

# A multiscale neural network based on hierarchical matrices \*

Yuwei Fan<sup>†</sup>, Lin Lin<sup>‡</sup>, Lexing Ying<sup>§</sup>, and Leonardo Zepeda-Núñez<sup>¶</sup>

**Abstract.** In this work we introduce a new multiscale artificial neural network based on the structure of  $\mathcal{H}$ -matrices. This network generalizes the latter to the nonlinear case by introducing a local deep neural network at each spatial scale. Numerical results indicate that the network is able to efficiently approximate discrete nonlinear maps obtained from discretized nonlinear partial differential equations, such as those arising from nonlinear Schrödinger equations and the Kohn-Sham density functional theory.

**Key words.**  $\mathcal{H}$ -matrix; multiscale neural network; locally connected neural network; convolutional neural network

**AMS subject classifications.**

**1. Introduction.** In the past decades, there has been a great combined effort in developing efficient algorithms to solve linear problems issued from discretization of integral equations (IEs), and partial differential equations (PDEs). In particular, multiscale methods such as multi-grid methods [7], the fast multipole method [17], wavelets [38], and hierarchical matrices [19, 6], have been strikingly successful in reducing the complexity for solving such systems. In several cases, for operators of pseudo-differential type, these algorithms can achieve linear or quasi-linear complexity. In a nutshell, these methods aim to use the inherent multiscale structure of the underlying physical problem to build efficient representations at each scale, thus compressing the information contained in the system. The gains in complexity stem mainly from processing information at each scale, and merging it in a hierarchical fashion.

Even though these techniques have been extensively applied to linear problems with outstanding success, their application to nonlinear problems has been, to the best of our knowledge, very limited. This is due to the high complexity of the solution maps. In particular, building a global approximation of such maps would normally require an extremely large amount of parameters, which in return, is often translated to algorithms with a prohibitive computational cost. The development of algorithms and heuristics to reduce the cost is an area of active research [4, 18, 37, 15, 14]. However, in general, each method is application-

---

\*Submitted to the editors July 4, 2018.

**Funding:** The work of Y.F. and L.Y. is partially supported by the U.S. Department of Energys Advanced Scientific Computing Research program under award DE-FC02-13ER26134/DE-SC0009409 and the National Science Foundation under award DMS-1521830. The work of L.L. and L. Z. is partially supported by the Department of Energy under Grant No. DE-SC0017867 and the CAMERA project.

<sup>†</sup>Department of Mathematics, Stanford University, Stanford, CA 94305 ([ywfan@stanford.edu](mailto:ywfan@stanford.edu), <http://web.stanford.edu/~ywfan/>)

<sup>‡</sup>Department of Mathematics, University of California, Berkeley, and Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720 ([linlin@math.berkeley.edu](mailto:linlin@math.berkeley.edu), <https://math.berkeley.edu/~linlin/>)

<sup>§</sup>Department of Mathematics and Institute for Computational and Mathematical Engineering, Stanford University, Stanford, CA 94305 ([lexing@stanford.edu](mailto:lexing@stanford.edu), <https://web.stanford.edu/~lexing/>)

<sup>¶</sup>Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720 ([lzepeda@lbl.gov](mailto:lzepeda@lbl.gov), <http://math.mit.edu/~lzepeda/>)

30 dependent, and requires a deep understanding of the physics behind.

31 On the other hand, the surge of interest in machine learning methods, in particular, deep  
 32 neural networks, has dramatically improved speech recognition [25], visual object recognition  
 33 [31], object detection, etc. This has fueled breakthroughs in many domains such as drug dis-  
 34 covery [36], genomics [33], and automatic translation [45], among many others. Several recent  
 35 reviews include [32] and [43]. Deep neural networks have empirically shown that it is possible  
 36 to obtain efficient representations of very high-dimensional functions. Even though the uni-  
 37 versality theorem holds for neural networks [27, 39, 11], *i.e.*, they can approximate arbitrarily  
 38 well any function with mild regularity conditions, how to build such approximation efficiently  
 39 remains an open question. In particular, the degree of approximation depends dramatically on  
 40 the architecture of the neural network, *i.e.* how the different layers are connected. In addition,  
 41 there is a fine balance between the number of parameters, the architecture, and the degree of  
 42 approximation [22, 23, 39].

43 This paper aims to combine the tools developed in deep neural networks with ideas from  
 44 multiscale methods. In particular, we extend hierarchical matrices ( $\mathcal{H}$ -matrices) to nonlinear  
 45 problems within the framework of neural networks. Let

$$46 \quad (1.1) \quad u = \mathcal{M}(v), \quad u, v \in \Omega \subset \mathbb{R}^n,$$

47 be a nonlinear map, issued from an underlying physical problem, described by either an  
 48 integral equation or a partial differential equation, where  $v$  can be considered as a parameter  
 49 in the equation,  $u$  is either the solution of the equation or a function of it, and  $n$  is the number  
 50 of variables.

51 We build a neural network with a novel multiscale structure inspired by hierarchical ma-  
 52 trices. To this end, we interpret the application of an  $\mathcal{H}$ -matrix to a vector using a neural  
 53 network structure as follows. We first reduce the dimensionality of the vector, or *restrict* it,  
 54 by multiplying it by a structured short and wide matrix. Then we *process* the encoded vector  
 55 by multiplying it by a structured square matrix. Then we bring the vector to its original size,  
 56 or *interpolate* it, by multiplying it with a structured tall and skinny matrix. Such operations  
 57 are performed at different spatial scales separately. The contributions from all spatial scales,  
 58 together with the near-field contribution represented by a near-diagonal matrix, are added  
 59 together to obtain the final approximation to the matrix-vector multiplication. Since every  
 60 step is linear, the overall operation is also a linear mapping. For nonlinear problems, such  
 61 an interpretation allows us to directly generalize the  $\mathcal{H}$ -matrix by replacing the structured  
 62 square matrix in the processing stage by a structured nonlinear network with several layers.  
 63 The resulting artificial neural network, which we call *multiscale neural network*, maintains a  
 64 relatively modest amount of parameters even for large problems.

65 We demonstrate the performance of the multiscale neural network by approximating the  
 66 solution to the nonlinear Schrödinger equation [40, 2], as well as the Kohn-Sham map [26, 30].  
 67 These mappings are highly nonlinear, and yet can be well approximated by the proposed  
 68 neural network, with a relative accuracy on the order of  $10^{-4} \sim 10^{-3}$ . Furthermore, we do  
 69 not observe overfitting even in the presence of a relatively small set of training samples.

70 **1.1. Related work.** Although machine learning and deep learning literature is extremely  
 71 vast, the application of deep learning to numerical analysis problems is relatively new, though

72 that seems to be rapidly changing. Several works have used deep neural networks to solve  
 73 PDEs [5, 42, 9, 13, 44, 34]. A combination of machine learning and multiscale methods was  
 74 applied by Chan, and Elsheikh using a multiscale finite elements with a response trained us-  
 75 ing a neural network [8]. For general applications of machine learning to nonlinear numerical  
 76 analysis problems the work of Raissi and Karniadakis uses machine learning, in particular,  
 77 Gaussian processes, to find parameters in nonlinear equations [41]; Khoo, Lu and Ying used  
 78 neural network in the context of uncertainty quantification [28]; Zhang et al used neural net-  
 79 work in the context of generating high-quality interatomic potentials for molecular dynamics  
 80 [51, 50]. Wang et al. used multiscale model reduction methods coupled with deep neural  
 81 networks [47].

82 **1.2. Organization.** The remainder of the paper is organized as follows. [Section 2](#) re-  
 83 views the  $\mathcal{H}$ -matrices and interprets the  $\mathcal{H}$ -matrices using the framework of neural networks.  
 84 [Section 3](#) extends the neural network representation of  $\mathcal{H}$ -matrices to the nonlinear case. [Sec-](#)  
 85 [tion 4](#) discusses the implementation details and demonstrates the accuracy of the architecture  
 86 in representing nonlinear maps, followed by the conclusion and future directions in [section 5](#).

87 **2. Neural network architecture for  $\mathcal{H}$ -matrices.** In this section, we aim to represent the  
 88 matrix-vector multiplication of  $\mathcal{H}$ -matrices within the framework of neural networks. For  
 89 the sake of clarity, we succinctly review the structure of  $\mathcal{H}$ -matrices for the one dimensional  
 90 case in [subsection 2.1](#). We interpret  $\mathcal{H}$ -matrices using the framework of neural networks in  
 91 [subsection 2.2](#), and then extend it to the multi-dimensional case in [subsection 2.3](#).

92 **2.1.  $\mathcal{H}$ -matrices.** Hierarchical matrices ( $\mathcal{H}$ -matrices) were first introduced by Hackbusch  
 93 et al. in a series of papers [19, 21, 20] as an algebraic framework for representing matrices with  
 94 a hierarchically off-diagonal low-rank structure. This framework provides efficient numerical  
 95 methods for solving linear systems arising from integral equations (IE) and partial differential  
 96 equations (PDE) [6]. In the sequel, we follow the notation in [35] to provide a brief introduction  
 97 to the framework of  $\mathcal{H}$ -matrices in a simple uniform and Cartesian setting. The interested  
 98 readers are referred to [19, 6, 35] for further details.

99 Consider the integral equation

$$100 \quad (2.1) \quad u(x) = \int_{\Omega} g(x, y)v(y) dy, \quad \Omega = [0, 1),$$

101 where  $u$  and  $v$  are periodic in  $\Omega$  and  $g(x, y)$  is smooth and numerically low-rank away from  
 102 the diagonal. A discretization with an uniform grid with  $N = 2^L m$  discretization points yields  
 103 the linear system given by

$$104 \quad (2.2) \quad u = Av,$$

105 where  $A \in \mathbb{R}^{N \times N}$ , and  $u, v \in \mathbb{R}^N$  are the discrete analogues of  $u(x)$  and  $v(x)$  respectively.

106 We introduce a hierarchical dyadic decomposition of the grid in  $L + 1$  levels as follows.  
 107 We start by the 0-th level of the decomposition, which corresponds to the set of all grid points  
 108 defined as

$$109 \quad (2.3) \quad \mathcal{I}^{(0)} = \{k/N : k = 0, \dots, N - 1\}.$$

110 At each level  $\ell$  ( $0 \leq \ell \leq L$ ), we decompose the grid in  $2^\ell$  disjoint *segments*.

111 Each segment is defined by  $\mathcal{I}_i^{(\ell)} = \mathcal{I}^{(0)} \cap [(i-1)/2^\ell, i/2^\ell]$  for  $i = 1, \dots, 2^\ell$ . Throughout  
 112 this manuscript,  $\mathcal{I}^{(\ell)}$  (or  $\mathcal{J}^{(\ell)}$ ) denote a generic segment of a given level  $\ell$ , and the superscript  
 113  $\ell$  will be omitted when the level is clear from the context. Moreover, following the usual  
 114 terminology in  $\mathcal{H}$ -matrices, we say that a segment  $\mathcal{J}^{(\ell)}$  ( $\ell \geq 1$ ) is the *parent* of a segment  
 115  $\mathcal{I}^{(\ell-1)}$  if  $\mathcal{I}^{(\ell-1)} \subset \mathcal{J}^{(\ell)}$ . Symmetrically,  $\mathcal{I}^{(\ell-1)}$  is called a *child* of  $\mathcal{J}^{(\ell)}$ . Clearly, each segment,  
 116 except those on level  $L$ , have two child segments.

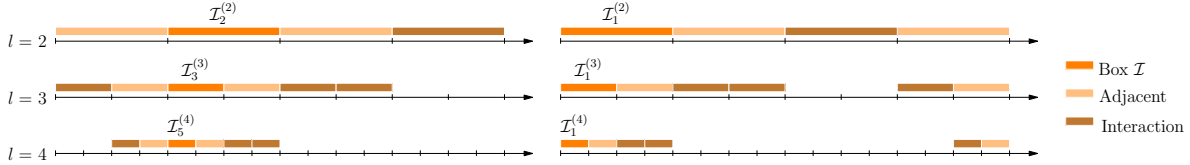


Figure 1: Illustration of the computational domain at level  $\ell = 2, 3, 4$ . The left and right figures represent an interior segment and a boundary segment and their adjacent and interaction list at different levels.

