



Certified Private Inference on Neural Networks via Lipschitz-Guided Abstraction Refinement*

Edoardo Manino¹, Bernardo Magri¹,
Mustafa A. Mustafa^{1,2}, and Lucas C. Cordeiro¹

¹ The University of Manchester, Department of Computer Science,
Oxford Road, Manchester M13 9PL, United Kingdom

² imec-COSIC, KU Leuven, Kasteelpark Arenberg 10, B-3001, Belgium

Abstract

Private inference on neural networks requires running all the computation on encrypted data. Unfortunately, neural networks contain a large number of non-arithmetic operations, such as ReLU activation functions and max pooling layers, which incur a high latency cost in their encrypted form. To address this issue, the majority of private inference methods replace some or all of the non-arithmetic operations with a polynomial approximation. This step introduces approximation errors that can substantially alter the output of the neural network and decrease its predictive performance. In this paper, we propose a Lipschitz-Guided Abstraction Refinement method (LiGAR), which provides strong guarantees on the global approximation error. Our method is iterative, and leverages state-of-the-art Lipschitz constant estimation techniques to produce increasingly tighter bounds on the worst-case error. At each iteration, LiGAR designs the least expensive polynomial approximation by solving the dual of the corresponding optimization problem. Our preliminary experiments show that LiGAR can easily converge to the optimum on medium-sized neural networks.

1 Introduction

The rise of very large neural network models has made it attractive to offload all the related computation to remote servers in the cloud [27]. While this configuration eases the computational burden on the client side, it opens the door to privacy concerns, as the user is forced to send their private data to a third-party machine [22]. An ideal solution would involve a private inference scheme (see Figure 1), where the server evaluates the neural network on encrypted data and the client decrypts the result afterwards.

In order to realise such vision, each operation performed by the neural network needs to be adapted for operating on ciphertext. Homomorphic encryption schemes enable this translation as follows. For a given encryption scheme $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec})$, we say that \mathcal{E} is additively

*This work is funded by the EPSRC grant EP/T026995/1 entitled “EnnCore: End-to-End Conceptual Guarding of Neural Architectures” under *Security for all in an AI enabled society*.

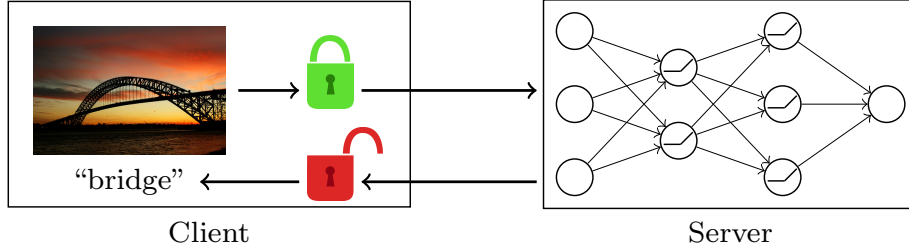


Figure 1: Private inference requires running a neural network on encrypted data.

homomorphic if for two ciphertexts $c_1 \leftarrow \text{Enc}(m_1)$ and $c_2 \leftarrow \text{Enc}(m_2)$, one can compute $c_3 = c_1 + c_2$ where $\text{Dec}(c_3) = m_1 + m_2$. One example of such scheme is due to Paillier [21]. Similarly, we say that \mathcal{E} is multiplicative homomorphic if for two ciphertexts $c_1 \leftarrow \text{Enc}(m_1)$ and $c_2 \leftarrow \text{Enc}(m_2)$, one can compute $c_3 = c_1 \cdot c_2$ where $\text{Dec}(c_3) = m_1 \cdot m_2$. The RSA encryption scheme [23] is of this flavor. A Fully Homomorphic Encryption (FHE) scheme is one that combines both operations. In particular, it allows for a client to compute an encryption of its secret input x and a server to homomorphically evaluate any circuit C over the ciphertext; as an output the server produces an encryption of the evaluation of $C(x)$ while it never learns anything about x . The first construction of FHE by Gentry [7] is based on the hardness of lattices problems and it is extremely inefficient. Since then, a lot of work has been done trying to bring FHE closer to the practical realm [2, 3, 15].

In this regard, neural networks contain a significant number of operations that cannot be represented as a combination of additions and multiplications: e.g. ReLU, hyperbolic tangent, sigmoid, softmax and maxpool [1]. In order to avoid the computational cost of directly executing these in FHE, existing private inference methods propose to approximate them with polynomials [4, 18, 16, 14, 13]. Such polynomial conversion introduces some approximation error at every layer of the neural network. For some inputs, the error can accumulate and drastically distort the output of the network [6]. While most existing works claim that the approximated network has similar predictive accuracy to the original one, some authors express the need for stronger guarantees on the approximation error [6]. To the best of our knowledge, this need has not been met so far. The only exception is the theoretical analysis in [14], which proves that the output error is bounded, but does not provide a constructive algorithm to compute it.

