$  $\mathcal{P} = \mathcal{P} \backslash \{\mathcal{S}'\} \cup SplitPre(\mathcal{S}', \mathcal{S})$
7     **return** Return $\mathcal{P}$ and $\mathcal{N}/\mathcal{P}$
8 **end**

---

Next, we show that Algorithm 1 returns the coarsest bisimulation, and hence, the reduced network is the smallest bisimulation equivalent network.

**Definition 14.** *A partition $\mathcal{P}$ of $\mathcal{N}$ is the coarsest bisimulation, if it is an NN bisimulation and it is coarser than every NN-bisimulation $\mathcal{P}'$ of $\mathcal{N}$.*

**Theorem 3.** *Algorithm 1 terminates and returns the coarsest bisimulation $\mathcal{P}$ of $\mathcal{N}$.*

*Proof.* Termination of the algorithm is straightforward, since, if there exists an inconsistent pair $(\mathcal{S}', \mathcal{S})$, then $SplitPre(\mathcal{S}', \mathcal{S})$ splits $\mathcal{S}'$ into at least two regions. Hence, the number of regions in $\mathcal{P}$ strictly increases. However, since, $\mathcal{N}$ has finitely many nodes, the number of regions in $\mathcal{P}$ is upper-bounded.

Next, we will argue that $\mathcal{P}$ that is returned is an NN-bisimulation. After the $SplitActBias(\mathcal{S}_i)$ operations, $\mathcal{P}$ only consists of regions which agree on the activation functions and biases. When the while loop terminates, there are no inconsistent pairs, that is, the pre-sum condition of the bisimulation definition is satisfied. Hence, the value of $\mathcal{P}$ when exiting the while loop is an NN-bisimulation.

To show that $\mathcal{P}$ is the coarsest bisimulation, it remains to show that $\mathcal{P}$ is coarser than any bisimulation of $\mathcal{N}$. Let $\mathcal{P}'$ be any bisimulation of $\mathcal{N}$. We will show that $\mathcal{P}'$ is finer than $\mathcal{P}$ at every stage of the algorithm.

Note that after exiting the for loop, $\mathcal{P}$ contains the maximal groups which agree on both the activation functions and biases. Every region of $\mathcal{P}'$ has to agree on the activation functions and biases, since it is a bisimulation. So, every region of $\mathcal{P}'$ is contained in some regions of $\mathcal{P}$, that is, $\mathcal{P}' \preceq \mathcal{P}$.

Next, we show that $\mathcal{P}' \preceq \mathcal{P}$ is an invariant for the while loop, that is, if it holds at the beginning of the loop, then it also holds at the end of the loop. So, when the while loop exits, we still have $\mathcal{P}' \preceq \mathcal{P}$. More precisely, we need to show that if $\mathcal{P}' \preceq \mathcal{P}$, then replacing $\mathcal{S}'$ by $SplitPre(\mathcal{S}', \mathcal{S})$ will still result in a partition that is coarser than $\mathcal{P}'$. In particular, we need to ensure that each region of $\mathcal{P}'$ that is contained in $\mathcal{S}'$ is not split by the $SplitPre(\mathcal{S}', \mathcal{S})$ operation. Suppose a region $\mathcal{S}'' \subseteq \mathcal{S}'$ of $\mathcal{P}'$ is split, then there exists $s_1'', s_2'' \in \mathcal{S}''$ such that $PreSum_i^{\mathcal{N}}(\mathcal{S}, s_1'') \neq PreSum_i^{\mathcal{N}}(\mathcal{S}, s_1'')$. But $\mathcal{S}$ is the disjoint union of some sets $\{\mathcal{S}_1'', \cdots, \mathcal{S}_l''\}$ of $\mathcal{P}'$. Hence, $PreSum_j^{\mathcal{N}}(\mathcal{S}_i'', s_1'') \neq PreSum_i^{\mathcal{N}}(\mathcal{S}_i'', s_1'')$ for some $i$, since $PreSum_i^{\mathcal{N}}(\mathcal{S}, s'') = \sum_i PreSum_i^{\mathcal{N}}(\mathcal{S}_i'', s'')$. However, this contradicts the fact that $\mathcal{P}'$ is an NN-bisimulation.

