# Approximation Capabilities of Feedforward Neural Networks with GELU Activations

Konstantin Yakovlev[*]        Nikita Puchkin[†]

## Abstract

We derive an approximation error bound that holds simultaneously for a function and all its derivatives up to any prescribed order. The bounds apply to elementary functions, including multivariate polynomials, the exponential function, and the reciprocal function, and are obtained using feedforward neural networks with the Gaussian Error Linear Unit (GELU) activation. In addition, we report the network size, weight magnitudes, and behavior at infinity. Our analysis begins with a constructive approximation of multiplication, where we prove the simultaneous validity of error bounds over domains of increasing size for a given approximator. Leveraging this result, we obtain approximation guarantees for division and the exponential function, ensuring that all higher-order derivatives of the resulting approximators remain globally bounded.

# Contents

[*]HSE University, Russian Federation, kdyakovlev@hse.ru
[†]HSE University, Russian Federation, npuchkin@hse.ru

1

# 1   Introduction

We investigate simultaneous approximation of multivariate functions and their higher-order derivatives using deep feedforward neural networks. Since approximating derivatives necessitates a smooth activation function, we employ networks with infinitely differentiable the Gaussian Error Linear Unit (GELU) [Hendrycks and Gimpel, 2016]. Our choice is motivated by the fact that higher-order derivatives of the GELU activation can be expressed in terms of Hermite polynomials, which enables simple and tractable bounds on their absolute values. In addition, this activation function is widely adopted in state-of-the-art large language models [Devlin et al., 2019, Raffel et al., 2020, Shoeybi et al., 2019], which further underscores its practical relevance.

The core of our constructive approach is the localized approximation of polynomials. While prior works have developed approximation theory for smooth functions and their derivatives on fixed compact sets [Yarotsky, 2017, De Ryck et al., 2021, Gühring and Raslan, 2021, Belomestny et al., 2023], they do not provide error bounds beyond the original domain nor offer simultaneous guarantees across a sequence of increasingly large domains. Despite recent advances in approximation of functions with noncompact domain presented in Schwab and Zech [2021], van Nuland [2024], Abdeljawad and Dittrich [2024], the results either do not focus on the simultaneous approximation of derivatives or impose strong assumptions on the weight function of the underlying weighted $L^p$ space [Abdeljawad and Dittrich, 2024]. We bridge this gap by providing explicit control on how approximation errors for fundamental operations (like multiplication) scale as the domain size grows. This allows us to construct an approximation of monomials with globally bounded higher-order derivatives. This properties are relevant in approximation of functions with unbounded domains including generative modelling [Oko et al., 2023, Tang and Yang, 2024, Azangulov et al., 2024, Yakovlev and Puchkin, 2025, Fukumizu et al., 2025] and physics-informed neural networks [Abdo et al., 2024, Alejo et al., 2024].

Our key technical innovations are twofold. First, we systematically employ a clipping operation on the network input. By clipping the neural network input, we ensure that the derivatives are globally bounded, since they are bounded on a compact domain. Second, we establish approximation guarantees for partition-of-unity functions in Sobolev seminorms, constructing functions with globally bounded derivatives and light tails.

As a consequence of these results, we provide approximation error bounds for the exponential function approximation and division approximation together with its derivatives, ensuring that their higher-order derivatives are globally bounded. Consequently, we extend the approximation results for elementary functions presented in Oko et al. [2023], Yakovlev and Puchkin [2025] to Sobolev seminorms on domains of increasing size.

**Paper structure.**   The remainder of the paper is structured as follows. Section 2 establishes necessary preliminaries and notations. In Section 3, we present our main result on approximation error bounds. Proofs not included in the main text are provided in the Appendix.

**Notation.**   The set of non-negative integers is denoted by $\mathbb{Z}_+ = \{0, 1, 2, \dots\}$. A multi-index $\mathbf{k} \in \mathbb{Z}_+^d$, where $d \in \mathbb{N}$, is denoted in bold. We also define $|\mathbf{k}| = k_1 + k_2 + \dots + k_d$, $\mathbf{k}! = k_1! \cdot k_2! \cdot \dots \cdot k_d!$ For a vector $v \in \mathbb{R}^d$ we define $v^{\mathbf{k}} = v_1^{k_1} v_2^{k_2} \dots v_d^{k_d}$. For a function $f$ of $d$ variables, its weak derivative with

respect to the multi-index $\mathbf{k} \in \mathbb{Z}_+^d$ is denoted as

$$\partial^{\mathbf{k}} f = \frac{\partial^{|\mathbf{k}|} f}{\partial x_1^{k_1} \partial x_2^{k_2} \ldots \partial x_d^{k_d}}.$$

Throughout the paper, we employ the notation $f \lesssim g$ to indicate that $f = \mathcal{O}(g)$. If $f \lesssim g$ and $g \lesssim f$, then we write $f \asymp g$. We frequently replace the expression for $\min\{a, b\}$ and $\max\{a, b\}$ with $a \vee b$ and $a \wedge b$, respectively. For any $x > 0$, we define $\log(x) = \ln(x \vee e)$.

