

NeurObfuscator: A Full-stack Obfuscation Tool to Mitigate Neural Architecture Stealing

Jingtao Li*, Zhezhi He[†], Adnan Siraj Rakin*, Deliang Fan*, Chaitali Chakrabarti*

*School of Electrical Computer and Energy Engineering, Arizona State University, Tempe, AZ, 85287

[†]Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai

*{jingtao1, asrakin, dfan, chaitali}@asu.edu; [†]{zhezhi.he}@sjtu.edu.cn

Abstract—Neural network stealing attacks have posed grave threats to neural network model deployment. Such attacks can be launched by extracting neural architecture information, such as layer sequence and dimension parameters, through leaky side-channels. To mitigate such attacks, we propose *NeurObfuscator*, a full-stack obfuscation tool to obfuscate the neural network architecture while preserving its functionality with very limited performance overhead. At the heart of this tool is a set of obfuscating knobs, including layer branching, layer widening, selective fusion and schedule pruning, that increase the number of operators, reduce/increase the latency, and number of cache and DRAM accesses. A genetic algorithm-based approach is adopted to orchestrate the combination of obfuscating knobs to achieve the best obfuscating effect on the layer sequence and dimension parameters so that the architecture information cannot be successfully extracted. Results on sequence obfuscation show that the proposed tool obfuscates a ResNet-18 ImageNet model to a totally different architecture (with 44 layer difference) without affecting its functionality with only 2% overall latency overhead. For dimension obfuscation, we demonstrate that an example convolution layer with 64 input and 128 output channels can be obfuscated to generate a layer with 207 input and 93 output channels with only a 2% latency overhead.

Index Terms—Neural Network, Side-channel attack, Architecture Stealing, Obfuscation

I. INTRODUCTION

The architecture information of a Deep Neural Network (DNN) model is very sensitive and should never be exposed. It is a valuable Intellectual Property (IP) that costs companies lots of time and resources. Knowledge of the exact architecture allows an adversary to build a more precise substitute model and such a model can be used to launch devastating adversarial attacks. For instance, it is shown in [1] that accurate architecture information enables the adversary to improve the attack success rate of input adversarial attack by almost 3 times.

Side-channel based DNN architecture stealing has been reported in several prior works [1], [2]. Even without access to the service, an outsider can extract the DNN architecture through side-channel information leakage, as shown in Fig. 1. Specifically, when the owner of the neural network IP hosts the application on a third-party cloud computing platform or on a local device with GPU support, it opens it up to architecture stealing through side-channel attacks [1]–[4]. A typical architecture stealing flow consists of profiling the target device, training sequence predictor (e.g., LSTM [5]), predicting layer sequence based on run-time trace of the target DNN model and

then extracting the dimension parameters of each layer. This is quite different from stealing through Machine-Learning-as-a-Service (MLaaS) [6], [7], where the attacker has access to the public prediction API and the confidence score of the labels.

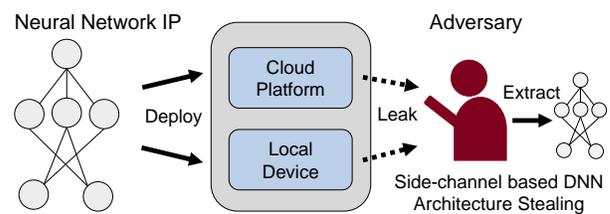


Fig. 1. Architecture stealing threatens intellectual property.

Previous efforts on preventing DNN architecture stealing have focused on hardware to eliminate information leakage. Oblivious Random Access Memory (ORAM) technology [8], [9] prevents memory access leakage by encrypting the memory address. [10] proposes re-design of the Miss Status Holding Registers (MSHR) to obfuscate GPU memory access to add a layer of randomness. Though hardware modifications are effective countermeasures, they are not beneficial to existing devices and have high performance overhead. Recently, [11] proposed a decision tree-based detection method against spy applications on GPU. However, it suffers from high false positive rate and is not practical. TVM [12] has also been proposed as a potential countermeasure. Nevertheless, as shown by our experiments (Fig. 6), standard TVM does not show enough randomness to be an effective countermeasure.

In this work, we propose *NeurObfuscator*, a full-stack tool which obfuscates neural network execution to effectively mitigate neural architecture stealing. Our obfuscating tool consists of 8 obfuscating knobs for two kinds of obfuscation, namely *sequence obfuscation* which obfuscates the layer depth and types and connection topologies between layers, and *dimension obfuscation* which obfuscates the dimension parameters of each layer, including the number of input and output channel, weight kernel size, etc. Function-preserving knobs such as layer branching, layer deepening, layer skipping followed by selective fusion at graph optimization step are used for sequence obfuscation, and layer widening, dummy addition, kernel widening and schedule modification in the back end are used for dimension obfuscation.

We use genetic algorithm to search for the best set of

obfuscation combinations for sequence and dimension obfuscation that achieve strong obfuscation for a given user-defined time budget. The obfuscation strength is measured by Layer Error Rate (LER) which represents normalized editing distance of extracted layer sequence given the ground-truth layer sequence in sequence obfuscation, and Dimension Error Rate (DER) which represents the normalized error of extracted dimension parameters in a layer in dimension obfuscation. Our contributions can be summarized as follows:

- This is the first work on mitigating the NN architecture stealing attack with pure-software obfuscations. We propose a total of 8 obfuscating knobs across the entire DNN execution stack to achieve sequence & dimension obfuscations and demonstrate the performance on state-of-the-art GPUs.
- We present an obfuscation tool backed by genetic algorithm to search for the best combination of obfuscations to obfuscate any neural network architecture with user-defined inference latency budget. Source code is available¹.
- For sequence obfuscation, our obfuscation tool can obfuscate a ResNet-18 architecture to have a 2.44 LER (which translates to a 44-layer editing distance [13]) against state-of-art LSTM-based sequence predictors with only 2% increase in overall latency.
- For dimension obfuscation, we show how a convolution layer with 64 input and 128 output channels can be obfuscated so that it is extracted as a layer with 207 input and 93 output channels with only 2% increase in layer-level latency.

II. BACKGROUND

A. Neural Network Notation

We summarize the neural network notation that is used throughout the paper, with focus on the most common Conv2D operator (represented in 4D by $k1, k2, c, j$) in Table I.

