

HierNet: A Hierarchical Deep Learning Model for SCA on Long Traces

Suvadeep Hajra* and Debdeep Mukhopadhyay*

**Indian Institute of Technology Kharagpur, India*

Abstract—In Side-Channel Analysis (SCA), statistical or machine learning methods are employed to extract secret information from power or electromagnetic (EM) traces. In many practical scenarios, raw power/EM traces can span hundreds of thousands of features, with relevant leakages occurring over only a few small segments. Consequently, existing SCAs often select a small number of features before launching the attack, making their success highly dependent on the feasibility of feature selection. However, feature selection may not always be possible, such as in the presence of countermeasures like masking or jitters.

Several recent works have employed deep learning (DL) methods to conduct SCA directly on long raw traces, thereby reducing reliance on feature selection steps. However, these methods often perform poorly in the presence of various jitter-based countermeasures. While some methods, such as the one proposed by Hajra et al. (TCHES 2024), have shown high robustness to several jitter-based countermeasures on relatively short traces, we demonstrate in this work that their performance deteriorates as trace lengths increase. Consequently, existing DL models perform poorly on long traces or in the presence of various jitter-based countermeasures.

To address these challenges, we propose HierNet, a hierarchical DL model designed for SCA on long traces. HierNet uses a two-level information assimilation process: at the base level, a shift-invariant DL model processes smaller trace segments; at the top level, another DL model aggregates the outputs from all segments to produce the final output. HierNet has been experimentally evaluated against various combinations of masking, random delay, and clock jitter countermeasures, using three publicly available SCA datasets with trace lengths up to 250K features. The results have been compared with four existing SCA benchmark models. They demonstrate HierNet’s superiority, particularly on long traces or in the presence of clock jitter countermeasures, showcasing its ability to reach guessing entropy 1 with fewer than or around 10 attack traces, while the benchmark models fail to do so even using 5K attack traces. HierNet also exhibits significantly better results in low training data scenarios.

Index Terms—SCA, Deep Learning, Shift-invariance, Transformer Network

1. Introduction

In modern days, cryptographic systems serve as the foundation of security and privacy in digital systems. However, these cryptographic systems can be vulnerable to attacks that exploit unintentional information leakage from physical implementations. Side-channel analysis (SCA) [1] exploits these unintended channels of information leakage, such as power consumption [2], electromagnetic (EM) emanations [3], or timing variations [1], arising from the physical implementation of cryptographic systems to compromise their security. This work concentrates on power/EM SCA.

A CMOS device’s power consumption, or EM emanation depends on the data being processed within the device. This characteristic is leveraged by SCAs to extract the secret key used in a cryptographic device. SCAs involve the collection of power/EM profiles (commonly referred to as power/EM traces) from one or more devices during their execution and recovering the secret key by analyzing these traces using statistical techniques or machine learning methods. Template attack [4] and stochastic attack [5] are two classical approaches to SCAs. More recently, various machine learning (ML) techniques [6], including deep learning (DL) [7], are introduced in SCA.

In many practical scenarios, raw power/EM traces can span hundreds of thousands of features, with the leakages of relevant cryptographic operations occurring over a few small segments, known as Points of Interest (PoIs). However, most of the existing works in SCA literature conduct attacks on some pre-selected high signal-to-noise ratio (SNR) features or on a pre-selected attack window containing the PoIs. For example, classical SCA like template and stochastic attacks select a few (in the order of 100) high SNR features prior to the actual attack. Similarly, most of the ML and DL-based methods are investigated on some pre-selected attack window of size in the order of 1K. However, in implementations protected by masking countermeasures [8], such feature selection requires the knowledge of additional intermediate variables of the cryptographic implementation, which might not always be available. Moreover, even with such knowledge, the feature selection might not be possible in the presence of jitter-based countermeasures (like random delay [9], clock jitter [10], and shuffling [11]). Therefore, the success of the most of the existing SCAs depends on the feasibility of feature selection.

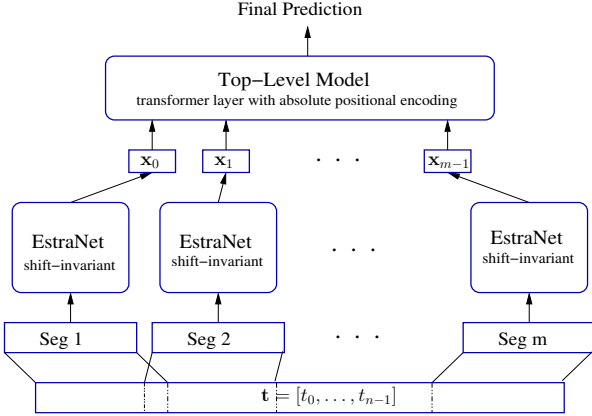


Figure 1. HierNet Architecture. All parameters except the last layer of the EstraNet models are shared.

Several deep learning (DL) based research efforts [12], [13], [14], [15], [16] have focused on improving the scalability and effectiveness of DL models for long attack windows or raw traces, thereby reducing their sensitivity to the attack window selection. Those models have demonstrated good performance even without selecting any smaller attack window. Moreover, several works, such as [13], [14], have demonstrated significantly improved performance when attacking long traces compared to attacking a smaller, selected attack window. However, the performance of most of those methods vastly deteriorates in the presence of countermeasures like random delay and clock jitter. For example, the investigation of [14] reveals that their methods might require more than 30 times more traces for the successful attack (to reach the guessing entropy 1) when subjected to desynchronized traces as compared to synchronized traces. Similarly, the results of [16] illustrate that the DL models from [13] may fail to reach guessing entropy 1 even when using up to $5K$ traces if the traces are misaligned using clock jitter. In contrast, they require only around 10 traces to reach the same in absence of clock jitter. In [17], Bursztein et al. have not shown the success of their DL model, GPAM, against jitter-based countermeasures, though they demonstrated its success on traces with lengths more than $10M$.

Recently, Hajra et al. [16] introduced a Transformer Network (TN)-based shift-invariant DL model, named EstraNet, for SCA on long traces, which has demonstrated state-of-the-art results against random delay and clock jitter countermeasures even when they are combined with masking. In this work, we closely analyze EstraNet and show that its performance against masking countermeasures rapidly deteriorates as trace lengths increase. Therefore, the existing DL models either perform poorly on long traces or in the presence of various jitter-based countermeasures.

In an effort to develop a DL model that is scalable to long traces while being robust against countermeasures such as masking and various kinds of jitter, we introduce HierNet, a hierarchical DL model for SCA. HierNet adopts a two-level information assimilation approach, as illustrated

in Figure 1. At the lower level, it processes smaller trace segments using a base DL model, summarizing the information extracted from each segment into distinct vectors. A top-level DL model then aggregates these outputs from all segments to produce the final result. We use EstraNet as the base model due to its shift-invariance and strong resilience against masking and jitter-based countermeasures on relatively smaller attack windows [16]. The top-level model is designed with a transformer layer that incorporates a novel self-attention mechanism with absolute positional encoding, allowing it to (a) generalize effectively across different trace lengths and segmentations and (b) exploit the positional information of the features. As a result, HierNet retains EstraNet’s robustness against masking, random delay, and clock jitter countermeasures while overcoming its limitations with longer traces.

Contributions:

The contributions of the work are the following:

- EstraNet, the shift-invariant DL model recently introduced in [16], demonstrated the state-of-the-art performance against various combinations of masking, random delay and clock jitter countermeasures. The model also has linear time and memory cost, making it computationally scalable to long traces. We investigate the attack efficacy of EstraNet on longer traces. We identify that the performance of EstraNet degrades with increasing attack window size or trace length, necessitating a) more effort in hyper-parameter tuning, b) potential overfitting issues, and c) diminished performance.
- We introduce HierNet, a hierarchical DL model for SCA on long traces. HierNet employs a two-level hierarchical DL architecture. At the first level, it utilizes a shift-invariant base DL model to assimilate information from smaller trace segments. At the top level, another DL model aggregates the outputs of the base model from all segments. The top-level DL model of HierNet employs a transformer layer, enabling HierNet to generalize effectively across different trace lengths and segmentations. Additionally, this transformer layer incorporates a novel self-attention operation with absolute positional encoding, exploiting positional information of the features.
- We experimentally evaluate HierNet against various combinations of masking, random delay, and clock jitter countermeasures using three publicly available datasets with trace lengths of up to $250K$ features. Its performance is compared with three benchmark DL models. The results show that HierNet can significantly outperform the benchmark models, particularly in the presence of clock jitter countermeasures, or on longer traces. HierNet can reach the guessing entropy 1 using less than or close to ten attack traces, whereas the benchmark models struggle to reach the same even using as many as $5K$ attack traces in these scenarios. HierNet has

also demonstrated significantly better performance on low training data scenarios.

- We perform an ablation study on various design choices of the top-level DL model of HierNet, underscoring the importance of using the transformer layer and absolute positional encoding in the model.

The paper is organized as follows. Section 2 introduces the notations and provides background information. In Section 3, we briefly describe the EstraNet model. Section 4 investigates the effectiveness of EstraNet on long traces. In Section 5, we introduce the proposed HierNet architecture. Section 6 presents the empirical evaluation of HierNet. In Section 7, we explore the future directions for the work. Finally, we conclude the paper in Section 8.

2. Preliminaries

2.1. Notations

The notational conventions used in the paper are as follows. A letter in the capital (like X), the corresponding small letter (like x), and the corresponding calligraphic letter (like \mathcal{X}) are respectively used to represent a random variable, an instantiation, and the domain of the random variable. Similarly, a capital letter in bold (like \mathbf{X}) and the corresponding small letter in bold (like \mathbf{x}) are respectively used to represent a random vector and its instantiation. We use a capital letter in Roman style (like M) to represent a matrix. The i -th element of a vector \mathbf{x} is represented by $\mathbf{x}[i]$. Similarly, the element of i -th row and j -th column of a matrix M is represented by $M[i, j]$. The notations $\mathbb{P}[\cdot]$ and $\mathbb{E}[\cdot]$ are respectively used to represent a random variable's probability distribution function and expectation.

2.2. Side-Channel Analysis

The power consumption and EM emanation of a semiconductor device are influenced by the values manipulated within the device. SCA takes advantage of this behavior in semiconductor devices to extract information about the sensitive variables in a cryptographic implementation, revealing the device's secret key. To achieve this, in an SCA, an attacker gains control of the target device, also known as Device Under Test (DUT). Then he/she collects power/EM measurements, known as traces, by repeatedly executing the encryption/decryption algorithm with various plaintexts/ciphertexts. Subsequently, the attacker uses statistical tests to deduce the device's secret key from the traces.

