

# Residual Matrix Product State for Machine Learning

Ye-Ming Meng<sup>1,2</sup>, Jing Zhang<sup>3</sup> Peng Zhang<sup>3</sup> Chao Gao<sup>2\*</sup> and Shi-Ju Ran<sup>4†</sup>

<sup>1</sup> State Key Laboratory of Precision Spectroscopy, East China Normal University, Shanghai 200062, China

<sup>2</sup> Department of Physics, Zhejiang Normal University, Jinhua, 321004, China

<sup>3</sup> School of Computer Science and Technology, Tianjin University, Tianjin, China

<sup>4</sup> Department of Physics, Capital Normal University, Beijing 100048, China

\* gaochao@zjnu.edu.cn † sjran@cnu.edu.cn

December 9, 2022

## 1 Abstract

2 Tensor network, which originates from quantum physics, is emerging as an efficient tool  
 3 for classical and quantum machine learning. Nevertheless, there still exists a consider-  
 4 able accuracy gap between tensor network and the sophisticated neural network models  
 5 for classical machine learning. In this work, we combine the ideas of matrix product  
 6 state (MPS), the simplest tensor network structure, and residual neural network and  
 7 propose the residual matrix product state (ResMPS). The ResMPS can be treated as a  
 8 network where its layers map the “hidden” features to the outputs (e.g., classifications),  
 9 and the variational parameters of the layers are the functions of the features of the sam-  
 10 ples (e.g., pixels of images). This is different from neural network, where the layers  
 11 map feed-forwardly the features to the output. The ResMPS can equip with the non-  
 12 linear activations and dropout layers, and outperforms the state-of-the-art tensor net-  
 13 work models in terms of efficiency, stability, and representation power. Besides, ResMPS  
 14 is interpretable from the perspective of polynomial expansion, where the factorization  
 15 and exponential machines naturally emerge. Our work contributes to connecting and  
 16 hybridizing neural and tensor networks, which is crucial to further enhance our under-  
 17 standing of the working mechanisms and improve the performance of both models.

18

## 19 Contents

20	<b>1 Introduction</b>	<b>2</b>
21	<b>2 Residual matrix product state</b>	<b>3</b>
22	2.1 Definition of residual matrix product state	3
23	2.2 The Architecture of ResMPS	4
24	2.3 The working mechanism of ResMPS	6
25	2.4 Benchmarking results	7
26	2.4.1 Classification accuracy	7
27	2.4.2 Redundancy of regular MPS	7
28	2.4.3 Sparse ResMPS and its representation power	9
29	2.5 Comparison with the recurrent neural network and transformer	9
30	<b>3 Properties of the residual structure</b>	<b>10</b>
31	3.1 Avoiding the gradient problems by residual terms	10
32	3.2 Relations to polynomial expansion	10

33	<b>4 Conclusion</b>	<b>12</b>
34	<b>A Training details</b>	<b>13</b>
35	<b>B Path independency</b>	<b>13</b>
36	<b>C Additional benchmarks</b>	<b>14</b>
37	<b>References</b>	<b>16</b>

---

38  
39

## 40 1 Introduction

41 The tensor network (TN), as a mathematical model widely used to describe quantum many-  
42 body states [1–4], has been successfully applied to machine learning (ML). For instance, TN is  
43 used in supervised and unsupervised image classification, natural language processing, etc. [5–  
44 11]. Several recent works also demonstrate TN’s ability of establishing the connection between  
45 physics and artificial intelligence [12, 13]. Nevertheless, despite the high interpretability of  
46 TN [14–16], there still exists a considerable performance gap between TN and neural network  
47 (NN) [7, 17].

48 In the context of quantum physics, TN is used to represent linear ansatz. While in machine  
49 learning, TN realizes a non-linear map from the features to the outputs, where there exists  
50 a local kernel function [5] that maps the features of the samples to the quantum states in  
51 Hilbert space. It is still an open issue to determine whether the NN techniques can enhance  
52 TN performance. Several recent works have explored different ways to combine TN and NN,  
53 which includes adopting the convolutional neural network (CNN) as a feature extractor in  
54 TN [7, 17, 18]; compressing the linear layers of deep NN by matrix product operators [19];  
55 and implementing the convolutional operations using TN [20], etc. These attempts further  
56 motivate us to investigate the possible hybridization of TN and NN.

57 In this work, we incorporate the information highways (also known as shortcuts) [21, 22],  
58 non-linear activations, and dropout [23] into TN (MPS in specific), and propose Residual MPS  
59 (ResMPS in short). The essential underlying idea of ResMPS is a delicate way of inputting  
60 data such that the variational parameters of the network layers are the functions of the data  
61 features. Such idea is inspired by the traditional feed-forward neural network (FNN), while  
62 in FNN, the data is input only in the initial step.

63 We provide two specific examples of ResMPS dubbed as simple and activated ResMPS. The  
64 simple version (sResMPS in short) is a multi-linear model that can exactly be written into a  
65 standard MPS, and the activated version (aResMPS in short) is a non-linear model equipped  
66 with NN layers. The results on Fashion-MNIST show that the simple ResMPS achieves the same  
67 accuracy as MPS while its parameter complexity is half of the MPS. For the activated ResMPS,  
68 we find a significant enhancement of efficiency and accuracy by introducing the non-linear  
69 activations and the dropout layers on the residual terms.

70 Furthermore, we determine the model interpretability of sResMPS by polynomial expres-  
71 sion. The truncated model achieves a high level of accuracy with only a few low-order terms  
72 of sResMPS. Surprisingly, the factorization [24] and exponential machines [25] have naturally  
73 emerged in this expansion scheme. ResMPS shows the underlying connections between TN  
74 and NN for ML and can shed light on novel possibilities and flexibility for developing powerful  
75 ML models beyond NN or TN.

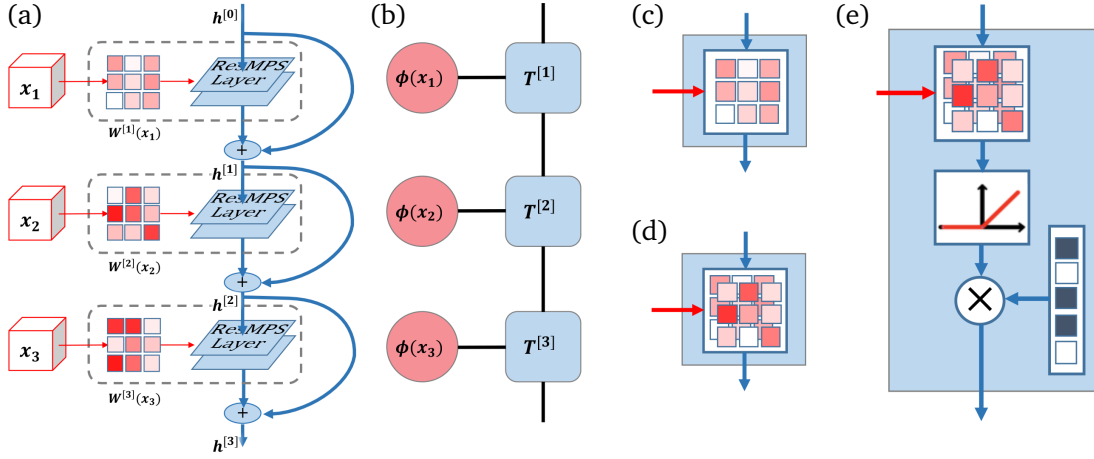


Figure 1: Illustrations of a typical ResMPS compared with a standard MPS. (a) An illustration of ResMPS containing a three-layer FNN in which the variational parameters are functions of the features,  $\mathbf{x}$ . (b) An illustration of a three-tensor MPS, which is contracted with the feature vectors [see Eq. (13)]. (c) An illustration of sResMPS, which is only parameterized by a single channel weight matrix. (d) An illustration of ResMPS, which is equivalent to the standard MPS. (e) An illustration of aResMPS, where the hidden feature will pass through a two-channel linear layer, ReLU activation, and dropout layer in sequence.

76 The remainder of this paper is organized as follows. We define ResMPS in Sec. 2.1 and  
 77 present its two typical architectures in Sec. 2.2. We decipher the working mechanism of  
 78 ResMPS in Sec. 2.3. The image recognition benchmarks compared with other networks are  
 79 exhibited in 2.4. Relation to recurrent neural network (RNN) and transformer is discussed in  
 80 Sec. 2.5. The contribution of residual terms to training stability is discussed in Sec. 3.1. We  
 81 establish equivalence between truncated ResMPS and factorization machine in Sec. 3.2. We  
 82 conclude the result in Sec. 4.

