

# U-NETS AS BELIEF PROPAGATION: EFFICIENT CLASSIFICATION, DENOISING, AND DIFFUSION IN GENERATIVE HIERARCHICAL MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

U-Nets are among the most widely used architectures in computer vision, renowned for their exceptional performance in applications such as image segmentation, denoising, and diffusion modeling. However, a theoretical explanation of the U-Net architecture design has not yet been fully established.

This paper introduces a novel interpretation of the U-Net architecture by studying certain generative hierarchical models, which are tree-structured graphical models extensively utilized in both language and image domains. With their encoder-decoder structure, long skip connections, and pooling and up-sampling layers, we demonstrate how U-Nets can naturally implement the belief propagation denoising algorithm in such generative hierarchical models, thereby efficiently approximating the denoising functions. This leads to an efficient sample complexity bound for learning the denoising function using U-Nets within these models. Additionally, we discuss the broader implications of these findings for diffusion models in generative hierarchical models. We also demonstrate that the conventional architecture of convolutional neural networks (ConvNets) is ideally suited for classification tasks within these models. This offers a unified view of the roles of ConvNets and U-Nets, highlighting the versatility of generative hierarchical models in modeling complex data distributions.

## 1 INTRODUCTION

U-Nets are one of the most prominent network architectures in computer vision, primarily employed for tasks such as image segmentation, denoising (Ronneberger et al., 2015; Zhou et al., 2018; Siddique et al., 2021; Oktay et al., 2018), and diffusion modeling (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song & Ermon, 2019; Song et al., 2020). These networks are structured as encoder-decoder convolutional neural networks equipped with long skip connections, and their input and output typically maintain the same dimensions. While U-Nets have demonstrated exceptional performance across a variety of applications, the theoretical foundations of their key components—including the encoder-decoder structure, long skip connections, and the pooling and up-sampling layers—remain inadequately understood. Notably, long skip connections have a significant impact on performance as shown in empirical studies (Drozdal et al., 2016; Wang et al., 2022). Existing explanations, often anecdotal, suggest their efficacy stems from improved information propagation and reduction of the vanishing gradient issue, but a thorough theoretical exploration is still lacking.

In this paper, we introduce a novel interpretation of the U-Net architecture, viewing it through the lens of neural network approximation. We posit that:

*U-Nets naturally approximate the belief propagation denoising algorithm in certain generative hierarchical models.*

The generative hierarchical model discussed herein is a tree-structured graphical model, which has been widely employed in language and image generative modeling (Chomsky, 1959; Lee, 1996; Allen-Zhu & Li, 2023; Li et al., 2000; Willsky, 2002; Jin & Geman, 2006). We detail the precise definition of such generative hierarchical models in Section 2. A series of recent work (Mossel, 2016; Sclocchi et al., 2024; Tomasini & Wyart, 2024; Petrini et al., 2023; Kadkhodaie et al., 2023a)

054 have pioneered the use of generative hierarchical models in studying classification tasks and dif-  
 055 fusion models. Kadkhodaie et al. (2023a) has empirically shown that U-Nets can effectively learn  
 056 the denoising function for these models. Furthermore, as noted by Sclocchi et al. (2024), the belief  
 057 propagation denoising algorithm, which computes the denoising function in these models exactly,  
 058 includes a downward process and an upward process, with the latter reusing the intermediate results  
 059 from the downward process. In Section 4, we demonstrate how the belief propagation algorithm  
 060 naturally induces the encoder-decoder structure, the long skip connections, and the pooling and up-  
 061 sampling operations of the U-Nets. This gives rise to an efficient sample complexity bound for  
 062 learning the denoising function in generative hierarchical models using U-Nets.

063 In addition to our main findings, in Section 3, we demonstrate that the standard architecture of con-  
 064 volutional neural networks (ConvNets) is well-suited for classification tasks within the same gener-  
 065 ative hierarchical model. We provide efficient sample complexity results to support this assertion.  
 066 This offers a unified perspective on the role of both ConvNets and U-Nets in image classification  
 067 and denoising tasks, and also highlights the versatility of generative hierarchical models in modeling  
 068 data distributions across language and image domains.

## 070 2 THE GENERATIVE HIERARCHICAL MODEL

071  
 072 To define the generative hierarchical model, we start by introducing some key notations. Consider  
 073 a tree  $\mathcal{T} = (\mathcal{V}, \mathcal{E})$  with a height of  $L$ , where we conventionally designate the root of the tree  $r$   
 074 as layer 0. For each node  $v \in \mathcal{V}$ , we denote  $\text{pa}(v)$  as the parent of  $v$ ,  $\mathcal{C}(v)$  as the children of  $v$ ,  
 075 and  $\mathcal{N}(v)$  as the siblings of  $v$ . We denote  $\mathcal{V}^{(\ell)}$  as the set of nodes at layer  $\ell$ . We assume that for  
 076 any  $v \in \mathcal{V}^{(\ell-1)}$ , the number of children is precisely  $m^{(\ell)}$  for  $\ell \in [L]$ . The leaf nodes  $v \in \mathcal{V}^{(L)}$   
 077 have no children. Additionally, for each  $v \in \mathcal{V}^{(\ell)}$ , we assume an ordering function (a bijection)  
 078  $\iota : \mathcal{C}(v) \rightarrow [m^{(\ell)}]$ , ensuring that any child  $v' \in \mathcal{C}(v)$  possesses a unique rank  $\iota(v') \in [m^{(\ell)}]$ . We  
 079 denote the number of nodes at layer  $\ell$  as  $d^{(\ell)}$ , and the number of nodes at layer  $L$  as  $d = d^{(L)}$ . We  
 080 further denote  $\underline{m} = (m^{(\ell)})_{\ell \in [L]}$ , and  $\|\underline{m}\|_1 = \sum_{\ell=1}^L m^{(\ell)}$ . By these definitions and assumptions,  
 081 we have  $d^{(\ell)} = \prod_{1 \leq s \leq \ell} m^{(s)}$ , and  $1 = d^{(0)} \leq d^{(1)} \leq \dots \leq d^{(L)} = d$ .

082  
 083 For each layer  $\ell = 0, \dots, L$ , every tree node  $v \in \mathcal{V}^{(\ell)}$  is associated with a variable  $x_v^{(\ell)} \in [S]$  for  
 084 some  $S \in \mathbb{N}_{\geq 2}$ . (For simplicity, we use the same variable space  $[S]$  across all layers, although our  
 085 framework can accommodate variations across different layers  $\ell$ .) We denote  $\mathbf{x}^{(\ell)} = (x_v^{(\ell)})_{v \in \mathcal{V}^{(\ell)}} \in$   
 086  $[S]^{d^{(\ell)}}$  as the variables at layer  $\ell$ . The variables at the leaves,  $\mathbf{x} = \mathbf{x}^{(L)} \in [S]^d$  are considered  
 087 the observed covariates, exemplified by the pixel representation of an image. Conversely, the root  
 088 node variable  $y = x_r^{(0)} \in [S]$  is treated as the associated label. Variables for the intermediate layers  
 089  $\{\mathbf{x}^{(\ell)}\}_{1 \leq \ell \leq L-1}$  remain unobserved.

090  
 091 **The generative hierarchical model.** We consider a specific type of generative hierarchical model  
 092 (GHM)<sup>1</sup>, which is a joint distribution  $\mu_{\star}$  over variables

$$093 \quad (y = x^{(0)} \in [S], \quad \mathbf{x}^{(1)} \in [S]^{d^{(1)}}, \quad \dots, \quad \mathbf{x}^{(L-1)} \in [S]^{d^{(L-1)}}, \quad \mathbf{x}^{(L)} = \mathbf{x} \in [S]^d),$$

094 associated with a set of functions  $\{\psi^{(\ell)} : [S] \times [S]^{m^{(\ell)}} \rightarrow \mathbb{R}_{\geq 0}\}_{\ell \in [L]}$ , defined as

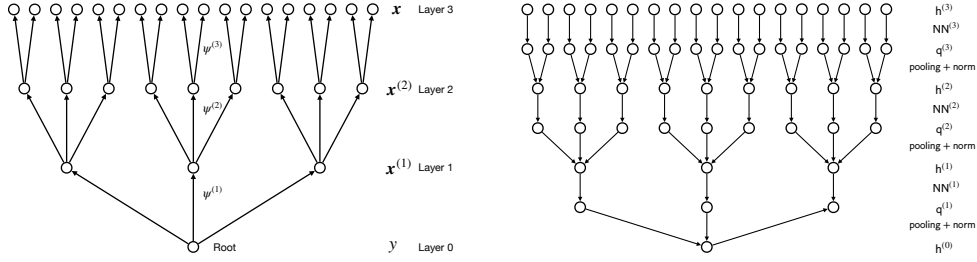
$$095 \quad \mu_{\star}(y, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(L-1)}, \mathbf{x})$$

$$096 \quad \propto \psi^{(1)}(y, \mathbf{x}^{(1)}) \cdot \left( \prod_{v \in \mathcal{V}^{(1)}} \psi^{(2)}(x_v^{(1)}, x_{\mathcal{C}(v)}^{(2)}) \right) \cdots \left( \prod_{v \in \mathcal{V}^{(L-1)}} \psi^{(L)}(x_v^{(L-1)}, x_{\mathcal{C}(v)}^{(L)}) \right). \quad \text{(GHM)}$$

097  
 098 The formula specifies that any two nodes  $v_1, v_2 \in \mathcal{V}^{(\ell)}$  within the same level uses the same function  
 099  $\psi^{(\ell)}$ , thereby embedding specific invariance properties into  $\mu_{\star}$ . Consequently, this GHM ensures  
 100 that  $(x_v)_{v \succeq v_1} \stackrel{d}{=} (x_v)_{v \succeq v_2}$  for any  $v_1, v_2 \in \mathcal{V}^{(\ell)}$ . The notation  $v \succeq v_1$  denotes that  $v$  is either

101  
 102  
 103  
 104  
 105 <sup>1</sup>We define “generative hierarchical models” as general probabilistic models with a hierarchical structure.  
 106 The specific model discussed in this paper is an instance of such generative hierarchical models. These mod-  
 107 els are also known by various other names, including “hierarchical generative models”, “latent hierarchical  
 models”, “Bayesian hierarchical models”, “hierarchical Markov random fields”, among others.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117



118 Figure 1: Left: The generative hierarchical model with 3 layers and  $m^{(1)} = 3$ ,  $m^{(2)} = 3$ , and  
119  $m^{(3)} = 2$  children in each layer. Right: A 3-layer convolutional neural network.

120  
121  
122

identical to or a descendant of  $v_1$ . We will explore later how this invariance property interacts with the convolutional structure of the neural networks to be introduced.

124  
125  
126  
127

Throughout the paper, we impose a specific assumption on the  $\psi$  functions in GHM, namely that each  $\psi$  function can be factorized into a product of functions that depend solely on the ordering of the child. This assumption is essential for enabling convolutional neural networks to approximate the associated belief propagation algorithm.

128

**Assumption 1** (Factorization of  $\psi$ ). For each layer  $\ell \in [L]$  and node  $v \in \mathcal{V}^{(\ell-1)}$ , we have

129

$$\psi^{(\ell)}(x_v^{(\ell-1)}, x_{\mathcal{C}(v)}^{(\ell)}) = \prod_{v' \in \mathcal{C}(v)} \psi_{\iota(v')}^{(\ell)}(x_v^{(\ell-1)}, x_{v'}^{(\ell)}). \quad (1)$$

130  
131

We also state a technical assumption concerning the boundedness of these  $\psi$  functions of GHM.

133  
134

**Assumption 2** (Boundedness of  $\psi$ ). For any layer  $\ell \in [L]$  and child node rank  $\iota \in [m^{(\ell)}]$ , the transition probabilities are bounded as follows:

135

$$1/K \leq \min_{x, x'} \psi_{\iota}^{(\ell)}(x, x') \leq \max_{x, x'} \psi_{\iota}^{(\ell)}(x, x') \leq K. \quad (2)$$

136  
137

It is helpful to think about the joint distribution  $\mu_*$  as a tree-structured Markov process, which admits the factorization

140

$$\begin{aligned} & \mu_*(y, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(L-1)}, \mathbf{x}) \\ &= \mu_*(y) \mu_*(\mathbf{x}^{(1)} | y) \cdot \left( \prod_{v \in \mathcal{V}^{(1)}} \mu_*(x_{\mathcal{C}(v)}^{(2)} | x_v^{(1)}) \right) \cdots \left( \prod_{v \in \mathcal{V}^{(L-1)}} \mu_*(x_{\mathcal{C}(v)}^{(L)} | x_v^{(L-1)}) \right), \end{aligned} \quad (3)$$

141  
142  
143

where we abuse the notation to denote  $\mu_*(x_{\mathcal{C}(v)}^{(\ell)} | x_v^{(\ell-1)})$  as the conditional probability of  $x_{\mathcal{C}(v)}^{(\ell)}$  given  $x_v^{(\ell-1)}$ , and  $\mu_*(y)$  as the marginal probability of  $y$ . Indeed, any graphical model specified as in Eq. (1) can be cast into the form of (3). Furthermore, It can be checked that Eq. (3) coincides with (GHM) upon taking  $\psi^{(\ell)}(x_v^{(\ell-1)}, x_{\mathcal{C}(v)}^{(\ell)}) = \mu_*(x_{\mathcal{C}(v)}^{(\ell)} | x_v^{(\ell-1)})$  for  $\ell \geq 1$ , and  $\psi^{(1)}(y, \mathbf{x}^{(1)}) = \mu_*(y) \mu_*(\mathbf{x}^{(1)} | y)$ . In this scenario, the factorization assumption (1) implies that  $(x_{v'}^{(\ell)})_{v' \in \mathcal{C}(v)}$  given  $x_v^{(\ell-1)}$  are conditionally independent. We include a schematic plot of the generative hierarchical model with 3 layers as in Figure 1(left).

152

153  
154  
155  
156  
157  
158  
159

**GHMs as natural models for languages and images.** In the field of linguistics, GHMs are very similar to context-free grammars (CFGs) (Chomsky, 1959; Lee, 1996; Allen-Zhu & Li, 2023). The generative process of a context-free grammar involves creating valid strings or sentences based on a given set of production rules (the  $\psi$  functions) that dictate how symbols can be extended to form new strings. Starting with an initial symbol (the label  $y$ ), the generation proceeds iteratively by applying production rules until all symbols belong to the terminal set, thus forming a complete sentence (the covariate  $\mathbf{x}$ ).

160  
161

In computer vision, GHMs are often utilized to model natural images (Li et al., 2000; Willsky, 2002; Jin & Geman, 2006), where they are sometimes referred to as multi-resolution Markov models. The hypothetical image generation process begins with a high-level concept of the image (represented

by the label  $y$ ), which is then iteratively refined using a production rule (the  $\psi$  function). This refinement continues through successive resolution levels until the detail reaches the pixel level, resulting in the final image (the covariate  $\mathbf{x}$ ).

We remark that a series of studies (Sclocchi et al., 2024; Tomasini & Wyart, 2024; Petrini et al., 2023; Kadkhodaie et al., 2023a) have provided both theoretical and empirical evidence supporting the efficacy of GHMs as powerful tools for modeling the properties of combinatorial data.

### 3 THE WARM-UP PROBLEM: CLASSIFICATION IN GHMS

In this section, we consider the warm-up problem of classification within the GHM. In the subsequent section, we will investigate the denoising task, where the results and intuitions will be similar and parallel to those presented in this section.

In the classification task, consider the scenario where we observe a set of iid samples  $\{(\mathbf{x}_i, y_i)\}_{i \in [n]}$  drawn from  $\mu_*$  under the GHM. Our objective is to learn a probabilistic classifier  $\hat{\mu}(y|\mathbf{x})$  from this dataset. With a suitable loss function, the optimal classifier is the Bayes classifier  $\mu_*(y|\mathbf{x})$ , which represents the true conditional probability of  $y$  given  $\mathbf{x}$ . We aim to examine the sample complexity of learning this classifier through empirical risk minimization over the class of convolutional networks.

**The ConvNet architecture.** We here introduce the convolutional neural network (ConvNet) architecture used for classification, represented as  $\mu_{\text{NN}}(\cdot|\mathbf{x}) \in \Delta([S])$  for input  $\mathbf{x} \in [S]^d$ . Initially, we set  $\mathbf{h}_v^{(L)} = x_v \in [S]$  for each node  $v \in \mathcal{V}^{(L)}$ . The operational flow of the network unfolds as follows:

$$\begin{aligned} \mathbf{q}_v^{(\ell)} &= \text{NN}_{\iota(v)}^{(\ell)}(\mathbf{h}_v^{(\ell)}) \in \mathbb{R}^S, & \ell \in [L], v \in \mathcal{V}^{(\ell)}, \\ \mathbf{h}_v^{(\ell-1)} &= \text{normalize}\left(\sum_{v' \in \mathcal{C}(v)} \mathbf{q}_{v'}^{(\ell)}\right) \in \mathbb{R}^S, & \ell \in [L], v \in \mathcal{V}^{(\ell-1)}, \\ \mu_{\text{NN}}(\cdot|\mathbf{x}) &= \text{softmax}(\mathbf{h}_r^{(0)}) \in \Delta([S]). \end{aligned} \quad (\text{ConvNet})$$

The functions  $\text{NN}_{\iota(v)}^{(\ell)} : \mathbb{R}^{S_{\text{in}}^{(\ell)}} \rightarrow \mathbb{R}^S$  (which adjust their input dimensionality  $S_{\text{in}}^{(\ell)}$  based on layer depth,  $S_{\text{in}}^{(\ell)} = (S+1) \cdot 1\{\ell \neq L\} + 2 \cdot 1\{\ell = L\}$ ) and  $\text{normalize} : \mathbb{R}^S \rightarrow \mathbb{R}^S$  are defined as follows:

$$\text{NN}_{\iota}^{(\ell)}(\mathbf{h}) = W_{1,\iota}^{(\ell)} \cdot \text{ReLU}(W_{2,\iota}^{(\ell)} \cdot \text{ReLU}(W_{3,\iota}^{(\ell)} \cdot [\mathbf{h}; 1])), \quad (4)$$

$$\text{normalize}(h)_s = h_s - \max_{s' \in [S]} h_{s'}. \quad (5)$$

The dimensions of the weight matrices within the ConvNet are specified as follows:

$$\mathbf{W} = \left\{ \left\{ W_{1,\iota}^{(\ell)} \in \mathbb{R}^{S \times D} \right\}_{\iota \in [m^{(\ell)}]}, \left\{ W_{2,\iota}^{(\ell)} \in \mathbb{R}^{D \times D} \right\}_{\iota \in [m^{(\ell)}]}, \left\{ W_{3,\iota}^{(\ell)} \in \mathbb{R}^{D \times S_{\text{in}}^{(\ell)}} \right\}_{\iota \in [m^{(\ell)}]} \right\}_{\ell \in [L]}. \quad (6)$$

Furthermore, we denote  $\mu_{\text{NN}}^{\mathbf{W}}$  as the ConvNet classifier  $\mu_{\text{NN}}$  parameterized by  $\mathbf{W}$ . A schematic illustration of a ConvNet with 3 layers is provided in Figure 1(right).

**Remark 1** (An explanation of the ‘‘ConvNet’’ architecture). *We remark that the neural network layers described in (ConvNet) are different from the ‘‘convolution operations’’ typically seen in practice. The convolution operations used in practice involve computing the inner products between convolutional filters and image patches, whereas in (ConvNet) and (4), a point-wise product is employed instead. Despite this, these layers are still referred to as convolutional layers because the mapping from  $\{\mathbf{h}_v^{(\ell)}\}_{v \in \mathcal{V}^{(\ell)}}$  to  $\{\mathbf{q}_v^{(\ell)}\}_{v \in \mathcal{V}^{(\ell)}}$ , as per the first line of (ConvNet), preserves the translation-invariance property. Specifically, we use the same function  $\text{NN}_{\iota}^{(\ell)}$  across different inputs  $\mathbf{h}_{v_1}^{(\ell)}$  and  $\mathbf{h}_{v_2}^{(\ell)}$  as long as  $\iota(v_1) = \iota(v_2) = \iota$ .*

*Additionally, the ‘‘normalization operator’’ defined in (5) differs from commonly used ones. We adopt this specific form for technical reasons, to effectively control the approximation error.*

*Despite these differences from standard convolutional networks, (ConvNet) represents an iterative composition of convolutional layers, pooling layers, and normalization layers, aligning closely with the architecture of convolutional networks used in practice. Figure 1(right) shows the sequence of these operations in detail.*

**The ERM estimator.** In the classification task, we employ empirical risk minimization over ConvNets as outlined in the following equation:

$$\widehat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}} \left\{ \widehat{\mathbf{R}}(\mu_{\text{NN}}^{\mathbf{W}}) = \frac{1}{n} \sum_{i=1}^n \text{loss}(y_i, \mu_{\text{NN}}^{\mathbf{W}}(\cdot | \mathbf{x}_i)) \right\}, \quad (7)$$

where, for simplicity of analysis, we opt for the square loss rather than the more commonly used cross-entropy loss:

$$\text{loss}(y, \mu_{\text{NN}}^{\mathbf{W}}(\cdot | \mathbf{x})) = \sum_{s=1}^S \left( 1\{y = s\} - \mu_{\text{NN}}^{\mathbf{W}}(s | \mathbf{x}) \right)^2. \quad (8)$$

The parameter space for the ConvNets is defined as:

$$\mathcal{W}_{d, \underline{m}, L, S, D, B} := \left\{ \mathbf{W} \text{ as defined in (6)} : \|\mathbf{W}\| := \max_{j \in [3]} \max_{\ell \in [L]} \max_{\iota \in [m^{(\ell)}]} \|W_{j, \iota}^{(\ell)}\|_{\text{op}} \leq B \right\}. \quad (9)$$

We anticipate that the empirical risk minimizer,  $\mu_{\text{NN}}^{\widehat{\mathbf{W}}}$ , could learn the Bayes classifier  $\mu_*$ , as the global minimizer of the population risk over all conditional distributions yields the Bayes classifier:

$$\mu_*(\cdot | \cdot) = \arg \min_{\mu} \left\{ \mathbf{R}(\mu) = \mathbb{E}[\text{loss}(y, \mu(\cdot | \mathbf{x}))] \right\}.$$

In our theoretical analysis, we measure the discrepancy between  $\mu_{\text{NN}}^{\widehat{\mathbf{W}}}$  and  $\mu_*$  using the squared Euclidean distance:

$$\mathbf{D}_2^2(\mu, \mu_*) = \mathbb{E}_{\mathbf{x} \sim \mu_*} \left[ \sum_{s=1}^S \left( \mu(s | \mathbf{x}) - \mu_*(s | \mathbf{x}) \right)^2 \right]. \quad (10)$$

**Sample complexity bound.** The subsequent theorem establishes the bound of the  $\mathbf{D}_2^2$ -distance between the ConvNet estimator  $\mu_{\text{NN}}^{\widehat{\mathbf{W}}}$  and the true Bayes classifier  $\mu_*$ .

