# Model design and parameter optimization of CNN for side-channel cryptanalysis

**Yun Lin Liu** [Corresp., 1] , **Yan Kai Chen** [1] , **Wei Xiong Li** [1] , **Yang Zhang** [1]

[1] Center of Equipment Simulation Training, Shijiazhuang campus of the Army Engineering University, Shijiazhuang, Hebei, China

Corresponding Author: Yun Lin Liu
Email address: llyun324@163.com

**Background.** The side-channel cryptanalysis method based on convolutional neural network (CNNSCA) can effectively carry out cryptographic attacks. The CNNSCA network models that achieve cryptanalysis mainly include CNNSCA based on the VGG variant (VGG-CNNSCA) and CNNSCA based on the Alexnet variant (Alex-CNNSCA). The learning ability and cryptanalysis performance of these CNNSCA models are not optimal, and the trained model has low accuracy, too long training time, and takes up more computing resources. In order to improve the overall performance of CNNSCA, the paper will improve CNNSCA model design and hyperparameter optimization.

**Methods.** The paper first studied the CNN architecture composition in the SCA application scenario, and derives the calculation process of the CNN core algorithm for side-channel leakage of one-dimensional data. Secondly, a new basic model of CNNSCA was designed by comprehensively using the advantages of VGG-CNNSCA model classification and fitting efficiency and Alex-CNNSCA model occupying less computing resources, in order to better reduce the gradient dispersion problem of error back propagation in deep networks , the SE (Squeeze-and-Excitation) module is newly embedded in this basic model , this module is used for the first time in the CNNSCA model, which forms a new idea for the design of the CNNSCA model.Then apply this basic model to a known first-order masked dataset from the side-channel leak public database (ASCAD). In this application scenario, according to the model design rules and actual experimental results, exclude non-essential experimental parameters. Optimize the various hyperparameters of the basic model in the most objective experimental parameter interval to improve its cryptanalysis performance, which results in a hyper-parameter optimization scheme and a final benchmark for the determination of hyper-parameters.

**Results.** Finally, a new CNNSCA model optimized architecture for attacking unprotected encryption devices is obtained——CNNSCAnew. Through comparative experiments, CNNSCAnew's guessing entropy evaluation results converged to 61. From model training to successful recovery of the key, the total time spent was shortened to about 30 minutes, and we obtained better performance than other CNNSCA models.

1   # Model design and parameter optimization of CNN for
2   # side-channel cryptanalysis

3

4   Yun Lin Liu, Yan Kai Chen, Wei Xiong Li, Yang Zhang
5   email address of the corresponding author(llyun324@163.com)
6   Center of Equipment Simulation Training,Shijiazhuang campus of the Army Engineering
7   University, Shijiazhuang, 050003, China

8

9   ## Abstract

10  **Background.** The side-channel cryptanalysis method based on convolutional neural network
11  (CNNSCA) can effectively carry out cryptographic attacks. The CNNSCA network models that
12  achieve cryptanalysis mainly include CNNSCA based on the VGG variant (VGG-CNNSCA) and
13  CNNSCA based on the Alexnet variant (Alex-CNNSCA). The learning ability and cryptanalysis
14  performance of these CNNSCA models are not optimal, and the trained model has low accuracy,
15  too long training time, and takes up more computing resources. In order to improve the overall
16  performance of CNNSCA, the paper will improve CNNSCA model design and hyperparameter
17  optimization.
18  **Methods.** The paper first studied the CNN architecture composition in the SCA application
19  scenario, and derives the calculation process of the CNN core algorithm for side-channel leakage
20  of one-dimensional data. Secondly, a new basic model of CNNSCA was designed by
21  comprehensively using the advantages of VGG-CNNSCA model classification and fitting
22  efficiency and Alex-CNNSCA model occupying less computing resources, in order to better
23  reduce the gradient dispersion problem of error back propagation in deep networks , the SE
24  (Squeeze-and-Excitation) module is newly embedded in this basic model , this module is used
25  for the first time in the CNNSCA model, which forms a new idea for the design of the CNNSCA
26  model.Then apply this basic model to a known first-order masked dataset from the side-channel
27  leak public database (ASCAD). In this application scenario, according to the model design rules
28  and actual experimental results, exclude non-essential experimental parameters. Optimize the
29  various hyperparameters of the basic model in the most objective experimental parameter
30  interval to improve its cryptanalysis performance, which results in a hyper-parameter
31  optimization scheme and a final benchmark for the determination of hyper-parameters.
32  **Results.** Finally, a new CNNSCA model optimized architecture for attacking unprotected
33  encryption devices is obtained——CNNSCAnew. Through comparative experiments,
34  CNNSCAnew's guessing entropy evaluation results converged to 61. From model training to
35  successful recovery of the key, the total time spent was shortened to about 30 minutes, and  we
36  obtained better performance than other CNNSCA models.

37

38  **Key words** Side-channel analysis,CNN,VGG, Alexnet,SEnet, Hyperparameter

39

## Introduction

40 **Introduction**

41 Side Channel Analysis[1](SCA) refers to bypassing the tedious analysis of encryption algorithms,
42 by using the information  (such as execution time, power consumption, electromagnetic radiation
43 , etc.)leaked by the hardware device embedded in the encryption algorithm during the calculation
44 process , combined with statistical analysis methods to attack cryptographic systems. The side-
45 channel cryptanalysis method is divided into profiling methods and non-profiling methods: non-
46 profiling methods include differential power attack[2] (DPA), correlation power attack[3] (CPA)
47 and mutual information attack[4] (MIA); profiling methods include template attack[5] (TA), side-
48 channel cryptography attack based on multi-layer perceptron (MLPSCA), and side-channel
49 cryptography attack based on convolutional neural networks (CNNSCA). Although the attack
50 method of the non-profiling method is simple and direct, weak side-channel signal or excessive
51 environmental noise can cause the attack to fail. The profiling method can effectively analyze the
52 characteristics of the side-channel signal when the encryption knowledge of the attacking device
53 is obtained in advance, so it is easier to crack the cryptogramme. In the case of an encrypted
54 implementation copy, the best cryptanalysis attack in the traditional SCA method is TA[5-8], but
55 TA has difficulties in statistical analysis when processing high-dimensional side-channel signals,
56 and cannot attack the implementation of protected encryption. With the rapid development of
57 supervised machine learning algorithms, it can effectively analyze one-dimensional data with
58 similar power consumption traces in other fields, and side-channel cryptanalysis based on
59 machine learning (MLSCA)[9-11] has begun to emerge. The new profiling method MLPSCA
60 surpasses the traditional profiling method in attack performance [11-13], and overcomes the
61 shortcomings of template attacks that cannot handle high-dimensional side-channel signals, but it
62 also loses effectiveness when attacking encryption with protection. Nowadays, with the
63 development of machine learning, deep learning techniques with excellent performance in image
64 classification and target recognition have become popular. Studies have shown that the
65 application of convolutional neural network algorithms under deep learning can produce better
66 encryption performance in side-channel analysis[12-16]. The deep network helps to mine the deep
67 features in the data, which can make the neural network have more powerful performance, which
68 makes CNNSCA can also attack the encryption implementation with protection. In the side-
69 channel analysis application scenario, deep learning eliminates the step of manually extracting
70 features from the workflow of model construction. For example, in the traditional bypass attack
71 method, the TA with better attack effect only selects 5 strong feature points, while the deep
72 learning model can select hundreds to thousands of feature points, select more features to
73 construct a template, it is extremely beneficial to the generalization and robustness of the side-
74 channel analysis model.

75

76 Analyze the above domestic and foreign documents, there are two main types of CNN structures
77 that have successfully used CNNSCA to achieve cryptanalysis, which are based on two variants
78 of Alexnet and VGGnet network structures [12,16-18]. Among them, the 2012 ILSVRC(ImageNet
79 Large Scale Visual Recognition Challenge) champion structure Alexnet[19], although successful

80    in the SCA application, but in fact, the training accuracy of CNNSCA based on this network
81    variant is not high, moreover, the Alex-CNNSCA network model in the literature [16] has a large
82    amount of training parameters and a long calculation time, which means that there is still room
83    for optimization of this network structure. The 2013 ILSVRC champion network ZFNet [20] has
84    not changed much from the 2012 first ILSVRC champion network Alexnet. The 2014 ILSVRC
85    runner-up structure VGGnet[21] also succeeded in breaking secrets in the SCA application. In the
86    literature [12,17-18], VGG-CNNSCA models with different parameters were proposed. Among
87    them, the best cryptanalysis performance is in the literature [12] proposed VGG-CNNSCA, but
88    its training accuracy is still not high. Obviously, there is still room for improvement in the
89    cryptanalysis performance. The 2014 ILSVRC champion network GoogLeNet[22] and the 2015
90    ILSVRC champion network ResNet[23] have also been used in SCA, but the effect is average.
91    This conclusion has been confirmed in the literature [12]. The last ILSVRC champion network in
92    2017 was the SEnet[24] proposed by Momenta and Oxford University. There is currently little
93    literature on applying this network to SCA scenarios.
94
95    Although CNNSCA overcomes the shortcomings of the previous profiling methods and
96    improves the cryptanalysis performance, the existing CNNSCA model learning ability and
97    cryptanalysis performance are not optimal. The disadvantages of these models are: low training
98    accuracy and excessive training time long, taking up too much computing resources, etc. The
99    reason is mainly affected by CNNSCA model design and hyperparameter optimization. In order
100   to improve the overall performance of CNNSCA, the paper will improve CNNSCA model
101   design and hyperparameter optimization, and has done the following work:
102       1. The composition of the CNN architecture in the SCA application scenario is studied, and
103          the calculation process of the CNN core algorithm for side-channel leakage of one-
104          dimensional data is deduced.
105       2. Taking advantage of the high efficiency of classification and fitting of the VGG-
106          CNNSCA model and the advantages of the Alex-CNNSCA model occupying less
107          computing resources, a new basic model of CNNSCA is designed to better reduce the
108          gradient dispersion of error back propagation in the deep network. The problem is that
109          the SE module is newly embedded in this basic model, so that the model basically
110          achieves the purpose of breaking the secrets, thereby solving the problem of constructing
111          the CNNSCA model.
112       3. Apply the above basic model to a known first-order mask data set of the side-channel
113          leak public database (ASCAD). In this application scenario, according to the model
114          design rules and actual experimental results, unnecessary experiments are maximized
115          parameter, optimize the various hyperparameters of the model in the most objective
116          experimental parameter interval to improve the breaking performance of the new
117          CNNSCA, which solves the problem of hyperparameter optimization, and gives the final
118          determination benchmark for hyperparameters. Finally, a new CNNSCA model

119        optimized architecture for attacking unprotected encryption devices-CNNSCAnew is

120        obtained.

121    4.  The performance verified on public data sets exceeds other profiling SCA methods.

122

123   The algorithms involved in the paper experiments are all programmed in the Python language,

124   and use the deep learning architecture Keras library[25] (version 2.4.3) or directly use the GPU

125   version of the Tensorflow library[26] (version 2.2.0). The experiment was carried out on an

126   ordinary computer equipped with 16 GB RAM and 8G GPU (Nvidia GF RTX 2060). All

127   experiments use side-channel leaking public data sets-known first-order mask data sets in the

128   ASCAD database, use 50,000 pieces of data from its training set to train the model, and

129   randomly select 1,000 pieces of data from its test set for testing. When testing the cryptanalysis

130   performance of the CNNSCA model, the guessing entropy index is used to evaluate the

131   cryptanalysis performance.

132

133   **Materials & Methods**

134   **Materials**

135   1 CNN

136   Convolutional Neural Network (CNN) is one of the most successful algorithms of artificial

137   intelligence, and it is a multi-layer neural network with a new structure. Its design is inspired by

138   the research on the optic nerve receptive field [27-28]. The core component of CNN, the

139   convolution kernel, is the structural embodiment of the local receptive field. It belongs to the

140   deep network of back propagation training. It uses the two-dimensional spatial relationship of the

141   data to reduce the number of parameters that need to be learned, and improves the training

142   performance of the BP algorithm(Error Back Propagation, which is used to calculate the gradient

143   of the loss function with respect to the parameters of the neural network) to a certain extent. The

144   main difference between CNN and MLP is the addition of the convolution block structure. In the

145   convolution block, a small part of the input data is used as the original input of the network

146   structure, and the data information is forwarded layer by layer in the network, and each layer

147   uses several convolution cores to extract features of the input data. Convolutional neural

148   networks have been successfully applied in computer vision, natural language processing,

149   disaster climate prediction and other fields, especially shine on ILSVRC [29]. ILSVRC is one of

150   the most popular and authoritative academic competitions in the field of machine vision in recent

151   years, representing the highest level in the field of imaging. The introduction of outstanding

152   CNNs in the image classification and target positioning projects of the ILSVRC competition

153   over the years is shown in Table 1(CNN with outstanding performance in previous ILSVRC

154   competitions).