In this paper, we address this research gap by designing neural network approximations with formal guarantees on the maximum output error. Essentially, we recognise that this research objective is fundamentally a problem of equivalence verification. As such, we leverage state-of-the-art verification techniques, and we adapt them to the polynomial approximation setting. More specifically, our contributions are the following:

- we propose an abstraction technique to represent the polynomial error;
- we adapt state-of-the-art reachability and Lipschitz estimation techniques to the architecture of our abstracted network;
- we derive a closed-form solution to the problem of optimal error allocation;
- we introduce our Lipschitz-Guided Abstraction Refinement (LiGAR) method to iteratively minimise the over-estimation caused by our abstraction;
- we show that LiGAR can easily converge on medium-sized neural networks;

- we discuss the practical limitations of a worst-case method like LiGAR.

All the code for the experiments in this paper is publicly available at [17].

2 Related Work

The literature on privacy-preserving machine learning is vast. Here, we only cover the main research trends on the topic of private inference.

Neural Architecture Search (NAS). Since the main bottleneck is executing non-arithmetic operations on encrypted data, some works propose to select neural architectures that contain the bare minimum of them. Delphi [18] uses NAS to select which ReLUs we can approximate with degree-2 polynomials. The rest is replaced with garbled circuits to avoid accuracy drop. Similarly, SAFENet [16] uses NAS to replace a subset of the ReLUs with polynomials of degree in $[0, 3]$. Other methods include a finetuning phase after altering the network architecture: DeepReDuce [9] uses pruning and distillation techniques, while SNL [5] replaces some ReLUs with the identity function before re-training.

Low-degree approximations. Earlier approaches aimed at replacing all non-arithmetic operations with low-degree polynomials. CryptoNets [8] replaces all activation functions with x^2 . Similarly, the authors of [4] use degree-2 polynomials but includes batch normalisation layers to keep the input distribution centred with the polynomial approximation. As discussed in [6], these low-degree approximation only allow for shallow neural networks, as it is very challenging to fine-tune the approximations without incurring in the gradient explosion problem.

High-degree approximations. An alternate approach is approximating non-arithmetic operations with high-degree polynomials [14]. This way, the approximated network does not require any fine-tuning. The major obstacle to such goal is obtaining stable high-degree approximations of the sign operator [12]. From it, a number of popular non-arithmetic operations can be derived, including ReLU and max-pooling. A recent example of this line of research is [13], which focuses on the challenge applying fully-homomorphic encryption to large-scale networks.

3 Preliminaries

3.1 Neural Network Definitions and Notation

In this paper, we focus on feedforward fully-connected neural networks with ReLU activation functions. This neural architecture encapsulates all the challenges in designing private inference methods. Extension of our method to general neural networks is possible, but we leave it for future work.

Definition 1 (ReLU Layer). *Define the Rectified Linear Unit (ReLU) as the function $\sigma(z_i) = \max(z_i, 0)$ with potential $z_i \in \mathbb{R}$. A ReLU layer is a mapping $h : \mathbb{R}^n \rightarrow \mathbb{R}^n$ that applies σ element-wise, i.e. $(h(z))_i = \sigma(z_i)$.*

Definition 2 (Affine Layer). *An affine layer is a mapping $g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ such that $g(z) = Wz + b$ where $W \in \mathbb{R}^{n \times m}$ is a matrix and $b \in \mathbb{R}^n$ is a column vector.*

Definition 3 (Feedforward Fully-Connected ReLU Neural Network). *A feedforward fully-connected ReLU neural network $f(x)$ is the composition of an alternating sequence of affine and ReLU layers:*

$$y = f(x) = g_k \circ h_{k-1} \circ g_{k-1} \circ \dots \circ h_1 \circ g_1(x) \tag{1}$$

3.2 Univariate Polynomial Approximation

The ReLU function is non-arithmetic and thus necessitates replacement. This is a classic problem of univariate polynomial approximation [26], with the goal of minimising the maximum error over a finite interval $[\underline{z}, \bar{z}]$:

$$p_d^* = \arg \min_{p_d} \left\{ \max_{z \in [\underline{z}, \bar{z}]} \{|p_d(z) - \sigma(z)|\} \right\} \tag{2}$$

where d is the polynomial degree. The solution of Equation 2 can be computed via Remez’s algorithm, of which many implementations exist (e.g. [20]); we present an example of it in Figure 3a. Crucially, univariate polynomial approximation becomes numerically unstable as the degree d grows into the hundreds. The authors of [12] address this issue by decomposing the minimax polynomial into a composition of lower-degree ones.

4 LiGAR: Lipschitz-Guided Abstraction Refinement

We define the problem of finding the optimal polynomial approximation of a neural network $f(x)$ as follows:

Definition 4 (Optimal Polynomial Approximation). *Given a (possibly infinite) set of inputs \mathcal{X} and a target output error δ_y , find the least computationally expensive polynomial approximations $p_{d_i}^* \approx \sigma(z_i)$ of each ReLU function i in $f(x)$, such that the maximum output error satisfies $\|f(x) - \hat{f}(x)\|_\infty \leq \delta_y$, where all the ReLUs in $\hat{f}(x)$ have been replaced with their respective approximations.*

Here, we interpret the least expensive computation objective in terms of minimising the sum of all degrees $d_{tot} = \sum_i d_i$. However, we note that other choices are possible (e.g. multiplication depth in [13]). These are correlated with the polynomial degree, but not equivalent to it.