**Theorem 1** (Learning to classify using ConvNets). *Let Assumption 1 and 2 hold. Let  $\mathcal{W}_{d, \underline{m}, L, S, D, B}$  be the set defined as in Eq. (9), where  $D \geq S^2 K^2 d \cdot 3^L$  and  $B = \text{Poly}(d, S, K, 3^L, D)$ . Let  $\widehat{\mathbf{W}}$  be the empirical risk minimizer as in Eq. (7). Then with probability at least  $1 - \eta$ , we have*

$$\mathbf{D}_2^2(\mu_{\text{NN}}^{\widehat{\mathbf{W}}}, \mu_*) \leq C \cdot \left( \frac{S^4 K^4 d^2 \cdot 3^{2L}}{D^2} + \sqrt{\frac{LD(D + 2S + 1) \|\underline{m}\|_1 \log(d \|\underline{m}\|_1 DS \cdot 3^L) + \log(1/\eta)}{n}} \right). \quad (11)$$

The proof of Theorem 1 is detailed in Section E.

**Remark 2.** *To ensure the  $\mathbf{D}_2^2$ -distance is less than  $\epsilon^2$ , Theorem 1 requires to take*

$$D = \Theta(S^2 K^2 d 3^L / \epsilon), \quad n = \tilde{\Theta}(LS^4 K^4 d^2 3^{2L} \|\underline{m}\|_1 / \epsilon^6), \quad (12)$$

where  $\tilde{\Theta}$  hides a logarithmic factor  $\log(d \|\underline{m}\|_1 S \cdot 3^L / (\eta \epsilon))$ . The dependency on any of the parameters  $(S, d, 3^L, K, \|\underline{m}\|_1, \epsilon)$  could potentially be refined by imposing additional assumptions on the  $\psi$  functions or through a more detailed analysis of approximation and generalization. This question of improving rates remains open for future work.

Consider a simplified scenario where  $S = 2$ ,  $K$  is constant, and  $m^{(\ell)} = m \geq 3$  for each  $\ell \in [L]$ , leading to  $d = m^L$ . In this setup, the sample complexity gives

$$n = \tilde{\Theta}(L^2 m \cdot d^{2+2 \log_m 3} / \epsilon^6) \leq \tilde{\Theta}(L^2 m \cdot d^4 / \epsilon^6), \quad (13)$$

exhibiting a polynomial dependence on  $d$  and  $1/\epsilon$ . Such polynomial scaling aligns with existing literature (Poggio et al., 2017; Malach & Shalev-Shwartz, 2020; Schmidt-Hieber, 2020; Allen-Zhu & Li, 2022; Petrini et al., 2023), which indicates that learning hierarchical models using multi-layer networks avoids the curse of dimensionality. Theorem 1 serves as a warm-up result in the classification context. In Section 4, we aim to extend similar methodologies to address denoising problems, employing analogous proof strategies.

**Proof strategy.** The proof strategy begins by decomposing the squared distance between the learned model and the true Bayes classifier into approximation and generalization error terms. The generalization error is bounded using a standard chaining argument, leading to a rate of  $\tilde{O}(\sqrt{d_p/n})$ , where  $d_p$  denotes the number of ConvNet parameters. The focus then shifts to controlling the approximation error. This is done by first introducing the belief propagation and message passing algorithm for computing the Bayes classifier, and subsequently showing that ConvNets can approximate this algorithm effectively. A detailed outline of the proof strategy is provided in Appendix A.1, with the complete proof presented in Appendix E.

## 4 DENOISING AND DIFFUSION IN GHMS

In this section, we consider the denoising task within the GHM. Consider the joint distribution of noisy and clean covariates  $(z, x)$ , generated from the following:  $x \sim \mu_*$  represents the clean covariates, and  $z = x + g$  where  $g \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  denotes the independent isotropic Gaussian noise. For simplicity in notation and with a slight abuse of notations, we continue to refer to the joint distribution of  $(z, x)$  as  $\mu_*$ .

We consider a scenario where a set of iid samples  $\{(z_i, x_i)\}_{i \in [n]} \sim_{iid} \mu_*$  is drawn from the distribution. Our objective is to learn a denoiser  $m(z)$  from this dataset. With a suitable loss function, the optimal denoiser is the Bayes denoiser  $m_*(z) = \mathbb{E}_{(x,z) \sim \mu_*}[x|z]$ , which calculates the posterior expectation of  $x$  given  $z$ . We aim to examine the sample complexity of learning this denoiser through empirical risk minimization over the class of U-Nets. The approaches and results of this section closely align with those discussed in Section 3 on the classification task.

**The U-Net architecture.** We here introduce the U-Net architecture used for denoising, represented as  $m_{\text{NN}}(z) \in \mathbb{R}^d$  for input  $z \in \mathbb{R}^d$ . Initially, we set  $h_{\downarrow,v}^{(L)} = -(x - z_v)^2/2)_{x \in [S]} \in \mathbb{R}^S$  for  $v \in \mathcal{V}^{(L)}$  for each node  $v \in \mathcal{V}^{(L)}$ . The operational flow of the network unfolds as follows:

$$\begin{aligned} \mathbf{q}_{\downarrow,v}^{(\ell)} &= \text{NN}_{\downarrow,\ell(v)}^{(\ell)}(\text{normalize}(\mathbf{h}_{\downarrow,v}^{(\ell)})) \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ \mathbf{h}_{\downarrow,v}^{(\ell-1)} &= \sum_{v' \in \mathcal{C}(v)} \mathbf{q}_{\downarrow,v'}^{(\ell)} \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell-1)}, \\ \mathbf{u}_{\uparrow,v}^{(\ell)} &= \mathbf{b}_{\uparrow,\text{pa}(v)}^{(\ell-1)} \in \mathbb{R}^S, \quad (\text{with } \mathbf{b}_{\uparrow,r}^{(0)} = \mathbf{h}_{\downarrow,r}^{(0)}) & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \quad (\text{UNet}) \\ \mathbf{b}_{\uparrow,v}^{(\ell)} &= \text{NN}_{\uparrow,\ell(v)}^{(\ell)}(\text{normalize}(\mathbf{u}_{\uparrow,v}^{(\ell)} - \mathbf{q}_{\downarrow,v}^{(\ell)})) + \mathbf{h}_{\downarrow,v}^{(\ell)} \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ \mathbf{m}_{\text{NN}}(z)_v &= \sum_{s \in [S]} s \cdot \text{softmax}(\mathbf{b}_{\uparrow,v}^{(L)})_s, & v \in \mathcal{V}^{(L)}. \end{aligned}$$

The functions  $\{\text{NN}_{\downarrow,\ell}^{(\ell)}, \text{NN}_{\uparrow,\ell}^{(\ell)} : \mathbb{R}^S \rightarrow \mathbb{R}^S\}_{\ell \in [L]}$  and  $\text{normalize} : \mathbb{R}^S \rightarrow \mathbb{R}^S$  are defined as follows:

$$\text{NN}_{\diamond,\ell}^{(\ell)}(h)_s = W_{1,\diamond,\ell}^{(\ell)} \cdot \text{ReLU}(W_{2,\diamond,\ell}^{(\ell)} \cdot \text{ReLU}(W_{3,\diamond,\ell}^{(\ell)} \cdot [h; 1])), \quad s \in [S], \diamond \in \{\downarrow, \uparrow\}, \ell \in [L], \ell \in [m^{(\ell)}], \quad (14)$$

$$\text{normalize}(h)_s = h_s - \max_{s' \in [S]} h_{s'}. \quad (15)$$

The dimensions of the weight matrices within the U-Net are specified as follows:

$$\mathbf{W} = \left\{ \{W_{1,\diamond,\ell}^{(\ell)} \in \mathbb{R}^{S \times D}\}_{\ell \in [m^{(\ell)}]}, \{W_{2,\diamond,\ell}^{(\ell)} \in \mathbb{R}^{D \times D}\}_{\ell \in [m^{(\ell)}]}, \{W_{3,\diamond,\ell}^{(\ell)} \in \mathbb{R}^{D \times (S+1)}\}_{\ell \in [m^{(\ell)}]} \right\}_{\ell \in [L], \diamond \in \{\downarrow, \uparrow\}} \quad (16)$$

Furthermore, we denote  $m_{\text{NN}}^{\mathbf{W}}$  as the U-Net denoiser  $m_{\text{NN}}$  parameterized by  $\mathbf{W}$ . A schematic illustration of a U-Net with 3 layers is provided in Figure 2.

**Remark 3** (An explanation of the ‘‘U-Net’’ architecture). *As noted in our discussion of the convolutional network as in Remark 1, the ‘‘convolutional layers’’ and the ‘‘normalization operator’’ described in (UNet) are different from practical implementations. However, we continue to use these terms because they retain core characteristics of their practical counterparts.*

*An important feature of (UNet) is its encoder-decoder architecture and the inclusion of long skip connections, which closely mirror practical implementations. Specifically, the encoder sequence in (UNet) progresses as  $h_{\downarrow}^{(L)} \rightarrow \mathbf{q}_{\downarrow}^{(L)} \rightarrow h_{\downarrow}^{(L-1)} \rightarrow \dots \rightarrow \mathbf{q}_{\downarrow}^{(1)} \rightarrow h_{\downarrow}^{(0)}$ , consisting of a series*

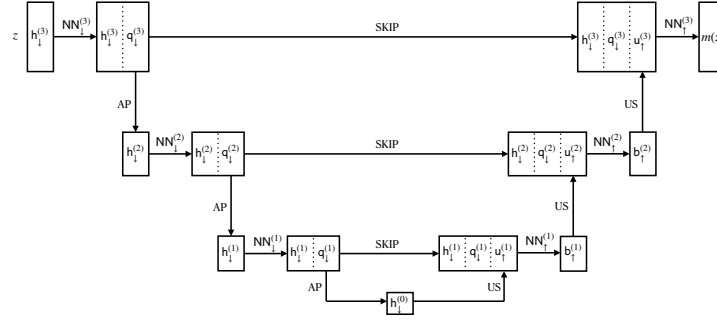


Figure 2: A U-Net with  $L = 3$ . “AP” stands for average-pooling. “US” stands for up-sampling. “SKIP” stands for long skip connections.

of convolutional, average pooling, and normalization layers. This part of the architecture is the same as (ConvNet) for the classification task. The decoder sequence ascends as  $b_{\uparrow}^{(0)} \rightarrow u_{\uparrow}^{(1)} \rightarrow b_{\uparrow}^{(1)} \rightarrow \dots \rightarrow b_{\uparrow}^{(L-1)} \rightarrow u_{\uparrow}^{(L)} \rightarrow \mathbf{m}_{\text{NN}}$ , consisting of a series of convolutional, up-sampling, and normalization layers. Moreover, the computation of  $b_{\uparrow}^{(\ell)}$  utilizes  $u_{\uparrow}^{(\ell)}$  from the upward sequence, and  $(h_{\downarrow}^{(\ell)}, q_{\downarrow}^{(\ell)})$  from the downward process, which is enabled by the long skip connections. This encoder-decoder architecture and the long skip connections are effectively visualized in Figure 2.

**The ERM estimator.** In the denoising task, we employ empirical risk minimization over U-Nets as outlined in the following equation:

$$\widehat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}} \left\{ \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) = \frac{1}{nd} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z}_i)\|_2^2 \right\}. \quad (17)$$

The parameter space for the U-Nets is defined as:

$$\mathcal{W}_{d, \underline{m}, L, S, D, B} := \left\{ \mathbf{W} \text{ as defined in (16)} : \|\mathbf{W}\| := \max_{\diamond \in \{\downarrow, \uparrow\}} \max_{j \in [3]} \max_{\ell \in [L]} \max_{\iota \in [m^{(\ell)}]} \|W_{j, \iota}^{(\ell)}\|_{\text{op}} \leq B \right\}. \quad (18)$$

We anticipate that the empirical risk minimizer,  $\mathbf{m}_{\text{NN}}^{\widehat{\mathbf{W}}}$ , could learn the Bayes denoiser  $\mathbf{m}_{\star}$ , as the global minimizer of the population risk over all functions yields the Bayes denoiser:

$$\mathbf{m}_{\star}(\cdot) = \arg \min_{\mathbf{m}} \left\{ \mathbf{R}(\mathbf{m}) = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mu_{\star}} [d^{-1} \|\mathbf{x} - \mathbf{m}(\mathbf{z})\|_2^2] \right\}.$$

In our theoretical analysis, we measure the discrepancy between  $\mathbf{m}_{\text{NN}}^{\widehat{\mathbf{W}}}$  and  $\mathbf{m}_{\star}$  using the squared Euclidean distance:

$$D_2^2(\mathbf{m}, \mathbf{m}_{\star}) = \mathbb{E}_{\mathbf{z} \sim \mu_{\star}} \left[ d^{-1} \|\mathbf{m} - \mathbf{m}_{\star}(\mathbf{z})\|_2^2 \right]. \quad (19)$$

**Sample complexity bound.** The subsequent theorem establishes the bound of the  $D_2^2$ -distance between the U-Net estimator  $\mathbf{m}_{\text{NN}}^{\widehat{\mathbf{W}}}$  and the true Bayes denoiser  $\mathbf{m}_{\star}$ .

**Theorem 2** (Learning to denoise using U-Nets). *Let Assumption 1 and 2 hold. Let  $\mathcal{W}_{d, \underline{m}, L, S, D, B}$  be the set defined as in Eq. (18), where  $B = \text{Poly}(d, S, K, 18^L, D)$ . Let  $\widehat{\mathbf{W}}$  be the empirical risk minimizer as in Eq. (17). Then with probability at least  $1 - \eta$ , we have*

$$D_2^2(\mathbf{m}_{\text{NN}}^{\widehat{\mathbf{W}}}, \mathbf{m}_{\star}) \leq C \cdot \left( \frac{S^6 K^4 d^2 \cdot 18^{2L}}{D^2} + S^2 \cdot \sqrt{\frac{LD(D + 2S + 1) \|\underline{m}\|_1 \log(d \|\underline{m}\|_1 DS \cdot 18^L) + \log(1/\eta)}{n}} \right).$$

The proof of Theorem 2 is detailed in Section H.

**Remark 4.** To ensure the  $D_2^2$ -distance is less than  $\epsilon^2$ , Theorem 2 requires to take

$$D = \Theta(S^3 K^2 d 18^L / \epsilon), \quad n = \tilde{\Theta}(LS^{10} K^4 d^2 18^{2L} \|\underline{m}\|_1 / \epsilon^6), \quad (20)$$

where  $\tilde{\Theta}$  hides a logarithmic factor  $\log(d \|\underline{m}\|_1 S \cdot 18^L / (\eta \epsilon))$ . Similar to Theorem 1, the dependency on any of the parameters  $(S, d, 18^L, K, \|\underline{m}\|_1, \epsilon)$  could potentially be refined by imposing additional assumptions on the  $\psi$  functions or through a more detailed analysis of approximation and generalization. This question of improving rates remains open for future work.

Consider a simplified scenario where  $S = 2$ ,  $K$  is constant, and  $m^{(\ell)} = m \geq 3$  for each  $\ell \in [L]$ , leading to  $d = m^L$ . In this setup, the sample complexity gives

$$n = \tilde{\Theta}(L^2 m \cdot d^{2+2 \log_m 18} / \epsilon^6) \leq \tilde{\Theta}(L^2 m \cdot d^8 / \epsilon^6), \quad (21)$$

exhibiting a polynomial dependence on  $d$  and  $1/\epsilon$ . Although the degree of the polynomial is substantial and potentially improvable, this gives the first polynomial sample complexity result for learning the Bayes denoiser in a hierarchical model using the U-Nets.

**Connection to the task of diffusion generative modeling.** The denoising task examined in this section is closely related to the diffusion model approach (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song & Ermon, 2019; Song et al., 2020) to generative modeling. Diffusion models involve learning a generative model from a dataset of  $n$  independent and identically distributed samples  $\{\mathbf{x}_i\}_{i=1}^n$ , drawn from an unknown data distribution  $\mu_* \in \mathcal{P}([S]^d)$ . The goal is to generate new samples  $\hat{\mathbf{x}} \sim \hat{\mu}$  that match the distribution  $\mu_*$ . Most diffusion model formulations involve a series of steps that are closely related to the denoising task described earlier. Here, we illustrate how they work using a variant of diffusion models, the stochastic localization process (Eldan, 2013; El Alaoui et al., 2022; Montanari & Wu, 2023; Celentano, 2022; Montanari, 2023):

- **Step 1.** Fit approximate denoising functions  $\hat{\mathbf{m}}_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$  for  $t \in [0, T]$ . This is done by minimizing the empirical risk over a class of neural networks  $\mathcal{F}$  (i.e., the denoising task discussed in this section):

$$\hat{\mathbf{m}}_t \equiv \arg \min_{\text{NN}_t \in \mathcal{F}} \frac{1}{nd} \sum_{i=1}^n \|\mathbf{x}_i - \text{NN}_t(t \cdot \mathbf{x}_i + \sqrt{t} \cdot \mathbf{g}_i)\|_2^2, \quad \mathbf{g}_i \sim_{iid} \mathcal{N}(\mathbf{0}, \mathbf{I}_d). \quad (\text{ERM})$$

- **Step 2.** Simulate a discretized version of a stochastic differential equation (SDE) starting from zero, whose drift term gives the approximate denoising function:

$$d\mathbf{z}_t = \hat{\mathbf{m}}_t(\mathbf{z}_t) \cdot dt + d\mathbf{B}_t, \quad t \in [0, T], \quad \mathbf{z}_0 = \mathbf{0}, \quad (\text{SDE})$$

and generate an approximate sample  $\hat{\mathbf{x}} = \mathbf{z}_T / T \in \mathbb{R}^d$  at the final time  $T$ .

Standard analysis shows that by replacing the fitted denoising functions  $\hat{\mathbf{m}}_t(\mathbf{z})$  with the true denoising functions  $\mathbf{m}_t(\mathbf{z})$  in Eq. (SDE) and allowing  $T \rightarrow \infty$ , we can effectively recover the original data distribution  $\mu_*$ . Consequently, the quality of samples generated from diffusion models hinges on two critical factors: (1) How well the fitted denoising functions  $\hat{\mathbf{m}}_t$  (ERM) approximate the true denoising functions  $\mathbf{m}_t$ ; (2) How accurately the SDE discretization scheme approximates the continuous process as in (SDE).

Recent work has made substantial progress in addressing these two theoretical questions: controlling the SDE discretization error, assuming a reliable denoising function estimator is available (Chen et al., 2022a; 2023a; Lee et al., 2023; Li et al., 2023; Benton et al., 2023); and controlling the denoising function approximation error through neural networks (Oko et al., 2023; Chen et al., 2023b; Mei & Wu, 2023). However, these works have not explained the benefit of employing U-Net in image diffusion modeling, which is the primary focus of the current work. Indeed, by integrating the sample complexity bounds for learning denoising functions, as established in Theorem 2, with standard SDE discretization error bounds, such as the result established in Benton et al. (2023), it is straightforward to derive an end-to-end error bound for the sampling process of diffusion models in GHMs, similar to the strategy of Oko et al. (2023); Chen et al. (2023b); Mei & Wu (2023)<sup>2</sup>.

<sup>2</sup>We note that the stochastic localization formulation is equivalent to the DDPM diffusion model, differing only in parametrization (Montanari, 2023). In the DDPM model, U-Nets serve to approximate the score function. The score function is a linear combination of the denoising function with an identity map, as per Tweedie’s formula. Consequently, Theorem 2 can be readily adapted to establish a sample complexity bound for learning this score function.



**Proof strategy.** The proof strategy for the denoising task parallels that of the classification task. The squared distance between the learned model and the true Bayes denoiser is decomposed into approximation and generalization error terms. The generalization error is bounded via a standard parameter counting argument. The approximation error is controlled by first introducing the belief propagation and message passing algorithm for computing the Bayes denoiser and then showing that U-Nets can effectively approximate this algorithm. A detailed outline of the proof strategy is provided in Appendix A.2, with the complete proof presented in Appendix H.

## 5 FURTHER RELATED WORK

**Generative hierarchical models.** Hierarchical modeling of data distributions has been proposed in a series of works (Mossel, 2016; Poggio et al., 2017; Malach & Shalev-Shwartz, 2020; Schmidt-Hieber, 2020; Allen-Zhu & Li, 2022; Petrini et al., 2023; Sclocchi et al., 2024; Tomasini & Wyart, 2024; Cagnetta & Wyart, 2024; Garnier-Brun et al., 2024; Kadkhodaie et al., 2023a;b). While the hierarchical models in Poggio et al. (2017); Malach & Shalev-Shwartz (2020); Schmidt-Hieber (2020); Allen-Zhu & Li (2022) remain deterministic, Mossel (2016); Petrini et al. (2023); Sclocchi et al. (2024); Tomasini & Wyart (2024); Cagnetta & Wyart (2024); Garnier-Brun et al. (2024) studied the generative version of hierarchical models. The diffusion model for multi-scale image distribution representations has been empirically examined in Kadkhodaie et al. (2023a;b), which demonstrated that U-Nets are effective in modeling denoising algorithms. The theoretical and empirical evidence presented in Sclocchi et al. (2024); Tomasini & Wyart (2024); Petrini et al. (2023) underscores the effectiveness of generative hierarchical models in capturing the combinatorial properties of image datasets. Given their significant relevance to this work, we delve deeper into these studies.