155

156   Table 1 sorts out the champion networks and individual runner-up networks of the last ILSVRC

157   classification task from 2012 to 2017, and briefly introduces their names, rankings, classification

158   results under the top1 and top5 indicators, and some remarks. Top1 refers to the largest

159  probability vector as the prediction result, if the classification is correct, it is correct. Top5 is
160  correct as long as there is a correct classification in the top five of the largest probability vectors.
161  Among them, the error rate of the classification results of the last champion network SEnet (2017)
162  under the top5 index is obviously the lowest, reaching 2.25%. Deep convolutional networks have
163  greatly promoted the development of various fields of deep learning.
164
165  2 CNNSCA model hyperparameters
166  Hyperparameters of neural network models are a concept often used in machine learning or deep
167  learning, including the structural parameters and training parameters of the network model. To
168  design a CNN model in SCA application scenarios, all the parameters that need to be set are as
169  follows:
170      1) Structural parameters
171  Define all the parameters of the neural network architecture, including the regular parameter
172  network layer activation function, classification function, loss function, and optimizer. In the
173  convolutional neural network, the network layer is subdivided into convolutional blocks (a
174  combination of different numbers of convolutional layers and pooling layers), convolutional
175  layers, full link layers, pooling layers, the number of convolution kernels , convolution kernel
176  size and fill.
177
178  The convolution block, convolution layer, pooling layer, number of convolution kernels,
179  convolution kernel size and padding in these parameters mainly control the scale and
180  performance of feature extraction in the feature extraction stage of the CNNSCA model. Full
181  link layer, activation function, classification function and loss function, these parameters
182  constitute the main body of the CNNSCA network, and perform feature learning and fitting
183  classification on side-channel leakage data.
184
185      2) Training parameters
186  Control the parameters of the network model training phase, including the number of iterations,
187  batch learning volume, and learning rate. When training a network model, a complete training set
188  is processed at one time, which is called complete batch learning. If a single training sample is
189  processed at a time, it is called random learning. In practice, in order to improve efficiency, a
190  compromise method is usually adopted, called small-batch learning, that is, small batches of
191  training samples are processed at one time during the model learning process. The batch size
192  depends on environmental factors[30] (such as network architecture, computer GPU performance,
193  the trade-off between network regularization effect and stability, etc.). The number of iterations
194  is an important parameter to be adjusted. A small value will cause the network model to underfit
195  (the model is too poor to capture the feature trend in the training data set), while a higher value
196  will cause the network model to overfit (the model is too Complex, perfectly fits the training data
197  set, but cannot generalize its prediction to other data sets). In addition, the variable that optimizes

198    the training effect of the network model-the learning rate (also called the step size), aims to
199    promote the gradient (ie, the error gradient) drop during the training process.
200
201    The number of iterations and the amount of batch learning affect the degree of model training,
202    and the optimizer and learning rate are used to control the gradient of the error. These parameters
203    all have an important impact on CNNSCA's cryptanalysis performance and need to be adjusted
204    according to specific attack scenarios.
205
206    3 Core algorithm and network structure of CNNSCA
207    3.1 CNN network structure for SCA
208    Combined with the side-channel cryptanalysis scenario, the CNN applied to the side-channel
209    attack mainly has six network layers stacked layer by layer and an embeddable SE module:
210    a)    Convolutional layers (Conv for short) are linear layers. The incomplete connection
211          between layers can avoid the two shortcomings of a fully connected network: training
212          weights requires a huge amount of calculation and model overfitting. The weights of
213          the same convolution kernel (also known as filters) in the same layer are shared,
214          allowing the convolution layer to extract constant displacement features while
215          reducing parameters. The convolutional layer can also use multiple convolution
216          kernels. Each convolution kernel extracts different abstract features from the input
217          vector. These abstract features are arranged side by side in an additional dimension
218          (the so-called depth), making the CNN resistant to time-domain distortion Vector
219          features[31]. The convolutional layer usually needs to set the padding mode , one is
220          valid padding , so that the dimension of the feature vector after convolution is smaller
221          than the original vector; the other is the same padding , so that the convolutional The
222          feature vector dimension is the same as the original vector.
223    b)    Batch Normalization layers[32] (BN for short), whose role is to reduce the deviation of
224          covariates in the two stages of training and prediction, which is conducive to the use
225          of a higher learning rate for the network model[33].
226    c)    Activation layers (ACT for short) are non-linear layers and consist of a single real
227          function, which acts on each coordinate of the input vector. The ReLU function is
228          currently the first choice in deep learning.
229    d)    Pooling layers (POOL for short) are non-linear layers. Use the pooling window to
230          slide on the input vector to extract salient feature points to reduce the feature
231          dimension. There is no weight in the pooling layer, which will not cause distortion of
232          the input signal.
233    e)    Fully-Connected layers (FC for short), the neurons between the layers are completely
234          connected, and these layers need to train a lot of weights. This layer is expressed by
235          an affine function as: D-dimensional x vector is the input, and Ax+B is the output.
236          Among them, $A \in RC \times D$ is the weight matrix and $B \in RC$ is the deviation vector.
237          These weights and deviations are the training parameters of the FC layer.

238    f)    Softmax layer (SOFT for short). In multi-classification tasks, softmax is usually used
239        as the activation function of the output layer. Here, softmax is used to represent the
240        output layer. This layer classifies the input, obtains the predicted value of each label,
241        and takes the label corresponding to the maximum value as the global classification
242        result.
243    g)    SE module, SEnet is a classic attention model structure, and it is also a required basic
244        network structure for fine-grained classification tasks. SEnet proposed the Squeeze-
245        and-Excitation (SE) module, which did not introduce a new spatial dimension, and
246        improved the representation ability of the model by displaying the channel correlation
247        between the features of the convolutional layer. The feature recalibration mechanism:
248        by using global information to selectively enhance informatized features and
249        compress those useless features at the same time. In deep network training, this
250        mechanism can effectively overcome the gradient dispersion problem in error back
251        propagation. The SE module is universal. Even if it is embedded in an existing
252        model, its parameters do not increase significantly. It is a relatively successful
253        attention module[24]. The structure of the SE module is shown in Figure 1 (SE
254        module).
255

256        In Figure 1, the SE module uses global pooling as a squeeze operation, and then uses
257        two FC layers to form an excitation structure to profile the correlation between
258        channels, and output and input the same number of feature channels weights. The
259        advantages of this are: 1) it has more nonlinearity and can better fit the complex
260        correlation between channels; 2) the amount of parameters and the amount of
261        calculation are greatly reduced. Then obtain the normalized weight between 0 and 1
262        through a sigmoid function, and then use a scale operation to weight the normalized
263        weight to the features of each channel[24]. Finally, the output of scale is superimposed
264        on the input $x$ before the SE module to generate a new vector $\tilde{x}$ .
265

266    3.2 Core algorithm of CNN for SCA
267    1)  Convolution calculation
268    Usually convolution operations in the field of computer vision are numerical operations on two-
269    dimensional image data. In the SCA application scenario, the dimensionality of the convolution
270    operation is adjusted, which is to slide the convolution kernel on the one-dimensional energy
271    trace data. The number of steps moved each time is called the step length, and the convolution
272    calculation is performed on each sliding to obtain a value. After one round of calculation is
273    completed, a feature vector representing the vector feature is obtained. The rule of numerical
274    operation is to multiply a one-dimensional convolution kernel with a value at the corresponding
275    position of a one-dimensional vector, and then sum. For example, there is a 1x3 convolution
276    kernel, which convolves a 1x6 one-dimensional vector with a step size of 1. The calculation
277    process is shown in Figure 2 (Convolution calculation process).

278

279    In Figure 2(a), the convolution kernel slides from the left side of the input vector. The first
280    numerical calculation is: 1x1+0x0+1x1=2, and the first value 2 of the new feature vector is
281    obtained. Then, the convolution kernel slides one step to the right to continue the numerical
282    calculation: 1x0+0x1+1x0=0, to get the second value 0 of the new feature vector, as shown in
283    Figure 2(b). Repeat this process until the convolution kernel slides to the far right of the input
284    vector, and the convolution calculation is complete.

285

286    2)  Pooling calculation
287    There are three ways of pooling: Max-Pooling, Mean-Pooling and Stochastic Pooling. Maximum
288    pooling is to extract the maximum value of the value in the pooling window, average pooling is
289    to extract the average value of the value in the pooling window, and random pooling is to
290    randomly extract the value in the pooling window. The original pooling operation of CNN is also
291    a numerical operation on two-dimensional image data. In the SCA application scenario, the
292    pooling calculation has also been dimensionally adjusted, and a pooling mode is selected for
293    calculation on the one-dimensional energy trace data. For example, the pooling window size is
294    1x2, and the maximum or average pooling operation is performed on a 1x6 one-dimensional
295    vector with a step size of 2. The pooling calculation is shown in Figure 3 (Pooling calculation
296    process).

297

298    In Figure 3(a), the maximum pooling starts from the left side of the input vector. Every two steps
299    of the pooling window, the maximum value of the two values in the window is selected as a
300    value of the new feature vector. The average pooling is shown in Figure 3(b). For every two
301    sliding steps of the pooling window, the average of the two values of the window class is
302    calculated as a value of the new feature vector. The pooling window slides to the right until the
303    rightmost of the input vector, and the pooling calculation is complete.

304

305    3)  softmax function
306    This function normalizes the output value and converts all output values into probabilities. The
307    sum of the probabilities is 1. The formula of softmax is:

308    $$softmax\left(x_i\right) = \frac{exp\left(x_i\right)}{\Sigma_j exp\left(x_j\right)} \qquad (1)$$

309    Here $x_i$ represents the input of the i-th neuron in the softmax layer, $x_j$ represents the input of
310    the j-th neurons in the softmax layer, and $\Sigma_j$ is the sum of calculations for $x_j$. The result of the
311    function is used as the fitting probability of the i-th neuron label.

312

313    4)  Principle of weight adjustment
314    Using the cost function and gradient descent algorithm[34], each time the network model is
315    trained, the weights are automatically adjusted in the direction of error reduction, so that the

316    training parameters are repeated until all iterations are over, and the weight adjustment is
317    completed.
318
319    5)  Evaluation of Cryptanalysis Performance
320    Generally, security officers consider two indicators when evaluating CNNSCA's cryptanalysis
321    performance: one is the training accuracy of the neural network model during modeling, the Acc
322    indicator[35], and the other is the security indicator guessing entropy of the key obtained in the
323    attack phase[36-37]. The guessing entropy index is commonly used to evaluate the SCA
324    cryptanalysis performance, and the guessing entropy is used to measure the efficiency of
325    decrypt.Guessing Entropy (GE) is obtained through a custom rank function $Rank(\cdot)$, which is
326    defined as:
327    $$Rank(\hat{g}, D_{train}, D_{test}, n) = \left| \left\{ k \in K \middle| d_n[k] \geq d_n[k^*] \right\} \right| \quad (2)$$
328
329    The adversary uses the modeling data set $D_{train}$ to establish a bypass analysis model $g$, and uses $n$
330    energy trace samples in the attack data set $D_{test}$ to perform $n$ attacks during the attack phase.
331    After each attack, the logarithm value of the distribution probability of 256 types of hypothetical
332    cryptograms is obtained, compose a vector $d_i = [d_i[1], d_i[2], L, d_i[k]]$, whose indexes are arranged in
333    the positive order of the hypothetical cryptogramme's key space (the index counts from zero),
334    where $i \in n$, $k \in K$, and $K$ is the key space of the hypothetical cryptogramme. The results of each
335    attack are accumulated. Then, the rank function $Rank(\cdot)$ sorts all the elements of the vector $d_i$ in
336    reverse order by value, and keeps the position of the corresponding index of each element in the
337    vector before and after sorting consistent with the position of the element, and obtains a new
338    ranking vector $D_i = [D_i[1], D_i[2], L, D_i[k]]$, where each the element $D_i[k]$ contains two values $k$ and
339    $d[k]$, and finally the index of the logarithmic element of the known cryptogramme $k^*$ probability
340    in $D_i$ is output, that is, the guessing entropy $GE(d[k^*])$. At the i-th attack, the higher the matching
341    rate of the energy trace model of the real cryptogramme, the higher the index ranking of its
342    $GE(d[k^*])$. Guessing entropy is the $GE(d[k^*])$ index ranking output of each attack——$rank$. In $n$
343    attacks, the better the performance of the cryptanalysis method and the higher the efficiency, the
344    faster the ranking of $GE(d[k^*])$ converge to zero. It shows that in the i-th attack, the guessing
345    entropy converges to zero and continues to converge in subsequent attacks. The adversary only
346    needs $i$ attacks to crack the cryptogramme, that is, only $i$ power consumption traces are needed to
347    break the secret. Equation (2) can be rewritten as (3):
348    $$GE_n(\hat{g}) = Every\left[ Rank(\hat{g}, D_{train}, D_{test}, n) \right] \quad (3)$$
349
350    4 Side-channel leaking public data sets
351    The newly published ASCAD database[12] aims to achieve AES-128 with first-order mask
352    protection, namely 8-bit AVR microcontroller (ATmega8515), in which the energy trace is the
353    data signal converted by the collected electromagnetic radiation. The adversary outputs the
354    collected signal for the third S-box of the first round of AES encryption, and launches an attack

