# STD-NET: Search of Image Steganalytic Deep-learning Architecture via Hierarchical Tensor Decomposition

Shunquan Tan, *Senior Member, IEEE,* Qiushi Li, Laiyuan Li, Bin Li, *Senior Member, IEEE,*
and Jiwu Huang, *Fellow, IEEE*

*Abstract*—Recent studies shows that the majority of existing deep steganalysis models have a large amount of redundancy, which leads to a huge waste of storage and computing resources. The existing model compression method cannot flexibly compress the convolutional layer in residual shortcut block so that a satisfactory shrinking rate cannot be obtained. In this paper, we propose STD-NET, an unsupervised deep-learning architecture search approach via hierarchical tensor decomposition for image steganalysis. Our proposed strategy will not be restricted by various residual connections, since this strategy does not change the number of input and output channels of the convolution block. We propose a normalized distortion threshold to evaluate the sensitivity of each involved convolutional layer of the base model to guide STD-NET to compress target network in an efficient and unsupervised approach, and obtain two network structures of different shapes with low computation cost and similar performance compared with the original one. Extensive experiments have confirmed that, on one hand, our model can achieve comparable or even better detection performance in various steganalytic scenarios due to the great adaptivity of the obtained network architecture. On the other hand, the experimental results also demonstrate that our proposed strategy is more efficient and can remove more redundancy compared with previous steganalytic network compression methods.

*Index Terms*—Steganalysis, steganography, deep learning, convolutional neural network, tensor decomposition.

## I. INTRODUCTION

STEGANALYSIS aims to reveal covert communication established via steganography. For steganalytic frameworks "into the wild", low memory/computational cost, as well as lightweight model size are just as important as high detection performance.

For steganography, digital image is the most commonly used cover medium, and is the main battleground of the war between steganography and steganalysis [1]. Over the past decade, the so-called embedding distortion minimizing

framework [2] has reigned supreme on both spatial and frequency domain for image steganography. Most state-of-the-art steganographic algorithms can be categorized as additive embedding distortion minimizing schemes, including HILL [3] and MiPOD [4] in spatial domain, UERD [5] in JPEG domain, as well as UNIWARD [6] (including S-UNIWARD in spatial domain, and J-UNIWARD in JPEG domain). Research on non-additive distortion functions has made progress both in spatial domain [7]–[9] and JPEG domain [10]–[12]. Further on, some pioneering works have been devoted to incorporating deep-learning and reinforcement learning models into embedding distortion minimizing framework [13]–[17].

In the arm race with steganography, steganalysis has envolved from the old-style "rich model" hand-crafted features family [18]–[22] equipped with an ensemble classifier [23] to deep-learning based solutions. Started from the work of Tan and Li [24], deep-learning based steganalysis has gradually overtaken "rich model" features family [25]–[29] in the last five years. In [28], Ye et al. proposed a deep steganalytic network equipped with a pre-processing layer SRM (spatial rich model) and a new activation function called Truncated Linear Unit (TLU), achieving significant improvement in spatial domain. In [30], Xu proposed a 20-layer deep residual steganalytic network. In [31], Zeng et al. proposed a generic hybrid deep-learning framework aiming at large-scale JPEG image steganalysis. Afterwards in [32], Zeng et al. proposed WISERNet, the wider separate-then-reunion network specifically designed for steganalysis of true-color images. Boroumand et al. proposed a deep residual steganalytic network SRNet [33], which is the first end-to-end steganalytic network. It becomes a popular deep-learning solution due to its good performance in both spatial and frequency domain steganalysis. Inspired by Yedroudj-Net [34], Zhang et al. proposed another deep-learning framework with distinct advantage in spatial-domain steganalysis in [35]. In [36], You et al. proposed SiaStegNet, a siamese CNN (Convolutional Neural Network) framework aiming at steganalysis to images of arbitrary size. All of the above mentioned frameworks more or less incorporate the domain knowledge behind the "rich model" hand-crafted features family, leading to dramatic computational cost and storage overheads in its residual extracting.

It is well known that aimlessly scaling up deep-learning network is not necessarily a guarantee of better performance [37]. The scalability of deep-learning steganalyzers is by and large relevant to their application scenarios [38].

Therefore, research of network architecture search [39] and model compression/pruning [40], [41] becomes a magnet in recent years. In [42], the authors proposed CALPA-NET, a channel-pruning-assisted deep residual network architecture search approach to shrink the network structure of existing deep-learning based steganalyzers. The obtained architecture can achieve comparative performance with less than 2% parameters compared to SRNet. However, CALPA-NET still has the following deficiencies:

- It is a supervised data-driven solution in which labeled stego samples are indispensable;
- It cannot shrink the bottom layers linked via direct shortcut connections, especially the noise residual extraction blocks in SRNet, with satisfying shrinking rates;
- Its performance might be mildly degraded (1%~2%) especially when aiming at J-UNIWARD steganography.

Many researchers have pointed out that tensor analysis and consequently hierarchical tensor decomposition can be used to dissect and optimize CNNs, the cornerstone of deep-learning techniques [43]–[46]. Specifically, tensor decomposition based techniques such as CP decomposition [44] and Tucker decomposition [45], which replace convolutional layers with their low-rank matrices approximations to achieve the purpose of effectively removing redundant information, has become a significant branch of model compression in deep learning. In this paper, the authors move a step further and propose STD-NET, an image steganalytic deep-learning architecture search approach via hierarchical tensor decomposition. In order to prevent the resulting framework from being affected by specific steganography algorithm, and to make our decomposition strategy more general, we proposed a normalized distortion threshold to evaluate the sensitivity of each involved convolutional layer of the target model and to guide STD-NET to compress it in an unsupervised way. Therefore with the help of hierarchical tensor decomposition, starting from SRNet, we obtain a novel image steganalytic deep-learning architecture via shrinking as well as deepening the structure of SRNet. Compared with prior works, STD-NET has the following novel improvements:

- It is an unsupervised data-driven solution in which only a few unlabeled samples are required, and hence a unified criterion can be applied to compress all involved convolutional layers;
- It can significantly "shrink" the noise residual extraction blocks in SRNet and consequently dramatically save the parameters and FLOPs.
- The extensive experiments conducted on large-scale public datasets have shown that STD-NET has shown fairly mild degradation in detection performance. Under certain scenarios its performance has even outperformed the original SRNet model.

The rest of the paper is organized as follows. Sect. II firstly gives a brief overview of convolution operation, hierarchical tensor decomposition of convolution operation as well as SRNet, the representative deep steganalytic model. Then STD-NET, our proposed image steganalytic deep-learning architecture search approach is described in detail. Results

of experiments conducted on large-scale public datasets and corresponding discussions are presented in Sect. III. Finally, we make a conclusion in Sect. IV.

## II. Our proposed STD-NET

Throughout the text, boldface capital letters, e.g., $\mathbf{X}$, denotes matrices, and boldface Euler script letters, e.g., $\mathcal{X}$, denotes high-order tensors. Given a tensor $\mathcal{X}$, its Euclidean norm is denoted as $\|\mathcal{X}\|$.

### A. Preliminaries

CNN is not only the cornerstone of deep-learning techniques, but also the infrastructure of existing deep-learning based steganalyzers. As the name implies, the most important components of a CNN is its convolutional layers. They take up the overwhelming majority of learnable parameters.

Given a convolutional layer $L_l$, its convolution operation takes an input tensor $\mathcal{I}^{l-1} \in \mathbb{R}^{J^{l-1} \times H^{l-1} \times W^{l-1}}$ which contains $J^{l-1}$ input channels with height $H^{l-1}$ and width $W^{l-1}$, and maps it into an output tensor $\mathcal{O}^l \in \mathbb{R}^{J^l \times H^l \times W^l}$ which contains $J^l$ output channels with height $H^l$ and width $W^l$. Let $\mathcal{K}^l \in \mathbb{R}^{J^{l-1} \times J^l \times D^l \times D^l}$ denote the kernel tensor. Elementwise:

$$\mathcal{O}^l_{i,h,w} = \sum_{j=1}^{J^{l-1}} \sum_{k=1}^{D^l} \sum_{s=1}^{D^l} \mathcal{K}^l_{j,i,k,s} \cdot \mathcal{I}^{l-1}_{j,h_k,w_s},$$
$$h_k = (h-1) \cdot \Delta + k - \Psi, \ w_s = (w-1) \cdot \Delta + s - \Psi \quad (1)$$

in which $\mathcal{O}^l_{i,h,w}$, $\mathcal{K}^l_{j,i,k,s}$, and $\mathcal{I}^{l-1}_{j,h_k,w_s}$ denotes the $(i,h,w)$-th, the $(j,i,k,s)$-th, and the $(j,h_k,w_s)$-th element of tensor $\mathcal{O}^l$, $\mathcal{K}^l$, and $\mathcal{I}^{l-1}$ respectively. $\Delta$ is stride and $\Psi$ is zero-padding length. In existing deep-learning based steganalyzers, the kernel filters are usually quite small, typically $3 \times 3$ or $5 \times 5$.

In the terminology of tensor analysis, *mode* is an alias of dimension. Given a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$ and a matrix $\mathbf{U} \in \mathbb{R}^{J \times I_n}$, the *n-mode product* of $\mathcal{X}$ and $\mathbf{U}$ is actually multiplying $\mathcal{X}$ by $\mathbf{U}$ along its $n$-th dimension, and is denoted as $\mathcal{X} \times_n \mathbf{U} \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times J \times I_{n+1} \times \cdots \times I_N}$. Elementwise:

$$(\mathcal{X} \times_n \mathbf{U})_{i_1 \cdots i_{n-1} j i_{n+1} \cdots i_N} = \sum_{i_n=1}^{I_n} \mathcal{X}_{i_1 \cdots i_n \cdots i_N} \cdot \mathbf{U}_{j i_n}$$

*1) Tensor decomposition of convolution operation:* Here a famed variant of tensor decomposition, Tucker decomposition which is a form of higher-order PCA (Principal Component Analysis) is adopted [43]. The Tucker decomposition of convolution operation of a given kernel tensor $\mathcal{K}^l$ can be defined as:

$$\mathcal{K}^l \approx \mathcal{G}^l \times_1 \underset{\triangle}{\mathbf{T}}^l \times_2 \overset{\triangle}{\mathbf{T}}^l \quad (2)$$

in which the so-called *core tensor* $\mathcal{G}^l \in \mathbb{R}^{I^l \times O^l \times D^l \times D^l}$, $\underset{\triangle}{\mathbf{T}}^l \in \mathbb{R}^{J^{l-1} \times I^l}$, $\overset{\triangle}{\mathbf{T}}^l \in \mathbb{R}^{O^l \times J^l}$. $\underset{\triangle}{\mathbf{T}}^l$ and $\overset{\triangle}{\mathbf{T}}^l$ are in orthogonal form after the decomposition. $\mathcal{G}^l$ can be regarded as a compressed version of $\mathcal{K}^l$ as long as $I^l < J^{l-1}$ and $O^l < J^l$. The mode 3 and 4 are not involved in Tucker decomposition due to the fact that they correspond to the size of the kernel filters and are quite small, and consequently their decomposition cannot lead in prominent cut down on computation costs.
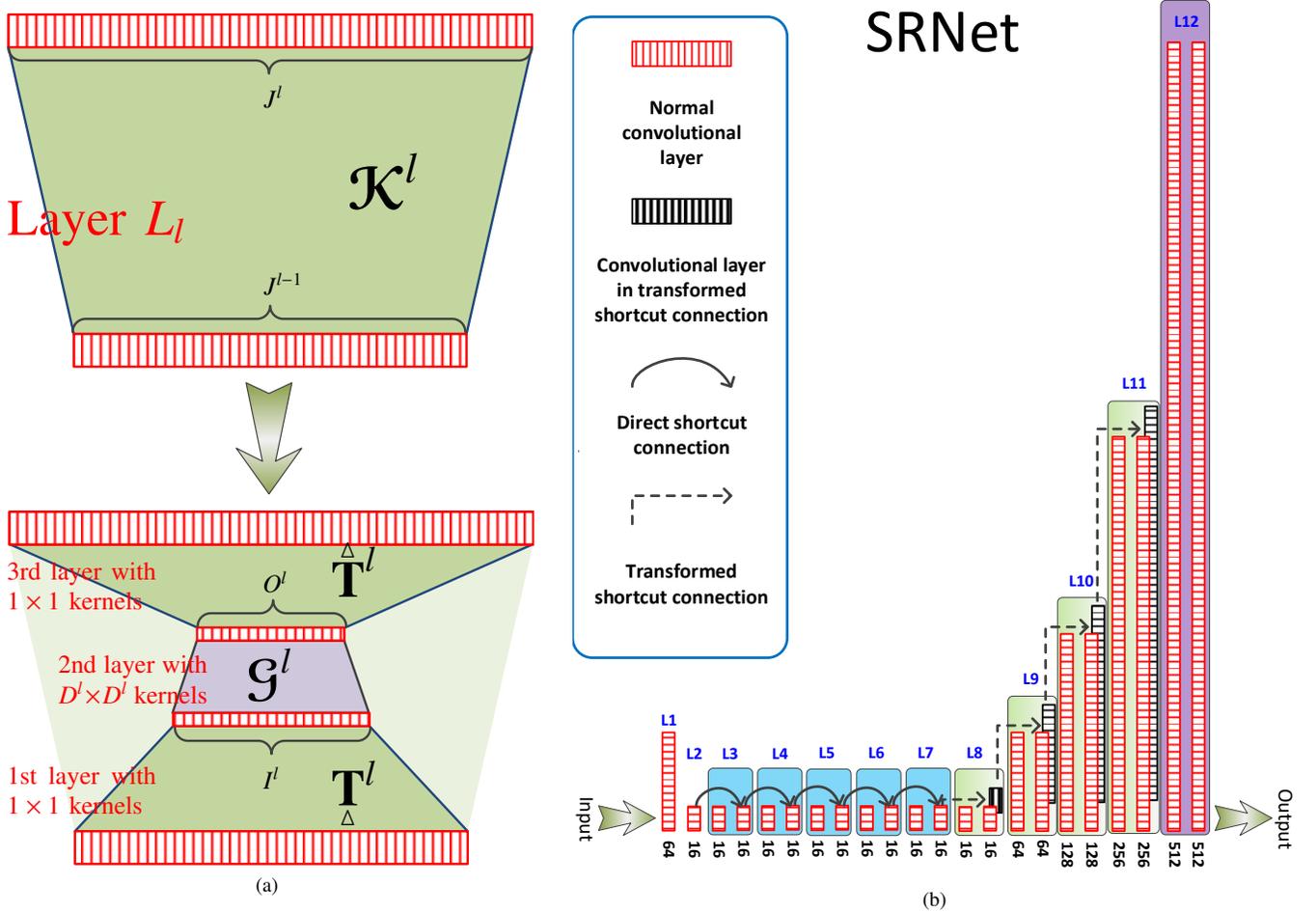
Fig. 1. (a) The conceptual map of Tucker decomposition of convolution operation. (b) The conceptual structure of SRNet.