**Contributions of Petrini et al. (2023); Sclocchi et al. (2024); Tomasini & Wyart (2024).** The series of works on hierarchical generative models (Petrini et al., 2023; Sclocchi et al., 2024; Tomasini & Wyart, 2024) inspired the current study. Sclocchi et al. (2024) first pointed out that the belief propagation denoising algorithm of hierarchical models consists of downward and upward processes. Through mean-field analysis on a random generative hierarchical model, they identified a phase transition phenomenon, aligning with empirical observations in diffusion models, thereby providing strong evidence of the efficacy of these models in handling combinatorial data properties. Petrini et al. (2023), on the other hand, first introduced these models in a classification context. Petrini et al. (2023); Tomasini & Wyart (2024) demonstrated that learning hierarchical models using multi-layer networks circumvents the curse of dimensionality. Specifically, they theoretically and empirically characterized the sample complexity, showing that it remains polynomial in dimension when learning convolutional networks under random generative rules. On the other hand, in the absence of correlations, they showed that the sample complexity is again exponential in the dimension, even for hierarchical generative models. This learning incapability is not captured by our analysis which does not consider optimization.

**ConvNets and U-Nets and their implicit bias.** Convolutional networks (LeCun et al., 1989; 1998; Krizhevsky et al., 2012; Szegedy et al., 2015; He et al., 2016) have become the state-of-the-art architecture for image classification and have been the backbone for many computer vision tasks. U-Nets (Ronneberger et al., 2015; Zhou et al., 2018; Siddique et al., 2021; Oktay et al., 2018) have been particularly well-suited for image segmentation and denoising tasks (Ronneberger et al., 2015), and have served as the backbone architecture for diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song & Ermon, 2019; Song et al., 2020). A series of theoretical works has explained the inductive bias of CNNs (Bruna & Mallat, 2013; Gunasekar et al., 2018; Bietti & Mairal, 2019b;a; Scetbon & Harchaoui, 2020; Li et al., 2020; Bietti, 2021; Mei et al., 2021; Misiakiewicz & Mei, 2022; Cagnetta et al., 2023; Favero et al., 2021; Bietti et al., 2021; Xiao, 2022; Petrini et al., 2023; Tomasini & Wyart, 2024; Wang & Wu, 2024). However, they mostly focused on the classification and regression setting and were not concerned with the role of U-Nets in denoising tasks. The implicit bias of U-Nets has been theoretically investigated in Williams et al. (2024); Falck et al. (2022), where they found that the U-Nets are conjugate to the ResNets. In contrast, we demonstrate that U-Nets can effectively approximate the belief propagation denoising algorithms of GHMs. We note that Cui et al. (2023) analyzed the learning dynamics for a simple U-Net in diffusion models.

**Neural networks approximation of algorithms.** A recent line of work has investigated the expressiveness of neural networks through an algorithm approximation viewpoint (Wei et al., 2022; Bai et al., 2024; Giannou et al., 2023; Liu et al., 2022a; Marwah et al., 2021; 2023; Lin et al., 2023; Mei & Wu, 2023). In particular, Wei et al. (2022); Bai et al. (2024); Giannou et al. (2023); Liu et al. (2022a); Lin et al. (2023) demonstrate that transformers can efficiently approximate several algorithm classes, such as gradient descent, reinforcement learning algorithms, and even Turing machines. In the context of diffusion models, Mei & Wu (2023) shows that ResNets can efficiently approximate the score function of high-dimensional graphical models by approximating the variational inference algorithm. Our work is closely related to Mei & Wu (2023), except that we study neural network approximation in a different statistical model and network architecture.

From a practical viewpoint, a line of work has focused on neural network denoising by unrolling iterative denoising algorithms into deep networks (Gregor & LeCun, 2010; Zheng et al., 2015; Zhang & Ghanem, 2018; Pappan et al., 2017; Ma et al., 2021; Chen et al., 2018; Borgerding et al., 2017; Monga et al., 2021; Yu et al., 2023a;b). While this literature has primarily focused on devising better denoising algorithms, our work leverages this perspective to develop neural network approximation theory and explain existing network architectures.

**Related theory of diffusion models.** In recent years, diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song & Ermon, 2019; Song et al., 2020) have emerged as a leading approach for generative modeling. Neural network-based score function approximation has been recently studied from the function approximation viewpoint in Oko et al. (2023); Chen et al. (2023b); Yuan et al. (2023); Shah et al. (2023); Biroli & Mézard (2023), and from the algorithm approximation viewpoint in Mei & Wu (2023). Theoretical studies of other aspects of diffusion models include Liu et al. (2022b); Li et al. (2023); Lee et al. (2023); Chen et al. (2022b; 2023d; 2022a; 2023c;a); Benton et al. (2023); El Alaoui et al. (2022); Montanari & Wu (2023); Celentano (2022); Ghio et al. (2023); Biroli & Mézard (2023); Biroli et al. (2024); Cui et al. (2023); Fu et al. (2024); Wu et al. (2024). For a comprehensive introduction to the theory of diffusion models, see the recent review (Chen et al., 2024).

## 6 CONCLUSIONS AND DISCUSSIONS

In this paper, we introduced a novel interpretation of the U-Net architecture through the lens of generative hierarchical models. We demonstrated that their belief propagation denoising algorithm naturally induces the encoder-decoder structure, the long skip connections, and the pooling and up-sampling operations of the U-Nets. We also provided an efficient sample complexity bound for learning the denoising function with U-Nets. Furthermore, we discussed the broader implications of these findings for diffusion models. We also showed that ConvNets are well-suited for classification tasks within these models. Our study offers a unified perspective on the roles of ConvNets and U-Nets, highlighting the versatility of generative hierarchical models in capturing complex data distributions across language and image domains.

The results presented in this paper offer considerable scope for enhancement. We initially assumed that the covariates  $x$  lie in the discrete space  $[S]^d$ , and extending these results to continuous spaces would be an intriguing direction for future research. Additionally, the dependencies of the sample complexity bound on  $d$  and  $1/\epsilon$  may be amenable to improvement through more careful analysis. Moreover, the convolution operations employed in this paper are different from those commonly employed in practical settings. It would be worthwhile to explore graphical models where the belief propagation algorithm aligns more naturally with ConvNets and U-Nets that utilize standard convolution operations.

On the practical side, our theoretical findings generated a hypothesis of the functionality of each layer of the U-Nets. Verifying these hypotheses in pre-trained U-Nets, such as those used in stable diffusion models, using interpretability methods, could yield valuable insights. Furthermore, extending these results to include conditional denoising functions represents an exciting direction for future research. Finally, we hope that the insights provided in this paper could guide the design of innovative network architectures.

## REFERENCES

- 540  
541  
542 Zeyuan Allen-Zhu and Yuanzhi Li. How can deep learning performs deep (hierarchical) learning.  
543 2022.
- 544 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar. *arXiv*  
545 *preprint arXiv:2305.13673*, 2023.
- 546  
547 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:  
548 Provable in-context learning with in-context algorithm selection. *Advances in neural information*  
549 *processing systems*, 36, 2024.
- 550 Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Linear convergence  
551 bounds for diffusion models via stochastic localization. *arXiv preprint arXiv:2308.03686*, 2023.
- 552  
553 Alberto Bietti. Approximation and learning with deep convolutional models: a kernel perspective.  
554 *arXiv preprint arXiv:2102.10032*, 2021.
- 555  
556 Alberto Bietti and Julien Mairal. Group invariance, stability to deformations, and complexity of  
557 deep convolutional representations. *Journal of Machine Learning Research*, 20(25):1–49, 2019a.
- 558  
559 Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. *Advances in Neural*  
*Information Processing Systems*, 32, 2019b.
- 560  
561 Alberto Bietti, Luca Venturi, and Joan Bruna. On the sample complexity of learning under geometric  
562 stability. *Advances in neural information processing systems*, 34:18673–18684, 2021.
- 563  
564 Giulio Biroli and Marc Mézard. Generative diffusion in very large dimensions. *arXiv preprint*  
*arXiv:2306.03518*, 2023.
- 565  
566 Giulio Biroli, Tony Bonnaire, Valentin de Bortoli, and Marc Mézard. Dynamical regimes of diffu-  
567 sion models. *arXiv preprint arXiv:2402.18491*, 2024.
- 568  
569 Mark Borgerding, Philip Schniter, and Sundeep Rangan. Amp-inspired deep networks for sparse  
linear inverse problems. *IEEE Transactions on Signal Processing*, 65(16):4293–4308, 2017.
- 570  
571 Joan Bruna and Stéphane Mallat. Invariant scattering convolution networks. *IEEE transactions on*  
572 *pattern analysis and machine intelligence*, 35(8):1872–1886, 2013.
- 573  
574 Francesco Cagnetta and Matthieu Wyart. Towards a theory of how the structure of language is  
acquired by deep neural networks. *arXiv preprint arXiv:2406.00048*, 2024.
- 575  
576 Francesco Cagnetta, Alessandro Favero, and Matthieu Wyart. What can be learnt with wide con-  
577 volutional neural networks? In *International Conference on Machine Learning*, pp. 3347–3379.  
578 PMLR, 2023.
- 579  
580 Michael Celentano. Sudakov-ferniqie post-amp, and a new proof of the local convexity of the tap  
free energy. *arXiv preprint arXiv:2208.09550*, 2022.
- 581  
582 Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling:  
583 User-friendly bounds under minimal smoothness assumptions. In *International Conference on*  
584 *Machine Learning*, pp. 4735–4763. PMLR, 2023a.
- 585  
586 Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estima-  
587 tion and distribution recovery of diffusion models on low-dimensional data. *arXiv preprint*  
*arXiv:2302.07194*, 2023b.
- 588  
589 Minshuo Chen, Song Mei, Jianqing Fan, and Mengdi Wang. An overview of diffusion models: Ap-  
590 plications, guided generation, statistical rates and optimization. *arXiv preprint arXiv:2404.07771*,  
591 2024.
- 592  
593 Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy  
as learning the score: theory for diffusion models with minimal data assumptions. *arXiv preprint*  
*arXiv:2209.11215*, 2022a.

- 594 Sitan Chen, Sinho Chewi, Holden Lee, Yuanzhi Li, Jianfeng Lu, and Adil Salim. The probability  
595 flow ode is provably fast. *arXiv preprint arXiv:2305.11798*, 2023c.
- 596
- 597 Sitan Chen, Giannis Daras, and Alex Dimakis. Restoration-degradation beyond linear diffusions:  
598 A non-asymptotic analysis for ddim-type samplers. In *International Conference on Machine*  
599 *Learning*, pp. 4462–4484. PMLR, 2023d.
- 600 Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin. Theoretical linear convergence of  
601 unfolded ista and its practical weights and thresholds. *Advances in Neural Information Processing*  
602 *Systems*, 31, 2018.
- 603
- 604 Yongxin Chen, Sinho Chewi, Adil Salim, and Andre Wibisono. Improved analysis for a proximal  
605 algorithm for sampling. In *Conference on Learning Theory*, pp. 2984–3014. PMLR, 2022b.
- 606
- 607 Noam Chomsky. On certain formal properties of grammars. *Information and control*, 2(2):137–167,  
608 1959.
- 609 Hugo Cui, Florent Krzakala, Eric Vanden-Eijnden, and Lenka Zdeborová. Analysis of learning a  
610 flow-based generative model from limited sample complexity. *arXiv preprint arXiv:2310.03575*,  
611 2023.
- 612
- 613 Michal Drozdal, Eugene Vorontsov, Gabriel Chartrand, Samuel Kadoury, and Chris Pal. The im-  
614 portance of skip connections in biomedical image segmentation. In *International workshop on*  
615 *deep learning in medical image analysis, international workshop on large-scale annotation of*  
616 *biomedical data and expert label synthesis*, pp. 179–187. Springer, 2016.
- 617 Ahmed El Alaoui, Andrea Montanari, and Mark Sellke. Sampling from the sherrington-kirkpatrick  
618 gibbs measure via algorithmic stochastic localization. In *2022 IEEE 63rd Annual Symposium on*  
619 *Foundations of Computer Science (FOCS)*, pp. 323–334. IEEE, 2022.
- 620
- 621 Ronen Eldan. Thin shell implies spectral gap up to polylog via a stochastic localization scheme.  
622 *Geometric and Functional Analysis*, 23(2):532–569, 2013.
- 623 Fabian Falck, Christopher Williams, Dominic Danks, George Deligiannidis, Christopher Yau,  
624 Chris C Holmes, Arnaud Doucet, and Matthew Willetts. A multi-resolution framework for u-  
625 nets with applications to hierarchical vaes. *Advances in Neural Information Processing Systems*,  
626 35:15529–15544, 2022.
- 627
- 628 Alessandro Favero, Francesco Cagnetta, and Matthieu Wyart. Locality defeats the curse of dimen-  
629 sionality in convolutional teacher-student scenarios. *Advances in Neural Information Processing*  
630 *Systems*, 34:9456–9467, 2021.
- 631 Hengyu Fu, Zhuoran Yang, Mengdi Wang, and Minshuo Chen. Unveil conditional diffusion models  
632 with classifier-free guidance: A sharp statistical theory. *arXiv preprint arXiv:2403.11968*, 2024.
- 633
- 634 Jérôme Garnier-Brun, Marc Mézard, Emanuele Moscatò, and Luca Saglietti. How transformers  
635 learn structured data: insights from hierarchical filtering. *arXiv preprint arXiv:2408.15138*, 2024.
- 636 Davide Ghio, Yatin Dandi, Florent Krzakala, and Lenka Zdeborová. Sampling with flows, diffusion  
637 and autoregressive neural networks: A spin-glass perspective. *arXiv preprint arXiv:2308.14085*,  
638 2023.
- 639
- 640 Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dim-  
641 itris Papailiopoulos. Looped transformers as programmable computers. *arXiv preprint*  
642 *arXiv:2301.13196*, 2023.
- 643 Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In *Proceedings of*  
644 *the 27th international conference on international conference on machine learning*, pp. 399–406,  
645 2010.
- 646
- 647 Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent  
on linear convolutional networks. *Advances in neural information processing systems*, 31, 2018.

- 648 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
649 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
650 770–778, 2016.
- 651 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
652 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 653 Ya Jin and Stuart Geman. Context and hierarchy in a probabilistic image model. In *2006 IEEE*  
654 *computer society conference on computer vision and pattern recognition (CVPR'06)*, volume 2,  
655 pp. 2145–2152. IEEE, 2006.
- 656 Zahra Kadkhodaie, Florentin Guth, Stéphane Mallat, and Eero P Simoncelli. Learning multi-scale  
657 local conditional probability models of images. *arXiv preprint arXiv:2303.02984*, 2023a.
- 658 Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization  
659 in diffusion models arises from geometry-adaptive harmonic representation. *arXiv preprint*  
660 *arXiv:2310.02557*, 2023b.
- 661 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-  
662 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 663 Yann LeCun, Bernhard Boser, John Denker, Donnie Henderson, Richard Howard, Wayne Hubbard,  
664 and Lawrence Jackel. Handwritten digit recognition with a back-propagation network. *Advances*  
665 *in neural information processing systems*, 2, 1989.
- 666 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
667 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 668 Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for  
669 general data distributions. In *International Conference on Algorithmic Learning Theory*, pp.  
670 946–985. PMLR, 2023.
- 671 Lillian Lee. Learning of context-free languages: A survey of the literature. 1996.
- 672 Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Towards faster non-asymptotic convergence for  
673 diffusion-based generative models. *arXiv preprint arXiv:2306.09251*, 2023.
- 674 Jia Li, Robert M Gray, and Richard A Olshen. Multiresolution image classification by hierarchical  
675 modeling with two-dimensional hidden markov models. *IEEE transactions on information theory*,  
676 46(5):1826–1841, 2000.
- 677 Zhiyuan Li, Yi Zhang, and Sanjeev Arora. Why are convolutional nets more sample-efficient than  
678 fully-connected nets? *arXiv preprint arXiv:2010.08515*, 2020.
- 679 Licong Lin, Yu Bai, and Song Mei. Transformers as decision makers: Provable in-context reinforce-  
680 ment learning via supervised pretraining. *arXiv preprint arXiv:2310.08566*, 2023.
- 681 Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers  
682 learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022a.
- 683 Xingchao Liu, Lemeng Wu, Mao Ye, and Qiang Liu. Let us build bridges: Understanding and  
684 extending diffusion generative models. *arXiv preprint arXiv:2208.14699*, 2022b.
- 685 Yao Ma, Xiaorui Liu, Tong Zhao, Yozen Liu, Jiliang Tang, and Neil Shah. A unified view on  
686 graph neural networks as graph signal denoising. In *Proceedings of the 30th ACM International*  
687 *Conference on Information & Knowledge Management*, pp. 1202–1211, 2021.
- 688 Eran Malach and Shai Shalev-Shwartz. The implications of local correlation on learning some deep  
689 functions. *Advances in Neural Information Processing Systems*, 33:1322–1332, 2020.
- 690 Tanya Marwah, Zachary Lipton, and Andrej Risteski. Parametric complexity bounds for approx-  
691 imating pdes with neural networks. *Advances in Neural Information Processing Systems*, 34:  
692 15044–15055, 2021.

- 702 Tanya Marwah, Zachary Chase Lipton, Jianfeng Lu, and Andrej Risteski. Neural network approxi-  
703 mations of pdes beyond linearity: A representational perspective. In *International Conference on*  
704 *Machine Learning*, pp. 24139–24172. PMLR, 2023.
- 705 Song Mei and Yuchen Wu. Deep networks as denoising algorithms: Sample-efficient learning of  
706 diffusion models in high-dimensional graphical models. *arXiv preprint arXiv:2309.11420*, 2023.
- 707  
708 Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Learning with invariances in random  
709 features and kernel models. In *Conference on Learning Theory*, pp. 3351–3418. PMLR, 2021.
- 710  
711 Marc Mezard and Andrea Montanari. *Information, physics, and computation*. Oxford University  
712 Press, 2009.
- 713 Theodor Misiakiewicz and Song Mei. Learning with convolution and pooling operations in kernel  
714 methods. *Advances in Neural Information Processing Systems*, 35:29014–29025, 2022.
- 715 Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep  
716 learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):18–44, 2021.
- 717  
718 Andrea Montanari. Sampling, diffusions, and stochastic localization. *arXiv preprint*  
719 *arXiv:2305.10690*, 2023.
- 720  
721 Andrea Montanari and Yuchen Wu. Posterior sampling from the spiked models via diffusion pro-  
722 cesses. *arXiv preprint arXiv:2304.11449*, 2023.
- 723 Elchanan Mossel. Deep learning and hierarchal generative models. *arXiv preprint*  
724 *arXiv:1612.09057*, 2016.
- 725  
726 Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribu-  
727 tion estimators. *arXiv preprint arXiv:2303.01861*, 2023.
- 728 Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa,  
729 Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention u-net:  
730 Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*, 2018.
- 731  
732 Vardan Papyan, Yaniv Romano, and Michael Elad. Convolutional neural networks analyzed via con-  
733 volutional sparse coding. *The Journal of Machine Learning Research*, 18(1):2887–2938, 2017.
- 734  
735 Judea Pearl. Reverend bayes on inference engines: A distributed hierarchical approach. In *Proba-  
736 bilistic and Causal Inference: The Works of Judea Pearl*, pp. 129–138. 1982.
- 737  
738 Leonardo Petrini, Francesco Cagnetta, Umberto M Tomasini, Alessandro Favero, and Matthieu  
739 Wyart. How deep neural networks learn compositional data: The random hierarchy model. *arXiv*  
740 *preprint arXiv:2307.02129*, 2023.
- 741  
742 Tomaso Poggio, Hrushikesh Mhaskar, Lorenzo Rosasco, Brando Miranda, and Qianli Liao. Why  
743 and when can deep-but not shallow-networks avoid the curse of dimensionality: a review. *Inter-  
744 national Journal of Automation and Computing*, 14(5):503–519, 2017.
- 745  
746 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-  
747 ical image segmentation. In *Medical image computing and computer-assisted intervention-  
748 MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceed-  
749 ings, part III 18*, pp. 234–241. Springer, 2015.
- 750  
751 Meyer Scetbon and Zaid Harchaoui. Harmonic decompositions of convolutional networks. In *Inter-  
752 national Conference on Machine Learning*, pp. 8522–8532. PMLR, 2020.
- 753  
754 Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation  
755 function. 2020.
- 756  
757 Antonio Sclocchi, Alessandro Favero, and Matthieu Wyart. A phase transition in diffusion models  
758 reveals the hierarchical nature of data. *arXiv preprint arXiv:2402.16991*, 2024.
- 759  
760 Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of gaussians using the ddpm objec-  
761 tive. *arXiv preprint arXiv:2307.01178*, 2023.