## 2 Preliminaries and notations

**Norms.** We denote the Euclidean norm of a vector $v$ as $\|v\|$, the maximal absolute value of its entries as $\|v\|_\infty$, and the number of its non-zero entries as $\|v\|_0$. Similarly, $\|A\|_\infty$ and $\|A\|_0$ represent the maximal absolute value of entries and the number of non-zero entries of matrix $A$, respectively. Finally, for a set $\Omega \subseteq \mathbb{R}^r$ and a function $f : \Omega \to \mathbb{R}^d$, we define

$$\|f\|_{L^\infty(\Omega)} = \operatorname*{esssup}_{x \in \Omega} \|f(x)\|.$$

**Smoothness spaces.** We introduce the Sobolev space to characterize the regularity of functions in our analysis, and its definition is provided below.

**Definition 2.1** (Sobolev space). *Let $\Omega \subseteq \mathbb{R}^r$ be an open set, and let $k \in \mathbb{Z}_+$. Then, the Sobolev space $W^{k,\infty}(\Omega)$ is defined as follows:*

$$W^{k,\infty}(\Omega) = \{f \in L^\infty(\Omega) : \partial^{\mathbf{k}} f \in L^\infty(\Omega) \quad \text{for every } \mathbf{k} \in \mathbb{Z}_+^r \text{ with } |\mathbf{k}| \leqslant k\}.$$

*Here, $L^\infty(\Omega)$ is the Lebesgue space. We define the Sobolev seminorm on $W^{k,\infty}(\Omega)$ as*

$$|f|_{W^{k,\infty}(\Omega)} = \max_{\mathbf{k} \in \mathbb{Z}_+^r,\, |\mathbf{k}|=k} \|\partial^{\mathbf{k}} f\|_{L^\infty(\Omega)}.$$

*Finally, we define the Sobolev norm on $W^{k,\infty}(\Omega)$ as*

$$\|f\|_{W^{k,\infty}(\Omega)} = \max_{0 \leqslant m \leqslant k} |f|_{W^{m,\infty}(\Omega)}.$$

**Neural networks.** In this paper, we focus on feed-forward neural networks employing the Gaussian Error Linear Unit (GELU) activation function:

$$\mathrm{GELU}(x) = x \cdot \Phi(x), \quad \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-t^2/2} \mathrm{d}t.$$

The choice of GeLU is motivated by its infinite smoothness and bounded derivatives (see Lemma B.2). For a vector $b = (b_1, \ldots, b_r) \in \mathbb{R}^r$, we define the shifted activation function $\mathrm{GELU}_b : \mathbb{R}^r \to \mathbb{R}^r$ as

$$\mathrm{GELU}_b(x) = (\mathrm{GELU}(x_1 - b_1), \ldots, \mathrm{GELU}(x_r - b_r)), \quad x = (x_1, \ldots, x_r) \in \mathbb{R}^r.$$

Given a network depth $L \in \mathbb{N}$ and a vector of layer sizes $W = (W_0, W_1, \ldots, W_L) \in \mathbb{N}^{L+1}$. Then, a neural network of depth $L$ and architecture $W$ is a function $f : \mathbb{R}^{W_0} \to \mathbb{R}^{W_L}$ such that

$$f(x) = -b_L + A_L \circ \text{GELU}_{b_{L-1}} \circ A_{L-1} \circ \text{GELU}_{b_{L-2}} \circ \cdots \circ A_2 \circ \text{GELU}_{b_1} \circ A_1 \circ x, \qquad (1)$$

where $A_j \in \mathbb{R}^{W_j \times W_{j-1}}$ is a weight matrix and $b_j \in \mathbb{R}^{W_j}$ is a bias vector for all $j \in 1, \ldots, L$. The maximum number of neurons of each layer is given by $\|W\|_\infty$ and is reffered to as the width of the neural network. We define the class $\text{NN}(L, W, S, B)$ of neural networks of the form (1) with at most $S$ non-zero weights and the weight magnitude $B$ as follows:

$$\text{NN}(L, W, S, B) = \left\{ f \text{ of the form } (1) : \sum_{j=1}^{L}(\|A_j\|_0 + \|b_j\|_0) \leqslant S, \ \max_{1 \leqslant j \leqslant L} \|A_j\|_\infty \vee \|b_j\|_\infty \leqslant B \right\}.$$

## 3 Main results

This section presents our main results. Specifically, Subsection 3.1 details approximation error bounds for elementary operations, including the identity function, partition of unity, and square operation. Subsection 3.2 elaborates on the approximation of monomials. Finally, Subsection 3.3 provides approximation error bounds for the exponentiation and division operations.