SCA can be categorized into two types: profiling SCA and non-profiling SCA. In profiling SCA, the attacker possesses a clone of the DUT under control. With this clone device, he/she can create a profile of the DUT's power consumption or EM emission characteristics and employs this profile to carry out the attack on the DUT. On the other hand, in non-profiling SCA, the attacker lacks a clone device and, therefore, cannot build a power/EM profile a priori. Instead, he/she attempts to recover the secret key solely from the traces of the DUT. This paper considers profiling SCA.

2.3. Deep Learning based Profiling SCA

Like any profiling SCA, DL-based profiling SCA also involves a two-phase process. In the initial phase, known as the profiling phase, the attacker configures a predefined key in the clone device and acquires a large number of traces by executing encryption (or decryption) operations on well-known plaintexts (or ciphertexts) using the device. For each trace, the attacker calculates the value of an intermediate secret variable denoted as $Z = F(X, K)$. Here, X represents a component of the random plaintext (or ciphertext), K represents a component of the known key, and $F(\cdot, \cdot)$ denotes a cryptographic primitive. Then, the attacker uses the traces and the values of the intermediate variable Z to train a DL model to predict Z from the \mathbf{L} , where \mathbf{L} represents a trace. More specifically, the DL model takes a trace \mathbf{L} as its input and outputs a probability distribution over the values of Z . Let $f(\cdot; \theta^*)$ denote the trained DL model, with θ^* denoting the model parameters. The output of the DL model for a trace, \mathbf{l} , can be given as $\mathbf{p} = f(\mathbf{l}; \theta^*)$ where $\mathbf{p} \in \mathbb{R}^{|\mathcal{Z}|}$ with $\mathbf{p}[j]$, for $j = 0, \dots, |\mathcal{Z}| - 1$, denoting the predicted probability of the intermediate variable $Z = j$.

In the attack phase, attacker collect multiple traces $\{\tilde{\mathbf{I}}^i\}_{i=0}^{T_a-1}$ by executing the DUT. Let $\{\tilde{\mathbf{p}}^i\}_{i=0}^{T_a-1}$ be the corresponding plaintexts (or ciphertexts). Then, he/she uses the set of attack trace-plaintext pairs $\{(\tilde{\mathbf{I}}^i, \tilde{\mathbf{p}}^i)\}_{i=0}^{T_a-1}$ to compute the score of each key $k \in \mathcal{K}$ as

$$\hat{\delta}_k = \sum_{i=0}^{T_a-1} \log \mathbf{p}^i[F(\tilde{\mathbf{p}}^i, k)] \quad (1)$$

where $\mathbf{p}^i = f(\tilde{\mathbf{I}}^i; \theta^*)$ is the output of the DL model for the i -th trace. The key with the highest score, i.e., $\hat{k} = \operatorname{argmax}_k \hat{\delta}_k$, is considered as the predicted key. The prediction is said to be correct if $\hat{k} = k^*$ holds. Alternatively, the rank of the correct key within the list of all possible keys, sorted based on their respective scores $\hat{\delta}_k$, serves as a metric to determine the success level of the attack. The rank of the correct key averaged over multiple repetitions of the attack, referred to as guessing entropy, is a widely used metric to evaluate the success of an attack.

2.4. Transformer Network

A Transformer Network (TN) is a kind of DL model composed of multiple transformer layers. The input trace is fed into the first transformer layer, and the output of each layer serves as the input for the subsequent layer. The final output of the last transformer layer is considered the output of the TN model.

Each transformer layer consists of a self-attention layer and a position-wise feed-forward layer. It takes a sequence $\mathbf{X} = [\mathbf{x}_0, \dots, \mathbf{x}_{n-1}]^T \in \mathbb{R}^{n \times d}$ of n feature vectors as input and outputs another sequence $\mathbf{Y} = [\mathbf{y}_0, \dots, \mathbf{y}_{n-1}]^T \in \mathbb{R}^{n \times d}$ where n is the sequence length (trace length in the context of SCA) and d is the dimension of the feature vectors. The output \mathbf{Y} is given by

$$\begin{aligned}\hat{Y} &= f_{SA}(X) + X \\ Y &= f_{PFF}(\hat{Y}) + \hat{Y}\end{aligned}$$

where $f_{SA} : \mathbb{R}^{n \times d} \mapsto \mathbb{R}^{n \times d}$ is the function representing the self-attention layer and $f_{PFF} : \mathbb{R}^{n \times d} \mapsto \mathbb{R}^{n \times d}$, is the function representing the position-wise feed-forward layer. The position-wise feed-forward layer independently applies a non-linear transformation to each feature vector, increasing the non-linearity (and, thus, its representation power) of the TN model. On the other hand, the self-attention layer captures the inter-dependencies among different features [15], [16], [18]. It is worth mentioning that the number of parameters in both the self-attention and position-wise feed-forward layer is independent of the sequence length n . Therefore, it can scale to a very long sequence without being prone to overfitting.

The self-attention layer is a key component of a Transformer Network (TN). Different variations of this layer lead to vastly different types of TN models, each varying significantly in performance and computational cost. In the next part of this section, we briefly describe the self-attention layer and some of its variations, setting the groundwork for the following discussion.

2.4.1. Self-attention Layer. Let the sequence $X = [\mathbf{x}_0, \dots, \mathbf{x}_{n-1}]^T \in \mathbb{R}^{n \times d}$ of n feature vectors is the input of a self-attention layer and $\hat{Y} = [\hat{\mathbf{y}}_0, \dots, \hat{\mathbf{y}}_{n-1}]^T \in \mathbb{R}^{n \times d}$ be the corresponding output sequence. Then, in a self-attention layer, each output feature vector $\hat{\mathbf{y}}_i$, for $i = 0, \dots, n-1$, is computed as

$$\hat{\mathbf{y}}_i = \sum_{j=0}^{n-1} \frac{k(f(\mathbf{x}_i), g(\mathbf{x}_j))}{\sum_{l=0}^{n-1} k(f(\mathbf{x}_i), g(\mathbf{x}_l))} \phi(\mathbf{x}_j) \quad (2)$$

where $f(\mathbf{x}_i)$, $g(\mathbf{x}_i)$, and $\phi(\mathbf{x}_i)$ are some vector valued functions of input feature \mathbf{x}_i , and $k : \mathbb{R}^{d_k} \times \mathbb{R}^{d_k} \mapsto \mathbb{R}^+$ is called the kernel function. The kernel function outputs a high score, $k(f(\mathbf{x}_i), g(\mathbf{x}_j))$, if it finds the features \mathbf{x}_i , and \mathbf{x}_j to be related and a low score otherwise. In summary, in self-attention, each output feature, $\hat{\mathbf{y}}_i$, is computed by taking the weighted sum of the input features $\phi(\mathbf{x}_0), \dots, \phi(\mathbf{x}_{n-1})$ where the weight for the j -th input feature is given by the normalized score $k(f(\mathbf{x}_i), g(\mathbf{x}_j)) / \sum_{l=0}^{n-1} k(f(\mathbf{x}_i), g(\mathbf{x}_l))$. Therefore, if the i -th feature is related to (or dependent on) the j -th feature, the weight of the j -th input feature in the i -th output feature can be significantly greater than weights of the other features, increasing the component of the j -th feature (in the input) in the i -th feature (in the output). In that way, self-attention layer can combine the information of the i and j -th features.

One of the most common kernel used in TN literature is obtained by setting the kernel as $k(f(\mathbf{x}_i), g(\mathbf{x}_j)) = \exp(\langle f(\mathbf{x}_i), g(\mathbf{x}_j) \rangle / \sqrt{d_k})$, where $\langle \cdot, \cdot \rangle$ represents the dot product between two vectors and d_k is the dimension of the outputs of $f(\cdot)$ and $g(\cdot)$. The above kernel leads to following self-attention known as softmax self-attention:

$$\hat{\mathbf{y}}_i = \sum_{j=0}^{n-1} \frac{\exp(\langle f(\mathbf{x}_i), g(\mathbf{x}_j) \rangle / \sqrt{d_k})}{\sum_{l=0}^{n-1} \exp(\langle f(\mathbf{x}_i), g(\mathbf{x}_l) \rangle / \sqrt{d_k})} \phi(\mathbf{x}_j) \quad (3)$$

$$= \sum_{j=0}^{n-1} \text{softmax}(\langle f(\mathbf{x}_i), g(\mathbf{x}_j) \rangle / \sqrt{d_k})_j \phi(\mathbf{x}_j) \quad (4)$$

In the self-attention mechanism of the form of Eq. (2), the attention score from feature \mathbf{x}_i to feature \mathbf{x}_j does not depend on their positions, i or j , in the input sequence. Therefore, the actual order of the features in the input sequence does not affect the attention scores. In other words, permuting the input sequence also permutes the output sequence in the same way.

In many tasks, we want the attention a feature gives to another feature to depend on their absolute positions or the distance between them. For instance, in the presence of random delay or clock jitter countermeasures, the distances between the informative features remain almost the same across all traces. These properties are exploited in the TN models proposed in [16], [18] by making the self-attention shift-invariant. This approach is described as follows:

Relative Positional Encoding. To make the self-attention shift-invariant, one can incorporate relative positional encoding into kernel function by generalizing it to $k_r(f(\mathbf{x}_i), g(\mathbf{x}_j), e(i-j))$, where $e(i-j)$ is a vector space encoding of the relative distance $i-j$. Therefore, the output of the resultant self-attention takes the form

$$\hat{\mathbf{y}}_i = \sum_{j=0}^{n-1} \frac{k_r(f(\mathbf{x}_i), g(\mathbf{x}_j), e(i-j))}{\sum_{l=0}^{n-1} k_r(f(\mathbf{x}_i), g(\mathbf{x}_l), e(i-l))} \phi(\mathbf{x}_j) \quad (5)$$

With relative positional encoding, the kernel function returns a high score not only based on the similarity between \mathbf{x}_i , and \mathbf{x}_j , but also considering their relative distance $i-j$.

Another scenario where incorporating position information into the attention scores is beneficial is when there is no jitter-based countermeasure in the implementation. In the absence of such countermeasures, the positions of the informative features in the input sequence remain consistent across all traces. To exploit this property, one approach is to make the attention scores dependent on the absolute positions of the features. In this work, we achieve this by incorporating absolute positional encoding into the self-attention layer. We introduce the self-attention with absolute positional encoding in Section 5.2.

The next section briefly describes the EstraNet architecture for analyzing its performance on longer traces.

3. EstraNet Model

We first provide a brief overview of EstraNet's overall architecture, followed by a discussion of its key component, the GaussiP attention.

3.1. The Overall Architecture

The overall EstraNet architecture is shown in Figure 2a. Being a transformer network, EstraNet consists of

multiple EstraNet layers stacked one after another. The output sequence of the top EstraNet layer (the sequence $\{\mathbf{y}_0, \dots, \mathbf{y}_{m-1}\}$, where each $\mathbf{y}_i \in \mathbb{R}^d$) is reduced to a single d -dimensional vector $\bar{\mathbf{y}}$ using a softmax attention layer. The output $\bar{\mathbf{y}}$ is fed to a classification layer for classification.

