

# Optics Letters

## Learning complex scattering media for optical encryption

LINA ZHOU,  YIN XIAO, AND WEN CHEN\*

Department of Electronic and Information Engineering, The Hong Kong Polytechnic University, Hong Kong, China

\*Corresponding author: owen.chen@polyu.edu.hk

Received 15 June 2020; revised 11 August 2020; accepted 12 August 2020; posted 13 August 2020 (Doc. ID 400174); published 15 September 2020

**Optical encryption has provided a new insight for securing information; however, it is always desirable that high security can be achieved to withstand the attacks. In this Letter, we propose a new method via learning complex scattering media for optical encryption. After the recordings through complex scattering media, a designed learning model is trained. The proposed method uses an optical setup with complex scattering media to experimentally record the ciphertexts and uses a learning model to generate security keys. During the decryption, the trained learning model with its parameters is applied as security keys. In addition, various parameters, e.g., virtual phase-only masks, can be flexibly applied to further enlarge key space. It is experimentally demonstrated that the proposed learning-based encryption approach possesses high security. The proposed method could open up a new research perspective for optical encryption.** © 2020 Optical Society of America

<https://doi.org/10.1364/OL.400174>

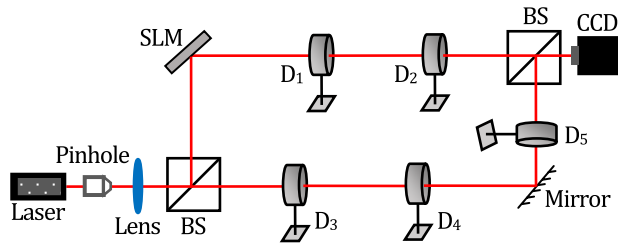
Optical encryption has become one of the advanced encoding methods owing to its outstanding properties [1–5]. Optical encoding technologies originated from the milestone work done by Refregier and Javidi, who put forward the idea of double random phase encoding (DRPE) [3]. Since then, there is an explosion of optical encryption schemes, such as the fractional Fourier transform domain and the Fresnel transform domain [6,7]. With a rapid development of optical techniques, different optical cryptosystems are accordingly proposed, such as interferometric imaging, diffractive imaging, and ghost imaging [8–12]. However, it has been found that optical encoding schemes cannot withstand some attacks [13–17]. For instance, Carnicer *et al.* proposed a groundbreaking philosophy of a chosen-ciphertext attack to vet the vulnerability of the DRPE scheme [13]. Similarly, chosen-plaintext, known-plaintext attack, and ciphertext-only attack provided an insight for the cryptanalysis of optical encryption [14–16]. The major objective in the developed optical cryptanalyses was to extract an estimated plaintext from the ciphertext by using the estimated security keys. Apart from the aforementioned attacking technologies, another method, called learning-based attack, has also been developed [17]. It can allow a direct extraction of unknown plaintexts from the given ciphertexts without

individual retrieval of various security keys or the usage of complex-phase retrieval algorithms. The recent progress in optical cryptanalysis becomes a serious threat to optical encryption schemes, which would request the advances in optical encryption methods or infrastructures.

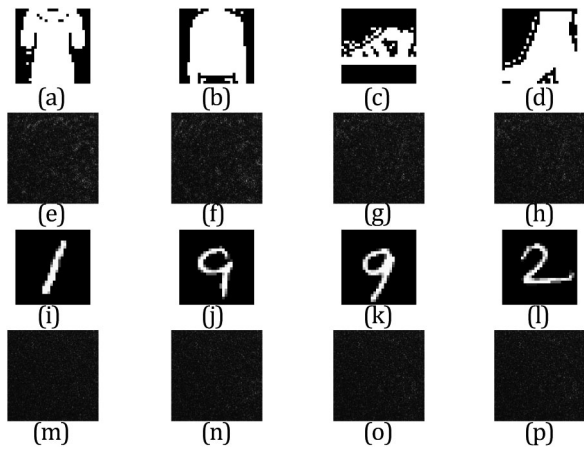
Inspired by remarkable characteristics of machine learning [18–22], we introduce machine learning into optical cryptography. In this Letter, learning complex scattering media for optical encryption is proposed for the first time to our knowledge. The proposed method uses an optical setup with complex scattering media to experimentally record the ciphertexts and uses a learning model to generate security keys. It is experimentally demonstrated that the proposed learning-based encryption is feasible and effective and possesses high security.