4.1 Approximation Error Abstraction

Substituting a ReLU activation functions with a polynomials introduces some input-dependent approximation error $\epsilon(z_i) = |\sigma(z_i) - p_{d_i}^*(z_i)|$. In order to facilitate the solution of the optimisation problem in Definition 4, we propose to remove the dependency of $\epsilon(z_i)$ on z_i by letting its value vary freely in the interval $\epsilon_i \in [-\delta_i, +\delta_i]$. In this way, we can isolate the contribution of each approximation $p_{d_i}^*$ on the output error (see Section 4.4).

We present a high-level picture of our abstraction in Figure 2. The abstracted network uses the same weights, biases and activations of the original network. However, a new input ϵ_i is added to the output of each activation function:

Definition 5 (Polynomial Neural Network Abstraction). *Call $\hat{f}(x)$ the polynomial approximation of neural network $f(x)$, such that $|\sigma(z_i) - p_{d_i}(z_i)| \leq \delta_i$ for each ReLU i with potentials*

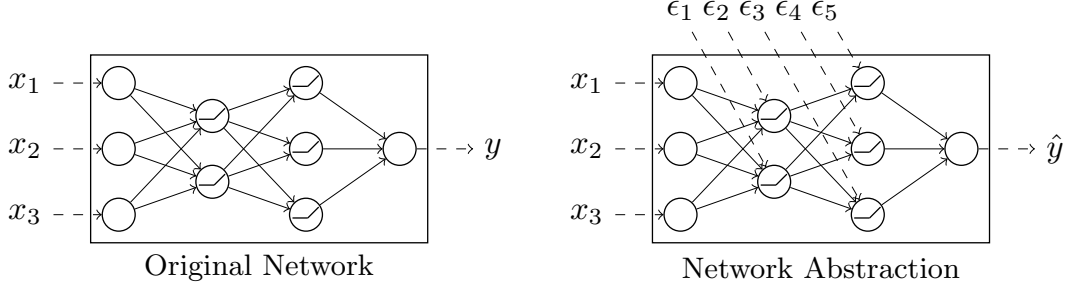


Figure 2: We abstract away the details of the polynomial approximation by introducing an additional error input ϵ_i to the output of each ReLU activation.

$z_i \in [\underline{z}_i, \bar{z}_i]$. Define $f(x, \epsilon)$ as the abstraction of $\hat{f}(x)$ if each ReLU layer h_ℓ of the abstraction takes the following form:

$$(\hat{h}_\ell(z))_i = \sigma(z_i) + \epsilon_i \quad (3)$$

where $\epsilon_i \in [-\delta_i, +\delta_i]$ and $z_i \in [\underline{z}_i, \bar{z}_i]$ for all i and $x \in \mathcal{X}$.

4.2 Potential Range Estimation

Note that estimating the range of the potentials z_i for each ReLU neuron i corresponds to the classic reachability problem in neural network verification [29, 25]. In our experiments, we use the FastLin bound propagation technique in [28] to compute linear bounds on the output of the affine layers g_ℓ . Differently from the standard definition therein, the bounds depend on both the original input vector x and the error vector ϵ :

$$\underline{S}_\ell x + \underline{U}_\ell \epsilon + \underline{v}_\ell \leq g_\ell(x, \epsilon) \leq \bar{S}_\ell x + \bar{U}_\ell \epsilon + \bar{v}_\ell \quad (4)$$

Note that the matrices \underline{U}_ℓ and \bar{U}_ℓ have non-zero entries only in the columns $j \in [1, \sum_{k=0}^{\ell-1} n_k]$, where n_k is the number of ReLU neurons in layer k . Furthermore, we derive the concrete range of the ReLU potentials $(g_\ell(x, \epsilon))_i \in [\underline{z}_i, \bar{z}_i]$ by maximising (resp. minimising) the right-hand side (resp. left-hand side) of Equation 4 over both $x \in \mathcal{X}$ and $\epsilon_i \in [-\delta_i, +\delta_i]$.

4.3 Local Lipschitz Constant Estimation

In order to compute the optimal polynomial approximation of each ReLU, we need both the potential range $[\underline{z}_i, \bar{z}_i]$ and the optimal approximation error ϵ_i (see Section 4.4). In this regard, we can inform our decision by estimating the impact of each ϵ_i on the output error. We achieve this goal by computing the local Lipschitz constant of ϵ_i :

Definition 6 (Lipschitz Constant). *The local Lipschitz constant of ϵ_i is:*

$$L_i^\infty = \max_{x, \epsilon, \epsilon'} \frac{|f(x, \epsilon) - f(x, \epsilon')|_\infty}{|\epsilon_i - \epsilon'_i|_\infty} \quad (5)$$

where $x \in \mathcal{X}$, $\epsilon_j, \epsilon'_j \in [-\delta_j, +\delta_j]$ for all j and $\epsilon_j = \epsilon'_j$ for all $j \neq i$.