TABLE I
NEURAL NETWORK NOTATION

Notation	Definition
$\mathbf{X}^{(i)}$	Inputs of i -th layer
$\mathbf{W}^{(i)}, \mathbf{U}^{(i)}, \mathbf{V}^{(i)}$	Weights of i -th layer
$\varphi(\cdot, \cdot)$	Activation function
$k1, k2$	Conv2D kernel sizes
c, j	Input/Output channel sizes
h_i, w_i	Height/Width of inputs
h_o, w_o	Height/Width of outputs

B. NN Execution Flow

Generally, an NN architecture is a topology of neural network layers with non-linear functions. Fig. 2 demonstrates

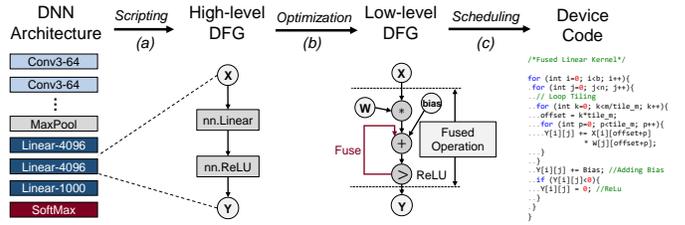


Fig. 2. DNN model execution flow. (a) Scripting (Coding) in Python on popular deep learning framework, (b) using TVM [12], a deep learning compiler for graph optimization via TVM Relay, and (c) Auto-TVM scheduling to generate device-level code.

a typical NN execution process, which consists of multiple steps. The first step is scripting (coding) of a DNN architecture using Python with popular frameworks such as Pytorch or Tensorflow. The scripting transforms the raw design into a high-level dataflow graph (aka. computational graph). Next, the high-level graph is ported to TVM for further optimization. One can also directly use TorchScript or Tensorflow XLA for graph optimization. For instance, in TVM, the graph optimization process is handled by Relay module, which provides handy options such as: 1) “FoldConstant()”, which evaluates expression involves only constants; 2) “EliminateCommonSubexpr()” which creates a shared variable for multiple expressions with same output to avoid the same expression being evaluated multiple times; and 3) “FuseOps()”, which fuses multiple expressions together. User can specify which optimizations to enable.

The last step is scheduling which optimizes the execution of operators on a given device. In TVM framework, a machine-learning based scheduling called “AutoTVM” [14] is used to generate optimized codes. For each operator in the optimized low-level graph, AutoTVM module uses Xgboost [15] to search for the best schedule within the predefined search space. Fig. 2 (c) shows a generic multi-level loop nest implementation of the linear operator. The search space for such a linear operator is defined by one single knob, $tile_m : [1, m]$, which determines the tiling parameter m for input X .

C. Architecture Stealing Attack Flow

Extracting the architecture sequence is not trivial. Since neural network execution goes through several steps of optimization as shown in Fig. 2, the intermediate steps bring in lots of variations in the final device code which directly affects the hardware trace. Prior works [1], [2] have adopted machine learning to extract the architecture from side-channel information. Both works successfully extract common architectures with very high accuracy. They share similar stealing attack methodologies as illustrated in Fig. 3 but differ in their prediction models.

The works in [1], [2] both use Long-Short-Term-Memory (LSTM) models to predict the layer sequence. First, massive profiling of randomly generated DNNs on the target devices is done offline. After proper labeling (an example is shown in Fig Fig. 9), the attacker acquires a trace-sequence dataset and uses it to train the LSTM model. At the time of the

¹Source Code: <https://github.com/zlijingtao/Neurofuscator>

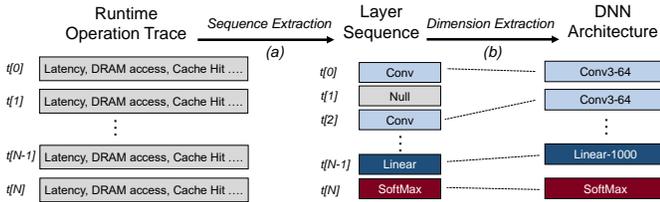


Fig. 3. Architecture stealing flow. (a): Layer Sequence Extraction (b): Layer Dimension Extraction.

attack, the attacker uses the LSTM predictor to perform the layer sequence extraction on the run-time trace of the target DNN, which is time-series data consisting of multiple features as shown in Fig. 3. The sequence prediction locates the *layer operator* in the run-time trace and classifies it by layer type.

Then, dimension extraction is done for each identified layer operator once its *time-step* (position) and *class* (layer type) is known. This is considered to be simpler than sequence extraction. Note that, dimension extraction can be done either manually [1] or automatically [2].

In summary, existing architecture stealing attacks heavily rely on the run-time trace, and so to mitigate such stealing attacks, our obfuscating tool changes the run-time trace as much as possible.

III. THREAT MODEL

We only consider architecture stealing on applications running on common GPU devices. For other devices such as FPGAs, CPUs and ASICs, we believe that the obfuscation methods proposed in this paper can also be used. We consider NN applications running in both remote and local settings.

In remote setting, we assume that the owner runs the NN application on a third-party cloud computing platform and the attacker acts as a normal user without system privilege on the same machine. Specifically, the attacker can perform “driver downgrading attack” to access the profiling API and thus conduct GPU profiling on target neural network applications at run-time. This is similar to the threat model presented in [2].

In local setting, we assume that the device is off-the-shelf and the attacker can do profiling on an identical device to train a predictor model. While the target application is running, the attacker can get access to the run-time hardware traces of the target neural network applications through side-channel attacks [1], [3], [16].

Depending on the attack scenario and capability, we categorize the attacker w.r.t the extent of information leakage. Table II describes three cases (from weakest to strongest):

- *Case-A*: Timing side-channel. The attacker can get accurate operator latency information for each time step. For example, the cycle information for each issued operator is acquired in [3], [4]. This case naturally includes Electromagnetic (EM) side-channel [16], as the EM reflects the cycle information of each operator.
- *Case-B*: DRAM side-channel. The attacker has access to the DRAM read/write information of each operator, as well as the latency through PCIE side-channel [1], [17].

- *Case-C*: Cache side-channel. The attacker enables context-switching side-channel [2] or exploits the collocation [18] side-channel. By profiling spy applications, the attacker samples the cache performance counters of the target applications and uses it to extract cache performance, DRAM transactions, latency of the target kernels, etc. The additional cache performance counters in case C include L1 cache and L2 cache utilization, hit rate and read and write data volumes.

In all cases, the attacker does massive profiling of DNN model’s run-time trace to steal the architecture.

TABLE II
DIFFERENT CASES OF INFORMATION LEAKAGE

	Latency	DRAM-Access	Cache-Counters
Case-A	✓	✗	✗
Case-B	✓	✓	✗
Case-C	✓	✓	✓

IV. TRACE OBFUSCATION

Architecture stealing is possible because neural network execution process is deterministic, as described in Fig. 2. To provide countermeasure against architecture stealing, we propose six obfuscating knobs in scripting: layer widening, layer branching, dummy addition, layer deepening, layer skipping and kernel widening. Next, we propose selective fusion under graph optimization and schedule modification in the backend.