A schematic experimental setup is shown in Fig. 1, and complex scattering media can be flexibly designed and applied in the optical setup. The collimated He–Ne laser beam with a wavelength of 633.0 nm propagates through a beam splitter to be separated into two beams. One beam illuminates a reflective spatial light modulator (SLM, Holoeye LC-R720). The input grayscale images (i.e., plaintexts) are sequentially embedded into the SLM that are fashion-product images from the fashion Modified National Institute of Standards and Technology (MNIST) database [18] or handwritten-digit images from the MNIST database [19]. Then, the modulated laser beam successively propagates through the diffusers  $D_1$  and  $D_2$ . The reference beam is sequentially diffracted by the diffusers  $D_3$ ,  $D_4$ , and  $D_5$ , and then interferes with the object beam. The interference patterns recorded by a CCD are used as ciphertexts (size of  $600 \times 600$  pixels) in this study. Figures 2(a)–2(p) show several plaintexts selected from the two databases [18,19] and their corresponding ciphertexts experimentally recorded by the CCD.

It has been demonstrated that various attacks [13–17] pose a great threat to optical cryptosystems. Here, we propose learning-based optical encryption. Using machine learning technologies [18–22], we generate a trained model with its parameters as security keys. The decryption process is not the directly reverse process of encryption. Meanwhile, parameters in the optical setup used to record the ciphertexts can be discarded. The ciphertexts and the plaintext-ciphertext pairs cannot be directly generated by the attackers. The proposed method uses an optical setup with complex scattering media to experimentally



**Fig. 1.** Schematic experimental setup with complex scattering media: SLM, spatial light modulator (Holoeye LC-R720);  $D_1$ ,  $D_2$ ,  $D_3$ ,  $D_4$ , and  $D_5$ , diffusers (Thorlabs, DG10-600); BS, beam splitter cube; CCD, charge-coupled device (Thorlabs, DCC1240C).



**Fig. 2.** (a)–(d) and (i)–(l) Typical input images (i.e., plaintexts) embedded into the SLM, and (e)–(h) and (m)–(p) the corresponding speckle patterns (i.e., ciphertexts) experimentally recorded by a CCD.

record the ciphertexts and uses a learning model to generate security keys.

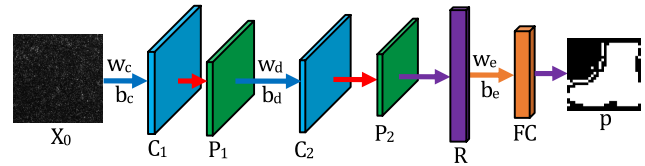
Figure 3 shows a designed convolutional neural network (CNN) for the encryption. Five thousand images are randomly selected from each database to serve as plaintexts, and their corresponding speckle patterns (i.e., ciphertexts) are recorded by the CCD shown in Fig. 1. The designed learning model is composed of two convolutional layers followed by one pooling layer after each convolutional layer, one reshaping layer, and one fully connected layer. The ciphertext is preprocessed by resizing it to reduce the dimensionality and computational complexity. The size of convolutional kernels depends on the dimension of the input, and the number of the kernels is determined accordingly. The first convolutional layer is labeled as  $C_1$  with weights and biases respectively denoted as  $w_c$  and  $b_c$ . The feature map for  $C_1$  can be described by

$$x_1 = \sigma[(w_c * x_0) + b_c], \quad (1)$$

where  $x_0$  denotes the ciphertext,  $*$  denotes convolution, and  $\sigma$  denotes the activation function used in  $C_1$ . After down sampling, the first pooling layer ( $P_1$ ) is processed by convolution, and the feature map for  $C_2$  is given by

$$x_2 = \sigma[(w_d * x_p) + b_d], \quad (2)$$

where  $w_d$  and  $b_d$ , respectively, represent weights and biases of the kernels used in  $C_2$ , and  $x_p$  denotes the feature map of  $P_1$ .