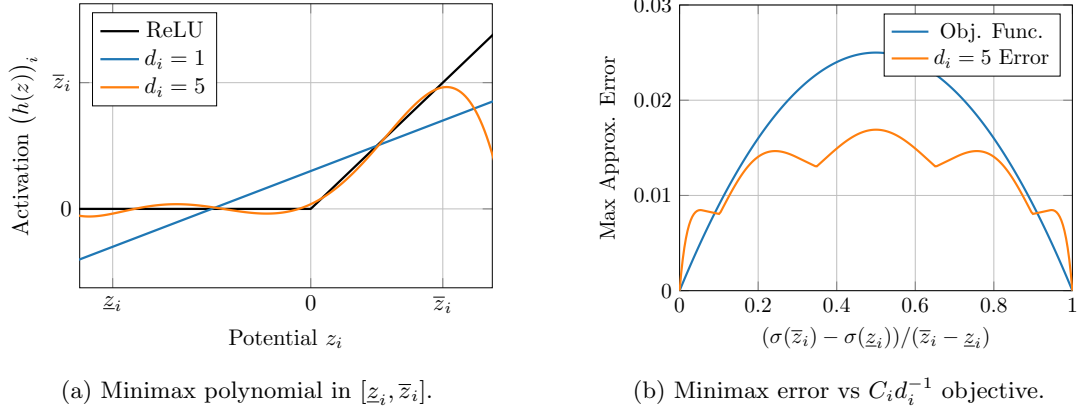


Figure 3: Approximating a ReLU activation with polynomials (3a) and the resulting error (3b). Note that the example on the left has $(\sigma(\bar{z}_i) - \sigma(z_i))/(\bar{z}_i - z_i) = 0.4$.

In other terms, Definition 6 quantifies the maximum rate of change in the output as we modify one single error input ϵ_i . Thanks to it, we can bound the total approximation error as follows:

$$\|f(x) - \hat{f}(x)\|_\infty \leq \max_\epsilon \{ \|f(x) - f(x, \epsilon)\|_\infty \} \leq \sum_i L_i^\infty \max_{\epsilon_i} \{ |\epsilon_i| \} \quad (6)$$

Computing the Lipschitz constant of a neural network is highly non-trivial, and it is currently an area of active research [24, 10]. In our experiments, we adapt the algorithm in [24] to the architecture of our abstracted neural network, by extracting the Jacobian of ϵ from the intermediate results of the backward pass.

4.4 Optimal Error Allocation

Recall that our objective is approximating the neural network $f(x)$ in such a way that the output error never exceed a given margin δ_y (see Definition 4). With the Lipschitz estimates from Section 4.3, we can decide how to optimally distribute the local approximation error ϵ_i as follows:

$$\text{Minimise } d_{tot} = \sum_i d_i(\epsilon_i, z_i, \bar{z}_i) \quad (7)$$

$$\text{Subject to } \sum_i L_i^\infty \epsilon_i \leq \delta_y \quad (8)$$

$$\text{And } 0 \leq \epsilon_i \leq \delta_i, \quad \forall i \quad (9)$$

where ϵ_i is a shorthand for $\max_{\epsilon_i} \{ |\epsilon_i| \}$ from Equation 6, i.e. the maximum approximation error of each ReLU function in $[z_i, \bar{z}_i]$.

Crucially, the objective function estimates the inference cost of replacing each activation i with a polynomial of degree d_i . We give an example of such polynomial in Figure 3a. Clearly, the degree is a monotonically increasing function in the size of the potential range $[z_i, \bar{z}_i]$, and monotonically decreasing in the required precision ϵ_i . Furthermore, we know that the minimax approximation $p_{d_i}^*(z_i)$ of degree d_i of any continuous non-differentiable function $\sigma(z_i)$

Algorithm 1 LiGAR

Input: network abstraction $f(x, \epsilon)$, finite input set \mathcal{X} , input domain δ_x , error domain δ_ϵ , output precision δ_y , convergence threshold h .