- 756 Nahian Siddique, Sidike Paheding, Colin P Elkin, and Vijay Devabhaktuni. U-net and its variants for  
757 medical image segmentation: A review of theory and applications. *Ieee Access*, 9:82031–82057,  
758 2021.
- 759 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
760 learning using nonequilibrium thermodynamics. In *International Conference on Machine Learn-*  
761 *ing*, pp. 2256–2265. PMLR, 2015.
- 762 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.  
763 *Advances in neural information processing systems*, 32, 2019.
- 764 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
765 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*  
766 *arXiv:2011.13456*, 2020.
- 767 Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Du-  
768 mitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In  
769 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- 770 Umberto Tomasini and Matthieu Wyart. How deep networks learn sparse and hierarchical data: the  
771 sparse random hierarchy model. *arXiv preprint arXiv:2404.10727*, 2024.
- 772 Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series  
773 in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. doi: 10.1017/  
774 9781108627771.
- 775 Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational  
776 inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- 777 Haonan Wang, Peng Cao, Jiaqi Wang, and Osmar R Zaiane. Uctransnet: rethinking the skip con-  
778 nections in u-net from a channel-wise perspective with transformer. In *Proceedings of the AAAI*  
779 *conference on artificial intelligence*, volume 36, pp. 2441–2449, 2022.
- 780 Zihao Wang and Lei Wu. Theoretical analysis of the inductive biases in deep convolutional networks.  
781 *Advances in Neural Information Processing Systems*, 36, 2024.
- 782 Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on  
783 approximating turing machines with transformers. *Advances in Neural Information Processing*  
784 *Systems*, 35:12071–12083, 2022.
- 785 Christopher Williams, Fabian Falck, George Deligiannidis, Chris C Holmes, Arnaud Doucet, and  
786 Saifuddin Syed. A unified framework for u-net design and analysis. *Advances in Neural Infor-*  
787 *mation Processing Systems*, 36, 2024.
- 788 Alan S Willsky. Multiresolution markov models for signal and image processing. *Proceedings of*  
789 *the IEEE*, 90(8):1396–1458, 2002.
- 790 Yuchen Wu, Minshuo Chen, Zihao Li, Mengdi Wang, and Yuting Wei. Theoretical insights for  
791 diffusion guidance: A case study for gaussian mixture models. *arXiv preprint arXiv:2403.01639*,  
792 2024.
- 793 Lechao Xiao. Eigenspace restructuring: a principle of space and frequency in neural networks. In  
794 *Conference on Learning Theory*, pp. 4888–4944. PMLR, 2022.
- 795 Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin D  
796 Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *arXiv preprint*  
797 *arXiv:2306.01129*, 2023a.
- 798 Yaodong Yu, Tianzhe Chu, Shengbang Tong, Ziyang Wu, Druv Pai, Sam Buchanan, and  
799 Yi Ma. Emergence of segmentation with minimalistic white-box transformers. *arXiv preprint*  
800 *arXiv:2308.16271*, 2023b.
- 801 Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Reward-directed  
802 conditional diffusion: Provable distribution estimation and reward improvement. *arXiv preprint*  
803 *arXiv:2307.07055*, 2023.

810 Jian Zhang and Bernard Ghanem. Ista-net: Interpretable optimization-inspired deep network for  
811 image compressive sensing. In *Proceedings of the IEEE conference on computer vision and*  
812 *pattern recognition*, pp. 1828–1837, 2018.

813

814 Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Da-  
815 long Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural net-  
816 works. In *Proceedings of the IEEE international conference on computer vision*, pp. 1529–1537,  
817 2015.

818 Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++:  
819 A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image*  
820 *Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop,*  
821 *DLMA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI*  
822 *2018, Granada, Spain, September 20, 2018, Proceedings 4*, pp. 3–11. Springer, 2018.

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863



864	CONTENTS	
865		
866		
867	<b>1 Introduction</b>	<b>1</b>
868		
869	<b>2 The generative hierarchical model</b>	<b>2</b>
870		
871		
872	<b>3 The warm-up problem: Classification in GHMs</b>	<b>4</b>
873		
874		
875	<b>4 Denoising and diffusion in GHMs</b>	<b>6</b>
876		
877		
878	<b>5 Further related work</b>	<b>9</b>
879		
880		
881	<b>6 Conclusions and Discussions</b>	<b>10</b>
882		
883	<b>A Proof strategy</b>	<b>18</b>
884		
885	A.1 Proof strategy: ConvNets approximate the belief propagation algorithm . . . . .	18
886	A.2 Proof strategy: U-Nets approximate the belief propagation algorithm . . . . .	19
887		
888		
889	<b>B Technical preliminaries</b>	<b>20</b>
890		
891		
892	<b>C Proof of Proposition 3</b>	<b>21</b>
893		
894	<b>D Proof of Theorem 4</b>	<b>21</b>
895		
896	D.1 Auxillary lemmas . . . . .	23
897		
898		
899	<b>E Proof of Theorem 1</b>	<b>28</b>
900		
901	E.1 Error decomposition . . . . .	28
902	E.2 Results on generalization . . . . .	28
903	E.3 Auxillary lemmas . . . . .	29
904		
905		
906	<b>F Proof of Proposition 5</b>	<b>30</b>
907		
908		
909	<b>G Proof of Theorem 6</b>	<b>30</b>
910		
911	G.1 Auxillary lemmas . . . . .	31
912		
913	<b>H Proof of Theorem 2</b>	<b>33</b>
914		
915	H.1 Error decomposition . . . . .	33
916	H.2 Results on generalization . . . . .	33
917	H.3 Auxillary lemmas . . . . .	34

## 918 A PROOF STRATEGY

### 919 A.1 PROOF STRATEGY: CONVNETS APPROXIMATE THE BELIEF PROPAGATION ALGORITHM

920 Lemma 17 introduces a decomposition of  $D_2^2(\mu_{\text{NN}}^{\widehat{\mathbf{W}}}, \mu_*)$  into two components, approximation error  
921 and generalization error:

$$922 \quad D_2^2(\mu_{\text{NN}}^{\widehat{\mathbf{W}}}, \mu_*) \leq \underbrace{\inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mu_{\text{NN}}^{\mathbf{W}}, \mu_*)}_{\text{approximation error}} + 2 \cdot \underbrace{\sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbf{R}}(\mu_{\text{NN}}^{\mathbf{W}}) - \mathbf{R}(\mu_{\text{NN}}^{\mathbf{W}}) \right|}_{\text{generalization error}}. \quad 923$$

924 The bound of generalization error follows a standard approach: employing a chaining argument, the  
925 error is controlled by  $\tilde{O}(\sqrt{d_p/n})$ , where  $d_p$  represents the number of parameters in the ConvNet  
926 class.

927 In the following, we describe our strategy to control the approximation error: we first present the  
928 belief propagation and message passing algorithm for computing the Bayes classifier  $\mu_*$ , and then  
929 demonstrate that ConvNets are capable of effectively approximating this message passing algorithm.  
930

931 **The belief propagation and message passing algorithm.** The belief propagation algorithm  
932 operates on input  $\mathbf{x} \in [S]^d$  and iteratively calculates the beliefs  $\{\nu_v^{(\ell)} \in \Delta([S])\}_{\ell \in \{0, \dots, L\}, v \in \mathcal{V}^{(\ell)}}$  as  
933 follows:  
934

$$935 \quad \begin{aligned} 936 \quad \nu_v^{(L)}(x_v^{(L)}) &= 1\{x_v^{(L)} = x_v\}, \\ 937 \quad \nu_v^{(\ell)}(x_v^{(\ell)}) &\propto \sum_{x_{\mathcal{C}(v)}^{(\ell+1)}} \prod_{v' \in \mathcal{C}(v)} \left( \psi_{v'}^{(\ell+1)}(x_v^{(\ell)}, x_{v'}^{(\ell+1)}) \nu_{v'}^{(\ell+1)}(x_{v'}^{(\ell+1)}) \right), \quad \ell = L-1, \dots, 0, \\ 938 \quad \mu_{\text{BP}}(y|\mathbf{x}) &= \nu_r^{(0)}(y). \end{aligned} \quad 939$$

(BP-CLS) 940

941 Classical results in graphical models verify that the belief propagation algorithm accurately com-  
942 putes the Bayes classifier in this tree graph.

943 **Lemma 1** (BP calculates the Bayes classifier exactly (Pearl, 1982; Wainwright et al., 2008; Mezard  
944 & Montanari, 2009)). *When applying the belief propagation algorithm (BP-CLS) starting with  $\mathbf{x} \in$*   
945  *$[S]^d$ , it holds that  $\mu_*(\cdot|\mathbf{x}) = \mu_{\text{BP}}(\cdot|\mathbf{x})$ .*

946 The belief propagation algorithm can be streamlined into a message passing algorithm, starting with  
947 the initialization  $h_v^{(L)} = x_v$  for each node  $v$  in the highest layer  $\mathcal{V}^{(L)}$ . The operations are defined as  
948 follows:

$$949 \quad \begin{aligned} 950 \quad q_v^{(\ell)} &= f_{v(v)}^{(\ell)}(h_v^{(\ell)}) \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ 951 \quad h_v^{(\ell-1)} &= \text{normalize} \left( \sum_{v' \in \mathcal{C}(v)} q_{v'}^{(\ell)} \right) \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell-1)}, \\ 952 \quad \mu_{\text{MP}}(y|\mathbf{x}) &= \text{softmax}(h_r^{(0)}). \end{aligned} \quad 953$$

(MP-CLS) 954

955 The functions  $f_v^{(L)} : [S] \rightarrow \mathbb{R}^S, \{f_v^{(\ell)} : \mathbb{R}^S \rightarrow \mathbb{R}^S\}_{\ell \in [L-1]}$  are defined as:

$$956 \quad \begin{aligned} 957 \quad f_v^{(L)}(x)_s &= \log \psi_v^{(L)}(s, x), & x \in [S], s \in [S], \\ 958 \quad f_v^{(\ell)}(h)_s &= \log \sum_{a \in [S]} \psi_v^{(\ell)}(s, a) e^{h_a}, & h \in \mathbb{R}^S, s \in [S], \ell \in [L-1]. \end{aligned} \quad 959$$

(22) 960

961 We note that the normalization operator in (MP-CLS) is non-essential and could be dropped; how-  
962 ever, we include it to ensure the formulation closely mirrors (ConvNet), offering a technical benefit.

963 The subsequent proposition affirms that message passing is essentially equivalent to belief propaga-  
964 tion:

965 **Proposition 3** (BP reduces to MP). *Consider the belief propagation algorithm and the message*  
966 *passing algorithm, both starting from  $\mathbf{x} \in [S]^d$ , as in Eq. (BP-CLS) and (MP-CLS). Then we have*  
967  *$\nu_v^{(\ell)}(\cdot) = \text{softmax}(h_v^{(\ell)})$  for all  $0 \leq \ell \leq L-1$  and  $v \in \mathcal{V}^{(\ell)}$ . In particular, we have  $\mu_{\text{BP}}(\cdot|\mathbf{x}) =$*   
968  *$\mu_{\text{MP}}(\cdot|\mathbf{x}) = \mu_*(\cdot|\mathbf{x})$ .*

969 The proof of Proposition 3 is presented in Section C.

**Approximating message passing with ConvNets.** By comparing the message passing algorithm (MP-CLS) alongside the ConvNet (ConvNet), the primary distinction lies in the nonlinear functions used:  $f_i^{(\ell)}$  versus  $\text{NN}^{(\ell)}$ . Given the expression  $f_i^{(\ell)}(h)_s = \log \sum_{a \in [S]} \psi_i^{(\ell)}(s, a) e^{h_a}$ , it becomes evident that approximating the logarithmic and exponential functions using one-hidden-layer networks enables  $f_i^{(\ell)}(h)$  to be effectively approximated by a two-hidden-layer neural network. This leads to the following theorem:

**Theorem 4** (ConvNets approximation of Bayes classifier). *Let Assumption 1 and 2 hold. For any  $\delta > 0$ , take*

$$D = 4\lceil S^2 K^2 d \cdot 3^L / \delta \rceil, \quad B = \text{Poly}(d, S, K, 3^L, 1/\delta).$$

*Then there exists  $\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}$  as in Eq. (9), such that defining  $\mu_{\text{NN}}^{\mathbf{W}}$  as in Eq. (ConvNet), we have*

$$\max_{y \in [S], \mathbf{x} \in [S]^d} \left| \log \mu_{\star}(y|\mathbf{x}) - \log \mu_{\text{NN}}^{\mathbf{W}}(y|\mathbf{x}) \right| \leq \delta.$$

The proof of Theorem 4 is detailed in Section D.

## A.2 PROOF STRATEGY: U-NETS APPROXIMATE THE BELIEF PROPAGATION ALGORITHM

The proof strategy for the denoising task closely aligns with that of the classification task as detailed in Section A.1.

Lemma 22 introduces a decomposition of  $D_2^2(\mathbf{m}_{\text{NN}}^{\widehat{\mathbf{W}}}, \mathbf{m}_{\star})$  into two components, approximation error and generalization error:

$$D_2^2(\mathbf{m}_{\text{NN}}^{\widehat{\mathbf{W}}}, \mathbf{m}_{\star}) \leq \underbrace{\inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mathbf{m}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}_{\star})}_{\text{approximation error}} + 2 \cdot \underbrace{\sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) - \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) \right|}_{\text{generalization error}}.$$

The bound of generalization error follows a standard approach: employing a chaining argument, the error is controlled by  $\widehat{O}(\sqrt{d_p/n})$ , where  $d_p$  represents the number of parameters in the U-Net class.

In the following, we describe our strategy to control the approximation error: we first present the belief propagation and message passing algorithm for computing the Bayes denoiser  $\mathbf{m}_{\star}$ , and then demonstrate that U-Nets are capable of effectively approximating this message passing algorithm.

**The belief propagation and message passing algorithm.** The belief propagation algorithm operates on input  $\mathbf{z} \in \mathbb{R}^d$  and iteratively calculates the beliefs  $\{\nu_{\downarrow, v}^{(\ell)}, \nu_{\uparrow, v}^{(\ell)} \in \Delta([S])\}_{\ell \in \{0, \dots, L\}, v \in \mathcal{V}^{(\ell)}}$  as follows:

$$\nu_{\downarrow, v}^{(L)}(x_v^{(L)}) = \psi^{(L+1)}(x_v^{(L)}, z_v),$$

$$\nu_{\downarrow, v}^{(\ell)}(x_v^{(\ell)}) \propto \sum_{x_{\mathcal{C}(v)}^{(\ell+1)}} \prod_{v' \in \mathcal{C}(v)} \left( \psi_{i(v')}^{(\ell+1)}(x_v^{(\ell)}, x_{v'}^{(\ell+1)}) \nu_{\downarrow, v'}^{(\ell+1)}(x_{v'}^{(\ell+1)}) \right), \quad \ell = L-1, \dots, 0,$$

$$\nu_{\uparrow, r}^{(0)}(x_r^{(0)}) \propto 1,$$

$$\nu_{\uparrow, v}^{(\ell)}(x_v^{(\ell)}) \propto \sum_{x_{\text{pa}(v)}^{(\ell-1)}, x_{\mathcal{N}(v)}^{(\ell)}} \psi^{(\ell)}(x_{\text{pa}(v)}^{(\ell-1)}, x_{\mathcal{C}(\text{pa}(v))}^{(\ell)}) \nu_{\uparrow, \text{pa}(v)}^{(\ell-1)}(x_{\text{pa}(v)}^{(\ell-1)}) \prod_{v' \in \mathcal{N}(v)} \nu_{\downarrow, v'}^{(\ell)}(x_{v'}^{(\ell)}), \quad \ell = 1, \dots, L,$$

$$\nu_v^{(L)}(x_v^{(L)}) \propto \nu_{\uparrow, L}^{(L)}(x_v^{(L)}) \psi^{(L+1)}(x_v^{(L)}, z_v),$$

$$\mathbf{m}_{\text{BP}}(\mathbf{z})_v = \sum_{x_v^{(L)}} x_v^{(L)} \nu_v^{(L)}(x_v^{(L)}), \quad (\text{BP-DNS})$$

where  $\psi^{(L+1)}(x_v^{(L)}, z_v) := \exp\{-(x_v^{(L)} - z_v)^2/2\}$ . Classical results in graphical models verify that the belief propagation algorithm accurately computes the Bayes denoiser in this tree graph.

**Lemma 2** (BP calculates the Bayes denoiser exactly (Pearl, 1982; Wainwright et al., 2008; Mezard & Montanari, 2009)). *When applying the belief propagation algorithm (BP-DNS) starting with  $\mathbf{z} \in \mathbb{R}^d$ , it holds that  $\mu_{\star}(x_v|\mathbf{z}) = \nu_v^{(L)}(x_v)$  for  $v \in \mathcal{V}^{(L)}$ , so that  $\mathbf{m}_{\text{BP}}(\mathbf{z}) = \mathbf{m}(\mathbf{z})$ .*

We remark that the downward-upward structure of belief propagation in generative hierarchical models has been pointed out in the literature (Sclocchi et al., 2024).

The belief propagation algorithm can be streamlined into a message passing algorithm, starting with the initialization  $h_{\downarrow,v}^{(L)} = -(x - z_v)^2/2)_{x \in [S]} \in \mathbb{R}^S$  for each node  $v$  in the highest layer  $\mathcal{V}^{(L)}$ . The operations are defined as follows:

$$\begin{aligned} q_{\downarrow,v}^{(\ell)} &= f_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(h_{\downarrow,v}^{(\ell)})) \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ h_{\downarrow,v}^{(\ell-1)} &= \sum_{v' \in \mathcal{C}(v)} q_{\downarrow,v'}^{(\ell)} \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell-1)}, \\ u_{\uparrow,v}^{(\ell)} &= b_{\uparrow,\text{pa}(v)}^{(\ell-1)} \in \mathbb{R}^S, \quad (\text{with } b_{\uparrow,r}^{(0)} = h_{\downarrow,r}^{(0)}) & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \quad (\text{MP-DNS}) \\ b_{\uparrow,v}^{(\ell)} &= f_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(u_{\uparrow,v}^{(\ell)} - q_{\downarrow,v}^{(\ell)})) + h_{\downarrow,v}^{(\ell)} \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ \mathbf{m}_{\text{MP}}(\mathbf{z})_v &= \sum_{s \in [S]} s \cdot \text{softmax}(b_{\uparrow,v}^{(L)})_s, & v \in \mathcal{V}^{(L)}. \end{aligned}$$

The functions  $\{f_{\downarrow,\iota}^{(\ell)}, f_{\uparrow,\iota}^{(\ell)} : \mathbb{R}^S \rightarrow \mathbb{R}^S\}_{\ell \in [L]}$  are defined as

$$\begin{aligned} f_{\downarrow,\iota}^{(\ell)}(h)_s &= \log \sum_{a \in [S]} \psi_{\iota}^{(\ell)}(s, a) e^{h_a}, & h \in \mathbb{R}^S, s \in [S], \ell \in [L], \\ f_{\uparrow,\iota}^{(\ell)}(h)_s &= \log \sum_{a \in [S]} \psi_{\iota}^{(\ell)}(a, s) e^{h_a}, & h \in \mathbb{R}^S, s \in [S], \ell \in [L]. \end{aligned} \quad (23)$$

We note that the normalization operator in (MP-DNS) is non-essential and could be dropped; however, we include it to ensure the formulation closely mirrors (UNet), offering a technical benefit.

The subsequent proposition affirms that message passing is essentially equivalent to belief propagation:

**Proposition 5** (BP reduces to MP). *Consider the belief propagation algorithm and the message passing algorithm, both with input  $\mathbf{z} \in \mathbb{R}^d$ , as in Eq. (BP-DNS) and (MP-DNS). Then we have  $\nu_{\downarrow,v}^{(\ell)}(\cdot) = \text{softmax}(h_{\downarrow,v}^{(\ell)})$  and  $\nu_{\uparrow,v}^{(\ell)}(\cdot) = \text{softmax}(b_{\uparrow,v}^{(\ell)} - h_{\downarrow,v}^{(\ell)})$ , and  $\nu_v^{(L)}(\cdot) = \text{softmax}(b_{\uparrow,v}^{(L)})$ . In particular,  $\mathbf{m}_{\text{MP}}(\mathbf{z}) = \mathbf{m}_{\text{BP}}(\mathbf{z}) = \mathbf{m}(\mathbf{z})$ .*

The proof of Proposition 5 is presented in Section F.

**Approximating message passing with ConvNets.** By comparing the message passing algorithm (MP-DNS) alongside the U-Net (UNet), the primary distinction lies in the nonlinear functions used:  $f_{\downarrow,\iota}^{(\ell)}$  versus  $\text{NN}_{\downarrow,\iota}^{(\ell)}$ . Notably,  $f_{\downarrow,\iota}^{(\ell)}$  entails a log-sum-exponential structure. This structure suggests that approximating the logarithmic and exponential functions with one-hidden-layer neural networks can allow  $f_{\downarrow,\iota}^{(\ell)}(h)$  to be effectively approximated by a two-hidden-layer neural network. This leads to the following theorem:

**Theorem 6** (U-Nets approximation of Bayes denoiser). *Let Assumption 1 and 2 hold. For any  $\delta > 0$ , take*

$$D = 4\lceil S^3 K^2 d \cdot 18^L / \delta \rceil, \quad B = \text{Poly}(d, S, K, 18^L, 1/\delta).$$

*Then there exists  $\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}$  as in Eq. (18), such that defining  $\mathbf{m}_{\text{NN}}^{\mathbf{W}}$  as in Eq. (UNet), we have*

$$\sup_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{m}(\mathbf{z}) - \mathbf{m}_{\text{NN}}(\mathbf{z})\|_{\infty} \leq \delta.$$

The proof of Theorem 6 is detailed in Section G.

## B TECHNICAL PRELIMINARIES

We here present a bound on the supremum of sub-Gaussian processes, whose proof was based on the chaining argument.

**Lemma 3** (Proposition A.4 of (Bai et al., 2024)). *Suppose that  $\{X_w\}_{w \in \Theta}$  is a zero-mean random process given by*

$$X_w \equiv \frac{1}{n} \sum_{i=1}^n f(z_i; w) - \mathbb{E}_z[f(z; w)],$$

*where  $z_1, \dots, z_n$  are i.i.d samples from a distribution  $\mathbb{P}_z$  such that the following assumption holds:*

- 1080 (a) The index set  $\Theta$  is equipped with a distance  $\rho$  and diameter  $B_p$ . Further, assume that for  
 1081 some constant  $A_p$ , for any ball  $\Theta'$  of radius  $r$  in  $\Theta$ , the covering number admits upper bound  
 1082  $\log N(\Delta; \Theta', \rho) \leq d_p \log(2A_p r / \Delta)$  for all  $0 < \Delta \leq 2r$ .  
 1083  
 1084 (b) For any fixed  $w \in \Theta$  and  $z$  sampled from  $\mathbb{P}_z$ , the random variable  $f(z; w) - \mathbb{E}_z[f(z; w)]$  is a  
 1085  $\sigma$ -sub-Gaussian random variable ( $\mathbb{E}[e^{\lambda[f(z; w) - \mathbb{E}_z[f(z; w)]]}] \leq e^{\lambda^2 \sigma^2 / 2}$  for any  $\lambda \in \mathbb{R}$ ).  
 1086  
 1087 (c) For any  $w, w' \in \Theta$  and  $z$  sampled from  $\mathbb{P}_z$ , the random variable  $f(z; w) - f(z; w')$  is a  
 1088  $\sigma' \rho(w, w')$ -sub-Gaussian random variable ( $\mathbb{E}[e^{\lambda[f(z; w) - f(z; w')]}] \leq e^{\lambda^2 (\sigma')^2 \rho^2(w, w') / 2}$  for any  
 1089  $\lambda \in \mathbb{R}$ ).