A. Obfuscation in Scripting

We realize that many of the function-preserving transformations that have been successfully used in evolution NAS [19] can be used in obfuscation. More specifically, we use layer widening, layer branching, layer deepening, layer skipping and kernel widening and dummy addition obfuscating knobs in this phase. Note that while many of these operators have been introduced before in the context of architecture evolution [20], [21], we are the first to use them in terms of side-channel countermeasures. Layer branching is redesigned, and dummy addition knob is added for dimension obfuscation.

1) **Layer Widening**: *Layer widening* increases output channel j of a Conv2D layer or a linear layer. Basically, the weights of the added output channels are duplicates of the weights of existing output channels. We allow the widening operator to take fractional numbers. For example, if the weight $\mathbf{W}_{k_1, k_2, c, j}^{(i)}$ takes a widening factor of $0.25 \times$ and results in $U_{k_1, k_2, c, 1.25j}^{(i)}$, then the first $0.5j$ of the output channels come from the duplication of the first $0.25j$ output channel of original $\mathbf{W}_{k_1, k_2, c, j}^{(i)}$.

To preserve the functionality, next layer’s weights need to be adjusted accordingly. In this example, the next layer $\mathbf{W}_{k_1, k_2, j, m}^{(i+1)}$ must increase its input channel size accordingly, resulting in $U_{k_1, k_2, 1.25j, m}^{(i+1)}$ to match the increased output

channels. The dimension parameters of $U^{(i+1)}$ for the first $0.5j$ input channels have to be adjusted for the duplicated input channels.

Purpose: Layer widening increases memory accesses for the current and the next layer by around $(N - 1)$ times for widening factor N . This results in increased number of input/output channels and affects dimension extraction.

2) **Layer Branching:** *Layer branching* breaks a single NN layer operator into smaller ones. For example, a Conv2D operator $W_{k_1, k_2, c, j}^{(i)}$ is branched into two parts (output-wise branching): $U_{k_1, k_2, c, j/2}^{(i)}$ and $V_{k_1, k_2, c, j/2}^{(i)}$ and the final output is the concatenation of the two partial convolutions:

$$\text{Concate}(U_{k_1, k_2, c, j/2}^{(i)} * X^{(i)}, V_{k_1, k_2, c, j/2}^{(i)} * X^{(i)}) \quad (1)$$

While the version in [19] only considers branching in the output channel dimension of Conv2D/linear layers, we also consider layer branching in the input channel dimension. A Conv2D layer of weight $W_{k_1, k_2, c, j}^{(i)}$ is branched into two (input-wise branching): $U_{k_1, k_2, c/2, j}^{(i)}$ and $V_{k_1, k_2, c/2, j}^{(i)}$, and the final result is the addition of the two:

$$\text{Add}(U_{k_1, k_2, c/2, j}^{(i)} * X^{(i)}, V_{k_1, k_2, c/2, j}^{(i)} * X^{(i)}) \quad (2)$$

Here, the activation input needs to be sliced into two as well to match the halved input channel dimension of two smaller convolutions. Various branching methods are feasible, for example, one can also separate it into more than two parts or even do unbalanced branching. Here we consider balanced branching into two or four parts, for both input-wise and output-wise branching.

Purpose: Layer branching increases the number of layer operators and changes the data volume that needs to be accessed for each operator. For input-wise branching, the input activation and weight volume are halved for each small kernel, and for output-wise branching, input activation is the same but weight and output activation volumes are halved. This knob can be used for both sequence and dimension obfuscation.

3) **Dummy Addition:** *Dummy addition* is simply adding zero to the activation results. We create a zero matrix of the same shape as the activation output X of current layer.

$$D_{b, j, h_o, w_o}^{(i)} = O_{b, j, h_o, w_o} \quad (3)$$

A dummy addition factor of N means that we create and add the dummy matrix to the output repeatedly N times.

Purpose: Addition operators are “fused” into previous layer operators in the fusion step in graph optimization (refer to Fig. 2) and so the extra cache accesses from the addition operator get added to the layer computation and affect the dimension extraction of that layer.

4) **Layer Deepening:** *Layer deepening* inserts an extra computational layer at the end of current layer’s activation function. The insertion of a deepening layer $U^{(i)}$ does not change the original result.

$$\varphi(U^{(i)} * \varphi(W^{(i)} * X^{(i)})) = \varphi(W^{(i)} * X^{(i)}) \quad (4)$$

For linear layers, the deepening layer $U^{(i)}$ is simply an identity matrix of the same size as its input. For Conv2D layer, layer $U^{(i)}$ of size (k_1, k_2, j, j) need to be initialized as:

$$U_{a, b, d, m}^{(i)} = \begin{cases} 0 & a = \frac{k_1-1}{2} \wedge b = \frac{k_2-1}{2} \wedge d = m \\ 1 & \text{otherwise} \end{cases} \quad (5)$$

We favor a kernel size of $k_1 = k_2 = 1$ which avoids too much extra computation. Notice that the correctness of Eq. (4) also depends on whether the activation function $\varphi(\cdot)$ results stay the same when it gets stacked $\varphi(\cdot) = \varphi(\varphi(\cdot))$. Fortunately, the most popular ReLU activation subscribes to this property. The same property does not hold for batch normalization, so the deepening layer must be added before batch normalization, as shown in Fig. 4 (a).

Purpose: Add an extra computational layer to the layer extraction result. This can be used for sequence obfuscation.

5) **Layer Skipping:** *Layer skipping* inserts an extra computational layer as illustrated in Fig. 4 (b). The additional layer, referred to as *skipping layer*, operates on the activation output of an existing layer and adds it to the original activation output. The skipping layer is initialized to zero and thus always have a zero output matrix.

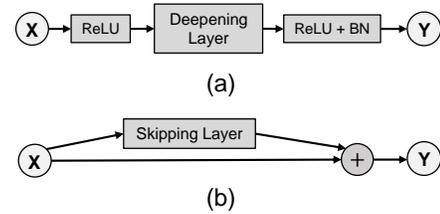


Fig. 4. (a) Illustration of Layer Deepening. (b) Illustration of Layer Skipping.

For an activation of size (b, j, h_o, w_o) , the skipping layer can be a Conv2D layer $U^{(i)}$ that has a shape of (k_1, k_2, j, j) and all entries are zero. The output of the skipping layer is:

$$U^{(i)} * X_{b, j, h_o, w_o} + X_{b, j, h_o, w_o} = X_{b, j, h_o, w_o} \quad (6)$$

Purpose: Add an extra computational layer to the layer extraction result. This can be used for sequence obfuscation.

6) **Kernel Widening:** *Kernel widening* increases the kernel size of a Conv2D layer. It is done by padding zeros to both the input and convolution kernels. A kernel widening of “+1” to a Conv2D layer of shape (k_1, k_2, c, j) will result in new weight of shape $(k_1 + 2, k_1 + 2, c, j)$ and input of shape $(b, c, h_i + 2, w_i + 2)$. We find this to be useful in particular for Conv2D layers that have kernel size of 1×1 . These small 1×1 kernels would then transform to 3×3 kernels after widening.