Output: optimal error allocation ϵ^* , potential ranges \underline{z}, \bar{z} .

```

1:  $t \leftarrow 0$ 
2:  $\delta_i(1) \leftarrow \delta_\epsilon, \forall i$  ▷ Largest error domain
3: repeat
4:    $t \leftarrow t + 1$ 
5:   for all  $x \in \mathcal{X}$  do
6:      $\underline{z}(x), \bar{z}(x) \leftarrow \text{ComputePotentials}(f, x, \delta_x, \delta_i(t))$  ▷ Forward pass
7:      $L^\infty(x) \leftarrow \text{ComputeLipschitz}(f, \underline{z}, \bar{z})$  ▷ Backward pass
8:   end for
9:    $C_i \leftarrow C(\min_x \{z_i(x)\}, \max_x \{\bar{z}_i(x)\}), \forall i$  ▷ Cost coefficients
10:   $L_i^\infty \leftarrow \max_x \{L_i^\infty(x)\}, \forall i$  ▷ Lipschitz estimates
11:   $\epsilon^*(t), v(t) \leftarrow \text{OptimalAllocation}(C, L^\infty, \delta_i(t), \delta_y)$  ▷ Solve dual problem
12:   $\delta_i \leftarrow \frac{1}{2}(\max_t \{\epsilon_i^*(t)\} + \min_{t: \epsilon_i^*(t) < \delta_i(t)} \{\epsilon_i^{max}(t)\}), \forall i$  ▷ Tighter error domain
13: until  $v(t-1) - v(t) \leq h$  ▷ Stop if no progress
14:  $t^* \leftarrow \arg \min_t \{v(t)\}$ 
15: return  $\epsilon^*(t^*), v(t^*)$  ▷ Return best allocation
    
```

guarantees a maximum approximation error $\epsilon_i \propto O(d^{-1})$ [26]. As a result, we can compute the exact relation between ϵ_i and $d_i = 1$ and extend it to any $d_i > 1$:

$$\epsilon_i \approx \frac{C_i}{d_i} \quad \text{where} \quad C_i = \frac{1}{2} \left(\sigma(\underline{z}_i) - \underline{z}_i \frac{\sigma(\bar{z}_i) - \sigma(\underline{z}_i)}{\bar{z}_i - \underline{z}_i} \right) \quad (10)$$

which is exact for $d_i = 1$ and overestimated for $d_i > 1$, unless the constant C_i is very close to 0 or 1 (see example for $d_i = 5$ in Figure 3b). We show how to derive the corresponding optimal solution in Appendix A.

4.5 Abstraction Refinement Cycle

Now, we can discuss the key idea behind our LiGAR method. When we compute the Lipschitz constants $L - i^\infty$ and potential ranges $[\underline{z}_i, \bar{z}_i]$, we need to specify a domain $\epsilon_i \in [-\delta_i, +\delta_i]$ for the approximation error. If we choose a domain that is too large, our estimates will be conservative. If we choose it too small, the values of ϵ_i will also be small, thus yielding a larger cost C_i/ϵ_i .

We solve this issue by making our abstraction-optimisation process iterative (see Algorithm 1). More specifically, we start with a large user-defined error domain (Line 2). Then, we estimate the potential ranges and Lipschitz estimates accordingly (Lines 5-10). Next, we solve the dual optimisation problem (Line 11), which gives us the optimal error allocation *given the current estimates*. With this result, we can tighten the error domain (Line 12) and repeat the process.

Note that the error domain is always kept larger than any of the solutions found so far, and smaller than all error domains that did not activate the constraint $\epsilon_i^* \leq \delta_i$. In our experiments, LiGAR converges after 20-30 iterations. Finally, we use an indirect definition for the input domain in Lines 5-10. That is, akin to many neural network verification settings [19], we select a finite set of concrete inputs \mathcal{X} , and consider a norm- ∞ ball of size δ_x around each of them.

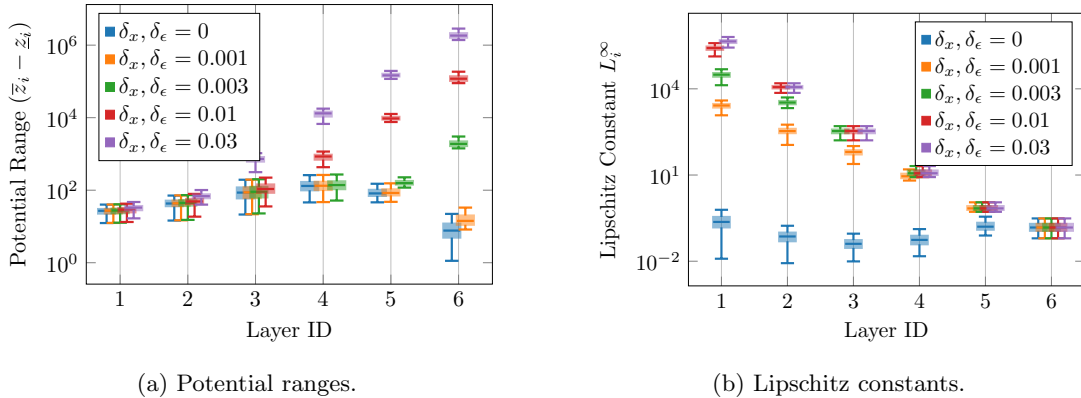


Figure 4: Potential ranges (Fig. 4a) and Lipschitz constants (Fig. 4b) on the MNIST (fc) network. The values become more conservative as δ_x and δ_ϵ increase.