355  against the first AES key byte. The database follows the MNIST database[38] rules and provides a
356  total of four data sets, each with 60,000 entries power consumption traces, of which 50,000
357  power consumption traces are used for analysis/training, and 10,000 power consumption traces
358  are used for testing/attack. The first three ASCAD data sets respectively represent the encryption
359  realization leakage with three different random delay protection countermeasures. The signal
360  offsets desync=0, desync=50, and desync=100 are used to represent these three data sets with
361  two strategies of mask and delay. All power consumption traces in the first three types of data
362  sets contain 700 feature points. These feature points are selected from the original energy trace
363  containing 100,000 feature points, and the selection basis is the position of the largest signal
364  peak. When the mask is known, the maximum signal-to-noise ratio of the data set can reach 0.8,
365  but it is almost 0 when the mask is unknown. The last ASCAD data set stores the original energy
366  trace.
367

## Methods

369  1 Design of CNNSCA base model
370  With reference to the advantages of the VGG-CNNSCA model with high classification and
371  fitting efficiency and the Alex-CNNSCA model occupying less computing resources, the paper
372  selects the same structural parameters from these two models, and some of the factors that
373  promote the high fitting efficiency of the two types of models parameter. These parameters
374  construct a new CNN simple model specifically for SCA scenarios, which is used to test the
375  impact of different hyperparameters on model performance. The convolution block of this simple
376  model consists of the Conv layer, BN layer, and ACT layer. After the block, a POOL layer is
377  usually added to reduce the feature dimension. The new convolution block is repeated n times in
378  the network model until it is reasonable. Until the output of the size. Then, introduce n FC layers,
379  use the softmax function in the last FC layer, and finally output the classification prediction
380  results. In addition, in order to improve the classification and recognition performance of the
381  CNNSCA model, the SE module is newly embedded in the simple model. Its main function is to
382  reduce the gradient dispersion problem in error back propagation. This is the first use in the
383  CNNSCA model. The SE module will be embedded between the convolutional layer and the
384  pooling layer of the convolutional block of the simple model, and the simple model containing
385  the SE module will be renamed to the CNNSCA base model-CNNSCAbase. The newly designed
386  CNNSCAbase structure is shown in Figure 4(Convolutional network structure in a side-channel
387  attack scenario).
388

389  The initial configuration basis and selection of CNNSCAbase are as follows: Find out the two
390  prototypes of Alex-CNNSCA and VGG-CNNSCA to set the same parameters, set these
391  parameters in CNNSCAbase in the same way, these parameters are as follows: 5 convolutional
392  blocks , 3 full connections , the padding modes of the convolutional layers are SAME, and the
393  activation functions of all layers before the last layer of the network select ReLU. In addition, in
394  most classification tasks, convolutional networks use softmax, crossentropy, and RMSprop as the

395  model's classification function, loss function, and optimizer[19-23]. Here, CNNSCAbase also
396  chooses to use these three activation functions . Since the side-channel leakage data belongs to
397  one-dimensional data, the processing complexity is less than that of two-dimensional data. Here,
398  the convolution layer of each convolution block is initialized to 1, and the number of convolution
399  kernels in the first convolution layer is 64 (choose the smaller number of Alexnet or VGGnet),
400  the size of the convolution kernel is 3x3, the step size is 1, the pooling mode is tentatively
401  averaged pooling mode, the pooling window size is 2, and the step size is 2. In addition, in the
402  initial setting of CNNSCAbase, a new SE module is embedded in the last four convolution
403  blocks. All initial structure parameters of CNNSCAbase are shown in Table 2(CNNSCAbase
404  Configuration).
405
406  Here we first verify and analyze the model training effect of CNNSCAbase with and without SE
407  module. Remove the SE module from the CNNSCAbase model, all other parameters remain
408  unchanged, and name this model CNNSCAnoSE. Train the models CNNSCAnoSE and
409  CNNSCAbase on the training set of the ASCAD dataset with known masks. The training results
410  of the two models are shown in Figure 5(Training effect of CNNSCAnoSE model and
411  CNNSCAbase model):
412
413  As shown in Figure 5, when the training iteration reaches 28 times, the accuracy of
414  CNNSCAbase is significantly higher than that of CNNSCAnoSE, which is about 96%. Continue
415  to train the CNNSCAnoSE model, and when the training iteration reaches 70 times, its accuracy
416  rate rises to about 90%. In addition, when the training accuracy of the two models is close, the
417  training time of the 28-iteration CNNSCAbase model is about 1393 seconds, which is
418  significantly less than the training time of the 70-iteration CNNSCAnoSE model, the training
419  time of the latter is about 2240 seconds. This proves that the SE module can promote the
420  improvement of the classification performance of the CNNSCAbase model and can reduce the
421  model training time.
422
423  Next, we discuss the hyperparameter optimization of the CNNSCA model. Model structure
424  parameters and training parameters are hyperparameters and need to be set in advance. Later, we
425  will design a set of experimental procedures to optimize these hyperparameters in specific
426  application scenarios. For example, we choose to determine the model parameters first rather
427  than the global training parameters, first determine the number of Conv layers, rather than first
428  determine the kernel size or the number of filters . The reason for this design is: currently,
429  Python's deep learning architecture library [25-26] is mainly used to program the CNN network.
430  When using these library methods, the CNN network structure is usually programmed first. The
431  order in which these parameters appear in the program code will be affected by the library
432  methods, and then the training parameters are designed according to the size of the network and
433  the size of the training set. It is precisely in consideration of the order in which the parameters

434 appear during programming, we have designed the order of the following experimental
435 procedures.
436
437 2 Selection and optimization of CNN structure parameters for side-channel cryptanalysis
438 2.1 Structural parameter selection rules
439 In section Methods 1, the base model CNNSCAbase is set, and the best model after parameter
440 optimization will be named CNNSCAnew later. In CNNSCAbase, in addition to the specific set
441 of structural parameters, the remaining structural parameters need to be customized. These
442 structural parameters include classification function, loss function, optimizer, the number of
443 convolutional layers in each convolution block, the number of convolution kernels in the
444 convolution layer, convolution kernel size, pooling layer pooling mode. When choosing these
445 custom structure parameters, you need to follow the classic rules of building a deep learning
446 network structure[20,22], which can reduce the number of unnecessary test parameters. The rules
447 are as follows:
448
449 Rule 1: Set the same parameters for the convolutional layers in the same convolutional block to
450 keep the amount of data generated by different layers unchanged.
451 Rule 2: The dimensionality of each pooling window is 2, and the window sliding step is also 2,
452 each operation reduces the dimensionality of the input data to half.
453 Rule 3: In the convolutional layer of the i-th block (starting from i=1), the number of convolution
454 kernels is n: $n_i = n_1 \times 2^{i-1}$, i≥2. This rule keeps the amount of data processed by different
455 convolution blocks as constant as possible. The network structure characteristics of VGG-16 in
456 this reference [21] are formulated.
457 Rule 4: The size of the convolution kernel of all convolution layers is the same.
458
459 2.2 Structural parameter optimization
460 Among the custom structure parameters, the structure parameters that need to be further adjusted
461 through experimental analysis are: the number of convolution layers in each convolution block,
462 the number of convolution kernels in the convolution layer, the size of the convolution kernel,
463 the pooling mode of the pooling layer, SE module. The experimental process of structural
464 parameter optimization is as follows:
465
466 1) Number of convolutional layers
467 In Section Methods 1, in the initial CNNSCAbase structure, the number of convolutional layers
468 for each convolution block is 1, and the convolutional structure is named Cnov1. Refer to the
469 number of convolutional layers of different convolutional blocks of the Alexnet and VGGnet16
470 prototypes. It is found that the minimum number is 1 and the maximum is 3, and the small
471 number is distributed in the front convolution block, and the large number is at the back. This is
472 also to build deep learning The common habit of the Internet. Therefore, the upper limit of the
473 number of convolutional layers of the CNNSCAbase convolution block is set to 3, and the

474　　baseline is Cnov1, and a certain convolutional layer parameter configuration can be obtained
475　　through two sets of necessary experiments. When training the CNNSCAbase model, the training
476　　iteration and batch parameters of the current optimal CNNSCA model[12] are used, which are 75
477　　and 200 (in all experiments in section Methods 2.2, unless otherwise specified, the iteration and
478　　batch parameters are used. experiment).
479
480　　Experiment 1: Set up a model in which the number of convolutional layers in 5 convolutional
481　　blocks is 2, and other parameters are consistent with CNNSCAbase, and the structure is named
482　　Cnov2. Then set the number of convolutional layers of the first 4 convolutional blocks to 2, and
483　　the convolutional layer of the last convolutional block to 3. Other parameters are consistent with
484　　CNNSCAbase, and the structure is named Cnov3. The specific settings of the number of
485　　convolutional layers of each convolution block of Cnov1~3 are shown in Table
486　　3(CNNSCAbase.Conv1-7 Configuration). The three structures constructed are trained and tested,
487　　and the results of experiment 1 are shown in Figure 6(Convergence of guessing entropy of
488　　Cnov1~3).
489
490　　From the results in Figure 6, it is found that when Cnov2 and Cnov1 attack the 750th energy
491　　trace, their guessing entropy basically converges to 0, while Cnov3 cannot converge in a finite
492　　number of (1000) attacks. When doing further analysis, if you set two or more convolutional
493　　blocks with 3 convolutional layers in the 5 convolutional blocks of the model, the calculation
494　　amount and parameter amount of model training will increase by several times, and the 8G GPU
495　　memory used in the experiment will be directly exhausted , unable to run the code, then this
496　　parameter setting method will have no practical significance. Therefore, the upper limit of the
497　　number of convolutional layers for each convolutional block is determined to be 2.
498
499　　Experiment 2: On the basis of the conclusion of Experiment 1, the convolutional layer parameter
500　　setting of each convolution block is further accurate. As shown in Figure 6, the convergence of
501　　the orange line representing the entropy of Cnov2's guess is more stable than that of Cnov1, but
502　　it is obvious that there are more convolutional layers, which means that the amount of model
503　　calculations and parameters are relatively large, which affects the overall performance of the
504　　model. Therefore, four structures of Cnov4~7 are set, and each structure sequentially reduces the
505　　number of convolutional layers in each convolution block of Cnov2 by one. The specific settings
506　　of the number of convolutional layers of each convolution block of Cnov4~7 are shown in Table
507　　3(CNNSCAbase.Conv1-7 Configuration). Train and test these constructed structures, and the
508　　results of experiment 2 are shown in Figure 7(Convergence of guessing entropy of
509　　Cnov1~2,4~7).
510
511　　The red curve representing the entropy of Cnov5 guessing in Figure 7 converges optimally.
512　　Finally, the number of convolutional layers of the Cnov5 structure is determined, and the
513　　benchmark is set for the number of convolutional layers of CNNSCAnew.

514

515   2)  The number of convolution kernels in the convolution layer

516   It is known that the number of convolution kernels of each convolutional layer of CNNSCAbase

517   is initially set according to Rule 3. The number of convolution kernels of the first convolutional

518   layer is 64. Usually increasing the number of convolution kernels means that more dimensional

519   feature extraction is performed on the input data, thereby improving the classification efficiency

520   of the convolutional network. But it will inevitably lead to an increase in the amount of

521   calculation and storage of the attack device, which will lead to an increase in the training time of

522   the model. Therefore, under the condition that the efficiency loss of the guarantee model is not

523   large, the model training time can be reduced by reducing the number of convolution kernels.

524   Since the number of convolution kernels in the later layer increases by a factor of 2 of the

525   number of convolution kernels in the first convolution layer, to determine the number of

526   convolution kernels as a benchmark, it is only necessary to test the number of convolution

527   kernels in the first convolution layer. At the same time, the CNNSCA model in reference [12],

528   the upper limit of the number of convolution kernels in the convolution layer is 512, which can

529   achieve the effect of breaking the density, so the paper also adjusts the upper limit of the number

530   of convolution kernels to 512.