Purpose: Change kernel size of the Conv2D operator, resulting in a completely different trace. This affects dimension extraction.

B. Obfuscation in Graph Optimization

Fusion is an important graph optimization technique in the TVM Relay module. It fuses subsequent injective operators (scaling or addition) in complex layer operators, such as

Conv2D, linear and max-pooling, and transforms the shape of the inputs completely. Fusion ensures execution efficiency as it improves the data reuse and avoids context switching overhead. As shown in Fig. 5, the fused operator is significantly faster than sum of the separate operators.



Fig. 5. Fusion saves time by reducing kernel switching overhead. Top: *winograd_conv2d_kernel2* is fused with add and ReLU kernel. Bottom: *winograd_conv2d_kernel2* is issued separately from add and ReLU kernels, resulting in significant increase in execution time.

7) **Selective Fusion:** *Selective fusion* is a controllable version of the generic fusion. While the generic fusion fuses successive injective operators greedily, the selective fusion allows N successive operators to fuse and forbids more operators to fuse. For example, by setting N to zero for a Conv2D operator shown in Fig. 5, the Conv2D operator will be issued separately, as shown in the lower part of the figure.

Purpose: Increase the number of operators. Setting N to a small value decreases the memory access and latency of a layer operator, and affects both sequence and dimension extraction.

C. Obfuscation in Scheduling

In the backend, AutoTVM handles the compilation and generates optimized code for a given device. It provides options such as the number of trials for tuning, etc. We investigated whether these options can be used to generate randomness in the final result and thereby help in obfuscation. We tried 3 rounds with different number of trials using the default Xgboost (XGB) tuner in AutoTVM for a Conv2D operator. All these trials generated the same schedule, which is understandable because the tuning is designed to optimize latency. The profiling results in Fig. 6 show that the cycle, DRAM read and L1 cache utilization are very similar for different number of tuning options (XGB-200, 400 and 800 denoted by bars 1-3), meaning AutoTVM derived schedule is deterministic and cannot be directly used for obfuscation.

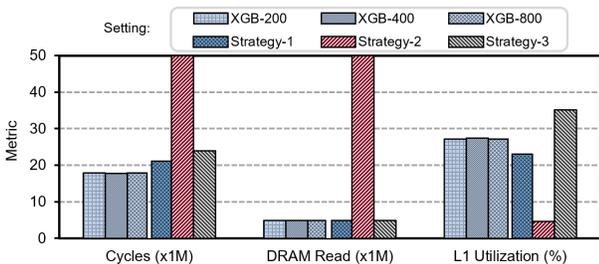


Fig. 6. Profiling results of AutoTVM derived schedules using Xgboost tuner with different number of trials (bars 1-3) and results using the proposed strategies (bars 4-6).

8) **Schedule Modification:** To generate schedules via AutoTVM with different outcomes, the search space has to be

modified. Actually changing the search space requires a time-consuming tuning (searching) process each time, and so we propose a simple approach that directly modifies the derived schedules with a small sacrifice on operator’s performance. For example, the schedule derived by Xgboost in Fig. 6 is $[-1, 4, 8, 4]$ for “Tile-Y” and $[-1, 2, 4, 2]$ for “Tile-X”. The first dimension of the tiling is for mini-batch so it is fixed as -1. We present a modification strategy by forcing each of the other dimensions to be 1. For example, the schedule for “Strategy-1” forces the second dimension to be 1, which produces $[-1, 1, 8, 16]$ for “Tile-Y” and $[-1, 1, 4, 4]$ for “Tile-X”. We set the values of other two dimensions by keeping the product be the same as original ($8 \times 16 = 4 \times 8 \times 4$) and let them be as close as possible. We derived three modified schedules using this method. Their profiling result is shown in bars 4-6 in Fig. 6. For latency and L1 cache utilization, all three schedules show noticeable difference. Strategy-2 is an example of bad modification, where the DRAM read and number of cycles explodes and L1 cache utilization is very poor. Strategy-1, on the other hand, helps achieve obfuscation without hurting the performance too much.

Purpose: Derive different schedules for the same operator that present differences in latency, DRAM access and cache performance. This affects dimension extraction.

V. NEUROBFUSCATOR TOOL FLOW

The NeuroObfuscator tool flow consists of two key steps: 1) sequence obfuscation which obfuscates the layer sequence including layer type and topology, and 2) dimension obfuscation which obfuscates the dimensions of individual layer operators. We summarize the role of each obfuscating knob in Fig. 7. Knobs with star superscripts affect both sequence and dimension extraction.

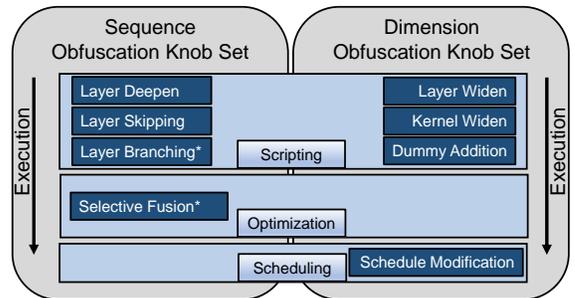


Fig. 7. Obfuscating knobs in NeuroObfuscator are separated into two sets.

Since the purpose of sequence and dimension obfuscation is orthogonal, we investigate them independently. To reduce search time, we take one more step in limiting the search space as follows.

Knob Partition. First, we partition the set of knobs, as shown in Fig. 7. We put selective fusion and layer branching together with layer deepening and layer skipping in the sequence obfuscation knob set. The remaining 4 knobs, namely, layer widening, kernel widening, dummy addition and schedule pruning are considered for dimension obfuscation. Among

i.e., a Turing GPU (GTX-1660) to profile models on CIFAR-10 dataset and an Ampere GPU (RTX-3090) to profile models on ImageNet dataset. We collect the number of cycles, DRAM and cache performance metrics for each issued operator of the model running in inference mode. After profiling, we collect three sets of features to match the three attack cases. The selected features in Nsight Compute profiling process are listed in Table III. In practice, the attacker gets noisy trace information through side-channels [1], [2]. To study the worst-case (i.e., strongest attack), we assume that the attacker that can obtain an accurate trace.