Please note that a *n-mode product* in Eq. (2) is equivalent to introducing a convolutional layer with 1×1 kernels. Therefore as illustrated in Fig. 1(a), after Tucker decomposition, the convolutional layer $L_l$ with kernel tensor $\mathcal{K}^l$ is replaced by three cascading layers:

- The first layer with $1 \times 1$ kernels which reduces the number of feature maps by passing from an input tensor of $\mathbb{R}^{J^{l-1} \times H^{l-1} \times W^{l-1}}$ to an output tensor of $\mathbb{R}^{I^l \times H^{l-1} \times W^{l-1}}$ since $I^l < J^{l-1}$;
- The second layer with $D^l \times D^l$ kernels representing the convolution by the *core tensor* $\mathcal{G}^l$, in which the feature map tensor goes from the size of $\mathbb{R}^{I^l \times H^{l-1} \times W^{l-1}}$ to the size of $\mathbb{R}^{O^l \times H^l \times W^l}$.
- The third layer again with 1×1 kernels which re-increases the number of feature maps to the final output tensor of $\mathbb{R}^{J^l \times H^l \times W^l}$ since $O^l < J^l$.

For a given convolutional layer $L_l$, the number of parameters and the number of FLOPs (FLoating-point OPerations) are determined as follows (let $| \bullet |$ denotes the number of the corresponding variable):

- $|\text{params}|=J^{l-1} \cdot J^l \cdot D^l \cdot D^l$; [1]
- $|\text{FLOPs}|=|\text{params}| \cdot H^l \cdot W^l$.

[1]The bias, as well as a few parameters in other types of layers are omitted.

After Tucker decomposition as defined in (2), the corresponding metrics of the decomposed $L_l$ can be calculated as:

- $|\text{params}|=J^{l-1} \cdot I^l + I^l \cdot O^l \cdot D^l \cdot D^l + O^l \cdot J^l$;
- $|\text{FLOPs}|=|\text{params}| \cdot H^l \cdot W^l$.

*2) Interior structure of SRNet:* SRNet [33] consists of three modules: the bottom module tries to suppress the image contents and boost SNR (Signal-to-Noise Ratio) [2], while the middle module aims at learning compact representative features and the top module is a simple binary "cover" vs. "stego" classifier.

As illustrated in Fig. 1(b), following the notations in [33], "L1" and "L2", two hierarchical convolutional layers and the subsequent five unpooled residuals blocks with direct shortcut connections (from "L3" to "L7") make up the bottom module of SRNet. From its "L8" up to "L12", the middle module gradually halves sizes of feature maps (256×256 → 128×128 → 64 × 64 → 32 × 32 → 16 × 16) as well as doubles and even quadruple (for "L8") numbers of output channels (16 → 64 → 128 → 256 → 512) layer by layer. Among the four middle blocks, "L8" to "L11" are with transformed shortcut connections. The top module is a standard fully connected

[2]In the literature of steganalysis, image content is "noise" while stego noise is "signal".

layer followed by a softmax node. Every convolutional layer is directly followed by a batch normalization layer. The design of the middle module as well as the top module of SRNet, like most of the state-of-the-art deep-learning based steganalyzers, just simply followed the effective recipes of the research field of computer vision and is therefore pretty redundant. In our proposed CALPA-NET [42], we demonstrated that the middle/top module of SRNet can be aggressively shrinked with comparative performance. On the contrary, as the key part of SRNet, the bottom module is with compact design and cannot be effectively shrinked even with CALPA-NET.

### B. Algorithm of our proposed STD-NET

*1) The overall procedure:* As mentioned in Sect. II-A, convolutional layer is the most important component of CNN, taking up the majority of learnable parameters, and there is a large amount of redundancy in it. The diagram of overall procedure is shown in Fig. 2. According to CALPA-NET [42], there is a lot of redundancy in the layers that can compact feature representation of SRNet (from "L9" to "L12"), so we set the number of output channels for these layers to 64, and named the intermediate model as SRNetC64. Firstly, the SRNetC64 model $N$ is trained using the original training protocol. Then the well-trained $N$ is traversed from bottom to top and the input/output channel number for the resulting core tensor of every involved convolutional layer is determined in an unsupervised data-driven manner (see Sect. II-B2). Two brand-new model structures (cylinder/ladder-shaped structure, see Sect. II-B3) are obtained by replacing model $N$'s kernel tensor with corresponding Tucker decomposed version. The resulting model can be trained in two different ways. The first one is preserving the model parameters after Tucker decomposition, and the other one is resetting all of the model weights (including those in all the $\mathcal{G}^l$, $\mathbf{T}^l$, and $\overset{\triangle}{\mathbf{T}}^l$) and then training it from scratch. Since the first layer of SRNet is more sensitive to compression, the following experiment will not decompose the first convolutional layer.

*2) The proposed unsupervised tensor decomposition criterion:* Since a direct residual shortcut requires the number of channels at both ends to be the same, the existing channel-based network pruning strategies cannot efficiently prune both ends of a given residual shortcut with a unified criterion. On the contrary, with Tucker decomposition, our proposed STD-NET can guarantee that the numbers of input and output channels of the resulting convolutional group are the same as the corresponding convolutional layer prior to the decomposition. In other words, the compression for every involved convolutional layer is independent and will not be affected by the direct residual connections in the network. All the convolutional layers can be treated equally, regardless of the residual modules they belong to. Therefore, our proposed STD-NET approach is more concise and efficient than CALPA-NET.

It should be noted that the determination of the input/output channel number of the core tensor directly affects the final efficiency of the network compression. Therefore, we define a
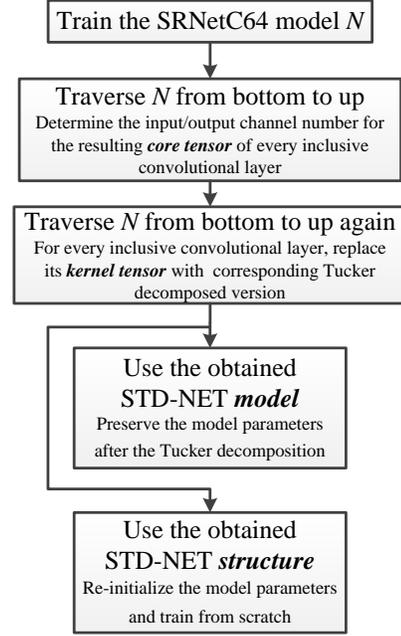


Fig. 2. The overall STD-NET diagram.

normalized distortion threshold:

$$\mathcal{T}_d^l = \frac{\|\widehat{\mathcal{O}}^l - \widehat{\mathcal{O}}_{\mathbf{T}}^l\|}{\|\widehat{\mathcal{O}}^l\|} \tag{3}$$

where $\widehat{\mathcal{O}}^l$ and $\widehat{\mathcal{O}}_{\mathbf{T}}^l$ are the batch-normalized version of the features maps of the $l$-th convolutional layer in SRNetC64 model $S$ and the batch-normalized output feature maps of the corresponding Tucker decomposed convolutional group, respectively. The $\|\cdot\|$ denotes Euclidean norm. The resulting normalized distortion threshold $\mathcal{T}_d^l$ can effectively reflect the sensitivity of each involved convolutional layer to Tucker decomposition and assist us in choosing the appropriate number of input and output channels in an unsupervised way.

For the $l$-th convolutional layer, determining the normalized distortion $\mathcal{T}_d^l$ can be summarized as 3 steps. Firstly, initialize the corresponding shrinking rate to 100%. Secondly, decrease the shrinking rate and calculate a sequence of normalized distortion values $\mathcal{T}^l$. Finally, choose an appropriate normalized distortion value from $\mathcal{T}^l$ and set it as $\mathcal{T}_d^l$. The detailed traversal algorithm for $\mathcal{T}^l$ is shown in Algorithm. 1. Within a certain shrinking rate change range, if the normalized distortion values changes more drastically, it indicates that the corresponding feature representation and even the final detection performance will suffer greater impact. Therefore, we need to make a trade-off between compression efficiency and model performance so as to obtain a more lightweight model without obvious performance drop. Through traversing from bottom to top, we can select an appropriate normalized distortion threshold $\mathcal{T}_d^l$ for every involved convolutional layer $L_l$. In Sect. III-B we provide relevant experimental data of determining the normalized distortion values. Please note that in real scenarios, since cover images are easy to obtain, we can only use cover images as samples for this step.

---

**Algorithm 1** Normalized distortion traversal algorithm for $\mathcal{T}^l$.

---

**Require:** A SRNetC64 model $N$ with the best validation accuracy and $\mathcal{K}^l$, kernel tensor of the layer $L_l$ being processed, a batch of images $\mathcal{B}$ (without label information) randomly selected from the training dataset. a pre-defined step $\epsilon$, a pre-defined lower bound of the shrinking rate $\tau$.

1: Initialize $I'$ and $O'$, the input and output channel number of $\mathcal{G}^l$ as: $I' = J^{l-1}$ and $O' = J^l$, the shrinking rate $\gamma = 100\%$. Initialize the resulting distortion threshold list $\mathcal{T}^l$ as an empty one-dimensional array.
2: Feed $N$ with $\mathcal{B}$, and then feedforward the input in cascaded layers bottom to up till $\widehat{\mathcal{O}}^l$ is obtained.
3: **repeat**
4:     Let $I' = \lfloor J^{l-1} \cdot \gamma \rfloor$ and $O' = \lfloor J^l \cdot \gamma \rfloor$.
5:     $\mathcal{T}_\gamma$=CalNormDistortion($I'$, $O'$).
6:     Append $\mathcal{T}_\gamma$ to $\mathcal{T}^l$.
7:     $\gamma = \gamma - \epsilon$.
8: **until** $\gamma < \tau$

9: **function** CalNormDistortion($I$, $O$)
10:     Apply Tucker decomposition to $\mathcal{K}^l$, and use $\mathcal{G}^l \times_1 \overset{\triangle}{\mathbf{T}}{}^l \times_2 \overset{\triangle}{\mathbf{T}}{}^l$ to replace $\mathcal{K}^l$ in $N$.
11:     Let $N'$ denotes the reconstructed model. Feed $\mathcal{B}$ into $N'$ again and then feedforward it to get $\widehat{\mathcal{O}}^l_\mathbf{T}$. Calculate the normalized distortion between $\widehat{\mathcal{O}}^l$ and $\widehat{\mathcal{O}}^l_\mathbf{T}$ as: $D = \frac{\|\widehat{\mathcal{O}}^l - \widehat{\mathcal{O}}^l_\mathbf{T}\|}{\|\widehat{\mathcal{O}}^l\|}$.
12:     **return** $D$.
13: **end function**

---

Eq. (3) actually reflects the sensitivity of each involved convolutional layer to Tucker decomposition. The choice of such a distortion threshold comes directly from optimization objective of Tucker decomposition, which minimizes changes to the output tensor using the *core tensor* with least rank.