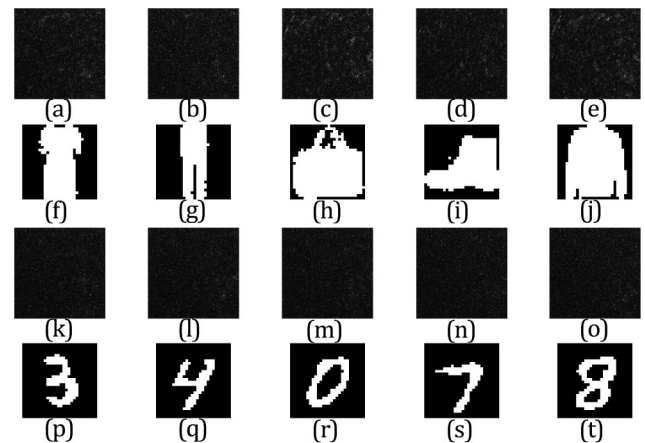
**Fig. 3.** Framework for the proposed learning-based encryption:  $x_0$ , ciphertext;  $C_1$  and  $C_2$ , the first and second convolutional layer;  $P_1$  and  $P_2$ , the first and second pooling layer; R, shaping layer; FC, fully connected layer; p, plaintext.

Subsequently, the pooling layer ( $P_2$ ) is formed by down sampling. After image resizing,  $P_2$  is reshaped to a one-dimensional vector ( $R$ ). To achieve a prediction of the plaintext,  $R$  is connected to a fully connected layer (FC) with the usage of weights ( $w_e$ ) and biases ( $b_e$ ). The feature map of FC is described by

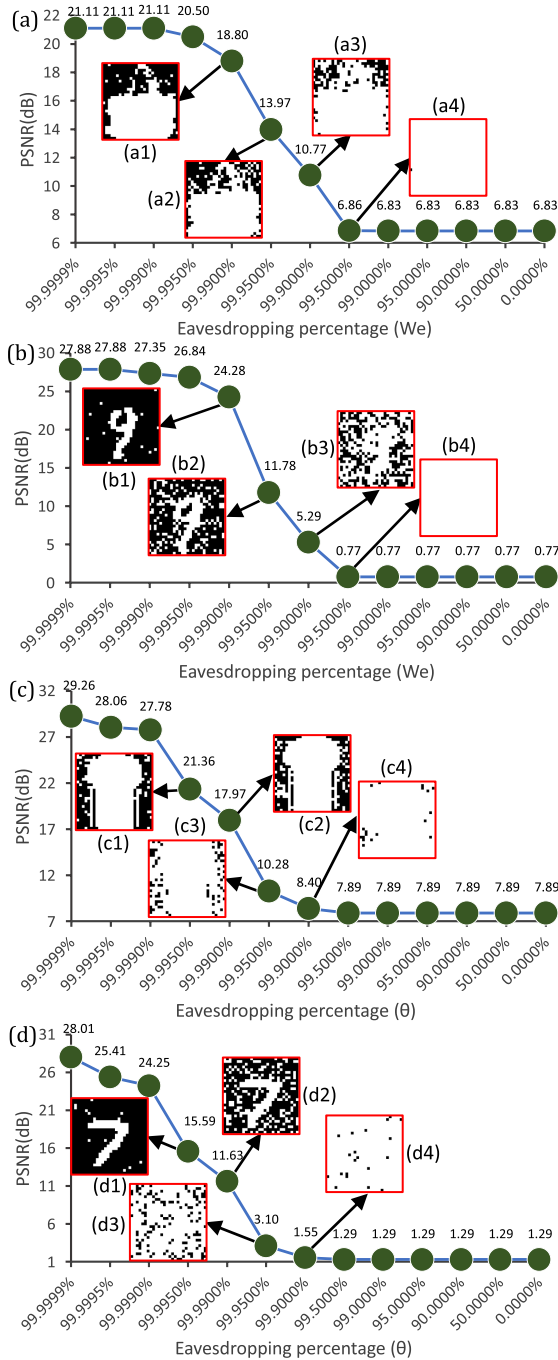
$$x_{FC} = (w_e * x_R) + b_e, \quad (3)$$

where  $x_R$  denotes feature map of the reshaping layer. Then, the one-dimensional vector of FC is reshaped to a two-dimensional vector, which is the ultimate prediction. Here,  $\theta$  is a combination of all the weights and biases.