### 3.1 Approximation of elementary operations

Passing the output of one layer to a non-adjacent layer is frequently beneficial. Note that ReLU activation allows for an exact identity mapping [Nakada and Imaizumi, 2020a]. However, in the case of GELU, an approximate mapping is guaranteed, as demonstrated by the following lemma, which provides an approximation error bound for a single-layer neural network.

**Lemma 3.1** (approximation of identity operation). *Let $m \in \mathbb{N}$ and let $\text{id}(x) = x$. Then, for any $\varepsilon \in (0, 1)$ there exists $\varphi_{id} \in \text{NN}(L, W, S, B)$ satisfying*

$$\|\varphi_{id} - \text{id}\|_{W^{m,\infty}([-C,C])} \leqslant C^2 \varepsilon, \quad \text{for all } C \geqslant 1.$$

*Furthermore, $L = 2$, $\|W\|_\infty = 1$, $S = 3$, and $\log B \lesssim \log(1/\varepsilon) + \log m$.*

*Proof.* The proof follows the same approach as outlined in Scarselli and Tsoi [1998]. We let

$$\varphi_{id}(x) = -\frac{R \cdot \text{GELU}(0)}{\partial^1 \text{GELU}(0)} + \frac{R}{\partial^1 \text{GELU}(0)} \text{GELU}\left(\frac{x}{R}\right),$$

where $R > 0$ and will be determined later in the proof. We also emphasize that the form of $\varphi_{id}$ is valid, since $\partial^1 \text{GELU}(0) = 1/2$. Taylor expansion suggests that for any $x \in [-C, C]$ it holds that

$$|\varphi_{id}(x) - x| \leqslant \frac{|\text{GELU}|_{W^{2,\infty}(\mathbb{R})} C^2}{2 \cdot \partial^1 \text{GELU}(0) R}.$$

Similarly, we deduce that

$$|\partial^1 \varphi_{id}(x) - 1| \leqslant \frac{|\text{GELU}|_{W^{2,\infty}(\mathbb{R})} C}{\partial^1 \text{GELU}(0) R}.$$

4

Additionally, for any $k \geqslant 2$ we find that

$$|\varphi_{id} - \mathrm{id}|_{W^{k,\infty}(\mathbb{R})} = |\varphi_{id}|_{W^{k,\infty}(\mathbb{R})} \leqslant \frac{|\mathrm{GELU}|_{W^{k,\infty}(\mathbb{R})}}{\partial^1 \mathrm{GELU}(0) R^{k-1}}.$$

Therefore, choosing

$$R = \max_{2 \leqslant k \leqslant m} \left( \frac{|\mathrm{GELU}|_{W^{k,\infty}(\mathbb{R})}}{\partial^1 \mathrm{GELU}(0)\varepsilon} \right)^{1/(k-1)} \vee 1,$$

ensures that for any $C \geqslant 1$

$$\max_{0 \leqslant k \leqslant 1} |\varphi_{id} - \mathrm{id}|_{W^{k,\infty}([-C,C])} \leqslant C^2 \varepsilon, \quad \max_{2 \leqslant k \leqslant m} |\varphi_{id} - \mathrm{id}|_{W^{k,\infty}(\mathbb{R})} \leqslant \varepsilon.$$

Therefore, it holds that

$$\|\varphi_{id} - \mathrm{id}\|_{W^{m,\infty}([-C,C])} \leqslant C^2 \varepsilon, \quad \text{for all } C \geqslant 1.$$

We next specify the configuration of $\varphi_{id}$. Clearly, $L = 2$, $\|W\|_\infty = 1$ and $S = 3$. As for the weight magnitude, we apply Lemma B.2, arriving at

$$\log B \lesssim \log(1/\varepsilon) + \max_{2 \leqslant k \leqslant m} \frac{\log \left( |\mathrm{GELU}|_{W^{k,\infty}(\mathbb{R})} \vee 1 \right)}{k-1} \lesssim \log(1/\varepsilon) + \max_{2 \leqslant k \leqslant m} \frac{\log\left( (k+1)! \right)}{k-1}.$$

Now Stirling's approximation implies that $\log((k+1)!) \lesssim k \log k$, and thus,

$$\log B \lesssim \log(1/\varepsilon) + \log m.$$

The proof is complete.

$\square$

The following lemma generalizes the result presented in Lemma 3.1 to the case of multiple layers.

**Lemma 3.2** (approximation of identity operation with multiple layers)**.** *Let $m \in \mathbb{N}$ and let* $\mathrm{id} : x \mapsto x$. *Then, for every $\varepsilon \in (0,1)$, every $L \in \mathbb{N}$ with $L \geqslant 2$, and every $K \geqslant 1$, there exists $\varphi_{id} \in \mathsf{NN}(L, W, S, B)$ such that*

$$
\begin{aligned}
(i) &\quad \|\varphi_{id} - \mathrm{id}\|_{W^{m,\infty}([-K,K])} \leqslant \varepsilon, \\
(ii) &\quad \|\varphi_{id}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m\log(m\log(1/\varepsilon)) + \log(2K))\}.
\end{aligned}
$$

*Moreover, it holds that*

$$\|W\|_\infty \lesssim 1, \quad S \lesssim L, \quad \log B \lesssim (m+L)\log m + \log(1/\varepsilon) + m\log(K).$$

The proof of Lemma 3.2 is moved to Appendix A.1. Next, we move to the approximation of partition of unity, a crucial component in the framework of localized Taylor polynomials [Gühring and Raslan, 2021, De Ryck et al., 2021]. First, we approximate the Heaviside step function, as presented in the following lemma.