## 83 2 Residual matrix product state

### 84 2.1 Definition of residual matrix product state

85 The traditional FNN, including the residual neural network, consists of multiple trainable lay-  
 86 ers [26]. For instance, in supervised learning, FNN maps the input sample  $\mathbf{x}$  to the output  $l$ ,  
 87 e.g., the corresponding classification. The typical form of one layer can be written as

$$\mathbf{h}^{[n]} = \sigma \left( F^{[n]} \left( \mathbf{h}^{[n-1]}; \mathbf{W}^{[n]} \right) + \mathbf{b}^{[n]} \right), \quad (1)$$

88 where  $\mathbf{h}^{[n-1]}$  denotes the hidden variables that are input to the  $n$ -th layer with  $\mathbf{h}^{[0]} = \mathbf{x}$ ,  $F^{[n]}$   
 89 denotes the mapping of the  $n$ -th layer (e.g., fully connected, convolution, or pooling layer).  
 90 Each layer may consist some variational parameters  $\mathbf{W}^{[n]}$  (weights) and  $\mathbf{b}$  (bias). Furthermore,  
 91  $\sigma$  denotes the activation function.

92 Inspired by the matrix product state [27, 28] and residual neural network [21, 22], here  
 93 we propose a novel machine learning architecture dubbed as residual matrix product state  
 94 (ResMPS). Different from FNN [see, Eq. (1)], ResMPS does not explicitly map the features  
 95 with a feed-forward network. Instead, it uses the features to parameterize FNN variational  
 96 parameters. This enables an FNN to map the hidden features to the expected outputs (see,

97 Fig. 1). In the ResMPS, the mapping of one layer is

$$\mathbf{h}^{[n]} = \mathbf{h}^{[n-1]} + \mathbf{v}^{[n]}(\mathbf{h}^{[n-1]}; \mathbf{W}^{[n]}(x_n), \mathbf{b}^{[n]}), \quad (2)$$

98 where the second term  $\mathbf{v}^{[n]}$  is the function of  $\mathbf{h}^{[n]}$  parametrized by weights  $\mathbf{W}^{[n]}$  and bias  $\mathbf{b}^{[n]}$ ,  
 99 and weights  $\mathbf{W}^{[n]}$  are further parameterized by the  $n$ -th feature  $x_n$ , which distinguish ResMPS  
 100 from FNNs; the initial hidden  $\mathbf{h}^{[0]}$  is initialized by ones for simplicity. Therefore, the depth of  
 101 ResMPS depends on the input size. Eq. (2) gives the general form of  $\mathbf{v}^{[n]}$ , in this work we limit  
 102  $\mathbf{v}^{[n]}$  in a simple form,

$$\mathbf{v}^{[n]} = \sigma(L^{[n]}(\mathbf{h}^{[n-1]}; \mathbf{W}^{[n]}(x_n)) + \mathbf{b}^{[n]}). \quad (3)$$

103 where  $L^{[n]}$  is a linear map, and  $\sigma$  is the activation. Similar to ResNet, the output of one layer is  
 104 the addition of the output of  $\mathbf{v}^{[n]}$ , and the input includes the hidden features. This is to form a  
 105 shortcut of the information flow, which can avoid the vanishing/explosion of the gradients (see  
 106 Sec. 3.1 for details). We further note that one obtains a standard FNN by adopting  $\mathbf{h}^{[0]} = \mathbf{x}$   
 107 and removing the dependence of  $\mathbf{W}$  on  $\mathbf{x}$ .

108 Note that the dimension of the hidden variable is not equal to the dimension of the label  
 109 index (i.e., the number of classes in the discriminative task). Therefore, one additional linear  
 110 layer without bias and activation should be added to map the final hidden  $\mathbf{h}^{[N]}$  to the output  
 111  $f^{(l)}$  as  $f^{(l)} = \sum_i L_{li} h_i^{[N]}$  where  $\dim(l)$  equals to the number of classes. The linear map can be  
 112 flexibly replaced by any other map as long as its output dimension matches the dimension of  
 113 the label.

## 114 2.2 The Architecture of ResMPS

115 In the following, we examine two instances of ResMPS, called simple ResMPS [sResMPS, see  
 116 Fig.1(c)] and activated ResMPS [aResMPS, see Fig.1(e)]. The sResMPS is a multi-linear model  
 117 that is equivalent to MPS. It achieves the same accuracy with only half of the parameter com-  
 118 plexity of the MPS. The aResMPS is a generalized version of sResMPS, in which the general-  
 119 ization efficiency is enhanced by introducing non-linear activation functions and dropout in  
 120 the FNN part. The map of one layer in the sResMPS is written as

$$h_j^{[n]} = h_j^{[n-1]} + \sum_i x_n W_{ij}^{[n]} h_i^{[n-1]}. \quad (4)$$

121 The weights of the layers in the FNN are linearly dependent on the features  $\mathbf{x}$ . The bias terms  
 122 are also disabled in this example.

123 sResMPS is equivalent to a restricted version of MPS and can achieve identical performance  
 124 with only a half parameter complexity of standard MPS. See Sec. 2.4.2 for details.

125 It is seen that MPS has a remarkable representation power. The training error is less than  
 126 1% [29]. However, a gap between the training and testing accuracy suggests an over-fitting  
 127 issue. To address the over-fitting issue, we propose the aResMPS by incorporating the non-  
 128 linear activation functions and dropout. This modification also enhances the generalization  
 129 power [30]. The map of each layer in the FNN of the aResMPS is more-or-less a fully-connected  
 130 layer with a shortcut, which reads

$$\mathbf{h}^{[n]} = \mathbf{h}^{[n-1]} + \sigma(L^{[n]}(\mathbf{h}^{[n-1]}) + \mathbf{b}^{[n]}), \quad (5)$$

131 where  $\sigma$  is an activation function. The map  $L^{[n]}$  relies on the feature  $x_n$  in a non-linear fashion

$$L^{[n]}(\mathbf{h}^{[n-1]})_j = \sum_{c=1,2} \left[ \xi^{[c]}(x_n) \sum_i W_{ij}^{[n,c]} h_i^{[n-1]} \right]. \quad (6)$$

132 The architecture of ResMPS is flexible, due to the choice of  $\xi^{[c]}(x_n)$  and the number of  
 133 channels  $\dim(c)$ . Here we choose  $\xi^{[1]}(x_n) = x_n$  and  $\xi^{[2]}(x_n) = 1 - x_n$ . We introduce  $\xi^{[c]}$  to  
 134 enhance the non-linearity of the aResMPS. It is worth mentioning that even sResMPS repre-  
 135 sents a non-linear map on the features  $\mathbf{x}$  (but a linear map on the hidden features).

136 For the aResMPS, the map on either the features or the hidden features is non-linear. In-  
 137 deed, the FNN embedded inside the aResMPS can be replaced by any NN. Here, we choose a  
 138 standard fully-connected network with two channels labeled by  $c$ .

139 Throughout this paper, we choose the ReLU activation function that can screen the nega-  
 140 tive inputs [31, 32]. Due to its piecewise linear characteristics, the gradient can directly pass  
 141 through without attenuation or enhancement. Therefore, the ReLU function is suitable for en-  
 142 hancing the non-linearity of the deep networks, which can improve its representation power  
 143 and avoid the vanishing/explosion of the gradient. Furthermore, we use a dropout layer com-  
 144 bined with the residual structure to improve the generalization ability of ResMPS. The dropout  
 145 layer effectively creates an ensemble of networks while avoiding the co-adaptation of inter-  
 146 mediate variables [23, 33, 34]. We impose dropout on the residual terms, i.e.  $\mathbf{h}^{[n]} = \mathbf{h}^{[n-1]}$   
 147  $+ \text{dropout}(\sigma(\dots))$ .