531

532   Experiment 3: Name the four structures tested as filter1, filter2, filter3, and filter4. The

533   convolution kernel values of the first convolution layer are 8, 16, 32, 64, and the number of

534   convolution kernels of the remaining four convolution blocks is also increased by a factor of 2

535   respectively. The upper limit of the number of convolution kernels is always 512. Other

536   structural parameters are the parameters of the current CNNSCAnew. Train and test the filter1~4

537   structure, and the result of experiment 3 is shown in Figure 8(Convergence of guessing entropy

538   of filter1~4 (epochs=75)).

539

540   Figure 8 shows that after 75 iterations of training, the guessing entropy of the filter4 structure

541   cannot converge. Although the guessing entropy of the filter1~3 structure converges, it fluctuates

542   all the time. When checking the training accuracy of the filter1~3 structure, it is found that the

543   accuracy of the three structures has reached more than 99%, or even reached 1. Obviously, the

544   model has an overfitting phenomenon, which is the most common problem in neural networks.

545   Therefore, the number of training iterations of the filter1~3 structure is reduced to 40, the three

546   structures are retrained, and then the test set is attacked again, and the result shown in Figure

547   9(Convergence of guessing entropy of filter1~3 (epochs=40))  is obtained. The guessing entropy

548   of the filter1~3 structure in Figure 9 all converge to rank 0, and filter3 converges to the position

549   of rank 0 earliest. In summary, the benchmark for selecting convolution kernel parameters is the

550   filter3 structure, and the convolution kernel parameters of the CNNSCAnew structure are

551   updated synchronously.

552

553   3)  Pooling mode of pooling layer

554    It is known that the initial setting mode of the pooling layer of CNNSCAbase is AveragePool,
555    and another common pooling mode is MaxPooling. According to rule 3, both the pooling
556    window and the pooling step size are still selected here.
557
558    Experiment 4: Will test the impact of two pooling modes AveragePool and MaxPooling on the
559    current CNNSCAnew structure. The result of experiment 4 is shown in Figure 10(Convergence
560    of guessing entropy of AveragePool and MaxPool structure).
561
562    In Figure 10, it is obvious that the guessing entropy convergence of the average pooling structure
563    is better than the maximum pooling structure, so the benchmark of the pooling layer pooling
564    mode is average pooling, and the pooling mode of the CNNSCAnew structure is set to average
565    pooling.
566
567    4)  Convolution kernel size
568    The size of the convolution kernel of each convolution layer in CNNSCAbase is initially set to
569    1x3, or 3 for short. In deep learning, people often reduce the size of the convolution kernel by
570    increasing the depth of the network, thereby reducing the computational complexity of the
571    network. In the VGG-CNNSCA structure, the convolution kernel uses a larger size of 11, and in
572    the Alex-CNNSCA structure, a small size of 3 is used.
573
574    Experiment 5: Test the attack effects of the models with the convolution kernel sizes of 3, 5, 7, 9,
575    and 11, and name these five structures as kernel3, kernel5, kernel7, kernel9, and kernel11. The
576    other parameters of these structures are compared with The current CNNSCAnew is the same.
577    The result of experiment 5 is shown in Figure 11(Convergence of guessing entropy of different
578    convolution kernel size structures).
579
580    In Figure 11, the convolution kernel size of the structure kernel3 is 3, which guesses that the
581    entropy convergence is better than other structures, so the size 3 is used as the setting reference
582    for the convolution kernel size. At the same time, the size of the convolution kernel of the
583    CNNSCAnew structure is set to 3.
584
585    5)  SE module
586    The attention mechanism in deep learning is essentially similar to the selective visual attention
587    mechanism of humans, and the core role is to select information that is more critical to the
588    current task goal from a large number of information[24]. The paper has initially added an SE
589    fixed module to the last four convolution blocks of CNNSCAbase. The initial setting of the
590    dimensional change ratio of the first full link layer of the SE module is 1/16, but this
591    conventional setting is in the SCA scene The suitability of the medium requires further
592    verification.
593

594 Experiment 6: Test the SE module of the model. The rate of the dimensional change of the first
595 full link layer is 1/4, 1/8, 1/16, and 1/32 respectively. The other parameters of the test model are
596 the same as CNNSCAnew. . Experiment 7: Test the attack effect of the SE module used 1, 2, and
597 3 times for the last four convolutional blocks in the current CNNSCAnew structure. The results
598 of experiment 6 are shown in Figure 12(Convergence of guessing entropy for different SE
599 dimension ratios).
600
601 As shown in Figure 12, when the dimensional ratio of the SE module is 1/8, the guessing entropy
602 convergence of the overall structure of the CNN is the best. On the basis of this dimensional ratio,
603 the result of Experiment 7 is shown in Figure 13(Convergence of guessing entropy for different
604 number of SE cycles). It is found that when the last four convolution blocks of CNNSCAnew use
605 the SE module twice, the guessing entropy converges fastest. Therefore, the dimensional ratio of
606 the SE module is 1/8 and the SE module is looped twice as a new benchmark for the parameters
607 of the SE module in the CNNSCAnew structure.
608
609 6)   Number of channels at the full link layer
610 The CNNSCA model in literature [12,16] uses 4096 channels in the fully connected layer, which
611 is similar to the number of channels in the fully connected layer of the original VGGnet and
612 Alexnet structures. Considering that the classification task of the ImageNet competition is 1000
613 classifications, and only 256 classifications are needed in the SCA scene, the number of channels
614 can be adjusted appropriately to reduce the training complexity of the model.
615
616 Experiment 8: Will test the model attack effect of the four cases where the number of channels in
617 the fully connected layer is 4096, 3072, 2048, and 1024. The other parameters of these test
618 models are the same as CNNSCAnew. The reason why the number of channels is not set lower
619 than 1024 is that from the convolutional layer to the fully connected layer, if the vector
620 dimension changes sharply, the feature points of the vector are greatly reduced, which will affect
621 the training effect of the model. The result of experiment 8 is shown in Figure 14(Convergence
622 of guessing entropy of the four channel number structure of FC layer).
623
624 It is found from Figure 14 that when the number of channels of the FC layer is 1024, the
625 guessing entropy of its structure converges fastest and continues to be stable. Therefore, 1024 is
626 selected as the reference for the number of channels in the FC layer of the CNNSCAnew
627 structure.
628
629 In summary, the parameter benchmark of the CNNSCAnew structure has been optimized. The
630 new structure parameters are shown in Table 4 (CNNSCAnew Configuration).
631
632 3 Selection and optimization of CNN training parameters for side-channel cryptanalysis

633　Almost all experiments in Section Methods 2.2 use the three parameters of 75 iterations, 200
634　batches of learning, and 1x10-4 learning rate for experiments. These training parameters have
635　little effect on the experimental effects of optimizing various structural parameters, but It is not
636　the optimal setting. The convolutional network of deep learning is applied to the side-channel
637　attack, and these training parameters should also be tuned according to the actual processed side-
638　channel signal data. The order of training parameter tuning is usually learning rate, batch
639　learning amount, and number of iterations[39-40]. In the experiment of training parameter
640　optimization, the current CNNSCAnew structure is used. The training parameter optimization
641　experiment process is as follows:

642

643　1) Learning rate
644　The learning rate is a hyperparameter that is artificially set. The learning rate is used to adjust the
645　size of the weight change, thereby adjusting the speed of model training. The learning rate is
646　generally between 0-1. The learning rate is too large, and the learning will be accelerated in the
647　early stage of model training, making it easier for the model to approach the local or global
648　optimal solution, but there will be large fluctuations in the later stage of the training, and even
649　the value of the loss function may hover around the minimum value, which is always difficult to
650　reach Optimal solution; the learning rate is too small, the model weight adjustment is too slow,
651　and the number of iterations is too much.

652

653　Experiment 9: Will test the impact of five commonly used learning rates on the model's
654　cryptanalysis effect, namely $lr1=1\times10^{-2}$, $lr2=1\times10^{-3}$, $lr3=1\times10^{-4}$, $lr4=1\times10^{-5}$, $lr5=1\times10^{-6}$ . The
655　result of experiment 9 is shown in Figure 15(Convergence of model guessing entropy under five
656　learning rates).

657

658　Figure 15 reflects that when the learning rate is lr2, the guessing entropy of CNNSCAnew
659　converges fastest and is the most stable. Therefore, $1\times10^{-3}$ is selected as the learning rate
660　benchmark of the CNNSCAnew structure.

661

662　2) Batch size
663　The appropriate batch size is more important for the optimization of the model. This parameter
664　does not need to be fine-tuned, just take a rough number, usually $2^n$ (GPU can play a better
665　performance for batches of the power of 2). A batch size that is too large will be limited by the
666　GPU memory, the calculation speed will be slow, and it cannot increase indefinitely (the training
667　set has 50,000 data); it cannot be too small, which may cause the algorithm to fail to converge.

668

669　Experiment 10: According to the size of the ASCAD data set in Section Materials 4, this
670　experiment selects the batch size values: 32, 64, 128, and 256 for the experiment. The result of
671　experiment 10 is shown in Figure 16(Convergence of model guessing entropy under four
672　batches).

673

674 It can be seen from Figure 16 that when the batch learning amount is 128, the guessing entropy
675 of CNNSCAnew converges fastest and is the most stable. Therefore, 128 is selected as the batch
676 size benchmark of CNNSCAnew structure.

677

678 3) Number of iterations (epoch)
679 The number of iterations is related to the fitting performance of the CNNSCA model. The model
680 has been fitted (the accuracy rate reaches 1), and there is no need to continue training; on the
681 contrary, if all epochs have been calculated, but the loss value of the model is still declining, and
682 the model is still optimizing, then the epoch is too small. Should increase. At the same time, the
683 number of iterations of model training also refers to the actual cryptanalysis effect of the model,
684 which is measured by guessing entropy.

685

686 Experiment 11: In experiments 1-10, almost all use iteration number 75 to train the CNNSCA
687 model. This experiment will center on iteration number 75, and train the CNNSCAnew model at
688 10 intervals in the upper and lower intervals to further optimize the model parameter epoch. The
689 interval number of 10 is chosen because the step interval is too small, and the error loss of model
690 training is not much different, so the setting is meaningless; the interval is too large, and repeated
691 experiments may be required to determine an appropriate number of iterations. Therefore,
692 Experiment 11 will test 8 iteration parameters epoch1=15, epoch2=25, epoch3=35, epoch4=45,
693 epoch5=55, epoch6=65, epoch7=75, epoch8=85. The current CNNSCAnew structure has
694 achieved higher training accuracy and breaking performance, in order to reduce model
695 calculation pressure and calculation time, lower iteration parameters are usually selected when
696 the model performs better. Therefore, the upper limit of the epoch test parameter is set to 85. The
697 result of experiment 11 is shown in Figure 17(Convergence of model guessing entropy under
698 eight epochs).

699

700 From the results in Figure 17, it is found that the model of epoch1-4 guesses that the entropy
701 does not converge. Separately recalculate the graph of epoch5-8 model. The result is shown in
702 Figure 18(Convergence of model guessing entropy under four epochs). It can be clearly seen that
703 the epoch5-8 model has a convergence trend. Among them, the epoch5 curve is closest to the
704 position of ranking 0, and the epoch8 curve first converges to ranking 0, but afterwards it
705 fluctuates more widely, and it is obviously over-fitting. The convergence of epoch6 and 7 is
706 similar, the curve begins to fluctuate greatly, and it is close to ranking 0 in the later period.

707

708 Experiment 12: Continue to debug the epoch parameters in a smaller range, and test the other
709 two iterations with an interval of only 5: epoch60=60, epoch70=70. The trained epoch60 and
710 epoch70 models and the previously trained epoch5, epoch6, and epoch7 models are
711 simultaneously attacked on the target set. The results of Experiment 12 are shown in Figure
712 19(Convergence of model guessing entropy under five epochs).

