

# Exploring explicit coarse-grained structure in artificial neural networks

Xi-Ci Yang,<sup>1</sup> Z. Y. Xie,<sup>2,\*</sup> and Xiao-Tao Yang<sup>1,†</sup>

<sup>1</sup>*College of Power and Energy Engineering, Harbin Engineering University, Harbin 150001, China*

<sup>2</sup>*Department of Physics, Renmin University of China, Beijing 100872, China*

We propose to employ the hierarchical coarse-grained structure in the artificial neural networks explicitly to improve the interpretability without degrading performance. The idea has been applied in two situations. One is a neural network called TaylorNet, which aims to approximate the general mapping from input data to output result in terms of Taylor series directly, without resorting to any magic nonlinear activations. The other is a new setup for data distillation, which can perform multi-level abstraction of the input dataset and generate new data that possesses the relevant features of the original dataset and can be used as references for classification. In both cases, the coarse-grained structure plays an important role in simplifying the network and improving both the interpretability and efficiency. The validity has been demonstrated on MNIST and CIFAR-10 datasets. Further improvement and some open questions related are also discussed.

## I. INTRODUCTION

In the past decade, machine learning has drawn great attention from almost all natural science and engineering communities, such as mathematics [1–3], physics [4–10], biology [11–13], and materials sciences [14–16], and has been widely used in various aspects of modern society, e.g., automatic driving systems, face recognition, fraud detection, expert recommendation system, speech enhancement, and natural language processing, etc. Especially, the deep learning techniques based on the artificial neural networks [17, 18] have become the most popular and dominant machine learning approaches progressively, and their interactions with many-body physics have been intensively explored in recent years. On the one hand, some typical neural networks, such as multi-layer perceptron [18], restricted boltzmann machine [19], autoencoder [20], convolutional neural network [21], and autoregressive network [22], have been successfully applied to the study of quantum magnetization [23–26], Fermi-Dirac statistics [27, 28], superconductivity [29, 30], statistical averages [31] and phase transitions [5, 32, 33] in physical systems. On the other hand, the ideas and techniques developed in physics are introduced into neural networks to improve the performance as well as interpretability [34, 35]. This approach might obtain a deeper insight of the neural networks and is sometimes referred to as the physics-inspired machine learning [36, 37]. A successful example is the introduction of tensor-network state into deep learning. It stems from quantum information and develops fastly in quantum many-body physics, and recently it has been used to realize the supervised learning [38–40], generative models [41–43], and network reconstruction [44–46], etc. Though there are some limitations and difficulties in the current stage, it can still be expected that the interplay between deep learning and many-body physics will continue to flourish in the next

decade.

Among the discussions about machine learning and many-body physics, the connection between deep neural networks and the renormalization group (RG) has been extensively studied in the literature recently [47–53]. This relevance probably stems from the essential similarity between the underlying hierarchical structure of the inference process in supervised learning [18] and the generated coarse-grained structure in the RG flow in physical systems [54], and can be seen more clearly in the context of tensor renormalization group, where the tensor-network structure and the RG-based techniques are combined together to study the many-body physics [55–59]. In fact, it shows that not only the hierarchical structure but also the backpropagation method employed in the training process of neural networks resemble those of the tensor networks very much [60], and this actually lays the foundation of the increasing interplays between the two fields.

In this work, we propose to explicitly employ the hierarchical coarse-grained structure, as generated in the RG process in tensor networks similarly, in neural networks, and apply it to image classification and data distillation [61–63] for better interpretability in both physics and mathematics. To be specific, in the classification task, we construct a neural network called TaylorNet, which expresses the mapping from the input data to the output label in terms of the Taylor series approximately without using any nonlinear activation functions. The network is simple and can be expressed as a polynomial manifestly in mathematics. This is very different from the ordinary neural networks whose explicit expressions are difficult to obtain, thus unveiling part of the mysteries of neural networks and providing clear direction for further improvement. In the second part, we design a multi-level distillation process which imitates the coarse-graining (CG) operations in the RG flow and displays the underlying hierarchical structure of the inference process explicitly. It shows that the data distilled from lower-level abstraction contains much more details than those distilled from higher-level abstraction, and the final data distilled from the highest-level layer can be used as good

\* qingtaoxie@ruc.edu.cn

† yangxiaotao@hrbeu.edu.cn

references for image classification directly. The results obtained in both tasks are rather satisfying, as demonstrated in the MNIST [64] and CIFAR-10 [65] datasets.

The rest of the paper is organized as follows. In Sec. II, we briefly review the coarse-grained structure generated in the RG process in the context of the matrix product operator. In Sec. III and Sec. IV, we introduce the TaylorNet and the new setup for data distillation, respectively, and demonstrate their validity in MNIST and CIFAR-10 datasets. In Sec. V, we summarize our work and discuss the possible improvement as well as promising extensions briefly.

## II. COARSE-GRAINED STRUCTURE GENERATED IN THE RG PROCESS

The RG is one of the most profound tools conceptually of theoretical physics [66–70], and its impact spans from high-energy to statistical and condensed matter physics [66, 71–73]. Essentially, the RG is a conceptual framework comprising various techniques, such as the original block spin approach [58, 67], functional RG [73], Monte Carlo RG [71], density matrix RG [72], and tensor network RG [55–57, 59], and so on. Though these schemes differ substantially in details, they share a same essential feature, namely the RG process aims to identify the relevant degrees of freedom (DOFs), integrate out the irrelevant ones iteratively, and eventually arrive at a low-energy effective theory. The extraction of the relevant information is realized by a set of RG transformations, which map the DOFs in a lower scale to those in its neighboring higher scale. A hierarchical coarse-grained structure is essentially generated during this kind of scale transformation. It can be seen more clearly in the real-space RG schemes [55, 56, 58, 67, 72], as exemplified in the following.

To show the coarse-grained structure mentioned above clearly, let's consider the real-space RG transformations of a matrix product operator [74, 75] defined on a one-dimensional lattice, which may represent a many-body Hamiltonian of a quantum system or a transfer matrix of a classical statistical model. In this context, through a series of scale transformations, the RG process aims to find a finite-dimensional representation of the operator, which can preserve the low-energy part of the Hamiltonian or the dominant-eigenvalue part of the transfer matrix approximately. For simplicity, let's assume the MPO is symmetric, and consider the simple case where a binary mapping is performed in each scale transformation. Then the RG process in such a system with length  $L = 8$  and open boundary condition can be illustrated in Fig. 1. At the beginning, the operator is expressed in terms of the variables  $\{\sigma^{(0)}\}$  sitting on the blue lines. The 1st scale transformation is composed of four isometries denoted as  $U^{(1)}$ s, each of which maps  $\{\sigma^{(0)}\}$  sitting on two neighboring blue lines to the variables  $\{\sigma^{(1)}\}$  sitting on the corresponding green lines. Similarly, the 2nd scale

transformation is composed of two isometries  $U^{(2)}$ s, and each  $U^{(2)}$  maps  $\{\sigma^{(1)}\}$  to variables  $\{\sigma^{(2)}\}$  sitting on the red lines.  $U^{(3)}$  constitutes the last scale transformation, and maps  $\{\sigma^{(2)}\}$  to variables  $\{\sigma^{(3)}\}$  sitting on the black lines. Eventually, the operator is represented in terms of  $\{\sigma^{(3)}\}$ , and this completes the full RG process.

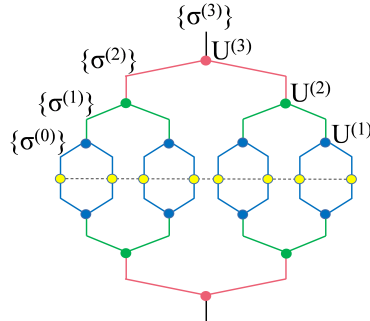


FIG. 1. Hierarchical coarse-grained structure generated in the RG process for a matrix product operator with open boundary condition. The hollow circles connected by a dashed line and colored by yellow denote the lattice sites where the operator is defined. As described in the main text, the local RG transformations are represented by the rank-3 tensors  $U$ s denoted by solid dots, and the DOFs reside on the links between the dots and are denoted as  $\{\sigma\}$ . For the sake of clarity, the various scales are distinguished by different colors.