148 If we discard the activation and the dropout layers of aResMPS [see Fig.1(e)], we will get  
 149 a standard two-channel MPS. For a standard MPS with physical bond dimension  $d = 2$ , the  
 150 map given by a local-tensor construction is [29]

$$h_j^{[n]} = \sum_{c=1,2} \left[ \xi^{[c]}(x_n) \sum_i T_{ij}^{[n,c]} h_i^{[n-1]} \right]. \quad (7)$$

151 If we introduce transformation  $T_{ij}^{[n,c]} = W_{ij}^{[n,c]} + \delta_{ij}$ , we can simply get

$$h_j^{[n]} = h_j^{[n-1]} \left( \sum_{c=1,2} \xi^{[c]}(x_n) \right) + \sum_{c=1,2} \left[ \xi^{[c]}(x_n) \sum_i W_{ij}^{[n,c]} h_i^{[n-1]} \right]. \quad (8)$$

152 By taking feature map with norm-1 normalization [29], i.e.  $\sum_{c=1,2} \xi^{[c]}(x_n) = 1$ , we get a  
 153 ResMPS with map

$$h_j^{[n]} = h_j^{[n-1]} + \sum_{c=1,2} \left[ \xi^{[c]}(x_n) \sum_i W_{ij}^{[n,c]} h_i^{[n-1]} \right]. \quad (9)$$

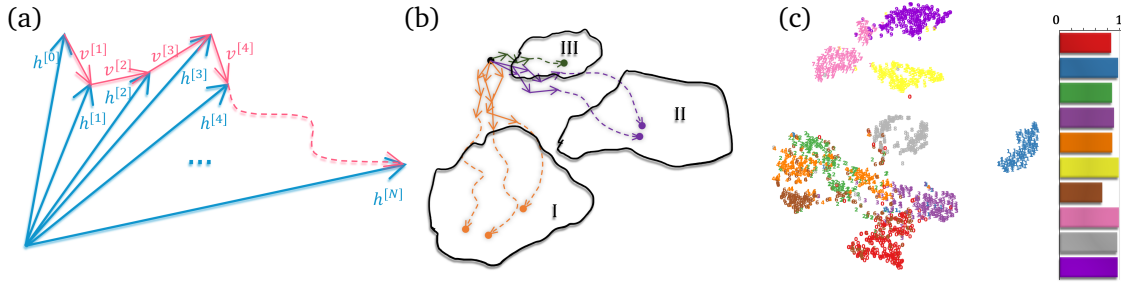


Figure 2: Encoding process of the ResMPS. (a) An illustration of a high-dimensional path of one sample. Blue arrows represent hidden features between different layers. Red arrows represent shift-vectors contributed by the residual part. (b) An illustration of the aggregation behavior of samples. The same color denotes samples belonging to the same class. (c) The two-dimensional data distribution generated by t-SNE on the endpoint dimension reduction of (b), the data points come from the Fashion-MNIST data set, and the corresponding accuracy is on the right. Note we reduce the feature dimension to 2 for illustration, which is far less than the original dimension of the hidden and weakened separations of the samples from different classes.

### 154 2.3 The working mechanism of ResMPS

155 We illustrate the path of the hidden state  $\mathbf{h}^{[i]}$  of ResMPS in the high-dimensional vector space  
 156 [as shown in Fig.2(a)]. Each layer of the ResMPS updates the state  $\mathbf{h}^{[i]}$  once to make it one  
 157 step forward with shift-vector  $\mathbf{v}^{[i+1]} = \mathbf{h}^{[i+1]} - \mathbf{h}^{[i]}$ . After passing through all layers, all shift-  
 158 vectors are connected into a continuous path, namely  $\sum_{i=1}^N \mathbf{v}^{[i]}$ . For the same ResMPS, differ-  
 159 ent features of the samples share the same initial point (i.e.,  $\mathbf{h}^{[0]}$ ). Since the parameter  $W$  of  
 160 shift-vector  $\mathbf{v}$  is a function of feature  $\mathbf{x}$ , the path encodes the information of samples. Besides,  
 161 Similar samples have close paths in the vector space [as shown in Fig.2(b)]. After training  
 162 convergence, samples of the same category will eventually gather together.

163 In order to show the behavior of paths of the hidden variables in the high-dimensional  
 164 space, we use the Fashion-MNIST dataset to train sResMPS, and use the t-SNE algorithm [35,  
 165 36] to embed the endpoints of the ResMPS to a two-dimensional plane after the network  
 166 converges. Note that before we apply t-SNE for dimensionality reduction, the original virtual  
 167 feature has 100 components, which is sufficient large for accuracy yet economical. Fig. 2(c)  
 168 illustrates the visualization of final hidden features  $h_t^{[N]}$  in the two-dimensional space. Samples  
 169 with better classification accuracy are relatively separated, while those with poor classification  
 170 accuracy overlap with other classifications. We also demonstrate the visualizations by t-SNE  
 171 for the intermediate hidden features in Fig. 3 to show how the hidden features gradually cluster  
 172 into different sub-regions. It shows that the initial hidden features are uniformly distributed,  
 173 and then gradually separated during the encoding process.

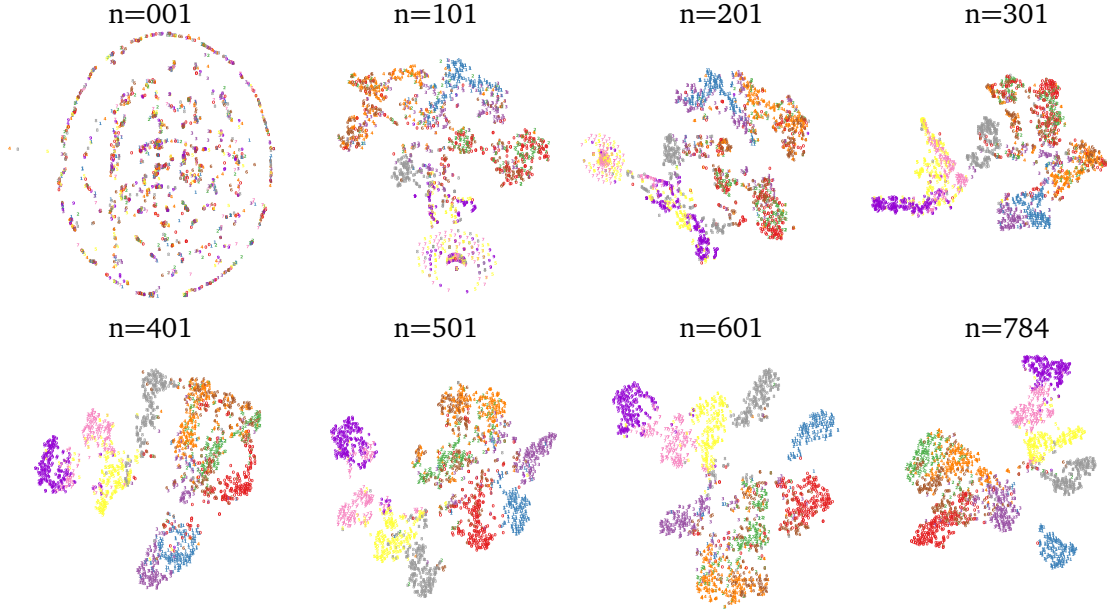


Figure 3: Intermediate hidden features  $\mathbf{h}^{[n]}$  ( $1 \leq n \leq 784$ ) visualized by t-SNE. It is shown that samples of the same classes gradually enter the same regions. This characterizes the encoding process of ResMPS.

## 174 2.4 Benchmarking results

### 175 2.4.1 Classification accuracy

176 For the MNIST and Fashion-MNIST datasets, Table 1 shows the accuracy of the sResMPS and  
 177 aResMPS, compared with several established NN [39] and TN models [7, 9, 15, 17, 29, 38].  
 178 The MPS and ResMPS models represent a high level of representation power, as indicated by  
 179 their high training accuracy. The aResMPS also surpasses the probabilistically interpretable  
 180 Bayesian [15] and other TN models, including the two-dimensional TN known as projected-  
 181 entangled pair state (PEPS) [17]. It also achieves a (slightly) better accuracy than that of  
 182 CNN-PEPS model, in which CNN is adopted as the feature extractor. This accuracy surpasses  
 183 the CNN without the stacking architecture, such as AlexNet [39]. The aResMPS still does not  
 184 overperform the ResNet which is formed by stacking multiple convolution layers.

### 185 2.4.2 Redundancy of regular MPS

186 To see the equivalence between the standard MPS and sResMPS mentioned in Sec.2.2, we  
 187 introduce the third-order tensors  $\mathbf{T}^{[n]}$  satisfying

$$T_{1,::}^{[n]} = \mathbf{I}, \quad T_{2,::}^{[n]} = \mathbf{W}^{[n]}. \quad (10)$$

188 The feature vectors  $\phi(x_n)$  are obtained by the feature map as  $\phi(x_n) = (1, x_n)$  [5, 9, 29].  
 189 Summing their joint index gives a linear mapping represented by a matrix

$$\begin{aligned} A_{a_{n-1}, a_n}^{[n]} &= \sum_{p_n} T_{p_n, a_{n-1}, a_n}^{[n]} \phi(x_{n-1})_{p_n} \\ &= \delta_{a_{n-1}, a_n} + x_n W_{a_{n-1}, a_n}^{[n]}. \end{aligned} \quad (11)$$

Table 1: Experimental results on MNIST [37] and Fashion-MNIST dataset. The first 6 models are pure TN architectures, which means they are multi-linear, and no neural structures like pooling, activation and convolution are introduced. AlexNet, ResNet, and CNN-PEPS are NN or TN-NN hybrid models. For aResMPS, we use ReLU as activation, and the dropout probability is set to be 0.6.