Like a transformer layer, each EstraNet layer consists of a self-attention layer and a position-wise feed-forward layer (shown in Figure 2b). However, instead of vanilla self-attention, the EstraNet layer uses a novel self-attention named GaussiP attention. The standard layer normalization operation of the vanilla transformer layer is also replaced by a novel layer-centering operation in EstraNet layer.

3.2. GaussiP Attention

The GaussiP attention is similar to self-attention with relative positional encoding given in Eq. (5). However, it uses a form of Gaussian kernel over the relative distance $i - j$ and special encoding to facilitate a) information flow from one feature to a distant feature, and b) controllable sparsity of the attention scores. Moreover, it uses a closed-form approximation for the denominator in Eq. (5). More precisely, the GaussiP attention's kernel takes the form:

$$\begin{aligned} k_{GPA}(i-j; s_p, c_p) &= \exp\left(-\frac{\|\mathbf{i}-\mathbf{j}\|_2^2}{2}\right) \\ &= \exp\left(-\frac{\beta_2^2 s_p^2 (i-j-c_p n)^2 \|\mathbf{W}_p \mathbf{p}\|_2^2}{2}\right) \\ &= \exp\left(-\frac{\hat{\beta}_2^2 s_p^2 (i/n-j/n-c_p)^2 \|\mathbf{W}_p \mathbf{p}\|_2^2}{2}\right), \end{aligned}$$

where $\mathbf{i} = \beta_2 s_p \mathbf{W}_p (\mathbf{b} + i \mathbf{p})$ and $\mathbf{j} = \beta_2 s_p \mathbf{W}_p (\mathbf{b} + j \mathbf{p} + c_p n \mathbf{p})$ be the encoding of the positional indices i and j , with $s_p \in \mathbb{R}^+$, $c_p \in (0, 1)$ being the parameters, $\mathbf{W}_p \in \mathbb{R}^{d_e \times d}$, $\mathbf{p} \in \mathbb{R}^d$ being two constants, and $\hat{\beta}_2 = n \beta_2 > 0$ being a hyper-parameter of the model. To make the time and memory complexity of the attention linear in the sequence length n , the above kernel is approximately factorized using the d_p -dimensional Fourier feature map defined as

$$\begin{aligned} \phi_{fr}(\mathbf{x}) &= \frac{1}{\sqrt{d_p}} [\sin(\mathbf{w}_0^T \mathbf{x}), \dots, \sin(\mathbf{w}_{d_p-1}^T \mathbf{x}), \\ &\quad \cos(\mathbf{w}_0^T \mathbf{x}), \dots, \cos(\mathbf{w}_{d_p-1}^T \mathbf{x})]^T, \end{aligned} \quad (6)$$

where $\mathbf{w}_0, \dots, \mathbf{w}_{d_p-1}$ are the i.i.d. (independent and identically distributed) samples from d_e dimensional Gaussian distribution with zero mean and identity covariance matrix ensuring

$$k_{GPA}(i-j; s_p, c_p) \approx \phi_{fr}(\mathbf{i})^T \phi_{fr}(\mathbf{j}) \quad (7)$$

to hold. Furthermore, the GaussiP attention approximates the term $\sum_{l=0}^{n-1} k_{GPA}(i-l; s_p, c_p)$ in the closed form $n/\hat{\beta}_2 s_p \|\mathbf{W}_p \mathbf{p}\|_2$. Therefore, the output feature vectors $\hat{\mathbf{y}}_i$ s, of the resultant GaussiP attention,

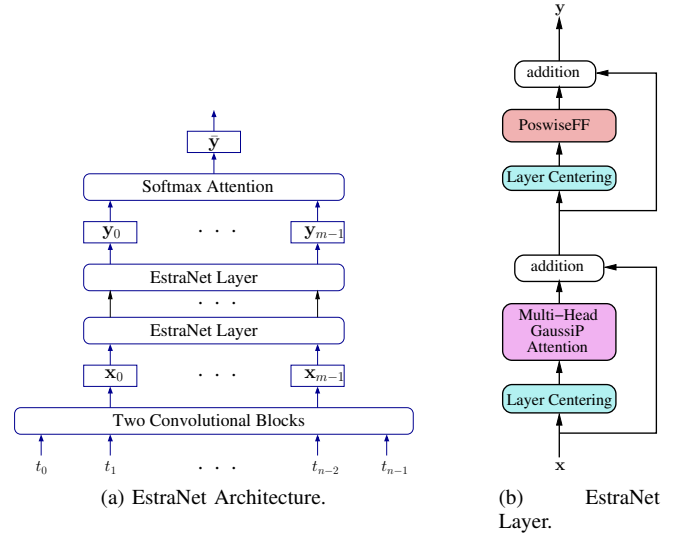


Figure 2. Schematic of EstraNet [16] Architecture.

$$\hat{\mathbf{y}}_i^T = \frac{\sum_{j=0}^{n-1} k_{GPA}(i-j; s_p, c_p) \phi(\mathbf{x}_j)^T}{\sum_{l=0}^{n-1} k_{GPA}(i-l; s_p, c_p)}, \text{ following Eq. (5),} \quad (8)$$

$$\begin{aligned} &\approx \frac{\hat{\beta}_2 s_p \|\mathbf{W}_p \mathbf{p}\|_2}{n} \sum_{j=0}^{n-1} (\phi_{fr}(\mathbf{i})^T \phi_{fr}(\mathbf{j})) \phi(\mathbf{x}_j)^T \\ &\approx \frac{\hat{\beta}_2 s_p \|\mathbf{W}_p \mathbf{p}\|_2}{n} \left[\phi_{fr}(\mathbf{i})^T \sum_{j=0}^{n-1} \phi_{fr}(\mathbf{j}) \phi(\mathbf{x}_j)^T \right], \end{aligned} \quad (9)$$

for all $i = 0, \dots, n-1$, can be computed in linear time and memory with respect to n .

One significant advantage of the GaussiP attention is that the i -th output feature puts high attention around $i - j - c_p n = 0$ or $j = i - c_p n$, allowing information flow from features near $i - c_p n$ to feature i . Since $i - c_p n$ can be far away from i (by setting $c_p \in (0, 1)$ properly), GaussiP attention allows information flow from one region of the traces to a distant region, enabling integration of the leakages of different regions of the traces.

3.3. Multi-head GaussiP Attention

In implementations protected by masking countermeasures, the intermediate secret variable is split into multiple shares such that any proper subset of the shares is independent of the unmasked secret. Therefore, to successfully attack such implementations, one needs to combine the information from the leakages of all the shares. However, in many masked implementations (e.g., software implementations), different secret shares leak at different regions (which may be long distances apart from each other) of the traces. Thus, to attack such implementations, the DL models must be able to combine the leakages from distant regions.

Although the vanilla GaussiP attention facilitates information flow from features near $i - c_p n$ to feature i , this

TABLE 1. THE T_{GE1} (LOWER IS BETTER) VALUES OBTAINED FROM THE ESTRANET MODEL ON TRACES WITH INCREASING TRACE LENGTH. FOR EACH EXPERIMENT, THE MODEL HAS BEEN TRAINED THRICE. THE TABLE SHOWS THE T_{GE1} FROM ALL THREE TRAINING RUNS.

Dataset	Window size 10K			Full-length Traces			Double-length Traces		
	run 1	run 2	run 3	run 1	run 2	run 3	run 1	run 2	run 3
ASCADf	13	9	15	6	13	30	> 5K	> 5K	> 5K
ASCADr	4	6	5	> 5K	> 5K	> 5K	> 5K	> 5K	> 5K
CHES20	11	4	7	5	6	5	> 5K	> 5K	> 5K

flow may be too restricted to effectively combine leakages from multiple regions. In EstraNet, Hajra et al. solve this problem using multi-head GaussiP attention. More precisely, an H -head GaussiP attention, where H is a positive integer, uses H parallel GaussiP attentions to compute the output sequence. Therefore, the output of a H -head GaussiP attention is given by

$$\hat{\mathbf{y}}_i = \sum_{h=0}^{H-1} W_o^{(h)} \hat{\mathbf{y}}_i^{(h)}, \quad \text{where}$$

$$\hat{\mathbf{y}}_i^{(h)} = \sum_{j=0}^{n-1} \frac{k_{GPA}(i-j; s_p^{(h)}, c_p^{(h)})}{\sum_{l=0}^{n-1} k_{GPA}(i-l; s_p^{(h)}, c_p^{(h)})} \phi(\mathbf{x}_j) \quad (10)$$

with $\hat{\mathbf{y}}_i^{(h)}$ being the i -th output feature and $s_p^{(h)} \in \mathbb{R}^+$, $c_p^{(h)} \in (0, 1)$, $W_o^{(h)} \in \mathbb{R}^{d \times d_o}$ being three parameters of h -th head and $\phi(\cdot)$ being a vector valued function with output dimension d_o . The h -th head, for $h = 0, 1, \dots, H-1$, of the attention propagates information from trace region near $(i - c_p^{(h)}n)$ to the i -th feature. Therefore, the H -head GaussiP attentions facilitates information flow from H different regions, $(i - c_p^{(0)}n), \dots, (i - c_p^{(H-1)}n)$, to the i -th feature. If any of these regions contain PoIs, information from those PoIs is combined near i -th output feature.