In this study, the input ciphertext is resized from  $600 \times 600$  pixels to  $100 \times 100$ . The resized ciphertext convolves with 20 convolutional kernels (size of  $5 \times 5$ ) forming the first convolutional layer with size of  $96 \times 96 \times 20$ . Weights  $w_c$  (dimension of  $5 \times 5 \times 20$ ) and biases  $b_c$  (dimension of  $20 \times 1$ ) in Eq. (1) are randomly initialized. Activation function adopted for each convolutional layer is sigmoid. After down sampling, the first pooling layer is generated with the size of  $48 \times 48 \times 20$ . Down-sampling size is 2 for each pooling layer. The pooled data is further processed by convolution with 20 kernels (size of  $5 \times 5$ ) forming the second convolutional layer with the size of  $44 \times 44 \times 20$ . Weights  $w_d$  (dimension of  $5 \times 5 \times 20$ ) and biases  $b_d$  (dimension of  $20 \times 1$ ) in Eq. (2) are randomly initialized. Next, down-sampling processing is adopted again to generate the second pooling layer with the size of  $22 \times 22 \times 20$ . Subsequently, the second pooling layer is reshaped from a three-dimensional vector to a one-dimensional vector with



**Fig. 4.** Decrypted images obtained by using all correct security keys: (a)–(e) and (k)–(o) ciphertexts; (f)–(j) and (p)–(t) the retrieved plaintexts. Peak signal-to-noise ratios (PSNR) for (f)–(j) and (p)–(t) are 25.22 dB, 26.76 dB, 18.61 dB, 22.04 dB, 28.34 dB, 28.14 dB, 23.69 dB, 25.26 dB, 29.20 dB, and 22.63 dB, respectively.



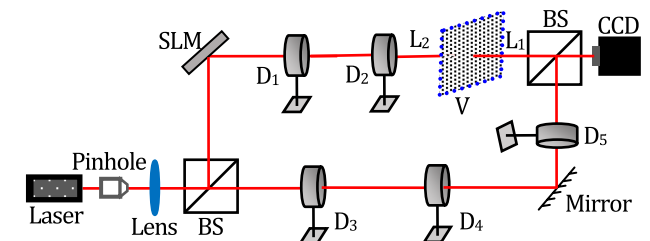
**Fig. 5.** Eavesdropping analysis of security keys  $w_e$  and  $\theta$ .

the size of  $1 \times 9680$ . The reshaped data is fully connected to the reshaped plaintext (size of  $1 \times 784$ ). Weights  $w_e$  (dimension of  $784 \times 9680$ ) and biases  $b_e$  (dimension of  $9680 \times 1$ ) in Eq. (3) are randomly initialized. Four thousand eight hundred speckle patterns and their corresponding plaintexts are used as the training data, and 200 other ciphertexts are used to test the learning model. To evaluate the difference between the retrieved plaintexts and original plaintexts, the mean squared error (MSE) is calculated. When the MSE value is higher than a preset threshold, the error is back-propagated, and then the weights and biases of each layer are updated by stochastic gradient descent [22]. Here, the training epoch is selected to be 5.

The momentum is set to be  $-9.5 \times 10^{-4}$ , and the learning rate is  $10^{-6}$ . The learning model is trained by using Matlab 2009 running on a PC with Intel Core i7 at 8 GHz, 64 GB RAM, and Nvidia GTX1080Ti. Total time taken for the model training is about 4.0 h. After the training, the unknown plaintext can be retrieved in real time by using the trained learning model with its parameters.