Model	MPS Train	MPS Test	Fashion- MNIST Train	Fashion- MNIST Test
MPS machine [29]	1.0000	0.9880	0.9988	0.8970
Unitary tree TN [9]	0.98	0.95	-	-
Tree curtain model [38]	-	-	0.9538	0.8897
Bayesian TN [15]	-	-	0.8950	0.8692
EPS-SBS [7]	-	0.9885	-	0.886
PEPS [17]	-	-	-	0.883
CNN-PEPS [17]	-	-	-	0.912
AlexNet [39]	-	-	-	0.8882
ResNet [39]	-	-	-	0.9339
sResMPS	1.0000	0.9873	0.9987	0.8909
aResMPS	1.0000	0.9907	0.9999	0.9142

190 Applying this mapping to  $\mathbf{h}^{[n-1]}$ , one gets

$$\begin{aligned}
 h^{[n]} &= \sum_{a_{n-1}} h_{a_{n-1}}^{[n-1]} A_{a_{n-1}, a_n}^{[n]} \\
 &= h_{a_n}^{[n-1]} + \sum_{a_{n-1}} h_{a_{n-1}}^{[n-1]} x_n W_{a_{n-1}, a_n}^{[n]}.
 \end{aligned} \tag{12}$$

191 This form is exactly the definition of sResMPS, shown in Eq. (4). Therefore, the sResMPS is  
 192 equivalent to the standard MPS formed by the following tensors

$$\mathcal{T} = \sum_{\{a_n\}} \prod_n T_{p_n a_n a_{n+1}}^{[n]} \tag{13}$$

193 as its tensor-train cores [40] [Fig. 1 (b)]. The numbers of the input and output hidden features  
 194 for different layers provide the two virtual bond dimensions of the MPS, i.e.,  $\{\dim(a_n)\}$ . In  
 195 this work, we fix  $\dim(a_n) = \chi$ ,  $\forall n$ . The physical dimension of the MPS should also match the  
 196 dimension of the feature vector, i.e.  $\dim(\phi(x_n)) = \dim(p_n)$ .

197 For  $\dim(p_n) = 2$ , the number of variational parameters in sResMPS is  $\sim O(N\chi^2)$  where  
 198  $N$  is total number of features. This is only half of that in the MPS which is  $\sim O(2N\chi^2)$ . Our  
 199 numerical simulations show that the accuracy of both models is almost the same. See the  
 200 training and testing accuracy versus epochs on Fashion-MNIST dataset [41] in Fig. 4 (a) with  
 201  $\chi = 40$ . This is because one of the two channels of each tensor in the MPS is much less  
 202 “activated”. The inset of Fig. 4 (a) shows the average norm of the two channels of different  
 203 tensors

$$q_p^{[n]} = \frac{1}{\chi^2} \sum_{j=1}^{\chi} \sum_{k=1}^{\chi} \left| T_{pjk}^{[n]} - \delta_{jk} \right|, \tag{14}$$

204 with  $p = 1, 2$  representing the channels. The main contribution to the output is from the  
 205 second channel. Therefore, one channel is sufficient to pass the information to the output.



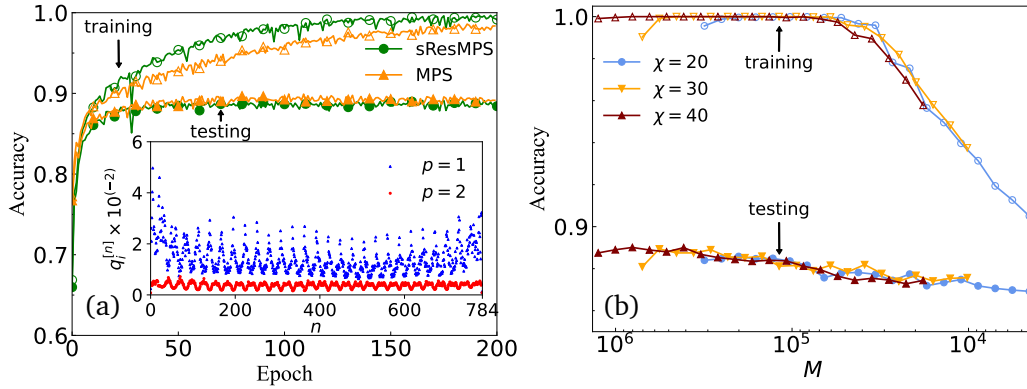


Figure 4: Numerical results of the simple ResMPS. (a) Training and testing accuracy of sResMPS (without dropout) and MPS versus epochs on the Fashion-MNIST dataset. The inset shows the average norm [Eq. (14)] of the two channels in the MPS for different tensors  $n$ . (b) Training and testing accuracy of the sResMPS versus the total number of unmasked weights in the sResMPS. The left end of each curve corresponds to the un-pruned result. It is also seen that the first few steps of pruning improve the accuracy. The total number of parameters of sResMPS with  $\chi = 20, 30,$  and  $40$  equals to about  $3 \times 10^5, 7 \times 10^5$  and  $13 \times 10^5$ , respectively.

### 206 2.4.3 Sparse ResMPS and its representation power

207 For sufficient large  $\chi$ , the accuracy achieves saturation, which shows the possibility to obtain  
 208 a sparse ResMPS by removing those less activated parameters, without reducing its representa-  
 209 tion power. To this end, a pruning process [14] is applied. We first train a dense ResMPS to  
 210 convergence and then mask its parameters with relatively small absolute values by multiplying  
 211 zeros. After that, we optimize unmasked parameters. The mask-retraining step can be applied  
 212 multiple times until one achieves a desired sparsity.

213 Fig. 4 (b) shows the accuracy values versus the number of unmasked parameters  $M$ , for  
 214 different virtual bond dimensions,  $\chi = 20, 30,$  and  $40$ . The training and testing accuracy are  
 215 consistent for different  $M$ . Hence, the representation power of sparse ResMPS is character-  
 216 ized by  $M$ . It's worth noting that to achieve a desired accuracy, the sparse ResMPS needs fewer  
 217 parameters than the dense one. For example, to achieve (approximately) 86.5% testing set  
 218 accuracy, the dense one needs  $\chi^2 L = 12544$  parameters ( $\chi = 4$ ), while the sparse one only  
 219 needs  $M = 4000$  parameters.

### 220 2.5 Comparison with the recurrent neural network and transformer

221 In the first sight the ResMPS is quite similar to the RNN and transformer, which are intensively  
 222 explored nowadays. Here we compare the ResMPS with the later two respectively.

223 The layer mapping of ResMPS is similar to recurrent neural network (RNN). However,  
 224 there are two essential differences between ResMPS and RNN. Firstly, ResMPS is non-uniform,  
 225 and can be reduced to the translation-invariant MPS named as uniform MPS [42, 43]. From  
 226 this point of view, RNN, in which the layers share the same variational parameters, is closer  
 227 to the uniform MPS. Secondly, the recurrent mapping (also called a unit) is constructed by  
 228 several linear operations and non-linear activations. Different constructions of the units define  
 229 different RNN, such as the long short-term memory model. For the ResMPS, it is formed by  
 230 tensor units, which are simply multi-linear mappings and can be analyzed by the established  
 231 methods such as tensor decompositions [2, 27, 40].

232 ResMPS is different from transformer networks. Firstly, the transformer network is es-

233 tablished on the attention mechanism, which is given by the inner product of “queries” and  
 234 “keys”, while ResMPS contains higher order interaction of input features (see its connection  
 235 to exponential machine in Sec. 3.2 for detail). Hence, the working mechanism of ResMPS, as  
 236 explained in Sec. 2.3, is different to transformer networks. Secondly, A general transformer  
 237 network consists of an encoder and a decoder. The encoder transforms a data sequence into a  
 238 vector, while the decoder transforms the vector into another sequence. In contrast, ResMPS,  
 239 and in general the MPS-based networks output a vector. Moreover, ResMPS does not need  
 240 positional encoding, see appendix B for more discussions.

## 241 3 Properties of the residual structure

### 242 3.1 Avoiding the gradient problems by residual terms

243 A typical MPS architecture designed for pattern recognition contains hundreds of tensor cores.  
 244 Such an architecture probably encounters gradient vanishing/exploding problems. To avoid  
 245 the gradient problems, some existing MPS schemes apply a DMRG-like algorithm where the  
 246 MPS takes the canonical form [5, 6, 10, 44]. In these attempts, however, the accuracy is sensitive  
 247 to the hidden features’ dimensions (virtual bonds). Recently, an MPS algorithm was proposed  
 248 based on the automatic gradient technique [29] that can achieve higher accuracy than the  
 249 previous ones, while its performance is not sensitive to the virtual dimensions. To find out  
 250 why such a deep network avoids the gradient problems, here we construct the tensor cores to  
 251 satisfy a specific form given by Eq. (10). The identity in  $T_{1,::}^{[n]}$  plays the role of “highway” to pass  
 252 the information from the previous tensor core directly to the latter ones. The components  $T_{2,::}^{[n]}$   
 253 represent the residual terms, which is  $\ll O(1)$ . The application of residual conditions implies  
 254 that each layer of ResMPS can easily express identity mapping. In other words, the architecture  
 255 of ResMPS satisfies the identity parameterization [21, 22, 45].