5 Empirical Analysis

Here, we present our preliminary experimental results. Our main goal is twofold: demonstrate the behaviour of LiGAR on medium-size networks; gather evidence on the practical utility of worst-case approximations for private inference.

5.1 Experimental Setting

Neural Network Model. We select the MNIST (fc) image classification network from VNN-COMP’21.¹ This network has six hidden layer containing 256 neurons each. The inputs are grey-scaled images with 784 pixels, whose values are scaled to the interval $[0, 1]$. The outputs are ten scores that represent the log-likelihood of the image containing a hand-written digit between zero and nine [11].

LiGAR Parameters. For reasons of time, we run LiGAR with \mathcal{X} containing only 1000 random samples from the MNIST training set [11]. We set the initial error domain to $\delta_\epsilon = 0.1$, the convergence threshold to $h = 0.01$ and the output precision to $\delta_y = 0.01$, unless otherwise specified.

5.2 Experimental Results

Potential Ranges and Lipschitz Constants. Computing worst-case guarantees on the potential ranges and Lipschitz constants incur a cost regarding over-estimation. In Figure 4, we show what happens when we increase the input δ_x and error δ_ϵ domains. As expected, the bound propagation yields fairly tight estimates for the potential range $\bar{z}_i - \underline{z}_i$ over the first layers of the network until the over-approximation caused by the FastLin abstraction technique begins to accumulate. Similarly, Lipschitz constants are tight for the final few layers but get increasingly worse as the backward estimation procedure reaches earlier layers. Note that LiGAR can iteratively reduce the over-estimation caused by δ_ϵ , but not the one caused by δ_x .

¹<https://github.com/stanleybak/vnncomp2021/tree/main/benchmarks>

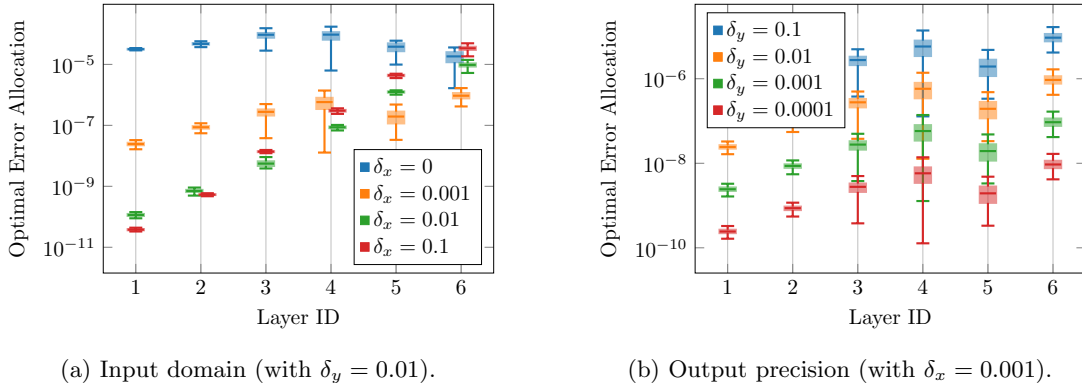


Figure 5: Optimal error allocation for different input domain sizes (Fig. 5a) and output precision requirements (Fig. 5b) on the MNIST (fc) network. ReLU neurons have been omitted in their linear state (always active or inactive).

LiGAR Error Allocation. We report the optimal error allocation computed by LiGAR in Figure 5. Note how increasing the input domain δ_x forces LiGAR to be more conservative in its estimates. Interestingly, this conservative behaviour is more pronounced for earlier layers. We believe this is caused by the corresponding over-estimation of the Lipschitz constants (see Figure 4b). In contrast, the error allocation is proportional to the output precision δ_y . This is not surprising, since δ_y appears only in one constraint of the optimization problem (see Equation 8).

Linearly Removable ReLU Activations. It has been claimed in the literature that many ReLU activations in a given network can be replaced by a linear function [5]. We test this claim by checking how many ReLUs are either always active or always inactive according to LiGAR’s potential range estimates. The results in Table 1 show that the number of ReLUs we can remove quickly falls to zero as we cover a larger portion of the input space by increasing δ_x .

MNIST (fc)	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Layer 6
$\delta_x = 0$	23	11	53	81	238	162
$\delta_x = 0.001$	23	11	53	76	161	0
$\delta_x = 0.01$	17	7	0	0	0	0
$\delta_x = 0.1$	0	0	0	0	0	0

Table 1: Number of ReLUs in the always-active or always-inactive state.

Inference-Time Output Error. LiGAR guarantees that the output approximation error is never greater than δ_y . Here, we estimate the practical performance of the approximated network by measuring both the maximum observed output error and the change in predictive accuracy. Note that due to numerical instability in existing univariate polynomial approximation libraries (see Section 3.2) we report only the result for the abstracted network. In Table 2 we show results for both the 1000-sample LiGAR train set \mathcal{X} (see Section 5.1), and the full 10000-sample MNIST test set. Note how the accuracy matches the original network, and the maximum error is at least

two orders of magnitude smaller than the LiGAR guarantee $\delta_y = 0.01$. This shows that using worst-case estimates in designing neural network approximations leads to very conservative solutions.