TABLE III
FEATURE MEASUREMENT ENABLED IN PROFILING

Feature	exists in Case
sm_cycles_active.sum	A, B, C
dram_sectors_read.sum, dram_sectors_write.sum	B, C
L1 transaction ^a , utilization ^b , hit rate ^c	C
L2 transaction ^d , utilization ^e , hit rate ^f	C

^al1tex_t_sectors_pipe_lsu_mem_global_op_ld/st.sum

^bl1tex_lsu_writeback_active.avg.pct_of_peak_sustained_active

^cl1tex_t_sector_hit_rate.pct

^dl2tex_t_sectors_op_read/write.sum

^el2tex_t_sectors.avg.pct_of_peak_sustained_elapsed

^fl2tex_t_sector_hit_rate.pct

LER metric. The LSTM-based predictor for the testbed is a single-layer LSTM-RNN model with a Connectionist Temporal Classification (CTC) decoder as adopted in Deep-sniffer [1]. The *Layer prediction Error Rate* (LER) is used to quantitatively measure the performance of a trained predictor. The LER has the form:

$$LER = \frac{ED(L, L^*)}{|L^*|} \quad (7)$$

where L is the predicted sequence and L^* is the ground-truth, ED denotes editing distance (Levenshtein distance [13]) and $|\cdot|$ denotes the length.

We derive three sets of LSTM predictors - one for each attack case. For each case, we set different number of hidden units of the LSTM network to 64, 96, 128, 256 and 512, resulting in a total of $3 \times 5 = 15$ LSTM predictors. We split each dataset into 4:1 for training and validation subsets, and train for 150 epochs. The final validation LERs for all the LSTM predictors are shown in Table IV. We observe an excellent layer sequence extraction performance on case C where all the latency, DRAM and cache features are considered for each time-step, and a comparative poor performance on case A where only the latency feature is considered.

Predictor Training. The evaluator uses the bagging approach [22] and provides the average LER of LSTM predictors for different input sizes, where we choose the input sizes to match that of CIFAR-10 and ImageNet datasets. The evaluator is shown in Fig. 10. The training needs to be done once for each new device.

TABLE IV
VALIDATION LER OF LSTM-BASED LAYER SEQUENCE PREDICTORS

LSTM unit	CIFAR-10			ImageNet		
	case A	case B	case C	case A	case B	case C
64-unit	0.095	0.100	0.001	0.178	0.027	0.002
96-unit	0.121	0.087	0.007	0.291	0.035	0.000
128-unit	0.126	0.045	0.008	0.292	0.044	0.001
256-unit	0.077	0.023	0.013	0.283	0.031	0.004
512-unit	0.098	0.074	0.000	0.303	0.036	0.003

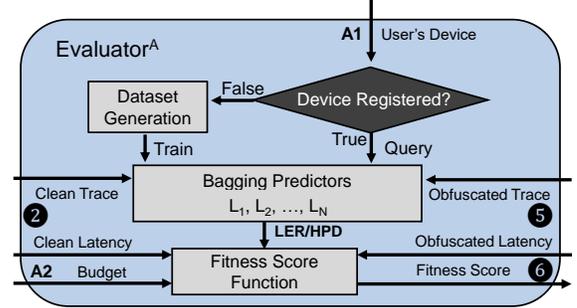


Fig. 10. Design of the **Evaluator** in the proposed framework. Bagging predictors supply LER (in sequence obfuscation) or DER (in dimension obfuscation) to calculate the fitness score function.

2) *GA-based Obfuscator*: Our next goal is to maximize the obfuscation given a user-defined latency budget, B . For instance, $B = 0.1$ means that the user can afford up to 10% extra inference latency. Then the optimization problem can be set up as a constrained discrete optimization problem that maximizes the average LER given the latency budget:

$$\begin{aligned} \min_S \quad & \frac{1}{N} \sum_{i=1}^N LER_i(S) \\ \text{s.t.} \quad & T \leq (1+B)T^* \end{aligned} \quad (8)$$

where, N is the number of predictors in bagging, S denotes the set of obfuscation options, T is the latency with obfuscation and T^* is the clean latency without obfuscation.

Genetic Algorithm. We choose to use the genetic algorithm (GA) to solve the discrete optimization problem. Since the optimization with *constraints* in Eq. (8) cannot be directly used in GA (as GA works well for unconstrained optimization problems [23]), the reward R (a.k.a. fitness score) for GA is designed as follows:

$$R = \frac{1}{N} \sum_{i=1}^N LER_i(S) \Big/ \left[\epsilon + \left(\frac{T - (1+B)T^*}{T^*} \right)^2 \right] \quad (9)$$

We replace the constraints in Eq. (8) with a penalty term, which penalizing the reward when latency T deviates from the total latency $(1+B)T^*$. This deviation is normalized and squared and a small offset term ϵ is added to avoid zero proximity. The block diagram of GA-based obfuscator is shown in Fig. 11.

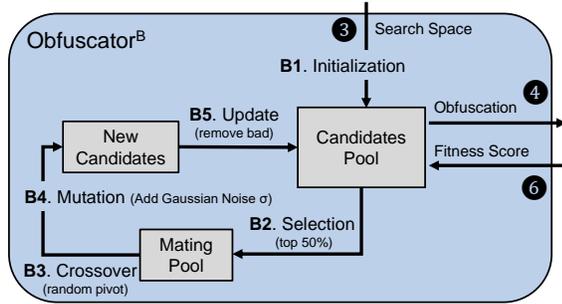


Fig. 11. Design of the GA-based **Obfuscator** in the proposed framework.

The initial value of each obfuscating knob is randomly generated based on the search space provided in step ③ in Fig. 8. For the mating process, we add top 50% of the population based on fitness score into the mating pool. The crossover process takes random pivot of two lists and produces the same number of offsprings. In the mutation process, we apply Gaussian noise with standard deviation σ on each of the offspring obfuscation sets with rounding and clipping to keep the value in legal format. Mutated offsprings are added into the candidate pool, and half of the candidates (including newly added offsprings) that have lowest fitness scores are removed from the candidate pool. The pool of candidates gradually improve over generations resulting in high fitness scores.

B. Dimension Parameter Extraction

For dimension obfuscation, we focus on the obfuscating knobs, namely, layer widening, kernel widening, dummy addition and schedule modification, that affect the dimension parameters the most. We only describe obfuscation on standard Conv2D operators as they appear most frequently in the DNN architectures under testing.

DER metric. To evaluate the prediction error of a layer’s dimension parameter, we define dimension parameter prediction error rate (DER) as a measure of the obfuscation effect similar to the LER metric. If the number of input/output channels of a Conv2D operator be (c, j) , the DER for a given prediction (c, j) on layer i is defined as:

$$DER(i) = \frac{|c - c^*|}{c^*} + \frac{|j - j^*|}{j^*} \quad (10)$$

where, c^* and j^* represent the original (without obfuscation) input/output channels.