256 To further demonstrate the role of identity parameterization in ResMPS, we use Gaussian  
 257 distributions with zero mean and standard deviation  $\varepsilon$  to randomly initialize the elements of  
 258  $T_{2,::}^{[n]}$ . Fig. 5 shows the testing accuracy at the 10-th, 20-th, and 50-th epochs. For a sufficiently  
 259 small  $\varepsilon$ , the accuracy is quickly and stably converged. However, for relatively large  $\varepsilon$  [e.g.,  
 260  $O(10^{-1})$ ] as illustrated by the green region, the gradients become unstable. Consequently, the  
 261 accuracy stays around 0.1 and cannot improve further through the training process. Note that  
 262 this may be unstable in most cases if instead of the identity parameterization, the entire  $T$  is  
 263 randomly initialized.

### 264 3.2 Relations to polynomial expansion

265 The forward propagation of the sResMPS (4) is fully linear on the hidden features. Applying  
 266 the maps to the initial hidden features  $\mathbf{h}_0$  in sequence, we can then rewrite the output hidden  
 267 features in an expansive form [Fig. 6 (b)] as

$$\begin{aligned} \mathbf{h}^{[N]} &= (\mathbf{I} + x_N \mathbf{W}^{[N]}) \dots (\mathbf{I} + x_2 \mathbf{W}^{[2]}) (\mathbf{I} + x_1 \mathbf{W}^{[1]}) \mathbf{h}^{[0]} \\ &= \sum_{k=0}^N \mathbf{M}^{[k]} \mathbf{h}^{[0]}, \end{aligned} \quad (15)$$

268 where  $N$  is the total number of features  $\mathbf{x}$ . The output  $\mathbf{h}^{[N]}$  is the stack of  $N + 1$  terms. The  
 269 zeroth term satisfies  $\mathbf{M}^{[0]} = \mathbf{I}$ , which is the result of the information highway from the first  
 270 input hidden features to the output. The term  $\mathbf{M}^{[1]} = \sum_{\alpha=1}^N x_{\alpha} \mathbf{W}^{[\alpha]}$  is the part in ResMPS  
 271 which is linear on the features  $\mathbf{x}$ . The  $k$ -th term contains the  $k$ -th order contributions from  $\mathbf{x}$ ,

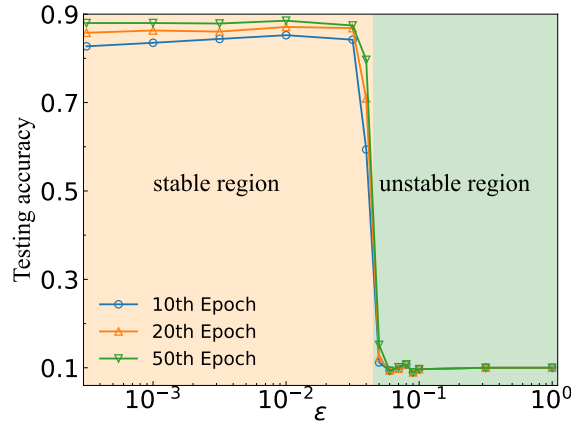


Figure 5: The testing accuracy of the sResMPS versus  $\varepsilon$  on Fashion-MNIST dataset. Here  $\varepsilon$  is the standard deviation of the initial residual part. We fix the number of epochs to be 10, 20, and 50. The network can be trained stably for small values of  $\varepsilon$ . Otherwise, the training process encounters gradient vanishing (or explosion) problems. The stable and unstable regions are illustrated by orange and green colors, respectively. Note that for the unstable region, the value of the network elements diverges. The network will give classifications randomly, thus, the accuracy tends to be 0.1.

272 i.e.,

$$\mathbf{M}^{[k]} = \sum_{\alpha_1 \dots \alpha_k = 1}^N G_{\alpha_1 \dots \alpha_k} x_{\alpha_1} \dots x_{\alpha_k} \mathbf{W}^{[\alpha_1]} \dots \mathbf{W}^{[\alpha_k]}, \quad (16)$$

$$G_{a_1, a_2, \dots, a_n} = \begin{cases} 1, & a_1 > a_2 > \dots > a_n \\ 0, & \text{otherwise.} \end{cases} \quad (17)$$

273 This formula is actually a specific form of the Exponential Machines [25]. Due to their essential  
 274 similarity, the algebraic properties of Exponential Machines are also valid for sResMPS. For  
 275 instance, the output feature  $\mathbf{h}^{[N]}$  is a linear mapping concerning the initial hidden feature  $h_0$ ,  
 276 and a multi-linear mapping concerning the feature  $\mathbf{x}$ .

277 Due to the residual condition [see Eq. (10) with  $|\mathbf{W}^{[n]}| \ll O(10^{-1})$ ], the contributions from  
 278 the higher-order terms of Eq. (16) should decay exponentially with  $k$ . Therefore, we can define  
 279 a set of lower-order effective models by retaining the first few terms. For instance, by keeping  
 280 the zeroth- and first-order terms in Eq. (16), we obtain a model in which the output features  
 281 are linear to both hidden and sample features. Keeping the zeroth, linear, and quadratic terms  
 282 gets a model

$$\mathbf{h}^{[N](2)} = \left( \mathbf{I} + \sum_{\alpha=1}^N x_{\alpha} \mathbf{W}^{[\alpha]} + \sum_{\alpha, \beta=1}^N G_{\alpha, \beta} x_{\alpha} x_{\beta} \mathbf{W}^{[\alpha]} \mathbf{W}^{[\beta]} \right) \mathbf{h}^{[0]}. \quad (18)$$

283 This model is similar to Factorization Machines [24] and polynomial NN [46].

284 Fig. 6 (a) shows the difference between the accuracy of several lower-order models and  
 285 the sResMPS. This implies that the significant improvement achieved by the sResMPS has its  
 286 root in a few lower-order terms, especially the linear term. As the order increases, the cost of  
 287 directly computing Eq. (15) is also exponentially increased. Therefore, truncating the order  
 288 of expansion is not economical. ResMPS adopts a different and efficient scheme for retaining  
 289 all higher-order interactions.

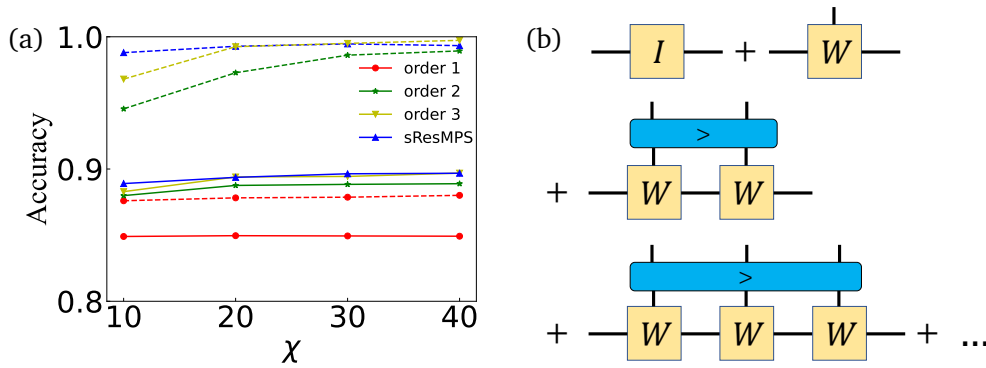


Figure 6: Polynomial expansion based on ResMPS. (a) Training (dashed) and testing (solid) accuracy versus  $\chi$  by taking different orders in the expansion form; (b) An illustration of the polynomial expansion picture of the sResMPS, see Eq. (16).

## 290 4 Conclusion

291 We propose ResMPS by incorporating MPS with the information highways, non-linear activa-  
 292 tions, and dropout. In contrast to FNN, the variational parameters in ResMPS are replaced by  
 293 adjustable functions. For FNN, features are input at the first layer of the network. For ResMPS,  
 294 however, features are divided and input into the weight matrices of each layer, which is in-  
 295 herited from MPS. Furthermore, the introduction of neural network structures in the ResMPS  
 296 brings more vital representation power than the usual MPS. For concreteness, we present two  
 297 specific versions of ResMPS.

298 The first derived architecture, sResMPS, is simply a linear version of ResMPS. By comparing  
 299 MPS' learning performance on the Fashion-MNIST dataset, we reveal the channel redundancy  
 300 of MPS. sResMPS thus discards the redundant channel. Consequently, it achieves consistent  
 301 accuracy while the parameter complexity is reduced by half.