**Lemma 3.3** (Approximation of Heaviside step function). *For every $\varepsilon \in (0,1)$, every $\varkappa \in (0,1)$, and every $m \in \mathbb{N}$, there exists a GELU network $\varphi_\varkappa \in \mathsf{NN}(L,W,S,B)$ such that*

$$(i) \quad \|\varphi_\varkappa\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\left\{\mathcal{O}(m\log(m\log(1/\varepsilon)/\varkappa))\right\},$$

$$(ii) \quad \|\varphi_\varkappa\|_{W^{m,\infty}((-\infty,-\varkappa])} \vee \|1 - \varphi_\varkappa\|_{W^{m,\infty}([\varkappa,+\infty))} \leqslant \varepsilon$$

*Moreover, $\varphi_\varkappa$ has $L = 2$, $\|W\|_\infty \vee S \lesssim 1$, and $\log B \lesssim m\log(m/\varkappa) + \log(1/\varepsilon)$.*

*Proof.* Let

$$\eta(x) = \frac{\mathrm{GELU}(x+\varepsilon_0) - \mathrm{GELU}(x-\varepsilon_0)}{2\varepsilon_0}, \quad x \in \mathbb{R},$$

where $\varepsilon_0 \in (0,1)$ will be determined later. Note that $\eta$ approximates $\partial^1\mathrm{GELU}$ is Sobolev norm. Formally, the Taylor expansion suggests that for any $k \in \mathbb{Z}_+$ and $x \in \mathbb{R}$ we have

$$|\partial^k\eta(x) - \partial^{k+1}\mathrm{GELU}(x)| = \frac{|\partial^k\mathrm{GELU}(x+\varepsilon_0) - \partial^k\mathrm{GELU}(x-\varepsilon_0) - 2\varepsilon\partial^{k+1}\mathrm{GELU}(x)|}{2\varepsilon_0}$$

$$\leqslant \frac{\varepsilon_0^2}{6}|\mathrm{GELU}|_{W^{k+3,\infty}(\mathbb{R})}.$$

Hence, it holds that

$$\|\eta - \partial^1\mathrm{GELU}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \frac{\varepsilon_0^2}{6}\max_{0\leqslant k\leqslant m}|\mathrm{GELU}|_{W^{k+3,\infty}(\mathbb{R})}. \tag{2}$$

Now let $\varphi_\varkappa(x) = \eta(\alpha x)$ for $\alpha \geqslant 1$ that will be optimized later. Therefore, triangle inequality yields

$$\|\varphi_\varkappa\|_{W^{m,\infty}([-\infty,-\varkappa])} \leqslant \alpha^m\|\eta\|_{W^{m,\infty}((-\infty,-\alpha\varkappa])}$$
$$\leqslant \alpha^m\left(\|\partial^1\mathrm{GELU} - \eta\|_{W^{m,\infty}((-\infty,-\alpha\varkappa])} + \|\partial^1\mathrm{GELU}\|_{W^{m,\infty}((-\infty,-\alpha\varkappa])}\right).$$

Thus, Lemma B.2 together with (2) implies that

$$\|\varphi_\varkappa\|_{W^{m,\infty}([-\infty,-\varkappa])} \leqslant \alpha^m\left(\frac{\varepsilon_0^2}{6}(m+4)\sqrt{(m+1)!} + 2e^{-\alpha^2\varkappa^2/4}\sqrt{(m+1)!}\right)$$

Setting $\alpha = 2\varkappa^{-1}\sqrt{2\log(1/\varepsilon_0)} \geqslant 1$ ensures that

$$\|\varphi_\varkappa\|_{W^{m,\infty}([-\infty,-\varkappa])} \leqslant \alpha^m\varepsilon_0^2(m+3)\sqrt{(m+1)!} \leqslant (8\varkappa^{-2}\log(1/\varepsilon_0))^{m/2}\varepsilon_0^2(m+3)\sqrt{(m+1)!}.$$

Using the fact that

$$\sup_{\varepsilon\in(0,1)}\varepsilon(\log(1/\varepsilon))^{m/2} \leqslant \left(\frac{m}{2e}\right)^{m/2} \leqslant m^{m/2}$$

we find that

$$\|\varphi_\varkappa\|_{W^{m,\infty}((-\infty,-\varkappa])} \leqslant (8\kappa^{-2}m)^{m/2}\varepsilon_0(m+3)\sqrt{(m+1)!}.$$