In the RG transformations described above, by isometry we mean there are fewer DOFs after the mapping, and the generated variables and DOFs are usually termed as coarse-grained. As clearly sketched in Fig. 1, the whole RG process generates a hierarchical coarse-grained structure with three levels. For a given level, the RG transformations identify the relevant DOFs from lower-level DOFs, and output the identified DOFs to a higher-level transformation for further extraction. This CG operation is an essential gradient of the RG process, and has been heavily employed in the original block spin numerical RG calculations [76–78] and the more recent tensor network RG proposals [57, 79, 80].

Without discussing the RG flow in the parameter space and the corresponding fixed-point properties, in this work, we just focus on the hierarchical coarse-grained structure described above and emphasize its similarity to the inference process in supervised learning tasks like image classification. Deep learning falls in the category of representation learning, whose central task is to extract high-level abstract features relevant to the final target from the raw data possessing many irrelevant details and variations [18], and it solves this problem by constructing higher-level representations out of simpler lower-level representations. More specifically, a representation in a given level is characterized by the output of the previous lower-level neural network layer, and is regarded as the input of a new hidden layer to construct the more abstract higher-level representation. The desired representation is eventually obtained by multistep abstraction,

each step of which is realized by a hidden layer and extracts increasingly abstract features from the original input data. This multistep abstraction process is very similar to the RG process discussed before and illustrated in Fig. 1, and also shows an underlying hierarchical coarse-grained structure. This similarity is more evident for the convolutional neural network, where the local structure is emphasized by convolution operations [21].

In the following sections, we introduce this hierarchical coarse-grained structure into neural networks manifestly, which makes the conceptual similarity described here more explicit, and the resulting networks are much more easier to understand in both physics and mathematics.

### III. TAYLORNET

The quintessential example of a deep learning model is the deep feedforward neural network, and sometimes is referred to as multilayer perceptron model [18]. It represents a mapping  $\mathcal{F}$  from the input data to the output result, and generally can be expressed as a composite function of many linear ( $\mathcal{L}$ ) and nonlinear ( $\mathcal{N}$ ) transformations ordered alternately. For example, a feedforward neural network with  $n$  layers can be represented as

$$\mathcal{F} = \mathcal{N}_n \mathcal{L}_n \cdots \mathcal{N}_2 \mathcal{L}_2 \mathcal{N}_1 \mathcal{L}_1 \quad (1)$$

where the linear mappings  $\mathcal{L}$ s contain many variational parameters that need to be determined, while the nonlinear mappings  $\mathcal{N}$ s contain almost no free parameters and are realized by some known operations called activations, such as rectified linear unit, logistic sigmoid, max-poolings, and so on. The nonlinear mappings  $\mathcal{N}$ s are indispensable to approximate a nonlinear  $\mathcal{F}$  [18]. Apparently, Eq. (1) seems oversimplified, but fortunately, when  $\mathcal{F}$  is Borel measurable, the validity is guaranteed by the universal approximation theorem [81, 82], as long as the neural network is sufficiently wide and at least one  $\mathcal{N}$  is squashing in some sense. Therefore, Eq. (1) provides a quite general belief to approximate an actual mapping in practical applications of neural networks. However, for a given mapping  $\mathcal{F}$ , generally, there is no clue on either how wide the neural network is or how we can obtain the desired  $\mathcal{L}$ s. In order to determine the parameters effectively, much effort has been devoted to designing special structures, and this greatly boosted the development of deep neural networks. Successful structures include the convolution operation [21], shortcut connection in residual network [83, 84], attention structure in the transformer model [85], etc. Nevertheless, the specific design of structures relies mainly on empirical experience, and there is no theoretical guidance generally, and this is why deep learning is usually regarded as a magic black box.

To reduce the mystery in the structure design, in this work, we propose to use another universal expansion, i.e., the multi-variable Taylor formula valid for an arbitrary

analytic function  $f$ . Expanded at a certain point, the expression can be written as

$$f(X) = f_0 + \sum_{i=1}^N a_i^{(1)} x_i + \sum_{i,j=1}^N a_{ij}^{(2)} x_i x_j + \sum_{i,j,k=1}^N a_{ijk}^{(3)} x_i x_j x_k + \cdots, \quad (2)$$

where  $X$  is the input vector with  $N$  elements denoted as  $\{x_1, x_2, \dots, x_N\}$ ,  $a^{(n)}$  is the coefficient related to the corresponding  $n$ -th order derivative, and  $f_0$  is a collected constant. Eq. (2) is also universal, since the nonanalytic functions encountered in our daily life are always expected to have a finite number of singular points and thus can be well approximated by Eq. (2) arguably. Hereafter, just for convenience, we simply refer to  $a$  and  $x_i x_j x_k \dots$  as the Taylor coefficient and Taylor term, respectively.

The validity of Eq. (2) can be directly verified by an experiment on the classification of MNIST dataset. In the experiment, we regard each image as a vector  $X$ , and assume  $f^{(\alpha)}(X^{(i)})$  can be expanded as Eq. (2), where  $f^{(\alpha)}(X^{(i)})$  denotes the probability of the  $i$ -th image belongs to the  $\alpha$ -th category. The whole neural network has only a linear layer that holds the coefficients, and the output can be expressed as

$$f^{(\alpha)}(X^{(i)}) = \sum_{j=1}^p W_{\alpha,j} \tilde{X}_j^{(i)}. \quad (3)$$

In the above equation,  $\tilde{X}^{(i)}$  is a vector containing  $p$  individual Taylor terms corresponding to  $X^{(i)}$  which are retained in Eq. (2), and  $W$  is a  $10 \times p$  weight matrix with element  $W_{\alpha,j}$  the Taylor coefficient corresponding to  $\tilde{X}_j^{(i)}$ . To make the calculation feasible, we resize the original  $28 \times 28$  images into  $7 \times 7$  through the well-established bilinear interpolation technique [86], do expansion up to the fourth order, and collect all possible terms in  $\tilde{X}$ .

The result is shown in Fig. 2. It is clear that the test accuracy can be systematically improved as the expansion order  $n$  is increased. When  $n = 4$ , the number of total terms retained is about 293 thousand, and the obtained accuracy is about 98%. As a comparison, on the same  $7 \times 7$  MNIST dataset, a residual network with 1.3 million parameters can achieve an accuracy of about 99%. The performance can be further improved by some detailed analysis. Especially, it shows that, though there are about totally 293 thousand terms retained in the expansion, the contribution of a great number of terms is very small. For example, the distribution of weight corresponding to the quadratic terms is displayed in Fig. 3. For a given term  $x_i x_j$ , Fig. 3(a) tells the dominant weight always consumes the least portion no matter how far  $x_i$  and  $x_j$  are separated, and Fig. 3(b) tells that all the weights have a preferred distribution as a function of the distance between  $x_i$  and  $x_j$ . This reminds us that there is much redundancy in  $W$ , and thus the number of parameters,  $N_0$ , can be greatly reduced by discarding the

small weights. In fact, experiments show that the accuracy remains unchanged when  $N_0$  is reduced by half, and drops only by 0.83 percent when  $N_0$  is reduced to 30%. Even if  $N_0$  is reduced to 10%, we can still obtain an accuracy of about 85%. This actually reflects the spirit of Taylor expansion, since it means the accuracy can be systematically improved by adding more subtle terms.