302 The second derived architecture, aResMPS, is equipped with activation and dropout layers.  
 303 We compare aResMPS with several TN and NN models on the Fashion-MNIST dataset. The  
 304 activation and dropout layer enhance the non-linearity and generalization ability of the model,  
 305 respectively. Therefore, aResMPS surpasses the state-of-the-art TN methods and AlexNet in  
 306 terms of accuracy, although it is still inferior to ResNet formed by stacking multiple convolution  
 307 layers. Going beyond present aResMPS to achieve higher accuracy, e.g., replacing the weight  
 308 matrices with convolution layers, is a valuable improvement direction of ResMPS.

309 The perspectives of the residual network derive the polynomial expansion of ResMPS. The  
 310 benefits are two-fold. Firstly, we give the condition of vanishing/explosion of the gradients  
 311 of ResMPS. This helps the feature design of MPS and ResMPS algorithms with stable conver-  
 312 gence. Secondly, it establishes the equivalence between MPS and polynomial networks such as  
 313 Factorization Machines and Exponential Machines. Further numerical evidence suggests that  
 314 the contribution of high-order terms is insignificant. This helps to better understand the MPS  
 315 and ResMPS.

316 Are other NN structures (e.g., convolution and pooling layers) compatible with ResMPS?  
 317 Is it possible to propose a ResMPS structure based on general NN structures (e.g., Tree TN or  
 318 Projected Entangled-Pair States)? These problems are worthy of investigation in the future.

## 319 Acknowledgements

320 **Funding information** Y.-M.M. and C.G. are supported by National Natural Science Founda-  
 321 tion of China (NSFC, Grant No. 1183501 and No. 12074342) and Zhejiang Provincial Natural  
 322 Science Foundation of China (Grant No. LR22A040001 and No. LY21A040004). S.-J.R. is  
 323 supported by NSFC (Grant No. 12004266 and No. 11834014), Beijing Natural Science Foun-  
 324 dation (No. 1192005 and No. Z180013), Foundation of Beijing Education Committees (No.  
 325 KM202010028013), and the Academy for Multidisciplinary Studies, Capital Normal Univer-  
 326 sity. J.Z. and P.Z. are supported by NSFC (Grant No. 61772363).

## 327 A Training details

328 All benchmarks of this paper are implemented on MNIST and Fashion-MNIST datasets, each  
 329 of which includes 60,000 training and 10,000 testing grayscale images with  $L = 10$  labels.  
 330 To fit inputs of ResMPS, we choose a specific path to reorder 2D pixels into a 1D squence  
 331  $\{(x_1, x_2, \dots, x_n; n = 784) | x_i \in [0, 1]\}$ , where 784 is the pixel number of one image. For single  
 332 channel ResMPS, e.g., the sResMPS, there is no need to feature map the original data. And  
 333 for double channel ResMPS, a linear feature map  $(x, 1 - x)$  is adapted to match the channel  
 334 dimension. The dimension of the hidden feature  $\chi$  is set to 100, which is sufficient large for  
 335 accuracy yet economical. The predicted classification is given by the largest component

$$f(\mathbf{x}) = \operatorname{argmax}_l f^{(l)}(\mathbf{x}), \quad (19)$$

336 where  $f^{(l)}(\mathbf{x})$  is the overall mapping of ResMPS with  $l = 1, \dots, L$  denotes classification. To train  
 337 ResMPS, we use the Stochastic Gradient Descent (SGD) method with Adam optimizer [47] and  
 338 learning rate  $10^{-4}$ . Data are divided into mini-batches of size 1000. We choose cross-entropy  
 339 as the loss function,

$$CE = - \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \log \operatorname{softmax}^{f(y_i)}(\mathbf{x}_i) \quad (20)$$

340 where

$$\operatorname{softmax}^{f(y_i)}(\mathbf{x}_i) = \frac{e^{f(y_i)(\mathbf{x}_i)}}{\sum_{l=0}^{L-1} e^{f(y_l)(\mathbf{x}_i)}}. \quad (21)$$

341 Here  $\mathcal{D}$  is the data set of a mini-batch, and  $\mathbf{x}_i$  and  $y_i$  are the  $i$ -th input and classification  
 342 respectively. If dropout is taken, the corresponding probability is set to 0.6. We take ReLU  
 343 for the non-linear case. The whole implementation is based on the PyTorch library, and the  
 344 relevant code is available on GitHub.

## 345 B Path independency

346 ResMPS naturally deals with sequential data, but general data like images are usually high-  
 347 dimensional. Therefore, one needs to choose a specific path to unfold high-dimensional data.  
 348 To study how path choice effects the performance, we here compare three typical paths —  
 349 zigzag path, Hilbert path, and random path. The first two are illustrated in Fig. (7). The zigzag  
 350 path, as the most popular one, arranges data row by row; the Hilbert path preserves more  
 351 neighboring information than others; as for the random path, the original position information  
 352 is completely dropped, i.e., all points are collected together, shuffled completely, and then  
 353 reordered in a random way. Note that the numbers of rows and columns of the Hilbert path  
 354 need to fit  $2^n$  with  $n$  a positive integer, so we need to extend the original image by using

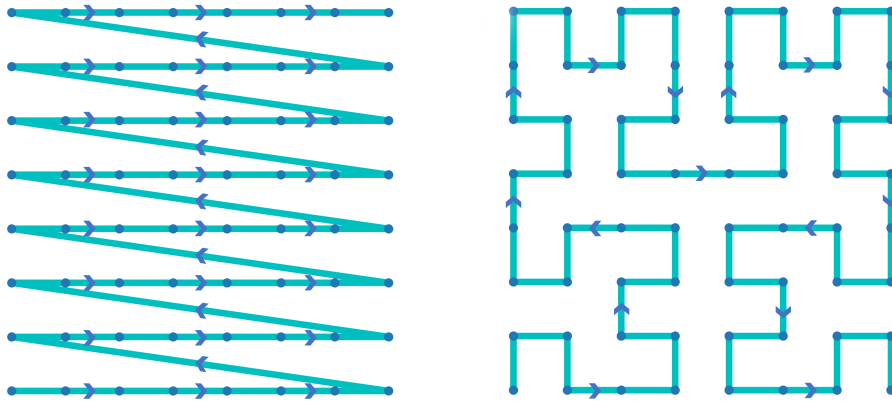


Figure 7: The zigzag path (left) and the Hilbert path (right).

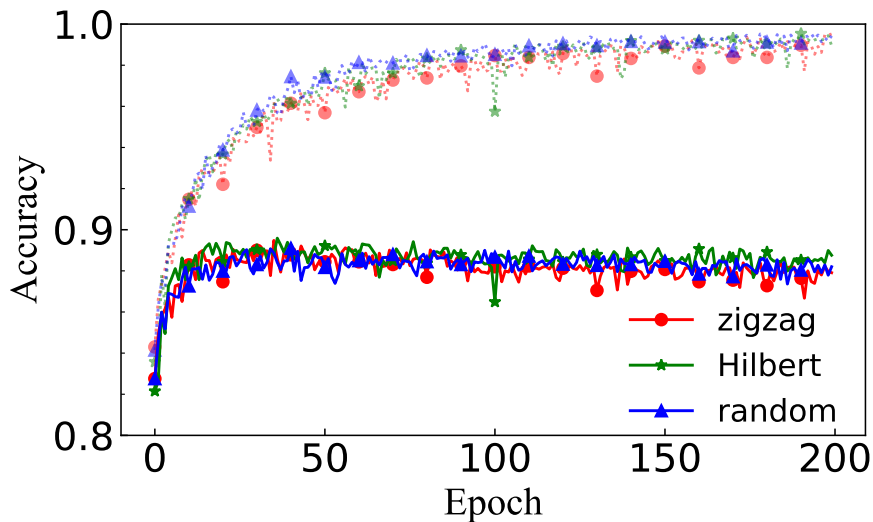


Figure 8: Convergence history of different paths on Fashion-MNIST dataset.

355 borders consisting of zeros. For the case of MNIST and Fashion-MNIST, image size is extended  
 356 from  $28 \times 28$  to  $32 \times 32$ .

357 We run benchmarks for sResMPS without activation and dropout on the Fashion-MNIST  
 358 dataset. The hidden dimension is set to be 40 and other training details are the same as  
 359 appendix A. The result is shown in Fig. 8. To our surprise, there is no significant performance  
 360 difference observed, even the random path still performs well.

361 Since the correlation decay exponentially in MPS [48], one would expect that ResMPS can  
 362 not capture long-range information, and hence is path-independent. However, the previous ex-  
 363 periment show that ResMPS is path-independent. We suppose that this inconsistency is mainly  
 364 due to the introduction of residual connection, which preserve the sensitivity of gradient in  
 365 the long-range, and presumably maintains the long-range correlation.