The vanilla Tucker decomposition is implemented with manual selection of the rank of the core tensor. Instead, as shown in Eq. (3),we adopt an unsupervised data-driven criterion (no sample labels are involved in the criterion) in our proposal. The idea comes from ThiNet [40], one of the networking pruning scheme utilized in our prior work CALPA-NET [42]. ThiNet greedily prunes those channels with smallest effect on the activation values of the next layer. With a similar scheme, in our proposal the convolutional layer $L_l$ is replaced by three cascading layers, and the rank of the core tensor is tuned to guarantee that the corresponding Tucker decomposition is with smallest effect on the activation values of the third layer with $1 \times 1$ kernels.

*3) Input/output channel configuration:* Based on the proposed unsupervised tensor decomposition criterion, we design an algorithm of input/output channel number determination for the core tensor $\mathcal{G}^l$.

As described in Algorithm. 2, for the $l$-th convolutional layer to be decomposed in a given SRNetC64 model $N$, the numbers of input and output channels are respectively denoted as $J^{l-1}$ and $J^l$.

---

**Algorithm 2** Input/output channel number determination algorithm for $\mathcal{G}^l$.

---

**Require:** A SRNetC64 model $N$ with the best validation accuracy and the kernel tensor $\mathcal{K}^l$ of the being processed layer $L_l$, a batch of images $\mathcal{B}$ (without label information) randomly selected from the training dataset, a pre-determined distortion threshold $\mathcal{T}^l_d$, and a distortion margin $\varsigma$.

1: Initialize $I'$ and $O'$, the input and output channel number of $\mathcal{G}^l$ as: $I' = J^{l-1}$ and $O' = J^l$.
2: Define the input/output channel number for a cylinder-shaped, and a ladder-shaped $\mathcal{G}^l$ as $I^l_{\text{cylinder}}/O^l_{\text{cylinder}}$, and $I^l_{\text{ladder}}/O^l_{\text{ladder}}$, respectively.
3: **if** $J^{l-1} < J^l$ **then**
4:     $\epsilon_{I'} = 1$, $\epsilon_{O'} = \lfloor O'/I' \rfloor$.
5: **else**
6:     $\epsilon_{I'} = \lfloor I'/O' \rfloor$, $\epsilon_{O'} = 1$.
7: **end if**
8: Feed $N$ with $\mathcal{B}$, and then feedforward the input in cascaded layers bottom to up till $\widehat{\mathcal{O}}^l$ is obtained.
9: **repeat**
10:     Let $I' = I' - \epsilon_{I'}$ and $O' = O' - \epsilon_{O'}$.
11:     $\mathcal{T}'_d$=CalNormDistortion($I'$, $O'$).
12: **until** $\mathcal{T}'_d > \mathcal{T}^l_d$
13: **if** $J^{l-1} \neq J^l$ **then**
14:     The input/output channel number for a ladder-shaped $\mathcal{G}^l$ has been determined as: $I^l_{\text{ladder}} = I'$ and $O^l_{\text{ladder}} = O'$. **exit**
15: **else**
16:     $I^l_{\text{cylinder}} = I'$ and $O^l_{\text{cylinder}} = O'$.
17: **end if**
18: Let $I'' = I''' = I'$ and $O'' = O''' = O'$.
19: **repeat**
20:     Let $I'' = I'' - 1$ and $O'' = O'' + 1$.
21:     $\mathcal{T}''_d$=CalNormDistortion($I''$, $O''$).
22: **until** $\mathcal{T}''_d > \mathcal{T}^l_d + \varsigma$
23: **repeat**
24:     Let $I''' = I''' + 1$ and $O''' = O''' - 1$.
25:     $\mathcal{T}'''_d$=CalNormDistortion($I'''$, $O'''$).
26: **until** $\mathcal{T}'''_d > \mathcal{T}^l_d + \varsigma$
27: **if** $I'' \cdot O'' < I'''' \cdot O'''$ **then**
28:     $I^l_{\text{ladder}} = I''$ and $O^l_{\text{ladder}} = O''$.
29: **else**
30:     $I^l_{\text{ladder}} = I'''$ and $O^l_{\text{ladder}} = O'''$.
31: **end if**

---

Firstly, the numbers of input and output channels of the core tensor $\mathcal{G}^l$ are initialized to $I' = J^{l-1}$, $O' = J^l$. In this initial stage $\mathcal{G}^l$ is actually with the same size as the original kernel tensor $\mathcal{K}^l$, as illustrated in the left-hand side of Fig. 3.

Secondly, as shown in line 3–17 of Algorithm 2, $I'$ and $O'$ are reduced in certain steps at the same time until the corresponding distortion threshold reaches a pre-defined normalized distortion. Then the reduced core tensor $\mathcal{G}^l$ is obtained. In this stage, $\epsilon_{I'}$ and $\epsilon_{O'}$, the reduction steps of $I'$ and $O'$ are all set
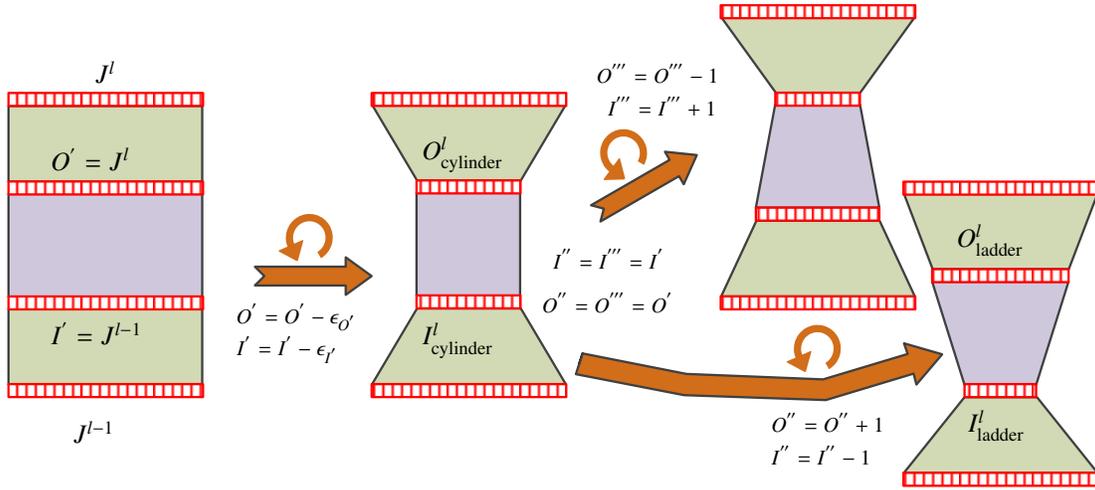
Fig. 3. Flow chart of input/output channel number determination algorithm for $\mathcal{G}^l$.

to 1 for those convolutional layers in which both input and output are with identical number of channels. $L_2$ and $L_{9-1}$ of SRNetC64 are two exceptions since their input and output are with different number of channels. For those exceptional layers, $\epsilon_{I'}$ and $\epsilon_{O'}$ are set to make $I'$ and $O'$ reduce in the same proportion.

Hence after this stage, for those convolutional layers in which both input and output are with identical number of channels, the shape of the reduced core tensor $\mathcal{G}^l$ is similar to a cylinder, as illustrated in the middle of Fig. 3. The number of input and output channels currently searched for $\mathcal{G}^l$ is denoted as $I^l_{cylinder}$ and $O^l_{cylinder}$, respectively.

Thirdly, as shown in line 18–31 of Algorithm. 2, we try to magnify the difference between the input and output channel numbers under the constraint that their sum is fixed, in order to further reduce the model parameters of the obtained three cascading layers after the Tucker decomposition. The shape of the further reduced core tensor $\mathcal{G}^l$ obtained in this way is similar to a ladder.

In this stage two iterative cycles are conducted, in which the first one decreases input channels and at the same time increases output channels ($I'' = I'' - 1$ and $O'' = O'' + 1$), and the second one does the opposite ($I''' = I''' + 1$ and $O''' = O''' - 1$). A distortion margin $\varsigma$ is introduced in this stage. Each of the two cycles stops once the normalized distortion exceeds the distortion threshold plus margin, namely $\mathcal{T}^l_d + \varsigma$.

After this stage, as illustrated in the right-hand side of Fig. 3, two reduced core tensor are obtained. The one with the biggest difference is selected, and the number of its input and output channel is denoted as $I^l_{ladder}$ and $O^l_{ladder}$, respectively.

Please note that the last stage is unnecessary for those layers whose input and output are with different number of channels (i.e. $L_2$ and $L_{9-1}$ of SRNetC64). They are already with a ladder-shaped $\mathcal{G}^l$, hence the number of the input and output channel of such a ladder-shaped $\mathcal{G}^l$ is denoted as $I^l_{ladder}$ and $O^l_{ladder}$, respectively.

After traversing the given SRNetC64 model from bottom to up, every original involved convolutional layer from the SRNetC64 model is replaced with Tucker-decomposed convo-

lutional group with a determined input/output channel number, and the corresponding STD-NET is obtained. We name the obtained STD-NET with cylinder-shaped core tensors (except $L_2$ and $L_{9-1}$) the cylinder-shaped STD-NET. Likewise, we name the one with ladder-shaped core tensors the ladder-shaped STD-NET.

*4) Theoretical analysis regarding to ladder-shaped configuration:* Here we prove that the ladder-shaped structure is more computational effective than the corresponding cylinder-shaped structure. From Sect. II-A1 we have known that for a given convolutional layer $L_l$, its parameters after Tucker decomposition as defined in (2) can be calculated as:

$$|\text{params}| = J^{l-1} \cdot I^l + I^l \cdot O^l \cdot D^l \cdot D^l + O^l \cdot J^l \quad (4)$$

If a cylinder-shaped configuration is adopted, we can get $I^l = I^l_{cylinder}$ and $O^l = O^l_{cylinder}$ for $\mathcal{G}^l$ in line 16 of Algorithm 2. Since $L_2$ and $L_{15}$ have been excluded, $I^l = O^l$ and $J^l = J^{l-1}$. From (4) we can get:

$$|\text{params}| = \underbrace{D^l \cdot D^l \cdot (I^l \cdot O^l)}_{①} + \underbrace{J^{l-1} \cdot (I^l + O^l)}_{②} \quad (5)$$

With a fixed total channel number, namely a fixed $\Delta = I^l_{cylinder} + O^l_{cylinder}$:

$$0 \le (I^l - O^l)^2$$
$$\Rightarrow 0 \le (I^l + O^l)^2 - 4 \cdot I^l \cdot O^l$$
$$\Rightarrow I^l \cdot O^l \le \frac{\Delta^2}{4} \quad (6)$$

In (6) the equality is hold when $I^l = O^l$, which means that ① obtains its maximum with $I^l = O^l$.

Since ② $= J^{l-1} \cdot \Delta$, it is fixed. Therefore $|\text{params}|$ reaches its maximum with a cylinder-shaped configuration.

Again from (6) we can get:

$$\frac{\Delta^2}{4} - I^l \cdot O^l = \frac{(I^l - O^l)^2}{4} \ge 0$$

Since $\frac{\Delta^2}{4}$ is fixed, the larger the difference between $I^l$ and $O^l$ is, the smaller $I^l \cdot O^l$ is, and consequently the smaller ①
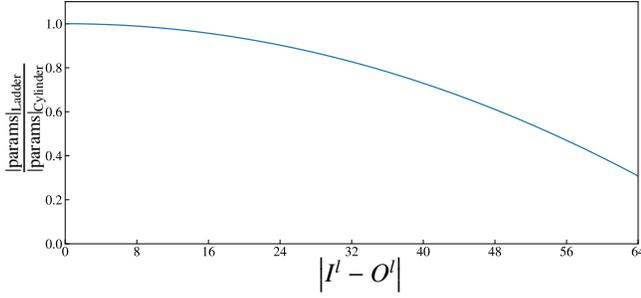
Fig. 4. The complexity proportion of a ladder-shaped core tensor to a cylinder-shaped one when the sums of their input and output channels are equal.

is. As a result, from (5) we can see the larger the difference between $I^l$ and $O^l$ is, the less the parameters are required for $L_l$. Consequently, on the basis of $I^l_{cylinder}$ and $O^l_{cylinder}$, keeping the sum of input and output channel numbers in the core tensor constant, magnifying the difference between them helps to further reduce the model parameters.

Fig. 4 provides a further demonstration. Assume that we apply Tucker decomposition to a convolutional layer with $J^{l-1} = J^l = 64$. we set $D^l \cdot D^l = 9$, and the shrinking rate for searching cylinder-shaped core tensor $\gamma = 50\%$ (which means that for the core tensor $I^l + O^l = 64$). We plot the complexity proportion of the resulting ladder-shaped core tensor to the corresponding cylinder-shaped one when the sums of their input and output channels are equal in Fig. 4. From Fig. 4 we can observe that the greater the difference between $I^l$ and $O^l$, the less the parameters of the ladder-shaped core tensor compared to its cylinder-shaped peer.