6

Hence, setting

$$\varepsilon_0 = \left( (8\kappa^{-2}m)^{m/2}(m+3)\sqrt{(m+1)!} \right)^{-1} \varepsilon \in (0,1) \tag{3}$$

ensures that

$$\|\varphi_{\varkappa}\|_{W^{m,\infty}((-\infty,-\varkappa])} \leqslant \varepsilon.$$

Using similar argument, we also deduce that

$$\|1 - \varphi_{\varkappa}\|_{W^{m,\infty}([\kappa,+\infty))} \leqslant \varepsilon.$$

We next note that in view of (2) and Lemma B.2, it holds that

$$\|\varphi_{\varkappa}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \alpha^m (\|\eta - \partial^1 \mathrm{GELU}\|_{W^{m,\infty}(\mathbb{R})} + \|\partial^1 \mathrm{GELU}\|_{W^{m,\infty}(\mathbb{R})})$$
$$\leqslant \alpha^m (\frac{\varepsilon_0^2}{6}(m+4)\sqrt{(m+1)!} + (m+2)\sqrt{(m-1)!}).$$

The choice of $\alpha$ indicates that

$$\|\varphi_{\varkappa}\|_{W^{m,\infty}(\mathbb{R})} \leqslant 2\alpha^m (m+2)\sqrt{(m+1)!} \leqslant 2(8\kappa^{-2}\log(1/\varepsilon_0))^{m/2}(m+2)\sqrt{(m+1)!}.$$

Now the choice of $\varepsilon_0$ from (3) yields

$$\|\varphi_{\varkappa}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\left\{ \mathcal{O}(m\log(m\log(1/\varepsilon)/\varkappa)) \right\}.$$

Finally, we specify the configuration of $\varphi_{\varkappa}$. Clearly, $L = 2$, $\|W\|_\infty \vee S \lesssim 1$, and

$$\log B \lesssim \log(\alpha) + \log(1/\varepsilon_0) \lesssim m\log(m/\varkappa) + \log(1/\varepsilon).$$

The proof is complete.

$\square$

Subsequently, we use Lemma 3.3 to approximate a partition of unity. Following Yakovlev and Puchkin [2025], we use non-uniform partition, a key element in approximating the division operation. The result is presented below.

**Lemma 3.4** (partition of unity approximation). *Define $a_i = 2^{-N+i}$ for each $i \in \{0, 1, \ldots, N\}$, where $N \in \mathbb{N}$ and $N \geqslant 3$. Then, for every $\varepsilon \in (0,1)$ and every $m \in \mathbb{N}$, there exist $\{\psi_i\}_{i=1}^N$, with $\psi_i \in \mathsf{NN}(L, W, S, B)$ for each $1 \leqslant i \leqslant N$, such that*

$$(i) \quad \sum_{i=1}^N \psi_i(x) = 1, \quad \text{for all } x \in \mathbb{R},$$

$$(ii) \quad \max_{1 \leqslant i \leqslant N} \|\psi_i\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(mN + m\log(m\log(1/\varepsilon)))\},$$

$$(iii) \quad \|\psi_N\|_{W^{m,\infty}(-\infty,a_{N-2}]} \vee \|\psi_1\|_{W^{m,\infty}([a_2,+\infty))} \vee \max_{2 \leqslant i \leqslant N-1} \|\psi_i\|_{W^{m,\infty}(\mathbb{R}\backslash(a_{i-2},a_{i+1}))} \leqslant \varepsilon,$$

*Furthermore, $L = 2$, $\|W\|_\infty \vee S \lesssim 1$ and $\log B \lesssim \log(1/\varepsilon) + mN + m\log m$.*

*Proof.* Next, for a Heaviside function approximation $\varphi_{a_0}$ from Lemma 3.3 formulated with accuracy parameter $\varepsilon/2$ and $\kappa = a_0$, we define

$$\psi_i(x) = \begin{cases} 1 - \varphi_{a_0}(x - a_1), & i = 1, \\ \varphi_{a_0}(x - a_{i-1}) - \varphi_{a_0}(x - a_i), & i \in \{2, \ldots, N-1\}, \\ \varphi_{a_0}(x - a_{N-1}), & i = N \end{cases}$$

It is clear that for all $x \in \mathbb{R}$

$$\sum_{i=1}^{N} \psi_i(x) = 1.$$

In other words, $\{\psi_i\}_{i=1}^{N}$ forms a partition of unity. Now derive the behavior of tails for each $\psi_i$. First, note that for each $1 \leqslant i \leqslant N$ we have that

$$\|\varphi_{a_0}(\cdot - a_i)\|_{W^{m,\infty}((-\infty, a_{i-1}])} \leqslant \|\varphi_{a_0}\|_{W^{m,\infty}((-\infty, -a_0])} \leqslant \varepsilon/2 \tag{4}$$

and similarly

$$\|1 - \varphi_{a_0}(\cdot - a_i)\|_{W^{m,\infty}([a_{i+1}, +\infty))} \leqslant \|1 - \varphi_{a_0}\|_{W^{m,\infty}([a_0, +\infty))} \leqslant \varepsilon/2. \tag{5}$$