## 366 C Additional benchmarks

367 Note that aResMPS can achieve better performance than sResMPS, while the former differs to  
 368 the later in three aspects: the channel number is doubled, the nonlinear action and dropout are  
 369 added. To explore the contribution of each aspect to the final performance, we use different

Table 2: Benchmarks for fine partitioned models. sResMPS and aResMPS corresponds to the first and the last row respectively. Other models filled the gap between sResMPS and aResMPS.

Model	MPS Train	MPS Test	Fashion- MNIST Train	Fashion- MNIST Test
1-channel, -ReLU, -dropout	1.0000	0.9873	0.9987	0.8909
1-channel, -ReLU, +dropout	1.0000	0.9889	0.9618	0.9022
1-channel, +ReLU, -dropout	1.0000	0.9864	1.0000	0.8957
1-channel, +ReLU, +dropout	1.0000	0.9885	0.9904	0.9108
2-channel, -ReLU, -dropout	1.0000	0.9880	0.9988	0.8970
2-channel, -ReLU, +dropout	1.0000	0.9896	0.9899	0.9081
2-channel, +ReLU, -dropout	1.0000	0.9873	1.0000	0.9102
2-channel, +ReLU, +dropout	1.0000	0.9907	0.9999	0.9142

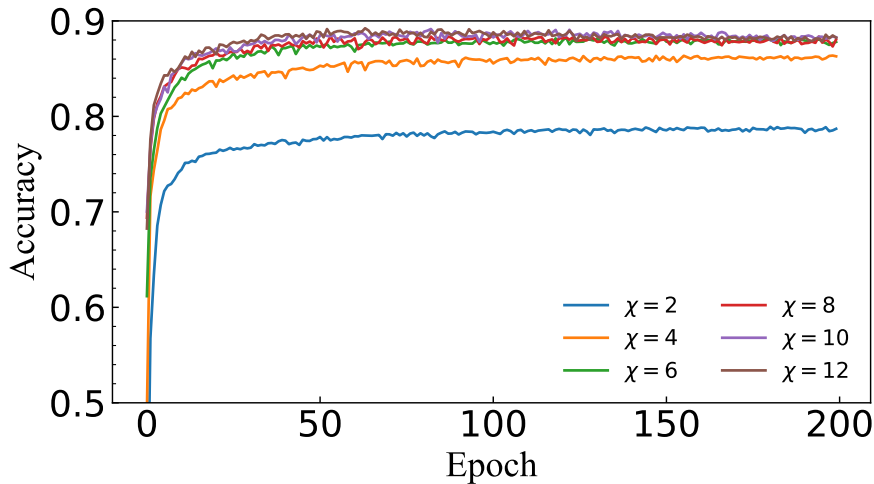


Figure 9: Testing accuracy for sResMPS on Fashion-MNIST dataset for various hidden dimension  $\chi$ . The performance will increase as  $\chi$  increase for small  $\chi$ , and achieves saturation after a critical point about  $\chi_c = 6$ , which is the same as the target space dimension.

370 combinations to give more fine partitioned models. The benchmark results are shown in Ta-  
 371 ble 2. From the experimental results, we can see that these components more or less improve  
 372 the performance, and aResMPS as the most complex one achieves maximum performance.

373 The result of the pruning test shows that the hidden dimension  $\chi$  does not affect the  
 374 accuracy, but this statement is only valid for sufficiently large  $\chi$ . If we gradually reduce the  
 375 value of  $\chi$ , a predictable result is that the dimensions of hidden space are too small, so that  
 376 different classes become hard to distinguish. We show convergence procedure of sResMPS for  
 377 small  $\chi$ s, see Fig. 9. The results show that increasing the bond dimension indeed helps to  
 378 improve the accuracy, but only for a small enough  $\chi$ . Once a certain threshold is exceeded,  
 379 increasing  $\chi$  has no benefit to accuracy growth. Another fact is that even a really small value  
 380 of hidden dimension  $\chi = 2$  has already achieved about 78% accuracy.

## References

- 381
- 382 [1] U. Schollwöck, *The density-matrix renormalization group in the age of matrix product*  
383 *states*, *Annals of Physics* **326**(1), 96 (2011), doi:[10.1016/j.aop.2010.09.012](https://doi.org/10.1016/j.aop.2010.09.012).
- 384 [2] S.-J. Ran, E. Tarrico, C. Peng, X. Chen, L. Tagliacozzo, G. Su and M. Lewenstein, *Tensor*  
385 *Network Contractions*, Springer International Publishing, doi:[10.1007/978-3-030-](https://doi.org/10.1007/978-3-030-34489-4)  
386 [34489-4](https://doi.org/10.1007/978-3-030-34489-4) (2020).
- 387 [3] F. Verstraete, V. Murg and J. Cirac, *Matrix product states, projected entangled pair states,*  
388 *and variational renormalization group methods for quantum spin systems*, *Advances in*  
389 *Physics* **57**(2), 143 (2008), doi:[10.1080/14789940801912366](https://doi.org/10.1080/14789940801912366).
- 390 [4] G. Evenbly and G. Vidal, *Tensor network states and geometry*, *Journal of Statistical Physics*  
391 **145**(4), 891 (2011), doi:[10.1007/s10955-011-0237-4](https://doi.org/10.1007/s10955-011-0237-4).
- 392 [5] E. Stoudenmire and D. J. Schwab, *Supervised learning with tensor networks*, In D. D.  
393 Lee, M. Sugiyama, U. V. Luxburg, I. Guyon and R. Garnett, eds., *Advances in Neural*  
394 *Information Processing Systems 29*, pp. 4799–4807. Curran Associates, Inc. (2016).
- 395 [6] Z.-Y. Han, J. Wang, H. Fan, L. Wang and P. Zhang, *Unsupervised generative*  
396 *modeling using matrix product states*, *Physical Review X* **8**, 031012 (2018),  
397 doi:[10.1103/PhysRevX.8.031012](https://doi.org/10.1103/PhysRevX.8.031012).
- 398 [7] I. Glasser, N. Pancotti and J. I. Cirac, *From probabilistic graphical models to gen-*  
399 *eralized tensor networks for supervised learning*, *IEEE Access* **8**, 68169 (2020),  
400 doi:[10.1109/access.2020.2986279](https://doi.org/10.1109/access.2020.2986279).
- 401 [8] S. Cheng, L. Wang, T. Xiang and P. Zhang, *Tree tensor networks for generative modeling*,  
402 *Physical Review B* **99**, 155131 (2019), doi:[10.1103/PhysRevB.99.155131](https://doi.org/10.1103/PhysRevB.99.155131).
- 403 [9] D. Liu, S.-J. Ran, P. Wittek, C. Peng, R. B. García, G. Su and M. Lewenstein, *Machine*  
404 *learning by unitary tensor network of hierarchical tree structure*, *New Journal of Physics*  
405 **21**(7), 073059 (2019), doi:[10.1088/1367-2630/ab31ef](https://doi.org/10.1088/1367-2630/ab31ef).
- 406 [10] Z.-Z. Sun, S.-J. Ran and G. Su, *Tangent-space gradient optimization of ten-*  
407 *sor network for machine learning*, *Physical Review E* **102**, 012152 (2020),  
408 doi:[10.1103/PhysRevE.102.012152](https://doi.org/10.1103/PhysRevE.102.012152).
- 409 [11] P. Zhang, Z. Su, L. Zhang, B. Wang and D. Song, *A quantum many-body wave function*  
410 *inspired language modeling approach*, In *Proceedings of the 27th ACM International Confer-*  
411 *ence on Information and Knowledge Management*. ACM, doi:[10.1145/3269206.3271723](https://doi.org/10.1145/3269206.3271723)  
412 (2018).
- 413 [12] J. Chen, S. Cheng, H. Xie, L. Wang and T. Xiang, *Equivalence of restricted boltz-*  
414 *mann machines and tensor network states*, *Physical Review B* **97**(8), 085104 (2018),  
415 doi:[10.1103/physrevb.97.085104](https://doi.org/10.1103/physrevb.97.085104).
- 416 [13] V. Khulkov, A. Novikov and I. Oseledets, *Expressive power of recurrent neural networks*,  
417 In *International Conference on Learning Representations* (2018).
- 418 [14] Y. Levine, O. Sharir, N. Cohen and A. Shashua, *Quantum entanglement*  
419 *in deep learning architectures*, *Physical Review Letters* **122**, 065301 (2019),  
420 doi:[10.1103/PhysRevLett.122.065301](https://doi.org/10.1103/PhysRevLett.122.065301).