Next, we present some complexity results on checking if a partition is a bisimulation/$\delta$-bisimulation, complexity of constructing reduced systems from a bisimulation/$\delta$-bisimulation and the complexity of computing the coarsest bisimulation.

**Theorem 4.** *Given an NN $\mathcal{N}$, a partition $\mathcal{P}$ and $\epsilon \geq 0$, checking if $\mathcal{P}$ is a bisimulation and checking if $\mathcal{P}$ is an $\delta$-NN-bisimulation both take time $O(m)$, where $m$ is the number of edges of $\mathcal{N}$. Further, constructing $\mathcal{N}/\mathcal{P}$ for some $\mathcal{N}' \in \mathcal{N}/_\delta \mathcal{P}$ takes time $O(m)$ as well.*

*Proof.* To check if $\mathcal{P}$ is a bisimulation, we can iterate over all the nodes in a region to check if they have same activation function, bias, and pre-sums with respect to every region of $\mathcal{P}_{i-1}$. In doing so, we need to access each node and each edge at most once, hence, the complexity is bounded by $O(m)$. For the $\delta$-bisimulation, we need to check if the biases and pre-sums are within $\epsilon$. We can compute the biases and pre-sums in one pass over the network in time

$O(m)$ as before. Then we can find the max and min values of the bias/pre-sum values within each region, and check if the max and min values are within $\epsilon$, this will take time $O(m)$. For constructing the reduced system, we need to find the activation functions and biases of all the nodes in the reduced system, and the weight on the edge between two groups. The total computation needs to access each edge at most once.

**Theorem 5.** *The minimization algorithm has a time complexity of $O(\hat{n}(m + n \log n))$, where $n$ is the number of nodes and $m$ is the number of edges of $\mathcal{N}$, and $\hat{n}$ is the number of nodes in the minimized neural network.*

*Proof. SplitActBias($\mathcal{S}_i$)* needs to sort the elements in every group by the activation function/bias values, hence, takes time $O(n \log n)$. Finding an inconsistent pair takes the same time as checking whether $\mathcal{P}$ is a bisimulation, that is, $O(m)$. *SplitPre($\mathcal{S}', \mathcal{S}$)* take time at most $O(m)$ to compute the pre-sums and $O(n \log n)$ to split. Replacing $\mathcal{S}'$ by *SplitPre($\mathcal{S}', \mathcal{S}$)* takes time at most $O(\hat{n})$ which is upper-bounded by $O(n)$. So, each loop takes time $O(m + n \log n)$. The number of iterations of the while loop is upper bounded by the number of regions in the minimized neural network, that is, $O(\hat{n})$. Hence, the minimization algorithm has a runtime of $O(\hat{n}(m + n \log n))$.

## 7   Conclusions

We presented the notions of bisimulation and approximate bisimulation for neural networks that provide semantic equivalence and semantic closeness, respectively, and are applicable to neural networks with a wide range of activation functions. These provide foundational theoretical tools for exploring the trade-off between the amount of reduction and the semantic deviation in an approximation based verification paradigm for neural networks. Our future work will focus on experimental analysis of this trade-off on large scale neural networks. The notions of bisimulation explored are syntactic in nature, and we will explore semantic notions in the future. We provide a minimization algorithm for finding the smallest NN that is bisimilar to a given neural network. Though a unique minimal network does not exist with respect to $\delta$-bisimulations, we will explore heuristics for constructing small networks that are $\delta$-bisimilar.