*5) Computational complexity analysis:* In order to show under what circumstances the computational complexity after Tucker decomposition is lower than the original one, we turn to solve the following inequality:

$$J^{l-1} \cdot J^l \cdot D^l \cdot D^l - J^{l-1} \cdot I^l - I^l \cdot O^l \cdot D^l \cdot D^l - O^l \cdot J^l \geq 0 \quad (7)$$

Since existing deep-learning based steganalyzers usually adopt fixed small kernel filters, such as $3 \times 3$ or $5 \times 5$, in Eq. (7) we treat $D^l \cdot D^l$ as a constant. $J^l$ is usually an integral multiple of $J^{l-1}$. Here for simplicity, we assume that $J^{l-1} = J^l$. The proofs extended to a general case are straightforward and are omitted for brevity. As for $I^l$ and $O^l$, we only consider the cylinder-shaped structure, namely $I^l = O^l$ since in Sect. II-B4 we have proved that the ladder-shaped structure is more computational effective than the corresponding cylinder-shaped structure.

To simplify the notation we set $D^l \cdot D^l = n$, $J^{l-1} = J^l = x$, and $I^l = O^l = \gamma x$, $\gamma \in [0, 1]$. So that Eq. (7) turns to:

$$nx^2 - \gamma x^2 - n\gamma^2 x^2 - \gamma x^2 \geq 0 \Rightarrow n\gamma^2 + 2\gamma - n \leq 0 \quad (8)$$

Please note that now the solution of Eq. (8) is actually irrelevant to $J^{l-1}$ as well as $J^l$. Since $n\gamma^2 + 2\gamma - n$ is a univariate quadratic function of $\gamma$ with roots $\frac{-1 \pm \sqrt{1+n^2}}{n}$, Eq. (8) holds in $\gamma \in [0, \frac{-1+\sqrt{1+n^2}}{n}]$ and its left side attains a minimum at $\gamma = 0$ (the actual minimum of the quadratic function arrives at $-\frac{1}{n}$ which is outside the value range of $\gamma$). Specifically, with kernel filter sizes $3 \times 3$ ($n = 9$) and $5 \times 5$ ($n = 25$), Eq. (8)

holds in $[0, 0.895]$ and $[0, 0.961]$, respectively. That is to say, for instance, even with the small $3 \times 3$ filters used in SRNet, Tucker decomposition brings in computational complexity reduction as long as the shrinking rate $\gamma < 89.5\%$. Even with a conservative shrinking rate $\gamma = 50\%$, Tucker decomposition can cut the computational complexity by 63.8%.

Please note that re-increasing the number of channels in STEP 3) is just a natural consequence of Tucker decomposition. However, by doing so, there is no longer any need to chain the modifications in the network as in CALPA-NET and the modifications for each layer can be made independently. Compared with the above-mentioned gain, the cost of the introduction of STEP 3) is negligible. Please note that in Eq. (7), the part corresponding to STEP 3) is $O^l \cdot J^l$. With the assumption and the simplified notation used here, we can get the proportion of the part corresponding to STEP 3) in the reduction of computational complexity (the left side of the inequality in Eq. (8)) as $\frac{\gamma}{n - n\gamma^2 - \gamma - \gamma}$ which is next to zero. For instance, it gets to 0.039 with $n = 9$ and a normal shrinking rate $\gamma = 30\%$.

## III. EXPERIMENTS

### A. Experiment setup

*1) Datasets:* The primary image dataset used in our experiments is the union of BOSSBase v1.01 [47] and BOWS2 [48], each of which contains 10,000 $512 \times 512$ grayscale spatial images. All of the images were resized to $256 \times 256$ using Matlab function *imresize*. The corresponding JPEG images were further generated with QFs (Quality Factors) 75 and 95. In our experiments, 10,000 BOWS2 images and 4,000 randomly selected BOSSBase images were used for training. Another 1,000 randomly selected BOSSBase images were for validation. The remaining 5,000 BOSSBase images were retained for testing.

ALASKA v2 [49] is another large-scale dataset with totally 80,005 images, introduced to evaluate the performance of STD-NET. The uncompressed and JPEG compressed (with quality factors: 75 and various QF) grayscale images datasets of size $256 \times 256$ were downloaded from ALASKA v2's official website[3]. We used training and validation sets of randomly selected 56,000 and 4,000 images, respectively, and the remaining images for testing. There are no identical images between the three sets.

*2) Steganography schemes:* Four representative steganographic schemes, UERD [5] and J-UNIWARD [6] for JPEG domain, and HILL [3] and S-UNIWARD [6] for spatial domain, were our attacking targets in the experiments. For JPEG steganographic algorithms, the embedding payloads were set to 0.2 and 0.4 bpnzAC (bits per non-zero AC DCT coefficient). For spatial domain steganographic algorithms, the embedding payloads were set to 0.2 and 0.4 bpp (bits per pixel).

As pointed out in [50], the performance of non-additive schemes varies dramatically on different cover sources. Therefore, we mainly focus on the detection of additive embedding distortion steganography algorithms. As for non-additive

---
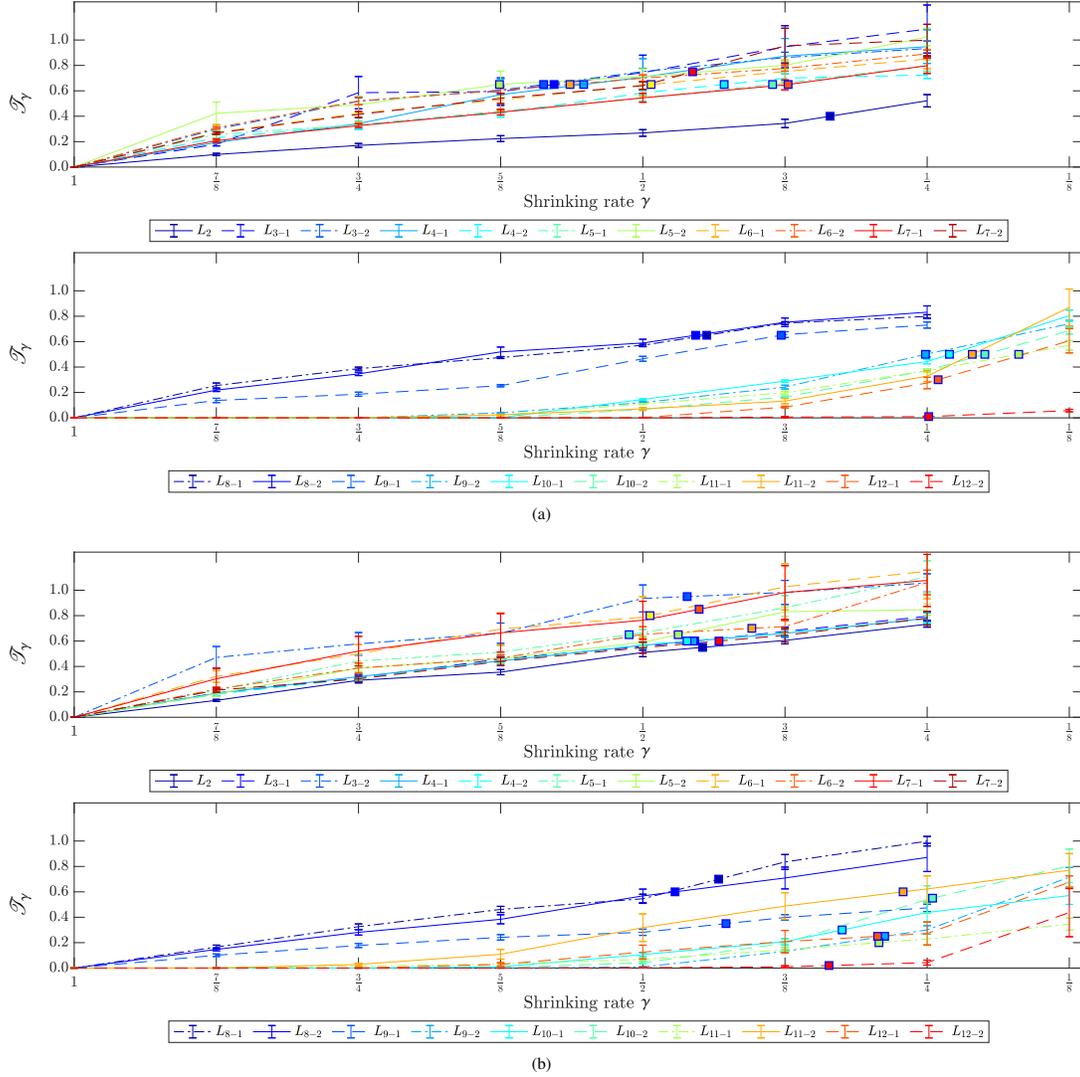
[3]https://alaska.utt.fr/#material

Fig. 5. (a) and (b) are the errorbars of normalized distortion vs. shrinking rate for every involved convolutional layer in the SRNetC64, conducted on JPEG-domain QF75 BOSSBase+BOWS2 dataset and JPEG-domain QF75 ALASKA v2 dataset, respectively. The square points with blue borders are the final determined $\mathscr{T}_d^l$.

schemes, the complementary experimental results can be found in Sect. III-F.

*3) Detectors:* SRNetC64 (a variant of SRNet, see Sect. II-B1) was selected as the initial architectures of our proposed STD-NET. Our implementation of STD-NET and its corresponding initial architectures are based on Tensorflow [51]. Unless otherwise specified, the initial architecture were trained with the hyperparameters mentioned in [33]. The batch size in the training procedure was set to 32 (namely 16 cover-stego pairs). The maximum numbers of iterations was set to $50 \times 10^4$ (about 571 epochs) on BOSSBase + BOWS2 dataset, and $80 \times 10^4$ (about 228 epochs) on ALASKA v2 dataset. The final searched STD-NET models adopted the same maximum number of iterations as the corresponding initial architecture. The optimizer used for the initial architectures and STD-NET was Adamax with initial learning rate 0.001. After $40 \times 10^4$ iterations the learning rate was reduced to 0.0001. In our experiments, the cylinder-shaped, ladder-shaped STD-NET architectures searched on JPEG-domain

QF75 BOSSBase+BOWS2 and cylinder-shaped STD-NET architecture searched on JPEG-domain QF75 ALASKA v2 are abbreviated as STD-BB-Cylinder, STD-BB-Ladder and STD-ALA-Cylinder, respectively.

The one with the best validation accuracy was evaluated on the corresponding testing set. All of the experiments were conducted on a GPU cluster with single NVIDIA Tesla P100 GPU card. Bounded by computational resources, every experiment was repeated three times, and the mean of the results on testing set were reported.

The source codes and auxiliary materials are available for download from GitHub. [4]

### B. Determination of the normalized distortion values

As described in Sect. II-B2, we need to determine the appropriate normalized distortion threshold $\mathscr{T}_d^l$ for every involved convolutional layer in the corresponding SRNetC64 model. In

[4]https://github.com/tansq/STD-NET

## TABLE I

Obtained structures of our proposed STD-NET under three different configurations and scenarios. Dimension parameters with ▮ (pink), ▮ (green) and ▮ (blue) backgrounds are from a cylinder-shaped (on BOSSBase+BOWS2), a ladder-shaped (on BOSSBase+BOWS2) and a cylinder-shaped configuration (on ALASKAv2/QF75/256×256/gray-scale dataset), respectively.

**SRNetC64**

| | $L_1$ | $L_2$ | $L_{3-1}$ | $L_{3-2}$ | $L_{4-1}$ | $L_{4-2}$ | $L_{5-1}$ | $L_{5-2}$ | $L_{6-1}$ | $L_{6-2}$ | $L_{7-1}$ | $L_{7-2}$ | $L_{8-1}$ | $L_{8-2}$ | $L_{9-1}$ | $L_{9-2}$ | $L_{10-1}$ | $L_{10-2}$ | $L_{11-1}$ | $L_{11-2}$ | $L_{12-1}$ | $L_{12-2}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\mathcal{K}^l$ | 1×64 | 64×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×16 | 16×64 | 64×64 | 64×64 | 64×64 | 64×64 | 64×64 | 64×64 | 64×64 |