- 421 [15] S.-J. Ran, *Bayesian tensor network with polynomial complexity for probabilistic machine*  
422 *learning* (2019), [1912.12923v2](https://arxiv.org/abs/1912.12923v2).
- 423 [16] J. Martyn, G. Vidal, C. Roberts and S. Leichenauer, *Entanglement and tensor networks for*  
424 *supervised image classification* (2020), [2007.06082v1](https://arxiv.org/abs/2007.06082v1).
- 425 [17] S. Cheng, L. Wang and P. Zhang, *Supervised learning with projected entangled pair states*,  
426 *Physical Review B* **103**(12), 125117 (2021), doi:[10.1103/physrevb.103.125117](https://doi.org/10.1103/physrevb.103.125117).
- 427 [18] D. Liu, Z. Yao and Q. Zhang, *Quantum-classical machine learning by hybrid tensor networks*  
428 (2020), [2005.09428v1](https://arxiv.org/abs/2005.09428v1).
- 429 [19] Z.-F. Gao, S. Cheng, R.-Q. He, Z. Y. Xie, H.-H. Zhao, Z.-Y. Lu and T. Xiang, *Compressing*  
430 *deep neural networks by matrix product operators*, *Physical Review Research* **2**, 023300  
431 (2020), doi:[10.1103/PhysRevResearch.2.023300](https://doi.org/10.1103/PhysRevResearch.2.023300).
- 432 [20] P. Blagoveschensky and A. H. Phan, *Deep convolutional tensor network* (2020), [2005.](https://arxiv.org/abs/2005.14506v1)  
433 [14506v1](https://arxiv.org/abs/2005.14506v1).
- 434 [21] K. He, X. Zhang, S. Ren and J. Sun, *Identity mappings in deep residual networks*,  
435 In *Computer Vision – ECCV 2016*, pp. 630–645. Springer International Publishing,  
436 doi:[10.1007/978-3-319-46493-0\\_38](https://doi.org/10.1007/978-3-319-46493-0_38) (2016).
- 437 [22] K. He, X. Zhang, S. Ren and J. Sun, *Deep residual learning for image recognition*, In *2016*  
438 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778. IEEE,  
439 doi:[10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90) (2016).
- 440 [23] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov, *Dropout:*  
441 *A simple way to prevent neural networks from overfitting*, *Journal of Machine Learning*  
442 *Research* **15**(1), 1929–1958 (2014).
- 443 [24] S. Rendle, *Factorization machines*, In *2010 IEEE International Conference on Data Mining*.  
444 IEEE, doi:[10.1109/icdm.2010.127](https://doi.org/10.1109/icdm.2010.127) (2010).
- 445 [25] A. Novikov, M. Trofimov and I. V. Oseledets, *Exponential machines*, In *5th International*  
446 *Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017,*  
447 *Workshop Track Proceedings*. OpenReview.net (2017).
- 448 [26] H. B. Demuth, M. H. Beale, O. De Jess and M. T. Hagan, *Neural Network Design*, Martin  
449 Hagan, Stillwater, OK, USA, 2nd edn., ISBN 0971732116 (2014).
- 450 [27] I. V. Oseledets, *Tensor-train decomposition*, *SIAM Journal on Scientific Computing* **33**(5),  
451 2295 (2011), doi:[10.1137/090752286](https://doi.org/10.1137/090752286).
- 452 [28] D. Perez-García, F. Verstraete, M. M. Wolf and J. I. Cirac, *Matrix product*  
453 *state representations*, *Quantum Information and Computation* **7**(5-6), 401 (2007),  
454 doi:[10.5555/2011832.2011833](https://doi.org/10.5555/2011832.2011833).
- 455 [29] S. Efthymiou, J. Hidary and S. Leichenauer, *Tensor network for machine learning* (2019),  
456 [1906.06329v1](https://arxiv.org/abs/1906.06329v1).
- 457 [30] J. Gao, L.-F. Qiao, Z.-Q. Jiao, Y.-C. Ma, C.-Q. Hu, R.-J. Ren, A.-L. Yang, H. Tang, M.-H.  
458 Yung and X.-M. Jin, *Experimental machine learning of quantum states*, *Physical Review*  
459 *Letters* **120**, 240501 (2018), doi:[10.1103/PhysRevLett.120.240501](https://doi.org/10.1103/PhysRevLett.120.240501).
- 460 [31] X. Glorot, A. Bordes and Y. Bengio, *Deep sparse rectifier neural networks*, In *14th Inter-*  
461 *national Conference on Artificial Intelligence and Statistics*, vol. 15, pp. 315–323 (2011).

- 462 [32] A. F. Agarap, *Deep learning using rectified linear units (relu)* (2018), [1803.08375v2](#).
- 463 [33] P. Baldi and P. J. Sadowski, *Understanding dropout*, In C. J. C. Burges, L. Bottou,  
464 M. Welling, Z. Ghahramani and K. Q. Weinberger, eds., *Advances in Neural Information*  
465 *Processing Systems*, vol. 26, pp. 2814–2822. Curran Associates, Inc. (2013).
- 466 [34] W. Zaremba, I. Sutskever and O. Vinyals, *Recurrent neural network regularization* [1409.](#)  
467 [2329v5](#).
- 468 [35] L. van der Maaten and G. Hinton, *Visualizing data using t-sne*, *Journal of Machine*  
469 *Learning Research* **9**(86), 2579 (2008).
- 470 [36] L. van der Maaten, *Learning a parametric embedding by preserving local structure*, In  
471 D. van Dyk and M. Welling, eds., *Proceedings of the Twelfth International Conference on*  
472 *Artificial Intelligence and Statistics*, vol. 5 of *Proceedings of Machine Learning Research*,  
473 pp. 384–391. PMLR, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA  
474 (2009).
- 475 [37] Y. LECUN, *The mnist database of handwritten digits*, <http://yann.lecun.com/exdb/mnist/>  
476 .
- 477 [38] E. M. Stoudenmire, *Learning relevant features of data with multi-scale tensor net-*  
478 *works*, *Quantum Science and Technology* **3**(3), 034003 (2018), doi:[10.1088/2058-](#)  
479 [9565/aaba1a](#).
- 480 [39] K. Meshkini, J. Platos and H. Ghassemian, *An analysis of convolutional neural network*  
481 *for fashion images classification (fashion-mnist)*, In S. Kovalev, V. Tarasov, V. Snasel and  
482 A. Sukhanov, eds., *Proceedings of the Fourth International Scientific Conference “Intelligent*  
483 *Information Technologies for Industry” (IITI’19)*, pp. 85–95. Springer International Pub-  
484 lishing, Cham, ISBN 978-3-030-50097-9, doi:[10.1007/978-3-030-50097-9\\_10](#) (2020).
- 485 [40] T. G. Kolda and B. W. Bader, *Tensor decompositions and applications*, *SIAM Review* **51**(3),  
486 455 (2009), doi:[10.1137/07070111x](#).
- 487 [41] H. Xiao, K. Rasul and R. Vollgraf, *Fashion-mnist: a novel image dataset for benchmarking*  
488 *machine learning algorithms* (2017), [1708.07747v2](#).
- 489 [42] H. N. Phien, G. Vidal and I. P. McCulloch, *Infinite boundary conditions for*  
490 *matrix product state calculations*, *Physical Review B* **86**(24), 245107 (2012),  
491 doi:[10.1103/physrevb.86.245107](#).
- 492 [43] J. Miller, G. Rabusseau and J. Terilla, *Tensor networks for probabilistic sequence modeling*,  
493 In A. Banerjee and K. Fukumizu, eds., *Proceedings of The 24th International Conference on*  
494 *Artificial Intelligence and Statistics*, vol. 130 of *Proceedings of Machine Learning Research*,  
495 pp. 3079–3087. PMLR (2021).
- 496 [44] S. R. White, *Density matrix formulation for quantum renormalization groups*, *Physical*  
497 *Review Letters* **69**, 2863 (1992), doi:[10.1103/PhysRevLett.69.2863](#).
- 498 [45] M. Hardt and T. Ma, *Identity matters in deep learning* (2016), [1611.04231v3](#).
- 499 [46] L.-L. Huang, A. Shimizu, Y. Hagihara and H. Kobatake, *Face detection from clut-*  
500 *tered images using a polynomial neural network*, *Neurocomputing* **51**, 197 (2003),  
501 doi:[10.1016/s0925-2312\(02\)00616-1](#).

- 502 [47] D. P. Kingma and J. Ba, *Adam: A method for stochastic optimization*,  
503 doi:[10.48550/ARXIV.1412.6980](https://doi.org/10.48550/ARXIV.1412.6980) (2014).
- 504 [48] J. Eisert, *Entanglement and tensor network states*, *Modeling and Simulation* 3, 520 (2013)  
505 (2013), [1308.3318](https://arxiv.org/abs/1308.3318).