117 In addition, for a segment  $\mathcal{I}$  on level  $\ell \geq 2$ , we define the following lists:

118  $\text{NL}(\mathcal{I})$  *neighbor list* of  $\mathcal{I}$ . List of the segments on level  $\ell$  that are adjacent to  $\mathcal{I}$  including  $\mathcal{I}$   
 119 itself;

120  $\text{IL}(\mathcal{I})$  *interaction list* of  $\mathcal{I}$ . If  $\ell = 2$ ,  $\text{IL}(\mathcal{I})$  contains all the segments on level 2 minus  $\text{NL}(\mathcal{I})$ .

121 If  $\ell > 2$ ,  $\text{IL}(\mathcal{I})$  contains all the segments on level  $\ell$  that are children of segments in  
 122  $\text{NL}(\mathcal{P})$  with  $\mathcal{P}$  being  $\mathcal{I}$ 's parent minus  $\text{NL}(\mathcal{I})$ .

123 **Figure 1** illustrates the dyadic partition of the computational domain and the lists on  
 124 levels  $\ell = 2, 3, 4$ . Clearly,  $\mathcal{J} \in \text{NL}(\mathcal{I})$  if and only if  $\mathcal{I} \in \text{NL}(\mathcal{J})$ , and  $\mathcal{J} \in \text{IL}(\mathcal{I})$  if and only if  
 125  $\mathcal{I} \in \text{IL}(\mathcal{J})$ .

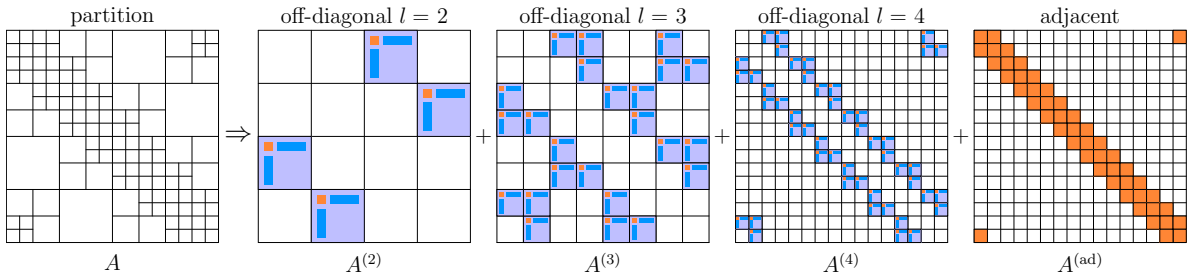


Figure 2: Partition of the matrix  $A$  for  $L = 4$  and nonzero blocks of  $A^{(\ell)}$ ,  $\ell = 2, 3, 4$  (colored blue) and  $A^{(\text{ad})}$  (colored orange). Nonzero blocks of  $A^{(\ell)}$  are approximated by low rank approximation and nonzero blocks of  $A^{(\text{ad})}$  are retained.

126 For a vector  $v \in \mathbb{R}^N$ ,  $v_{\mathcal{I}}$  denotes the elements of  $v$  indexed by  $\mathcal{I}$  and for a matrix  $A \in$   
 127  $\mathbb{R}^{N \times N}$ ,  $A_{\mathcal{I}, \mathcal{J}}$  represents the submatrix given by the elements of  $A$  indexed by  $\mathcal{I} \times \mathcal{J}$ . The  
 128 dyadic partition of the grid and the interaction lists induce a multilevel decomposition of  $A$

129 as follows

$$130 \quad (2.4) \quad A = \sum_{\ell=2}^L A^{(\ell)} + A^{(\text{ad})}, \quad \begin{aligned} A_{\mathcal{I},\mathcal{J}}^{(\ell)} &= \begin{cases} A_{\mathcal{I},\mathcal{J}}, & \mathcal{I} \in \text{IL}(\mathcal{J}); \\ 0, & \text{otherwise,} \end{cases} \quad \mathcal{I}, \mathcal{J} \text{ at level } \ell, 2 \leq \ell \leq L, \\ A_{\mathcal{I},\mathcal{J}}^{(\text{ad})} &= \begin{cases} A_{\mathcal{I},\mathcal{J}}, & \mathcal{I} \in \text{NL}(\mathcal{J}); \\ 0, & \text{otherwise,} \end{cases} \quad \mathcal{I}, \mathcal{J} \text{ at level } L. \end{aligned}$$

131 In a nutshell,  $A^{(\ell)}$  considers the interaction at level  $\ell$  between a segment and its interaction  
132 list, and  $A^{(\text{ad})}$  considers all the interactions between adjacent segments at level  $L$ .

133 **Figure 2** illustrates the block partition of  $A$  induced by the dyadic partition, and the  
134 decomposition induced by the different interaction lists at each level that follows (2.4).

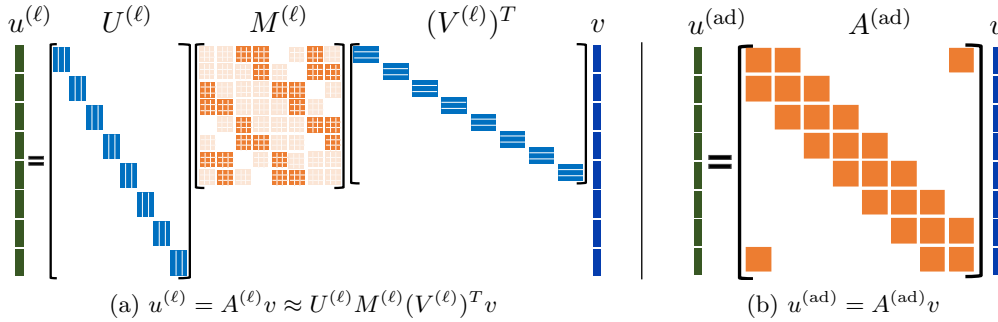


Figure 3: Diagram of matrix-vector multiplication (2.7) of the low-rank part and adjacent part of  $\mathcal{H}$ -matrices. The blocks of  $M^{(\ell)}$  colored by pale orange are zero blocks, and if we treat these blocks as non-zero blocks, the matrices  $M^{(\ell)}$  are block cyclic band matrices.

135 The key idea behind  $\mathcal{H}$ -matrices is to approximate the nonzero blocks  $A^{(\ell)}$  by a low rank  
136 approximation (see [29] for a thorough review). This idea is depicted in **Figure 2**, in which  
137 each non-zero block is approximated by a tall and skinny matrix, a small squared one and a  
138 short and wide one, respectively. In this work, we focus on the uniform  $\mathcal{H}$ -matrices [16], and,  
139 for simplicity, we suppose that each block has a fixed rank at most  $r$ , *i.e.*

$$140 \quad (2.5) \quad A_{\mathcal{I},\mathcal{J}}^{(\ell)} \approx U_{\mathcal{I}}^{(\ell)} M_{\mathcal{I},\mathcal{J}}^{(\ell)} (V_{\mathcal{J}}^{(\ell)})^T, \quad U_{\mathcal{I}}^{(\ell)}, V_{\mathcal{J}}^{(\ell)} \in \mathbb{R}^{N/2^\ell \times r}, \quad M_{\mathcal{I},\mathcal{J}}^{(\ell)} \in \mathbb{R}^{r \times r}.$$

141 where  $\mathcal{I}$ , and  $\mathcal{J}$  are any interacting segments at level  $\ell$ .

142 The main observation is that it is possible to factorize each  $A^{(\ell)}$  as  $A^{(\ell)} \approx U^{(\ell)}M^{(\ell)}(V^{(\ell)})^T$   
143 depicted in **Figure 3**.  $U^{(\ell)}$  is a block diagonal matrix with diagonal blocks  $U_{\mathcal{I}}^{(\ell)}$  for  $\mathcal{I}$  at level  
144  $\ell$ ,  $V^{(\ell)}$  is a block diagonal matrix with diagonal blocks  $V_{\mathcal{J}}^{(\ell)}$  for  $\mathcal{J}$  at level  $\ell$ , and finally  $M^{(\ell)}$   
145 aggregates all the blocks  $M_{\mathcal{I},\mathcal{J}}^{(\ell)}$  for all interacting segments  $\mathcal{I}, \mathcal{J}$  at level  $\ell$ . This factorization  
146 induces a decomposition of  $A$  given by

$$147 \quad (2.6) \quad A = \sum_{\ell=2}^L A^{(\ell)} + A^{(\text{ad})} \approx \sum_{\ell=2}^L U^{(\ell)}M^{(\ell)}(V^{(\ell)})^T + A^{(\text{ad})}.$$

148 Thus the matrix-vector multiplication (2.2) can be expressed as

$$149 \quad (2.7) \quad u \approx \sum_{\ell=2}^L u^{(\ell)} + u^{(\text{ad})} = \sum_{\ell=2}^L U^{(\ell)} M^{(\ell)} (V^{(\ell)})^T v + A^{(\text{ad})} v,$$

150 as illustrated in Figure 3, which constitutes the basis for writing  $\mathcal{H}$ -matrices as a neural  
151 network.

152 In addition, the matrices  $\{U^{(\ell)}, V^{(\ell)}, M^{(\ell)}\}_{\ell=2}^L$  and  $A^{(\text{ad})}$  have the following properties.

153 **Property 1.** *The matrices*

- 154 1.  $U^{(\ell)}$  and  $V^{(\ell)}$ ,  $\ell = 2, \dots, L$  are block diagonal matrices with block size  $N/2^\ell \times r$ ;
- 155 2.  $A^{(\text{ad})}$  is a block cyclic tridiagonal matrix with block size  $m \times m$ ;
- 156 3.  $M^{(\ell)}$ ,  $\ell = 2, \dots, L$  are block cyclic band matrices with block size  $r \times r$  and band size  
157  $n_b^{(\ell)}$ , which is 2 for  $\ell = 2$  and 3 for  $\ell \geq 3$ , if we treat all the pale orange colored blocks  
158 of  $M^{(\ell)}$  in Figure 3a as nonzero blocks.

159 We point out that the band sizes  $n_b^{(\ell)}$  and  $n_b^{(\text{ad})}$  depend on the definitions of NL and IL. In  
160 this case, the list were defined using the *strong admissible condition* in  $\mathcal{H}$ -matrices. Other  
161 conditions can be certainly used, such as the *weak admissibility condition*, leading to similar  
162 structures.

163 **2.2. Matrix-vector multiplication as a neural network.** An artificial neural network,  
164 in particular, a feed-forward network, can be thought of the composition of several simple  
165 functions, usually called *propagation functions*, in which the intermediate one-dimensional  
166 variables are called *neurons*, which in return, are organized in vector, or tensor, variables  
167 called *layers*. For example, an artificial feed-forward neural network

$$168 \quad (2.8) \quad u = \mathcal{F}(v), \quad u, v \in \mathbb{R}^n$$

169 with  $K + 1$  layer can be written using the following recursive formula

$$170 \quad (2.9) \quad \begin{aligned} \xi^{(0)} &= v, \\ \xi^{(k)} &= \phi(W^{(k)} \xi^{(k-1)} + b^{(k)}), \\ u &= \xi^{(K)}, \end{aligned}$$

171 where for each  $k = 1, \dots, K$  we have that  $\xi^{(k)}, b^{(k)} \in \mathbb{R}^{n_k}$ ,  $W^{(k)} \in \mathbb{R}^{n_k \times n_{k-1}}$ . Following the  
172 terminology of machine learning,  $\phi$  is called the *activation function* that is applied entry-wise,  
173  $W^{(k)}$  are the *weights*,  $b^{(k)}$  are the *biases*, and  $\xi^{(k)}$  is the  $k$ -th layer containing  $n_k$  number  
174 of *neurons*. Typical choices for the activation function are linear function, the rectified-  
175 linear unit (ReLU), or sigmoid function. In addition, (2.9) can easily be rewritten using  
176 tensors by replacing the matrix-vector multiplication by the more general tensor contraction.  
177 We point out that representing layers with tensors or vectors is equivalent up to reordering  
178 and reshaping. The main advantage of using the former is that layers, and its propagating  
179 functions, can be represented in a more compact fashion. Therefore, in what follows we  
180 predominantly use a tensor representation.