Therefore,

$$\|\psi_N\|_{W^{m,\infty}(-\infty, a_{N-2}]} \vee \|\psi_1\|_{W^{m,\infty}([a_2, +\infty))} \leqslant \varepsilon.$$

Next, for any $2 \leqslant i \leqslant N - 1$ it holds that

$$\|\psi_i\|_{W^{m,\infty}(\mathbb{R} \setminus (a_{i-2}, a_{i+1}))} = \|\psi_i\|_{W^{m,\infty}((-\infty, a_{i-2}])} \vee \|\psi_i\|_{W^{m,\infty}([a_{i+1}, +\infty))}$$

First, from (4) we find that

$$\|\psi_i\|_{W^{m,\infty}((-\infty, a_{i-2}])} \leqslant \|\varphi_{a_0}(\cdot - a_{i-1})\|_{W^{m,\infty}((-\infty, a_{i-2}])} + \|\varphi_{a_0}(\cdot - a_i)\|_{W^{m,\infty}((-\infty, a_{i-2}])} \leqslant \varepsilon.$$

Second, (5) implies that

$$\|\psi_i\|_{W^{m,\infty}([a_{i+1}, +\infty))} \leqslant \|1 - \varphi_{a_0}(\cdot - a_{i-1})\|_{W^{m,\infty}([a_{i+1}, +\infty))} + \|1 - \varphi_{a_0}(\cdot - a_i)\|_{W^{m,\infty}([a_{i+1}, +\infty))} \leqslant \varepsilon.$$

Thus, we arrive at

$$\|\psi_N\|_{W^{m,\infty}(-\infty, a_{N-2}]} \vee \|\psi_1\|_{W^{m,\infty}([a_2, +\infty))} \vee \max_{2 \leqslant i \leqslant N-1} \|\psi_i\|_{W^{m,\infty}(\mathbb{R} \setminus (a_{i-2}, a_{i+1}))} \leqslant \varepsilon.$$

Now we focus on the behavior of each $\psi_i$ on the real line. Formally, Lemma 3.3 suggests that for any $1 \leqslant i \leqslant N$

$$\|\psi_i\|_{W^{m,\infty}(\mathbb{R})} \leqslant 2\|\varphi_{a_0}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\left\{\mathcal{O}(m \log(m \log(1/\varepsilon)/a_0))\right\}.$$

Recall that $a_0 = 2^{-N}$. Hence, it holds that

$$\|\psi_i\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(mN + m \log(m \log(1/\varepsilon)))\}.$$

8

We now specify the configuration for each $\psi_i$. Using the configuration of $\varphi_{a_0}$ outlined in Lemma 3.3 and parallelization argument from Lemma B.6, we conclude that

$$L = 2, \quad \|W\|_\infty \vee S \lesssim 1,$$
$$\log B \lesssim \log(1/\varepsilon) + m \log(m/a_0) \lesssim \log(1/\varepsilon) + mN + m \log m.$$

The proof is finished.

$\square$

Next, we aim to approximate the clipping operation, which is essential for controlling the Sobolev norm at infinity of the approximator. The following lemma demonstrates the existence of a shallow GELU network for approximating clipping.

**Lemma 3.5** (approximation of clipping operation). *For every $A \geqslant 1$, every $\varepsilon \in (0,1)$, and every $m \in \mathbb{N}$, there exists $\varphi_{clip} \in \mathsf{NN}(L, W, S, B)$ such that*

$(i)$    $\|\varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-A,A])} \leqslant \varepsilon,$

$(ii)$    $\|\varphi_{clip} + A + 1/2\|_{W^{m,\infty}((-\infty,-A-1])} \vee \|\varphi_{clip} - A - 1/2\|_{W^{m,\infty}([A+1,+\infty))} \leqslant \varepsilon,$

$(iii)$    $\|\varphi_{clip}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m \log m + m \log \log(1/\varepsilon) + \log(2A))\},$

$(iv)$    $\|\varphi_{clip}\|_{W^{0,\infty}(\mathbb{R})} \leqslant A + 5/2,$

$(v)$    $\|\varphi_{clip} + A + 1/2\|_{W^{0,\infty}((-\infty,-A])} \vee \|\varphi_{clip} - A - 1/2\|_{W^{0,\infty}([A,+\infty))} \leqslant \varepsilon + 1,$

$(vi)$    $|\varphi_{clip}|_{W^{k,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(k \log m + k \log \log(1/\varepsilon))\}.$

*Moreover, $\varphi_{clip}$ has $L = \|W\|_\infty = 2$, $S = 7$ and $\log B \lesssim \log(Am/\varepsilon)$.*