**Corresponding STD-NET structures**

| | $L_1$ | $L_2$ | $L_{3-1}$ | $L_{3-2}$ | $L_{4-1}$ | $L_{4-2}$ | $L_{5-1}$ | $L_{5-2}$ | $L_{6-1}$ | $L_{6-2}$ | $L_{7-1}$ | $L_{7-2}$ | $L_{8-1}$ | $L_{8-2}$ | $L_{9-1}$ | $L_{9-2}$ | $L_{10-1}$ | $L_{10-2}$ | $L_{11-1}$ | $L_{11-2}$ | $L_{12-1}$ | $L_{12-2}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\mathbf{T}^l_\triangle$ | ~ | 64 64 64 × 32 32 24 | 16 16 16 × 7 7 9 | 16 16 16 × 11 7 4 | 16 16 16 × 7 6 9 | 16 16 16 × 8 8 12 | 16 16 16 × 8 11 9 | 16 16 16 × 9 11 9 | 16 16 16 × 8 4 6 | 16 16 16 × 11 14 9 | 16 16 16 × 9 5 8 | 16 16 16 × 10 14 9 | 16 16 16 × 9 7 10 | 16 16 16 × 6 5 9 | 16 16 16 × 3 3 9 | 64 64 64 × 13 11 27 | 64 64 64 × 13 9 24 | 64 64 64 × 11 12 17 | 64 64 64 × 10 7 23 | 64 64 64 × 13 11 25 | 64 64 64 × 15 10 25 | 64 64 64 × 16 23 19 |
| $\mathcal{G}^l$ | ~ | 32 32 24 × 8 8 6 | 7 7 9 × 7 7 9 | 11 7 4 × 11 15 4 | 7 6 9 × 7 8 9 | 8 8 12 × 8 8 12 | 8 11 9 × 8 5 9 | 9 11 9 × 9 7 9 | 8 4 6 × 8 12 6 | 11 14 9 × 11 8 9 | 9 5 8 × 9 13 8 | 10 14 9 × 10 6 9 | 9 7 10 × 9 11 10 | 6 5 9 × 6 7 9 | 3 3 9 × 12 12 36 | 13 11 27 × 13 15 27 | 13 9 24 × 13 17 24 | 11 12 17 × 11 10 17 | 10 7 23 × 10 13 23 | 13 11 25 × 13 15 25 | 15 10 25 × 15 20 25 | 16 23 19 × 16 9 19 |
| $\mathbf{T}^l_\triangle$ | ~ | 8 8 6 × 16 16 16 | 7 7 9 × 16 16 16 | 11 15 4 × 16 16 16 | 7 8 9 × 16 16 16 | 8 8 12 × 16 16 16 | 8 5 9 × 16 16 16 | 9 7 9 × 16 16 16 | 8 12 6 × 16 16 16 | 11 8 9 × 16 16 16 | 9 13 8 × 16 16 16 | 10 6 9 × 16 16 16 | 9 11 10 × 16 16 16 | 6 7 9 × 16 16 16 | 12 12 36 × 64 64 64 | 13 15 27 × 64 64 64 | 13 17 24 × 64 64 64 | 11 10 17 × 64 64 64 | 10 13 23 × 64 64 64 | 13 15 25 × 64 64 64 | 15 20 25 × 64 64 64 | 16 9 19 × 64 64 64 |
| | "L1" | "L2" | "L3" | | "L4" | | "L5" | | "L6" | | "L7" | | "L8" | | "L9" | | "L10" | | "L11" | | "L12" | |

## TABLE II

Comparison of parameters and FLOPs of our proposed STD-NET under three different configurations and scenarios. Those for CALPA-SRNet (with $\varsigma = 5\%$) are listed as well for comparison. The percentages of parameters and FLOPs of STD-NETs compared to CALPA-SRNet are shown in parentheses. Datas with ▮ (pink), ▮ (green) and ▮ (blue) backgrounds are from a cylinder-shaped (on BOSSBase+BOWS2), a ladder-shaped (on BOSSBase+BOWS2) and a cylinder-shaped configuration (on ALASKAv2/QF75/256×256/gray-scale dataset), respectively.

**CALPA-SRNet (with $\varsigma = 5\%$)**

**Parameters (×10⁴)**

| 6.93 | 0.05 | 0.92 | 0.18 | 0.18 | 0.14 | 0.14 | 0.05 | 0.05 | 0.14 | 0.14 | 0.20 | 0.20 | 0.14 | 0.13 | 0.37 | 0.55 | 0.11 | 0.06 | 0.41 | 0.85 | 0.56 | 1.14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**FLOPs (×10⁸)**

| 19.70 | 0.37 | 6.03 | 1.22 | 1.22 | 0.94 | 0.94 | 0.37 | 0.37 | 0.94 | 0.94 | 1.32 | 1.32 | 0.94 | 0.88 | 0.61 | 0.90 | 0.04 | 0.02 | 0.04 | 0.08 | 0.01 | 0.02 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Corresponding STD-NET structures for SRNet**

**Parameters (×10⁴)**

| Overall | $L_1$ | $L_2$ | $L_{3-1}$ | $L_{3-2}$ | $L_{4-1}$ | $L_{4-2}$ | $L_{5-1}$ | $L_{5-2}$ | $L_{6-1}$ | $L_{6-2}$ | $L_{7-1}$ | $L_{7-2}$ | $L_{8-1}$ | $L_{8-2}$ | $L_{9-1}$ | $L_{9-2}$ | $L_{10-1}$ | $L_{10-2}$ | $L_{11-1}$ | $L_{11-2}$ | $L_{12-1}$ | $L_{12-2}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4.97 (71.7%) | 0.05 | 0.44 | 0.06 | 0.14 | 0.06 | 0.08 | 0.08 | 0.10 | 0.08 | 0.14 | 0.10 | 0.12 | 0.10 | 0.05 | 0.11 | 0.31 | 0.31 | 0.24 | 0.21 | 0.31 | 0.39 | 0.43 |
| 4.79 (69.0%) | 0.05 | 0.44 | 0.06 | 0.12 | 0.06 | 0.08 | 0.07 | 0.09 | 0.06 | 0.13 | 0.08 | 0.10 | 0.09 | 0.05 | 0.11 | 0.31 | 0.30 | 0.24 | 0.20 | 0.31 | 0.37 | 0.39 |
| 8.40 (121.2%) | 0.05 | 0.29 | 0.10 | 0.02 | 0.10 | 0.16 | 0.10 | 0.10 | 0.05 | 0.10 | 0.08 | 0.10 | 0.12 | 0.10 | 0.53 | 1.00 | 0.82 | 0.47 | 0.77 | 0.88 | 0.88 | 0.56 |

**FLOPs (×10⁸)**

| Overall | $L_1$ | $L_2$ | $L_{3-1}$ | $L_{3-2}$ | $L_{4-1}$ | $L_{4-2}$ | $L_{5-1}$ | $L_{5-2}$ | $L_{6-1}$ | $L_{6-2}$ | $L_{7-1}$ | $L_{7-2}$ | $L_{8-1}$ | $L_{8-2}$ | $L_{9-1}$ | $L_{9-2}$ | $L_{10-1}$ | $L_{10-2}$ | $L_{11-1}$ | $L_{11-2}$ | $L_{12-1}$ | $L_{12-2}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 12.00 (60.9%) | 0.37 | 2.93 | 0.43 | 0.94 | 0.43 | 0.54 | 0.54 | 0.66 | 0.54 | 0.94 | 0.66 | 0.79 | 0.66 | 0.33 | 0.18 | 0.52 | 0.13 | 0.10 | 0.02 | 0.03 | 0.01 | 0.01 |
| 11.42 (57.9%) | 0.37 | 2.93 | 0.43 | 0.85 | 0.42 | 0.54 | 0.49 | 0.64 | 0.45 | 0.89 | 0.54 | 0.70 | 0.64 | 0.33 | 0.18 | 0.51 | 0.12 | 0.10 | 0.02 | 0.03 | 0.01 | 0.01 |
| 13.32 (67.6%) | 0.37 | 1.91 | 0.66 | 0.17 | 0.66 | 1.10 | 0.66 | 0.66 | 0.33 | 0.66 | 0.54 | 0.66 | 0.79 | 0.66 | 0.87 | 1.64 | 0.33 | 0.19 | 0.07 | 0.09 | 0.02 | 0.01 |

| | $L_1$ | $L_2$ | $L_{3-1}$ | $L_{3-2}$ | $L_{4-1}$ | $L_{4-2}$ | $L_{5-1}$ | $L_{5-2}$ | $L_{6-1}$ | $L_{6-2}$ | $L_{7-1}$ | $L_{7-2}$ | $L_{8-1}$ | $L_{8-2}$ | $L_{9-1}$ | $L_{9-2}$ | $L_{10-1}$ | $L_{10-2}$ | $L_{11-1}$ | $L_{11-2}$ | $L_{12-1}$ | $L_{12-2}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Overall** | "L1" | "L2" | "L3" | | "L4" | | "L5" | | "L6" | | "L7" | | "L8" | | "L9" | | "L10" | | "L11" | | "L12" | |

our experiments, the pre-defined step of shrinking rate was set to 5%, and the pre-defined lower bound of the shrinking rate $\tau$ was set to 20%. We start searching from the shrinking rate 50% for the bottom module, and 30% for the middle module and "L12" of SRNetC64.

In Fig. 5(a), we show how the normalized distortion $\mathscr{T}_\gamma$ changes with successive decreasing shrinking rates on JPEG-domain QF75 BOSSBase+BOWS2. Similarly, the change curves on JPEG-domain QF75 ALASKA v2 datasets are shown in Fig. 5(b). It can be seen that in most cases, as the shrinking rate decreases, the growth rate of $\mathscr{T}_\gamma$ increases. This trend is particularly obvious in the top module, and it also indicates that there is more redundancy in the top module. From Fig. 5, we can see that the finally determined distortion

threshold $\mathscr{T}_d^l$ of the upper layers are generally smaller than those of the lower layers.

### C. Compactness of STD-NET

In this section, we analyze the compactness and effectiveness of our proposed STD-NET. Based on the determined normalized distortion, we follow Algorithm. 2 to search for the corresponding number of input and output channels for the corresponding core tensor $\mathcal{G}^l$. Tab. I shows the three obtained structures, including STD-BB-Cylinder, STD-BB-Ladder and STD–ALA-Cylinder.

Tab. II shows the comparison of the parameters and FLOPs of the three different STD-NET structures with CALPA-SRNet. Compared with CALPA-SRNet, our proposed STD-

NET can compress SRNet more effectively. Parameters and FLOPs of the STD-NETs are further reduced by about 70% and 60%, respectively, except for the parameters of STD-ALA-Cylinder. Specifically, all the involved convolutional layers in STD-BB-Cylinder and STD-BB-Ladder get better compression ratios compared with CALPA-NET, except those in "L5" and "L10". As for STD-ALA-Cylinder, it has more parameters than other comparison models, since the ALASKA v2 dataset is more complex than BOSSBase+BOWS2 and requires more model parameters to maintain detection performance.

It is highlighted that compared with the original SRNet, parameters and FLOPs of STD-BB-Cylinder, STD-BB-Ladder, and STD-ALA-Cylinder are further reduced by 1.04% and 20.16%, 1.00% and 19.18%, and 1.76% and 22.38%, respectively.

### D. Detection performance of STD-NETs

In Tab. III, we compare the detection performance of the three STD-NETs, the corresponding CALPA-SRNet and the original SRNet on eighteen different scenarios.

In the JPEG domain, the detection performance of the STD-BB-Cylinder and STD-BB-Ladder on almost all involved scenarios outperforms CALPA-SRNet. In particular, the cylinder-shaped STD-NETs even exceed the original SRNet on JPEG-domain ALASKA v2 dataset with various quality factors.

In the spatial domain, the two variants of STD-NETs (STD-BB-Cylinder and STD-BB-Ladder), originally searched on JPEG domain dataset (QF75 BOSSBase+BOWS2 aiming at J-UNIWARD 0.4 bpnzAC), still achieve similar or even better detection performance compared with CALPA-SRNet, as well as original SRNet.

In the following experiments, unless otherwise specified, we only report the results for those STD-NET models with the structures obtained on the JPEG-domain QF75 BOSS-Base+BOWS2 dataset.

Fig. 6 shows the performance comparison of the STD-BB-Cylinder, STD-BB-Ladder STD-NETs and original SRNet, on JPEG-domain QF75 BOSSBase+BOWS2 dataset aiming at J-UNIWARD/0.4 bpnzAC. We can see that the cylinder and ladder-shaped STD-NETs show similar validation and testing performance compared with original SRNet. It is worth noting that the larger the model is, the bigger the gap between training accuracies and validation accuracies for the corresponding model is. The gaps of SRNet, STD-BB-Cylinder, and STD-BB-Ladder are clearly reduced in accordance with their model scale, which has been sorted from the largest to the smallest.

From Tab. III and Fig. 6, we can see our STD-NET architectures, especially STD-BB-Cylinder and STD-BB-Ladder, can obtain competitive detection performance in most scenarios. For instance, compared with original SRNet, the two STD-NET architectures can achieve competitive accuracy even if it is trained from scratch (see Fig. 6). Therefore, we can make a conclusion that our STD-NET has a great adaptivity.

Fig. 7 shows comparison of training accuracies and validation accuracies vs. training iterations for the original SR-Net and two STD-NET models with two different training strategy. From Fig. 7 we can see that both training strategy,

namely "training-decomposing-finetuning" and "trained-from-scratch", guarantee the STD-NET models achieve similar performance when compared with SRNet on validation set. The STD-NET model with "training-decomposing-finetuning" pipeline converges faster than other models on the training set.

All in all, on one hand, because our STD-NET is an adaptive architecture, we can directly train from scratch without using an existing model. On the other hand, the STD-NET model that preserves the model parameters after Tucker decomposition has a certain expression and discriminative ability, which helps to converge faster on varying steganalysis tasks.

### E. Detection performance on non-additive steganography

To further investigate our STD-NET models' detection ability, we evaluate its detection performance on BOSS-Base+BOWS2, targeting at non-additive steganography. In spatial domain, we adopt HILL as basic additive cost function with the state-of-the-art non-additive steganographic framework CMD [8] at payload of 0.4 bpp. In JPEG domain, J-UNIWARD is chosen as basic additive steganography with non-additive steganographic frameworks BBM [12] and BBC-BBM (the combination of BBC [10] and BBM).