181 Within this context, training a network refers to finding the weights and biases, whose  
182 entries are collectively called *parameters*, in order to approximate a given map. This is usually  
183 done by minimizing a loss function using a stochastic optimization algorithm.

184 **2.2.1. Locally connected networks.** We interpret the structure of  $\mathcal{H}$ -matrices (2.6) using  
 185 the framework of neural networks. The different factors in (2.7) possess a distinctive structure,  
 186 which we aim to exploit by using locally connected (LC) network. LC networks are propagating  
 187 functions whose weights have a block-banded constraint. For the one-dimensional example,  
 188 we also treat  $\xi$  as a 2-tensor of dimensions  $\alpha \times N_x$ , where  $\alpha$  is the *channel dimension* and  
 189  $N_x$  is the *spatial dimension*, and  $\zeta$  be a 2-tensor of dimensions  $\alpha' \times N'_x$ . We say that  $\xi$  is  
 190 connected to  $\zeta$  by a LC networks if

$$191 \quad (2.10) \quad \zeta_{c',i} = \phi \left( \sum_{j=(i-1)s+1}^{(i-1)s+w} \sum_{c=1}^{\alpha} W_{c',c;i,j} \xi_{c,j} + b_{c',i} \right), \quad i = 1, \dots, N'_x, \quad c' = 1, \dots, \alpha',$$

192 where  $w$  and  $s \in \mathbb{N}$  are the *kernel window size* and *stride*, respectively. In addition, we say  
 193 that  $\zeta$  is a *locally connected* (LC) layer if it satisfies (2.10).

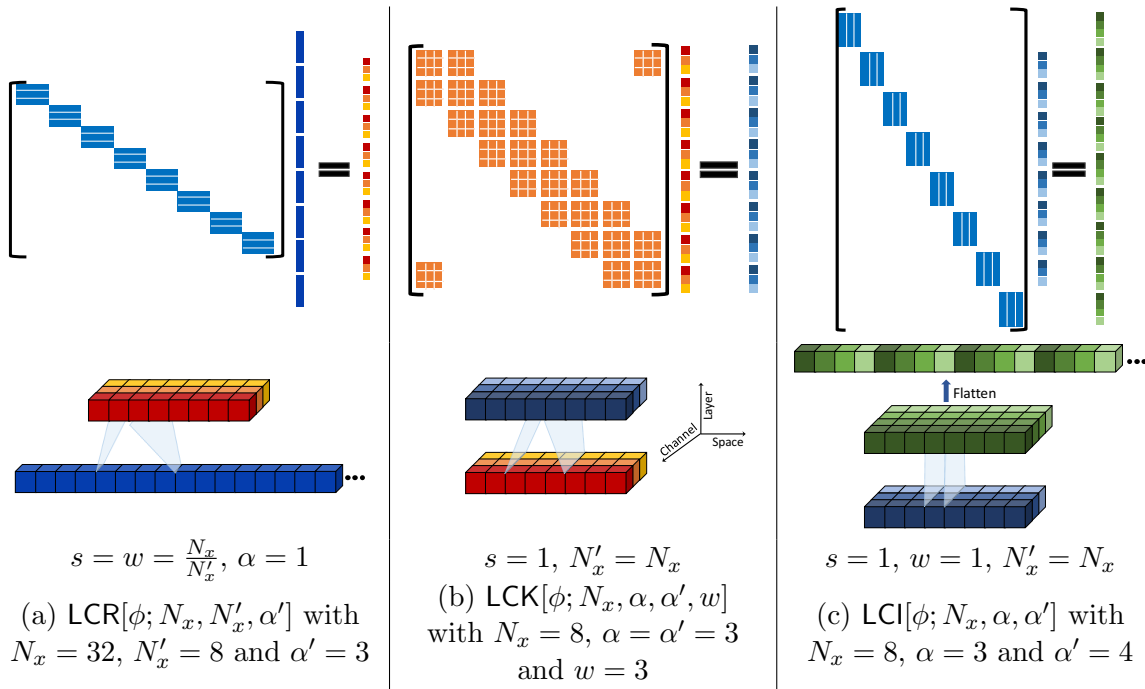


Figure 4: Three instances of locally connected networks used the represent the matrix-vector multiplication in (2.7). The upper portions of each column depicts the patterns of the matrices and the lower portions are their respective analogues using locally connect networks.

194 Each LC network requires 6 parameters,  $N_x, \alpha, N'_x, \alpha', w$  and  $s$  to be characterized. Next,  
 195 we define three types of LC network by specifying some of their parameters,

196 **LCR Restriction network:** we set  $s = w = \frac{N_x}{N'_x}$  and  $\alpha = 1$  in LC. This network represents the  
 197 multiplication of a block diagonal matrix with block sizes  $\alpha' \times s$  and a vector with size  
 198  $N_x \alpha$ , as illustrated by Figure 4 (a). We denote this network using  $\text{LCR}[\phi; N_x, N'_x, \alpha']$ .  
 199 The application of  $\text{LCR}[\text{linear}; 32, 8, 3]$  is depicted in Figure 4 (a).

200 LCK *Kernel* network: we set  $s = 1$  and  $N'_x = N_x$ . This network represents the multipli-  
 201 cation of a cyclic block band matrix of block size  $\alpha' \times \alpha$  and band size  $\frac{w-1}{2}$  times  
 202 a vector of size  $N_x\alpha$ , as illustrated by the upper portion of [Figure 4b](#). To account  
 203 for the periodicity we pad the input layer  $\xi_{c,j}$  on the spatial dimension to the size  
 204  $(N_x + w - 1) \times \alpha$ . We denote this network by  $\text{LCK}[\phi; N_x, \alpha, \alpha', w]$ . This network has  
 205 two steps: the periodic padding of  $\xi_{c,j}$  on the spatial dimension, and the application  
 206 of (2.10). The application of  $\text{LCK}[\text{linear}; 8, 3, 3, 3]$  is depicted in [Figure 4 \(b\)](#).  
 207 LCI *Interpolation* network: we set  $s = w = 1$  and  $N'_x = N_x$  in LC. This network represents  
 208 the multiplication of a block diagonal matrix with block size  $\alpha' \times \alpha$ , times a vector of  
 209 size  $N_x\alpha$ , as illustrated by the upper figure in [Figure 4 \(c\)](#). We denote the network  
 210  $\text{LCI}[\phi; N_x, \alpha, \alpha']$ , which has two steps: the application of (2.10), and the reshaping of  
 211 the output to a vector by column major indexing. The application of  $\text{LCI}[\text{linear}; 8, 3, 4]$   
 212 is depicted in [Figure 4 \(c\)](#).

213 **2.2.2. Neural network representation.** Following (2.7), in order to construct the neural  
 214 network (NN) architecture for (2.7), we need to represent the following four operations:

$$215 \quad (2.11a) \quad \xi^{(\ell)} = (V^{(\ell)})^T v,$$

$$216 \quad (2.11b) \quad \zeta^{(\ell)} = M^{(\ell)} \xi^{(\ell)},$$

$$217 \quad (2.11c) \quad u^{(\ell)} = U^{(\ell)} \zeta^{(\ell)},$$

$$218 \quad (2.11d) \quad u^{(\text{ad})} = A^{(\text{ad})} v.$$

220 From [Property 1.1](#) and the definition of LCR and LCI, the operations (2.11a) and (2.11c) are  
 221 equivalent to

$$222 \quad (2.12) \quad \xi^{(\ell)} = \text{LCR}[\text{linear}; N, 2^\ell, r](v), \quad u^{(\ell)} = \text{LCI}[\text{linear}; 2^\ell, r, \frac{N}{2^\ell}](\zeta^{(\ell)}),$$

223 respectively. Analogously, [Property 1.3](#) indicates that (2.11b) is equivalent to

$$224 \quad (2.13) \quad \zeta^{(\ell)} = \text{LCK}[\text{linear}; 2^\ell, r, r, 2n_b^{(\ell)} + 1](\xi^{(\ell)}).$$

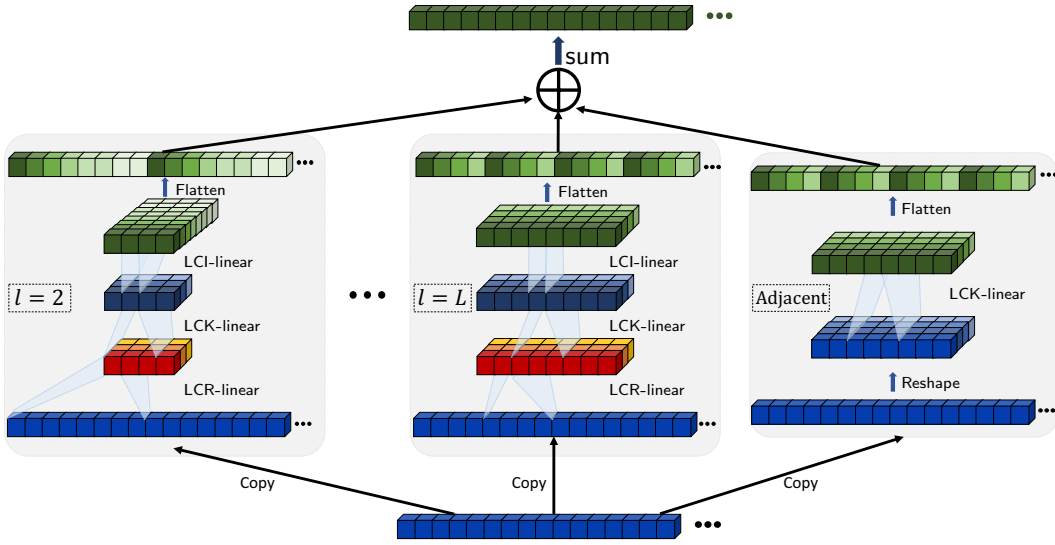
225 We point out that  $\xi^{(\ell)}$  is a vector in (2.11c) but a 2-tensor in (2.12) and (2.13). In principle,  
 226 we need to flatten  $\xi^{(\ell)}$  in (2.12) to a vector and reshape it back to a 2-tensor before (2.13).  
 227 These operations do not alter the algorithmic pipeline, so they are omitted.

228 Given that  $v, u^{(\text{ad})}$  are vectors, but LCK is defined for 2-tensors, we explicitly write the  
 229 reshape and flatten operations. Denote as  $\text{Reshape}[n_1, n_2]$  the map that reshapes a vector of  
 230 size  $n_1 n_2$  into a 2-tensor of size  $n_1 \times n_2$  by column major indexing, and  $\text{Flatten}$  is defined as  
 231 the inverse of  $\text{Reshape}$ . Using [Property 1.2](#), we can write (2.11d) as

$$232 \quad (2.14) \quad \tilde{v} = \text{Reshape}[m, 2^L](v), \quad \tilde{u}^{(\text{ad})} = \text{LCK} \left[ \text{linear}; 2^L, m, m, 2n_b^{(\text{ad})} + 1 \right] (\tilde{v}), \quad u^{(\text{ad})} = \text{Flatten}(\tilde{u}^{(\text{ad})}).$$

233 Combining (2.12), (2.13) and (2.14), we obtain [Algorithm 1](#), whose architecture is illus-  
 234 trated in [Figure 5](#). In particular, [Figure 5](#) is the translation to the neural network framework  
 235 of (2.7) (see [Figure 3](#)) using the building blocks depicted in [Figure 4](#).




 Figure 5: Neural network architecture for  $\mathcal{H}$ -matrices.