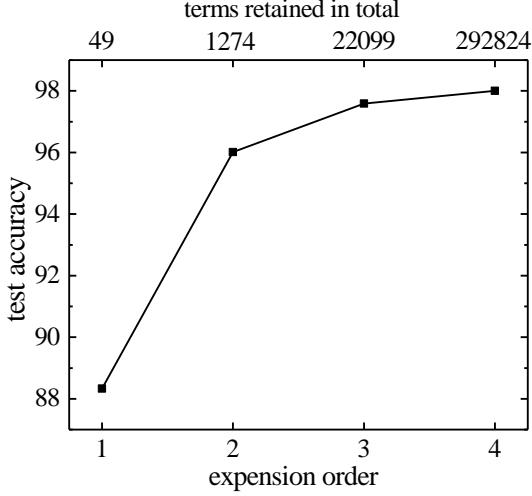


FIG. 2. Test accuracy of the experiment on direct Taylor expansion, as expressed in Eq. (3). The data is obtained on the MNIST dataset with resized  $7 \times 7$  images.

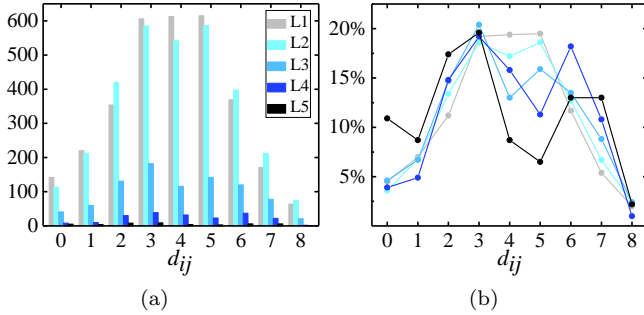


FIG. 3. Distribution of the obtained weight corresponding to the quadratic terms  $x_i x_j$ , as a function of distance  $d_{ij}$  between the two pixels in images of the  $7 \times 7$  MNIST dataset. Weights are ordered by value in ascending order, and equally divided into five groups that are referred to as level-1 (L1) to level-5 (L5), respectively. (a) Weight distribution at each distance  $d_{ij}$ . (b) Weight ratio distribution for each level as a function of  $d_{ij}$ . The details of the distribution can be found in Fig. 13 in App. D.

To extend Eq. (2) to large scale computer tasks, and to make the above procedure more practical and efficient, we propose the TaylorNet, which realizes Eq. (2) in a multistep manner by utilizing the hierarchical coarse-grained structure described in Sec. II. Based on the assumption that the dominant parts in the Taylor series mainly correspond to the product of  $x$ s in local clusters, as has

been partially evidenced in Fig. 3, we introduce a series of intermediate variables  $x_\alpha$  with the new index  $\alpha$  indicating the hierarchical levels. The variables at a higher level are to be expressed in terms of Taylor series with respect to the variables at the neighboring lower level, and constitute the Taylor expansions of the variables at the neighboring higher level. To be specific, if expanded to the second order, a local cluster indexed as  $\alpha$  with  $m$   $n$ -th level variables denoted as  $\{x_1^{(n)}, x_2^{(n)}, \dots, x_m^{(n)}\}$  is mapped to a variable  $x_\alpha^{(n+1)}$  at the  $(n+1)$ -th level, i.e.,

$$x_\alpha^{(n+1)} = c^{(n+1,0)} + \sum_{i=1}^m c_i^{(n+1,1)} x_i^{(n)} + \dots + \sum_{i,j=1}^m c_{ij}^{(n+1,2)} x_i^{(n)} x_j^{(n)} \quad (4)$$

where  $c^{(n,\sigma)}$  denotes the coefficients introduced in the  $\sigma$ -th order terms in the expansion of variables at the  $n$ -th level, and  $x_i^{(0)}$  is defined as the original input data  $x_i$ . Hereafter, Eq. (4) is referred to as a CG operation, and it is illustrated in Fig. 4, in which the usual convolution operation is also illustrated for comparison. Suppose we are considering the simplest case, i.e., the sizes of the local cluster is  $2 \times 2$ , and there is no overlap between the clusters. In the language of neural networks, this means the size of both the kernel and the stride is  $2 \times 2$ . In Fig. 4(a) and (b), the variables at two neighboring levels are denoted as dots and squares, respectively, and the variables associated with a single local mapping are indicated by the same color. In a convolution operation, a square is a linear combination of four dots, which corresponds to the linear terms in Eq. (4). While in the CG operation, a square is a nonlinear combination of the same four dots, which corresponds to Eq. (4) exactly. To indicate the nonlinear feature, an oval plate is added to distinguish from the convolution, as shown in Fig. 4(b).

The local introduction of the nonlinear mapping has a great advantage over the initial proposals of Taylor expansion, as expressed in Eq. (2) and Eq. (3). On the one hand, the locality puts a strong constraint on the distance among the variables showing up in the Taylor terms retained in the expansion. This greatly reduces the number of parameters that need to be determined, and also removes some unnecessary redundancy, since the contribution from the terms involving variables faraway separated is expected to be small, as partially evidenced in Fig. 3(b). On the other hand, the higher-order terms, as well as the terms involving variables belonging to different local clusters, can emerge naturally in the next several CG operations, which can be seen from Fig. 4(c) explicitly. For example, more complex terms like  $x_i^3$ ,  $x_i^4$ ,  $x_1 x_3$ ,  $x_1 x_2 x_3$ ,  $x_1 x_2 x_3 x_4$  show up in the expression of  $z$ , though  $y_1$  and  $y_2$  are only expanded to the second order locally. The full expression of  $z$  can be found in App. A. Thus by introducing the nonlinear terms locally, we can generate long-range and very complicated terms in the actual expansions of the variables at the highest level,



and there is no need to invoke any magic activations at all.

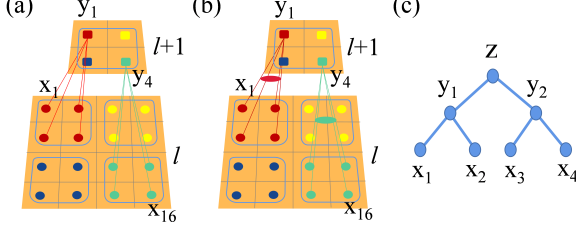


FIG. 4. Illustration of the CG operation in TaylorNet. (a) Convolution operation. (b) CG operation. Clearly, the product of  $x_1$  and  $x_{16}$  will emerge in the next CG after the next CG operation. (c) Two successive CG operations on four variables divided into two local clusters. Very complicated terms of  $x$  emerge in the expression of  $z$ , as shown in Fig. 10 in App. A.

In our experiment on MNIST, we use a TaylorNet with four CG layers, each of which maps a  $2 \times 2$  cluster to a single variable according to second-order Taylor expansion, and then use a linear layer that maps the resulting  $2 \times 2$  variables to a vector with 10 elements representing the probabilities. The detailed TaylorNet structure is shown in Fig. 5. The obtained accuracy is about 99.2%, which is quite satisfying. On the resized  $7 \times 7$  MNIST dataset, we can obtain an accuracy of about 97.9% with about 248 thousand parameters in total, which is much less than the parameters in both the original Taylor expansions (1.46 million) and the residual network (1.3 million) described before. As to the CIFAR-10 dataset, we can obtain an accuracy of about 71.7% with only 1.2 million parameters, and this is also more efficient than the recent MLP-Mixer proposal [87] without pre-training process, which combines the information from the inter-cluster and intra-cluster variables in a similar way. The detailed TaylorNet structure for CIFAR-10 dataset can be found in Fig. 12 in App. C. Furthermore, it shows that if we replace the CG operations with the convolution operations in the whole neural network, the accuracy will drop by about 7.3% and 30% immediately as expected, on MNIST and CIFAR-10 datasets, respectively. This clearly demonstrates the power of CG operations in the representations of nonlinear mappings.