Predictor Training. We adopt the Random Forest (RF) model, as a bagging version of decision trees for the dimension parameter extraction testbed. We collected around 50,000 traces for Conv2D operators with different input channel and output channel parameters (c, j) , which are the two most important dimension parameters. We neglect stride, kernel size and padding features because they rarely change. We train RF regression model with different number of trees (30, 50, 100 and 200) to predict c and j separately. The training and validation ratio is set to 4:1. We record the average DER of the validation dataset (20% of the data) for ImageNet and CIFAR-10 for the three attack cases. The results in Table V show that

the dimension extraction has negligible error for cases B and C and comparably high error for case A because it has only latency feature. Furthermore, the number of trees do not affect the prediction performance much.

TABLE V
AVERAGE DER OF RANDOM FOREST REGRESSION MODEL FOR DIMENSION EXTRACTION.

Number of Trees	CIFAR-10			ImageNet		
	case A	case B	case C	case A	case B	case C
30-tree	0.467	0.060	0.014	0.160	0.041	0.023
50-tree	0.465	0.060	0.014	0.160	0.041	0.023
100-tree	0.462	0.059	0.014	0.160	0.040	0.023
200-tree	0.464	0.059	0.014	0.159	0.040	0.023

The dimension obfuscation framework is similar to that of sequence obfuscation framework shown in Fig. 8. The user needs to specify a budget that limits the latency increase in dimension obfuscation for each layer. The evaluator (Fig. 10) uses RF regression model as the “Bagging Predictor” and average DER of all three attack cases to compute the fitness score. Note that the GA-based obfuscator (Fig. 11) for dimension obfuscation has a different search space because a different set of obfuscating knobs is considered.

VI. EVALUATION

A. Sequence Obfuscation Performance

We evaluate the performance of our obfuscation tool on a series of standard models [24]–[26]. Specifically, we select VGG-11, VGG-13, ResNet-20, ResNet-32 models on CIFAR-10 running on a Turing GPU (GTX-1660). We select VGG-19, ResNet-18 and MobileNet-V2 models on ImageNet running on an Ampere GPU (RTX-3090). For the GA, we set the population size to be 16 and run it till the fitness score stabilizes which occurs around 20 generations. The standard deviation σ for the the mutation step is set to a high value (i.e. $\sigma = 8.0$) at the beginning and gets halved after every 4 generations. To eliminate the randomness, for each data point reported here, we choose the average of 3 runs.

Effect of Individual Knobs. First, we investigate the effect of individual knobs on stand-alone Conv2D operators with different dimension parameters. We list the latency overhead for different dimension parameters in Table VI. We found layer branching introduces extra operators with a low latency cost. For example, output-wise layer branching into four adds 3 extra Conv2D operators and 1 concatenate operator with at most 49% latency increase. Selective fusion increases latency by around 15% but it only introduces one extra ReLU operator and BN operator. In contrast, deepening layer and skipping layer introduce an extra Conv2D operator at a lot higher latency cost, and is thus not effective.

Since the latency overhead due to application of an obfuscation knob on a single operator is large, the obfuscation knobs have to be applied selectively to only certain layers. Next, we demonstrate the contribution of individual knobs on

a full model using the GA-based obfuscator. We let only one obfuscating knob be available at a time during the GA search, and keep a fixed budget of $B = 0.02$. The results are shown in Fig. 12. Among all four sequence obfuscating knobs, layer branching and selective fusion have higher LER for the same latency budget and are clearly better choices. The selective combination of 4 knobs by NeurObfuscator achieves stronger obfuscation than any single knob, as expected.

TABLE VI
EFFECT OF INDIVIDUAL SEQUENCE OBFUSCATING KNOBS

Knobs	Extra Operator	Latency Overhead
Branching (output-wise by 2)	$1 \times \text{Conv2D}$, $1 \times \text{Concat}$	21% ~ 27%
Branching (output-wise by 4)	$3 \times \text{Conv2D}$, $1 \times \text{Concat}$	38% ~ 49%
Selective Fusion (N=0)	$1 \times \text{ReLU}$, $1 \times \text{BN}$	14% ~ 15%
Deepen	$1 \times \text{Conv2D}$ (1x1 kernel)	39% ~ 89%
Skipping	$1 \times \text{Conv2D}$	70% ~ 130%

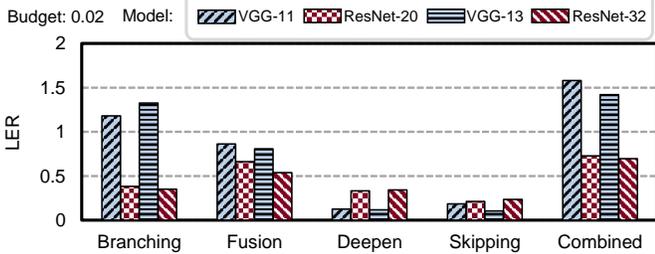


Fig. 12. Individual contribution of each sequence obfuscating knob in GA-obfuscator, followed by combination effect. Time budget is selected at 0.02. VGG-11, VGG-13, ResNet-20 and ResNet-32 are on CIFAR-10 dataset.

NeurObfuscator - Sequence Obfuscation. We demonstrate the performance of NeurObfuscator on CIFAR-10 and ImageNet datasets. For VGG-11, VGG-13 and ResNet-32 running on CIFAR-10, we use bagging of all 15 LSTM predictors. The LER results under different latency budgets are shown in Fig. 13. We notice that the LER absolute value is high for VGG-11 and VGG-13 while low for ResNet-20 and ResNet-32. This is because LER is the layer editing distance divided by total number of layers of the vanilla architecture (without obfuscation) and the sequence obfuscation affects the absolute editing distance directly rather than the relative editing distance (i.e., LER).

For VGG-19, ResNet-18 and MobileNet-V2 on ImageNet dataset, case A and case B LSTM predictors struggle to get good extraction performance, i.e., provide low LER for the baseline architecture. So we use bagging of three “elite” LSTM predictors (number of units of 128, 256, 512 LSTM predictors in case C), which have near-zero clean LER. The results are shown in Fig. 14. Moreover, we observe that LER increases sub-linearly with increasing latency budget. This is because the search space is kept fixed, and the most effective knobs with low latency overhead are chosen up front and so increasing the budget only allows knobs that are not as effective to get added to the obfuscation set.

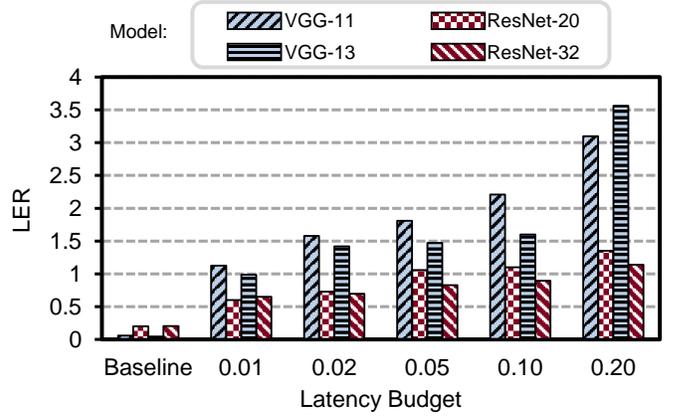


Fig. 13. Sequence obfuscation results on typical architectures on CIFAR-10 dataset including VGG-11, VGG-13, ResNet-20 and ResNet-32.