1090 Then with probability at least  $1 - \eta$ , it holds that

$$1091 \sup_{w \in \Theta} |X_w| \leq C \sigma \sqrt{\frac{d_p \cdot \log(2A_p(1 + B_p \sigma' / \sigma)) + \log(1/\eta)}{n}},$$

1092 where  $C$  is a universal constant.

1093 We next present a simple inequality used in the proof of Theorem 1.

1094 **Lemma 4** (From log ratio bound to square distance bound). *Let  $p$  and  $q$  be two probability measures on  $\Delta([S])$  such that*

$$1095 \max_{y \in [S]} \left| \log p(y) - \log q(y) \right| \leq \delta.$$

1096 Then we have

$$1097 \sum_{s=1}^S (p(s) - q(s))^2 \leq (e^\delta - 1)^2.$$

1098 *Proof of Lemma 4.* The lemma is by the fact that  $|p(y) - q(y)| \leq (\exp\{|\log p(y) - \log q(y)|\} - 1) \cdot p(y)$ .  $\square$

## 1100 C PROOF OF PROPOSITION 3

1101 *Proof of Proposition 3.* By Eq. (MP-CLS) and (22), defining  $\nu_v^{(\ell)}(\cdot) = \text{softmax}(h_v^{(\ell)})$ , then for  $\ell \leq L - 2$ , we get

$$1102 \nu_v^{(\ell)}(x_v^{(\ell)}) \propto \prod_{v' \in \mathcal{C}(v)} \left( \sum_{a \in [S]} \psi_{\iota(v')}^{(\ell+1)}(x_v^{(\ell)}, a) e^{(h_v^{(\ell+1)})_a} \right)$$

$$1103 \propto \sum_{x_{\mathcal{C}(v)}^{(\ell+1)}} \prod_{v' \in \mathcal{C}(v)} \left( \psi_{\iota(v')}^{(\ell+1)}(x_v^{(\ell)}, x_{v'}^{(\ell+1)}) \nu_{v'}^{(\ell+1)}(x_{v'}^{(\ell+1)}) \right).$$

1104 This coincides with the update rule in Eq. (BP-CLS). For  $\ell = L - 1$ , we get

$$1105 \nu_v^{(L-1)}(x_v^{(L-1)}) \propto \prod_{v' \in \mathcal{C}(v)} \left( \sum_{a \in [S]} \psi_{\iota(v')}^{(L)}(x_v^{(L-1)}, a) 1\{a = x_{v'}\} \right) \propto \prod_{v' \in \mathcal{C}(v)} \psi_{\iota(v')}^{(L)}(x_v^{(L-1)}, x_{v'}).$$

1106 This again coincides with the update rule in Eq. (BP-CLS). This finishes the proof of Proposition 3.  $\square$

## 1107 D PROOF OF THEOREM 4

1108 *Proof of Theorem 4.* By Lemma 13, take  $M_1^{(\ell)} = \lceil 2SKd3^L / \delta \rceil + 1$  and  $M_2^{(\ell)} = \lceil 2SK^2d3^L / \delta \rceil + 1$ . Then there exists  $\{(a_j, w_j, b_j)\}_{j \in [M_1^{(\ell)}]}$  and  $\{(\bar{a}_j, \bar{w}_j, \bar{b}_j)\}_{j \in [M_2^{(\ell)}]}$  with

$$1109 \sup_j |a_j| \leq 2SK, \quad \sup_j |w_j| \leq 1, \quad \sup_j |b_j| \leq SK, \quad \sup_j |\bar{a}_j| \leq 2, \quad \sup_j |\bar{w}_j| \leq 1, \quad \sup_j |\bar{b}_j| \leq \log(4 \cdot 3^L S d K^2 / \delta),$$

1110 such that defining

$$1111 \log_{\delta \star}(x) = \sum_{j=1}^{M_1^{(\ell)}} a_j \cdot \text{ReLU}(w_j x + b_j), \quad \exp_{\delta \star}(x) = \sum_{j=1}^{M_2^{(\ell)}} \bar{a}_j \cdot \text{ReLU}(\bar{w}_j x + \bar{b}_j),$$

and defining  $f_\ell^{(\ell)}(h), \bar{f}_\ell^{(\ell)}(h) \in \mathbb{R}^S$  by

$$f_\ell^{(\ell)}(h)_i = \log \sum_{j=1}^S \psi_\ell^{(\ell)}(i, j) \exp(h_j), \quad \bar{f}_\ell^{(\ell)}(h)_i = \log_{\delta^*} \sum_{j=1}^S \psi_\ell^{(\ell)}(i, j) \exp_{\delta^*}(h_j), \quad \forall i \in [S],$$

we have

$$\sup_{\max_i h_i=0} \|f_\ell^{(\ell)}(h) - \bar{f}_\ell^{(\ell)}(h)\|_\infty \leq \delta/(d3^L). \quad (24)$$

In addition, by Lemma 15, take  $M_1^{(L)} = \lceil 3^L dK/\delta \rceil + 1$  and  $M_2^{(L)} = 3$ . Then there exists  $\{(a_j^{(L)}, w_j^{(L)}, b_j^{(L)})\}_{j \in [M_1^{(L)}]}$  and  $\{(\bar{a}_j^{(L)}, \bar{w}_j^{(L)}, \bar{b}_j^{(L)})\}_{j \in [M_2^{(L)}]}$  with

$$\sup_j |a_j^{(L)}| \leq 2K, \quad \sup_j |w_j^{(L)}| \leq 1, \quad \sup_j |b_j^{(L)}| \leq SK, \quad \sup_j |\bar{a}_j^{(L)}| \leq 4, \quad \sup_j |\bar{w}_j^{(L)}| \leq 1, \quad \sup_j |\bar{b}_j^{(L)}| \leq 1,$$

such that defining

$$\log_\delta(x) = \sum_{j=1}^{M_1^{(L)}} a_j^{(L)} \cdot \text{ReLU}(w_j^{(L)} x + b_j^{(L)}), \quad \text{Ind}(x) = \sum_{j=1}^{M_2^{(L)}} \bar{a}_j^{(L)} \cdot \text{ReLU}(\bar{w}_j^{(L)} x + \bar{b}_j^{(L)}),$$

and defining  $f_\ell^{(L)}(h), \bar{f}_\ell^{(L)}(h) \in \mathbb{R}^S$  by

$$f_\ell^{(L)}(x)_i = \log \sum_{j=1}^S \psi_\ell^{(L)}(i, j) 1(x=j), \quad \bar{f}_\ell^{(L)}(x)_i = \log_\delta \sum_{j=1}^S \psi_\ell^{(L)}(i, j) \text{Ind}(x-j), \quad \forall i \in [S],$$

we have

$$\sup_{x \in [S]} \|f_\ell^{(L)}(x) - \bar{f}_\ell^{(L)}(x)\|_\infty \leq \delta/(d3^L). \quad (25)$$

By Eq. (24) and (25) and Lemma 5, taking  $h_r^{(0)} \in \mathbb{R}^S$  to be as defined in Eq. (MP-CLS) and  $\bar{h}_r^{(0)} \in \mathbb{R}^S$  to be as defined in Eq. (A-MP-CLS) with  $\{\bar{f}_\ell^{(\ell)}\}_{\ell \in [L], \ell \in [m^{(\ell)}]}$  as defined above, we have

$$\|h_r^{(0)} - \bar{h}_r^{(0)}\|_\infty \leq [\delta/(d3^L)] \times \prod_{1 \leq \ell \leq L} (2m^{(\ell)} + 1) \leq \delta.$$

As a consequence, we just need to show that the approximate version of message passing algorithm as in Eq. (A-MP-CLS) could be cast as a neural network.

Indeed, by Lemma 14, there exist two-hidden-layer neural networks (for  $\ell \in [L-1]$  and  $\iota \in [m^{(\ell)}]$ )

$$\text{NN}_{W_{1,\iota}^{(\ell)}, W_{2,\iota}^{(\ell)}, W_{3,\iota}^{(\ell)}}(h) = W_{1,\iota}^{(\ell)} \cdot \text{ReLU}(W_{2,\iota}^{(\ell)} \cdot \text{ReLU}(W_{3,\iota}^{(\ell)} \cdot [h; 1])),$$

with  $W_{1,\iota}^{(\ell)} \in \mathbb{R}^{S \times SM_1^{(\ell)}}$ ,  $W_{2,\iota}^{(\ell)} \in \mathbb{R}^{SM_1^{(\ell)} \times (SM_2^{(\ell)} + 1)}$ ,  $W_{3,\iota}^{(\ell)} \in \mathbb{R}^{(SM_2^{(\ell)} + 1) \times (S+1)}$ , and

$$\|W_{1,\iota}^{(\ell)}\|_{\max} \leq 2SK, \quad \|W_{2,\iota}^{(\ell)}\|_{\max} \leq \text{Poly}(SK M_1^{(\ell)} M_2^{(\ell)}), \quad \|W_{3,\iota}^{(\ell)}\|_{\max} \leq \log(4SK^2/\delta).$$

such that

$$\text{NN}_{W_{1,\iota}^{(\ell)}, W_{2,\iota}^{(\ell)}, W_{3,\iota}^{(\ell)}}(h) = \bar{f}_\ell^{(\ell)}(h), \quad \forall h \in \mathbb{R}^S \text{ such that } \max_j h_j = 0.$$

Furthermore, by Lemma 16, there exist two-hidden-layer neural networks (for  $\iota \in [m^{(L)}]$ )

$$\text{NN}_{W_{1,\iota}^{(L)}, W_{2,\iota}^{(L)}, W_{3,\iota}^{(L)}}(x) = W_{1,\iota}^{(L)} \cdot \text{ReLU}(W_{2,\iota}^{(L)} \cdot \text{ReLU}(W_{3,\iota}^{(L)} \cdot [x; 1])),$$

with  $W_{1,\iota}^{(L)} \in \mathbb{R}^{S \times SM_1^{(L)}}$ ,  $W_{2,\iota}^{(L)} \in \mathbb{R}^{SM_1^{(L)} \times (SM_2^{(L)} + 1)}$ ,  $W_{3,\iota}^{(L)} \in \mathbb{R}^{(SM_2^{(L)} + 1) \times 2}$ ,

$$\|W_{1,\iota}^{(L)}\|_{\max} \leq 2K, \quad \|W_{2,\iota}^{(L)}\|_{\max} \leq \text{Poly}(SK M_1^{(L)} M_2^{(L)}), \quad \|W_{3,\iota}^{(L)}\|_{\max} \leq 1.$$

such that

$$\text{NN}_{W_{1,\iota}^{(L)}, W_{2,\iota}^{(L)}, W_{3,\iota}^{(L)}}(x) = \bar{f}_\ell^{(L)}(x), \quad \forall x \in [S].$$

This proves that the approximate version of message passing as in Eq. (A-MP-CLS) can be cast into the convolutional neural network as in Eq. (ConvNet) with proper choice of dimension

$$D \geq \max_{\ell \in [L]} \{SM_1^{(\ell)}, SM_2^{(\ell)} + 1\} = S \times \left( \lceil 2SK^2 d3^L/\delta \rceil + 1 \right) + 1,$$

and norm of the weights. This finishes the proof of Theorem 4.  $\square$

## D.1 AUXILLARY LEMMAS

**Lemma 5** (Error propagation of the approximate version of message passing in classification). *Assume we have functions  $f_i^{(\ell)}$  and  $\bar{f}_i^{(\ell)}$  such that*

$$\begin{aligned} \|f_i^{(L)}(x) - \bar{f}_i^{(L)}(x)\|_\infty &\leq \delta, \quad \forall x \in [S], \\ \|f_i^{(\ell)}(h) - \bar{f}_i^{(\ell)}(h)\|_\infty &\leq \delta, \quad \forall h \in \mathbb{R}^S \text{ such that } \max_{j \in [S]} h_j = 0, \quad \forall \ell \leq L-1. \end{aligned} \quad (26)$$

*Furthermore, consider the following approximate version of message passing algorithm with initialization  $\bar{h}_v^{(L)} = x_v$  for  $v \in \mathcal{V}^{(L)}$ :*

$$\begin{aligned} \bar{q}_v^{(\ell)} &= \bar{f}_{i(v)}^{(\ell)}(\bar{h}_v^{(\ell)}) \in \mathbb{R}^S, & \ell \in [L-1], \quad v \in \mathcal{V}^{(\ell)}, \\ \bar{h}_v^{(\ell-1)} &= \text{normalize}\left(\sum_{v' \in \mathcal{C}(v)} \bar{q}_{v'}^{(\ell)}\right) \in \mathbb{R}^S, & \ell \in [L-1], \quad v \in \mathcal{V}^{(\ell-1)}. \end{aligned} \quad (\text{A-MP-CLS})$$

*Taking  $h_r^{(0)} \in \mathbb{R}^S$  to be as defined in Eq. (MP-CLS) and  $\bar{h}_r^{(0)} \in \mathbb{R}^S$  to be as defined in Eq. (A-MP-CLS), we have*

$$\|h_r^{(0)} - \bar{h}_r^{(0)}\|_\infty \leq \delta \times \prod_{1 \leq \ell \leq L} (2m^{(\ell)} + 1).$$

*Proof of Lemma 5.* We prove this lemma by induction, aiming to show that for any  $\ell \in [L-1]$  we have

$$\|\bar{h}_v^{(\ell)} - h_v^{(\ell)}\|_\infty \leq 2m^{(\ell+1)} \prod_{k=\ell+2}^L (2m^{(k)} + 1)\delta, \quad \forall v \in \mathcal{V}^{(\ell)}. \quad (27)$$

To prove the formula for  $\ell = L-1$ , since  $h_v^{(L)} = \bar{h}_v^{(L)}$ , by Eq. (26), we get

$$\|\bar{q}_v^{(L)} - q_v^{(L)}\|_\infty \leq \delta, \quad \forall v \in \mathcal{V}^{(L)}.$$

By Lemma 7, we get

$$\|h_v^{(L-1)} - \bar{h}_v^{(L-1)}\|_\infty \leq 2 \left\| \sum_{v' \in \mathcal{C}(v)} (\bar{q}_{v'}^{(L)} - q_{v'}^{(L)}) \right\|_\infty \leq 2m^{(L)}\delta, \quad \forall v \in \mathcal{V}^{(L-1)}.$$

This proves the formula (27) for  $\ell = L-1$ .

Assuming that (27) holds at the layer  $\ell$ , by the update formula, we have

$$\begin{aligned} \|\bar{q}_v^{(\ell)} - q_v^{(\ell)}\|_\infty &= \|\bar{f}_{i(v)}^{(\ell)}(\bar{h}_v^{(\ell)}) - f_{i(v)}^{(\ell)}(h_v^{(\ell)})\|_\infty \\ &\leq \|\bar{f}_{i(v)}^{(\ell)}(\bar{h}_v^{(\ell)}) - f_{i(v)}^{(\ell)}(\bar{h}_v^{(\ell)})\|_\infty + \|f_{i(v)}^{(\ell)}(\bar{h}_v^{(\ell)}) - f_{i(v)}^{(\ell)}(h_v^{(\ell)})\|_\infty \\ &\leq \delta + 2m^{(\ell+1)} \prod_{k=\ell+2}^L (2m^{(k)} + 1)\delta \leq \prod_{k=\ell+1}^L (2m^{(k)} + 1)\delta, \end{aligned}$$

where the middle inequality is by the assumption of  $f_i^{(\ell)}$  and by Lemma 6. By Lemma 7, we get

$$\|\bar{h}_v^{(\ell-1)} - h_v^{(\ell-1)}\|_\infty \leq 2m^{(\ell)} \prod_{k=\ell+1}^L (2m^{(k)} + 1)\delta, \quad \forall v \in \mathcal{V}^{(\ell)}.$$

This proves Lemma 5 by the induction argument.  $\square$

**Lemma 6** (Non-expansiveness of log-sum-exponential). *For  $h \in \mathbb{R}^S$  and  $\Psi \in \mathbb{R}^{S \times S}$ , define  $f(h) \in \mathbb{R}^S$  by*

$$f(h)_i = \log \sum_{j=1}^S \Psi_{ij} \exp(h_j), \quad \forall i \in [S].$$

*Then for  $h_1, h_2 \in \mathbb{R}^S$ , we have*

$$\|f(h_1) - f(h_2)\|_\infty \leq \|h_1 - h_2\|_\infty.$$

*Proof of Lemma 6.* Fix  $i \in [S]$ . We have

$$\nabla_h f(h)_i = \left( \frac{\Psi_{ij} \exp(h_j)}{\sum_{k \in [S]} \Psi_{ik} \exp(h_k)} \right)_{j \in [S]},$$

so that  $\|\nabla_h f(h)_i\|_1 = 1$ . By intermediate value theorem, we have

$$|f(h_1)_i - f(h_2)_i| = |\nabla_h f(\xi)_i^\top (h_1 - h_2)| \leq \|\nabla_h f(\xi)_i\|_1 \|h_1 - h_2\|_\infty = \|h_1 - h_2\|_\infty.$$

This proves Lemma 6.  $\square$

**Lemma 7** (Lipschitzness of the normalization operator). For  $h \in \mathbb{R}^S$ , define  $\text{normalize}(h) \in \mathbb{R}^S$  by

$$\text{normalize}(h)_i = h_i - \max_j h_j, \quad \forall i \in [S].$$

Then for  $h_1, h_2 \in \mathbb{R}^S$ , we have

$$\|\text{normalize}(h_1) - \text{normalize}(h_2)\|_\infty \leq 2\|h_1 - h_2\|_\infty.$$

*Proof of Lemma 7.* Note we have the following inequality

$$|\max_j h_{1,j} - \max_j h_{2,j}| \leq \|h_1 - h_2\|_\infty,$$

so that

$$|\text{normalize}(h_1)_i - \text{normalize}(h_2)_i| \leq \|h_1 - h_2\|_\infty + |\max_j h_{1,j} - \max_j h_{2,j}| \leq 2\|h_1 - h_2\|_\infty.$$

This completes the proof of Lemma 7.  $\square$

**Lemma 8** (ReLU approximation of the exponential function). For any  $\delta > 0$ , take  $M = \lceil 1/\delta \rceil + 1 \in \mathbb{N}$ . Then there exists  $\{(a_j, w_j, b_j)\}_{j \in [M]}$  with

$$\sup_j |a_j| \leq 2, \quad \sup_j |w_j| \leq 1, \quad \sup_j |b_j| \leq \log M, \quad (28)$$

such that defining  $\exp_\delta : \mathbb{R} \rightarrow \mathbb{R}$  by

$$\exp_\delta(x) = \sum_{j=1}^M a_j \cdot \text{ReLU}(w_j x + b_j),$$

we have  $\exp_\delta$  is non-decreasing on  $(-\infty, 0]$ , and

$$\sup_{x \in (-\infty, 0]} |\exp(x) - \exp_\delta(x)| \leq \delta, \quad \exp_\delta(0) = 1.$$

*Proof of Lemma 8.* Define  $e_j = j/(M-1)$ ,  $b_j = -\log(e_j)$  for  $j \in [M-1]$ ,  $a_1 = (e_2 - e_1)/(b_2 - b_1)$  and  $a_j = (e_{j+1} - e_j)/(b_{j+1} - b_j) - (e_j - e_{j-1})/(b_j - b_{j-1})$  for  $2 \leq j \leq M-2$ . Furthermore, define

$$\exp_\delta(x) = \sum_{j=1}^{M-2} a_j \text{ReLU}(x + b_j) + \text{ReLU}(-x + e_1) - \text{ReLU}(-x).$$

Then we have  $\exp_\delta(-b_j) = e_j$  for  $j \in [M-1]$ , and  $\exp_\delta$  is piece-wise linear and non-decreasing on  $(-\infty, 0]$ . Note that we also have  $\exp(-b_j) = e_j$  for  $j \in [M-1]$ , and  $\exp$  is increasing on  $(-\infty, 0]$ . This proves that  $\sup_{x \in (-\infty, 0]} |\exp(x) - \exp_\delta(x)| \leq \delta$ . Furthermore, it is easy to see that  $\exp_\delta(0) = 1$  and  $\exp_\delta$  is non-decreasing on  $(-\infty, 0]$ . Finally, since  $\exp$  is 1-Lipschitz, it is easy to see that  $\sup_{j \in [M]} |a_j| \leq 2$ . It is also easy to see the other parts of Eq. (28) are satisfied, and this proves Lemma 8.  $\square$

**Lemma 9** (ReLU approximation of the logarithm function). For any  $A > 0$ ,  $\delta > 0$ , take  $M = \lceil 2A/\delta \rceil + 1 \in \mathbb{N}$ . Then there exists  $\{(a_j, w_j, b_j)\}_{j \in [M]}$  with

$$\sup_j |a_j| \leq 2A, \quad \sup_j |w_j| \leq 1, \quad \sup_j |b_j| \leq A, \quad (29)$$

such that defining  $\log_\delta : \mathbb{R} \rightarrow \mathbb{R}$  by

$$\log_\delta(x) = \sum_{j=1}^M a_j \cdot \text{ReLU}(w_j x + b_j),$$

we have  $\log_\delta$  is non-decreasing on  $[1/A, A]$ , and

$$\sup_{x \in [1/A, A]} |\log(x) - \log_\delta(x)| \leq \delta.$$

*Proof of Lemma 9.* The proof of Lemma 9 is similar to Lemma 8.  $\square$



1296 **Lemma 10** (ReLU approximation of indicator function). *Define*

$$1297 \text{Ind}(x) = 2\text{ReLU}(x - 1/2) + 2\text{ReLU}(x + 1/2) - 4\text{ReLU}(x),$$