In Tab. IV, we compare the detection performance of the STD-NET models STD-BB-Cylinder and STD-BB-Ladder with the original SRNet. Please note the two STD-NET models' parameters are finetuned after Tucker decomposition. From Tab. IV we can see the fine-tuned STD-NET model gets similar or even better performance compared with SRNet. The performance of the two models, STD-BB-Cylinder and STD-BB-Ladder, is very close. Compared with SRNet, their accuracy losses do not exceed 0.2%. Further on, targeting at BBM, the performance of our two STD-NET models is even 0.11% and 0.05% higher than that of SRNet model, respectively.

### F. Generalization of STD-NET

To evaluate the generalization of STD-NET, we further implement our STD-NET strategy for SiaStegNet [36]. We trained the original network and the searched versions on BOSSBase + BOWS2 dataset aiming at S-UNIWARD 0.4 bpp for 500 epochs with batch size 32 (namely 16 cover-stego pairs). The optimizer used for the original architecture and its corresponding STD-NET versions was Adamax with initial learning rate 0.001 and weight decay 0.0001. All other hyperparameters and the training schedule were same as described in [36]. The two final obtained architectures are denoted as STD-Sia-Cylinder and STD-Sia-Ladder, standing for the cylinder-shaped and ladder-shaped STD-NET version of SiaStegNet, respectively.

The comparison results of model size and performance before and after decomposition are shown in Tab. V. It can be seen that the number of parameters and FLOPs of STD-Sia-Cylinder and STD-Sia-Ladder are around 17% and 25% of the original respectively. Those fine-tuned STD-NET models with decomposed parameters retained can even gain better performance compared with the original SiaStegNet model. Besides, though those STD-NET models trained from scratch

TABLE III

Comparison of detection performance of our proposed STD-NET under three different configurations and scenarios, and the corresponding CALPA-SRNet and the original SRNet. Datas with pink, green and blue backgrounds are from a cylinder-shaped (on BOSSBase+BOWS2), a ladder-shaped (on BOSSBase+BOWS2) and a cylinder-shaped configuration (on ALASKAv2/QF75/256×256/gray-scale dataset), respectively. Datas with blue frames are from the corresponding STD-NET model following curriculum learning method.

| Scenarios | | Targets | | STD-NETs for SRNet | | CALPA-SRNet | | SRNet | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | $P_{MD}$(5%) | wAUC | $P_{MD}$(5%) | wAUC | $P_{MD}$(5%) | wAUC |
| BOSSBase+BOWS2 | JPEG (QF75) | J-UNIWARD | 0.4 bpnzAC | 10.9% | 98.3% (↓0.4%) | 10.3% | 98.5% (↓0.14%) | 9.80% | 98.7% |
| | | | | 12.4% | 98.4% (↓0.3%) | | | | |
| | | | 0.2 bpnzAC | 47.8% | 91.1% (↓1.2%) | 43.9% | 92.3% (↑0.04%) | 41.2% | 92.3% |
| | | | | 47.2% | 90.7% (↓1.6%) | | | | |
| | | UERD | 0.4 bpnzAC | 1.6% | 99.7% (↓0.0%) | 1.9% | 99.62% (↓0.08%) | 1.7% | 99.7% |
| | | | | 2.0% | 99.7% (↓0.0%) | | | | |
| | | | 0.2 bpnzAC | 16.0% | 97.9% (↓0.0%) | 17.0% | 97.7% (↓0.2%) | 15.4% | 97.9% |
| | | | | 19.8% | 97.3% (↓0.6%) | | | | |
| | Spatial | MiPOD | 0.4 bpp | 34.8% | 94.8% (↓1.0%) | 50.0% | 94.4% (↓1.4%) | 27.9% | 95.8% |
| | | | | 33.0% | 95.1% (↓0.7%) | | | | |
| | | | 0.2 bpp | 57.8% | 87.1% (↑1%) | 57.0% | 86.0% (↓0.1%) | 56.9% | 86.1% |
| | | | | 57.4% | 87.0% (↑0.9%) | | | | |
| | | HILL | 0.4 bpp | 31.4% | 95.3% (↓0.7%) | 28.0% | 96.1% (↑0.13%) | 27.7% | 96.0% |
| | | | | 30.6% | 95.8% (↓0.2%) | | | | |
| | | | 0.2 bpp | 53.8% | 88.6% (↓1.1%) | 49.4% | 89.5% (↓0.15%) | 50.6% | 89.7% |
| | | | | 53.3% | 89.3% (↓0.4%) | | | | |
| ALASKA v2 | JPEG (QF75) | J-UNIWARD | 0.4 bpnzAC | 33.9% | 94.9% (↓1.2%) | 37.8% | 93.1% (↓3.0%) | 26.6% | 96.1% |
| | | | | 33.2% | 94.9% (↓1.2%) | | | | |
| | | | 0.2 bpnzAC | 67.1% | 84.0% (↓4.0%) | 76.5% | 80.5% (↓7.5%) | 84.1% | 88.0% |
| | | | | 68.8% | 83.1% (↓4.9%) | | | | |
| | | UERD | 0.4 bpnzAC | 16.4% | 96.4% (↓0.3%) | 50.8% | 92.8% (↓3.9%) | 16.0% | 96.7% |
| | | | | 20.5% | 96.2% (↓0.5%) | | | | |
| | | | 0.2 bpnzAC | 43.7% | 92.5% (↑1.5%) | 65.3% | 72.1% (↓18.9%) | 45.0% | 91.0% |
| | | | | 47.9% | 90.5% (↓0.5%) | | | | |
| | JPEG (various) | J-UNIWARD | 0.4 bpnzAC | 64.5% | 84.8% (↑0.3%) | 73.5% | 83.2% (↓1.3%) | 65.0% | 84.5% |
| | | | | 65.3% | 84.4% (↓0.1%) | | | | |
| | | | 0.2 bpnzAC | 86.1% | 71.7% (↑1.7%) | 93.1% | 65.2% (↓4.8%) | 89.6% | 70.0% |
| | | | | 93.6% | 61.6% (↓8.4%) | | | | |
| | | UERD | 0.4 bpnzAC | 50.1% | 90.5% (↓0.5%) | 57.7% | 87.0% (↓4.0%) | 49.0% | 91.0% |
| | | | | 49.1% | 89.6% (↓1.4%) | | | | |
| | | | 0.2 bpnzAC | 71.5% | 80.9% (↑0.9%) | 87.6% | 75.0% (↓5.0%) | 80.1% | 80.0% |
| | | | | 93.3% | 61.4% (↓18.6%) | | | | |

suffer losses of detection accuracy, the losses are still within a tolerable range.

Generally speaking, our proposed STD-NET approach is available for all the deep-learning steganalyzers equipped with convolutional layers, and its applicative scenarios are rich. Besides the results dealing with QF75 JPEG images, we have demonstrated its effectiveness in following scenarios:

- JPEG-domain QF90 256 × 256 scenario;
- JPEG-domain QF75 512 × 512 scenario.

As shown in Tab. VI, even in the above two harder scenarios, the detection performance of the cylinder-shaped as well as the ladder-shaped STD-NETs searched on those scenarios can be close to the original SRNet, and outperforms CALPA-SRNet.

We have compared our proposed unsupervised data-driven criterion with manually selected core-tensor ranks used in vanilla Tucker decomposition. We name the obtained cylinder-shaped STD-NET architecture for JPEG-domain QF75 BOSS-Base +BOWS2 aiming at J-UNIWARD/0.4 bpnzAC, with
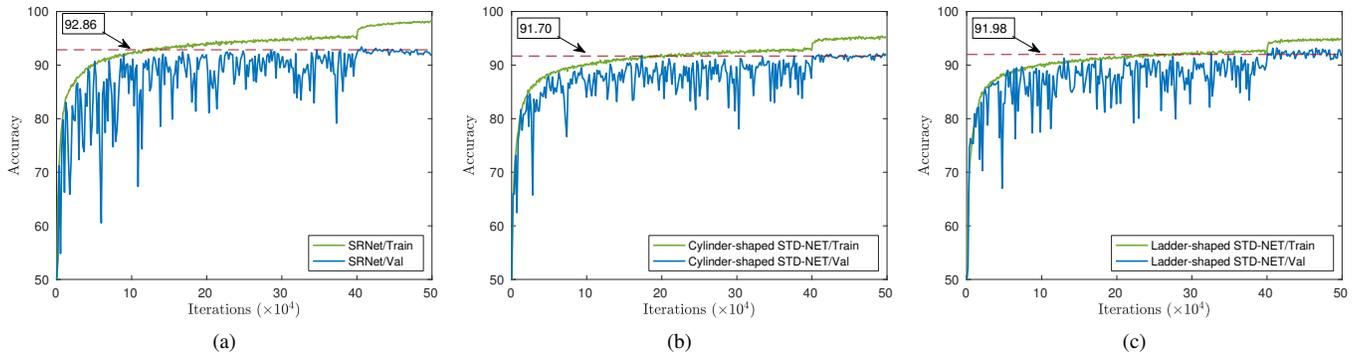
Fig. 6. Comparison of the training accuracies and validation accuracies vs. training iterations for SRNet and the corresponding STD-NETs (STD-BB-Cylinder and STD-BB-Ladder). The dashed line in every sub-figure marks the testing accuracy achieved by the trained model with highest validation accuracy. (a), (b), and (c) is the result for original SRNet, cylinder-shaped STD-NET, and ladder-shaped STD-NET, respectively. The experiments were conducted on JPEG-domain QF75 BOSSBase+BOWS2 dataset, aiming at J-UNIWARD/0.4 bpnzAC.
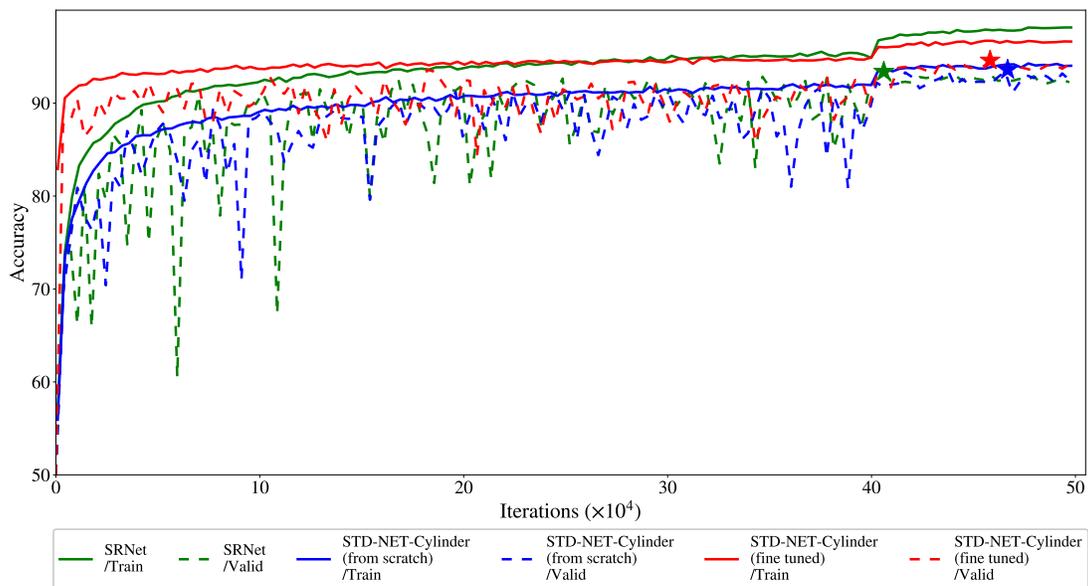


Fig. 7. Comparison of the training accuracies and validation accuracies vs. training iterations for the original SRNet, corresponding STD-NET trained from scratch, and STD-NET with "training-decomposing-finetuning" pipeline. The experiment was conducted on BOSSBase+BOWS2 dataset aiming at J-UNIWARD/0.4 bpnzAC. "★" denotes the point with highest value at the corresponding plot.

TABLE IV
COMPARISON OF DETECTION ACCURACY FOR SRNET AND THE CORRESPONDING STD-NET MODELS (WITH THE "TRAINING-DECOMPOSING-FINETUNING" PIPELINE) SEARCHED ON JPEG-DOMAIN QF75 BOSSBASE+BOWS2 AT NON-ADDITIVE STEGANOGRAPHY.

| Model | HILL 0.4 bpp | J-UNIWARD 0.4 bpnzAC | |
|---|---|---|---|
| | CMD | BBM | BBC-BBM |
| SRNet | 81.99 | 91.63 | 91.40 |
| STD-BB-Cylinder (fine-tuned) | 81.92 ($\downarrow$0.07) | 91.74 ($\uparrow$0.11) | 91.22 ($\downarrow$0.18) |
| STD-BB-Ladder (fine-tuned) | 81.94 ($\downarrow$0.05) | 91.68 ($\uparrow$0.05) | 91.20 ($\downarrow$0.20) |

vanilla Tucker decomposition as STD-BB-FIXED, in which the shrinking rate for the bottom module is fixed to 50%, and the shrinking rate for the middle module and "L12" is fixed to 30%.