713
714   Figure 19 shows that the guessed entropy of the epoch70 model converges best, and its guessed
715   entropy converges fastest and is the most stable. Therefore, 70 is selected as the training iteration
716   benchmark of the CNNSCAnew structure.
717

## Results

719   1 Get a new model CNNSCAnew for attacking ASCAD data set with known first-order mask
720   protection.
721   According to the 12 sets of experiments in Section Methods 2 and Section Methods 3, the best
722   benchmarks for CNNSCAnew structure parameters and training parameters are demonstrated.
723   The CNNSCAnew model contains 5 convolutional blocks, 8 convolutional layers, and 3 fully
724   connected layers. The size of the convolution kernel of each convolution layer is 3, the activation
725   function is ReLU, and the padding is Same. Each convolutional block is equipped with a pooling
726   layer, the pooling layer selects the average pooling mode, and the pooling window is (2, 2). The
727   number of output channels of the convolution layer in the convolution block 1-5 starts from 32
728   and increases by a multiple of 2 in turn. Two SE modules are added after the convolution layer
729   of each convolution block in the convolution block 2-4, and the dimension ratio of the SE is set
730   to 1/8. In the first two fully connected layers, set the number of output channels to 1024 and the
731   activation function to ReLU. The output channel number of the third fully connected layer is the
732   target classification number 256, and the classification function is Soft-max. The global
733   configuration loss function is crossentropy, the optimization method is RMSprop, the number of
734   training iterations is 70, the learning rate is $1\times10^{-3}$, and the batch learning volume is 128. All
735   parameters of the newly obtained CNNSCAnew are shown in Table 5(CNNSCAnew
736   Configuration).
737

738   2 The CNNSCA model design method and the convolutional network hyperparameter
739   optimization scheme for side-channel attack are refined.
740   The CNNSCA model design method is refined: comprehensively utilize the advantages of VGG-
741   CNNSCA model classification and fitting efficiency and Alex-CNNSCA model occupy less
742   computing resources, while using SEnet's SE module to reduce the gradient dispersion problem
743   of error back propagation in deep neural networks to save calculation time, a new basic model of
744   CNNSCA was designed, named CNNSCAbase.
745

746   At the same time, the hyperparameter optimization scheme of the convolutional network used for
747   side-channel attacks is refined: design the structural parameter optimization experiment and the
748   training parameter optimization experiment, and use CNNSCAbase to implement the attack
749   training. According to parameter selection rules, common sense of parameter optimization of
750   CNN model, and data characteristics of actual application scenarios, the test parameters of each
751   experiment are designed, and unnecessary test parameters are excluded. Each time, according to
752   the cryptanalysis results of the experiment, the parameters that make CNNSCAbase's

753    cryptanalysis effect better are selected. Relying on two sets of experimental processes, a

754    hyperparameter optimization scheme is formed, and the hyperparameters finally determined by

755    the experiment are used as the parameters of the new model CNNSCAnew.

756

## Discussion

758    Comparative analysis of CNNSCAnew and other profiling side-channel attack methods

759    1 Comparative analysis of CNNSCAnew, classic template attack and MLPSCA

760    Experiment 13: Compare the cryptanalysis's performance of CNNSCAnew with the HW-based

761    TA[1] and MLPSCA method proposed by Benadjila et al[12]. TA and MLPSCA are the profiling

762    methods that performed better in the early traditional profiling methods and the later new

763    profiling methods, respectively. Experiment 13 carried out an attack on the ASCAD data set with

764    a known mask, which represents the realization of the encryption in an unprotected state. The

765    result of experiment 13 is shown in Figure 20(TA, MLPSCA, CNNSCAnew guessing entropy

766    convergence).

767

768    It can be seen from Figure 20 that CNNSCAnew's guessing entropy convergence is significantly

769    better than TA and MLPSCA.

770

771    2 Comparative analysis with other existing CNNSCA

772    Experiment 14: Compare the breaking performance of CNNSCAnew model with VGG-

773    CNNSCA[12] and Alex-CNNSCA[16]. The latter two methods are the profiling methods with

774    better performance among the latest profiling methods. Among them, VGG-CNNSCA in [12]

775    uses the ASCAD public data set, and Alex-CNNSCA in [16] uses a self-collected data set.

776    Experiment 14 carried out an attack on the ASCAD data set with a known mask, which

777    represents the realization of the encryption in an unprotected state. The result of experiment 14 is

778    shown in Figure 21(CNNSCAnew, VGG-CNNSCA, Alex-CNNSCA guessing entropy

779    convergence).

780

781    It can be seen from Figure 21 that the CNNSCAnew proposed in this paper has a better guessing

782    entropy convergence than other CNNSCAs. In [12], the guessing entropy of VGG-CNNSCA

783    requires at least 650 power consumption traces to converge to rank zero, and the model training

784    time takes 37 minutes. The CNNSCAnew method constructed in this paper only requires 61

785    power consumption traces, and the model training time only needs about 28 minutes. The

786    training time of CNNSCAnew and VGG-CNNSCA in this paper are shown in Figure

787    22(CNNSCAnew and VGG-CNNSCA training time).

788

789    After comparing CNNSCAnew with VGG-CNNSCA and Alex-CNNSCA, the model

790    comparison analysis and the cryptanalysis performance comparison analysis, the results are

791    summarized in Table 6(Comparative analysis of CNNSCAnew, VGG-CNNSCA and Alex-

792    CNNSCA) to show.

793

## Conclusions

795 Among the profiling side-channel cryptography attack methods, the most popular one is
796 CNNSCA, a side-channel attack method combined with deep learning convolutional neural
797 network algorithms. Its cryptanalysis performance is significantly better than traditional profiling
798 methods. Among the existing CNNSCA methods, the CNNSCA network models that achieve
799 cryptanalysis mainly include CNNSCA based on the VGG variant (VGG-CNNSCA) and
800 CNNSCA based on the Alexnet variant (Alex-CNNSCA). The learning capabilities and
801 cryptanalysis performance of these CNNSCA models it is not optimal. The paper aims to explore
802 effective methods to obtain the performance gains of the new side-channel attack method
803 CNNSCA.

804

805 After studying the related knowledge, necessary structure and core algorithm of CNNSCA, the
806 paper found that CNNSCA model design and hyperparameter optimization can be used to
807 improve the overall performance of CNNSCA. In terms of CNNSCA model design, the
808 advantages of VGG-CNNSCA model classification and fitting efficiency and the Alex-CNNSCA
809 model occupying less computing resources can be used to design a new CNNSCA basic model.
810 In order to better reduce the gradient dispersion problem of error back propagation in the deep
811 network, it is a very effective method to embed the SE module in this basic model; in terms of
812 the hyperparameter optimization of the CNNSCA model, the above basic model is applied to
813 side-channel leakage A known first-order mask data set in the public database (ASCAD). In this
814 specific application scenario, according to the model design rules and actual experimental
815 results, unnecessary experimental parameters can be excluded to the greatest extent. Various
816 hyperparameters of the model are optimized within the parameter interval to improve the
817 performance of the new CNNSCA, and the final determination benchmark for each
818 hyperparameter is given. Finally, a new CNNSCA model optimized architecture for attacking
819 unprotected encryption devices is obtained——CNNSCAnew. The paper also verified through
820 experimental comparison that CNNSCAnew's cryptanalysis effect is completely superior to
821 traditional profiling methods and the new profiling methods in literature [12,16]. In the literature
822 [12,16], the results of CNNSCA's guessing entropy are: convergence to 650 and oscillation. The
823 result of CNNSCAnew's guessing entropy proposed in this paper is to converge to a minimum of
824 61. Under the same experimental environment and experimental equipment conditions, literature
825 [12] took 40 minutes from model training to attacking the key, while the total calculation time of
826 CNNSCAnew was shortened to 30 minutes.

827

828 It should be noted that in practice, the results of each training of the CNNSCAnew model will
829 have a slight deviation. This is a normal phenomenon during neural network training and will not
830 affect the average performance of the model. While proposing the new CNNSCA method, the
831 paper also provides a more comprehensive and detailed design plan and optimization method for
832 the side-channel cryptanalysis researchers who need to design the CNNSCA model. In the

833  future, we can use these design schemes and optimization methods to continue to explore the
834  CNNSCA model that is more suitable for attacking protected equipment to achieve efficient
835  attacks on encrypted equipment with protection, which is of great significance to information
836  security and encryption protection.
837

## Acknowledgements

847

## References

849  [1] Mangard S, Oswald E, Popp T. Energy analysis attack [M]. Translated by Feng Dengguo,
850  Zhou Yongbin, and Liu Jiye. Beijing: Science Press, 2010.
851  [2] Kocher P, Jaffe J, Jun B. Differential power analysis[C]. Annual International Cryptology
852  Conference. Springer, Berlin, Heidelberg, 1999: 388-397.
853  [3] Brier E, Clavier C, Olivier F. Correlation power analysis with a leakage
854  model[C].International workshop on cryptographic hardware and embedded systems.
855  Springer,Berlin, Heidelberg, 2004: 16-29.
856  [4] Gierlichs B, Batina L, Tuyls P, Preneel B. Mutual information analysis[C]. International
857  Workshop on Cryptographic Hardware and Embedded Systems. Springer, Berlin,Heidelberg,
858  2008: 426-442.
859  [5] Chari S, Rao J R, Rohatgi P. Template attacks[C]. International Workshop on Cryptographic
860  Hardware and Embedded Systems. Springer, Berlin, Heidelberg, 2002:13-28.
861  [6] Lerman, L., Bontempi, G., Markowitch, O.: Power analysis attack: An approachbased on
862  machine learning. Int. J. Appl. Cryptol. 3(2) (June 2014) 97–115
863  [7] Picek, S., Heuser, A., Guilley, S.: Template attack versus Bayes classifier. Journal of
864  Cryptographic Engineering 7(4) (Nov 2017) 343–351
865  [8] Choudary, O., Kuhn, M.G.: Efficient template attacks. In Francillon, A., Rohatgi,P., eds.:
866  Smart Card Research and Advanced Applications - 12th International Conference, CARDIS
867  2013, Berlin, Germany, November 27-29, 2013. Revised Selected Papers. Volume 8419 of
868  LNCS., Springer (2013) 253–270
869  [9] Lerman, L., Poussier, R., Bontempi, G., Markowitch, O., Standaert, F.: Template Attacks vs.
870  Machine Learning Revisited (and the Curse of Dimensionality in Side-Channel Analysis). In:
871  COSADE 2015, Berlin, Germany, 2015. Revised Selected Papers. (2015) 20–33

872    [10] Lerman, L., Bontempi, G., Markowitch, O.: A machine learning approach against a masked
873    AES - Reaching the limit of side-channel attacks with a learning model.J. Cryptographic
874    Engineering 5(2) (2015) 123–139
875    [11] Picek, S., Heuser, A., Jovic, A., Legay, A.: Climbing down the hierarchy: Hierarchical
876    classification for machine learning side-channel attacks. In Joye, M., Nitaj,A., eds.: Progress in
877    Cryptology - AFRICACRYPT 2017: 9th International Conference on Cryptology in Africa,
878    Dakar, Senegal, May 24-26, 2017, Proceedings,Cham, Springer International Publishing (2017)
879    61–78
880    [12] Benadjila R, Prouff E, Strullu R, Cagli E, Dumas C. Study of deep learning techniques for
881    side-channel analysis and introduction to ASCAD database[J]. ANSSI, France & CEA,LETI,
882    MINATEC Campus, France.2018, 22: 2018.
883    [13] Maghrebi, H., Portigliatti, T., Prouff, E.: Breaking cryptographic implementations using
884    deep learning techniques. In: Security, Privacy, and Applied Cryptography Engineering 6th
885    International Conference, SPACE 2016, Hyderabad, India, December 14-18, 2016, Proceedings.
886    (2016) 3–26
887    [14] Picek, S., Samiotis, I.P., Heuser, A., Kim, J., Bhasin, S., Legay, A.: On the performance of
888    convolutional neural networks for side-channel analysis. Cryptology ePrint Archive, Report
889    2018/004 (2018) https://eprint.iacr.org/2018/004.
890    [15] Cagli E, Dumas C, Prouff E. Convolutional neural networks with data augmentation against
891    jitter-based countermeasures[C]. International Conference on Cryptographic Hardware and
892    Embedded Systems. Springer, Cham, 2017: 45-68.
893    [16] Guo Dongxin, Chen Kaiyan, Zhang Yang, Hu Xiaoyang, Wei Yanhai. A new method for
894    attacking encrypted chip templates based on Alexnet convolutional neural network. Computer
895    Measurement and Control, 2018. 26(10): Pages 246-249+254.
896    [17] Guo Dongxin, Chen Kaiyan, Zhang Yang, Zhang Xiaoyu, Li Jianlong. A new method of
897    attacking encrypted chip templates based on VGGNet convolutional neural network. Computer
898    Application Research, 2019. 36(09): Page 2809-2812+2855.
899    [18] Kim J, Picek S, Heuser A, Bhasin S,  Hanjalic A. Make some noise. unleashing the power of
900    convolutional neural networks for profiled side-channel analysis[J]. IACR Transactions on
901    Cryptographic Hardware and Embedded Systems, 2019: 148-179.
902    [19] Krizhevsky,A.,Sutskever,I.,Hinton,G.E.:Imagenet classification with deep convolutional
903    neural networks. Commun. ACM 60(6),84–90 (2017). https://doi.org/10.1145/3065386
904    [20] Zeiler,M.D.,Fergus,R.:Visualizing and understanding convolutional networks. In:European
905    Conference on Computer Vision,pp.818–833. Springer (2014)
906    [21] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image
907    recognition (2014). arXiv preprintarXiv:1409.1556
908    [22] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D.,Erhan, D., Vanhoucke,
909    V., Rabinovich, A.: Going deeper with convolutions. In: Proceedings of the IEEE Conference on
910    Computer Vision and Pattern Recognition, pp. 1–9 (2015)