**Algorithm 1** Application of the NN representation of an  $\mathcal{H}$ -matrix to a vector  $v \in \mathbb{R}^N$ .

- |   |   |
|---|---|
| 1: $u = 0$ ;<br>2: <b>for</b> $\ell = 2$ to $L$ <b>do</b><br>3: $\xi = \text{LCR}[\text{linear}; N, 2^\ell, r](v)$ ;<br>4: $\zeta = \text{LCK}[\text{linear}; 2^\ell, r, r, 2n_b^{(\ell)} + 1](\xi)$ ;<br>5: $u = u + \text{LCI}[\text{linear}; 2^\ell, r, \frac{N}{2^\ell}](\zeta)$ ;<br>6: <b>end for</b> | 7: $\tilde{v} = \text{Reshape}[m, 2^L](v)$ ;<br>8: $\tilde{u}^{(\text{ad})} = \text{LCK}[\text{linear}; 2^L, m, m, 2n_b^{(\text{ad})} + 1](\tilde{v})$ ;<br>9: $u^{(\text{ad})} = \text{Flatten}(\tilde{u}^{(\text{ad})})$ ;<br>10: $u = u + u^{(\text{ad})}$ ;<br> |
|---|---|

236 Moreover, the memory footprints of the neural network architecture and  $\mathcal{H}$ -matrices are  
 237 asymptotically the same with respect the spatial dimension of  $u$  and  $v$ . This can be readily  
 238 shown by computing the total number of parameters. For the sake of simplicity, we only count  
 239 the parameters in the weights, ignoring those in the biases. A direct calculation yields the  
 240 number of parameters in LCR, LCK and LCI:

$$241 \quad (2.15) \quad N_p^{\text{LCR}} = N_x \alpha', \quad N_p^{\text{LCK}} = N_x \alpha \alpha' w, \quad N_p^{\text{LCI}} = N_x \alpha \alpha',$$

242 respectively. Hence, the number of parameters in [Algorithm 1](#) is

$$243 \quad (2.16) \quad \begin{aligned} N_p^{\mathcal{H}} &= \sum_{\ell=2}^L \left( N r + 2^\ell r^2 (2n_b^{(\ell)} + 1) + 2^\ell r \frac{N}{2^\ell} \right) + 2^L m^2 (2n_b^{(\text{ad})} + 1) \\ &\leq 2LNr + 2^{L+1} r^2 (2 \max_{\ell=2}^L n_b^{(\ell)} + 1) + Nm(2n_b^{(\text{ad})} + 1) \\ &\leq 2N \log(N)r + 3Nm(2n_b + 1) \sim O(N \log(N)), \end{aligned}$$

244 where  $n_b = \max(n_b^{(\text{ad})}, n_b^{(\ell)}, \ell = 2, \dots, L)$ , and  $r \leq m$  is used.

245 **2.3. Multi-dimensional case.** Following the previous section, the extension of [Algorithm 1](#)  
 246 to the  $d$ -dimensional case can be easily deduced using the tensor product of one-dimensional  
 247 cases. Consider  $d = 2$  below for instance, and the generalization to higher dimensional case  
 248 will be straight-forward. Suppose that we have an IE in 2D given by

$$249 \quad (2.17) \quad u(x) = \int_{\Omega} g(x, y)v(y) dy, \quad \Omega = [0, 1) \times [0, 1),$$

250 we discretize the domain  $\Omega$  with a uniform grid with  $n = N^2$  ( $N = 2^L m$ ) discretization points,  
 251 and let  $A$  be the resulting matrix obtained from discretizing (2.17). We denote the set of all  
 252 grid points as

$$253 \quad (2.18) \quad \mathcal{I}^{(d,0)} = \{(k_1/N, k_2/N) : k_1, k_2 = 0, \dots, N-1\}.$$

254 Clearly  $\mathcal{I}^{(d,0)} = \mathcal{I}^{(0)} \otimes \mathcal{I}^{(0)}$ , where  $\mathcal{I}^{(0)}$  is defined in (2.3), and  $\otimes$  is tensor product. At each  
 255 level  $\ell$  ( $0 \leq \ell \leq L$ ), we decompose the grid in  $4^\ell$  disjoint boxes as  $\mathcal{I}_i^{(d,\ell)} = \mathcal{I}_{i_1}^{(\ell)} \otimes \mathcal{I}_{i_2}^{(\ell)}$  for  
 256  $i_1, i_2 = 1, \dots, 2^\ell$ . The definition of the lists NL and IL can be easily extended. For each box  $\mathcal{I}$ ,  
 257 NL( $\mathcal{I}$ ) contains 3 boxes for 1D case,  $3^2$  boxes for 2D case. Similarly, the decomposition (2.4)  
 258 on the matrix  $A$  can easily be extended for this case. Following the structure of  $\mathcal{H}$ -matrices,  
 259 the off-diagonal blocks of  $A^{(\ell)}$  can be approximated as

$$260 \quad (2.19) \quad A_{\mathcal{I},\mathcal{J}}^{(\ell)} \approx U_{\mathcal{I}}^{(\ell)} M_{\mathcal{I},\mathcal{J}}^{(\ell)} (V_{\mathcal{J}}^{(\ell)})^T, \quad \mathcal{I}, \mathcal{J} \in \mathcal{I}^{(\ell)}, \quad U_{\mathcal{I}}^{(\ell)}, V_{\mathcal{J}}^{(\ell)} \in \mathbb{R}^{(N/2^\ell)^2 \times r}, M_{\mathcal{I},\mathcal{J}}^{(\ell)} \in \mathbb{R}^{r \times r}.$$

261 As mentioned before, we can describe the network using tensor or vectors. In what follows  
 262 we will switch between representations in order to illustrate the concepts in a compact fashion.  
 263 We denote an entry of a tensor  $T$  by  $T_{i,j}$ , where  $i$  is 2-dimensional index  $i = (i_1, i_2)$ . Using  
 264 the tensor notations,  $U^{(\ell)}, V^{(\ell)}$  in (2.7), can be treated as 3-tensors of dimension  $N \times N \times r$ .  
 265 We generalize the notion of band matrix to *band tensors*. A band tensor  $T$  satisfies that

$$266 \quad (2.20) \quad T_{i,j} = 0, \quad \text{if } |i_1 - j_1| > n_{b,1} \quad \text{or} \quad |i_2 - j_2| > n_{b,2},$$

267 where  $n_b = (n_{b,1}, n_{b,2})$  is called the band size for tensor. Thus [Property 1](#) can be generalized  
 268 to tensors yielding the following properties.

- 269 **Property 2.** 1. The 3-tensor  $U^{(\ell)}$  and  $V^{(\ell)}$ ,  $\ell = 2, \dots, L$  are block diagonal tensors  
 270 with block size  $N/2^\ell \times N/2^\ell \times r$ ;  
 271 2. the 4-tensor  $A^{(\text{ad})}$  is a block band cyclic tensor with block size  $m \times m \times m \times m$  and  
 272 band size  $n_b^{(\text{ad})} = (1, 1)$ ;  
 273 3. the 4-tensors  $M^{(\ell)}$ ,  $\ell = 2, \dots, L$  are block band cyclic tensor with block size  $r \times r$  and  
 274 band size  $n_b^{(\ell)}$ , which is (2, 2) for  $\ell = 2$  and (3, 3) for  $\ell \geq 3$ .

275 Next, we characterize LC networks for the 2D case. An NN layer for 2D can be represented  
 276 by a 3-tensor of size  $\alpha \times N_{x,1} \times N_{x,2}$ , in which  $\alpha$  is the channel dimension and  $N_{x,1}, N_{x,2}$  are  
 277 the spatial dimensions. If a layer  $\xi$  with size  $\alpha \times N_{x,1} \times N_{x,2}$  is connected to a locally connected  
 278 layer  $\zeta$  with size  $\alpha' \times N'_{x,1} \times N'_{x,2}$ , then

$$279 \quad (2.21) \quad \zeta_{c',i} = \phi \left( \sum_{j=(i-1)s+1}^{(i-1)s+w} \sum_{c=1}^{\alpha} W_{c',c;i,j} \xi_{c,j} + b_{c',i} \right), \quad i_1 = 1, \dots, N'_{x,1}, i_2 = 1, \dots, N'_{x,2}, c' = 1, \dots, \alpha',$$

280 where  $(i-1)s = ((i_1-1)s_1, (i_2-1)s_2)$ . As in the 1D case, the channel dimension corresponds  
 281 to the rank  $r$ , and the spatial dimensions correspond to the grid points of the discretized  
 282 domain. Analogously to the 1D case, we define the LC networks LCR, LCK and LCI and  
 283 use them to express the four operations in (2.11) which constitute the building blocks of the  
 284 neural network. The extension is trivial, the parameters  $N_x$ ,  $s$  and  $w$  in the one-dimensional  
 285 LC networks are replaced by their 2-dimensional counterpart  $N_x = (N_{x,1}, N_{x,2})$ ,  $s = (s_1, s_2)$   
 286 and  $w = (w_1, w_2)$ , respectively. We point out that  $s = w = \frac{N_x}{N_x}$  for the 1D case is replaced by  
 287  $s_j = w_j = \frac{N_{x,j}}{N_{x,j}'} , j = 1, 2$  for the 2D case in the definition of LC.

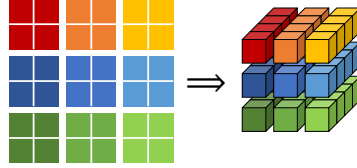


Figure 6: Diagram of Reshape[2<sup>2</sup>, 3, 3] in Algorithm 2.

288 Using the notations above we extend Algorithm 1 to the 2D case in Algorithm 2. We  
 289 crucially remark that the Reshape[ $r^2, n_1, n_2$ ] function in Algorithm 2 is not the usual major  
 290 column based reshaping. It reshapes a 2-tensor  $T$  with size  $rn_1 \times rn_2$  to a 3-tensor  $S$  with  
 291 size  $r^2 \times n_1 \times n_2$ , by treating the former as a block tensor with block size  $r \times r$ , and reshaping  
 292 each block as a vector following the formula  $S(k, i, j) = T((i-1)r + k_1, (j-1)r + k_2)$  with  
 293  $k = (k_1 - 1)r + k_2$ , for  $k_1, k_2 = 1, \dots, r$ ,  $i = 1, \dots, n_1$  and  $j = 1, \dots, n_2$ . Figure 6 provides an  
 294 example for the case Reshape[2<sup>2</sup>, 3, 3]. The Flatten is its inverse.

---

**Algorithm 2** Application of NN architecture for  $\mathcal{H}$ -matrices on a vector  $v \in \mathbb{R}^{N^2}$ .

---

- 1:  $u = 0$ ;
  - 2: **for**  $\ell = 2$  to  $L$  **do**
  - 3:    $\xi = \text{LCR}[\text{linear}; (N, N), (2^\ell, 2^\ell), r](v)$ ;
  - 4:    $\zeta = \text{LCK}[\text{linear}; (2^\ell, 2^\ell), r, r, (2n_{b,1}^{(\ell)} + 1, 2n_{b,2}^{(\ell)} + 1)](\xi)$ ;
  - 5:    $u = u + \text{LCI}[\text{linear}; (2^\ell, 2^\ell), r, (\frac{N}{2^\ell})^2](\zeta)$ ;
  - 6: **end for**
  - 7:  $\tilde{v} = \text{Reshape}[m^2, 2^L, 2^L](v)$ ;
  - 8:  $\tilde{u}^{(\text{ad})} = \text{LCK}[\text{linear}; (2^L, 2^L), m^2, m^2, (2n_{b,1}^{(\text{ad})} + 1, 2n_{b,2}^{(\text{ad})} + 1)](\tilde{v})$ ;
  - 9:  $u^{(\text{ad})} = \text{Flatten}(\tilde{u}^{(\text{ad})})$ ;
  - 10:  $u = u + u^{(\text{ad})}$ ;
- 

295 **3. Multiscale neural network.** In this section, we extend the aforementioned NN archi-  
 296 tecture to represent a general nonlinear mapping of the form

$$297 \quad (3.1) \quad u = \mathcal{M}(v), \quad u, v \in \mathbb{R}^{N^d},$$

298 Due to its multiscale structure, we refer to the resulting NN architecture as the *multiscale*  
 299 *neural network* (MNN). We consider the one-dimensional case below for simplicity, and the

300 generalization to higher dimensions follows directly as in [subsection 2.3](#).