As shown in Fig. VI, It is no doubt that the cylinder-shaped as well as the ladder-shaped STD-NETs searched

on those scenarios outperforms the corresponding STD-NETs with vanilla Tucker decomposition.

### G. Further discussions

*1) Impact of no pair constraint:* Experimental results reported in [52]–[54] indicate that training with pair constraint

TABLE V
THE COMPARISON OF MODEL SIZE AND PERFORMANCE BEFORE AND AFTER DECOMPOSITION FOR SIASTEGNET THROUGH OUR PROPOSED STD-NET .

| Model | Params (M) | FLOPs (M) | ACC (%) |
|---|---|---|---|
| SiaStegNet | 1.41 | 7159.9 | 88.83 |
| STD-Sia-Cylinder (fine-tuned) | 0.25 (17.86%) | 1874.46 (26.17%) | 89.42 ($\uparrow$0.59) |
| STD-Sia-Cylinder (from scratch) | 0.25 (17.86%) | 1874.46 (26.17%) | 88.58 ($\downarrow$0.25) |
| STD-Sia-Ladder (fine-tuned) | 0.22 (16.01%) | 1739.79 (24.29%) | 89.32 ($\uparrow$0.49) |
| STD-Sia-Ladder (from scratch) | 0.22 (16.01%) | 1739.79 (24.29%) | 86.84 ($\downarrow$1.99) |

TABLE VI
COMPARISON OF DETECTION PERFORMANCE FOR SRNET AND THE CORRESPONDING STD-NET MODELS (WITH THE "TRAIN-FROM-SCRATCH" PIPELINE) AIMING AT J-UNIWARD 0.4 BPNZAC. THE DATASET IS BOSSBASE+BOWS2. THE CYLINDER-SHAPED AS WELL AS THE LADDER-SHAPED STD-NETS ARE SEARCHED ON THE CORRESPONDING SCENARIOS.

| Model | QF90/256 × 256 | | QF75/512 × 512 | |
|---|---|---|---|---|
| | $P_{MD}(5\%)$ | wAUC | $P_{MD}(5\%)$ | wAUC |
| SRNet | 20.0 | 97.4 | 9.8 | 98.1 |
| CALPA-NET | 24.60 | 96.73 ($\downarrow$**0.67**) | 9.36 | 97.69 ($\downarrow$**0.41**) |
| STD-NET cylinder-shaped | 21.91 | 96.94 ($\downarrow$**0.46**) | 5.4 | 98.5 ($\uparrow$**0.4**) |
| STD-NET ladder-shaped | 21.24 | 96.97 ($\downarrow$**0.43**) | 13.7 | 96.6 ($\downarrow$**1.5**) |
| STD-BB-FIXED | 30.2 | 95.4 ($\downarrow$**2**) | 13.6 | 97.2 ($\downarrow$**0.9**) |

TABLE VII
COMPARISON OF DETECTION PERFORMANCE FOR SRNET AND THE CORRESPONDING STD-NET MODELS (WITH THE "TRAIN-FROM-SCRATCH" PIPELINE SEARCHED ON JPEG-DOMAIN QF75 BOSSBASE+BOWS2) WITHOUT PAIR CONSTRAINT IN THE TRAINING PROCEDURE. THE RESULTS FOR IMAGENET PRETRAINED EFFICIENTNETB0 ARE REPORTED AS WELL. THE UNDERLINED RESULTS ARE OBTAINED WITH CURRICULUM LEARNING.

| Model | BB-QF75-JU04 | | BB-QF75-UE04 | | BB-QF90-JU04 | | ALA2-QF75-JU04 | |
|---|---|---|---|---|---|---|---|---|
| | $P_{MD}(5\%)$ | wAUC | $P_{MD}(5\%)$ | wAUC | $P_{MD}(5\%)$ | wAUC | $P_{MD}(5\%)$ | wAUC |
| SRNet | 12.16 | 98.57 | 2.60 | 99.64 | 25.08 | 96.71 | 26.60 | 96.10 |
| EfficientNet-B0 | 66.6 | 86.63 | 58.20 | 92.10 | <u>86.26</u> | <u>75.46</u> | 29.01 | 95.12 |
| STD-BB-Cylinder | 10.84 | 98.53 | 2.11 | 99.66 | 29.20 | 96.01 | 33.91 | 94.91 |
| STD-BB-Ladder | 16.55 | 97.89 | 2.00 | 99.63 | 35.32 | 94.52 | 33.2 | 94.93 |

removed in the later stage may improve the detection performance. An experiment has been conducted on BOSS-Base+BOWS2 dataset as well as ALASKA v2 dataset to verify the impact of no pair constraint on the performance of our proposed STD-NET approach. In the experiment all the involved models were trained with 300 pair-constraint epochs and then 100 no-pair-constraint epochs. The results for a ImageNet pretrained EfficientNet-B0 model are reported as well.

compared the results in Tab. VII with those in Tab. III, we can see that on BOSSBase+BOWS2 dataset training without pair constraint does actually not introduce any relative advantage for the more complex SRNet models. Compared with the results reported in Tab.III of the revision, additional no-pair-constraint training epochs do mildly degrade the detection performance of SRNet models, as well as our proposed STD-NET models. Anyhow, on BOSSBase+BOWS2 dataset the extent

of the impact of no-pair-constraint training on SRNet and our proposed STD-NET is similar. Since no obvious impact can be observed, we have sticked to using the more straightforward pair-constraint training on BOSSBase+BOWS2 dataset, , and reported obtained results in Tab. III.

On the contrary, in the experiments we have observed that pair constraint does introduce obvious detection performance improvements for SRNet models on ALASKA v2 dataset. Now the huge gap between training and validation for SRNet models trained with pair constraint is markedly narrowed. Accordingly, we have removed pair constraint in the later stage of the training procedure on ALASKA v2 dataset, and report those obtained results in Tab. III.

However, please note that as reported in Tab. III, on ALASKA v2 dataset, even with pair constraint removed in the later stage of the training procedure, by and large the detection performances of our proposed STD-NETs are still comparable

to those of the original SRNets.

*2) STD-NET and EfficientNet family:* Please note that though as reported in [54], the EfficientNet family [37] pretrained on ImageNet can achieve state-of-the-art performance for JPEG steganalysis, STD-NET cannot be directly applied to EfficientNet to effectively shrink it. This is due to the fact that EfficientNets are mainly composed of *depthwise separable convolutional layers*, which cannot be handled with existing tensor decomposition techniques. Depthwise separable convolution is one of the most important components of EfficientNet family. However, depthwise separable convolutions sharply reduce the complexity margin which Tucker decomposition can cut off.

Given a block with depthwise separable convolution, it takes an input tensor $\mathcal{I}^{l-1} \in \mathbb{R}^{J^{l-1} \times H^{l-1} \times W^{l-1}}$ and maps it into an output tensor $\mathcal{O}^l \in \mathbb{R}^{J^l \times H^l \times W^l}$ with a depthwise convolution followed by a pointwise $1 \times 1$ convolution. Let $\mathcal{K}_d^l \in \mathbb{R}^{J^l \times D^l \times D^l}$ denotes the depthwise convolution kernel tensor and $\mathcal{K}_p^l \in \mathbb{R}^{J^{l-1} \times J^l}$ denotes the pointwise convolution kernel tensor. Element-wise (definitions of $h_k$ and $w_s$ are the same as those in Eq. (1), and are omitted here for clarity):

$$\mathcal{O}_{i,h,w}^l = \sum_{j=1}^{J^{l-1}} \mathcal{K}_{p\,j,i}^l \underbrace{\sum_{k=1}^{D^l} \sum_{s=1}^{D^l} \mathcal{K}_{d\,j,k,s}^l \cdot \mathcal{I}_{j,h_k,w_s}^{l-1}}_{\text{depthwise conv}} \qquad (9)$$

$$\underbrace{\phantom{\sum}}_{\text{pointwise conv}}$$

Measured with the number of parameters, the computational complexity of the block turns to:

$$J^{l-1} \cdot D^l \cdot D^l + J^{l-1} \cdot J^l \qquad (10)$$

Assuming that we apply Tucker decomposition to $\mathcal{K}_d^l$, the resulting complexity turns to:

$$J^{l-1} \cdot I^l + I^l \cdot O^l \cdot D^l \cdot D^l + O^l + J^{l-1} \cdot J^l \qquad (11)$$

Then subtract Eq. (11) from Eq. (10), we turn to solve the following inequality:

$$J^{l-1} \cdot D^l \cdot D^l - J^{l-1} \cdot I^l - I^l \cdot O^l \cdot D^l \cdot D^l - O^l \geq 0 \qquad (12)$$

Again we set $D^l \cdot D^l = n$, $J^{l-1} = J^l = x$, and $I^l = O^l = x\gamma$, $\gamma \in [0, 1]$. Eq. (12) turns to:

$$nx - x^2\gamma - nx^2\gamma^2 - x\gamma \geq 0 \Rightarrow nx\gamma^2 + (x+1)\gamma - n \leq 0 \qquad (13)$$

Regarding the left side of Eq. (13) as a univariate quadratic function of $\gamma$, we can get that Eq. (13) holds in $\gamma \in [0, \frac{-(x+1)+\sqrt{(x+1)^2+4n^2x}}{2nx}]$. Since $\sqrt{(x+1)^2 + 4n^2x} \leq (x+1) + 2n\sqrt{x}$, we can get $\frac{-(x+1)+\sqrt{(x+1)^2+4n^2x}}{2nx} \leq \frac{1}{\sqrt{x}}$. As a result, $\gamma$ actually lies in a quite narrow range included in $[0, \frac{1}{\sqrt{x}}]$. For instance, with $n = 3 \times 3 = 9$ and $x = 64$, Tucker decomposition can only bring in computational complexity reduction when the shrinking rate $\gamma < 8.1\%$. Since EfficientNet family has been well designed to be very compact given a number of parameters, With such an aggressive shrinking rate, Tucker decomposition will certainly severely degrade the performance of the resulting model.

However, according to the experimental results reported in [52]–[54], the EfficientNet family is still far from the best framework for image steganalysis. Firstly, it even cannot converge without pre-trained model weights obtained in other large-scale tasks. Secondly, even with pre-trained model weights, vanilla EfficientNet is inferior to SRNet in gray-scale image steganalysis as well as in aggressively down-sampled images. Thirdly, though designed for compactness, EfficientNet family is actually not superior in complexity even compared with SRNet, not to mention our proposed STD-NETs.

The experimental results in Tab. VII provide additional evidences. Tab. VII shows that the ImageNet pretrained EfficientNet-B0 model is inferior even with additional no-pair-constraint training epochs on BOSSBase+BOWS2 dataset. On ALASKA v2 dataset its performance cannot surpass that of SRNet. Certainly, the reported results are obtained for vanilla EfficientNet on gray-scale ALASKA v2 images. As reported in [54], specific renovations can improve the detection performance of EfficientNet in image steganalysis. However, all of those renovations increase its complexity by a wide margin.

To summarize, in a laboratory setting state-of-the-art deep-learning steganalyzers can detect stego images with satisfactory accuracies in quite a few scenarios. Therefore we are well past the stage that the scholars were only chasing performance improvements. Besides detection performance, computational complexity is another dimension to measure the overall virtue of a deep-learning steganalyzer. Our proposed STD-NET approach can achieve state-of-the-art performance on top of various notable architectures with even an order of magnitude smaller complexity. It can be highly complementary to the EfficientNet family in the field of image steganalysis.

## IV. Concluding remarks

In this paper, we propose STD-NET, which is aiming at searching an efficient image steganalytic deep-learning architecture to save the memory cost as well as the computational cost. The major contributions of this work are as follows:

- We have proposed a hierarchical tensor decomposition strategy, which can greatly reduce model parameters and FLOPs through Tucker decomposition. Unlike CALPA-NET, the STD-NET will not be restricted by various residual shortcut connections.
- We have proposed a normalized distortion threshold, which guide us to decompose involved convolutional layers on the basis of the original SRNet model in an unsupervised way, so as to search for an efficient and adaptive deep-learning image steganalysis architecture.
- The extensive experiments conducted on de-facto bench-marking image datasets show that our STD-NET models achieve comparative detection performance whether it is finetuned with the decomposed parameters or trained from scratch.

In the future, we will mainly focus on two aspects: 1) explore fully automatic deep steganalytic neural architecture

search strategies; 2) explore broader applications of our proposed tensor decomposition strategy in other areas of media security and forensics.

## References

[1] R. Böhme, *Advanced Statistical Steganalysis*, 1st ed. Springer Publishing Company, Incorporated, 2010.

[2] T. Filler and J. Fridrich, "Gibbs construction in steganography," *IEEE Transactions on Information Forensics and Security*, vol. 5, no. 4, pp. 705–720, 2010.