	LiGAR Train Set		Test Set	
	$ f(x) - f(x, \epsilon) $	Accuracy	$ f(x) - f(x, \epsilon) $	Accuracy
Original $f(x)$	-	0.980	-	0.969
$\delta_x = 0$	2.64×10^{-4}	0.980	7.26×10^{-4}	0.969
$\delta_x = 0.001$	4.90×10^{-6}	0.980	5.67×10^{-6}	0.969
$\delta_x = 0.01$	3.36×10^{-5}	0.980	4.00×10^{-5}	0.969
$\delta_x = 0.1$	1.26×10^{-4}	0.980	1.30×10^{-4}	0.969

Table 2: Output error and inference accuracy of the abstracted network.

6 Conclusions

In this paper, we propose the Lipschitz-Guided Abstraction Refinement (LiGAR) method, which can compute the optimal neural network approximation with a given worst-case output error. We show that LiGAR converges quickly on medium-sized networks and could in principle scale to larger ones. At the same time, the worst-case guarantees come at the price of conservative solutions. We believe that this is a major obstacle for the adoption of worst-case methods like LiGAR for the purpose of private inference.

References

- [1] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006.
- [2] Zvika Brakerski and Vinod Vaikuntanathan. Efficient fully homomorphic encryption from (standard) LWE. In Rafail Ostrovsky, editor, *IEEE 52nd Annual Symposium on Foundations of Computer Science, FOCS 2011, Palm Springs, CA, USA, October 22-25, 2011*, pages 97–106. IEEE Computer Society, 2011.
- [3] Zvika Brakerski and Vinod Vaikuntanathan. Fully homomorphic encryption from ring-lwe and security for key dependent messages. In Phillip Rogaway, editor, *Advances in Cryptology - CRYPTO 2011 - 31st Annual Cryptology Conference, Santa Barbara, CA, USA, August 14-18, 2011. Proceedings*, volume 6841 of *Lecture Notes in Computer Science*, pages 505–524. Springer, 2011.
- [4] Hervé Chabanne, Amaury de Wargny, Jonathan Milgram, Constance Morel, and Emmanuel Prouff. Privacy-preserving classification on deep neural network. Cryptology ePrint Archive, Paper 2017/035, 2017. <https://eprint.iacr.org/2017/035>.
- [5] Minsu Cho, Ameya Joshi, Brandon Reagen, Siddharth Garg, and Chinmay Hegde. Selective network linearization for efficient private inference. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 3947–3961. PMLR, 17–23 Jul 2022.
- [6] Karthik Garimella, Nandan Kumar Jha, and Brandon Reagen. Sisyphus: A cautionary tale of using low-degree polynomial activations in privacy-preserving deep learning, 2021.
- [7] Craig Gentry. Fully homomorphic encryption using ideal lattices. In Michael Mitzenmacher, editor, *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009, Bethesda, MD, USA, May 31 - June 2, 2009*, pages 169–178. ACM, 2009.