---

**Algorithm 3** Application of multiscale neural network to a vector  $v \in \mathbb{R}^N$ .

---

1: $u = 0$ ; 2: <b>for</b> $\ell = 2$ to $L$ <b>do</b> 3: $\xi_0 = \text{LCR}[\text{linear}; N, 2^\ell, r](v)$ ; 4: <b>for</b> $k = 1$ to $K$ <b>do</b> 5: $\xi_k = \text{LCK}[\phi; 2^\ell, r, r, 2n_b^{(\ell)} + 1](\xi_{k-1})$ ; 6: <b>end for</b> 7: $u = u + \text{LCI}[\text{linear}; 2^\ell, r, \frac{N}{2^\ell}](\xi_K)$ ; 8: <b>end for</b>	9: $\xi_0 = \text{Reshape}[m, 2^L](v)$ ; 10: <b>for</b> $k = 1$ to $K$ <b>do</b> 11: $\xi_k = \text{LCK}[\phi; 2^L, m, m, 2n_b^{(\text{ad})} + 1](\xi_{k-1})$ ; 12: <b>end for</b> 13: $u^{(\text{ad})} = \text{Flatten}(\xi_K)$ ; 14: $u = u + u^{(\text{ad})}$ ; 
---	---

---

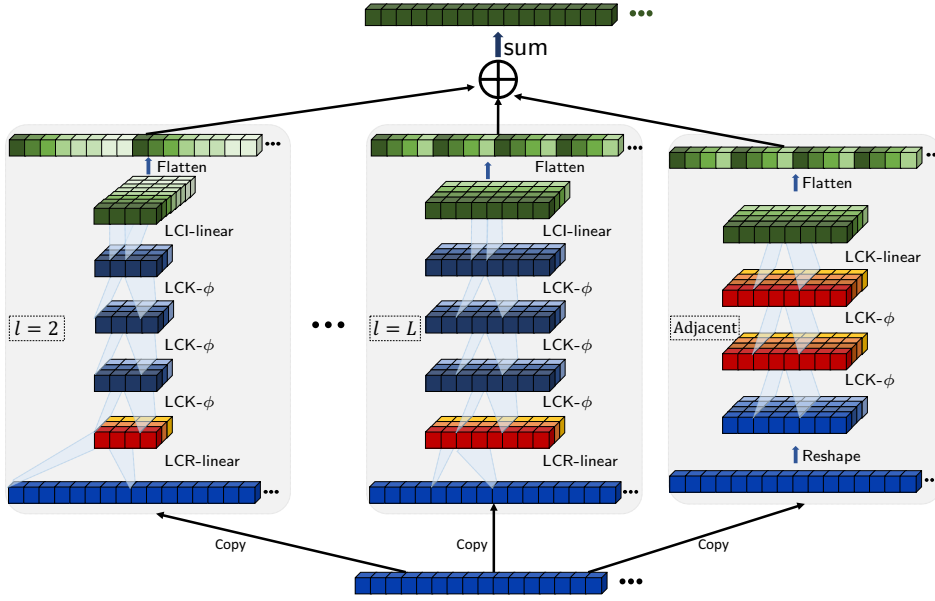


Figure 7: Multiscale neural network architecture for nonlinear mappings, which is an extension of the neural network architecture for  $\mathcal{H}$ -matrices [Figure 5](#).  $\phi$  is an activation function.

301 **3.1. Algorithm and architecture.** NN can represent nonlinearity by choosing the acti-  
302 vation function,  $\phi$ , to be nonlinear, such as ReLU or sigmoid. The range of the activation  
303 function also imposes constraints on the output of the NN. For example, the range of “ReLU”  
304 in  $[0, \infty)$  and the range of the sigmoid function is  $[0, 1]$ . Thus, the last layer is often chosen  
305 to be a linear layer to relax such constraint. [Algorithm 1](#) is then revised to [Algorithm 3](#), and  
306 the architecture is illustrated in [Figure 7](#). We remark that the nonlinear activation function  
307 is only used in the LCK network. The LCR and LCI networks in [Algorithm 1](#) are still treated  
308 as restriction and interpolation operations between coarse grid and fine grid, respectively, so  
309 we use linear activation functions in these layers. Particularly, we also use linear activation  
310 function for the last layer of the adjacent part, *i.e.* the  $\phi$  in line 11 in [Algorithm 3](#) is linear  
311 when  $k = K$ .

312 As in the linear case, we calculate the number of parameters of MNN and obtain (neglecting  
313 the number of parameters in  $b$  in (2.10))

$$\begin{aligned}
314 \quad (3.2) \quad N_p^{\text{MNN}} &= \sum_{\ell=2}^L \left( Nr + K2^\ell r^2(2n_b^{(\ell)} + 1) + 2^\ell r \frac{N}{2^\ell} \right) + K2^L m^2(2n_b^{(\text{ad})} + 1) \\
&\leq 2LNr + K2^{L+1}r^2(2 \max_{\ell=2}^L n_b^{(\ell)} + 1) + NKm(2n_b^{(\text{ad})} + 1) \\
&\leq 2N \log(N)r + 3NKm(2n_b + 1).
\end{aligned}$$

315 **3.2. Translation-invariant case.** For the linear case (2.1), if the kernel is *translation in-*  
316 *variant*, i.e.  $g(x, y) = g(x - y)$ , then the matrix  $A$  is a Toeplitz matrix. Then the matrices  $M^{(\ell)}$   
317 and  $A^{(\text{ad})}$  are Toeplitz matrices and all matrix blocks of  $U^{(\ell)}$  (resp.  $V^{(\ell)}$ ) can be represented  
318 by the same matrix. This leads to the *convolutional neural network* (CNN) as

$$319 \quad (3.3) \quad \zeta_{c', i} = \phi \left( \sum_{j=(i-1)s+1}^{(i-1)s+w} \sum_{c=1}^{\alpha} W_{c', c; j} \zeta_{c, j} + b_{c'} \right), \quad i = 1, \dots, N'_x, \quad c' = 1, \dots, \alpha'.$$

320 Compared to the LC network, the only difference is that the parameters  $W$  and  $b$  are indepen-  
321 dent of  $i$ . Hence, inheriting the definition of LCR, LCK and LCI, we define the layers CR, CK  
322 and CI, respectively. By replacing the LC layers in Algorithm 1 by the corresponding CNN  
323 layers, we obtain the neural network architecture for the translation invariant kernel.

324 For the nonlinear case, the translation invariant kernel for the linear case can be extended  
325 to kernels that are *equivariant to translation*, i.e. for any translation  $\mathcal{T}$ ,

$$326 \quad (3.4) \quad \mathcal{T}\mathcal{M}(v) = \mathcal{M}(\mathcal{T}v).$$

327 For this case, all the LC layers in Algorithm 3 can be replaced by its corresponding CNN  
328 layers. The number of parameters of CR, CK and CI are

$$329 \quad (3.5) \quad N_p^{\text{CR}} = \frac{N_x}{N'_x} \alpha', \quad N_p^{\text{CK}} = \alpha \alpha' w, \quad N_p^{\text{CI}} = \alpha \alpha'.$$

330 Thus, the number of parameters in Algorithm 3 using CNN is

$$\begin{aligned}
331 \quad (3.6) \quad N_{p, \text{CNN}}^{\text{MNN}} &= \sum_{\ell=2}^L \left( r \frac{N}{2^\ell} + Kr^2(2n_b^{(\ell)} + 1) + r \frac{N}{2^\ell} \right) + Km^2(2n_b^{(\text{ad})} + 1) \\
&\leq rN + (r^2 \log(N) + m^2)(2n_b + 1)K.
\end{aligned}$$

332 **4. Numerical results.** In this section we discuss the implementation details of MNN.  
333 We demonstrate the accuracy of the MNN architecture using two nonlinear problems: the  
334 nonlinear Schrödinger equation (NLSE), and the Kohn-Sham map (KS map) in the Kohn-  
335 Sham density functional theory.

336 **4.1. Implementation.** Our implementation of MNN uses Keras [10], a high-level applica-  
 337 tion programming interface (API) running, in this case, on top of TensorFlow [1] (a library  
 338 of toolboxes for training neural networks). The loss function is chosen as the mean squared  
 339 relative error, in which the relative error is defined with respect to  $\ell^2$  norm as

$$340 \quad (4.1) \quad \epsilon = \frac{\|u - u_{NN}\|_{\ell^2}}{\|u\|_{\ell^2}},$$

341 where  $u$  is the target solution generated by a numerical discretization of the PDEs and  $u_{NN}$   
 342 is the predicted solution by MNN. The optimization is performed using the NAdam optimizer  
 343 [12]. The weights and biases in MNN are initialized randomly from the normal distribution  
 344 and the batch size is always set between 1/100th and 1/50th of the number of train samples.

345 In all the tests, the band size is chosen as  $n_{b,ad} = 1$  and  $n_b^{(l)}$  is 2 for  $l = 2$  and 3 otherwise.  
 346 The nonlinear activation function is chosen as ReLU. All the test are run on GPU with data  
 347 type `float32`. All the numerical results are the best results by repeating the training a few  
 348 times, using different random seeds. The selection of parameters  $r$  (number of channels),  $L$   
 349 ( $N = 2^L m$ ) and  $K$  (number of layers in Algorithm 3) are problem dependent.

350 **4.2. NLSE with inhomogeneous background potential.** The nonlinear Schrödinger equa-  
 351 tion (NLSE) is widely used in quantum physics to describe the single particle properties of  
 352 the Bose-Einstein condensation phenomenon [40, 2]. Here we study the NLSE with inhom-  
 353 geneous background potential  $V(x)$

$$354 \quad (4.2) \quad \begin{aligned} & -\Delta u(x) + V(x)u(x) + \beta u(x)^3 = Eu(x), \quad x \in [0, 1]^d, \\ & \text{s.t. } \int_{[0,1]^d} u(x)^2 dx = 1, \text{ and } \int_{[0,1]^d} u(x) dx > 0, \end{aligned}$$

355 with period boundary condition, to find its ground state  $u_G(x)$ . We take a strongly nonlinear  
 356 case  $\beta = 10$  in this work and thus consider a defocusing cubic Schrödinger equation. Due  
 357 to the cubic term, an iterative method is required to solve (4.2) numerically. We employ the  
 358 method in [3] for the numerical solution, which solves a time-dependent NLSE by a normalized  
 359 gradient flow. The MNN is used to learn the map from the background potential to the ground  
 360 state

$$361 \quad (4.3) \quad V(x) \rightarrow u_G(x).$$

362 This map is equivariant to translation, and thus MNN is implemented using the CNN layers.  
 363 In the following, we study MNN on 1D and 2D cases, respectively.

364 The potential  $V$  is chosen as

$$365 \quad (4.4) \quad V(x) = - \sum_{i=1}^{n_g} \sum_{j_1, \dots, j_d = -\infty}^{\infty} \frac{\rho^{(i)}}{\sqrt{2\pi T}} \exp\left(-\frac{|x - j - c^{(i)}|^2}{2T}\right),$$

366 where the periodic summation imposes periodicity on the potential, and the parameters  $\rho^{(i)} \sim$   
 367  $\mathcal{U}(1, 4)$ ,  $c^{(i)} \sim \mathcal{U}(0, 1)^d$ ,  $i = 1, \dots, n_g$  and  $T \sim \mathcal{U}(2, 4) \times 10^{-3}$ .

368 **4.2.1. One-dimensional case.** For the one-dimensional case, the number of discretization  
 369 points is  $N = 320$ , and we set  $L = 6$  and  $m = \frac{N}{2L} = 5$ . In all the tests, the number of test  
 370 samples is the same as that the number of train samples if not properly specified. We perform  
 371 numerical experiments to study the behavior of MNN for different number of channels  $r$ ,  
 372 different number of CK layers  $K$ , different number of Gaussians  $n_g$  and different number of  
 373 training samples  $N_{\text{samples}}^{\text{train}}$ .