*Proof.* Define

$$\varphi_{clip}(x) = \alpha^{-1}\mathrm{GELU}(\alpha(x + A + 1/2)) - \alpha^{-1}\mathrm{GELU}(\alpha(x - A - 1/2)) - A - 1/2, \quad x \in \mathbb{R}, \quad (6)$$

where $\alpha \geqslant 1$ will be determined later in the proof. Therefore, Lemma B.2 implies that

$$\|\varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-A,A])} \leqslant \alpha^m \|\mathrm{GELU} - \mathrm{id}\|_{W^{m,\infty}([\alpha/2,+\infty))} + \alpha^m \|\mathrm{GELU}\|_{W^{m,\infty}((-\infty,-\alpha/2])}$$
$$\leqslant 4\alpha^m \exp(-\alpha^2/16).$$

We next note that

$$\sup_{\alpha > 0} \alpha^m \exp(-\alpha^2/32) \leqslant \exp\{\mathcal{O}(m \log m)\}, \quad (7)$$

which implies that

$$\|\varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-A,A])} \leqslant \exp\{\mathcal{O}(m \log m)\} \exp\{-\alpha^2/32\}.$$

Therefore, setting

$$\alpha \asymp \sqrt{m \log m + \log(1/\varepsilon)} \quad (8)$$

9

guarantees that

$$\|\varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-A,A])} \leqslant \varepsilon. \tag{9}$$

We next focus on the behavior of tails of $\varphi_{clip}$

$$\|\varphi_{clip}\|_{W^{m,\infty}(\mathbb{R}\setminus(-A-1,A+1))} = \|\varphi_{clip}\|_{W^{m,\infty}((-\infty,-A-1])} \vee \|\varphi_{clip}\|_{W^{m,\infty}([A+1,+\infty))}.$$

We also note that

$$\|\varphi_{clip} + A + 1/2\|_{W^{m,\infty}((-\infty,-A-1])} \leqslant 2\alpha^m \|\mathrm{GELU}\|_{W^{m,\infty}((-\infty,-\alpha/2])}$$

and, similarly,

$$\|\varphi_{clip} - A - 1/2\|_{W^{m,\infty}([A+1,+\infty))} \leqslant 2\alpha^m \|\mathrm{GELU} - \mathrm{id}\|_{W^{m,\infty}([\alpha/2,+\infty))}.$$

From Lemma B.2 we find that

$$\|\varphi_{clip} + A + 1/2\|_{W^{m,\infty}((-\infty,-A-1])} \vee \|\varphi_{clip} - A - 1/2\|_{W^{m,\infty}([A+1,+\infty))} \leqslant 4\alpha^m \exp(-\alpha^2/16) \leqslant \varepsilon,$$

where the last inequality uses (7) and (8). Therefore, due to the triangle inequality we have that

$$\|\varphi_{clip}\|_{W^{m,\infty}(\mathbb{R}\setminus(-A-1,A+1))} \leqslant \varepsilon + A + 1/2. \tag{10}$$

We next derive the Sobolev norm of $\varphi_{clip}$ on the real line. Using (10), we have that

$$\|\varphi_{clip}\|_{W^{m,\infty}(\mathbb{R})} = \|\varphi_{clip}\|_{W^{m,\infty}(\mathbb{R}\setminus(-A-1,A+1))} \vee \|\varphi_{clip}\|_{W^{m,\infty}([-A-1,A+1])}$$
$$\leqslant (\varepsilon + A + 1/2) \vee \|\varphi_{clip}\|_{W^{m,\infty}([-A-1,A+1])}.$$

Lemma B.2 together with (8) implies that

$$\|\varphi_{clip}\|_{W^{m,\infty}([-A-1,A+1])}$$
$$\leqslant 2(\alpha^m \|\partial^1 \mathrm{GELU}\|_{W^{m-1,\infty}(\mathbb{R})} \vee \alpha^{-1} \|\mathrm{GELU}\|_{W^{0,\infty}([-2A-3/2,2A+3/2])}) + A + 1/2$$
$$\leqslant \exp\{\mathcal{O}(m \log m + m \log\log(1/\varepsilon) + \log(2A))\}.$$

Therefore,

$$\|\varphi_{clip}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m \log m + m \log\log(1/\varepsilon) + \log(2A))\}.$$

Similarly, for any $1 \leqslant k \leqslant m$, Lemma B.2 in conjunction with (6) and (8) yields

$$|\varphi_{clip}|_{W^{k,\infty}(\mathbb{R})} \leqslant 2\alpha^{k-1} |\mathrm{GELU}|_{W^{k,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(k \log m + k \log\log(1/\varepsilon))\}.$$

In addition, from Lemma B.2 we find that

$$\|\varphi_{clip} - A - 1/2\|_{W^{0,\infty}([A,A+1])}$$
$$\leqslant \|\mathrm{GELU} - \mathrm{id}\|_{W^{0,\infty}([\alpha(2A+1/2),+\infty))} + \|\mathrm{id} - A - 1/2\|_{W^{0,\infty}([A,A+1])} + \alpha^{-1} \|\mathrm{GELU}\|_{W^{0,\infty}([-\alpha/2,\alpha/2])}$$
$$\leqslant \varepsilon + 1.$$