1299 *we have*

$$1300 1(x = j) = \text{Ind}(x - j), \quad \forall j, x \in \mathbb{Z}.$$

1303 *Proof of Lemma 10.* The lemma holds by direct calculation.  $\square$

1305 **Lemma 11** (Log-sum-exponential approximation). *Assume  $\log_{\delta_1} : \mathbb{R} \rightarrow \mathbb{R}$  and  $\exp_{\delta_2} : \mathbb{R} \rightarrow \mathbb{R}$  are such that,*

$$1307 \sup_{x \in [1/K, SK]} |\log(x) - \log_{\delta_1}(x)| \leq \delta_1, \quad \sup_{x \in (-\infty, 0]} |\exp(x) - \exp_{\delta_2}(x)| \leq \delta_2,$$

$$1309 \exp_{\delta_2}(0) = 1, \quad \exp_{\delta_2} \text{ is non-decreasing on } (-\infty, 0].$$

1311 *Assume that  $1/K \leq \min_{ij} \Psi_{ij} \leq \max_{ij} \Psi_{ij} \leq K$ . Define  $f(h), f_{\delta_1, \delta_2}(h) \in \mathbb{R}^S$  by*

$$1312 f(h)_i = \log \sum_{j=1}^S \Psi_{ij} \exp(h_j), \quad f_{\delta_1, \delta_2}(h)_i = \log_{\delta_1} \sum_{j=1}^S \Psi_{ij} \exp_{\delta_2}(h_j), \quad \forall i \in [S].$$

1314 *Then we have*

$$1315 \sup_{\max_i h_i = 0} \|f(h) - f_{\delta_1, \delta_2}(h)\|_{\infty} \leq \delta_1 + SK^2 \delta_2.$$

1318 *Proof of Lemma 11.* We have

$$1320 |f(h)_i - f_{\delta_1, \delta_2}(h)_i| = |\log \langle \Psi_{i\cdot}, \exp(h) \rangle - \log_{\delta_1} \langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle|$$

$$1321 \leq |\log \langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle - \log_{\delta_1} \langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle| + |\log \langle \Psi_{i\cdot}, \exp(h) \rangle - \log \langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle|.$$

1323 For the first term, since  $\exp_{\delta_2}(h) \leq 1$  for all  $h \leq 0$ , we have  $\langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle \leq S \max_{ij} \Psi_{ij} \leq SK$ .  
1324 Furthermore, since  $\max_i h_i = 0$  and  $\exp_{\delta_2}(0) = 1$ , we have  $\langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle \geq \min_{ij} \Psi_{ij} \geq 1/K$ .  
1325 As a consequence, by assumption, we have

$$1326 |\log \langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle - \log_{\delta_1} \langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle| \leq \delta_1.$$

1328 For the second term, since both  $\langle \Psi_{i\cdot}, \exp(h) \rangle$  and  $\langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle$  are within  $[1/K, SK]$  on which  
1329 log function has Lipschitz constant  $K$ , we have

$$1330 |\log \langle \Psi_{i\cdot}, \exp(h) \rangle - \log \langle \Psi_{i\cdot}, \exp_{\delta_2}(h) \rangle| \leq SK \max_{ij} \Psi_{ij} \cdot \|\exp(h) - \exp_{\delta_2}(h)\|_{\infty} \leq SK^2 \delta_2.$$

1333 This finishes the proof of Lemma 11.  $\square$

1334 **Lemma 12** (Log-Psi approximation). *Assume  $\log_{\delta} : \mathbb{R} \rightarrow \mathbb{R}$  is such that,*

$$1336 \sup_{x \in [1/K, K]} |\log(x) - \log_{\delta}(x)| \leq \delta.$$

1338 *Assume that  $1/K \leq \min_{ij} \Psi_{ij} \leq \max_{ij} \Psi_{ij} \leq K$ . For  $x \in [S]$ , define  $f(x), f_{\delta}(x) \in \mathbb{R}^S$  by*

$$1339 f(x)_i = \log \sum_{j=1}^S \Psi_{ij} 1(x = j), \quad f_{\delta}(x)_i = \log_{\delta} \sum_{j=1}^S \Psi_{ij} 1(x = j), \quad \forall i \in [S].$$

1342 *Then we have*

$$1343 \sup_{x \in [S]} \|f(x) - f_{\delta}(x)\|_{\infty} \leq \delta.$$

1346 *Proof of Lemma 12.* For any fixed  $x \in [S]$  and  $i \in [S]$ , we have

$$1347 |f(x)_i - f_{\delta}(x)_i| = |\log \Psi_{ix} - \log_{\delta} \Psi_{ix}| \leq \delta,$$

1349 where the last inequality is by assumption. This proves Lemma 12.  $\square$

**Lemma 13** (ReLU approximation of log-sum-exponential). *Assume that  $1/K \leq \min_{ij} \Psi_{ij} \leq \max_{ij} \Psi_{ij} \leq K$ . For any  $\delta > 0$ , take  $M_1 = \lceil 2SK/\delta \rceil + 1$  and  $M_2 = \lceil 2SK^2/\delta \rceil + 1$ . Then there exists  $\{(a_j, w_j, b_j)\}_{j \in [M_1]}$  and  $\{(\bar{a}_j, \bar{w}_j, \bar{b}_j)\}_{j \in [M_2]}$  with*

$$\sup_j |a_j| \leq 2SK, \quad \sup_j |w_j| \leq 1, \quad \sup_j |b_j| \leq SK, \quad \sup_j |\bar{a}_j| \leq 2, \quad \sup_j |\bar{w}_j| \leq 1, \quad \sup_j |\bar{b}_j| \leq \log(4SK^2/\delta),$$

such that defining

$$\log_{\delta_\star}(x) = \sum_{j=1}^{M_1} a_j \cdot \text{ReLU}(w_j x + b_j), \quad \exp_{\delta_\star}(x) = \sum_{j=1}^{M_2} \bar{a}_j \cdot \text{ReLU}(\bar{w}_j x + \bar{b}_j),$$

and defining  $f(h), f_\delta(h) \in \mathbb{R}^S$  by

$$f(h)_i = \log \sum_{j=1}^S \Psi_{ij} \exp(h_j), \quad f_\delta(h)_i = \log_{\delta_\star} \sum_{j=1}^S \Psi_{ij} \exp_{\delta_\star}(h_j), \quad \forall i \in [S],$$

we have

$$\sup_{\max_i h_i=0} \|f(h) - f_\delta(h)\|_\infty \leq \delta.$$

*Proof of Lemma 13.* Lemma 9 implies that taking  $M_1 = \lceil 2SK/\delta \rceil + 1$ , there exists  $\{(a_j, w_j, b_j)\}_{j \in [M_1]}$  with

$$\sup_j |a_j| \leq 2SK, \quad \sup_j |w_j| \leq 1, \quad \sup_j |b_j| \leq SK,$$

such that defining

$$\log_{\delta_\star}(x) = \sum_{j=1}^{M_1} a_j \cdot \text{ReLU}(w_j x + b_j),$$

we have

$$\sup_{1/(SK) \leq x \leq SK} |\log(x) - \log_{\delta_\star}(x)| \leq \delta/2.$$

Lemma 8 implies that taking  $M_2 = \lceil 2SK^2/\delta \rceil + 1$ , there exists  $\{(\bar{a}_j, \bar{w}_j, \bar{b}_j)\}_{j \in [M_2]}$  with

$$\sup_j |\bar{a}_j| \leq 2, \quad \sup_j |\bar{w}_j| \leq 1, \quad \sup_j |\bar{b}_j| \leq \log(4SK^2/\delta),$$

such that defining

$$\exp_{\delta_\star}(x) = \sum_{j=1}^{M_2} \bar{a}_j \cdot \text{ReLU}(\bar{w}_j x + \bar{b}_j),$$

we have  $\exp_{\delta_\star}$  is non-decreasing on  $(-\infty, 0]$ , and

$$\sup_{x \in (-\infty, 0]} |\exp(x) - \exp_{\delta_\star}(x)| \leq \delta/(2SK^2), \quad \exp_{\delta_\star}(0) = 1.$$

As a consequence, the condition of Lemma 11 is satisfied with  $\delta_1 = \delta/2$  and  $\delta_2 = \delta/(2SK^2)$ , so that we have

$$\sup_{\max_i h_i=0} \|f(h) - f_\delta(h)\|_\infty \leq \delta_1 + SK^2 \delta_2 = \delta.$$

This finishes the proof of Lemma 13.  $\square$

**Lemma 14** (Existence of ReLU network approximating log-sum-exponential). *Let  $f_\delta$  be the function as defined in Lemma 13. Then there exists a two-hidden-layer neural network*

$$\text{NN}_{W_1, W_2, W_3}(h) = W_1 \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_3 \cdot [h; 1])),$$

with  $W_1 \in \mathbb{R}^{S \times SM_1}$ ,  $W_2 \in \mathbb{R}^{SM_1 \times (SM_2+1)}$ ,  $W_3 \in \mathbb{R}^{(SM_2+1) \times (S+1)}$ , and

$$\|W_1\|_{\max} \leq 2SK, \quad \|W_2\|_{\max} \leq \text{Poly}(SKM_1M_2), \quad \|W_3\|_{\max} \leq \log(4SK^2/\delta).$$

such that

$$\text{NN}_{W_1, W_2, W_3}(h) = f_\delta(h), \quad \forall h \in \mathbb{R}^S \text{ such that } \max_j h_j = 0.$$

1404 *Proof of Lemma 14.* Define

$$1405 \bar{W}_2 = \begin{bmatrix} \bar{a}_{1:M_2}^\top & 0 & \cdots & 0 & 0 \\ 1406 0 & \bar{a}_{1:M_2}^\top & \cdots & 0 & 0 \\ 1407 \cdots & \cdots & \cdots & \cdots & \cdots \\ 1408 0 & 0 & \cdots & \bar{a}_{1:M_2}^\top & 0 \\ 1409 0 & 0 & \cdots & 0 & 1 \end{bmatrix} \in \mathbb{R}^{(S+1) \times (SM_2+1)},$$

$$1410 W_3 = \begin{bmatrix} \bar{w}_{1:M_2} & 0 & \cdots & 0 & \bar{b}_{1:M_2} \\ 1411 0 & \bar{w}_{1:M_2} & \cdots & 0 & \bar{b}_{1:M_2} \\ 1412 \cdots & \cdots & \cdots & \cdots & \cdots \\ 1413 0 & 0 & \cdots & \bar{w}_{1:M_2} & \bar{b}_{1:M_2} \\ 1414 0 & 0 & \cdots & 0 & 1 \end{bmatrix} \in \mathbb{R}^{(SM_2+1) \times (S+1)},$$

1415 then we have

$$1416 [\exp_{\delta^\star}(h); 1] = \bar{W}_2 \cdot \text{ReLU}(W_3 \cdot [h; 1]).$$

1417 Define

$$1418 W_1 = \begin{bmatrix} a_{1:M_1}^\top & 0 & \cdots & 0 \\ 1419 0 & a_{1:M_1}^\top & \cdots & 0 \\ 1420 \cdots & \cdots & \cdots & \cdots \\ 1421 0 & 0 & \cdots & a_{1:M_1}^\top \end{bmatrix} \in \mathbb{R}^{S \times SM_1},$$

$$1422 \tilde{W}_2 = \begin{bmatrix} w_{1:M_1} & 0 & \cdots & 0 & b_{1:M_1} \\ 1423 0 & w_{1:M_1} & \cdots & 0 & b_{1:M_1} \\ 1424 \cdots & \cdots & \cdots & \cdots & \cdots \\ 1425 0 & 0 & \cdots & w_{1:M_1} & \bar{b}_{1:M_1} \end{bmatrix} \in \mathbb{R}^{SM_1 \times (S+1)},$$

1426 then we have

$$1427 \log_{\delta^\star}(h) = W_1 \cdot \text{ReLU}(\tilde{W}_2 \cdot [h; 1]).$$

1428 As a consequence, we have

$$1429 f_\star(h) = \log_{\delta^\star} \Psi \exp_{\delta^\star}(h) = W_1 \cdot \text{ReLU}(\tilde{W}_2 \cdot \text{diag}(\Psi, 1) \cdot \bar{W}_2 \cdot \text{ReLU}(W_3 \cdot [h; 1])) = W_1 \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_3 \cdot [h; 1])),$$

1430 where we define  $W_2 = \tilde{W}_2 \cdot \text{diag}(\Psi, 1) \cdot \bar{W}_2$ . It is also direct to upper bound  $\|W_1\|_{\max}$ ,  $\|W_2\|_{\max}$ ,  
1431 and  $\|W_3\|_{\max}$ . This finishes the proof of Lemma 14.  $\square$

1432 **Lemma 15** (ReLU approximation of log-Psi). *Assume that  $1/K \leq \min_{ij} \Psi_{ij} \leq \max_{ij} \Psi_{ij} \leq K$ .  
1433 For any  $\delta > 0$ , take  $M_1 = \lceil K/\delta \rceil + 1$  and  $M_2 = 3$ . Then there exists  $\{(a_j, w_j, b_j)\}_{j \in [M_1]}$  and  
1434  $\{(\bar{a}_j, \bar{w}_j, \bar{b}_j)\}_{j \in [M_2]}$  with*

$$1435 \sup_j |a_j| \leq 2K, \quad \sup_j |w_j| \leq 1, \quad \sup_j |b_j| \leq SK, \quad \sup_j |\bar{a}_j| \leq 4, \quad \sup_j |\bar{w}_j| \leq 1, \quad \sup_j |\bar{b}_j| \leq 1,$$

1436 such that defining

$$1437 \log_\delta(x) = \sum_{j=1}^{M_1} a_j \cdot \text{ReLU}(w_j x + b_j), \quad \text{Ind}(x) = \sum_{j=1}^{M_2} \bar{a}_j \cdot \text{ReLU}(w_j x + b_j),$$

1438 and defining  $f(x), f_\delta(x) \in \mathbb{R}^S$  by

$$1439 f(x)_i = \log \sum_{j=1}^S \Psi_{ij} 1(x = j), \quad f_\delta(x)_i = \log_\delta \sum_{j=1}^S \Psi_{ij} \text{Ind}(x - j), \quad \forall i \in [S],$$

1440 we have

$$1441 \sup_{x \in [S]} \|f(x) - f_\delta(x)\|_\infty \leq \delta.$$

1442 *Proof of Lemma 15.* The proof of the lemma is similar to the proof of Lemma 13, using a combina-  
1443 tion of Lemma 10, 9, and 12.  $\square$

1444 **Lemma 16** (Existence of ReLU network approximating log-Psi). *Let  $f_\delta$  be the function as defined  
1445 in Lemma 15. Then there exists a two-hidden-layer neural network*

$$1446 \text{NN}_{W_1, W_2, W_3}(x) = W_1 \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_3 \cdot [x; 1])),$$

1447 with  $W_1 \in \mathbb{R}^{S \times SM_1}$ ,  $W_2 \in \mathbb{R}^{SM_1 \times (SM_2+1)}$ ,  $W_3 \in \mathbb{R}^{(SM_2+1) \times (S+1)}$ , and

$$1448 \|W_1\|_{\max} \leq 2K, \quad \|W_2\|_{\max} \leq \text{Poly}(SKM_1M_2), \quad \|W_3\|_{\max} \leq 1.$$

1449 such that

$$1450 \text{NN}_{W_1, W_2, W_3}(x) = f_\delta(x), \quad \forall x \in [S].$$

1451 *Proof of Lemma 16.* The proof of Lemma 16 is similar to the proof of Lemma 14.  $\square$

## 1458 E PROOF OF THEOREM 1

1459  
1460 *Proof of Theorem 1.* By Lemma 17, we have the error decomposition

$$1461 \quad D_2^2(\mu_{\widehat{\mathbf{W}}}, \mu_*) \leq \inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mu_{\mathbf{W}}, \mu_*) + 2 \cdot \sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbf{R}}(\mu_{\mathbf{W}}) - \mathbf{R}(\mu_{\mathbf{W}}) \right|.$$

1462  
1463 To control the first term (the approximation error), by Theorem 4, there exists  $\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}$  as  
1464 in Eq. (9) with norm bound  $B = \text{Poly}(d, S, K, 3^L, D)$ , such that defining  $\mu_{\mathbf{W}}^{\mathbf{W}}$  as in Eq. (ConvNet),  
1465 we have

$$1466 \quad \max_{y \in [S], \mathbf{x} \in [S]^d} \left| \log \mu_*(y|\mathbf{x}) - \log \mu_{\mathbf{W}}^{\mathbf{W}}(y|\mathbf{x}) \right| \leq C \cdot \frac{S^2 K^2 d \cdot 3^L}{D}.$$

1467  
1468 Furthermore, by Lemma 4, when  $D \geq S^2 K^2 d \cdot 3^L$ , we have

$$1469 \quad \inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mu_{\mathbf{W}}, \mu_*) \leq C \left( e^{\frac{S^2 K^2 d \cdot 3^L}{D}} - 1 \right)^2 \leq C \frac{S^4 K^4 d^2 \cdot 3^{2L}}{D^2}.$$

1470  
1471 To control the second term (the generalization error), by Proposition 7, with probability at least  $1 - \eta$ ,  
1472 we have

$$1473 \quad \sup_{\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}} \left| \widehat{\mathbf{R}}(\mu_{\mathbf{W}}) - \mathbf{R}(\mu_{\mathbf{W}}) \right| \leq C \cdot \sqrt{\frac{LD(D + 2S + 1) \|\underline{m}\|_1 \log(d \|\underline{m}\|_1 D S B \cdot 3^L) + \log(1/\eta)}{n}}.$$

1474  
1475 Combining the above two equations proves Theorem 1.  $\square$

### 1481 E.1 ERROR DECOMPOSITION

1482  
1483 **Lemma 17.** *Consider the setting of Theorem 1. We have decomposition*

$$1484 \quad D_2^2(\mu_{\widehat{\mathbf{W}}}, \mu_*) \leq \inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mu_{\mathbf{W}}, \mu_*) + 2 \cdot \sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbf{R}}(\mu_{\mathbf{W}}) - \mathbf{R}(\mu_{\mathbf{W}}) \right|.$$

1485  
1486 *Proof of Lemma 17.* We have that for any conditional distribution  $\mu_1(\cdot|\cdot)$ , there is decomposition

$$1487 \quad D_2^2(\mu_1, \mu_*) = \mathbb{E}_{\mathbf{x} \sim \mu_*} \left[ \sum_{s=1}^S \left( \mu_1(s|\mathbf{x}) - \mu_*(s|\mathbf{x}) \right)^2 \right]$$

$$1488 \quad = \mathbb{E}_{(\mathbf{x}, y) \sim \mu_*} \left[ \sum_{s=1}^S (\mu_1(s|\mathbf{x}) - 1\{y = s\})^2 \right] - \mathbb{E}_{(\mathbf{x}, y) \sim \mu_*} \left[ \sum_{s=1}^S (\mu_*(s|\mathbf{x}) - 1\{y = s\})^2 \right] = \mathbf{R}(\mu_1) - \mathbf{R}(\mu_*).$$

1489  
1490 Define

$$1491 \quad \mathbf{W}_* = \arg \min_{\mathbf{W} \in \mathcal{W}} \mathbf{R}(\mu_{\mathbf{W}}) = \arg \min_{\mathbf{W} \in \mathcal{W}} D_2^2(\mu_{\mathbf{W}}, \mu_*).$$

1492  
1493 Then we have

$$1494 \quad D_2^2(\mu_{\widehat{\mathbf{W}}}, \mu_*) = \mathbf{R}(\mu_{\widehat{\mathbf{W}}}) - \mathbf{R}(\mu_*)$$

$$1495 \quad = \mathbf{R}(\mu_{\widehat{\mathbf{W}}}) - \widehat{\mathbf{R}}(\mu_{\widehat{\mathbf{W}}}) + \widehat{\mathbf{R}}(\mu_{\widehat{\mathbf{W}}}) - \widehat{\mathbf{R}}(\mu_{\mathbf{W}_*}) + \widehat{\mathbf{R}}(\mu_{\mathbf{W}_*}) - \mathbf{R}(\mu_{\mathbf{W}_*}) + \mathbf{R}(\mu_{\mathbf{W}_*}) - \mathbf{R}(\mu_*)$$

$$1500 \quad \leq 2 \cdot \sup_{\mathbf{W} \in \mathcal{W}} \left| \mathbf{R}(\mu_{\mathbf{W}}) - \widehat{\mathbf{R}}(\mu_{\mathbf{W}}) \right| + D_2^2(\mu_{\mathbf{W}_*}, \mu_*)$$

1501  
1502 This proves Lemma 17.  $\square$

### 1506 E.2 RESULTS ON GENERALIZATION

1507  
1508 **Proposition 7** (Generalization error of the classification problem). *Let  $\mathcal{W}_{d, \underline{m}, L, S, D, B}$  be the set*  
1509 *defined as in Eq. (9). Then, with probability at least  $1 - \eta$ , we have*

$$1510 \quad \sup_{\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}} \left| \widehat{\mathbf{R}}(\mu_{\mathbf{W}}) - \mathbf{R}(\mu_{\mathbf{W}}) \right| \leq C \cdot \sqrt{\frac{LD(D + 2S + 1) \|\underline{m}\|_1 \log(d \|\underline{m}\|_1 D S B \cdot 3^L) + \log(1/\eta)}{n}}.$$

*Proof of Proposition 7.* In Lemma 3, we can take  $z = (y, \mathbf{x})$ ,  $w = \mathbf{W}$ ,  $\Theta = \mathcal{W}_{d, \underline{m}, L, S, D, B}$ ,  $\rho(w, w') = \|\mathbf{W} - \mathbf{W}'\|$ , and  $f(z_i; w) = \text{loss}(y, \mu_{\text{NN}}^{\mathbf{W}}(\cdot | \mathbf{x}))$ . Therefore, to show Proposition 7, we just need to apply Lemma 3 by checking (a), (b), (c).

**Check (a).** We note that the index set  $\Theta := \mathcal{W}_{d, \underline{m}, L, S, D, B}$  equipped with  $\rho(w, w') := \|\mathbf{W} - \mathbf{W}'\|$  has diameter  $B_p := 2B$ . Further note that  $\mathcal{W}_{d, \underline{m}, L, S, D, B}$  has a dimension bounded by  $d_p := D(D + 2S + 1)\|\underline{m}\|_1$ . According to Example 5.8 of (Wainwright, 2019), it holds that  $\log N(\Delta; \mathcal{W}_{d, \underline{m}, L, S, D, B}, \|\cdot\|) \leq d_p \cdot \log(1 + 2r/\Delta)$  for any  $0 < \Delta \leq 2r$ . This verifies (a).