4. Performance Degradation of EstraNet with Increasing Trace Lengths

4.1. Limitation of GaussiP Attention

The performance of GaussiP attention deteriorates drastically as the size of the attack window increases gradually. To illustrate the point, we recall from Eq. (8) that an output feature of the h -th head of an H -head GaussiP attention is computed as

$$\hat{\mathbf{y}}_i^{(h)} \approx \sum_{j=0}^{n-1} \frac{k_{GPA}(i-j; s_p^{(h)}, c_p^{(h)})}{\sum_{l=0}^{n-1} k_{GPA}(i-l; s_p^{(h)}, c_p^{(h)})} W_o^{(h)} \phi(\mathbf{x}_j)$$

$$= \sum_{j=0}^{n-1} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j) \quad (11)$$

where $w_j^{(h)} \triangleq \frac{k_{GPA}(i-j; s_p^{(h)}, c_p^{(h)})}{\sum_{l=0}^{n-1} k_{GPA}(i-l; s_p^{(h)}, c_p^{(h)})}$ is the weight for the j -th input feature \mathbf{x}_j and $\phi^{(h)}(\cdot) \triangleq W_o^{(h)} \phi(\cdot)$. Note that since k_{GPA} is a Gaussian kernel, it puts high attention to a small contiguous segment, say $\mathcal{R}^{(h)} = [a, b]$, of the input sequence where $0 \leq a < b < n$. Within this high-attention segment $\mathcal{R}^{(h)}$, a smaller segment, say

$\mathcal{S}^{(h)} \subseteq \mathcal{R}^{(h)}$, corresponds to some high SNR features. Thus, we can re-write Eq. (11) as

$$\hat{\mathbf{y}}_i^{(h)} \approx \sum_{j \in \mathcal{S}^{(h)}} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j) + \sum_{j \in \mathcal{R}^{(h)} \setminus \mathcal{S}^{(h)}} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j) + \sum_{j \notin \mathcal{R}^{(h)}} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j) \quad (12)$$

When $\hat{\beta}_2$ hyper-parameter of k_{GPA} remains constant, the number of high attention features, i.e., $|\mathcal{R}^{(h)}|$, increases almost linearly with the increase of the attack window size n (as explained in Appendix A). However, among this $\mathcal{R}^{(h)}$ features, the high SNR features $\mathcal{S}^{(h)}$ remain almost the same. Consequently, the norm of the signal part $\|\sum_{j \in \mathcal{S}^{(h)}} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j)\|_2$ within $\hat{\mathbf{y}}_i^{(h)}$ remains almost the same while the norm of the noise part, i.e., $\|\sum_{j \in \mathcal{R}^{(h)} \setminus \mathcal{S}^{(h)}} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j) + \sum_{j \notin \mathcal{R}^{(h)}} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j)\|_2$ increases with the attack window size n , reducing the effective SNR of $\hat{\mathbf{y}}_i^{(h)}$. One way to alleviate the above problem is to increase the value of $\hat{\beta}_2$ with the increase in the attack window size, reducing the size of high attention features $|\mathcal{R}^{(h)}|$ close to the number of high SNR features $|\mathcal{S}^{(h)}|$. However, increasing the $\hat{\beta}_2$ hyper-parameter reduces the size of $|\mathcal{R}^{(h)}|$ relative to the attack window size n , making it difficult for the H -head GaussiP attention to combine information from multiple PoIs. Alternatively, one can increase both $\hat{\beta}_2$ and the number of heads, H , of the multi-head GaussiP attention. The larger value of $\hat{\beta}_2$ helps to prevent the deterioration of the effective SNR at the outputs of the GaussiP attention, and the larger value of H increases the flow of information through the EstraNet layers, increasing the model's ability to combine leakages from multiple PoIs. However, the approach has several drawbacks:

- 1) EstraNet might require significantly more tuning of several hyper-parameters like $\hat{\beta}_2$ and H , increasing compute cost of training significantly.
- 2) Increasing the hyper-parameter H increases the model's number of parameters, increasing the model's training and inference cost. Moreover, having more parameters increases the risk of overfitting, significantly increasing the required training traces to train the model effectively.
- 3) Since the GaussiP attention is only an approximation of the Gaussian kernel, the weights, $w_j^{(h)}$ s, for the features in the set $\{j : j \notin \mathcal{R}^{(h)}\}$ in Eq. (12) are not exactly zeros, but some small values. When the size of the attack window is very large and the attentions are focused on very small

regions (by setting $\hat{\beta}_2$ to a high value), the set $\{j : j \notin \mathcal{R}^{(h)}\}$ becomes very large, increasing the component $\sum_{j \notin \mathcal{R}^{(h)}} w_j^{(h)} \phi^{(h)}(\mathbf{x}_j)$ in Eq. (12) to a significant value. That, in turn, reduces the effective SNR of the GaussiP attention’s outputs, making the training of EstraNet more difficult for the longer attack window sizes.

Next section empirically verifies the above arguments.

4.2. Experimental Evaluation of EstraNet on Longer Traces

To experimentally verify the above arguments, we investigate the effectiveness of EstraNet with increasing trace lengths. Towards that goal, we trained and evaluated EstraNet on three datasets: ASCAD fixed key (ASCADf), ASCAD random key (ASCADr), and CHES CTF 2020 (CHES20) datasets (see Section 6.1). These datasets respectively contain $100K$, $250K$, and $62.5K$ features. We conducted three sets of experiments on each of the datasets. In the first set, we selected an attack window of $10K$ from each dataset. In the second set, we trained and evaluated the model on the full-length traces. In the third set, we doubled the length of each trace by concatenating the trace with itself once, and then performed experiments on these double-length traces. The results are summarized in Table 1. Each experiment is repeated three times. The number of attack traces required to reach guessing entropy 1, denoted by T_{GE1} , in the three training runs of each experiment are shown in the table. From the table, it can be observed that for the attack window of $10K$, EstraNet requires close to 10 attack traces to reach the guessing entropy 1 across all three datasets. When attacks are performed on the full-length traces, EstraNet fails to reach the same using as many as $5K$ attack traces in one dataset (ASCADr dataset). However, when attacks are performed on the double-length traces, EstraNet is unable to reach the same using less than $5K$ traces in all datasets.

Following the arguments of the previous section, we tune both $\hat{\beta}_2$ and H hyper-parameters to improve EstraNet’s performance on longer traces. Therefore, we trained it on the ASCADf dataset using double-length traces, varying combinations of $\hat{\beta}_2$ and H . The results are summarized in Table 2. From the table, we observe that the best performance from EstraNet is achieved with $\hat{\beta}_2 = 300$ and $H = 16$. For this combination, the model can reach the guessing entropy 1 using less than $5K$ attack traces in only one out of the three training runs. This indicates that even with substantial hyper-parameter tuning, EstraNet may not consistently achieve good performance on longer traces.

It is worth noting that although EstraNet struggles on longer traces, it performs well with smaller attack windows. However, in many practical applications, such as those involving implementations protected by masking or jitter-based countermeasures, selecting sufficiently smaller attack windows from full-length traces may not be feasible. In such scenarios, applying EstraNet can be challenging. These

limitations are addressed by HierNet. Table 3 demonstrates consistent performance of HierNet even for double-length traces. Next section details the HierNet architecture.

5. HierNet Architecture

The core idea behind HierNet is to process a long trace in segmented manner using a shift-invariant model like EstraNet and then combine the output from all segments using a top-level DL model (as shown in Figure 1). This approach has several advantages. First, shift-invariant models like EstraNet may not effectively work on the entire trace but can perform very well on relatively smaller segments. Therefore, applying such models in a segmented manner retains its effectiveness against countermeasures like masking, random delay, and clock jitter while addressing its limitations with increasing trace length. Second, the top-level DL model can be designed differently from EstraNet, considering different criteria. For instance, shift-invariance of EstraNet enables it to better mitigate the effects of random delay and clock jitter compared to other models, as demonstrated in [16]. Therefore, using a shift-invariant model as the base model can enable HierNet to deal with such countermeasures better. However, once each segment of the traces is processed through the base model in HierNet, other positional encodings (e.g., absolute positional encoding) can be incorporated into the top-level DL model to exploit information like segment ordering. In fact, absolute positional encoding is integrated into the top-level DL model of HierNet, leveraging the fact that the random shifts of features caused by countermeasures like random delay and clock jitter are bounded (more details in Section 5.2).

We summarize the approach used in HierNet more precisely. In HierNet, the input trace \mathbf{t} is first segmented into multiple relatively smaller overlapping segments, denoted as $\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_{m-1}$. Each segment \mathbf{s}_i is independently processed by EstraNet models, producing outputs $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{m-1}$, where $\mathbf{x}_i \in \mathbb{R}^d$ represents the output of EstraNet applied to segment \mathbf{s}_i . Next, a top-level DL model (a single transformer layer) generates the final output by taking $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{m-1}\}$ as its input. Several design factors in HierNet are crucial, including the segmentation strategy, the design of the top-level model, and parameter sharing among EstraNet models for different segments. We discuss these factors in detail below.

5.1. Segmentation

HierNet inherits the robustness of EstraNet against random delay and clock jitter countermeasures by employing EstraNet models to process each trace segment. However, segmenting the input trace in HierNet can potentially compromise this robustness. For example, misalignments in traces may cause some PoI near a segment boundary to oscillate between segments. In one trace, the PoI might belong to one segment, while in another trace, it could be shifted to a different one, resulting in the PoI not being processed in any segment. Such issue can be mitigated by

TABLE 2. THE T_{GE1} (LOWER IS BETTER) VALUES OBTAINED FROM THE ESTRANET MODEL ON THE ASCADF DATASET WITH DOUBLE-LENGTH TRACES FOR VARIOUS COMBINATIONS OF H AND $\hat{\beta}_2$. FOR EACH HYPER-PARAMETER, THE MODEL HAS BEEN TRAINED THRICE. THE TABLE SHOWS THE T_{GE1} FROM ALL THREE TRAINING RUNS.

	$H = 8$			$H = 12$			$H = 16$		
	run 1	run 2	run 3	run 1	run 2	run 3	run 1	run 2	run 3
$\hat{\beta}_2 = 150$	> 5K	> 5K	> 5K	> 5K	> 5K	> 5K	> 5K	> 5K	> 5K
$\hat{\beta}_2 = 300$	> 5K	> 5K	> 5K	> 5K	32	> 5K	> 5K	> 5K	7

TABLE 3. THE T_{GE1} (LOWER IS BETTER) VALUES OBTAINED FROM HIERNET MODEL ON TRACES WITH INCREASING TRACE LENGTH.

Dataset	Full-length Traces			Double-length Traces		
	run 1	run 2	run 3	run 1	run 2	run 3
ASCADf	2	2	2	2	2	2
ASCADr	3	2	3	3	5	5
CHES20	2	2	2	3	3	3

making the segment length reasonably large. Additionally, we ensure significant overlap between consecutive segments. Specifically, if one segment ends at feature t , the next segment starts from feature $t - o$, ensuring that the span from $t - o$ to t belongs to both segments. It’s important to note that if the displacements are bounded by length o , each PoI will always remain within one of the two consecutive segments. However, if displacements exceed o , increasing its value (by setting a hyper-parameter) can enhance HierNet’s robustness against larger desynchronizations, albeit at the cost of increased computational overhead.

5.2. Designing the Top-level DL Model

When designing the top-level DL model, we primarily consider three criteria. First, we aim to ensure that the model’s parameterization remains independent of the trace lengths and segmentations. The number of segments, denoted as m , can vary depending on the trace length and chosen segment length, thereby affecting the input length to the top-level model. Therefore, we select a model whose parameterization is independent of the input length. Second, the model should effectively capture long-distance dependencies, enabling it to combine leakages from all input elements using a minimal number of layers. Finally, we require the model to easily incorporate absolute positional encoding. Based on these criteria, we selected a transformer layer as the top-level DL model, as it meets all three requirements. This model includes a self-attention and a position-wise feed-forward layer, as described in Section 2.4. We used a softmax self-attention operation and incorporated absolute positional encoding into it. A brief description of this operation is provided in the following section.