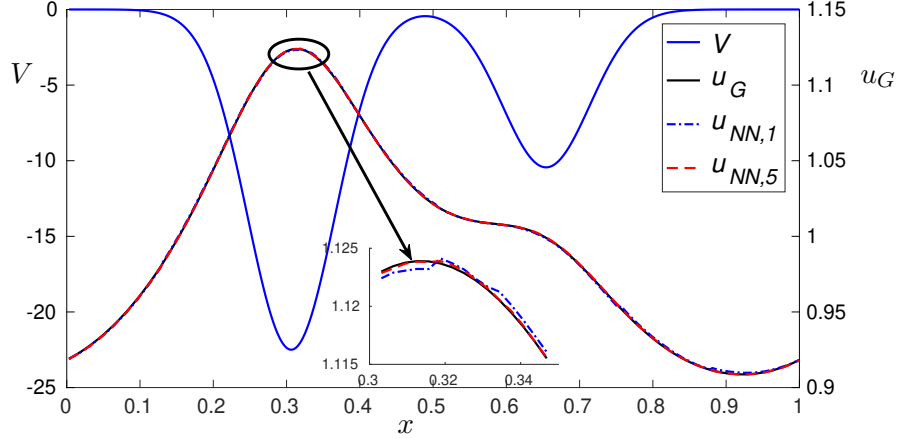


Figure 8: An example of the potential  $V$  and its corresponding solution  $u_G$  and predicted solution  $u_{NN,K}$  by MNN with  $r = 6$ ,  $n_g = 2$ .

$N_{\text{samples}}^{\text{train}}$	$N_{\text{samples}}^{\text{test}}$	Training error	Validation error
500	5000	2.1e-4	2.4e-4
1000	5000	1.8e-4	2.0e-4
5000	5000	1.4e-4	1.5e-4
20000	20000	1.5e-4	1.5e-4

Table 1: Relative error in approximating the ground state of NLSE for different number of samples  $N_{\text{samples}}^{\text{train}}$  for 1D case with  $r = 6$ ,  $K = 5$  and  $n_g = 2$ .

$r$	$N_{\text{params}}$	Training error	Validation error
2	1895	4.0e-4	4.0e-4
4	4555	2.0e-4	2.0e-4
6	8535	1.4e-4	1.5e-4

Table 2: Relative error in approximating the ground state of NLSE for different number of channels  $r$  for 1D case with  $K = 5$ ,  $n_g = 2$  and  $N_{\text{samples}}^{\text{train}} = 5000$ .

$K$	$N_{\text{params}}$	Training error	Validation error
1	3343	7.5e-4	7.5e-4
3	5939	2.1e-4	2.1e-4
5	8535	1.4e-4	1.5e-4

Table 3: Relative error in approximating the ground state of NLSE for different number of CK layers  $K$  for 1D case with  $r = 6$ ,  $n_g = 2$  and  $N_{\text{samples}}^{\text{train}} = 5000$ .

$n_g$	Training error	Validation error
2	1.4e-4	1.5e-4
4	2.6e-4	2.7e-4
6	2.6e-4	2.6e-4
8	2.8e-4	2.9e-4

Table 4: Relative error in approximating the ground state of NLSE for different number of Gaussians  $n_g$  for 1D case with  $K = 5$ ,  $r = 6$  and  $N_{\text{samples}}^{\text{train}} = 5000$ .

374 Usually, the number of samples should be greater than that of parameters to avoid over-  
375 fitting. But in neural network, it has been consistently found that the number of samples  
376 can be less than that of parameters [48, 49]. We present the numerical results for different  
377  $N_{\text{samples}}^{\text{train}}$  with  $K = 5$ ,  $r = 6$  and  $n_g = 2$  in Table 1. In this case, the number of parameters  
378 is  $N_{\text{params}} = 8535$ . There is no overfitting even  $N_{\text{samples}}^{\text{train}} = 500$ , and the error is only slightly  
379 larger than that when  $N_{\text{samples}}^{\text{train}} = 20000$ . For the case  $N_{\text{samples}}^{\text{train}} = 5000$ , the error is close to that  
380 for  $N_{\text{samples}}^{\text{train}} = 20000$ . This allows us to train MNN with  $N_{\text{samples}}^{\text{train}} < N_{\text{params}}$ . This feature is  
381 quite useful for high-dimensional case, because for high-dimensional case,  $N_{\text{params}}$  is usually  
382 very large and generating samples is expensive.

383 Table 2 presents the numerical results for different number of channels (*i.e.* the rank of  
384 the  $\mathcal{H}$ -matrix)  $r$  with  $K = 5$ ,  $n_g = 2$  and  $N_{\text{samples}}^{\text{train}} = 5000$ . As  $r$  increases, we find that the  
385 error first consistently decreases and then stagnates. We use  $r = 6$  for the 1D NLSE below to  
386 balance between efficiency and accuracy.

387 Similarly, Table 3 presents the numerical results for different number of CK layers  $K$  with  
388  $r = 6$ ,  $n_g = 2$  and  $N_{\text{samples}}^{\text{train}} = 5000$ . The error consistently decreases with respect to the  
389 increase of  $K$ , as NN can represent increasingly more complex functions with respect to the  
390 depth of the network. However, after a certain threshold, adding more layers provides very  
391 marginal gains in accuracy. In practice,  $K = 5$  is a good choice for the NLSE for 1D case.

392 Table 4 presents the numerical results for different number of Gaussians in the potential  
393  $V$ , for fixed  $K = 5$ ,  $r = 6$  and  $N_{\text{samples}}^{\text{train}} = 5000$ . Table 4 shows that within the class of input  
394 functions considered here, MNN is not sensitive to the complexity of the input. In particular,  
395 it shows that for the cases considered, increasing the number of wells only marginally increases  
396 the error.

397 Across the results in Tables 2 to 4, the validation errors are very close to the corresponding



398 training errors, and no overfitting is observed. Figure 8 presents a sample for the potential  
 399  $V$  and its corresponding solution and prediction solution by MNN. The prediction solution  
 400 agrees with the target solution very well.

401 **4.2.2. Two-dimensional case.** For two-dimensional case, the number of discretization  
 402 points is  $n = 80 \times 80$ , and we set  $L = 4$  and  $m = 5$ . In all the tests, the number of test  
 403 data is same as that of the train data. We perform simulation to study the behavior of MNN  
 404 for different number of channels  $r$ , different number of CK layers  $K$  and different number of  
 405 Gaussians  $n_g$ . As is discussed in 1D case, MNN allows  $N_{\text{samples}}^{\text{train}} < N_{\text{params}}$ . In all the tests, the  
 406 number of samples is always 20000 and numerical test shows no overfitting for all the test.

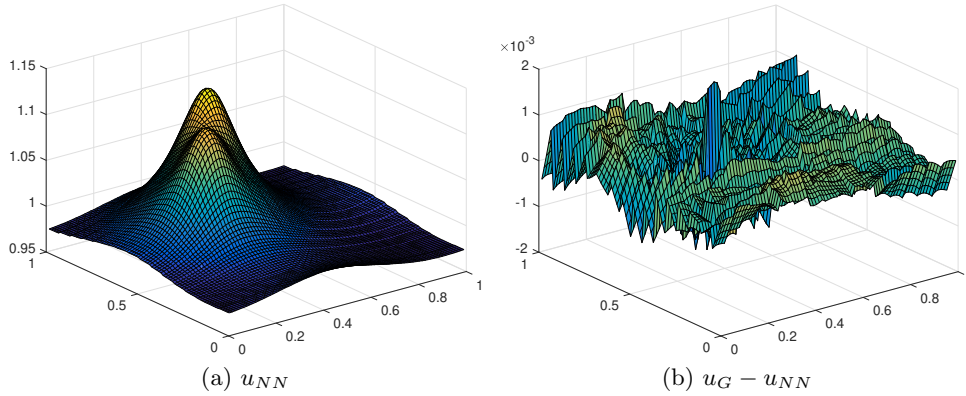


Figure 9: Prediction solution of MNN for  $K = 5$ ,  $r = 6$  and  $n_g = 2$  and its error with respect to the reference solution.

$r$	$N_{\text{params}}$	Training error	Validation error
2	33371	4.9e-4	4.9e-4
6	57323	1.5e-4	1.5e-4
10	100955	1.4e-4	1.4e-4

Table 5: Relative error in approximating the ground state of NLSE for different number of channels  $r$  for 2D case with  $K = 5$ ,  $n_g = 2$  and  $N_{\text{samples}}^{\text{train}} = 20000$ .

407 Tables 5 and 6 present the numerical results for different number of channels  $r$  and different  
 408 number of CK layers  $K$ , respectively. We find that similar to the 1D case, the choice of  
 409 parameters  $r = 6$  and  $K = 5$  also yield accurate results in the 2D case. Table 7 presents the  
 410 numerical results for different number of Gaussians in the potential  $V$  for  $K = 5$ ,  $r = 6$  and  
 411  $N_{\text{samples}}^{\text{train}} = 20000$ . We can find that MNN is also not sensitive to the complexity of the number  
 412 of Gaussians. Figure 9 presents the prediction of a sample and its corresponding error with  
 413 respect to the reference.

$K$	$N_{\text{params}}$	Training error	Validation error
1	16939	6.9e-4	7.0e-4
3	37131	2.1e-4	2.1e-4
5	57323	1.5e-4	1.5e-4
7	77515	1.4e-4	1.4e-4

Table 6: Relative error in approximating the ground state of NLSE for different number of CK layers  $K$  for 2D case with  $r = 6$ ,  $n_g = 2$  and  $N_{\text{samples}}^{\text{train}} = 20000$ .

$n_g$	Training error	Validation error
2	1.5e-4	1.5e-4
4	1.5e-4	1.5e-4
6	2.4e-4	2.4e-4
8	2.2e-4	2.2e-4

Table 7: Relative error in approximating the ground state of NLSE for different number of Gaussians  $n_g$  for the 2D case with  $r = 6$ ,  $K = 5$  and  $N_{\text{samples}}^{\text{train}} = 20000$ .

414 **4.3. Kohn-Sham map.** Kohn-Sham density functional theory [26, 30] is the most widely  
 415 used electronic structure theory. It requires the solution of the following set of nonlinear  
 416 eigenvalue equations (real arithmetic assumed for all quantities):

$$\begin{aligned}
 & \left( -\frac{1}{2}\Delta + V[\rho](x) \right) \psi_i(x) = \varepsilon_i \psi_i(x), \quad x \in \Omega = [-1, 1]^d \\
 & \int_{\Omega} \psi_i(x) \psi_j(x) \, dx = \delta_{ij}, \quad \rho(x) = \sum_{i=1}^{n_e} |\psi_i(x)|^2.
 \end{aligned}
 \tag{4.5}$$

418 Here  $n_e$  is the number of electrons (spin degeneracy omitted),  $d$  is the spatial dimension,  
 419 and  $\delta_{ij}$  stands for the Kronecker delta. In addition, all eigenvalues  $\{\varepsilon_i\}$  are real and ordered  
 420 non-decreasingly, and  $\rho(x)$  is the electron density, which satisfies the constraint

$$\rho(x) \geq 0, \quad \int_{\Omega} \rho(x) \, dx = n_e.
 \tag{4.6}$$

422 The Kohn-Sham equations (4.5) need to be solved self-consistently, which can also viewed as  
 423 solving the following fixed point map

$$\rho = \mathcal{F}_{\text{KS}}[V[\rho]].
 \tag{4.7}$$

425 Here the mapping  $\mathcal{F}_{\text{KS}}[\cdot]$  from  $V$  to  $\rho$  is called the Kohn-Sham map, which for a fixed potential  
 426 is reduced to a linear eigenvalue problem, and it constitutes the most computationally intensive  
 427 step for solving (4.5). We seek to approximate the Kohn-Sham map using a multiscale neural  
 428 network, whose output was regularized so it satisfies (4.6).

$r$	$N_{\text{params}}$	Training error	Validation error
2	2117	6.7e-4	6.7e-4
4	5183	3.3e-4	3.4e-4
6	9833	2.8e-4	2.8e-4
8	16067	3.3e-4	3.3e-4
10	33013	1.8e-4	1.9e-4

Table 8: Relative error on the approximation of the Kohn-Sham map for different  $r$ , with  $K = 6$ ,  $N_{\text{samples}}^{\text{train}} = 16000$ , and  $N_{\text{samples}}^{\text{test}} = 4000$ .