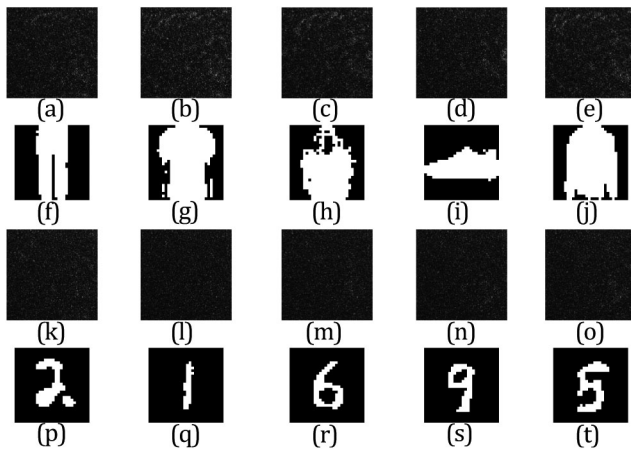
The trained learning model with its parameters, e.g., the size of kernels, the number of kernels, activation function,  $w_e$ ,  $w_d$ ,  $w_c$ , and  $\theta$  can be used as security keys. The security keys can also be updated with a new training iteration. It needs to be pointed out that security keys can be obtained by training the designed model using multiple databases. To save time, security keys for each database can be obtained individually. When all security keys are correct, the plaintexts can be fully retrieved as typically shown in Fig. 4. The security keys are further analyzed and demonstrated by using peak SNR (PSNR) values. Figures 5(a) and 5(b) show the performance of parameters  $w_e$  used in the decryption process, respectively, for the two different databases. In this case, all other security keys are assumed to be correct. When the eavesdropping percent for parameters  $w_e$  is lower than 99.90%, the decoded images cannot visually render useful information about the plaintexts. Moreover, the performance of  $\theta$  on PSNR values of the decrypted images has also been studied as shown in Figs. 5(c) and 5(d), respectively, for the two different databases. In this case, all other security keys are also assumed to be correct. It is demonstrated in Figs. 5(c) and 5(d) that when the eavesdropping percent for  $\theta$  is lower than 99.95%, the decoded images cannot visually render useful information about the plaintexts. For the sake of brevity, the eavesdropping analysis of other security keys is not presented here. These eavesdropping analyses demonstrate that the security of the proposed learning-based optical cryptosystem can be fully guaranteed. Without accurate knowledge about security keys, the plaintexts cannot be effectively extracted. The proposed method uses an optical setup with complex scattering media to experimentally record the ciphertexts, and uses a learning model to generate security keys. Therefore, high security can be achieved.

The higher security can be flexibly achieved by using virtual phase-only masks, which serve as supplementary security keys for the decryption. One virtual phase-only mask is further used here and shown in Fig. 6. The axial distance ( $L_1$ ) between the CCD and virtual phase-only mask is 4.0 cm and that ( $L_2$ ) between virtual phase-only mask and the  $D_2$  plane is 2.0 cm. The ciphertexts recorded by the CCD are back-propagated to the  $D_2$  plane when a virtual phase-only mask is used, and amplitude-only patterns obtained in the  $D_2$  plane rather than the recorded ciphertexts are used as the inputs for the designed

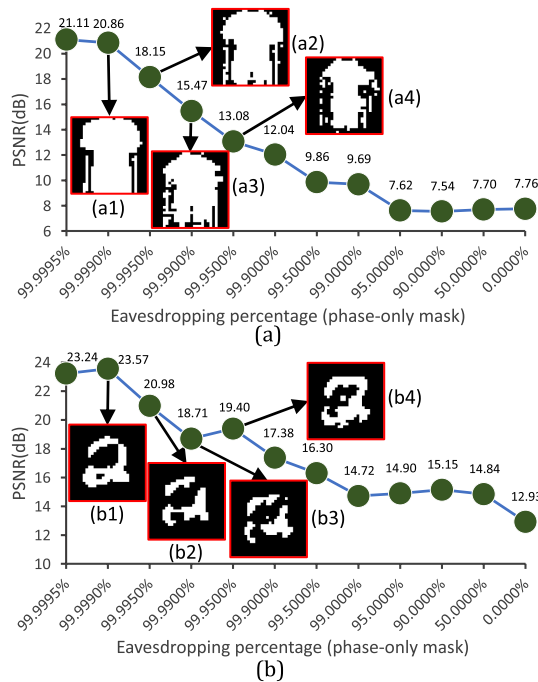


**Fig. 6.** Schematic experimental setup with complex scattering media:  $L_1$  and  $L_2$ , axial distances; V, virtual phase-only mask. A virtual phase-only mask is further used as supplementary security key for the decryption.





**Fig. 7.** Decrypted images obtained by using all correct security keys: (a)–(e) and (k)–(o) ciphertexts; (f)–(j) and (p)–(t) the retrieved plaintexts. The PSNR values for (f)–(j) and (p)–(t) are 30.65 dB, 24.63 dB, 16.09 dB, 29.26 dB, 25.42 dB, 23.08 dB, 34.36 dB, 24.63 dB, 29.59 dB, and 22.04 dB, respectively.