911 [23] He,K.,Zhang,X.,Ren,S.,Sun,J.:Deep residual learning for image recognition. In: Proceedings
912 of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)
913 [24] Jie H, Li S, Gang S, Albanie S. Squeeze-and-Excitation Networks[J]. IEEE Transactions on
914 Pattern Analysis and Machine Intelligence, 2017.
915 [25] Eldeeb A, Bursztein E, Chollet F, Haifeng Jin, Qianli Scott Zhu.: Keras (2015).
916 https://github.com/fchollet/keras
917 [26] Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado G.S., Davis A, Dean J,
918 Devin M, Ghemawat S,Goodfellow I, Harp A, Irving G, Isard M, Jia Y, Jozefowicz R, Kaiser L,
919 Kudlur M, Levenberg J, Mané D, Monga R, Moore S, Murray D, Olah C, Schuster M, Shlens J,
920 Steiner B, Sutskever I, Talwar K, Tucker P, Vanhoucke V, Vasudevan V, Viégas F, Vinyals O,
921 Warden P, Wattenberg M, Wicke M, Yu Y, Zheng X.: TensorFlow: Large-scale machine
922 learning on heterogeneous systems (2015). https://www.tensorflow.org/. Software available from
923 tensorflow.org.
924 [27] D. H. Hubel T. N. Wiesel "Receptive Fields And Functional Aechitecture of Monkey Striate
925 Cortex"[J]. Physiol. (1968), 195, pp. 215-243 With 3 plates and 14 text figures Printed in Great
926 Britain.
927 [28] Lecun Y, Bengio Y. Convolutional Networks for Images, Speech, and Time Series[J].  1998.
928 [29]Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A,
929 Bernstein M. :Imagenet large scale visual recognition challenge.Int.J.Comput.Vis.115(3),211–
930 252 (2015).
931 [30] Kim, T., Lee, J., Nam, J.: Sample-level CNN architectures for music auto-tagging using raw
932 waveforms. In: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing,
933 ICASSP 2018, Calgary, AB, Canada, April 15-20, 2018.(2018) 366–370
934 [31] Choi K, Fazekas G, Sandler M, Cho K. Convolutional Recurrent Neural Networks for
935 Music Classification[J]. IEEE, 2016.
936 [32] Ioffe, S., Szegedy, C.: Batch normalization: accelerating deep network training by reducing
937 internal covariate shift. CoRR (2015).arXiv:1502.03167.
938 [33] Goodfellow, I.J., Bengio, Y., Courville, A.C.: Deep Learning.Adaptive Computation and
939 Machine Learning. MIT Press, Cambridge (2016).
940 [34] Han Liqun, Kang Qian. Artificial Neural Network Theory, Design and Application——
941 Nerve Cells, Neural Networks and Neural System[J].Journal of Beijing Technology and
942 Business University: Natural Science Edition, 2005, 23(1): 52-52.
943 [35] Hawkins D M. The problem of overfitting[J]. Journal of chemical information and computer
944 sciences, 2004, 44(1): 1-12.
945 [36] Standaert F X, Malkin T G, Yung M. A unified framework for the analysis of side-channel
946 key recovery attacks[C]. Annual international conference on the theory and applications of
947 cryptographic techniques. Springer, Berlin, Heidelberg, 2009: 443-461.
948 [37] Masure L, Dumas C, Prouff E. A comprehensive study of deep learning for side-channel
949 analysis[J]. IACR Transactions on Cryptographic Hardware and Embedded Systems, 2020: 348-
950 375.

951   [38] Lecun Y, Cortes C. The mnist database of handwritten digits[J].

952   http://www.research.att.com/∽yann/ocr/mnist/, 2010.

953   [39] Smith S L, Kindermans P J, Ying C, Le Q V. Don't Decay the Learning Rate, Increase the
954   Batch Size[J]. 2017.

955   [40] Soltanolkotabi M, Javanmard A, Lee J D. Theoretical insights into the optimization
956   landscape of over-parameterized shallow neural networks[J]. IEEE Transactions on Information
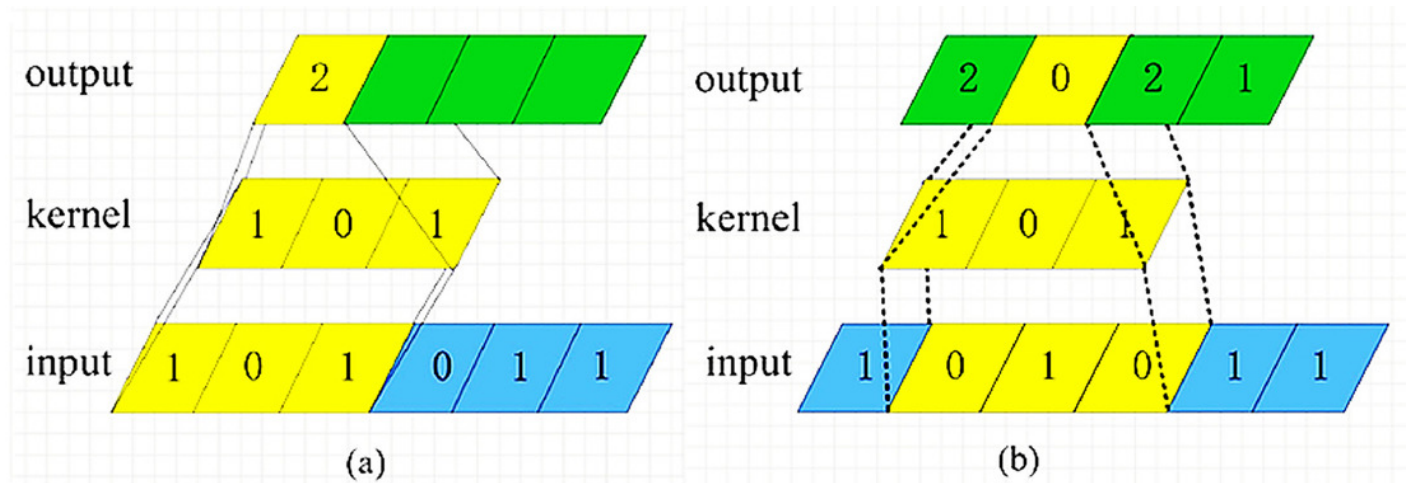957   Theory, 2017, PP(2):742-769.

# Figure 1

SEnet module

$X$

$D \mathrm{x} 1 \mathrm{x} W$

Global pooling    $D \mathrm{x} 1 \mathrm{x} 1$

FC    $D/rate \ \mathrm{x} 1 \mathrm{x} 1$

FC    $D \mathrm{x} 1 \mathrm{x} 1$

Sigmoid    $D \mathrm{x} 1 \mathrm{x} 1$

Scale    $D \mathrm{x} 1 \mathrm{x} W$

$+$   $\tilde{X}$    $D \mathrm{x} 1 \mathrm{x} W$

# Figure 2

Convolution_calculation_process

# Figure 3

Pooling_calculation_process

# Figure 4

Convolutional network structure in a side-channel attack scenario

# Figure 5

Training_effect_of_CNNSCAnoSE_model_and_CNNSCAbase_model

# Figure 6

Convergence_of_guessing_entropy_of_Cnov1~3

Each curve represents the convergence trend of guessing entropy under the six model structures of Cnov1~3 . The abscissa represents the number of energy traces used in the attack, and the ordinate represents the ranking of the guessing entropy.

# Figure 7

Convergence_of_guessing_entropy_of_Cnov1~2,4~7

Each curve represents the convergence trend of guessing entropy under the six model structures of Cnov1~2,4~7 . The abscissa represents the number of energy traces used in the attack, and the ordinate represents the ranking of the guessing entropy.

# Figure 8
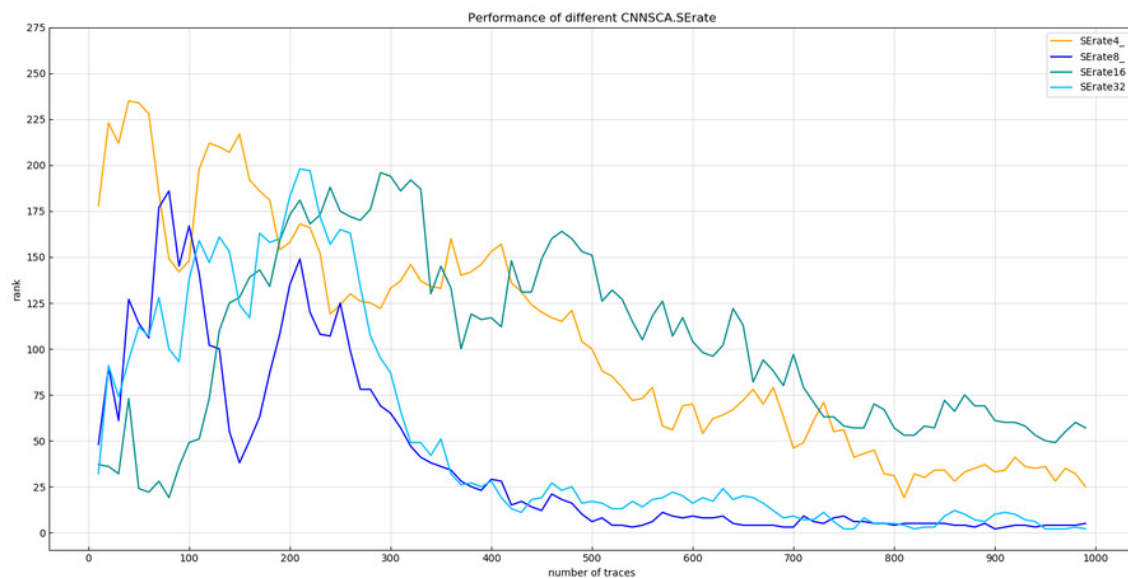
Convergence_of_guessing_entropy_of_filter1~4_(epochs=75)

Each curve represents the model guessing entropy of the four convolution kernel sizes. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 9

Convergence_of_guessing_entropy_of_filter1~3_(epochs=40)

Each curve represents the model guessing entropy of the three convolution kernel sizes. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 10

Convergence_of_guessing_entropy_of_AveragePool_and_MaxPool_structure
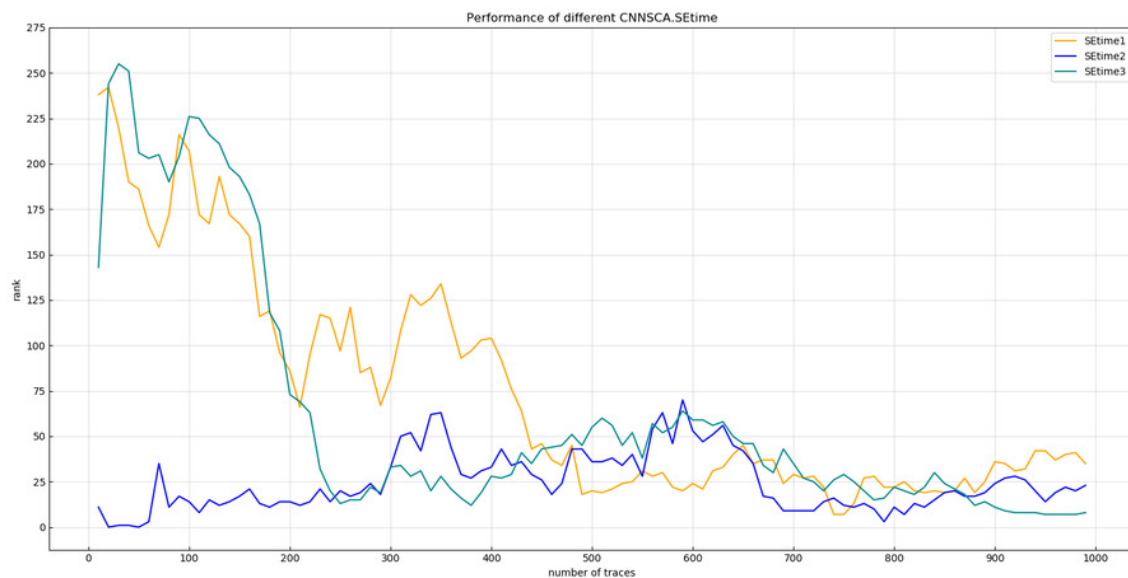
Each curve represents the model guessing entropy of two pooling methods. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 11