[3] B. Li, M. Wang, J. Huang, and X. Li, "A new cost function for spatial image steganography," in *Proc. IEEE 2014 International Conference on Image Processing, (ICIP'2014)*, 2014, pp. 4206–4210.

[4] V. Sedighi, R. Cogranne, and J. Fridrich, "Content-adaptive steganography by minimizing statistical detectability," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 2, pp. 221–234, 2016.

[5] L. Guo, J. Ni, W. Su, C. Tang, and Y. Q. Shi, "Using statistical image model for JPEG steganography: Uniform embedding revisited," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 12, pp. 2669–2680, 2015.

[6] V. Holub, J. Fridrich, and T. Denemark, "Universal distortion function for steganography in an arbitrary domain," *EURASIP Journal on Information Security*, vol. 2014, no. 1, pp. 1–13, 2014.

[7] T. Denemark and J. Fridrich, "Improving steganographic security by synchronizing the selection channel," in *Proc. 3rd ACM Information Hiding and Multimedia Security Workshop (IH&MMSec' 2015)*, 2015, pp. 5–14.

[8] B. Li, M. Wang, X. Li, S. Tan, and J. Huang, "A strategy of clustering modification directions in spatial image steganography," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 9, pp. 1905–1917, 2015.

[9] W. Zhang, Z. Zhang, L. Zhang, H. Li, and N. Yu, "Decomposing joint distortion for adaptive steganography," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 27, no. 10, pp. 2274–2280, 2017.

[10] W. Li, W. Zhang, K. Chen, W. Zhou, and N. Yu, "Defining joint distortion for JPEG steganography," in *Proc. 6th ACM Information Hiding and Multimedia Security Workshop (IH&MMSec' 2018)*, 2018, pp. 5–16.

[11] Y. Wang *et al.*, "BBC++: Enhanced block boundary continuity on defining non-additive distortion for jpeg steganography," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 5, pp. 2082–2088, 2021.

[12] Y. Wang, W. Zhang, W. Li, and N. Yu, "Non-additive cost functions for JPEG steganography based on block boundary maintenance," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 1117–1130, 2021.

[13] W. Tang, S. Tan, B. Li, and J. Huang, "Automatic steganographic distortion learning using a generative adversarial network," *IEEE Signal Processing Letters*, vol. 24, no. 10, pp. 1547–1551, 2017.

[14] W. Tang, B. Li, S. Tan, M. Barni, and J. Huang, "CNN-based adversarial embedding for image steganography," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 8, pp. 2074–2087, 2019.

[15] J. Yang, D. Ruan, J. Huang, X. Kang, and Y. Shi, "An embedding cost learning framework using GAN," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 839–851, 2020.

[16] M. Yedroudj, F. Comby, and M. Chaumont, "Steganography using a 3-player game," *Journal of Visual Communication and Image Representation*, vol. 72, p. 102910, 2020.

[17] X. Mo, S. Tan, B. Li, and J. Huang, "MCTSteg: A Monte Carlo tree search-based reinforcement learning framework for universal non-additive steganography," *IEEE Transactions on Information Forensics and Security*, pp. 1–1, 2021.

[18] J. Fridrich and J. Kodovský, "Rich models for steganalysis of digital images," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 3, pp. 868–882, 2012.

[19] T. Denemark, V. Sedighi, V. Holub, R. Cogranne, and J. Fridrich, "Selection-channel-aware rich model for steganalysis of digital images," in *Proc. 6th IEEE International Workshop on Information Forensic and Security (WIFS'2014)*, 2014, pp. 48–53.

[20] W. Tang, H. Li, W. Luo, and J. Huang, "Adaptive steganalysis based on embedding probabilities of pixels," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 4, pp. 734–745, 2016.

[21] S. Tan, H. Zhang, B. Li, and J. Huang, "Pixel-decimation-assisted steganalysis of synchronize-embedding-changes steganography," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 7, pp. 1658–1670, 2017.

[22] B. Li, Z. Li, S. Zhou, S. Tan, and X. Zhang, "New steganalytic features for spatial image steganography based on derivative filters and threshold LBP operator," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 5, pp. 1242–1257, 2018.

[23] J. Kodovský and J. Fridrich, "Ensemble classifiers for steganalysis of digital media," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 2, pp. 432–444, 2012.

[24] S. Tan and B. Li, "Stacked convolutional auto-encoders for steganalysis of digital images," in *Proc. Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA'2014)*, 2014, pp. 1–4.

[25] Y. Qian, J. Dong, W. Wang, and T. Tan, "Deep learning for steganalysis via convolutional neural networks," in *Proc. IS&T/SPIE Electronic Imaging 2015 (Media Watermarking, Security, and Forensics)*, 2015, pp. 94 090J–1–94 090J–10.

[26] L. Pibre, P. Jérôme, D. Ienco, and M. Chaumont, "Deep learning is a good steganalysis tool when embedding key is reused for different images, even if there is a cover source-mismatch," in *Proc. Media Watermarking, Security, and Forensics, Part of IS&T International Symposium on Electronic Imaging (EI'2016)*, 2016, pp. 1–11.

[27] G. Xu, H. Z. Wu, and Y. Q. Shi, "Structural design of convolutional neural networks for steganalysis," *IEEE Signal Processing Letters*, vol. 23, no. 5, pp. 708–712, 2016.

[28] J. Ye, J. Ni, and Y. Yi, "Deep learning hierarchical representations for image steganalysis," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 11, pp. 2545–2557, 2017.

[29] M. Chen, V. Sedighi, M. Boroumand, and J. Fridrich, "JPEG-phase-aware convolutional neural network for steganalysis of JPEG images," in *Proc. 5th ACM Information Hiding and Multimedia Security Workshop (IH&MMSec'2017)*, 2017, pp. 75–84.

[30] G. Xu, "Deep convolutional neural network to detect J-UNIWARD," in *Proc. 5th ACM Information Hiding and Multimedia Security Workshop (IH&MMSec'2017)*, 2017, pp. 67–73.

[31] J. Zeng, S. Tan, B. Li, and J. Huang, "Large-scale JPEG steganalysis using hybrid deep-learning framework," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 5, pp. 1242–1257, 2018.

[32] J. Zeng, S. Tan, G. Liu, B. Li, and J. Huang, "WISERNet: Wider separate-then-reunion network for steganalysis of color images," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 10, pp. 2735–2748, 2019.

[33] M. Boroumand, M. Chen, and J. Fridrich, "Deep residual network for steganalysis of digital images," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 5, pp. 1181–1193, 2019.

[34] M. Yedroudj, F. Comby, and M. Chaumont, "Yedroudj-net: An efficient CNN for spatial steganalysis," in *Proc. 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'2018)*, 2018, pp. 2092–2096.

[35] R. Zhang, F. Zhu, J. Liu, and G. Liu, "Depth-wise separable convolutions and multi-level pooling for an efficient spatial CNN-based steganalysis," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 1138–1150, 2020.

[36] W. You, H. Zhang, and X. Zhao, "A siamese CNN for image steganalysis," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 291–306, 2021.

[37] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proc. International Conference on Machine Learning (ICML' 2019)*, 2019, pp. 6105–6114.

[38] H. Ruiz *et al.*, "Analysis of the scalability of a deep-learning network for steganography "into the wild"," in *Proc. 2021 International Conference on Pattern Recognition (ICPR'2021), Worshop on MultiMedia FORensics in the WILD (MMForWILD'2021)*, 2021, pp. 439–452.

[39] C. Liu *et al.*, "Progressive neural architecture search," in *Proc. 2018 European Conference on Computer Vision (ECCV'2018)*, 2018, pp. 19–34.

[40] J.-H. Luo, J. Wu, and W. Lin, "ThiNet: A filter level pruning method for deep neural network compression," in *Proc. 2017 IEEE international conference on computer vision, (ICCV'2017)*, 2017, pp. 5058–5066.

[41] Z. Liu, M. Sun, T. Zhou, G. Huang, and T. Darrell, "Rethinking the value of network pruning," in *Proc. 2019 International Conference on Learning Representations, (ICLR'2019)*, 2019.

[42] S. Tan *et al.*, "CALPA-NET: Channel-pruning-assisted deep residual network for steganalysis of digital images," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 131–146, 2021.

[43] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM Review*, vol. 51, no. 3, p. 455–500, 2009.

[44] N. Cohen *et al.*, "Analysis and design of convolutional networks via hierarchical tensor decompositions," *arXiv preprint arXiv:1705.02302*, 2017.

[45] Y. Kim *et al.*, "Compression of deep convolutional neural networks for fast and low power mobile applications," in *Proc. 2016 International Conference on Learning Representations, (ICLR'2016)*, 2016.

[46] I. Glasser, R. Sweke, N. Pancotti, J. Eisert, and J. I. Cirac, "Expressive power of tensor-network factorizations for probabilistic modeling," in *Proc. 2019 Advances in neural information processing systems, (NIPS'2019)*, 2019, pp. 1498–1510.

[47] P. Bas, T. Filler, and T. Pevný, ""Break Our Steganographic System": The ins and outs of organizing boss," in *Information Hiding*, T. Filler, T. Pevný, S. Craver, and A. Ker, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 59–70.

[48] P. Bas and T. Furon, "BOWS-2," http://bows2.ec-lille.fr, accessed: 2019-4-6.

[49] R. Cogranne, Q. Giboulot, and P. Bas, "ALASKA#2," https://alaska.utt.fr, accessed: 2020-10-06.

[50] V. Sedighi, J. Fridrich, and R. Cogranne, "Toss that BOSSbase, Alice!" in *Proc. Media Watermarking, Security, and Forensics, Part of IS&T International Symposium on Electronic Imaging (EI'2016)*, San Francisco, CA, USA, 14-18 February 2016.

[51] M. Abadi *et al.*, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, software available from tensorflow.org. [Online]. Available: https://www.tensorflow.org/

[52] Y. Yousfi, J. Butora, E. Khvedchenya, and J. Fridrich, "ImageNet pretrained CNNs for JPEG steganalysis," in *Proc. 2020 IEEE International Workshop on Information Forensics and Security (WIFS' 2020)*, 2020.

[53] J. Butora, Y. Yousfi, and J. Fridrich, "How to pretrain for steganalysis," in *Proc. 2021 ACM Workshop on Information Hiding and Multimedia Security (IH&MMSec' 2021)*, 2021, pp. 143–148.

[54] Y. Yousfi, J. Butora, J. Fridrich, and C. Fuji Tsang, "Improving EfficientNet for JPEG steganalysis," in *Proc. 2021 ACM Workshop on Information Hiding and Multimedia Security (IH&MMSec' 2021)*, 2021, pp. 149–157.
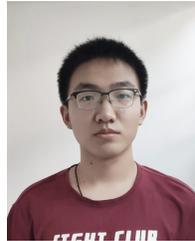
**Laiyuan Li** is currently a master student in Shenzhen University majoring in computer technology. His current research interests include multimedia forensics, Image tampering and deep learning.



**Bin Li (S'07-M'09-SM'17)** received the B.E. degree in communication engineering and the Ph.D. degree in communication and information system from Sun Yat-sen University, Guangzhou, China, in 2004 and 2009, respectively.

He was a Visiting Scholar with the New Jersey Institute of Technology, Newark, NJ, USA, from 2007 to 2008. He is currently a Professor with Shenzhen University, Shenzhen, China, where he joined in 2009. He is also the Director with the Guangdong Key Lab of Intelligent Information Processing and the Director with the Shenzhen Key Laboratory of Media Security. He is an Associate Editor of the IEEE Transactions on Information Forensics and Security. His current research interests include multimedia forensics, image processing, and deep machine learning.



**Jiwu Huang (M'98–SM'00–F'16)** received the B.S. degree from Xidian University, Xi'an, China, in 1982, the M.S. degree from Tsinghua University, Beijing, China, in 1987, and the Ph.D. degree from the Institute of Automation, Chinese Academy of Science, Beijing, in 1998. He is currently a Professor with the College of Electronics and Information Engineering, Shenzhen University, Shenzhen, China. Before joining Shenzhen University, he has been with the School of Information Science and Technology, Sun Yat-sen University, Guangzhou, China, since 2000. His current research interests include multimedia forensics and security. He is an Associate Editor of the IEEE Transactions on Information Forensics and Security.



**Shunquan Tan (M'10–SM'17)** received the B.S. degree in computational mathematics and applied software and the Ph.D. degree in computer software and theory from Sun Yat-sen University, Guangzhou, China, in 2002 and 2007, respectively.

He was a Visiting Scholar with New Jersey Institute of Technology, Newark, NJ, USA, from 2005 to 2006. He is currently an Associate Professor with College of Computer Science and Software Engineering, Shenzhen University, China, which he joined in 2007. He is the Vice Director with the Shenzhen Key Laboratory of Media Security. His current research interests include multimedia security, multimedia forensics, and machine learning.



**Qiushi Li** received the B.S. degree in information and computing science from Harbin University of Science and Technology, Harbin, China, in 2018. He is currently pursuing the Ph.D. degree in information and communication engineering with Shenzhen University, Shenzhen, China. His current research interests include multimedia forensics and machine learning.