Similar to the convolution operation, the above CG operation can be performed in slightly different manners. Firstly, the size of the local clusters can be different, and there can be overlaps between different clusters. Moreover, the translation and/or scale invariance of the CG kernels can be employed, that is, the Taylor coefficients in Eq. (4) for CG operations performed at different clusters and/or in different scales can be assumed identical. Secondly, the expansion order in Eq. (4) can be larger than 2, and this depends on the strength of the nonlinearity in the expansion of the actual  $\mathcal{F}$ .

It is worth mentioning that, in Fig. 5, we have em-

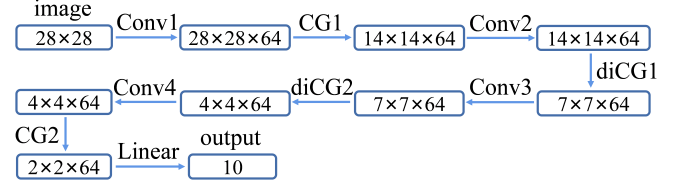


FIG. 5. Sketch of the TaylorNet used in the classification task on MNIST dataset. Hereafter, the numbers in the box denote the representation form of the data, e.g.,  $28 \times 28 \times 64$  denotes 64 feature maps with size  $28 \times 28$ , and the operations sit on the arrows correspond to different neural network layers, e.g.,  $\text{Conv}(l_1, l_2, c, s_1, s_2, p)$  means convolutional layer with kernel size  $l_1 \times l_2$ , channels  $c$ , stride size  $s_1 \times s_2$ , and padding number  $p$  (default 0), similar for CG operation and dilated CG operation as discussed in the main text and Fig. 6. Here, all the four convolutional layers have structure  $\text{Conv}(3, 3, 64, 1, 1, 1)$ , both CG layers have structure  $\text{CG}(2, 2, 64, 2, 2)$ ,  $\text{diCG1}$  and  $\text{diCG2}$  have structures  $\text{CG}(2, 2, 64, 1, 1, 7)$  and  $\text{CG}(2, 2, 64, 1, 1, 3)$ , respectively. The action of a multi-channel convolution operation is illustrated in Fig. 11 in App. B, and more details can also be found in Ref. [18].

ployed two simple ways to further enlarge the effective receptive fields [88] of the CG operations, without changing the size of the local clusters manifestly. One is the introduction of the dilated CG operation, as shown in Fig. 6. In the dilated CG, the variables need to be coarse grained scatter separately in different clusters instead of aggregating locally, which is very similar to the structure of the dilated convolution [89]. The other is performing convolution before the CG operation. This is easy to understand, since the convolution turns each dot in Fig. 4(b) as the linear combination of several dots nearby before the CG operation, and thus each square is effectively expressed in terms of more dots actually. These two techniques might be advantageous in some situations where the inter-cluster product is more important in Eq. (3), and are also considered in Ref. [87] similarly.

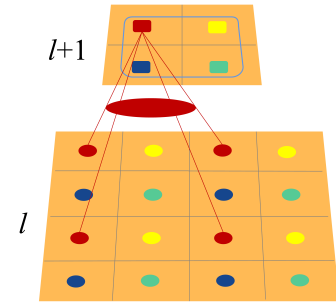


FIG. 6. Sketch of a dilated CG operation, used in Fig. 5. The structure shown in this figure is denoted as  $\text{CG}(2, 2, n, 1, 1, 2)$  which means the kernel size is  $2 \times 2$ , number of channels is  $n$ , stride size is  $1 \times 1$ , and the dilation is 2 in both directions.

#### IV. DATA DISTILLATION

The concept of knowledge distillation was originally proposed by Hinton *et al.* [61], and it aims to train a simpler neural network called student, which is expected to have the same performance approximately to a complex model referred to as teacher. Later, data distillation is proposed to train a smaller dataset from a larger dataset, and expect that the obtained distilled dataset can be used to efficiently train a neural network that has a similar performance to that of a neural network trained from the original larger dataset [62, 63]. Though the process of distillation is somewhat complicated, the idea is very simple and reasonable, namely the neural-network-based deep learning is believed to be able to extract some essential features from the original dataset.

In order to make the above idea clearer, and display the inference or abstraction process more explicitly, we propose a new setup of data distillation. Utilizing the hierarchical coarse-grained structure, the new proposal aims to extract the essential features through a multistep process, in which the abstraction is performed progressively from lower levels to higher levels. This is actually the essential spirit of deep learning [18], as discussed in Sec. II. It shows that the distilled dataset can be indeed used as references to perform classification task directly.

For concreteness, in the following we describe the distillation strategy applied to the MNIST dataset. The original dataset is composed of ten classes, each of which contains 6000 images, and is denoted as  $D^{(0)}(10, 6000)$  hereafter. Firstly we divide each class into 600 groups equally, and select one group from each class to constitute a subset which contains 10 images for each class. Thus totally we obtain 600 subsets, and for simplicity, the  $i$ -th subset is denoted as  $D_i^{(0)}(10, 10)$  and has 100 images in total. Then perform the usual distillation process on each subset  $D_i^{(0)}$  by neural networks, as will be described later, and distill 10 images corresponding to the 10 classes out of each subset. This completes the first level distillation procedure, from which 600 images for each class are distilled, and we denote the distilled dataset as  $D^{(1)}(10, 600)$  as a whole. Similarly, we further divide  $D^{(1)}(10, 600)$  into 100 subsets  $D_i^{(1)}(10, 6)$  on each of which the distillation process is performed and 10 distilled images are obtained, and then we obtain the distilled dataset  $D^{(2)}(10, 100)$  which contains 100 images for each class at the second level. Repeat this divide-and-conquer strategy, we can obtain dataset  $D^{(3)}(10, 20)$ ,  $D^{(4)}(10, 4)$ , and finally reach the highest distilled dataset  $D^{(5)}(10, 1)$ , which contains only a single distilled image for each class and can be regarded as the typical representatives hosting the essential features of that class. The whole process is illustrated in Fig. 7.

In this work, to perform the distillation procedure on each subsets  $D_i^{(\alpha)}(10, m)$ , we employ the distribution matching method and a similar neural network architecture proposed in Ref. [63], as is illustrated in Fig. 8

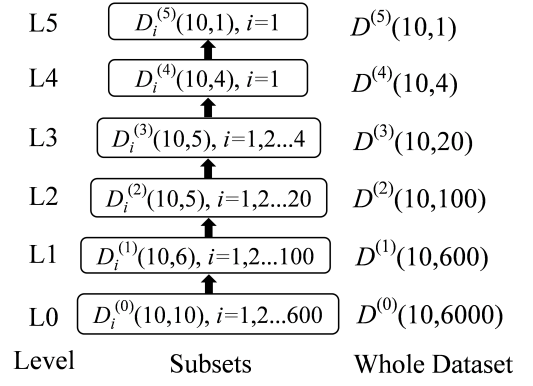


FIG. 7. The data distillation strategy described in the main text, designed for MNIST dataset. The whole distillation process has a five-level structure. The correspondence in the five distillation procedures can be represented as 10-to-1, 6-to-1, 5-to-1, 5-to-1, and 4-to-1 mappings, respectively.

in detail. The goal of the procedure is to determine  $n$  (distilled) images, denoted as  $Y$ s, which minimize a lost function defined in the following

$$L = \sum_{\alpha=1}^n \left( \lambda d_{\alpha, \alpha} - \sum_{\beta \neq \alpha}^n d_{\alpha, \beta} \right) \quad (5)$$

$$d_{\alpha, \beta} \equiv \sum_{i=1}^{m_{\beta}} |f(Y_{\alpha}) - f(X_{\beta, i})|^2$$

where  $Y_{\alpha}$  denotes the desired image for the  $\alpha$ -th class,  $X_{\alpha, i}$  denotes the  $i$ -th image in the  $\alpha$ -th class of the dataset,  $f$  denotes the nonlinear mapping represented by the neural network which produces the embedding vector for any given image, as illustrated in Fig. 8. In Eq. (5),  $n$  is the number of classes in the dataset,  $m_{\alpha}$  is the number of images belonging to the  $\alpha$ -th class, and  $\lambda$  is a hyperparameter to balance the two terms in the bracket, which is set to be 19 in our calculations. In essence,  $d_{\alpha, \beta}$  measures the Euclidean distance between the reference  $Y_{\alpha}$  and the images belonging to the  $\beta$ -th class of the original dataset, in the space where the embedding vectors are defined. Therefore, physically the lost function means that each desired reference is required to resemble the images in the same class to the greatest extent, and at the meanwhile differ from the images in the other classes as much as possible.