**Check (b).** Since  $f(z_i; w) = \text{loss}(y, \mu_{\text{NN}}^{\mathbf{W}}(\cdot | \mathbf{x}))$  is 2-bounded. As a consequence,  $f(z, w) - \mathbb{E}_z[f(z, w)]$  is a sub-Gaussian random variable with the sub-Gaussian parameter to be a universal constant.

**Check (c).** Lemma 20 implies that

$$|f(z; w_1) - f(z; w_2)| \leq \sigma' \cdot \|\mathbf{W}_1 - \mathbf{W}_2\|, \quad \sigma' := 12\|\underline{m}\|_1(3B^3)^L \cdot d \cdot S^{3/2} \cdot (S + D).$$

As a consequence,  $f(z; w_1) - f(z; w_2)$  is  $\sigma' \rho(w_1, w_2) = \sigma' \|\mathbf{W}_1 - \mathbf{W}_2\|$  sub-Gaussian.

Therefore, we apply Lemma 3 to conclude the proof of Proposition 7.  $\square$

### E.3 AUXILIARY LEMMAS

**Lemma 18** (Norm bound in the chain rule in classification settings). *Consider the ConvNet as in Eq. (ConvNet). Assume that  $\|\mathbf{W}\| \leq B$ . Then for any  $\ell, v, \iota$ , and  $\star \in [3]$ , we have*

$$\begin{aligned} \|\mathbf{h}_v^{(L)}\|_2 &\leq S^{3/2}, \\ \|\mathbf{q}_v^{(\ell)}\|_2 &\leq B^3 \cdot (\|\mathbf{h}_v^{(\ell)}\|_2 + 1), \\ \|\mathbf{h}_v^{(\ell-1)}\|_2 &\leq 2m^{(\ell)} \cdot \max_{v' \in \mathcal{C}(v)} \|\mathbf{q}_{v'}^{(\ell)}\|_2, \end{aligned}$$

$$\max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(k)}} \mathbf{q}_{v, i}^{(\ell)}\|_{\text{op}} \leq B^3 \cdot \max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(k)}} \mathbf{h}_{v, i}^{(\ell)}\|_{\text{op}}, \quad \forall k \geq \ell + 1,$$

$$\max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(\ell)}} \mathbf{q}_{v, i}^{(\ell)}\|_{\text{op}} \leq B^2 \cdot \|\mathbf{h}_v^{(\ell)}\|_2,$$

$$\max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(k)}} \mathbf{h}_{v, i}^{(\ell-1)}\|_{\text{op}} \leq 2m^{(\ell)} \cdot \max_{v' \in \mathcal{C}(v)} \max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(k)}} \mathbf{q}_{v', i}^{(\ell)}\|_{\text{op}}, \quad \forall k \geq \ell,$$

$$\max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(\ell)}} \text{softmax}(\mathbf{h}_r^{(0)})_i\|_{\text{op}} \leq \max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(\ell)}} \mathbf{h}_{r, i}^{(0)}\|_{\text{op}}.$$

*Proof of Lemma 18.* The proof of the lemma uses the chain rule, the 1-Lipschitzness of ReLU, the 2-Lipschitzness of normalize, and the 1-Lipschitzness of softmax.  $\square$

**Lemma 19.** *Consider the ConvNet as in Eq. (ConvNet). Assume that  $\|\mathbf{W}\| \leq B$ . Then we have*

$$\max_{\mathbf{x} \in [S]^d} \max_{\star \in [3]} \max_{\ell \in [L]} \max_{\iota \in [m^{(\ell)}]} \max_{i \in [S]} \|\nabla_{W_{\star, \iota}^{(\ell)}} \text{softmax}(\mathbf{h}_r^{(0)})_i\|_{\text{op}} \leq (3B^3)^L \cdot d \cdot S^{3/2}.$$

*Proof of Lemma 19.* This lemma is implied by Lemma 18 and an induction argument.  $\square$

**Lemma 20.** *Consider the ConvNet as in Eq. (ConvNet). Assume that  $\|\mathbf{W}\| \leq B$ . Then we have*

$$\max_{y, \mathbf{x}} \left| \mu_{\text{NN}}^{\mathbf{W}}(y | \mathbf{x}) - \mu_{\text{NN}}^{\overline{\mathbf{W}}}(y | \mathbf{x}) \right| \leq 3\|\underline{m}\|_1(3B^3)^L \cdot d \cdot S^{3/2} \cdot (S + D) \cdot \|\mathbf{W} - \overline{\mathbf{W}}\|.$$

Therefore, we have

$$\left| \text{loss}(y, \mu_{\text{NN}}^{\mathbf{W}}(\cdot | \mathbf{x})) - \text{loss}(y, \mu_{\text{NN}}^{\overline{\mathbf{W}}}(\cdot | \mathbf{x})) \right| \leq 12\|\underline{m}\|_1(3B^3)^L \cdot d \cdot S^{3/2} \cdot (S + D) \cdot \|\mathbf{W} - \overline{\mathbf{W}}\|.$$

*Proof of Lemma 20.* The first inequality is by the fact that

$$\begin{aligned} &\max_{y, \mathbf{x}} \left| \mu_{\text{NN}}^{\mathbf{W}}(y | \mathbf{x}) - \mu_{\text{NN}}^{\overline{\mathbf{W}}}(y | \mathbf{x}) \right| \\ &\leq \sum_{\star \in [3]} \sum_{\ell \in [L]} \sum_{\iota \in [m^{(\ell)}]} \min\{\text{nrow}(W_{\star, \iota}^{(\ell)}), \text{ncol}(W_{\star, \iota}^{(\ell)})\} \|\nabla_{W_{\star, \iota}^{(\ell)}} \text{softmax}(\mathbf{h}_r^{(0)})_i\|_{\text{op}} \|W_{\star, \iota}^{(\ell)} - \overline{W}_{\star, \iota}^{(\ell)}\|_{\text{op}}, \end{aligned}$$

where we have used the inequality that  $\text{trace}(A^T B) \leq \{\text{nrow}(A), \text{ncol}(A)\} \|A\|_{\text{op}} \|B\|_{\text{op}}$ .

To prove the second equation, we have

$$\begin{aligned} & \left| \text{loss}(y, \mu_{\text{NN}}^{\mathbf{W}}(\cdot|\mathbf{x})) - \text{loss}(y, \mu_{\text{NN}}^{\overline{\mathbf{W}}}(\cdot|\mathbf{x})) \right| = \left| \sum_{s=1}^S \left( 1\{y = s\} - \mu_{\text{NN}}^{\mathbf{W}}(s|\mathbf{x}) \right)^2 - \sum_{s=1}^S \left( 1\{y = s\} - \mu_{\text{NN}}^{\overline{\mathbf{W}}}(s|\mathbf{x}) \right)^2 \right| \\ & \leq \sum_{s=1}^S \left| 1\{y = s\} - \mu_{\text{NN}}^{\mathbf{W}}(s|\mathbf{x}) + 1\{y = s\} - \mu_{\text{NN}}^{\overline{\mathbf{W}}}(s|\mathbf{x}) \right| \cdot \left| \mu_{\text{NN}}^{\mathbf{W}}(s|\mathbf{x}) - \mu_{\text{NN}}^{\overline{\mathbf{W}}}(s|\mathbf{x}) \right| \leq 4 \cdot \max_s \left| \mu_{\text{NN}}^{\mathbf{W}}(s|\mathbf{x}) - \mu_{\text{NN}}^{\overline{\mathbf{W}}}(s|\mathbf{x}) \right|. \end{aligned}$$

This completes the proof of Lemma 20.  $\square$

## F PROOF OF PROPOSITION 5

*Proof of Proposition 5.* By Eq. (MP-DNS) and (23), defining  $\nu_{\downarrow, v}^{(\ell)}(\cdot) = \text{softmax}(h_{\downarrow, v}^{(\ell)})$ , then for  $\ell \leq L - 1$ , we get

$$\begin{aligned} \nu_{\downarrow, v}^{(\ell)}(x_v^{(\ell)}) & \propto \prod_{v' \in \mathcal{C}(v)} \left( \sum_{a \in [S]} \psi_{i(v')}^{(\ell+1)}(x_v^{(\ell)}, a) e^{(h_{\downarrow, v}^{(\ell+1)})_a} \right) \\ & \propto \sum_{x_{\mathcal{C}(v)}^{(\ell+1)}} \prod_{v' \in \mathcal{C}(v)} \left( \psi_{i(v')}^{(\ell+1)}(x_v^{(\ell)}, x_{v'}^{(\ell+1)}) \nu_{\downarrow, v'}^{(\ell+1)}(x_{v'}^{(\ell+1)}) \right). \end{aligned}$$

This coincides with the update rule of  $\nu_{\downarrow, v}^{(\ell)}$  as in Eq. (BP-DNS).

Furthermore, defining  $\nu_{\uparrow, v}^{(\ell)}(\cdot) = \text{softmax}(b_{\uparrow, v}^{(\ell)} - h_{\downarrow, v}^{(\ell)})$ , then for  $\ell = 0, 1, \dots, L$ , we get

$$\begin{aligned} \nu_{\uparrow, v}^{(\ell)}(x_v^{(\ell)}) & \propto \sum_{b \in [S]} \psi_{i(v)}^{(\ell+1)}(b, x_v^{(\ell)}) \nu_{\uparrow, \text{pa}(v)}^{(\ell+1)}(b) \prod_{v' \in \mathcal{N}(v)} \left( \sum_{a \in [S]} \psi_{i(v')}^{(\ell+1)}(b, a) e^{(h_{\downarrow, v}^{(\ell+1)})_a} \right) \\ & \propto \sum_{x_{\text{pa}(v)}^{(\ell-1)}, x_{\mathcal{N}(v)}^{(\ell)}} \psi^{(\ell)}(x_{\text{pa}(v)}^{(\ell-1)}, x_{\mathcal{C}(\text{pa}(v))}^{(\ell)}) \nu_{\uparrow, \text{pa}(v)}^{(\ell-1)}(x_{\text{pa}(v)}^{(\ell-1)}) \prod_{v' \in \mathcal{N}(v)} \nu_{\downarrow, v'}^{(\ell)}(x_{v'}^{(\ell)}). \end{aligned}$$

This coincides with the update rule of  $\nu_{\uparrow, v}^{(\ell)}$  as in Eq. (BP-DNS).

Finally, defining  $\nu_v^{(L)}(\cdot) = \text{softmax}(b_{\uparrow, v}^{(L)})$ , we get

$$\begin{aligned} \nu_v^{(L)}(x_v^{(L)}) & \propto \sum_{b \in [S]} \psi_{i(v)}^{(L)}(b, x_v^{(L)}) \nu_{\uparrow, \text{pa}(v)}^{(L)}(b) \prod_{v' \in \mathcal{N}(v)} \left( \sum_{a \in [S]} \psi_{i(v')}^{(L)}(b, a) e^{(h_{\downarrow, v}^{(L)})_a} \right) \times \nu_{\downarrow, v}^{(L)}(x_v^{(L)}) \\ & \propto \nu_{\uparrow, v}^{(L)}(x_v^{(L)}) \psi^{(L+1)}(x_v^{(L)}, z_v). \end{aligned}$$

This coincides with the formula of  $\nu_v^{(L)}$  as in Eq. (BP-DNS).

This finishes the proof of Proposition 5.  $\square$

## G PROOF OF THEOREM 6

*Proof of Theorem 6.* By Lemma 13, take  $M_1 = \lceil 2S^2 K d 18^L / \delta \rceil + 1$  and  $M_2 = \lceil 2S^2 K^2 d 18^L / \delta \rceil + 1$ . Then there exists  $\{(a_j, w_j, b_j)\}_{j \in [M_1]}$  and  $\{(\bar{a}_j, \bar{w}_j, \bar{b}_j)\}_{j \in [M_2]}$  with

$$\sup_j |a_j| \leq 2SK, \quad \sup_j |w_j| \leq 1, \quad \sup_j |b_j| \leq SK, \quad \sup_j |\bar{a}_j| \leq 2, \quad \sup_j |\bar{w}_j| \leq 1, \quad \sup_j |\bar{b}_j| \leq \log(4 \cdot 18^L S^2 d K^2 / \delta),$$

such that defining

$$\log_{\delta^*}(x) = \sum_{j=1}^{M_1} a_j \cdot \text{ReLU}(w_j x + b_j), \quad \exp_{\delta^*}(x) = \sum_{j=1}^{M_2} \bar{a}_j \cdot \text{ReLU}(\bar{w}_j x + \bar{b}_j),$$

and defining  $f_{\diamond,\iota}^{(\ell)}(h), \bar{f}_{\diamond,\iota}^{(\ell)}(h) \in \mathbb{R}^S$  for  $\diamond \in \{\downarrow, \uparrow\}$  by

$$\begin{aligned} f_{\downarrow,\iota}^{(\ell)}(h)_i &= \log \sum_{j=1}^S \psi_i^{(\ell)}(i, j) \exp(h_j), & \bar{f}_{\downarrow,\iota}^{(\ell)}(h)_i &= \log_{\delta^*} \sum_{j=1}^S \psi_i^{(\ell)}(i, j) \exp_{\delta^*}(h_j), & \forall i \in [S], \\ f_{\uparrow,\iota}^{(\ell)}(h)_i &= \log \sum_{j=1}^S \psi_i^{(\ell)}(j, i) \exp(h_j), & \bar{f}_{\uparrow,\iota}^{(\ell)}(h)_i &= \log_{\delta^*} \sum_{j=1}^S \psi_i^{(\ell)}(j, i) \exp_{\delta^*}(h_j), & \forall i \in [S], \end{aligned}$$

we have

$$\sup_{\max_i h_i=0} \|f_{\diamond,\iota}^{(\ell)}(h) - \bar{f}_{\diamond,\iota}^{(\ell)}(h)\|_{\infty} \leq \delta / (d18^L S). \quad (30)$$

By Eq. (30) and Lemma 21, taking  $b_{\uparrow,v}^{(L)} \in \mathbb{R}^S$  to be as defined in Eq. (MP-DNS) and  $\bar{b}_{\uparrow,v}^{(L)} \in \mathbb{R}^S$  to be as defined in Eq. (A-MP-DNS) with  $\{\bar{f}_{\diamond,\iota}^{(\ell)}\}_{\ell \in [L], \iota \in [m^{(\ell)}]}$  as defined above, we have

$$\|b_{\uparrow,v}^{(L)} - \bar{b}_{\uparrow,v}^{(L)}\|_{\infty} \leq [\delta / (d18^L S)] \times 18^L d = \delta / S,$$

which gives

$$\sup_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{m}(\mathbf{z}) - \mathbf{m}_{\text{NN}}(\mathbf{z})\|_{\infty} \leq S \cdot \|b_{\uparrow,v}^{(L)} - \bar{b}_{\uparrow,v}^{(L)}\|_{\infty} \leq \delta.$$

As a consequence, we just need to show that the approximate version of message passing algorithm as in Eq. (A-MP-DNS) could be cast as a neural network.

Indeed, by Lemma 14, there exist two-hidden-layer neural networks (for  $\diamond \in \{\downarrow, \uparrow\}$ ,  $\ell \in [L]$ , and  $\iota \in [m^{(\ell)}]$ )

$$\text{NN}_{W_{1,\diamond,\iota}^{(\ell)}, W_{2,\diamond,\iota}^{(\ell)}, W_{3,\diamond,\iota}^{(\ell)}}(h) = W_{1,\diamond,\iota}^{(\ell)} \cdot \text{ReLU}(W_{2,\diamond,\iota}^{(\ell)} \cdot \text{ReLU}(W_{3,\diamond,\iota}^{(\ell)} \cdot [h; 1])),$$

with  $W_{1,\diamond,\iota}^{(\ell)} \in \mathbb{R}^{S \times SM_1}$ ,  $W_{2,\diamond,\iota}^{(\ell)} \in \mathbb{R}^{SM_1 \times (SM_2+1)}$ ,  $W_{3,\diamond,\iota}^{(\ell)} \in \mathbb{R}^{(SM_2+1) \times (S+1)}$ , and

$$\|W_{1,\diamond,\iota}^{(\ell)}\|_{\max} \leq 2SK, \quad \|W_{2,\diamond,\iota}^{(\ell)}\|_{\max} \leq \text{Poly}(SKM_1M_2), \quad \|W_{3,\diamond,\iota}^{(\ell)}\|_{\max} \leq \log(4S^2d18^L K^2/\delta).$$

such that

$$\text{NN}_{W_{1,\diamond,\iota}^{(\ell)}, W_{2,\diamond,\iota}^{(\ell)}, W_{3,\diamond,\iota}^{(\ell)}}(h) = \bar{f}_{\diamond,\iota}^{(\ell)}(h), \quad \forall h \in \mathbb{R}^S \text{ such that } \max_j h_j = 0.$$

This proves that the approximate version of message passing as in Eq. (A-MP-DNS) coincides with the U-Net as in Eq. (UNet) with proper choice of dimension

$$D \geq \max_{\ell \in [L]} \{SM_1^{(\ell)}, SM_2^{(\ell)} + 1\} = S \times \left( \lceil 2S^2K^2d18^L/\delta \rceil + 1 \right) + 1$$

and norm of the weights. This finishes the proof of Theorem 6.  $\square$

## G.1 AUXILLARY LEMMAS

**Lemma 21** (Error propagation of the approximate version of message passing in denoising). *Assume we have functions  $\{f_{\downarrow,\iota}^{(\ell)}, f_{\uparrow,\iota}^{(\ell)}\}$  and  $\{\bar{f}_{\downarrow,\iota}^{(\ell)}, \bar{f}_{\uparrow,\iota}^{(\ell)}\}$  such that*

$$\|f_{\diamond,\iota}^{(\ell)}(h) - \bar{f}_{\diamond,\iota}^{(\ell)}(h)\|_{\infty} \leq \delta, \quad \forall h \in \mathbb{R}^S \text{ such that } \max_{j \in [S]} h_j = 0, \quad \diamond \in \{\downarrow, \uparrow\}. \quad (31)$$

Furthermore, consider the following approximate version of message passing algorithm with initialization  $h_{\downarrow,v}^{(L)} = -(x - z_v)^2/2)_{x \in [S]} \in \mathbb{R}^S$  for  $v \in \mathcal{V}^{(L)}$ , defined as below

$$\begin{aligned} \bar{q}_{\downarrow,v}^{(\ell)} &= \bar{f}_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{h}_{\downarrow,v}^{(\ell)})) \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ \bar{h}_{\downarrow,v}^{(\ell-1)} &= \sum_{v' \in \mathcal{C}(v)} \bar{q}_{\downarrow,v'}^{(\ell)} \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell-1)}, \\ \bar{u}_{\uparrow,v}^{(\ell)} &= \bar{b}_{\uparrow,\text{pa}(v)}^{(\ell-1)} \in \mathbb{R}^S, \quad (\text{with } \bar{b}_{\uparrow,r}^{(0)} = \bar{h}_{\downarrow,r}^{(0)}) & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ \bar{b}_{\uparrow,v}^{(\ell)} &= \bar{f}_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{u}_{\uparrow,v}^{(\ell)} - \bar{q}_{\downarrow,v}^{(\ell)})) + \bar{h}_{\downarrow,v}^{(\ell)} \in \mathbb{R}^S, & \ell \in [L], \quad v \in \mathcal{V}^{(\ell)}, \\ \bar{\mathbf{m}}_{\text{MP}}(\mathbf{z})_v &= \sum_{s \in [S]} s \cdot \text{softmax}(\bar{b}_{\uparrow,v}^{(L)})_s, & v \in \mathcal{V}^{(L)}. \end{aligned} \quad (\text{A-MP-DNS})$$

Taking  $b_{\uparrow,v}^{(L)} \in \mathbb{R}^S$  to be as defined in Eq. (MP-DNS) and  $\bar{b}_{\uparrow,v}^{(L)} \in \mathbb{R}^S$  to be as defined in Eq. (A-MP-DNS), we have

$$\max_{v \in \mathcal{V}^{(L)}} \|b_{\uparrow,v}^{(L)} - \bar{b}_{\uparrow,v}^{(L)}\|_{\infty} \leq \delta \times 18^L \cdot d.$$

1674 *Proof of Lemma 21.*

1675 **Step 1. Downward induction.** In the first step, we aim to show that for any  $\ell \in [L - 1]$  we have

$$1677 \quad \|\bar{h}_{\downarrow,v}^{(\ell)} - h_{\downarrow,v}^{(\ell)}\|_{\infty} \leq m^{(\ell+1)} \prod_{k=\ell+2}^L (2m^{(k)} + 1)\delta, \quad \forall v \in \mathcal{V}^{(\ell)}. \quad (32)$$

1678 To prove the formula for  $\ell = L - 1$ , since  $\text{normalize}(h_{\downarrow,v}^{(L)}) = \text{normalize}(\bar{h}_{\downarrow,v}^{(L)})$ , by Eq. (31), we get

$$1681 \quad \|\bar{q}_{\downarrow,v}^{(L)} - q_{\downarrow,v}^{(L)}\|_{\infty} \leq \delta, \quad \forall v \in \mathcal{V}^{(L)}.$$

1682 Hence we get

$$1683 \quad \|h_{\downarrow,v}^{(L-1)} - \bar{h}_{\downarrow,v}^{(L-1)}\|_{\infty} = \left\| \sum_{v' \in \mathcal{C}(v)} (\bar{q}_{\downarrow,v'}^{(L)} - q_{\downarrow,v'}^{(L)}) \right\|_{\infty} \leq m^{(L)}\delta, \quad \forall v \in \mathcal{V}^{(L-1)}.$$