5.2.1. Incorporating Absolute Positional Encoding into Self-attention Layer. In the presence of random delay or clock jitter countermeasure, the displacements of PoIs across the traces are likely to remain bounded within a segment. Therefore, if the leakage related to some intermediate variable appears in the i -th segment for some trace, its position

is unlikely to change for the other traces. In other words, the positions of the informative features of the input sequence $\{\mathbf{x}_0, \dots, \mathbf{x}_{m-1}\}$ will remain the same across all traces. To exploit such phenomenon, we incorporate the absolute positional encoding in the self-attention layer of the top-level DL model.

More precisely, in vanilla self-attention (the softmax self-attention in Eq. 4), the i -th output feature $\hat{\mathbf{y}}_i$, $i = 0, 1, \dots, m - 1$, is computed as $\hat{\mathbf{y}}_i = \sum_{j=0}^{m-1} p_{ij} \phi(\mathbf{x}_j)$ with p_{ij} , the attention weight from the i -th feature to the j -th feature, is given as

$$p_{ij} = \text{softmax} \left(\frac{\langle f(\mathbf{x}_i), g(\mathbf{x}_j) \rangle}{\sqrt{d_k}} \right)_j = \frac{\exp(\langle f(\mathbf{x}_i), g(\mathbf{x}_j) \rangle / \sqrt{d_k})}{\sum_{l=0}^{n-1} \exp(\langle f(\mathbf{x}_i), g(\mathbf{x}_l) \rangle / \sqrt{d_k})},$$

where $f(\cdot), g(\cdot)$ and $\phi(\cdot)$ are three vector valued functions of the input feature \mathbf{x}_i ($i = 0, \dots, m - 1$), $\langle \cdot, \cdot \rangle$ denotes the dot product between two vectors and d_k be the dimension of the outputs of f (or g). We incorporate the absolute positional encoding into the computation of attention probabilities as follows:

$$p_{ij}^a = (1 - \gamma) \text{softmax} \left(\frac{\langle f(\mathbf{x}_i), g(\mathbf{x}_j) \rangle}{\sqrt{d_k}} \right)_j + \gamma \text{softmax} \left(\frac{\langle \mathbf{v}_i, \mathbf{v}_j \rangle}{\sqrt{d_k}} \right)_j \quad (13)$$

$$= (1 - \gamma) \frac{\exp(\langle f(\mathbf{x}_i), g(\mathbf{x}_j) \rangle / \sqrt{d_k})}{\sum_{l=0}^{n-1} \exp(\langle f(\mathbf{x}_i), g(\mathbf{x}_l) \rangle / \sqrt{d_k})} + \gamma \frac{\exp(\langle \mathbf{v}_i, \mathbf{v}_j \rangle / \sqrt{d_k})}{\sum_{l=0}^{n-1} \exp(\langle \mathbf{v}_i, \mathbf{v}_l \rangle / \sqrt{d_k})}, \quad (14)$$

where $\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{m-1}$ are m distinct vectors and $\gamma \in [0, 1]$ is a hyper-parameter. The values of vectors $\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{m-1}$ are learned along with other parameters of the model during training. Note that the first component of the R.H.S. of Eq. 13 or 14 computes the attention probability based on the similarity (or inter-dependency) between the two features, \mathbf{x}_i and \mathbf{x}_j . The second component computes the probability based on their absolute positions or indices i, j . The hyper-parameters γ controls the relative weights of the two components in the attention probability p_{ij}^a .

5.3. Parameter Sharing among EstraNet Models

To keep HierNet’s model size small and to prevent it from overfitting, the parameters of all EstraNet models are shared. However, during our initial experiments, we observed that sharing the parameters of the last softmax layer of the EstraNet models resulted in a deceleration of training convergence for HierNet. We hypothesize that the parameter-sharing of the softmax attention layer was forcing each segment’s output towards a unified representation, thereby slowing the training progress. Hence, we opted not to share the parameters of last layer of the EstraNet models.

The overall HierNet architecture is shown in Figure 1. Next section experimentally evaluates HierNet.

6. Experimental Results

In this section, HierNet has been evaluated on three SCA datasets of masked implementations, with its performance compared to three benchmark models. We also present ablation studies of HierNet and compare its training time with that of the benchmark models.

6.1. Datasets Used

To carry out the experimental evaluation, we used two datasets of first-order masked AES implementation, namely ASCAD Fixed Key¹ (ASCADf) and ASCAD Random Key² (ASCADr) datasets, and a second-order masked implementation of Clyde-128 Tweakable Block Cipher, namely CHES 2020 CTF dataset³ (CHES20). The first two datasets were collected from an 8-bit ATmega8515 microcontroller, whereas the third dataset was collected from ARM Cortex-M0 microcontroller. We target first round’s third sbox for the first two datasets and 17th column from the right of the state matrix after the first subbyte operation for the third dataset. The datasets’s details are given in Appendix B. Table 4 summarizes the dataset statistics.

TABLE 4. STATISTICS OF THE DATASETS

	ASCADf	ASCADr	CHES20
# of Profiling traces	50K	200K	200K
# of Validation traces	5K	10K	10K
# of Attack traces	5K	10K	10K
Trace Length	100K	250K	62.5K
Mask order	1	1	2

6.2. Benchmark Models

We compare the performance of HierNet with three benchmark models, referred to as EffCNN [19], PolyCNN [12], and LSTMNet [13]. The EffCNN models are CNN models with a flattening layer, whereas the PolyCNN model is a CNN model with a global average-pooling layer. On the other hand, LSTMNet uses softmax attention after a bi-directional LSTM layer to accumulate information spread over a long trace. Further details of the models can be found in Appendix C. Since EstraNet is already found to perform poorly on longer traces in Section 4.2, we do not compare it with HierNet in this section.

6.3. Details of Hyper-parameters and Training

HierNet uses EstraNet model as the first-level DL model in its architecture. The hyper-parameters related to EstraNet in HierNet are adopted from [16], with a few exceptions.

1. https://github.com/ANSSI-FR/ASCAD/tree/master/ATMEGA_AES_v1/ATM_AES_v1_fixed_key
2. https://github.com/ANSSI-FR/ASCAD/tree/master/ATMEGA_AES_v1/ATM_AES_v1_variable_key
3. <https://ctf.spook.dev/>

Notably, following the recommendation of [16], the $\hat{\beta}_2$ hyper-parameter of the EstraNet model has been tuned for each dataset. In the top-level DL model, we use absolute positional encoding with γ hyper-parameter set to 1. We set the number of heads of the self-attention layers to 8, though setting it to some other value like 1 also provides similar results. HierNet introduces two new hyper-parameters: segment length and the overlap length between two consecutive segments. Across all datasets, the segment length has been set to $5K$, and the overlap length is set to $2K$.

We train HierNet for $4M$ training steps utilizing the Adam optimizer. Following the previous literature [16], [18] of TN for SCA, it has been trained using a cosine decay with a linear warmup learning rate schedule. More precisely, the learning rate is linearly increased from 0 to the maximum value of $2.5e-4$ over t_{warmup} training steps. Subsequently, it is gradually reduced to $0.004 \times 2.5e-4$ over the remaining $4M - t_{warmup}$ steps, following a cosine learning rate schedule. The default value for t_{warmup} is set to $1M$; however, it was observed that, sometimes, the model did not train effectively for this default value. In such cases, t_{warmup} was increased to either of $2M$ or $4M$.

6.4. Comparison with Benchmark Models in the Presence of Random Delay on Long Traces

In this section, we compare HierNet with the benchmark models in the presence of random delay on long traces. To achieve this, we conduct experiments on full-length traces from the three datasets, which are $100K$, $250K$, and $62.5K$ features long, respectively, with added random delay. The comparisons are carried out in two parts. In the first part, we use all available training traces from the datasets to train the models. In the second part, we compare their performance when fewer training traces are used to train them.

Experimental Setup. To introduce random delay effect, we added random shifts in the traces. In the profiling traces, the random delay has been introduced using a profiling desync of 600. In other words, we shifted each profiling traces by a random displacement in the range $[0, 600)$. The data augmentation has been performed using a desync of 400. Evaluations of the trained models were conducted for two attack desyncs: 600 and 1000. We use the minimum number of attack traces required to reach the guessing entropy 1, denoted by T_{GE1} , as the performance metric for the DL models.

Results with All Training Traces. We begin by comparing HierNet with the benchmark models, utilising all profiling traces from the respective datasets to train the model. Therefore, we utilize $50K$ traces from the ASCADf dataset, and $200K$ traces each from the ASCADr and CHES20 datasets for training. The results are detailed in Table 5. From the table, it becomes apparent that HierNet exhibits notable performance, requiring approximately 2 to 3 traces to reach the guessing entropy 1 across all datasets. Conversely, performance of the PolyCNN model is poor across the datasets, with relatively better performance on ASCADr dataset. Similarly, the EffCNN model exhibits

TABLE 5. THE T_{GE1} VALUES (LOWER IS BETTER) FOR DIFFERENT DL MODELS AGAINST THE RANDOM DELAY COUNTERMEASURE ON FULL-LENGTH TRACES. THE COLUMNS LABELED *Best*, *Med.*, AND *Avg.* RESPECTIVELY DISPLAY THE BEST, MEDIAN, AND AVERAGE T_{GE1} VALUES COMPUTED FROM THREE INDEPENDENT TRAINING TRIALS FOR EACH MODEL. ENTRIES MARKED WITH ‘-’ INDICATE THAT THE AVERAGE VALUE IS UNAVAILABLE BECAUSE THE MODEL FAILED TO REACH GUESSING ENTROPY 1 USING 5K ATTACK TRACES IN SOME OF THE TRAINING RUNS.

Dataset	Model	Attack Desync 600			Attack Desync 1000		
		Best	Med.	Avg.	Best	Med.	Avg.
ASCADf	PolyCNN	2626	> 5K	–	2102	> 5K	–
	EffCNN	2854	> 5K	–	2812	> 5K	–
	LSTMNet	157	208	1135.7	146	248	972.7
	HierNet	2	2	2.3	2	2	2.0
ASCADr	PolyCNN	5	286	711.7	4	326	445.7
	EffCNN	–	–	–	–	–	–
	LSTMNet	4	7	8.3	5	6	9.0
	HierNet	2	2	2.3	2	3	2.7
CHES20	PolyCNN	689	> 5K	–	1946	> 5K	–
	EffCNN	–	–	–	–	–	–
	LSTMNet	2	4	3.7	2	2	4.0
	HierNet	1	2	2.0	2	2	2.0

poor performance across all datasets. Specifically, our investigation revealed that the training of the model on the ASCADr and CHES20 datasets progresses too slowly to reduce the training loss significantly below the level of random guessing, even after 60 hours of training (as illustrated in Appendix D). In contrast, HierNet models require only 24 to 40 hours for the complete training. The LSTMNet model performs significantly worse on the ASCADf dataset, although its performance is comparable to HierNet on the ASCADr and CHES20 datasets.