The distilled images at different levels are sketched in Fig. 9(a). As expected, it seems that the obtained images obtained from lower-level distillations contain more details and are clearer. When the level of distillation goes up, the details gradually blur and only some indescribable features remain. This tendency is more evident for CIFAR-10 dataset, as is shown in Fig. 9(b).

It is reasonable to regard the remaining features in the final distilled images as essential ones which characterize, or even define the dataset in the perfect case. To check

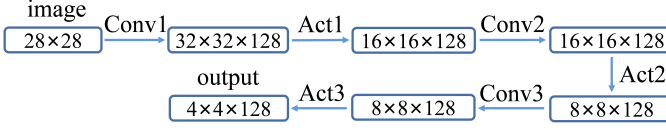


FIG. 8. The neural network structure used in the distillation procedure on each subset of MNIST. Conv1 has structure Conv(3,3,128,1,1,3), both Conv2 and Conv3 have structure Conv(3,3,128,1,1,1). Act1, Act2, Act3 are three activation layers, each of which is composed of instance normalization, relu, and avgpooling. The avgpooling is performed with kernel size  $2 \times 2$  and stride  $2 \times 2$ . Each image is mapped to an embedded space, and the result is a vector with length 16 and channels 128.

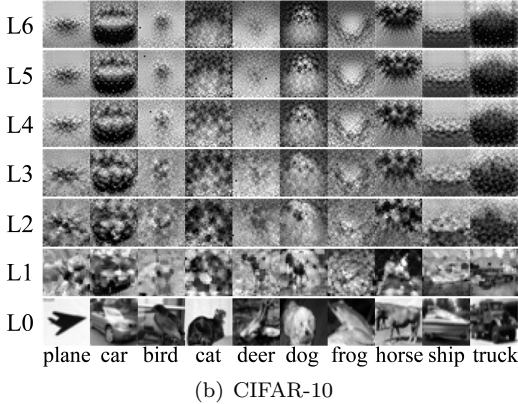
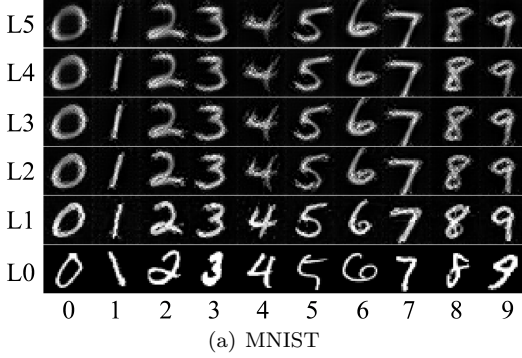


FIG. 9. Distilled figures at each abstraction level. For comparison, one sample of each class in the original dataset is shown in  $L_0$ .

this, firstly we train two residual networks, i.e., ResNet18 and ResNet50 on MNIST and CIFAR-10 datasets, respectively, through the usual classification task, and then use the trained models to collect the output embedding vectors [18] of both the test and distilled images. Without resorting to neural networks further, the final classification is performed by directly comparing the similarity between the embedding vectors of an test image and that of the distilled images, according to the angles in between, and the image is classified into a category whose distilled reference image has the highest similarity to it. It shows that this direct comparison can al-

ready give test accuracies of about 98.7% and 86.9% for MNIST and CIFAR-10, respectively. This confirms that the above distillation process can indeed capture some essential features of the original dataset, from lower-level abstraction to higher-level abstraction gradually, and this reflects exactly the spirit of deep learning.

The performance can be further improved in several ways. Firstly, in each distillation procedure, the lost function plays an important role and can be designed more smartly. In this work, the distance between two images is defined as the Euclidean distance in the embedded space, and one can use other measures, such as the Arcface loss [90], which emphasizes the angular separation and is frequently used in facial-recognition tasks. A preliminary experiment utilizing this lost function can produce an accuracy of about 99% for MNIST, and 89% for CIFAR-10. Secondly, the partition of the dataset, as well as the choice of the hyperparameter  $\lambda$ , might affect the result of the whole distillation. An optimal choice should consider the balance between performance and efficiency in a better way, while in this work, we just adopt the most convenient choice.

## V. SUMMARY

To summarize, inspired by the similarity between the RG flow in physical systems and the inference process in deep neural networks, we introduce the hierarchical coarse-grained structure into the artificial neural networks manifestly to improve the interpretability without degrading performance. To be specific, in the first part, we propose the TaylorNet by introducing the CG operation locally and hierarchically, which extends the linear convolution operation to nonlinear polynomial combinations. It approximates the mapping from the input signal to output result by Taylor expansions effectively, without resorting to the activation functions, and achieves satisfying results in the classification experiments on MNIST and CIFAR-10 datasets. In the distillation task, we propose a setup with a hierarchical coarse-grained structure, and make the inference process from lower levels to higher levels more transparent. It seems that the multistep distillation process is able to capture some essential features of the original dataset, and the distilled images possess less irrelevant details and can be used as reference images in classification tasks. In both cases, the resulting processes represented by the neural networks are more understandable, and the performance are very acceptable compared to the conventional neural networks.

Besides the specific issues discussed separately in Sec. III and Sec. IV, there are some other aspects that can be explored to further improve the performance. For example, we can use more than one TaylorNets to approximate a single mapping, and put the orthogonality constraint on these networks appropriately for higher efficiency. This might be achieved by adding penalties in the lost functions, training in momentum space, or using

other orthogonal complete sets like spherical functions. All these topics are interesting and have been discussed in physical systems, but are far beyond the scope of this paper, and we would like to leave them as pursuits in the near future.

As to the TaylorNet, it is also worth mentioning that our proposal is very different from the previous work in the literature [91–95], which have also explored the Taylor series in neural networks. Most of them have specific motivations and work in different frameworks, and there is no explicit hierarchical structure employed there. For example, Chen *et al.* [91] used a single-layer neural network similar to Eq. (3) to approximate the Taylor expansion of a single-variable function. Montavon *et al.* [92] explored the role of Taylor coefficients as derivatives in an ordinary neural network to analyze the importance of a single pixel in the classification task. Tong *et al.* [94] used the Taylor series to approximate the quadratic form of a Hermitian matrix. Rao *et al.* [95] expressed part of the nonlinearity in terms of a direct-product operation and applied it to the study of partial differential equations. The closest one to our work is probably Ref. [93], where Novikov *et al.* used the neural network to approximate the Taylor expansions; however, in their work both the Taylor terms and Taylor coefficients are represented as compact matrix product operators approximately; thus there is no coarse-grained structure emphasized in this work at all.

At last, though the fixed-point is not discussed at all in this work, the introduction of the hierarchical coarse-grained structure does provide the possibility of its exis-

tence. It is possible to study the fixed point of the scale transformations introduced in the TaylorNet, as well as the fixed point of the iterative distillation procedure introduced in Sec. IV. Whether the scale invariance can be related to some interesting critical phenomena in this context, as explored in Ref. [47], is an open question deserving attention.