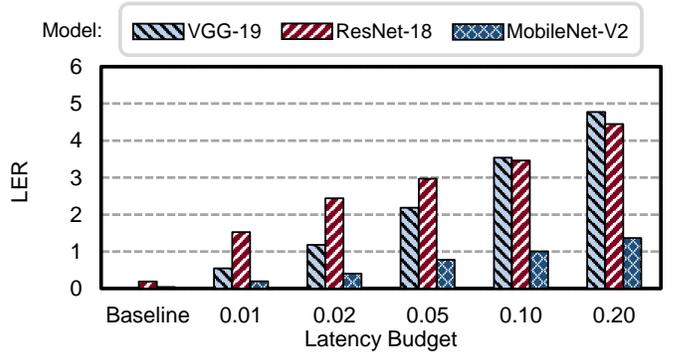


Fig. 14. Sequence obfuscation results on typical architectures on ImageNet dataset including VGG-19, ResNet-18 and MobileNet-V2.

Summary 1: We demonstrate the performance of four knobs, namely, layer deepening, layer skipping, layer branching and selective fusion, for sequence obfuscation. While layer branching and selective fusion have relatively strong performance, combination of all four knobs by GA in NeurObfuscator results in the strongest performance. We evaluated our tool on multiple models taking CIFAR-10 and ImageNet datasets as input data. On a ResNet-18 ImageNet model, we achieved a 2.44 LER (translates to 44 layers’ difference) with a mere 2% inference latency overhead.

B. Dimension Obfuscation Performance

For dimension obfuscation, we use the RF regression testbed and DER metric (Eq. (10)) to evaluate the effect of obfuscation. We pick a Conv2D layer (C2) with 64 input channels and 128 output channels with 3×3 kernel from VGG-19 network as an example. The top subfigure in Fig. 15 shows the layer operators marked by sequence obfuscation. Here C1 is the input Conv2D layer with 3 input channels and 64 output channels. In this example, the job of dimension extraction is to correctly predict the number of input channels and output channels of C2. We use the predictor to predict the output channel of C1 and input/output channel of C2. Here the ground-truth 64 is predicted twice: once as output channel of

C1 and once as input channel of C2. The average is taken if the two predictions do not match. Next, the effect of individual obfuscating knobs is evaluated. The DER and latency overhead for each knob are shown in Fig. 15.

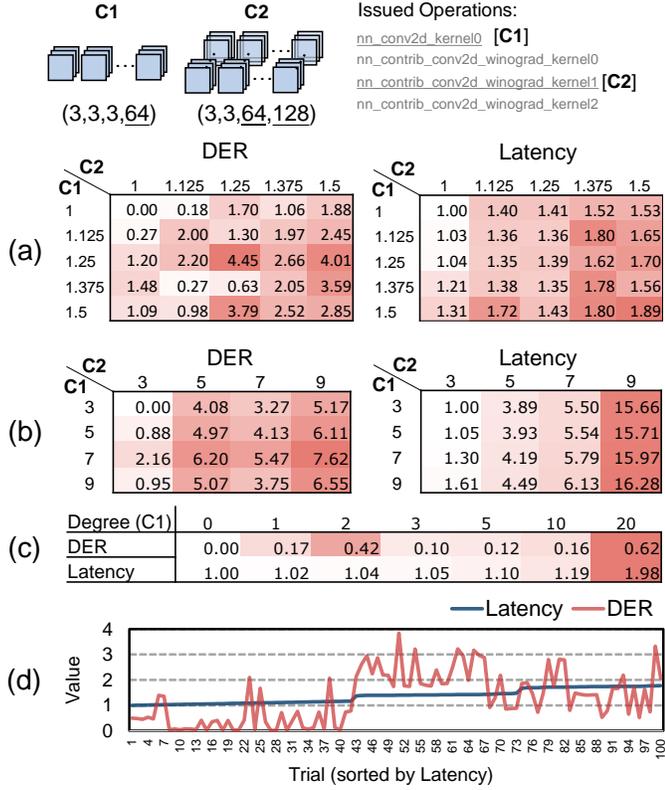


Fig. 15. Dimension Obfuscation on a Conv2D layer with 64 input channel and 128 output channel. Application of (a) layer widening to C1/C2, (b) kernel widening to C1/C2, (c) dummy addition to only C2, and (d) random schedule modification that results in highest DER for a given latency overhead.

Layer Widening. We use grid-search on applying widening factor of 1 to 1.5 (3/2) for C1 and C2. As shown in Fig. 15 (a), generally, applying higher widening factor increases the DER, and increases the latency. We found a sweet point where increasing the C1 output channel size by 1.25 \times can achieve a 1.20 DER with 1.04 \times latency.

Kernel Widening. Kernel widening affects both types of Conv2D operator. However, as shown in Fig. 15 (b), in most cases the large overhead makes it an expensive option to use in practice. The exception is that increasing kernel size of C1 from 3 \times 3 to 5 \times 5 results in 0.88 DER with 1.05 \times latency.

Dummy Addition. Dummy addition does not affect the dimension parameters of C2, because dummy operator are issued after “winograd_kernel2” and will not be fused into kernel1. However for a standard Conv2D such as C1, dummy addition has a dramatic effect. As shown in Fig. 15 (c), DER increases with increasing dummy addition factor and reaches a sweet point when dummy addition factor is 2; the corresponding DER is 0.42 and latency is 1.04 \times .

Schedule Modification. For the schedule modification knob, we target the schedules of two templates (plain-Conv2D and winograd-Conv2D), with a total of 13 distinct tunable

parameters. Since the search space is very large, we perform 100 trials of random choices. Fig. 15 (d) plots DER as a function of increasing latency. We see that there are DER spikes (value larger than 1.0) at trials 7, 23 and 25, even when the increase in latency is 1%. Thus, schedule modification is by far the most effective knob in dimension obfuscation.

NeurObfuscator - Dimension Obfuscation. We demonstrate the performance of NeurObfuscator on dimension parameter obfuscation. Using the same GA setting as in sequence obfuscation, and replacing the LER with DER, we obtain the results shown in Table VII. The final results are significantly better than when individual obfuscation knobs are used! A high DER of 2.51 is achieved with only 0.02 latency budget, i.e. a 2% increase in inference latency. This corresponds to the case where the original layer with 64 input c and 128 output channels is extracted to 207 input and 93 output channels.