Convergence_of_guessing_entropy_of_different_convolution_kernel_size_structures

Each curve represents the model guessing entropy of 5 convolution kernel sizes structures. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 12

Convergence_of_guessing_entropy_for_different_SE_dimension_ratios

Each curve represents the model guessing entropy of the four SE dimension ratios . The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 13

Convergence_of_guessing_entropy_for_different_number_of_SE_cycles

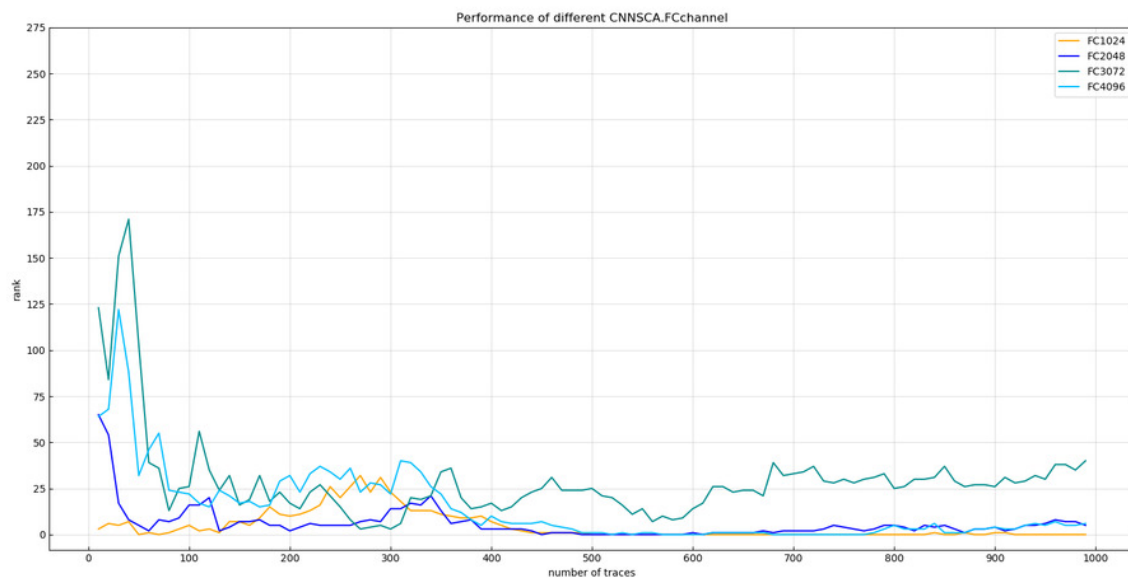Each curve represents the model guessing entropy of the three SE cycles . The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 14

Convergence_of_guessing_entropy_of_the_four_channel_number_structure_of_FC_layer
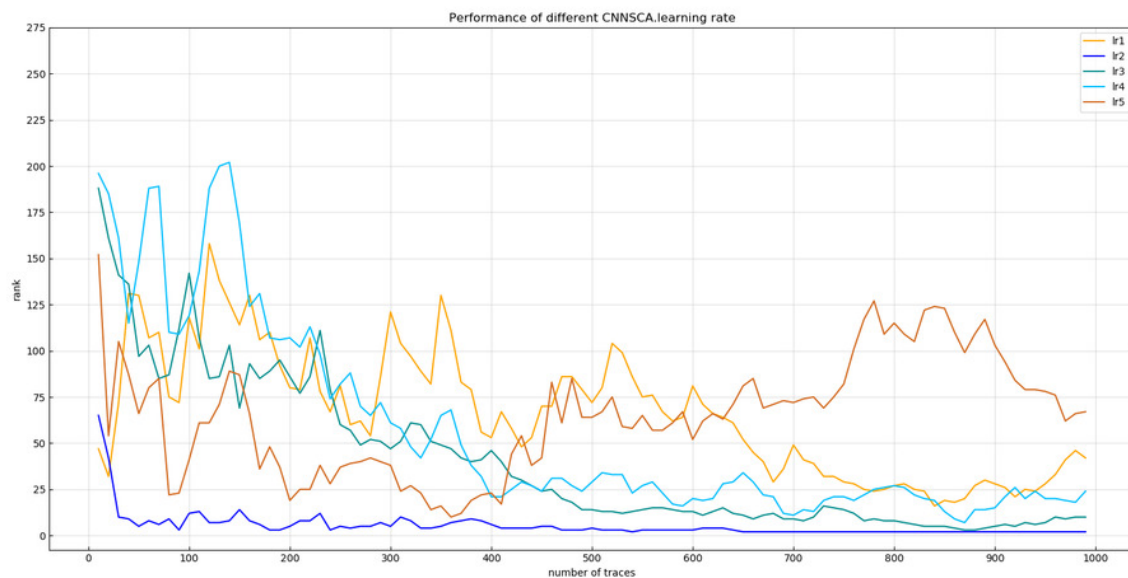
Each curve represents the model guessing entropy of the four FC layer structures. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 15

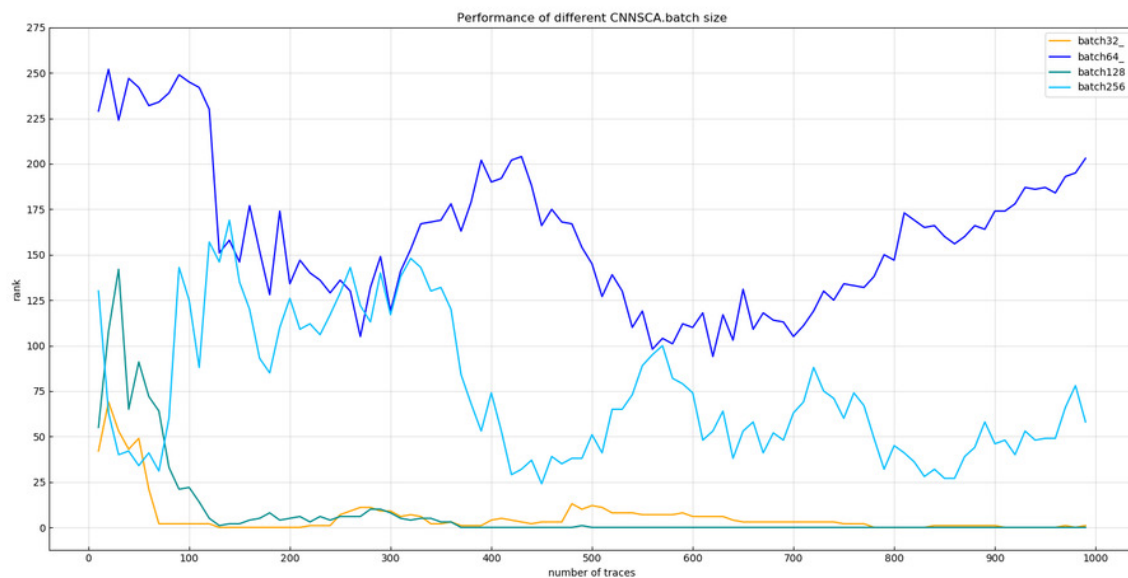Convergence_of_model_guessing_entropy_under_five_learning_rates

Each curve represents the model guessing entropy of five learning rates . The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 16

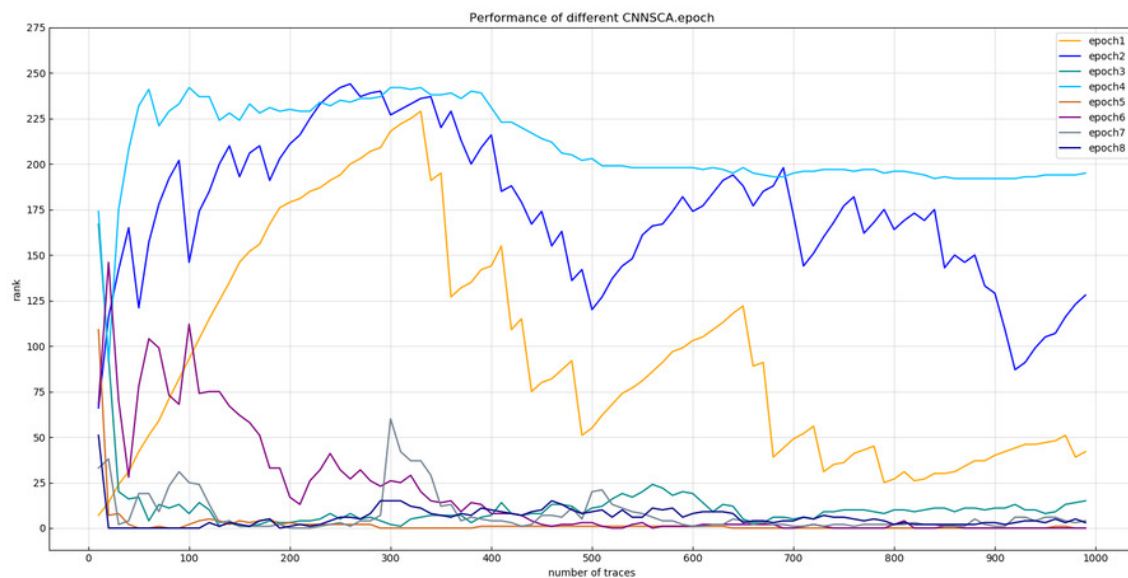Convergence_of_model_guessing_entropy_under_four_batches

Each curve represents the model guessing entropy of the four training batches. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 17

Convergence_of_model_guessing_entropy_under_eight_epochs

Each curve represents the model guessing entropy of eight training epochs . The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 18
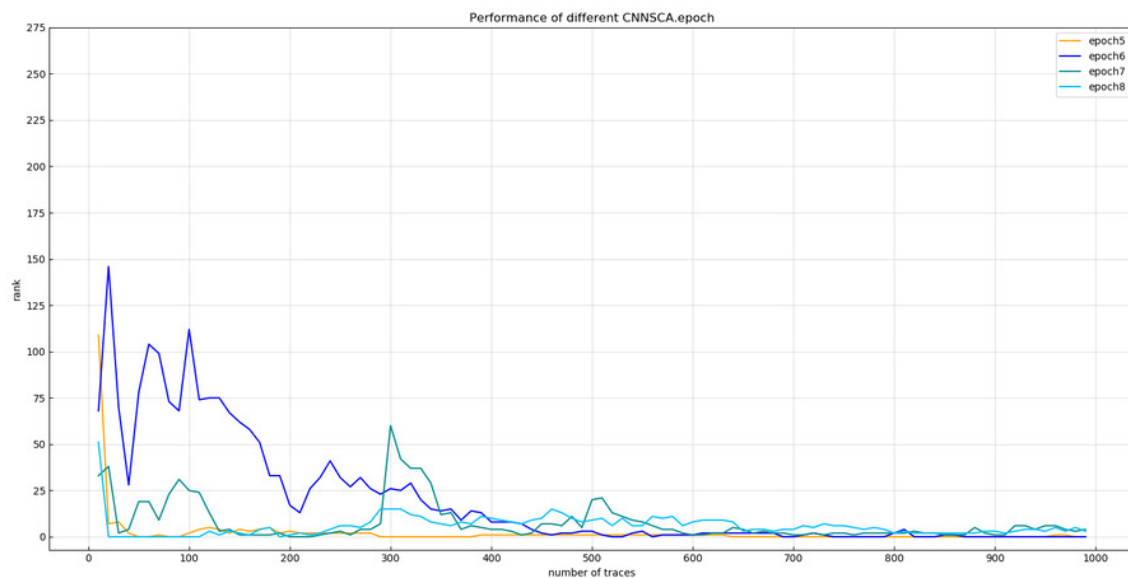
Convergence_of_model_guessing_entropy_under_four_epochs

Each curve represents the model guessing entropy of four training epochs . The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 19