429 In the following numerical experiments the potential,  $V$ , is given by

430 (4.8) 
$$V(x) = - \sum_{i=1}^{n_g} \sum_{j \in \mathbb{Z}^d} c_i \exp \left( - \frac{(x - r_i - 2j)^2}{2\sigma^2} \right), \quad x \in [-1, 1]^d,$$

431 where  $d$  is the dimension and  $r_i \in [-1, 1]^d$ . We set  $\sigma = 0.05$  for 1D and  $\sigma = 0.2$  for 2D case.  
 432 The coefficients  $c_i$  are randomly chosen following the uniform distribution  $\mathcal{U}([0.8, 1.2])$ , and the  
 433 centers of the Gaussian wells  $r_i$ , are chosen randomly under the constraint that  $|r_i - r_{i'}| > 2\sigma$ .  
 434 The Kohn-Sham map is discretized using a pseudo-spectral method [46], and solved by a  
 435 standard eigensolver.

436 **4.3.1. One-dimensional case.** We generated 7 data sets using different number of wells,  
 437  $n_g$ , which in this case is also equal to the number of electrons  $n_e$ , ranging from 2 to 8. The  
 438 number of discretization points is  $N = 320$ . We trained the architecture defined in section 3  
 439 for each  $n_g$ , setting the number of levels  $L = 6$ , using different values for  $r$  and  $K$ .

440 Table 8 shows that there is no overfitting, even at this level of accuracy and number of  
 441 parameters. This behavior is found in all the numerical examples, thus we only report the  
 442 test error in what follows.

443 From Table 9 we can observe that as we increase  $r$  the error decreases sharply. Figure 10  
 444 depict this behavior. In Figure 10 we have that if  $r = 2$ , then the network output  $\rho_{NN}$ , fails to  
 445 approximate  $\rho$  accurately; however, by modestly increasing  $r$ , the network is able to properly  
 446 approximate  $\rho$ .

447 However, the accuracy of the network stagnates rapidly. In fact, increasing  $r$  beyond 10  
 448 does not provide any considerable gains. In addition, Table 9 shows that the accuracy of the  
 449 network is agnostic to the number of Gaussian wells present in the system.

450 In addition, we studied the relation between the quality of the approximation and  $K$ . We  
 451 fixed  $r = 6$ , and we trained several networks using different values of  $K$ , ranging from 2, *i.e.*,  
 452 a very shallow network, to 10. The results are summarized in Table 10. We can observe that  
 453 the error decreases sharply as the depth of the network increases, and then stagnates as  $K$   
 454 becomes large.

455 **4.3.2. Two-dimensional case.** The discretization is the standard extension to 2D using  
 456 tensor products, using a  $64 \times 64$  grid. In this case we only used  $n_g = 2$  and we followed the  
 457 same number of training and test samples as that in the 1D case. We fixed  $K = 6$ ,  $L = 4$ ,

$n_g \backslash r$	2	4	6	8	10
2	6.7e-4	3.3e-4	2.8e-4	3.3e-4	1.8e-4
3	7.3e-4	3.6e-4	3.0e-4	2.0e-4	2.8e-4
4	8.5e-4	4.1e-4	2.9e-4	3.8e-4	2.4e-4
5	9.0e-4	5.9e-4	4.0e-4	3.6e-4	3.6e-4
6	6.3e-4	4.8e-4	3.8e-4	4.0e-4	4.2e-4
7	8.6e-4	5.5e-4	3.9e-4	3.7e-4	3.5e-4
8	1.2e-3	5.1e-4	3.7e-4	4.5e-4	3.7e-4

Table 9: Relative test error on the approximation of the Kohn-Sham map for different ranks  $r$ , with fixed  $K = 6$  and  $N_{\text{samples}}^{\text{train}} = 16000$ .

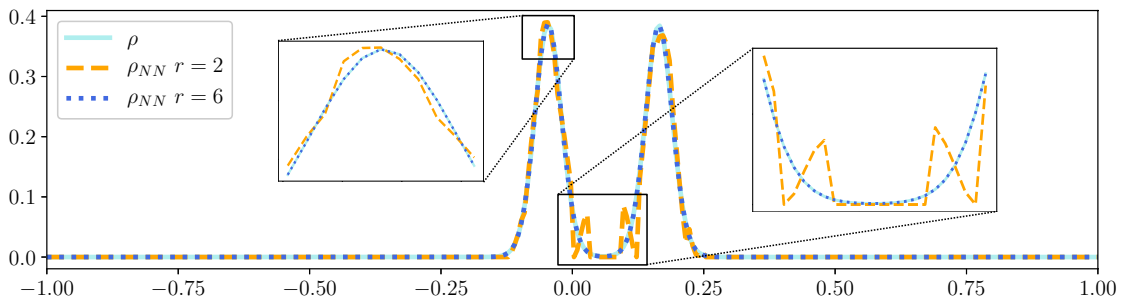


Figure 10: Estimation using two different multiscale networks with  $r = 2$ , and  $r = 6$ ; with  $K = 6$ , and  $L = 5$  fixed.

458 and we trained the network for different number of channels,  $r$ . The results are displayed in  
 459 [Table 11](#), which shows the same behavior as for the 1D case, in which the error decays sharply  
 460 and then stagnates, and there is no overfitting. In particular, the network is able to effectively  
 461 approximate the Kohn-Sham map as shown in [Figure 11](#). [Figure 11a](#) shows the output of  
 462 neural network for a test sample and [Figure 11b](#) shows the approximation error with respect  
 463 to the reference.

464 **5. Conclusion.** We have developed a multiscale neural network (MNN) architecture for  
 465 approximating nonlinear mappings, such as those arising from the solution of integral equa-  
 466 tions (IEs) or partial differential equations (PDEs). In order to control the number of param-  
 467 eters, we first rewrite the widely used hierarchical matrix into the form of a neural network,  
 468 which mainly consists of three sub-networks: restriction network, kernel network, and interpo-  
 469 lation network. The three sub-networks are all linear, and correspond to the components of a  
 470 singular value decomposition. We demonstrate that such structure can be directly generalized  
 471 to nonlinear problems, simply by replacing the linear kernel network by a multilayer kernel  
 472 network with nonlinear activation functions. Such “nonlinear singular value decomposition  
 473 operation” is performed at different spatial scales, which can be efficiently implemented by  
 474 a number of locally connected (LC) networks, or convolutional neural networks (CNN) when  
 475 the mapping is equivariant to translation. Using the parameterized nonlinear Schrödinger

$n_g \backslash K$	2	4	6	8	10
2	1.4e-3	3.1e-4	2.8e-4	3.5e-4	2.3e-4
3	2.0e-3	5.4e-4	3.0e-4	5.6e-4	5.3e-4
4	1.9e-3	5.8e-4	2.8e-4	6.0e-4	7.1e-4
5	1.8e-3	7.2e-4	4.0e-4	8.0e-4	7.4e-4
6	2.1e-3	7.3e-4	3.8e-4	6.7e-4	6.7e-4
7	2.2e-3	7.9e-4	3.8e-4	7.4e-4	5.8e-4
8	2.0e-3	8.8e-4	3.7e-4	6.7e-4	6.8e-4

Table 10: Relative test error on the approximation of the Kohn-Sham map for different  $K$  and fixed rank  $r = 6$ , and  $N_{\text{samples}}^{\text{train}} = 16000$ .

$r$	Training error	Validation error
4	5.2e-3	5.2e-3
6	1.6e-3	1.7e-3
8	1.2e-3	1.1e-3
10	9.1e-4	9.3e-4

Table 11: Relative errors on the approximation of the Kohn-Sham map for 2D case for different  $r$  and  $K = 6$ ,  $N_{\text{samples}}^{\text{train}} = 16000$  and  $N_{\text{samples}}^{\text{test}} = 4000$ .

476 equation and the Kohn-Sham map as examples, we find that MNN can yield accurate ap-  
 477 proximation to such nonlinear mappings. When the mapping has  $N$  degrees of freedom, the  
 478 complexity of MNN is only  $O(N \log N)$ . Thus the resulting MNN can be further used to  
 479 accelerate the evaluation of the mapping, especially when a large number of evaluations are  
 480 needed within a certain range of parameters.

481 In this work, we only provide one natural architecture of multiscale neural network based  
 482 on hierarchical matrices. The architecture can be altered depending on the target application.  
 483 Some of the possible modifications and extensions are listed below. 1) In this work, the neural  
 484 network is inspired by a hierarchical matrix with a special case of strong admissible condition.  
 485 One can directly construct architectures for  $\mathcal{H}$ -matrices with the weak admissible condition,  
 486 as well as other structures such as the fast multiple methods,  $\mathcal{H}^2$ -matrices and wavelets. 2)  
 487 The input,  $u$ , and output,  $v$ , in this work are periodic. The network can be directly extended  
 488 to the non-periodic case, by replacing the periodic padding in LCK by some other padding  
 489 functions. One may also explore the mixed usage of LC networks and CNNs in different  
 490 components of the architecture. 3) The matrices  $A^{(\text{ad})}$  and  $M^{(\ell)}$  can be block-partitioned in  
 491 different ways, which would result in different setups of parameters in the LCK layers. 4) The  
 492 LCR and LCI networks in Algorithm 3 can involve nonlinear activation functions as well and  
 493 can be extended to networks with more than one layer. 5) The LCK network in Algorithm  
 494 3 can be replaced by other architectures. In principle, for each scale, these LCK layers can  
 495 be altered to any network, for example the sum of two parallel subnetworks, or the ResNet  
 496 structure[24]. 6) It is known that  $\mathcal{H}$ -matrices can well approximate smooth kernels but become

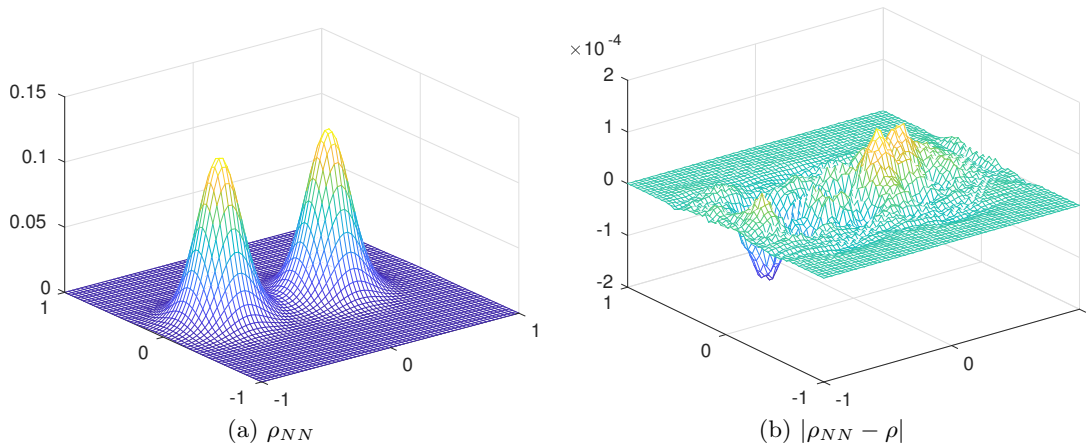


Figure 11: (a) Output of the trained network on a test sample for  $K = 6$ , and  $\alpha = 10$ ; (b) error with respect to the reference solution.

497 less efficient for highly oscillatory kernels, such as those arising from the Helmholtz operator  
 498 in the high frequency regime. The range of applicability of the MNN remains to be studied  
 499 both theoretically and numerically.

500 **Acknowledgements.** The authors thank Yuehaw Khoo for constructive discussions.