Comparison with Reduced Number of Training Traces. We perform further evaluation using a reduced number of profiling traces to train the models. Specifically, we compare HierNet and LSTMNet on the ASCADr and CHES20 datasets, utilizing a random subset of 50K profiling traces from the datasets’ total of 200K profiling traces to train the models. Note that, due to their poor performance on the full training data, PolyCNN and EffCNN models are not included in this comparison. Additionally, since LSTMNet performs poorly compared to HierNet with the full training data of the ASCADf dataset, we do not compare them again using a lesser number of training traces from that dataset. The results are presented in Table 6. From the table, it is evident that HierNet models can still reach the guessing entropy 1 using only three or fewer attack traces in both datasets. In contrast, the LSTMNet models necessitate over 2K attack traces on the ASCADr dataset to reach the guessing entropy 1. On the CHES20 dataset, though LSTMNet model reaches guessing entropy 1 using fewer than 5 attack traces in one training run, it fails to reach the same using less than 5K traces in the remaining two runs.

In summary, PolyCNN and EffCNN exhibit poor performance on longer traces. While LSTMNet shows more resilience to increasing trace length, achieving comparable performance to HierNet may require significantly more training data, which clearly demonstrate HierNet’s advantage over the benchmark models. The following section further compares HierNet with the benchmark models against the clock jitter countermeasure.

6.5. Comparison with Benchmark Models in the Presence of Clock Jitter Effect

An advantage in attacking countermeasures involving global trace desynchronizations, such as in many random delay countermeasures, lies in the ease with which these effects can be simulated during the training phase. Consequently, data augmentation can be readily applied, increasing the effective training data size and enhancing the robustness of the DL models to such misalignments. In contrast, the clock jitter countermeasure introduces small, local desynchronization throughout the full traces, leading to significant variation in distances between peaks within and across traces. While various data augmentation techniques are available for simulating the clock jitter effect [10], these approaches result in significant information loss due to the removal of informative features during the processing. Consequently, relying solely on data augmentation during training of the DL models to perform effectively against the clock jitter countermeasure may be infeasible.

This section assesses HierNet’s performance against clock jitter by conducting attacks after adding the clock jitter effect into traces. Note that prior studies, such as [20], demonstrated successful attacks against clock jitter effect under the assumption of a stronger adversary who has access to both clean version and noisy version of the traces. However, as in [10], [16], we consider a weaker adversary who lacks the access to the clean version of the noisy traces.

Generation of Clock Jitter Effect. We employed the methodology outlined in [10], [16], [20] to introduce the clock jitter effect into the traces. Specifically, we generated a noisy trace for each trace of the dataset using the following procedure. For every feature in the original trace, we performed one of the following three actions with equal probabilities: a) excluding the feature from the noisy trace, b) including the feature in the noisy trace, and c) including the feature along with an additional feature, where the additional feature is the average of the current and following features in the source trace. Notably, this method

TABLE 6. THE T_{GE1} VALUES (LOWER IS BETTER) FOR THE LSTMNET, AND HIERNET MODELS WHEN THE MODELS ARE TRAINED USING 50K PROFILING TRACES. THE COLUMNS LABELED *Best*, *Med.*, AND *Avg.* RESPECTIVELY DISPLAY THE BEST, MEDIAN, AND AVERAGE T_{GE1} VALUES COMPUTED FROM THREE INDEPENDENT TRAINING TRIALS FOR EACH MODEL. ENTRIES MARKED WITH ‘-’ INDICATE THAT THE AVERAGE VALUE IS UNAVAILABLE BECAUSE THE MODEL FAILED TO REACH THE GUESSING ENTROPY 1 USING 5K ATTACK TRACES IN SOME TRAINING RUNS.

Dataset	Model	Attack Desync 600			Attack Desync 1000		
		Best	Med.	Avg.	Best	Med.	Avg.
ASCADr	LSTMNet	2987	4893	–	2615	2812	–
	HierNet	2	2	2.3	2	3	2.7
CHES20	LSTMNet	2	> 5K	–	4	> 5K	–
	HierNet	1	2	1.7	2	2	2.0

aligns with the application of the clock jitter effect using the algorithm proposed by Wu et al. [20], with the parameter `clock_jitters_level` set to 1. It is pertinent to mention that Hajra et al. [16] also employed the same algorithm for introducing the clock jitter effect. However, prior to applying the clock jitter effect, they doubled the trace length by duplicating each feature in the source trace twice to avoid the loss of any informative features. In contrast, we do not perform such preprocessing. Therefore, our attack setup is more challenging compared to that of [16].

TABLE 7. SELECTED ATTACK WINDOWS FOR THE EXPERIMENTS WITH ADDED CLOCK JITTER EFFECT.

Attack Window	Dataset		
	ASCADf	ASCADr	CHES20
[40K, 50K]	[78K, 88K]	[46K, 56K]	

Increasing the Number of Traces in the Training Data. In Section 6.4, it was observed that the performance of the DL models improves with the increase of the number of training traces, although HierNet achieves its optimal performance with a significantly smaller number of training traces. To investigate whether the DL models can also be robust to clock jitter by using more training traces, we conduct experiments with an increased number of training traces. More precisely, the number of training traces for a dataset was increased to x -fold as follows: for each trace in the clean dataset, x noisy traces were added to the noisy dataset. Each of the x noisy traces was independently created by introducing clock jitter effects to the corresponding clean trace. Since the process adds x noisy traces in the noisy dataset for each trace in the clean dataset, it increases the size of the noisy dataset x -fold. Note that increasing the number of training traces in this manner does not increase the information content of the dataset, as multiple traces are generated from the same clean traces. However, creating multiple traces by adding clock jitter effects independently to the same clean trace should help the DL models become robust to the variations introduced by the clock jitter effect.

Selection of Attack Window. To reduce the compute cost, we did not conduct experiments on the full-length traces. Instead, the experiments were carried out on attack windows of size 10K. It should be noted that, earlier, Hajra et al. [16] also conducted attacks on selected attack windows with a length of 10K. We adopt the same attack windows for the experiments in this section. The details of these

TABLE 8. THE T_{GE1} VALUES (LESSER IS BETTER) FOR DIFFERENT DL MODELS ON THE THREE DATASETS WITH ADDED CLOCK JITTER EFFECT. WE HAVE PERFORMED THE ATTACKS ON A SELECTED ATTACK WINDOW OF SIZE 10K FOR THESE EXPERIMENTS.

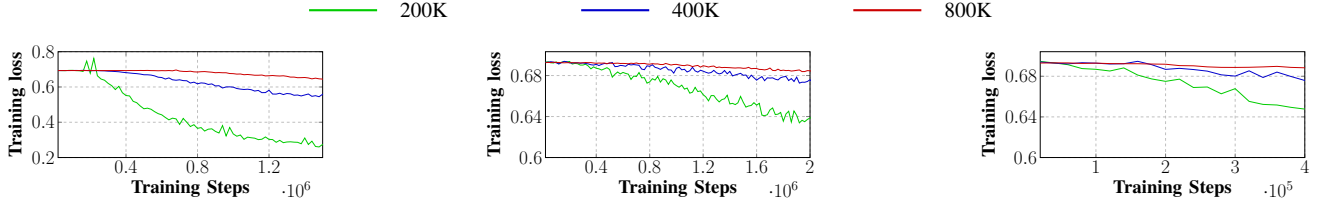
Dataset	Training Data Size	Model			
		PolyCNN	EffCNN	LSTM	HierNet
ASCADf	100K	> 5K	> 5K	> 5K	55
	200K	544	384	> 5K	23
	400K	442	60	197	22
	800K	262	34	160	17
ASCADr	200K	> 5K	> 5K	> 5K	20
	400K	> 5K	88	> 5K	16
	800K	25	40	23	13
CHES20	200K	> 5K	> 5K	> 5K	> 5K
	400K	> 5K	> 5K	> 5K	10
	800K	> 5K	> 5K	> 5K	12

attack windows are outlined in Table 7. It is important to note that the selection of the attack windows assumes the feasibility of SNR-based window selection, which may not hold in many practical attack scenarios. However, we emphasize that window selection was performed to reduce compute cost. We hope that the conclusions drawn from these experiments will also apply to full-length traces, albeit with a different scaling of the number of training traces.

Training Details. All benchmark models have been trained as outlined in Appendix C. The HierNet models were trained according to the procedure detailed in Section 6.3. As in Sections 6.4, data augmentation was applied to all models during training. Specifically, each training trace was shifted by a random displacement within the range $[0, 200)$ on-the-fly during the training.

Results. The attack results of the DL models under different numbers of training traces are presented in Table 8. Similar to Sections 4.2 and 6.4, the table uses T_{GE1} – the minimum number of attack traces required to reach the guessing entropy 1 – as the metric for assessing the effectiveness of the models. In these experiments also, each DL model was independently trained three times, and the table displays the median value of the three results.

In the case of the ASCADr dataset, the results of the DL models are presented for three distinct sizes of the training data: 200K, 400K, and 800K. As evident from the table, with 200K training traces, HierNet reaches guessing entropy 1 with around 20 traces, whereas other models fail to reach the same even using 5K attack traces. As the number of training traces increases, the other models can reach



(a) Training loss of the LSTMNet model

(b) Training loss of the PolyCNN model

(c) Training loss of the EffCNN model

Figure 3. Improvement in the training loss of the LSTMNet, PolyCNN, and EffCNN models during their training on the CHES20 dataset with the added clock jitter effect. The green, blue, and red lines, respectively, show the progress of the training loss for the models’ training with 200K, 400K, and 800K training traces.

guessing entropy 1. However, they require four times more training traces to achieve performance close to HierNet.

For the ASCADf dataset, the presented table displays attack results across four distinct training data sizes: 100K, 200K, 400K, and 800K. With 100K training traces, only HierNet can reach the guessing entropy 1 with fewer than 5K attack traces. As the number of training traces increases to 800K, the performance of EffCNN approaches that of HierNet. Conversely, the PolyCNN and LSTMNet models fail to reach performance comparable to HierNet, even with eight times more training traces.

The CHES20 dataset poses a challenging scenario for the benchmark models, though HierNet performs well. Specifically, HierNet achieves its optimal performance (reaching the guessing entropy 1 using approximately 10 attack traces) with only 400K training traces. However, the CNN and LSTM-based models fail to reach the same even using 5K attack traces. Furthermore, even when the models are trained with more training traces, their performance does not improve. In fact, in our experiments, we observe that as the number of training traces increases, these models become challenging to train, as indicated in Figure 3. We speculate that these models struggle to learn meaningful information from the traces during the training. Therefore, when trained with a smaller number of training traces, the models tend to overfit the noise, leading to a relatively faster decrease in training loss. However, more training traces mitigate overfitting, but the model remains untrainable.