## References

1. Ashok, P., Hashemi, V., Kretinsky, J., Mohr, S.: Deepabstract: Neural network abstraction for accelerating verification (2020)
2. Baier, C., Katoen, J.P.: Principles of Model Checking (Representation and Mind Series). The MIT Press (2008)
3. Bunel, R., Turkaslan, I., Torr, P.H.S., Kohli, P., Kumar, M.P.: Piecewise linear neural network verification: A comparative study. CoRR (2017)
4. Cheng, Y., Wang, D., Zhou, P., Zhang, T.: A survey of model compression and acceleration for deep neural networks. CoRR (2017)

5. Deng, L., Li, G., Han, S., Shi, L., Xie, Y.: Model compression and hardware acceleration for neural networks: A comprehensive survey. Proceedings of the IEEE **108**(4), 485–532 (2020)
6. Dutta, S., Jha, S., Sankaranarayanan, S., Tiwari, A.: Output range analysis for deep feedforward neural networks. In: Dutle, A., Muñoz, C., Narkawicz, A. (eds.) NASA Formal Methods. pp. 121–138. Springer International Publishing (2018)
7. Elboher, Y.Y., Gottschlich, J., Katz, G.: An abstraction-based framework for neural network verification. In: Lahiri, S.K., Wang, C. (eds.) Computer Aided Verification - 32nd International Conference, CAV 2020, Los Angeles, CA, USA, July 21-24, 2020, Proceedings, Part I. Lecture Notes in Computer Science, vol. 12224, pp. 43–65. Springer (2020)
8. Girard, A., Julius, A.A., Pappas, G.J.: Approximate simulation relations for hybrid systems. Discrete Event Dynamic Systems **18**(2), 163–179 (2008)
9. Girard, A., Pappas, G.J.: Approximate bisimulation relations for constrained linear systems. Automatica **43**(8), 1307–1317 (2007)
10. Girard, A., Pola, G., Tabuada, P.: Approximately bisimilar symbolic models for incrementally stable switched systems. In: Proceedings of the International Conference on Hybrid Systems: Computation and Control. pp. 201–214 (2008)
11. Huang, X., Kroening, D., Kwiatkowska, M.Z., Ruan, W., Sun, Y., Thamo, E., Wu, M., Yi, X.: Safety and trustworthiness of deep neural networks: A survey. CoRR **abs/1812.08342** (2018)
12. Katz, G., Barrett, C.W., Dill, D.L., Julian, K., Kochenderfer, M.J.: Reluplex: An efficient SMT solver for verifying deep neural networks. CoRR (2017)
13. Larsen, K.G., Skou, A.: Bisimulation through probabilistic testing. Inf. Comput. **94**(1), 1–28 (1991)
14. Milner, R.: Communication and Concurrency. Prentice-Hall, Inc (1989)
15. Prabhakar, P., Afzal, Z.R.: Abstraction based output range analysis for neural networks (2019)
16. Pulina, L., Tacchella, A.: An abstraction-refinement approach to verification of artificial neural networks. In: Touili, T., Cook, B., Jackson, P. (eds.) Computer Aided Verification. pp. 243–257. Springer Berlin Heidelberg, Berlin, Heidelberg (2010)
17. Sotoudeh, M., Thakur, A.: Abstract neural networks. In: Static Analysis Symposium (2020)
18. Virmaux, A., Scaman, K.: Lipschitz regularity of deep neural networks: analysis and efficient estimation. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) Advances in Neural Information Processing Systems 31, pp. 3835–3844. Curran Associates, Inc. (2018)
19. Xiang, W., Musau, P., Wild, A.A., Lopez, D.M., Hamilton, N., Yang, X., Rosenfeld, J.A., Johnson, T.T.: Verification for machine learning, autonomy, and neural networks survey. CoRR **abs/1810.01989** (2018)
20. Xiang, W., Tran, H., Johnson, T.T.: Output reachable set estimation and verification for multi-layer neural networks. CoRR **abs/1708.03322** (2017)