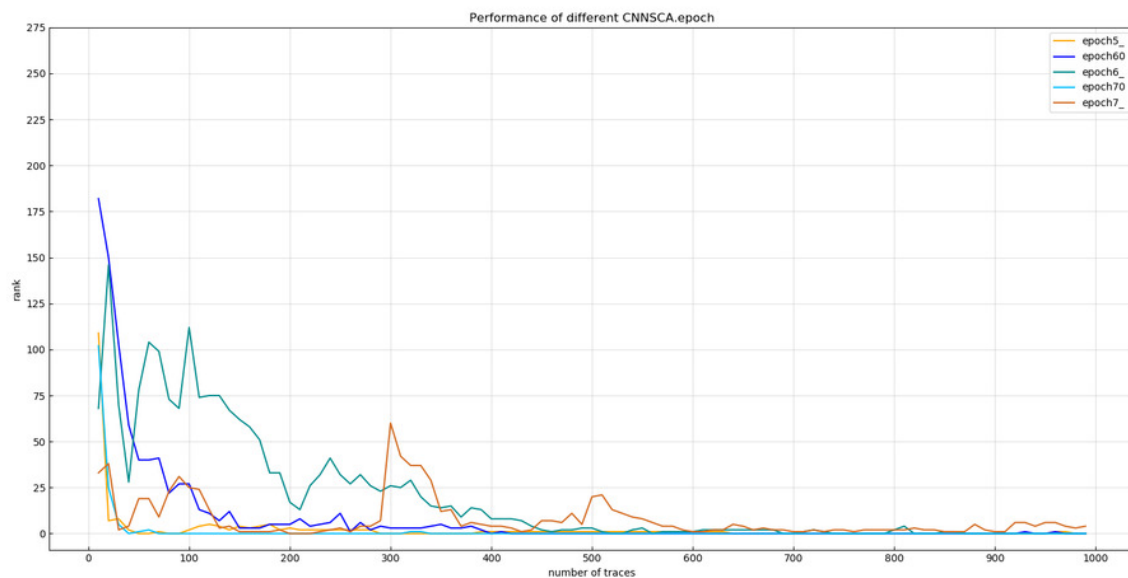Convergence_of_model_guessing_entropy_under_five_epochs

Each curve represents the model guessing entropy of five training epochs . The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 20

TA,_MLPSCA,_CNNSCAnew_guessing_entropy_convergence

Each curve represents the model guessing entropy of TA, MLPSCA, and CNNSCAnew respectively. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 21

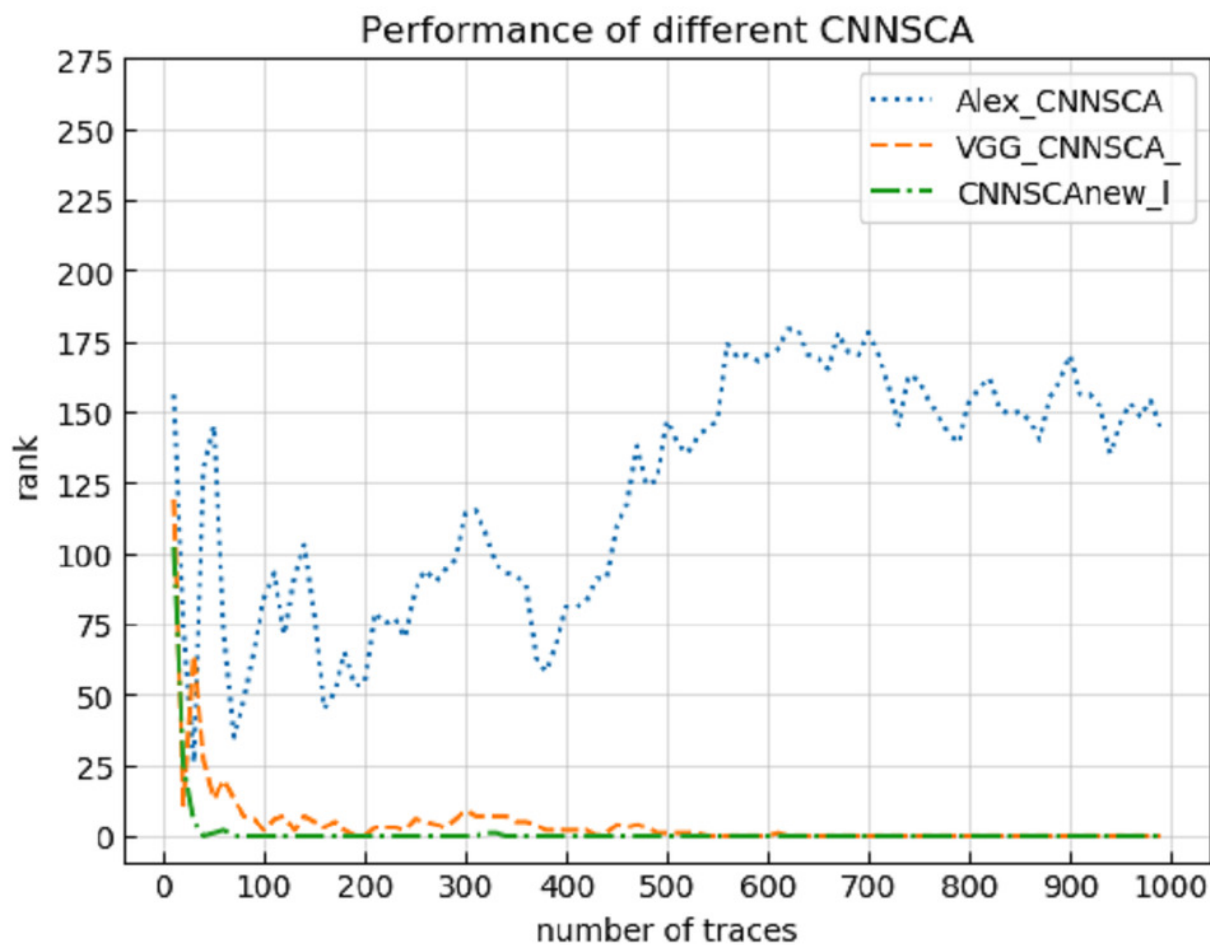CNNSCAnew,_VGG-CNNSCA,_Alex-CNNSCA_guessing_entropy_convergence

Each curve represents the model guessing entropy of CNNSCAnew, VGG-CNNSCA and Alex-CNNSCA respectively. The abscissa represents the number of energy trajectories used in the attack, and the ordinate represents the order of guessing entropy.

# Figure 22

CNNSCAnew_and_VGG-CNNSCA_training_time

**CNNSCAnew**    Running time:1653.920090106 Seconds

**VGG-CNNSCA**    Running time:2233.5491708040004 Seconds

**Table 1**(on next page)

CNN with outstanding performance in previous ILSVRC competitions

Table 1 CNN with outstanding performance in previous ILSVRC competitions

| Year | Network / Ranking | val top-1 | val top-5 | test top-5 | Remarks |
|---|---|---|---|---|---|
| 2012 | Alexnet（Champion） | 36.7% | 15.4% | 15.32% | 7CNNs、Used data from 2011 |
| 2013 | ZFnet（Champion） | -- | -- | 13.51% | The result on the ZFNet paper is 14.8 |
| 2014 | VGG（Runner-up） | 23.7% | 6.8% | 6.8% | Post-race、2 nets |
| 2014 | Googlenet v4（Champion） | 16.5% | 3.1% | 3.08% | Post-race、v4+Inception-Res-v2 |
| 2015 | Resnet（Champion） | -- | -- | 3.57% | 6 models |
| 2016 | Trimps-Soushen（Champion） | -- | -- | 2.99% | Public Security III (additional data) |
| 2017 | SEnet（Champion） | -- | -- | 2.25% | Momenta and Oxford University |

**Table 2**(on next page)

CNNSCAbase Configuration

Table2 CNNSCAbase Configuration

| ConvNet Configuration | | | |
|---|---|---|---|
| Input(1x700 vector) | | | |
| Block1 | (Conv3-64)x1 | Same\ReLU | AveragePool (2,2) |
| Block2 | (Conv3-128)x1 SE | Same\ReLU | AveragePool (2,2) |
| Block3 | (Conv3-256)x1 SE | Same\ReLU | AveragePool (2,2) |
| Block4 | (Conv3-512)x1 SE | Same\ReLU | AveragePool (2,2) |
| Block5 | (Conv3-1024)x1 SE | Same\ReLU | AveragePool (2,2) |
| (FC-4096)x2，ReLU | | | |
| (FC-256)x1，Soft-max | | | |
| Model compile（crossentropy、RMSprop） | | | |

**Table 3**(on next page)

CNNSCAbase.Conv1-7 Configuration

Table3 CNNSCAbase.Conv1-7 Configuration

| | Conv Configuration | | | | | | |
|---|---|---|---|---|---|---|---|
| Block | Conv1 | Conv2 | Conv3 | Conv4 | Conv5 | Conv6 | Conv7 |
| Block1 | (Conv3-64)x1 | (Conv3-64)x2 | (Conv3-64)x2 | (Conv3-64)x1 | (Conv3-64)x1 | (Conv3-64)x1 | (Conv3-64)x1 |
| Block2 | (Conv3-128)x1 SE | (Conv3-128)x2 SE | (Conv3-128)x2 SE | (Conv3-128)x2 SE | (Conv3-128)x1 SE | (Conv3-128)x1 SE | (Conv3-128)x1 SE |
| Block3 | (Conv3-256)x1 SE | (Conv3-256)x2 SE | (Conv3-256)x2 SE | (Conv3-256)x2 SE | (Conv3-256)x2 SE | (Conv3-256)x1 SE | (Conv3-256)x1 SE |
| Block4 | (Conv3-512)x1 SE | (Conv3-512)x2 SE | (Conv3-512)x2 SE | (Conv3-512)x2 SE | (Conv3-512)x2 SE | (Conv3-512)x2 SE | (Conv3-512)x1 SE |
| Block5 | (Conv3-1024)x1 SE | (Conv3-1024)x2 SE | (Conv3-1024)x3 SE | (Conv3-1024)x2 SE | (Conv3-1024)x2 SE | (Conv3-1024)x2 SE | (Conv3-1024)x2 SE |

**Table 4**(on next page)

CNNSCAnew Configuration

Table4 CNNSCAnew Configuration

| ConvNet Configuration | | | |
|---|---|---|---|
| Input(1x700 vector) | | | |
| Block1 | (Conv3-32)x1 | Same\ReLU | AveragePool (2,2) |
| Block2 | (Conv3-64)x1 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| Block3 | (Conv3-128)x2 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| Block4 | (Conv3-256)x2 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| Block5 | (Conv3-512)x2 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| (FC-1024)x2，ReLU | | | |
| (FC-256)x1，Soft-max | | | |
| Model compile（crossentropy、RMSprop） | | | |

**Table 5**(on next page)

CNNSCAnew Configuration

Table5 CNNSCAnew Configuration

| ConvNet Configuration | | | |
|---|---|---|---|
| Input(1x700 vector) | | | |
| Block1 | (Conv3-32)x1 | Same\ReLU | AveragePool (2,2) |
| Block2 | (Conv3-64)x1 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| Block3 | (Conv3-128)x2 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| Block4 | (Conv3-256)x2 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| Block5 | (Conv3-512)x2 SEx2（1/8） | Same\ReLU | AveragePool (2,2) |
| (FC-1024)x2，ReLU | | | |
| (FC-256)x1，Soft-max | | | |
| Model compile（crossentropy、RMSprop） | | | |
| Training parameters（$1x10^{-3}$、128、70） | | | |

**Table 6**(on next page)

Comparative analysis of CNNSCAnew, VGG-CNNSCA and Alex-CNNSCA

Table 6 Comparative analysis of CNNSCAnew, VGG-CNNSCA and Alex-CNNSCA

| Three CNNSCA | | CNNSCAnew | VGG-CNNSCA | Alex-CNNSCA |
|---|---|---|---|---|
| CNNSCA Configuration | Convol block1 | (Conv3-32)x1 Same\ReLU AveragePool (2,2) | (Conv11-64)x1 Same\ReLU AveragePool (2,2) | (Conv11-96)x1 Same\ReLU MaxPool (2,2) |
| | Convol block2 | (Conv3-64)x1 SEx2（1/8） Same\ReLU AveragePool (2,2) | (Conv11-128)x1 Same\ReLU AveragePool (2,2) | (Conv5-256)x1 Same\ReLU MaxPool (2,2) |
| | Convol block3 | (Conv3-128)x2 SEx2（1/8） Same\ReLU AveragePool (2,2) | (Conv11-256)x1 Same\ReLU AveragePool (2,2) | (Conv3-384)x1 Same\ReLU |
| | Convol block4 | (Conv3-256)x2 SEx2（1/8） Same\ReLU AveragePool (2,2) | (Conv11-512)x1 Same\ReLU AveragePool (2,2) | (Conv3-384)x1 Same\ReLU |
| | Convol block5 | (Conv3-512)x2 SEx2（1/8） Same\ReLU AveragePool (2,2) | (Conv11-512)x1 Same\ReLU AveragePool (2,2) | (Conv3-256)x1 Same\ReLU MaxPool (3,3) |
| | dense layer | (FC-1024)x2，ReLU (FC-256)x1，Soft-max | (FC-4096)x2，ReLU (FC-256)x1，Soft-max | (FC-4096)x2，ReLU (FC-256)x1，Soft-max |
| CNNSCA Performance | Learning rate | $1 \times 10^{-3}$ | $10^{-5}$ | $10^{-2}$ |
| | Batch size | 128 | 200 | 10 |
| | epoch | 70 | 75 | 20 |
| | calculating time | 28 minute | 37 minute | 2 hour |
| | Guess entropy | 61 | 650 | Did not converge |