501

#### REFERENCES

- 502 [1] M. ABADI, P. BARHAM, J. CHEN, Z. CHEN, A. DAVIS, J. DEAN, M. DEVIN, S. GHEMAWAT, G. IRVING,  
 503 M. ISARD, ET AL., *Tensorflow: A system for large-scale machine learning.*, in OSDI, vol. 16, 2016,  
 504 pp. 265–283.
- 505 [2] J. R. ANGLIN AND W. KETTERLE, *Bose–Einstein condensation of atomic gases*, Nature, 416 (2002),  
 506 p. 211.
- 507 [3] W. BAO AND Q. DU, *Computing the ground state solution of Bose–Einstein condensates by a normalized*  
 508 *gradient flow*, SIAM Journal on Scientific Computing, 25 (2004), pp. 1674–1697.
- 509 [4] M. BARRAULT, Y. MADAY, N. C. NGUYEN, AND A. T. PATERA, *An ‘empirical interpolation’ method:*  
 510 *application to efficient reduced-basis discretization of partial differential equations*, Comptes Rendus  
 511 Mathématique, 339 (2004), pp. 667 – 672, <https://doi.org/10.1016/j.crma.2004.08.006>.
- 512 [5] J. BERG AND K. NYSTRÖM, *A unified deep artificial neural network approach to partial differential*  
 513 *equations in complex geometries*, arXiv preprint arXiv:1711.06464, (2017).
- 514 [6] S. BÖRM, L. GRASEDYCK, AND W. HACKBUSCH, *Introduction to hierarchical matrices with applications*,  
 515 Engineering analysis with boundary elements, 27 (2003), pp. 405–422.
- 516 [7] A. BRANDT, *Multi-level adaptive solutions to boundary-value problems*, Mathematics of Computation, 31  
 517 (1977), pp. 333–390.
- 518 [8] S. CHAN AND A. H. ELSHEIKH, *A machine learning approach for efficient uncertainty quantification using*  
 519 *multiscale methods*, Journal of Computational Physics, 354 (2018), pp. 493 – 511, [https://doi.org/10.](https://doi.org/10.1016/j.jcp.2017.10.034)  
 520 [1016/j.jcp.2017.10.034](https://doi.org/10.1016/j.jcp.2017.10.034).
- 521 [9] P. CHAUDHARI, A. OBERMAN, S. OSHER, S. SOATTO, AND G. CARLIER, *Partial differential equations for*  
 522 *training deep neural networks*, in 2017 51st Asilomar Conference on Signals, Systems, and Computers,  
 523 Oct 2017, pp. 1627–1631, <https://doi.org/10.1109/ACSSC.2017.8335634>.
- 524 [10] F. CHOLLET ET AL., *Keras*. <https://keras.io>, 2015.

- 525 [11] N. COHEN, O. SHARIR, AND A. SHASHUA, *On the expressive power of deep learning: A tensor analysis*,  
526 arXiv preprint arXiv:1603.00988, (2018).
- 527 [12] T. DOZAT, *Incorporating nesterov momentum into adam*, 2015.
- 528 [13] W. E, J. HAN, AND A. JENTZEN, *Deep learning-based numerical methods for high-dimensional parabolic*  
529 *partial differential equations and backward stochastic differential equations*, Communications in Math-  
530 ematics and Statistics, 5 (2017), pp. 349–380, <https://doi.org/10.1007/s40304-017-0117-6>.
- 531 [14] Y. EFENDIEV, J. GALVIS, G. LI, AND M. PRESHO, *Generalized multiscale finite element methods.*  
532 *nonlinear elliptic equations*, Communications in Computational Physics, 15 (2014), pp. 733–755,  
533 <https://doi.org/10.4208/cicp.020313.041013a>.
- 534 [15] Y. EFENDIEV AND T. HOU, *Multiscale finite element methods for porous media flows and their applica-*  
535 *tions*, Applied Numerical Mathematics, 57 (2007), pp. 577 – 596, [https://doi.org/10.1016/j.apnum.](https://doi.org/10.1016/j.apnum.2006.07.009)  
536 [2006.07.009](https://doi.org/10.1016/j.apnum.2006.07.009). Special Issue for the International Conference on Scientific Computing.
- 537 [16] M. FENN AND G. STEIDL, *FMM and  $\mathcal{H}$ -matrices: a short introduction to the basic idea*, Tech. Report  
538 TR-2002-008, Department for Mathematics and Computer Science, University of Mannheim, 2002,  
539 <https://ub-madoc.bib.uni-mannheim.de/744/1/TR-02-008.pdf>.
- 540 [17] L. GREENGARD AND V. ROKHLIN, *A fast algorithm for particle simulations*, Journal of computational  
541 physics, 73 (1987), pp. 325–348.
- 542 [18] M. A. GREPL, Y. MADAY, N. C. NGUYEN, AND A. T. PATERA, *Efficient reduced-basis treatment of non-*  
543 *affine and nonlinear partial differential equations*, ESAIM: Mathematical Modelling and Numerical  
544 Analysis, 41 (2007), p. 575–605, <https://doi.org/10.1051/m2an:2007031>.
- 545 [19] W. HACKBUSCH, *A sparse matrix arithmetic based on  $\mathcal{H}$ -matrices. part I: Introduction to  $\mathcal{H}$ -matrices*,  
546 Computing, 62 (1999), pp. 89–108.
- 547 [20] W. HACKBUSCH, L. GRASEDYCK, AND S. BÖRM, *An introduction to hierarchical matrices*, Math. Bohem.,  
548 127 (2002).
- 549 [21] W. HACKBUSCH AND B. N. KHOROMSKIJ, *A sparse  $\mathcal{H}$ -matrix arithmetic: general complexity estimates*,  
550 Journal of Computational and Applied Mathematics, 125 (2000), pp. 479–501.
- 551 [22] K. HE AND J. SUN, *Convolutional neural networks at constrained time cost*, 2015 IEEE Conference on  
552 Computer Vision and Pattern Recognition (CVPR), (2015), pp. 5353–5360.
- 553 [23] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, 2016 IEEE  
554 Conference on Computer Vision and Pattern Recognition (CVPR), (2016), pp. 770–778.
- 555 [24] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in Proceedings of the  
556 IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- 557 [25] G. HINTON, L. DENG, D. YU, G. E. DAHL, A. R. MOHAMED, N. JAITLY, A. SENIOR, V. VANHOUCKE,  
558 P. NGUYEN, T. N. SAINATH, AND B. KINGSBURY, *Deep neural networks for acoustic modeling in*  
559 *speech recognition: The shared views of four research groups*, IEEE Signal Processing Magazine, 29  
560 (2012), pp. 82–97, <https://doi.org/10.1109/MSP.2012.2205597>.
- 561 [26] P. HOHENBERG AND W. KOHN, *Inhomogeneous electron gas*, Physical review, 136 (1964), p. B864.
- 562 [27] K. HORNIK, *Approximation capabilities of multilayer feedforward networks*, Neural Networks, 4 (1991),  
563 pp. 251–257, [https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T).
- 564 [28] Y. KHOO, J. LU, AND L. YING, *Solving parametric PDE problems with artificial neural networks*, arXiv  
565 preprint arXiv:1707.03351, (2017).
- 566 [29] N. KISHORE KUMAR AND J. SCHNEIDER, *Literature survey on low rank approximation of matrices*, Linear  
567 and Multilinear Algebra, 65 (2017), pp. 2212–2244.
- 568 [30] W. KOHN AND L. J. SHAM, *Self-consistent equations including exchange and correlation effects*, Physical  
569 review, 140 (1965), p. A1133.
- 570 [31] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convolutional neu-*  
571 *ral networks*, in Proceedings of the 25th International Conference on Neural Information Processing  
572 Systems - Volume 1, NIPS’12, USA, 2012, Curran Associates Inc., pp. 1097–1105.
- 573 [32] Y. LECUN, Y. BENGIO, AND G. HINTON, *Deep learning*, Nature, 521 (2015).
- 574 [33] M. K. K. LEUNG, H. Y. XIONG, L. J. LEE, AND B. J. FREY, *Deep learning of the tissue-regulated splicing*  
575 *code*, Bioinformatics, 30 (2014), pp. i121–i129, <https://doi.org/10.1093/bioinformatics/btu277>.
- 576 [34] Y. LI, X. CHENG, AND J. LU, *Butterfly-Net: Optimal function representation based on convolutional*  
577 *neural networks*, arXiv preprint arXiv:1805.07451, (2018).
- 578 [35] L. LIN, J. LU, AND L. YING, *Fast construction of hierarchical matrix representation from matrix–vector*

- 579 *multiplication*, Journal of Computational Physics, 230 (2011), pp. 4071–4087.
- 580 [36] J. MA, R. P. SHERIDAN, A. LIAW, G. E. DAHL, AND V. SVETNIK, *Deep neural nets as a method*  
581 *for quantitative structure–activity relationships*, Journal of Chemical Information and Modeling, 55  
582 (2015), pp. 263–274, <https://doi.org/10.1021/ci500747n>. PMID: 25635324.
- 583 [37] Y. MADAY, O. MULA, AND G. TURINICI, *Convergence analysis of the generalized empirical interpolation*  
584 *method*, SIAM Journal on Numerical Analysis, 54 (2016), pp. 1713–1731, [https://doi.org/10.1137/](https://doi.org/10.1137/140978843)  
585 [140978843](https://doi.org/10.1137/140978843).
- 586 [38] S. MALLAT, *A wavelet tour of signal processing: the sparse way*, in *A wavelet tour of signal processing: the*  
587 *sparse way*, Academic press, Boston, third ed., 2008, [https://doi.org/10.1016/B978-0-12-374370-1.](https://doi.org/10.1016/B978-0-12-374370-1.50001-9)  
588 [50001-9](https://doi.org/10.1016/B978-0-12-374370-1.50001-9).
- 589 [39] H. MHASKAR, Q. LIAO, AND T. POGGIO, *Learning functions: When is deep better than shallow*, arXiv  
590 preprint arXiv:1603.00988, (2018).
- 591 [40] L. PITAEVSKII, *Vortex lines in an imperfect Bose gas*, Sov. Phys. JETP, 13 (1961), pp. 451–454.
- 592 [41] M. RAISSI AND G. E. KARNIADAKIS, *Hidden physics models: Machine learning of nonlinear partial*  
593 *differential equations*, Journal of Computational Physics, 357 (2018), pp. 125 – 141, [https://doi.org/](https://doi.org/10.1016/j.jcp.2017.11.039)  
594 [10.1016/j.jcp.2017.11.039](https://doi.org/10.1016/j.jcp.2017.11.039).
- 595 [42] K. RUDD, G. D. MURO, AND S. FERRARI, *A constrained backpropagation approach for the adaptive solu-*  
596 *tion of partial differential equations*, IEEE Transactions on Neural Networks and Learning Systems,  
597 25 (2014), pp. 571–584, <https://doi.org/10.1109/TNNLS.2013.2277601>.
- 598 [43] J. SCHMIDHUBER, *Deep learning in neural networks: An overview*, Neural Networks, 61 (2015), pp. 85 –  
599 117, <https://doi.org/10.1016/j.neunet.2014.09.003>.
- 600 [44] K. SPILIOPOULOS AND J. SIRIGNANO, *Dgm: A deep learning algorithm for solving partial differential*  
601 *equations*, arXiv preprint arXiv:1708.07469, (2018).
- 602 [45] I. SUTSKEVER, O. VINYALS, AND Q. V. LE, *Sequence to sequence learning with neural networks*, in  
603 *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes,  
604 N. D. Lawrence, and K. Q. Weinberger, eds., Curran Associates, Inc., 2014, pp. 3104–3112.
- 605 [46] L. TREFETHEN, *Spectral Methods in MATLAB*, Society for Industrial and Applied Mathematics, 2000,  
606 <https://doi.org/10.1137/1.9780898719598>.
- 607 [47] Y. WANG, C. W. SIU, E. T. CHUNG, Y. EFENDIEV, AND M. WANG, *Deep multiscale model learning*,  
608 arXiv preprint arXiv:1806.04830, (2018).
- 609 [48] C. ZHANG, S. BENGIO, M. HARDT, B. RECHT, AND O. VINYALS, *Understanding deep learning requires*  
610 *rethinking generalization*, arXiv:1611.03530, (2016).
- 611 [49] K. ZHANG, W. ZUO, Y. CHEN, D. MENG, AND L. ZHANG, *Beyond a Gaussian denoiser: Residual learning*  
612 *of deep CNN for image denoising*, IEEE Transactions on Image Processing, 26 (2017), pp. 3142–3155.
- 613 [50] L. ZHANG, J. HAN, H. WANG, R. CAR, AND W. E, *Deepcg: constructing coarse-grained models via deep*  
614 *neural networks*, arXiv preprint arXiv:1802.08549, (2018).
- 615 [51] L. ZHANG, H. WANG, AND W. E, *Adaptive coupling of a deep neural network potential to a classical force*  
616 *field*, arXiv preprint arXiv:1806.01020, (2018).