In conclusion, the results suggest that HierNet demonstrates robust performance against clock jitter countermeasures across all three datasets. In contrast, the benchmark models might fail to reach to a comparable performance even when they are trained using eight times more training traces. In some scenarios, HierNet excels whereas all the benchmark models fail to be trained effectively, even with additional training data, highlighting the need for further research to enhance their performance in such contexts.

6.6. Ablation Studies

Ablation Study of Transformer Layer in Top-level DL Model. In this section, we first conduct an ablation study on the transformer layer used in the top-level DL model. To do this, we replace the transformer layer with a Feed-Forward Network (FFN) in the model. The results are

TABLE 9. THE T_{GE1} VALUES (LOWER IS BETTER) ARE SHOWN FOR THE HIERNET MODEL WITH ITS TOP-LEVEL DL MODEL’S TRANSFORMER LAYER REPLACED BY A FEED-FORWARD NETWORK (FFN). THE ATTACK DESYNC 1000 IS USED FOR THE EXPERIMENTS. REMAINING EXPERIMENTAL SETUP IS SAME AS SECTION 6.4.

Dataset	run 1	run 2	run 3
ASCADf	6	2	7
ASCADr	> 5K	> 5K	> 5K
CHES20	> 5K	> 5K	> 5K

shown in Table 9. From the table, we observe that while the modified model performs comparably to the original HierNet model on the ASCADf dataset, it fails to reduce the T_{GE1} below 5K traces on the remaining two datasets. In fact, the training loss of the model did not decrease significantly below the level of random guessing, even after training for the same duration as the original HierNet model. It is worth mentioning that due to computational constraints, we were unable to extensively tune the hyperparameters of the FFN model. Nonetheless, the results in Table 9 suggest that replacing the transformer layer in HierNet’s top-level DL model with another DL model while maintaining strong performance across datasets can be a non-trivial task.

Ablation Study of Absolute Positional Encoding.

The self-attention layer in the top-level DL model of HierNet’s transformer layer incorporates absolute positional encoding. In this section, we examine the performance of HierNet when the absolute positional encoding is ablated from the self-attention layer. The results are shown in Table 10. From the table, we observe that the modified HierNet performs similarly to the original version when trained with a sufficient number of training traces. However, in low training data scenarios, i.e., when trained with 50K training traces, the performance of the modified HierNet becomes significantly unstable. For example, when trained with 50K training traces on the CHES20 dataset, the modified HierNet may require as many as 74 attack traces to reach the guessing entropy 1 in some training runs, whereas the original HierNet with absolute positional encoding achieves this with only around two attack traces (as shown in Table 6).

6.7. Training Time

In our experimental environment (Google Colab TPU), the training time of LSTMNet can be more than five times

TABLE 10. THE T_{GE1} VALUES (LOWER IS BETTER) ARE SHOWN FOR THE HIERNET MODEL WITH THE ABSOLUTE POSITIONAL ENCODING REMOVED FROM THE SELF-ATTENTION LAYER OF THE TOP-LEVEL DL MODEL. THE ATTACK DESYNC 1000 IS USED FOR THE EXPERIMENTS. REMAINING EXPERIMENTAL SETUP IS SAME AS SECTION 6.4.

Dataset	# of Train Traces 200K			# of Train Traces 50K		
	run 1	run 2	run 3	run 1	run 2	run 3
ASCADf	not available			6	2	5
ASCADr	2	3	3	28	11	5
CHES20	3	2	2	3	74	2

slower for traces with lengths around 10K features, while it is about two times slower for longer traces (with lengths in the order of 100K features). The training time of the PolyCNN and EstraNet models varies between 3 to 12 times compared to HierNet on shorter traces, but their training becomes much slower on longer traces. In some cases, the training loss of the EffCNN model did not decrease significantly below the level of random guessing, even after training for more than two times the total duration of HierNet (see Appendix D). While an improved learning rate schedule might accelerate their training on longer traces, we still expect them to be significantly slower than HierNet.

While comparing the training time of HierNet with EstraNet, it’s worth noting that due to the overlapping segmentation in HierNet, its per-step training time is slightly higher than EstraNet’s. Specifically, in our experiments, HierNet’s per-step training time was between 1 and 1.5 times that of EstraNet across the three datasets. However, HierNet requires fewer training steps to converge, making the total training time comparable to EstraNet.

7. Future Works

HierNet against Random Delay Interrupt Countermeasure. The Random Delay Interrupt (RDI) countermeasure inserts random delays into multiple intermediate positions of traces, causing drastic variation in the distances between the PoIs across the traces. Since HierNet segments the traces into multiple segments and it is shift-invariant in each segment, it assumes that the intermediate distances between the PoIs within each segment remain approximately constant. Such an assumption may not be true in the presence of RDI. Therefore, HierNet may not be directly applicable against RDI countermeasures. Notably, various previously proposed pre-processing techniques [20], [21], [22] can be used to mitigate the RDI effects from the traces before training the DL models. It should also be noted that these pre-processing techniques may not completely eliminate the desynchronizations introduced by RDI. The pre-processed traces may still exhibit misalignment, albeit to a lesser extent. Therefore, HierNet, which has demonstrated its effectiveness against traces with bounded desynchronizations, could be a suitable choice for attacking the pre-processed traces. Investigating this aspect could be a future research direction.

HierNet against Shuffling Countermeasure. In the shuffling countermeasure [11], the order of execution of

multiple parallel operations of a cryptographic algorithm (e.g., sub-byte operations of different bytes in an AES encryption) is randomly shuffled, resulting in different permutations in the positions of leakages of the respective operations across traces. This shuffling might cause the relative distances between the PoIs to vary significantly, posing a challenge to the usefulness of DL models like EstraNet on such traces. However, HierNet, which segments entire traces into multiple segments and processes each segment independently, can still be employed against shuffling countermeasures. To illustrate this approach, consider targeting the sub-byte operations of the first round in an encryption scheme. If the approximate length of each sub-byte operation is known, the traces can be segmented accordingly, and HierNet can be applied based on that segmentation. However, utilizing HierNet in this manner necessitates knowledge of the approximate length of the operations. Such knowledge might be deduced when the source code of the implementation is available, following the procedures proposed in [23].

8. Conclusions

DL-based SCA has shown significant success in overcoming various countermeasures. However, the proper selection of the attack window is critical for its effectiveness. Consequently, several studies have explored methods to enhance the performance of DL models on longer attack windows or long raw traces, thereby reducing the dependency on the attack window selection. Despite these efforts, the performance of most of the existing DL models deteriorates in the presence of jitter-based countermeasures or as trace length increases. Therefore, developing DL models that perform well on long raw traces while remaining robust to jitter-based countermeasures remains an important task.

EstraNet, a TN-based model, was recently introduced to perform SCA on long traces. The model has demonstrated state-of-the-art performance against various combinations of masking, random delay, and clock jitter countermeasures. In this work, we closely analyze EstraNet and reveal that its performance gradually deteriorates as trace lengths increase. To address this limitation, we propose HierNet, a two-level hierarchical DL model for SCA. At the bottom level, HierNet employs EstraNet models to process smaller segments of the traces. At the top level, another DL model is used to combine information from outputs of the EstraNet models.

Extensive evaluations of HierNet have been conducted against various combinations of masking, random delay, and clock jitter countermeasures using several SCA datasets. The performance of HierNet has been compared with three existing benchmark DL models. The results indicate that HierNet outperforms the benchmark models on longer traces or in the presence of clock jitter countermeasure, showcasing its ability to reach the guessing entropy 1 using fewer than or close to 10 attack traces while the benchmark models are unable to reach the same using as many as 5K attack traces. Additionally, HierNet demonstrates significantly better performance in low training data scenarios. An ablation study related to the top-level DL model of HierNet highlights the

importance of its design choices. In conclusion, HierNet’s success establishes it as a promising method for SCA.