**Fig. 8.** Eavesdropping analysis of virtual phase-only mask. All other security keys are assumed to be correct.

learning model. After the training, the virtual phase-only mask, axial distances, the wavelength, and the trained learning model with its parameters can be used as security keys. The speckle patterns recorded by the CCD are still used as ciphertexts. Figures 7(a)–7(t) show the plaintexts retrieved from their ciphertexts by using all correct security keys. Performance of the virtual phase-only mask in the decryption process is demonstrated in Fig. 8. As shown in Figs. 8(a) and 8(b), when the

eavesdropping percent for the virtual phase-only mask is lower than 99.99%, the decoded images cannot visually render useful information about the plaintexts. For the sake of brevity, the eavesdropping analysis of other security keys is not presented here. It is expected that multiple virtual phase-only masks used for the decryption would further enhance the security.

In conclusion, we have proposed a new method for optical encryption by learning complex scattering media. Instead of directly using the parameters in an optical setup as security keys, the proposed method uses the trained learning model with its parameters as security keys. In addition to the trained model with its parameters, other parameters, e.g., a virtual phase-only mask, can be flexibly used to enlarge the key space. The proposed method uses an optical setup with complex scattering media to experimentally record the ciphertexts and uses a learning model to generate security keys. Therefore, high security is achieved in the proposed method. The proposed method can be theoretically and experimentally implemented, and is not limited to the typical optical setup and the typical number of diffusers presented in this study. The proposed learning-based encryption might open up a new research perspective for optical encryption.

**Funding.** Hong Kong Research Grants Council (25201416, C5011-19G).

**Disclosures.** The authors declare no conflicts of interest.

## REFERENCES

- B. Javidi, *Phys. Today* **50**(3), 27 (1997).
- O. Matoba, T. Nomura, E. Perez-Cabre, M. S. Millan, and B. Javidi, *Proc. IEEE* **97**, 1128 (2009).
- P. Refregier and B. Javidi, *Opt. Lett.* **20**, 767 (1995).
- A. Alfalou and C. Brosseau, *Adv. Opt. Photon.* **1**, 589 (2009).
- W. Chen, B. Javidi, and X. Chen, *Adv. Opt. Photon.* **6**, 120 (2014).
- G. Unnikrishnan, J. Joseph, and K. Singh, *Opt. Lett.* **25**, 887 (2000).
- O. Matoba and B. Javidi, *Opt. Lett.* **24**, 762 (1999).
- W. Chen, X. Chen, and C. J. R. Sheppard, *Opt. Lett.* **35**, 3817 (2010).
- P. Clemente, V. Durán, V. Torres-Company, E. Tajahuerce, and J. Lancis, *Opt. Lett.* **35**, 2391 (2010).
- Y. Zhang and B. Wang, *Opt. Lett.* **33**, 2443 (2008).
- E. G. Johnson and J. D. Brasher, *Opt. Lett.* **21**, 1271 (1996).
- Y. Shi, T. Li, Y. Wang, Q. Gao, S. Zhang, and H. Li, *Opt. Lett.* **38**, 1425 (2013).
- A. Carnicer, M. Montes-Usategui, S. Arcos, and I. Juvells, *Opt. Lett.* **30**, 1644 (2005).
- X. Peng, H. Wei, and P. Zhang, *Opt. Lett.* **31**, 3261 (2006).
- X. Peng, P. Zhang, H. Wei, and B. Yu, *Opt. Lett.* **31**, 1044 (2006).
- M. Liao, W. He, D. Lu, and X. Peng, *Sci. Rep.* **7**, 41789 (2017).
- L. Zhou, Y. Xiao, and W. Chen, *Opt. Express* **27**, 26143 (2019).
- H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms,” arXiv preprint arXiv:1708.07747 (2017).
- L. Deng, *IEEE Signal Process. Mag.* **29**(6), 141 (2012).
- Y. LeCun, Y. Bengio, and G. Hinton, *Nature* **521**, 436 (2015).
- K. Zhang, W. Zuo, S. Gu, and L. Zhang, *IEEE Conference on Computer Vision and Pattern Recognition* (2017), pp. 3929–3938.
- I. Sutskever, J. Martens, G. E. Dahl, and G. E. Hinton, *30th International Conference on Machine Learning (PMLR)* (2013), Vol. **28**, pp. 1139–1147.