## ACKNOWLEDGEMENT

We thank Jing Zhang for her contribution in the early stage of this work, and thank Tao Xiang and Ze-Feng Gao for helpful discussions. We are supported by the National R&D Program of China (Grants No. 2017YFA0302900 and No. 2016YFA0300503), the National Natural Science Foundation of China (Grants No. 52176064, 12274458 and No. 11774420), and by the Research Funds of Renmin University of China (Grants No. 20XNLG19).

## Appendix A: Complete expression of $z$

As illustrated in Fig. 4(c), if each local mapping is expanded to the second order, then the complete expression of  $z$  in terms of  $x_i$ , with  $i = 1, 2, 3, 4$ , can be written as

$$z = \mathbf{a} \mathbf{W} \mathbf{b}^T, \quad (\text{A1})$$

where the two vectors  $\mathbf{a}$  and  $\mathbf{b}$  are defined as

$$\begin{aligned} \mathbf{a} &= [1 \ x_1 \ x_1^2 \ x_1^3 \ x_1^4 \ x_2 \ x_2^2 \ x_2^3 \ x_2^4 \ x_1 x_2 \ x_1 x_2^2 \ x_1^2 x_2 \ x_1^2 x_2^2 \ x_1 x_2^3 \ x_1^3 x_2] \\ \mathbf{b} &= [1 \ x_3 \ x_3^2 \ x_3^3 \ x_3^4 \ x_4 \ x_4^2 \ x_4^3 \ x_4^4 \ x_3 x_4 \ x_3 x_4^2 \ x_3^2 x_4 \ x_3^2 x_4^2 \ x_3 x_4^3 \ x_3^3 x_4] . \end{aligned} \quad (\text{A2})$$

The weight  $W$  can be represented as a 15×15 matrix, whose nonzero elements are denoted as 1 in Fig. 10, just for clarity.

## Appendix B: The action of multi-channel convolutions

For concreteness, the convolutional layer with four three-channel kernels will turn the three-channel input data into four-channel output, as illustrated in Fig. 11, where the + sign means equal-weight superposition of the three dot-product results.

## Appendix C: TaylorNet structure used in the classification on CIFAR-10 dataset

The detailed structure of the TaylorNet used in Sec. III in the classification on CIFAR-10 dataset, is shown in Fig. 12.

## Appendix D: Weight details of the quadratic terms in the experiment on MNIST dataset

As to the direct experiment on the Taylor expansion, Eq. (3), on MNIST dataset in Sec. III, the trained weight of the quadratic terms are sketched in detail in Fig. 13. The statistics are shown in Fig. 3.

[1] N. Lei, Z. Luo, S. T. Yau, and D. X. Gu, arXiv:1805.10451.

[2] N. Lei, K. Su, S. T. Yau, and D. X. Gu, Comput. Aided Geometric Design **68**, 1 (2019).



$$\begin{array}{c}
1 \\
x_3 \\
x_3^2 \\
x_3^3 \\
x_3^4 \\
x_4 \\
x_4^2 \\
x_4^3 \\
x_4^4 \\
x_3x_4 \\
x_3x_4^2 \\
x_3^2x_4 \\
x_3^2x_4^2 \\
x_3x_4^3 \\
x_3^3x_4
\end{array}
\begin{pmatrix}
1 & x_1 & x_1^2 & x_1^3 & x_1^4 & x_2 & x_2^2 & x_2^3 & x_2^4 & x_1x_2 & x_1x_2^2 & x_1^2x_2 & x_1^2x_2^2 & x_1x_2^3 & x_1^3x_2 \\
\boxed{1} & \boxed{1} & 1 & 1 & 1 & \boxed{1} & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
\boxed{1} & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & & & & & \\
1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & & & & & \\
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & & & & & \\
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & & & & & \\
\boxed{1} & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & & & & & \\
1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & & & & & \\
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & & & 0 & & \\
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & & & & & \\
1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & & & & & \\
1 & & & & & & & & & & & & & & \\
1 & & & & & & & & & & & & & & \\
1 & & & & & & & & & & & & & & \\
1 & & & & & & & & & & & & & & \\
1 & & & & & & & & & & & & & & \\
1 & & & & & & & & & & & & & &
\end{pmatrix}$$

FIG. 10. Weight  $W$  appearing in Eq. (A1), corresponding to the coefficients of  $z$  illustrated in Fig. 4(c). For simplicity, the nonzero elements are denoted as 1. The terms in the red box are the ones showing up in the convolution operations.

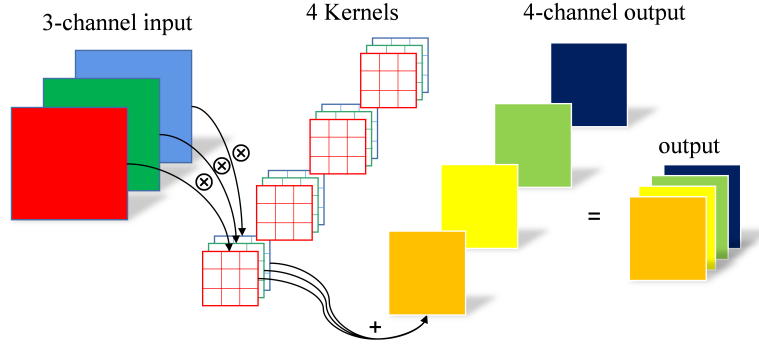


FIG. 11. The action of a convolutional layer with four three-channel kernels, as mentioned in Sec. III.

- [3] A. Fawzi, M. Balog, A. Huang, T. Hubert, B. R. Paredes, M. Barekatin, A. Novikov, F. J. R. Ruiz, J. Schrittwieser, G. Swirszcz, D. Silver, D. Hassabis, and P. Kohli, *Nature* **610**, 47 (2022).
- [4] X. Gao and L. M. Duan, *Nat. Commun.* **8**, 662 (2017).
- [5] J. Carrasquilla and R. G. Melko, *Nat. Phys.* **13**, 431 (2017).
- [6] D. Wu, L. Wang, and P. Zhang, *Phys. Rev. Lett.* **122**, 080602 (2019).
- [7] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, *Rev. Mod. Phys.* **91**, 045002 (2019).
- [8] E. Bedolla, L. C. Padierna, and R. C. Priego, *J. Phys.: Condens. Matter* **33**, 053001 (2021).
- [9] B. Font, G. Weymouth, V. T. Nguyen, and O. R. Tutty, *J. Comput. Phys.* **434**, 110199 (2021).
- [10] D. D. Sante, M. Medvidovic, A. Toschi, G. Sangiovanni, C. Franchini, A. M. Sengupta, and A. J. Millis, *Phys. Rev. Lett.* **129**, 136402 (2022).
- [11] M. Helmstaedter, K. L. Briggman, S. C. Turaga, V. Jain, H. S. Seung, and W. DenkHelmstaedter, *Nature* **500**, 168 (2013).
- [12] S. Webb, *Nature* **554**, 555 (2018).
- [13] J. Dauparas, I. Anishchenko, N. Bennett, H. Bail, R. J. Ragotte, L. F. Milles, B. I. M. Wicky, A. Courbet, R. J. deHaas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan, B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N.P.King, and D. Baker, *Science* **378**, 49 (2022).
- [14] V. Stanev, C. Oses, A. G. Kusne, E. Rodriguez, J. Paglione, S. Curtarolo, and I. Takeuchi, *npj Comput. Materials* **4**, 29 (2018).
- [15] A. Y. Wang, R. J. Murdock, S. K. Kauwe, A. O. Oliynyk, A. Gurlo, J. Brgoch, K. A. Persson, and T. D. Sparks, *Chem. Mater.* **32**, 4954 (2020).