1684 This proves the formula (32) for  $\ell = L - 1$ .

1685 Assuming that (32) holds at the layer  $\ell$ , by the update formula, we have

$$1686 \quad \begin{aligned} & \|\bar{q}_{\downarrow,v}^{(\ell)} - q_{\downarrow,v}^{(\ell)}\|_{\infty} = \|\bar{f}_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{h}_{\downarrow,v}^{(\ell)})) - f_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(h_{\downarrow,v}^{(\ell)}))\|_{\infty} \\ & \leq \|\bar{f}_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{h}_{\downarrow,v}^{(\ell)})) - f_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{h}_{\downarrow,v}^{(\ell)}))\|_{\infty} + \|f_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{h}_{\downarrow,v}^{(\ell)})) - f_{\downarrow,\iota(v)}^{(\ell)}(\text{normalize}(h_{\downarrow,v}^{(\ell)}))\|_{\infty} \\ & \leq \delta + 2m^{(\ell+1)} \prod_{k=\ell+2}^L (2m^{(k)} + 1)\delta \leq \prod_{k=\ell+1}^L (2m^{(k)} + 1)\delta, \end{aligned}$$

1687 where the middle inequality is by the assumption of  $f_{\downarrow,\iota}^{(\ell)}$  and by Lemma 6 and Lemma 7. Hence we get

$$1688 \quad \|\bar{h}_{\downarrow,v}^{(\ell-1)} - h_{\downarrow,v}^{(\ell-1)}\|_{\infty} \leq m^{(\ell)} \prod_{k=\ell+1}^L (2m^{(k)} + 1)\delta, \quad \forall v \in \mathcal{V}^{(\ell)}.$$

1689 This proves Eq. (32) by the induction argument.

1690 **Step 2. Upward induction.** The downward induction argument proves that, for  $\Gamma = \prod_{k=1}^L (2m^{(k)} + 1)$ , we have

$$1691 \quad \|\bar{q}_{\downarrow,v}^{(\ell)} - q_{\downarrow,v}^{(\ell)}\|_{\infty}, \|\bar{h}_{\downarrow,v}^{(\ell)} - h_{\downarrow,v}^{(\ell)}\|_{\infty} \leq \Gamma\delta, \quad \forall \ell = 0, 1, \dots, L, \quad \forall v \in \mathcal{V}^{(\ell)}.$$

1692 In this step, we aim to show that for any  $\ell = 0, 1, \dots, L$ , we have

$$1693 \quad \|\bar{b}_{\uparrow,v}^{(\ell)} - b_{\uparrow,v}^{(\ell)}\|_{\infty} \leq 6^{\ell} \cdot \Gamma \cdot \delta, \quad \forall v \in \mathcal{V}^{(\ell)}. \quad (33)$$

1694 To prove this formula for  $\ell = 0$ , note that  $b_{\uparrow,r}^{(0)} = h_{\downarrow,r}^{(0)}$  and  $\bar{b}_{\uparrow,r}^{(0)} = \bar{h}_{\downarrow,r}^{(0)}$ , we have

$$1695 \quad \|\bar{b}_{\uparrow,r}^{(0)} - b_{\uparrow,r}^{(0)}\|_{\infty} = \|\bar{h}_{\downarrow,r}^{(0)} - h_{\downarrow,r}^{(0)}\|_{\infty} \leq \Gamma \cdot \delta.$$

1696 This proves the formula (33) for  $\ell = 0$ .

1697 Assuming that (33) holds at layer  $\ell - 1$ , by the update formula, we have

$$1698 \quad \begin{aligned} & \|\bar{b}_{\uparrow,v}^{(\ell)} - b_{\uparrow,v}^{(\ell)}\|_{\infty} \\ & \leq \|f_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(b_{\uparrow,\text{pa}(v)}^{(\ell-1)} - q_{\downarrow,v}^{(\ell-1)})) - \bar{f}_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{b}_{\uparrow,\text{pa}(v)}^{(\ell-1)} - \bar{q}_{\downarrow,v}^{(\ell-1)}))\|_{\infty} + \|h_{\downarrow,v}^{(\ell)} - \bar{h}_{\downarrow,v}^{(\ell)}\|_{\infty} \\ & \leq \|f_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(b_{\uparrow,\text{pa}(v)}^{(\ell-1)} - q_{\downarrow,v}^{(\ell-1)})) - f_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{b}_{\uparrow,\text{pa}(v)}^{(\ell-1)} - \bar{q}_{\downarrow,v}^{(\ell-1)}))\|_{\infty} \\ & \quad + \|f_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{b}_{\uparrow,\text{pa}(v)}^{(\ell-1)} - \bar{q}_{\downarrow,v}^{(\ell-1)})) - \bar{f}_{\uparrow,\iota(v)}^{(\ell)}(\text{normalize}(\bar{b}_{\uparrow,\text{pa}(v)}^{(\ell-1)} - \bar{q}_{\downarrow,v}^{(\ell-1)}))\|_{\infty} + \|h_{\downarrow,v}^{(\ell)} - \bar{h}_{\downarrow,v}^{(\ell)}\|_{\infty} \\ & \leq \|\text{normalize}(b_{\uparrow,\text{pa}(v)}^{(\ell-1)} - q_{\downarrow,v}^{(\ell-1)}) - \text{normalize}(\bar{b}_{\uparrow,\text{pa}(v)}^{(\ell-1)} - \bar{q}_{\downarrow,v}^{(\ell-1)})\|_{\infty} + \delta + \Gamma \cdot \delta \\ & \leq 4 \cdot 6^{\ell-1} \cdot \Gamma \cdot \delta + \delta + \Gamma \cdot \delta \leq 6^{\ell} \cdot \Gamma \cdot \delta. \end{aligned}$$

1699 This proves Eq. (33) by the induction argument. This proves the Lemma 21 by observing that  $\Gamma \leq 3^L \prod_{\ell=1}^L m^{(\ell)} = 3^L \cdot d$ .  $\square$



## 1728 H PROOF OF THEOREM 2

1729 *Proof of Theorem 2.* By Lemma 22, we have the error decomposition

$$1730 D_2^2(\widehat{\mathbf{m}}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}) \leq \inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mathbf{m}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}) + 2 \cdot \sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) - \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) \right|.$$

1731 To control the first term (the approximation error), by Theorem 6, there exists  $\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}$  as  
1732 in Eq. (18) with norm bound  $B = \text{Poly}(d, S, K, 18^L, D)$ , such that defining  $\mathbf{m}_{\text{NN}}^{\mathbf{W}}$  as in Eq. (UNet),  
1733 we have

$$1734 \sup_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{m}(\mathbf{z}) - \mathbf{m}_{\text{NN}}(\mathbf{z})\|_{\infty} \leq C \cdot \frac{S^3 K^2 d \cdot 18^L}{D}.$$

1735 Therefore, we have

$$1736 \inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mathbf{m}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}) \leq \sup_{\mathbf{z}} \|\mathbf{m}(\mathbf{z}) - \mathbf{m}_{\text{NN}}(\mathbf{z})\|_{\infty}^2 \leq C \cdot \frac{S^6 K^4 d^2 \cdot 18^{2L}}{D^2}.$$

1737 To control the second term (the generalization error), by Proposition 8, with probability at least  $1 - \eta$ ,  
1738 we have

$$1739 \sup_{\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}} \left| \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) - \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) \right| \leq C \cdot S^2 \cdot \sqrt{\frac{LD(D + 2S + 1) \|\underline{m}\|_1 \log(d \|\underline{m}\|_1 DSB \cdot 18^L) + \log(1/\eta)}{n}}.$$

1740 Combining the above two equations proves Theorem 2.  $\square$

### 1741 H.1 ERROR DECOMPOSITION

1742 **Lemma 22.** *Consider the setting of Theorem 2. We have decomposition*

$$1743 D_2^2(\widehat{\mathbf{m}}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}) \leq \inf_{\mathbf{W} \in \mathcal{W}} D_2^2(\mathbf{m}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}) + 2 \cdot \sup_{\mathbf{W} \in \mathcal{W}} \left| \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) - \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) \right|.$$

1744 *Proof of Lemma 22.* We have that for any conditional expectation  $\mathbf{m}_1(\mathbf{z})$ , there is decomposition

$$1745 D_2^2(\mathbf{m}_1, \mathbf{m}) = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mu_{\star}} \left[ d^{-1} \|\mathbf{m}_1(\mathbf{z}) - \mathbf{m}(\mathbf{z})\|_2^2 \right]$$

$$1746 = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mu_{\star}} \left[ d^{-1} \|\mathbf{m}_1(\mathbf{z}) - \mathbf{x}\|_2^2 \right] - \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mu_{\star}} \left[ d^{-1} \|\mathbf{m}_1(\mathbf{z}) - \mathbf{x}\|_2^2 \right] = \mathbf{R}(\mathbf{m}_1) - \mathbf{R}(\mathbf{m}).$$

1747 Define

$$1748 \mathbf{W}_{\star} = \arg \min_{\mathbf{W} \in \mathcal{W}} \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) = \arg \min_{\mathbf{W} \in \mathcal{W}} D_2^2(\mathbf{m}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}).$$

1749 Then we have

$$1750 D_2^2(\widehat{\mathbf{m}}_{\text{NN}}^{\mathbf{W}}, \mathbf{m}) = \mathbf{R}(\widehat{\mathbf{m}}_{\text{NN}}^{\mathbf{W}}) - \mathbf{R}(\mathbf{m})$$

$$1751 = \mathbf{R}(\widehat{\mathbf{m}}_{\text{NN}}^{\mathbf{W}}) - \widehat{\mathbf{R}}(\widehat{\mathbf{m}}_{\text{NN}}^{\mathbf{W}}) + \widehat{\mathbf{R}}(\widehat{\mathbf{m}}_{\text{NN}}^{\mathbf{W}}) - \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}_{\star}}) + \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}_{\star}}) - \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}_{\star}}) + \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}_{\star}}) - \mathbf{R}(\mathbf{m})$$

$$1752 \leq 2 \cdot \sup_{\mathbf{W} \in \mathcal{W}} \left| \mathbf{R}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) - \widehat{\mathbf{R}}(\mathbf{m}_{\text{NN}}^{\mathbf{W}}) \right| + D_2^2(\mathbf{m}_{\text{NN}}^{\mathbf{W}_{\star}}, \mathbf{m})$$

1753 This proves Lemma 22.  $\square$

### 1754 H.2 RESULTS ON GENERALIZATION

1755 **Proposition 8** (Generalization error of the denoising problem). *Let  $\mathcal{W}_{d, \underline{m}, L, S, D, B}$  be the set defined  
1756 as in Eq. (9). Then, with probability at least  $1 - \eta$ , we have*

$$1757 \sup_{\mathbf{W} \in \mathcal{W}_{d, \underline{m}, L, S, D, B}} \left| \widehat{\mathbf{R}}(\mu_{\text{NN}}^{\mathbf{W}}) - \mathbf{R}(\mu_{\text{NN}}^{\mathbf{W}}) \right| \leq C \cdot S^2 \cdot \sqrt{\frac{LD(D + 2S + 1) \|\underline{m}\|_1 \log(d \|\underline{m}\|_1 DSB \cdot 18^L) + \log(1/\eta)}{n}}.$$

*Proof of Proposition 8.* In Lemma 3, we can take  $z = (z, \mathbf{x})$ ,  $w = \mathbf{W}$ ,  $\Theta = \mathcal{W}_{d, \underline{m}, L, S, D, B}$ ,  $\rho(w, w') = \|\mathbf{W} - \mathbf{W}'\|$ , and  $f(z_i; w) = \|\mathbf{x} - \mathbf{m}_{\text{NN}}^{\mathbf{W}}(z)\|_2^2$ . Therefore, to show Proposition 8, we just need to apply Lemma 3 by checking (a), (b), (c).

**Check (a).** We note that the index set  $\Theta := \mathcal{W}_{d, \underline{m}, L, S, D, B}$  equipped with  $\rho(w, w') := \|\mathbf{W} - \mathbf{W}'\|$  has diameter  $B_p := 2B$ . Further note that  $\mathcal{W}_{d, \underline{m}, L, S, D, B}$  has a dimension bounded by  $d_p := 2D(D + 2S + 1)\|\underline{m}\|_1$ . According to Example 5.8 of (Wainwright, 2019), it holds that  $\log N(\Delta; \mathcal{W}_{d, \underline{m}, L, S, D, B}, \|\cdot\|) \leq d_p \cdot \log(1 + 2r/\Delta)$  for any  $0 < \Delta \leq 2r$ . This verifies (a).

**Check (b).** Since  $f(z_i; w) = d^{-1}\|\mathbf{x} - \mathbf{m}_{\text{NN}}^{\mathbf{W}}(z)\|_2^2$  is  $S^2$ -bounded. As a consequence,  $f(z, w) - \mathbb{E}_z[f(z, w)]$  is a sub-Gaussian random variable with the sub-Gaussian parameter to be  $C \cdot S^2$ .

**Check (c).** Lemma 25 implies that

$$|f(z; w_1) - f(z; w_2)| \leq L_p \cdot \|\mathbf{W}_1 - \mathbf{W}_2\|, \quad L_p := 12\|\underline{m}\|_1 18^L B^{6L} \cdot d \cdot S^3 \left( \sum_{v \in \mathcal{V}^{(L)}} (S + |z_v|) \right) \cdot (S + D).$$

Since  $z \stackrel{d}{=} \mathbf{x} + \mathbf{g}$  where  $(\mathbf{x}, \mathbf{g}) \sim \mu_* \times \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , and  $\|\mathbf{x}\|_1 \leq Sd$ ,  $\|\mathbf{g}\|_1$  is  $Cd$ -sub-Gaussian. Hence  $\|z\|_1$  is  $CSd$ -sub-Gaussian, and hence  $f(z; w_1) - f(z; w_2)$  is  $\sigma' \rho(w_1, w_2)$  sub-Gaussian with

$$\sigma' = C\|\underline{m}\|_1 18^L B^{6L} \cdot d^2 \cdot S^4 \cdot (S + D).$$

Therefore, we apply Lemma 3 to conclude the proof of Proposition 8.  $\square$

### H.3 AUXILLARY LEMMAS

**Lemma 23** (Norm bound in the chain rule in denoising settings). *Consider the U-Net as in Eq. (UNet) with modified input  $\mathbf{h}_{\downarrow, v}^{(L)} = (-x^2/2 + xz_v)_{x \in [S]} \in \mathbb{R}^S$  (Since we will immediately normalize the input, this input is effectively the same as the input  $\mathbf{h}_{\downarrow, v}^{(L)} = (-(x - z_v)^2/2)_{x \in [S]} \in \mathbb{R}^S$ ). Assume that  $\|\mathbf{W}\| \leq B$ . Then for any  $\ell, v, \iota$ , and  $\star \in \{\downarrow, \uparrow\}$ , we have*

$$\begin{aligned} \|\mathbf{h}_{\downarrow, v}^{(L)}\|_2 &\leq S^3 + S^2|z_v|, \\ \|\mathbf{q}_{\downarrow, v}^{(\ell)}\|_2 &\leq B^3 \cdot (2 \cdot \|\mathbf{h}_{\downarrow, v}^{(\ell)}\|_2 + 1), \\ \|\mathbf{h}_{\downarrow, v}^{(\ell-1)}\|_2 &\leq m^{(\ell)} \cdot \max_{v' \in \mathcal{C}(v)} \|\mathbf{q}_{\downarrow, v'}^{(\ell)}\|_2, \\ \|\mathbf{b}_{\uparrow, \iota}^{(0)}\|_2 &= \|\mathbf{h}_{\downarrow, \iota}^{(0)}\|_2, \\ \|\mathbf{b}_{\uparrow, v}^{(\ell)}\|_2 &\leq B^3 \cdot (2\|\mathbf{b}_{\uparrow, \text{pa}(v)}^{(\ell-1)}\|_2 + 2\|\mathbf{q}_{\downarrow, v}^{(\ell)}\|_2 + 1) + \|\mathbf{h}_{\downarrow, v}^{(\ell)}\|_2, \\ \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{q}_{\downarrow, v, i}^{(\ell)}\|_{\text{op}} &\leq 2B^3 \cdot \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{h}_{\downarrow, v, i}^{(\ell)}\|_{\text{op}}, \quad \forall k \geq \ell + 1, \\ \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(\ell)}} \mathbf{q}_{\downarrow, v, i}^{(\ell)}\|_{\text{op}} &\leq 2B^2 \cdot \|\mathbf{h}_{\downarrow, v}^{(\ell)}\|_2, \\ \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{h}_{\downarrow, v, i}^{(\ell-1)}\|_{\text{op}} &\leq m^{(\ell)} \cdot \max_{v' \in \mathcal{C}(v)} \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{q}_{\downarrow, v', i}^{(\ell)}\|_{\text{op}}, \quad \forall k \geq \ell, \\ \max_{i \in [S]} \|\nabla_{W_{\star, \uparrow, \iota}^{(k)}} \mathbf{b}_{\uparrow, v, i}^{(\ell)}\|_{\text{op}} &\leq 2B^3 \cdot \max_{i \in [S]} \|\nabla_{W_{\star, \uparrow, \iota}^{(k)}} \mathbf{b}_{\uparrow, \text{pa}(v), i}^{(\ell-1)}\|_{\text{op}}, \quad \forall k \geq \ell + 1, \\ \max_{i \in [S]} \|\nabla_{W_{\star, \uparrow, \iota}^{(\ell)}} \mathbf{b}_{\uparrow, v, i}^{(\ell)}\|_{\text{op}} &\leq 2B^2 \cdot (\|\mathbf{b}_{\uparrow, \text{pa}(v)}^{(\ell-1)}\|_2 + \|\mathbf{q}_{\downarrow, v}^{(\ell)}\|_2), \\ \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{b}_{\uparrow, v, i}^{(\ell)}\|_{\text{op}} &\leq 2B^3 \cdot \left( \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{b}_{\uparrow, \text{pa}(v), i}^{(\ell-1)}\|_{\text{op}} \right. \\ &\quad \left. + \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{q}_{\downarrow, v, i}^{(\ell)}\|_{\text{op}} \right) + \max_{i \in [S]} \|\nabla_{W_{\star, \downarrow, \iota}^{(k)}} \mathbf{h}_{\downarrow, v, i}^{(\ell)}\|_{\text{op}}, \quad \forall k \in [L], \\ \|\nabla_{W_{\star, \diamond, \iota}^{(k)}} \mathbf{m}_{\text{NN}}^{\mathbf{W}}(z)_v\|_{\text{op}} &\leq S \cdot \max_{i \in [S]} \|\nabla_{W_{\star, \diamond, \iota}^{(k)}} \mathbf{b}_{\uparrow, v, i}^{(L)}\|_{\text{op}}, \quad \forall k \in [L], \diamond \in \{\downarrow, \uparrow\}. \end{aligned}$$

*Proof of Lemma 23.* The proof of the lemma uses the chain rule, the 1-Lipschitzness of ReLU, the 2-Lipschitzness of normalize, and the 1-Lipschitzness of softmax.  $\square$

1836 **Lemma 24.** Consider the U-Net as in Eq. (UNet). Assume that  $\|\mathbf{W}\| \leq B$ . Then for any  $v \in \mathcal{V}^{(L)}$ ,  
 1837 we have

$$1838 \max_{\diamond \in \{\downarrow, \uparrow\}} \max_{\mathbf{z} \in \mathbb{R}^d} \max_{\star \in [3]} \max_{\ell \in [L]} \max_{\iota \in [m^{(\ell)}]} \|\nabla_{W_{\star, \diamond, \iota}^{(\ell)}} \mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z})_v\|_{\text{op}} \leq 18^L B^{6L} \cdot d \cdot S^3 (S + |z_v|).$$

1841 *Proof of Lemma 24.* This lemma is implied by Lemma 23 and an induction argument.  $\square$

1842 **Lemma 25.** Consider the ConvNet as in Eq. (UNet). Assume that  $\|\mathbf{W}\| \leq B$ . Then for any  
 1843  $v \in \mathcal{V}^{(L)}$ , we have

$$1844 \max_{\mathbf{z}} \left| \mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z})_v - \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})_v \right| \leq 6 \|\underline{m}\|_1 18^L B^{6L} \cdot d \cdot S^3 (S + |z_v|) \cdot (S + D) \cdot \|\mathbf{W} - \overline{\mathbf{W}}\|.$$

1845 Therefore, we have

$$1846 \max_{\mathbf{z}} d^{-1} \left| \|\mathbf{x} - \mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z})\|_2^2 - \|\mathbf{x} - \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})\|_2^2 \right| \leq 12 \|\underline{m}\|_1 18^L B^{6L} \cdot d \cdot S^3 \left( \sum_{v \in \mathcal{V}^{(L)}} (S + |z_v|) \right) \cdot (S + D) \cdot \|\mathbf{W} - \overline{\mathbf{W}}\|.$$

1847 *Proof of Lemma 25.* The first inequality is by the fact that

$$1848 \max_{\mathbf{z}} \left| \mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z})_v - \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})_v \right|$$

$$1849 \leq \sum_{\diamond \in \{\downarrow, \uparrow\}} \sum_{\star \in [3]} \sum_{\ell \in [L]} \sum_{\iota \in [m^{(\ell)}]} \min\{\text{nrow}(W_{\star, \diamond, \iota}^{(\ell)}), \text{ncol}(W_{\star, \diamond, \iota}^{(\ell)})\} \|\nabla_{W_{\star, \diamond, \iota}^{(\ell)}} \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})_v\|_{\text{op}} \|W_{\star, \diamond, \iota}^{(\ell)} - \overline{W}_{\star, \diamond, \iota}^{(\ell)}\|_{\text{op}},$$

1850 where we have used the inequality that  $\text{trace}(A^T B) \leq \{\text{nrow}(A), \text{ncol}(A)\} \|A\|_{\text{op}} \|B\|_{\text{op}}$ .

1851 To prove the second inequality, we have

$$1852 \max_{\mathbf{z}} \left| \|\mathbf{x} - \mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z})\|_2^2 - \|\mathbf{x} - \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})\|_2^2 \right|$$

$$1853 \leq \max_{\mathbf{z}} \|2\mathbf{x} - \mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z}) - \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})\|_{\infty} \|\mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z}) - \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})\|_1$$

$$1854 \leq 2d \max_{\mathbf{z}} \|\mathbf{m}_{\text{NN}}^{\mathbf{W}}(\mathbf{z}) - \mathbf{m}_{\text{NN}}^{\overline{\mathbf{W}}}(\mathbf{z})\|_1.$$

1855 This completes the proof of Lemma 25.  $\square$