References

- [1] P. C. Kocher, “Timing attacks on implementations of Diffie-Hellman, RSA, DSS, and other systems,” in *CRYPTO ’96, California, USA*, ser. LNCS, N. Koblitz, Ed., vol. 1109, 1996, pp. 104–113.
- [2] P. C. Kocher, J. Jaffe, and B. Jun, “Differential power analysis,” in *CRYPTO, USA*, ser. LNCS, vol. 1666. Springer, 1999, pp. 388–397.
- [3] J. Quisquater and D. Samyde, “Electromagnetic analysis (EMA): Measures and counter-measures for smart cards,” in *E-smart 2001, France, Proceedings*, ser. LNCS, I. Attali and T. P. Jensen, Eds., vol. 2140. Springer, 2001, pp. 200–210.
- [4] S. Chari, J. R. Rao, and P. Rohatgi, “Template attacks,” in *CHES, USA, 2002*, ser. LNCS, B. S. K. Jr., Ç. K. Koç, and C. Paar, Eds., vol. 2523. Springer, 2002, pp. 13–28.
- [5] W. Schindler, K. Lemke, and C. Paar, “A stochastic model for differential side channel cryptanalysis,” in *CHES, UK, 2005*, ser. LNCS, J. R. Rao and B. Sunar, Eds., vol. 3659. Springer, 2005, pp. 30–46.
- [6] S. Picek, A. Heuser, A. Jovic, S. Bhasin, and F. Regazzoni, “The curse of class imbalance and conflicting metrics with machine learning for side-channel evaluations,” *TCHES*, vol. 2019, no. 1, pp. 209–237, 2019.
- [7] H. Maghrebi, T. Portigliatti, and E. Prouff, “Breaking cryptographic implementations using deep learning techniques,” in *SPACE, India, 2016*, ser. LNCS, vol. 10076. Springer, 2016, pp. 3–26.
- [8] S. Chari, C. S. Jutla, J. R. Rao, and P. Rohatgi, “Towards sound approaches to counteract power-analysis attacks,” in *CRYPTO, USA, 1999*, ser. LNCS, vol. 1666. Springer, 1999, pp. 398–412.
- [9] J. Coron and I. Kizhvatov, “An efficient method for random delay generation in embedded software,” in *CHES, Switzerland, 2009*, ser. LNCS, vol. 5747. Springer, 2009, pp. 156–170.
- [10] E. Cagli, C. Dumas, and E. Prouff, “Convolutional neural networks with data augmentation against jitter-based countermeasures - profiling attacks without pre-processing,” in *CHES, Taiwan, 2017*, ser. LNCS, vol. 10529. Springer, 2017, pp. 45–68.
- [11] C. Herbst, E. Oswald, and S. Mangard, “An AES smart card implementation resistant to power analysis attacks,” in *ACNS, Singapore, 2006*, ser. LNCS, J. Zhou, M. Yung, and F. Bao, Eds., vol. 3989, 2006, pp. 239–252.
- [12] L. Masure, N. Belleville, E. Cagli, M. Cornelié, D. Couroussé, C. Dumas, and L. Maingault, “Deep learning side-channel analysis on large-scale traces - A case study on a polymorphic AES,” in *ESORICS, UK, 2020*, ser. LNCS, vol. 12308. Springer, 2020, pp. 440–460.
- [13] X. Lu, C. Zhang, P. Cao, D. Gu, and H. Lu, “Pay attention to raw traces: A deep learning architecture for end-to-end profiling attacks,” *TCHES*, vol. 2021, no. 3, pp. 235–274, 2021.
- [14] G. Perin, L. Wu, and S. Picek, “Exploring feature selection scenarios for deep learning-based side-channel analysis,” *Cryptology ePrint Archive*, 2021.
- [15] E. Bursztein, L. Invernizzi, K. Král, D. Moghimi, J. M. Picod, and M. Zhang, “Generic attacks against cryptographic hardware through long-range deep learning,” *CoRR*, vol. abs/2306.07249, 2023.
- [16] S. Hajra, S. Chowdhury, and D. Mukhopadhyay, “EstraNet: An efficient shift-invariant transformer network for side-channel analysis,” *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2024, no. 1, p. 336–374, 2023.
- [17] E. Bursztein, L. Invernizzi, K. Král, D. Moghimi, J. M. Picod, and M. Zhang, “Generalized power attacks against crypto hardware using long-range deep learning,” *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2024, no. 3, pp. 472–499, 2024.
- [18] S. Hajra, S. Saha, M. Alam, and D. Mukhopadhyay, “TransNet: Shift invariant transformer network for side channel analysis,” in *AFRICACRYPT 2022, Fes, Morocco, 2022, Proceedings*, ser. LNCS, L. Batina and J. Daemen, Eds. Springer Nature Switzerland, 2022, pp. 371–396.
- [19] G. Zaid, L. Bossuet, A. Habrard, and A. Venelli, “Methodology for efficient CNN architectures in profiling attacks,” *TCHES*, vol. 2020, no. 1, pp. 1–36, 2020.
- [20] L. Wu and S. Picek, “Remove some noise: On pre-processing of side-channel measurements with autoencoders,” *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2020, no. 4, pp. 389–415, 2020.
- [21] F. Durvaux, M. Renauld, F. Standaert, L. van Oldeneel tot Oldenzeel, and N. Veyrat-Charvillon, “Efficient removal of random delays from embedded software implementations using hidden markov models,” in *CARDIS, Austria*, ser. Lecture Notes in Computer Science, S. Mangard, Ed., vol. 7771. Springer, 2012, pp. 123–140.
- [22] G. Chiari, D. Galli, F. Lattari, M. Matteucci, and D. Zoni, “A deep-learning technique to locate cryptographic operations in side-channel traces,” 2024.
- [23] L. Masure and R. Strullu, “Side channel analysis against the ANSSI’s protected AES implementation on ARM,” *IACR Cryptol. ePrint Arch.*, p. 592, 2021.
- [24] R. Benadjila, E. Prouff, R. Strullu, E. Cagli, and C. Dumas, “Deep learning for side-channel analysis and introduction to ASCAD database,” *J. Cryptogr. Eng.*, vol. 10, no. 2, pp. 163–188, 2020.
- [25] D. Bellizia, F. Berti, O. Bronchain, G. Cassiers, S. Duval, C. Guo, G. Leander, G. Leurent, I. Levi, C. Momin, O. Pereira, T. Peters, F. Standaert, B. Udvarhelyi, and F. Wiemer, “Spook: Sponge-based leakage-resistant authenticated encryption with a masked tweakable block cipher,” *IACR Trans. Symmetric Cryptol.*, vol. 2020, no. S1, pp. 295–349, 2020.

Appendix A.

Relation between $|\mathcal{R}^{(h)}|$ and n

Let us formally define $\mathcal{R}^{(h)}$ (introduced in Section 4.1) to be the set $\{j : k_{GPA}(i - j; s_p^{(h)}, c_p^{(h)}) > \delta\}$ for some $\delta \in (0, 1)$ and some $i \in \{0, 1, \dots, n - 1\}$. Then

$$\begin{aligned}
 |\mathcal{R}^{(h)}| &= \left| \{j : k_{GPA}(i - j; s_p^{(h)}, c_p^{(h)}) > \delta\} \right| \\
 &= \left| \left\{ j : \hat{\beta}_2 s_p^{(h)} \left(i/n - j/n - c_p^{(h)} \right) \|W_p \mathbf{p}\|_2 < \sqrt{2 \log(1/\delta)} \right\} \right| \\
 &= \left| \left\{ j : \left(i/n - j/n - c_p^{(h)} \right) < \frac{\sqrt{2 \log(1/\delta)}}{\hat{\beta}_2 s_p^{(h)} \|W_p \mathbf{p}\|_2} \right\} \right| \\
 &= \left| \{j : (i/n - j/n - c_p^{(h)}) < C\} \right|
 \end{aligned}$$

where $C = \frac{\sqrt{2 \log(1/\delta)}}{\hat{\beta}_2 s_p^{(h)} \|W_p \mathbf{p}\|_2}$. The parameters $c_p^{(h)}$ and $s_p^{(h)}$ are independent of n or $\hat{\beta}_2$ and $\|W_p \mathbf{p}\|_2$ is constant. Therefore, when $\hat{\beta}_2$ remains the same, $|\mathcal{R}^{(h)}|$ increases almost linearly with n .

Appendix B.

Details of Dataset

This section provides a detailed description of the datasets used.

ASCAD Fixed Key (ASCADf). The ASCAD fixed key dataset, referred to as the ASCADf dataset, is a collection of EM traces collected from a masked implementation of AES executed on the 8-bit microcontroller ATMega8515. This dataset comprises a total of 60K traces. Out of these, 50K traces were used for training, while two sets of 5K traces each were used for validation and testing. Each trace in the dataset is 100K features long. In line with earlier literature [13], [16], [19], [24], we aim to attack the third key byte of the first-round key and use the identity leakage model.

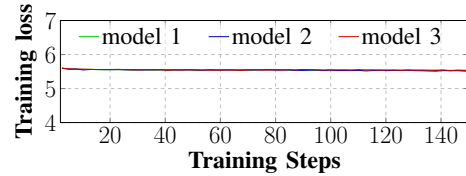
ASCAD Random Key (ASCADr). Similar to the ASCADf dataset, the ASCAD random key dataset, referred to as the ASCADr dataset, also corresponds to a masked implementation of AES executed on an 8-bit ATMega8515 microcontroller. The dataset consists of 200K profiling traces and 100K attack traces. Unless stated otherwise, all profiling traces have been utilized for training the DL models, while two disjoint random samples of 10K attack traces were used for validation and testing, respectively. Each trace in the dataset spans 250K features. To reduce memory usage, the traces were compressed to a length of 125K by replacing each consecutive pair of features of the original traces with their average. Therefore, the attacks were conducted on the compressed traces of length 125K. Similar to the ASCADf dataset, we attack the third key byte of the first round key using identity leakage model.

CHES 2020 CTF SW3 (CHES20). The CHES 2020 CTF presents a SCA challenge targeting the masked implementation of the Clyde-128 Tweakable Block Cipher (TBC) used in the Spook AEAD encryption scheme [25]. The challenge encompasses various datasets derived from distinct implementations of the cipher. In this study, we utilize the dataset corresponding to a second-order masked implementation executed on an ARM Cortex-M0 microcontroller. The dataset comprises 200K profiling traces and 500K attack traces. Unless stated otherwise, all profiling traces are employed for training the DL models, while for validation and testing, two separate and randomly selected sets of attack traces, each containing 10K traces, are utilized. It should be noted that Clyde-128 employs a (4×32) -bit state, with 4 representing the size of the non-linear S-box and 32 indicating the size of the linear L-box. Consistent with the approach in [16], we attack the four bits of the 17th column (from the right) of the state after the first round sbox operation.

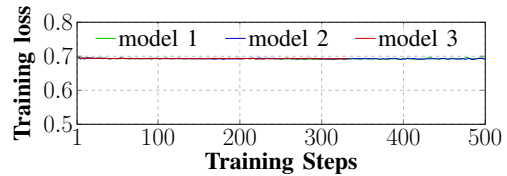
Appendix C. Details of the Benchmark Models

This section provides short descriptions of the benchmark DL models.

LSTMNet. LSTMNet [13] is a family of DL models that incorporate LSTM layers to aggregate information from multiple Points of Interest (PoIs). These models exhibit linear scalability with respect to input length, making those well-suited for processing very long traces. In their work



(a) The ASCADr Dataset.



(b) The CHES20 Dataset.

Figure 4. Training loss of three independent training of EffCNN model on the ASCADr and CHES20 datasets.

[13], Lu et al. demonstrated LSTMNet’s proficiency in attacking available full-length raw traces. For the ASCADf and ASCADr datasets, we trained the respective models available in their online repository⁴ to conduct the attack. Since no model is available for the CHES20 dataset, following [16], we used the model designed for the ASCADr dataset to train on this dataset. The models were trained for a maximum of 4K epochs, employing Adam optimizer with a learning rate of $1e-4$.

PolyCNN. The PolyCNN model [12], introduced by Masure et al. [12], is designed for attacking implementations protected by code polymorphism. The model is comprised of multiple convolutional blocks succeeded by a global average-pooling layer. [12] has demonstrated that the model can perform very well on very long traces. Consequently, we employ PolyCNN as a benchmark model in our study. The model underwent training for a maximum of 5K epochs, utilizing the Adam optimizer with a learning rate of $1e-5$.

EffCNN. In their work [19], Zaid et al. proposed a methodology for constructing CNN-based models for SCA. They successfully demonstrated the applicability of their methodology in creating compact CNN models suitable for both synchronized and desynchronized trace scenarios. However, it is worth noting that when applying their methodology to very long traces with large trace desynchronizations, the resultant models become large and exhibit slow training speeds. Despite these considerations, we utilize models constructed through their methodology as benchmark models in our study. It is important to emphasize that distinct models have been constructed based on their methodology for different datasets. The training process involved a maximum of 2K epochs, utilizing the Adam optimizer with a learning rate of $2.5e-5$.

4. https://github.com/lxj-sjtu/TCHES2021_Pay_attention_to_the_raw_traces

Appendix D.

Plots of Training Loss of EffCNN Model on the ASCADr and CHES20 datasets

In Figure 4, we illustrate the EffCNN model’s training loss convergence in the experiments with the ASCADr and CHES20 datasets. Figure 4a displays the training loss for the ASCADr dataset, indicating negligible improvement even after 150 epochs. Notably, training for 150 epochs on this dataset took approximately 50 hours, whereas HierNet models achieved full convergence within 24 to 40 hours. Similarly, Figure 4b shows that even after 500 epochs (equivalent to 80 hours of training), the training loss did not significantly decrease compared to random guessing.