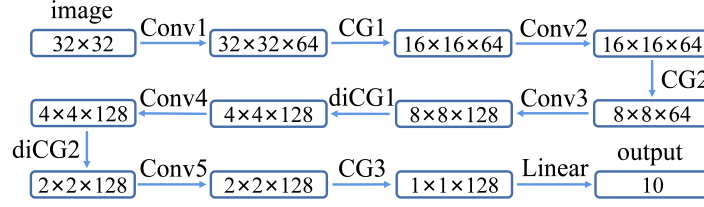


FIG. 12. The TaylorNet used in the classification on CIFAR-10 dataset, as mentioned in Sec. III. Convolutions Conv1 and Conv2 have structure Conv(3,3,64,1,1,1), Conv3, Conv4, and Conv5 have structure Conv(3,3,128,1,1,1). CG operations CG1 and CG2 have structure CG(2,2,64,2,2), CG3 have structure CG(2,2,128,1,1). Dilated CG operations diCG1 and diCG2 have structures diCG(2,2,128,1,1,4) and diCG(2,2,128,1,1,2), respectively.

- [16] R. Batra, L. Song, and R. Ramprasad, Nat. Rev. Mater. **6**, 655 (2021).
- [17] Y. LeCun, Y. Bengio, and G. Hinton, Nature **521**, 436 (2015).
- [18] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning* (MIT press, 2016).
- [19] P. Smolensky, *Information processing in dynamical systems: Foundations of harmony theory*, in *Parallel Distributed Processing* (MIT Press, 1986).
- [20] H. Bourslard and Y. Kamp, Biol. Cybern. **59**, 291 (1988).
- [21] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, Proc. IEEE **86**, 2278 (1998).
- [22] Y. Bengio, and S. Bengio, Adv. NIPS **12**, 400 (2000).
- [23] G. Carleo, and M. Troyer, Science **355**, 602 (2017).
- [24] Z. Cai, and J. Liu, Phys. Rev. B **97**, 035116 (2018).
- [25] X. Liang, W. Y. Liu, P. Z. Lin, G. C. Guo, Y. S. Zhang, L. He, Phys. Rev. B **98**, 104426 (2018).
- [26] Y. Nomura, and M. Imada, Phys. Rev. X **11**, 031034 (2021).
- [27] D. Pfau, J. S. Spencer, and A. G. D. G. Matthews, Phys. Rev. Research **2**, 033429 (2020).
- [28] J. Hermann, Z. Schatzle, and F. Noe, Nat. Chem. **12**, 891 (2020).
- [29] S. Zeng, Y. Zhao, G. Li, R. Wang, X. Wang, and J. Ni, npj Comput. Materials **5**, 84 (2019).
- [30] T. Konno, H. Kurokawa, F. Nabeshima, Y. Sakishita, R. Ogawa, I. Hosako, and A. Maeda, Phys. Rev. B **103**, 014509 (2021).
- [31] D. Wu, L. Wang, and P. Zhang, Phys. Rev. Lett. **122**, 080602 (2019).
- [32] E. P. L. van Nieuwenburg, Y. H. Liu, and S. D. Huber, Nat. Phys. **13**, 435 (2017).
- [33] S. J. Wetzel, Phys. Rev. E **96**, 022140 (2017).
- [34] L. G. Valiant, Commun. ACM **27**, 1134 (1984).
- [35] Y. Bahri, J. Kadmon, J. Pennington, S. S. Schoenholz, J. S. Dickstein, and S. Ganguli, Annu. Rev. Condens. Matter Phys. **11**, 501 (2020).
- [36] Nat. Mach. Intell. **3**, 925 (2021).
- [37] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and Liu Yang, Nat. Rev. Phys. **3**, 422 (2021).
- [38] E. Stoudenmire, and D. J. Schwab, Adv. NIPS **29**, 4799 (2016).
- [39] I. Glasser, N. Pancotti, and J. I. Cirac, IEEE Access **8**, 68619 (2020).
- [40] S. Cheng, L. Wang, and P. Zhang, Phys. Rev. B **103**, 125117 (2021).
- [41] Z. Y. Han, J. Wang, H. Fan, L. Wang, and P. Zhang, Phys. Rev. X **8**, 031012 (2018).
- [42] S. Cheng, L. Wang, T. Xiang, and P. Zhang, Phys. Rev. B **99**, 155131 (2019).
- [43] T. Vieijra, L. Vanderstraeten, and F. Verstraete, arXiv:2202.08177.
- [44] D. Liu, S. J. Ran, P. Wittek, C. Peng, R. B. Garcia, G. Su, and M. Lewenstein, New J. Phys. **21**, 073059 (2019).
- [45] Z. F. Gao, S. Cheng, R. Q. He, Z. Y. Xie, H. H. Zhao, Z. Y. Lu, and T. Xiang, Phys. Rev. Research **2**, 023300 (2020).
- [46] B. Zunkovic, Quantum Mach. Intell. **4**, 21 (2022).
- [47] S. Saremi, and T. J. Sejnowski, PNAS **110**, 3071 (2013).
- [48] C. Beny, arXiv:1301.3124.
- [49] H. W. Lin, M. Tegmark, D. Rolnick, J. Stat. Phys. **168**, 1223 (2017).
- [50] P. Mehta, D. J. Schwab, arXiv:1410.3831; arXiv:1609.0354.
- [51] M. K. Janusz, and Z. Ringel, Nat. Phys. **14**, 578 (2018).
- [52] S. H. Li, L. Wang, Phys. Rev. Lett. **121**, 260601 (2018).
- [53] L. De M. Koch, R. De M. Koch, and L. Cheng, IEEE Access **8**, 106487 (2020).
- [54] J. Cardy, *Scaling and Renormalization in Statistical Physics* (Cambridge University Press, 1996); M. Kardar, *Statistical Physics of Fields* (Cambridge University Press, 2007).
- [55] M. Levin, and C. P. Nave, Phys. Rev. Lett. **99**, 120601 (2007).
- [56] Z. Y. Xie, H. C. Jiang, Q. N. Chen, Z. Y. Weng, and T. Xiang, Phys. Rev. Lett. **103**, 160601 (2009); H. H. Zhao, Z. Y. Xie, Q. N. Chen, Z. C. Wei, J. W. Cai, and T. Xiang, Phys. Rev. B **81**, 174411 (2010).
- [57] Z. Y. Xie, J. Chen, M. P. Qin, J. W. Zhu, L. P. Yang, and T. Xiang, Phys. Rev. B **86**, 045139 (2012).
- [58] E. Efrati, Z. Wang, A. Kolan, L. P. Kadanoff, Rev. Mod. Phys. **86**, 647 (2014).
- [59] Y. Meurice, R. Sakai, J. U. Yockey, Rev. Mod. Phys. **94**, 025005 (2022).
- [60] B. B. Chen, Y. Gao, Y. B. Guo, Y. Z. Liu, H. H. Zhao, H. J. Liao, L. Wang, T. Xiang, W. Li, and Z. Y. Xie, Phys. Rev. B **101**, 220409(R) (2020).
- [61] G. Hinton, O. Vinyals, and J. Dean, arXiv:1503.02531.
- [62] T. Wang, J. Y. Zhu, A. Torralba, and A. A. Efros, arXiv:1811.10959.
- [63] B. Zhao, K. R. Mopuri, and H. Bilen, arXiv:2006.05929, ICLR (2021); B. Zhao, H. Bilen, arXiv:2110.04181, ICML (2021).
- [64] The official website of MNIST is available at <http://yann.lecun.com/exdb/mnist>.