TABLE VII
GA RESULTS FOR DIMENSION OBFUSCATION

Budget	0.00	0.01	0.02	0.05	0.10	0.20
DER	0.00	2.05	2.51	2.80	3.24	3.43
Prediction	(64, 128)	(177, 91)	(207, 93)	(225, 92)	(176, 319)	(141, 413)

Summary 2: We demonstrate the performance of four dimension obfuscating knobs, namely, layer widening, kernel widening, dummy addition and schedule modification. While schedule modification has the strongest performance among all four, NeurObfuscator achieves the best dimension obfuscation, as expected. On an example Conv2D layer with 64 input channels and 128 output channels, RF-based dimension extraction achieves 2.05 DER and 2.51 DER under 1% and 2% inference latency overhead, respectively.

VII. CONCLUSIONS AND FUTURE WORK

To mitigate the neural architecture stealing on GPU devices, we propose NeurObfuscator, a NN obfuscating tool that provides both sequence obfuscation and dimension obfuscation. We propose to use a total of eight obfuscating knobs across scripting, optimization and scheduling phases of a neural network model execution. Application of these knobs affect the number of computations, latency and number of memory accesses, thus altering the execution trace. To achieve the best obfuscation performance for a user-defined latency overhead, we leverage the genetic algorithm to identify the best combination of obfuscation knobs. On a ResNet-18 ImageNet model, sequence obfuscation helps achieve a 2.44 LER (which translates to 44 layers’ difference) with merely 2% latency overhead. Similarly, with 2% latency overhead, dimension obfuscation can achieve 2.51 DER corresponding to the case when a 64 input channel, 128 output channel Conv2D gets extracted as 207 input channel and 93 output channel.

While the proposed methodology has been designed for GPU devices, we plan to extend this to other hardware substrates, such as FPGAs and ASICs. Furthermore, we plan to

evaluate the end-to-end performance of such a system. Examples include accuracy evaluation if the obfuscated architecture is used to train a new model, and the attack success rate of the transfer adversarial attack if the obfuscated architecture is used to build a surrogate model.

REFERENCES

- [1] X. Hu, L. Liang, S. Li, L. Deng, P. Zuo, Y. Ji, X. Xie, Y. Ding, C. Liu, T. Sherwood *et al.*, “DeepSniffer: A dnn model extraction framework based on learning architectural hints,” in *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2020, pp. 385–399.
- [2] J. Wei, Y. Zhang, Z. Zhou, Z. Li, and M. A. Al Faruque, “Leaky dnn: Stealing deep-learning model secret with gpu context-switching side-channel,” in *2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 2020, pp. 125–137.
- [3] V. Duddu, D. Samanta, D. V. Rao, and V. E. Balas, “Stealing neural networks via timing side channels,” *arXiv preprint arXiv:1812.11720*, 2018.
- [4] S. Liu, Y. Wei, J. Chi, F. H. Shezan, and Y. Tian, “Side channel attacks in computation offloading systems with gpu virtualization,” in *2019 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2019, pp. 156–161.
- [5] H. Sak, A. Senior, and F. Beaufays, “Long short-term memory recurrent neural network architectures for large scale acoustic modeling,” in *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [6] F. Tramèr, F. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart, “Stealing machine learning models via prediction apis,” in *25th {USENIX} Security Symposium ({USENIX} Security 16)*, 2016, pp. 601–618.
- [7] B. Wang and N. Z. Gong, “Stealing hyperparameters in machine learning,” in *2018 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2018, pp. 36–52.
- [8] C. Liu, A. Harris, M. Maas, M. Hicks, M. Tiwari, and E. Shi, “Ghostrider: A hardware-software system for memory trace oblivious computation,” *ACM SIGPLAN Notices*, vol. 50, no. 4, pp. 87–101, 2015.
- [9] E. Stefanov, M. Van Dijk, E. Shi, C. Fletcher, L. Ren, X. Yu, and S. Devadas, “Path oram: an extremely simple oblivious ram protocol,” in *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, 2013, pp. 299–310.
- [10] E. Karimi, Y. Fei, and D. Kaeli, “Hardware/software obfuscation against timing side-channel attack on a gpu,” in *2020 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 2020, pp. 122–131.
- [11] Q. Xu, H. Naghibijouybari, S. Wang, N. Abu-Ghazaleh, and M. Annavaram, “Gpuguard: mitigating contention based side and covert channel attacks on gpus,” in *Proceedings of the ACM International Conference on Supercomputing*, 2019, pp. 497–509.
- [12] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, H. Shen, M. Cowan, L. Wang, Y. Hu, L. Ceze *et al.*, “{TVM}: An automated end-to-end optimizing compiler for deep learning,” in *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*, 2018, pp. 578–594.
- [13] G. Navarro, “A guided tour to approximate string matching,” *ACM computing surveys (CSUR)*, vol. 33, no. 1, pp. 31–88, 2001.
- [14] T. Chen, L. Zheng, E. Yan, Z. Jiang, T. Moreau, L. Ceze, C. Guestrin, and A. Krishnamurthy, “Learning to optimize tensor programs,” *arXiv preprint arXiv:1805.08166*, 2018.
- [15] T. Chen and C. Guestrin, “Xgboost: A scalable tree boosting system,” in *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, 2016, pp. 785–794.
- [16] H. Yu, H. Ma, K. Yang, Y. Zhao, and Y. Jin, “Deepem: Deep neural networks model recovery through em side-channel information leakage,” in *2020 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 2020, pp. 209–218.
- [17] X. Hu, L. Liang, L. Deng, S. Li, X. Xie, Y. Ji, Y. Ding, C. Liu, T. Sherwood, and Y. Xie, “Neural network model extraction attacks in edge devices by hearing architectural hints,” *arXiv preprint arXiv:1903.03916*, 2019.
- [18] H. Naghibijouybari, A. Neupane, Z. Qian, and N. Abu-Ghazaleh, “Rendered insecure: Gpu side channel attacks are practical,” in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, 2018, pp. 2139–2153.
- [19] H. Zhu, Z. An, C. Yang, K. Xu, E. Zhao, and Y. Xu, “Eena: efficient evolution of neural architecture,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [20] T. Chen, I. Goodfellow, and J. Shlens, “Net2net: Accelerating learning via knowledge transfer,” *arXiv preprint arXiv:1511.05641*, 2015.
- [21] M. Wistuba, “Deep learning architecture search by neuro-cell-based evolution with function-preserving mutations,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2018, pp. 243–258.
- [22] L. Breiman, “Bagging predictors,” *Machine learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [23] S. Rajeev and C. Krishnamoorthy, “Discrete optimization of structures using genetic algorithms,” *Journal of structural engineering*, vol. 118, no. 5, pp. 1233–1250, 1992.
- [24] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [25] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [26] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, “Mobilenets: Efficient convolutional neural networks for mobile vision applications,” *arXiv preprint arXiv:1704.04861*, 2017.