- [8] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin Lauter, Michael Naehrig, and John Wernsing. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 201–210, New York, New York, USA, 20–22 Jun 2016. PMLR.
- [9] Nandan Kumar Jha, Zahra Ghodsi, Siddharth Garg, and Brandon Reagen. Deepreduce: Relu reduction for fast private inference. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4839–4849. PMLR, 18–24 Jul 2021.
- [10] Jacob Laurel, Rem Yang, Shubham Ugare, Robert Nagel, Gagandeep Singh, and Sasa Misailovic. A general construction for abstract interpretation of higher-order automatic differentiation. *Proc. ACM Program. Lang.*, 6(OOPSLA2), oct 2022.
- [11] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [12] Eunsang Lee, Joon-Woo Lee, Jong-Seon No, and Young-Sik Kim. Minimax approximation of sign function by composite polynomial for homomorphic comparison. *IEEE Transactions on Dependable and Secure Computing*, 19(6):3711–3727, 2022.
- [13] Joon-Woo Lee, Hyungchul Kang, Yongwoo Lee, Woosuk Choi, Jieun Eom, Maxim Deryabin, Eunsang Lee, Junghyun Lee, Donghoon Yoo, Young-Sik Kim, and Jong-Seon No. Privacy-preserving machine learning with fully homomorphic encryption for deep neural network. *IEEE Access*, 10:30039–30054, 2022.
- [14] Junghyun Lee, Eunsang Lee, Joon-Woo Lee, Yongjune Kim, Young-Sik Kim, and Jong-Seon No. Precise approximation of convolutional neural networks for homomorphically encrypted data, 2021.
- [15] Wei-Kai Lin, Ethan Mook, and Daniel Wicks. Doubly efficient private information retrieval and fully homomorphic ram computation from ring lwe. *Cryptology ePrint Archive*, 2022.
- [16] Qian Lou, Yilin Shen, Hongxia Jin, and Lei Jiang. {SAFEN}et: A secure, accurate and fast neural network inference. In *International Conference on Learning Representations*, 2021.
- [17] Edoardo Manino, Bernardo Magri, Mustafa A. Mustafa, and Lucas Cordeiro. Certified Private Inference on Neural Networks via Lipschitz-Guided Abstraction Refinement, May 2023.
- [18] Pratyush Mishra, Ryan Lehmkuhl, Akshayaram Srinivasan, Wenting Zheng, and Raluca Ada Popa. Delphi: A cryptographic inference service for neural networks. In *29th USENIX Security Symposium (USENIX Security 20)*, pages 2505–2522. USENIX Association, August 2020.
- [19] Mark Niklas Müller, Christopher Brix, Stanley Bak, Changliu Liu, and Taylor T Johnson. The third international verification of neural networks competition (vnn-comp 2022): Summary and results, 2022.
- [20] Ricardo Pachón and Lloyd N. Trefethen. Barycentric-remez algorithms for best polynomial approximation in the chebfun system. *BIT Numerical Mathematics*, 49(4):721–741, 2009.
- [21] Pascal Paillier. Public-key cryptosystems based on composite degree residuosity classes. In Jacques Stern, editor, *Advances in Cryptology — EUROCRYPT ’99*, pages 223–238. Springer Berlin Heidelberg, 1999.
- [22] Bernardo Pulido-Gaytan, Andrei Tchernykh, Jorge M. Cortés-Mendoza, Mikhail Babenko, Gleb Radchenko, Arutyun Avetisyan, and Alexander Yu Drozdov. Privacy-preserving neural networks with homomorphic encryption: Challenges and opportunities. *Peer-to-Peer Networking and Applications*, 14(3):1666–1691, 2021.
- [23] Ronald L. Rivest, Adi Shamir, and Leonard M. Adleman. A method for obtaining digital signatures and public-key cryptosystems. *Commun. ACM*, 21(2):120–126, 1978.
- [24] Zhouxing Shi, Yihan Wang, Huan Zhang, J. Zico Kolter, and Cho-Jui Hsieh. Efficiently computing local lipschitz constants of neural networks via bound propagation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 2350–2364. Curran Associates, Inc., 2022.

- [25] Hoang-Dung Tran, Diago Manzanas Lopez, Patrick Musau, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang, and Taylor T. Johnson. Star-based reachability analysis of deep neural networks. In Maurice H. ter Beek, Annabelle McIver, and José N. Oliveira, editors, *Formal Methods – The Next 30 Years*, pages 670–686, Cham, 2019. Springer International Publishing.
- [26] Lloyd N Trefethen. *Approximation Theory and Approximation Practice, Extended Edition*. SIAM, 2019.
- [27] Joost Verbraeken, Matthijs Wolting, Jonathan Katzy, Jeroen Kloppenburg, Tim Verbelen, and Jan S. Rellermeyer. A survey on distributed machine learning. *ACM Comput. Surv.*, 53(2), mar 2020.
- [28] Lily Weng, Huan Zhang, Hongge Chen, Zhao Song, Cho-Jui Hsieh, Luca Daniel, Duane Boning, and Inderjit Dhillon. Towards fast computation of certified robustness for ReLU networks. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5276–5285. PMLR, 10–15 Jul 2018.
- [29] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural network robustness certification with general activation functions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

A Closed-Form Dual Solution

The objective function $d_{tot} \approx \sum_i C_i/\epsilon_i$ is convex, and it is possible to derive its global minimum in closed form as follows. First, let us write the Lagrangian of our optimization problem:

$$\mathcal{L}(\epsilon, \mu, \lambda) = \sum_i \frac{C_i}{\epsilon_i} + \sum_i \mu_i(\epsilon_i - \delta_i) + \lambda(\sum_i L_i^\infty \epsilon_i - \delta_y) \tag{11}$$

where we omit the constraints $\epsilon_i \geq 0$ for simplicity. Now, we can derive the minimum of the Lagrangian over ϵ by equating its derivative to zero:

$$\epsilon_i^* = \arg \min_{\epsilon} \{\mathcal{L}(\epsilon, \mu, \lambda)\} = \sqrt{\frac{C_i}{\mu_i + \lambda L_i^\infty}} \tag{12}$$

for all i , which already satisfies the implicit constraint $\epsilon_i \geq 0$. Finally, the values of the dual variables $\mu_i \geq 0$ and $\lambda \geq 0$ can be found by maximizing ϵ_i^* . This operations requires satisfying the following system of equations:

$$\sum_i \sqrt{\frac{C_i}{\mu_i + \lambda L_i^\infty}} L_i^\infty = \delta_y \tag{13}$$

$$\mu_i = \max\left(\frac{C_i}{(\epsilon_i^*)^2} - \lambda L_i^\infty, 0\right) \tag{14}$$

which can be done iteratively by first solving Equation 13 with a binary search, and then computing Equation 14. In our experience, very few iterations are required until convergence.