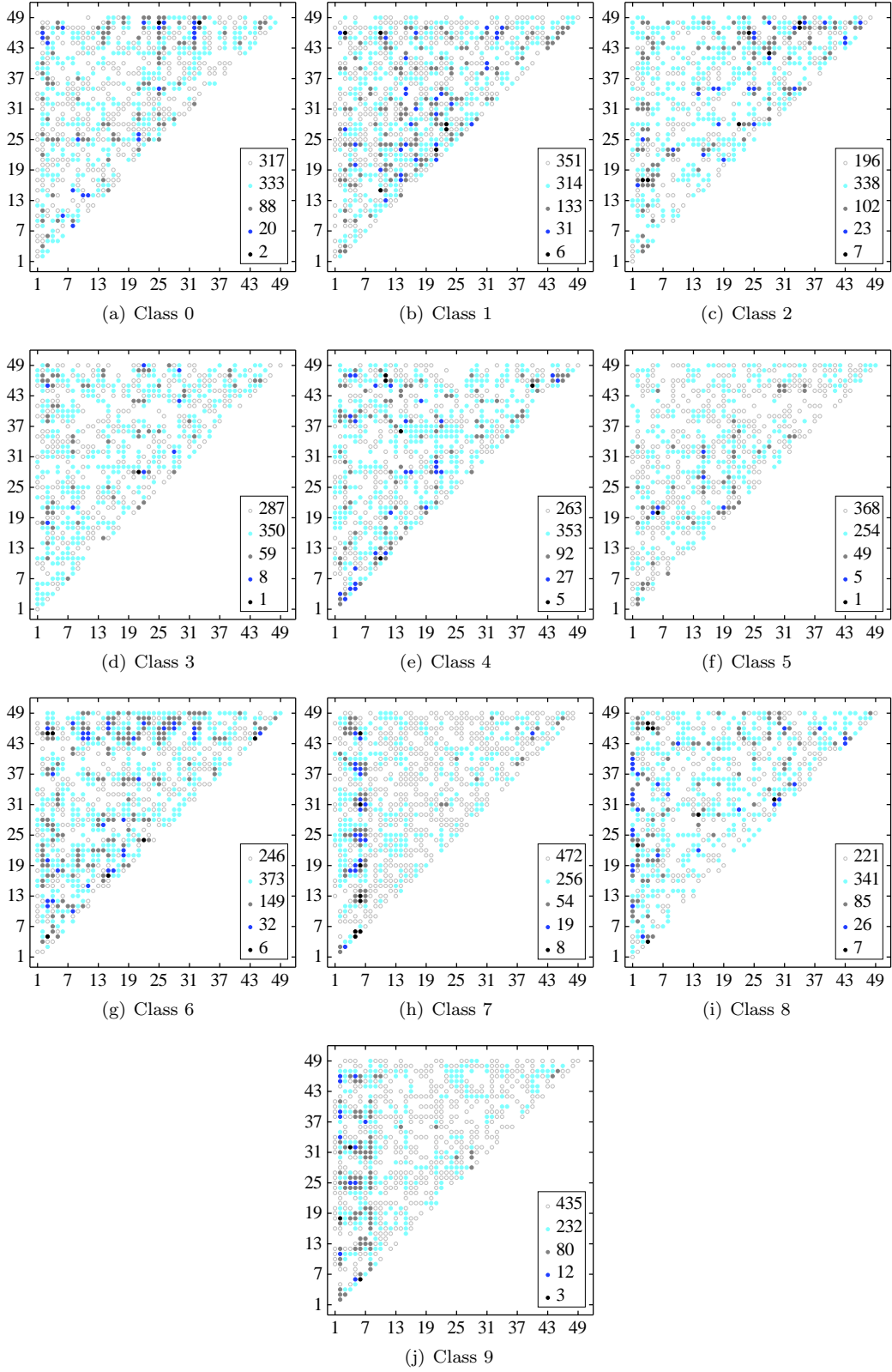


FIG. 13. The weight corresponding to quadratic terms  $x_i x_j$  for each class, for experiment of Eq. (3) on the  $7 \times 7$  MNIST dataset. The abscissa and ordinate represent the linear coordinates  $i$  and  $j$  in the input image, respectively. The weight is divided into five levels according to the magnitude of the value, with darker colors representing larger weights. The inset counts the number of weights belonging to different levels.

- [65] The official website of CIFAR-10 is available at <https://www.cs.toronto.edu/~kriz/cifar.html>.
- [66] M. Gell-Mann and F. E. Low, Phys. Rev. **95**, 1300 (1954).
- [67] L. P. Kadanoff, Physics 2, **263** (1966).
- [68] K. G. Wilson, Phys. Rev. B **4**, 3174 (1971); Phys. Rev. B **4**, 3184 (1971).
- [69] K. G. Wilson, Rev. Mod. Phys. **47**, 773 (1975).
- [70] L. P. Kadanoff, Phys. Rev. Lett. **34**, 1005 (1975).
- [71] R. H. Swendsen, Phys. Rev. Lett. **42**, 859 (1979).
- [72] S. R. White, Phys. Rev. Lett. **69**, 2863 (1992).
- [73] C. Wetterich, Phys. Lett. B **301**, 90 (1993).
- [74] F. Verstraete, J. J. García-Ripoll, and J. I. Cirac, Phys. Rev. Lett. **93**, 207204 (2004).
- [75] B. Pirvu, V. Murg, J. I. Cirac and F. Verstraete, New J. of Phys. **12**, 025012 (2010).
- [76] J. W. Bray and S. T. Chui, Phys. Rev. B **19**, 4876 (1979).
- [77] C. Y. Pan and X. Y. Chen, Phys. Rev. B **36**, 8600 (1987).
- [78] M. D. Kovarik, Phys. Rev. B **41**, 6889 (1990).
- [79] L. P. Yang, Y. Z. Liu, H. Y. Zou, Z. Y. Xie, and Y. Meurice, Phys. Rev. E **93**, 012138 (2016).
- [80] B. B. Chen, L. Chen, Z. Chen, W. Li, and A. Weichselbaum, Phys. Rev. X **8**, 031082 (2018).
- [81] K. Hornik, M. Stinchcombe, and H. White, Neural Networks **2**, 359 (1989).
- [82] G. Cybenko, Math. Control Signals Sys. **2**, 303 (1989).
- [83] K. M. He, X. Y. Zhang, S. Q. Ren, and J. Sun, Deep residual learning for image recognition, CVPR (2017), pp. 4700–4708, see [http://openaccess.thecvf.com/content\\_cvpr\\_2016/html/He\\_Deep\\_Residual\\_Learning\\_CVPR\\_2016\\_paper.html](http://openaccess.thecvf.com/content_cvpr_2016/html/He_Deep_Residual_Learning_CVPR_2016_paper.html).
- [84] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, Densely connected convolutional networks, CVPR (2017), see [http://openaccess.thecvf.com/content\\_cvpr\\_2017/html/Huang\\_Densely\\_Connected\\_Convolutional\\_CVPR\\_2017\\_paper.html](http://openaccess.thecvf.com/content_cvpr_2017/html/Huang_Densely_Connected_Convolutional_CVPR_2017_paper.html).
- [85] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, Adv. NIPS **31** (2017).
- [86] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical recipes: The art of scientific computing* (Cambridge University Press, 2007).
- [87] I. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit, M. Lucic, and A. Dosovitskiy, arXiv:2105.01601.
- [88] W. Luo, Y. Li, R. Urtasun, and R. Zemel, arXiv:1701.04128.
- [89] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, G. Cottrell, arxiv:1702.08502.
- [90] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, arXiv:1801.07698, CVPR 4685 (2019).
- [91] X. Chen, Q. Ma, and T. Alkharobi, IEEE **2**, 291 (2009).
- [92] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, and K. R. Muller, Pattern Recognit. **65**, 211 (2017).
- [93] A. Novikov, M. Trofimov, and I. Oseledets, arXiv:1605.03795, ICLR (2017).
- [94] Y. Tong, S. Xiong, X. He, G. Pan, and B. Zhu, J. Comput. Phys. **437**, 110325 (2021).
- [95] C. Rao, H. Sun, and Y. Liu, arXiv:2106.04781.