Similarly,

$$\|\varphi_{clip} + A + 1/2\|_{W^{0,\infty}([-A-1,-A])} \leqslant \varepsilon + 1.$$

Therefore, from (9) and (10) we deduce that

$$\|\varphi_{clip}\|_{W^{0,\infty}(\mathbb{R})} \leqslant \varepsilon + A + 3/2 \leqslant A + 5/2.$$

Finally, we specify the configuration of $\varphi_{clip}$. The choice of $\alpha$ from (8) in conjunction with (6) suggest that

$$L = \|W\|_\infty = 2, \quad S = 7, \quad \log B \lesssim \log(\alpha \vee \alpha^{-1} \vee A) \lesssim \log(Am/\varepsilon).$$

The proof is complete.

$\square$

## 3.2 Approximation of monomials

Now, we focus on approximating polynomials. A key starting point is the approximation of the square operation, as demonstrated in the following lemma.

**Lemma 3.6** (approximation of square operation). *Define $f_{sq} : x \mapsto x^2$. Then, for every $\varepsilon \in (0,1)$ and every $m \in \mathbb{N}$, there exists $\varphi_{sq} \in \mathsf{NN}(L, W, S, B)$ such that*

$$\|\varphi_{sq} - f_{sq}\|_{W^{m,\infty}([-C,C])} \leqslant C^3 \varepsilon, \quad \text{for all } C \geqslant 1.$$

*Furthermore, $L = \|W\|_\infty = 2$, $S = 6$ and $\log B \lesssim \log(1/\varepsilon) + \log m$.*

*Proof.* Inspired by [Scarselli and Tsoi, 1998, Theorem 2], we let

$$\varphi_{sq}(x) = \frac{R^2}{\partial^2 \mathrm{GELU}(0)} \left( \mathrm{GELU}\left(\frac{2x}{R}\right) - 2 \cdot \mathrm{GELU}\left(\frac{x}{R}\right) \right), \tag{11}$$

where $R > 0$. We also highlight that $\partial^2 \mathrm{GELU}(0) = \sqrt{2/\pi}$. Using Taylor expansion it can be shown for any $x \in [-C, C]$ that

$$\left|\varphi_{sq}(x) - x^2\right| = \frac{|x|^3 \cdot |4\partial^3 \mathrm{GELU}(\xi) - \partial^3 \mathrm{GELU}(\zeta)|}{3R \cdot \partial^2 \mathrm{GELU}(0)} \leqslant \frac{5C^3 |\mathrm{GELU}|_{W^{3,\infty}(\mathbb{R})}}{3R \cdot \partial^2 \mathrm{GELU}(0)},$$

and similarly for the derivatives

$$\left|\partial^1 \varphi_{sq}(x) - 2x\right| = \frac{2|x|^2 \cdot |2\partial^3 \mathrm{GELU}(\widetilde{\xi}) - \partial^3 \mathrm{GELU}(\widetilde{\zeta})|}{R \cdot |\partial^2 \mathrm{GELU}(0)|} \leqslant \frac{6C^2 |\mathrm{GELU}|_{W^{3,\infty}(\mathbb{R})}}{R \cdot \partial^2 \mathrm{GELU}(0)}$$

and the second derivatives

$$\left|\partial^2 \varphi_{sq}(x) - 2\right| = \frac{|x| \cdot |8\partial^3 \mathrm{GELU}(\xi^\circ) - 2\partial^3 \mathrm{GELU}(\zeta^\circ)|}{R \cdot \partial^2 \mathrm{GELU}(0)} \leqslant \frac{10C |\mathrm{GELU}|_{W^{3,\infty}(\mathbb{R})}}{R \cdot \partial^2 \mathrm{GELU}(0)},$$

11

where $\xi, \widetilde{\xi}, \xi^\circ$ and $\zeta, \widetilde{\zeta}, \zeta^\circ$ are all within the interval defined by the origin and $x$. To proceed, we derive an explicit bound for the derivatives of order $k$ with $k \geqslant 3$ and $x \in \mathbb{R}$ as

$$|\partial^k \varphi_{sq}(x)| = \frac{1}{R^{k-2}\partial^2 \mathrm{GELU}(0)} \left| 2^k \partial^k \mathrm{GELU}\left(\frac{2x}{R}\right) - 2 \cdot \partial^k \mathrm{GELU}\left(\frac{x}{R}\right) \right| \leqslant \frac{(2^k+2)|\mathrm{GELU}|_{W^{k,\infty}(\mathbb{R})}}{R^{k-2}\partial^2 \mathrm{GELU}(0)}.$$

Therefore, from Lemma B.2 we find that

$$|\partial^k \varphi_{sq}(x)| \leqslant \frac{(k+1)(2^k+2)}{R^{k-2}\partial^2 \mathrm{GELU}(0)} \sqrt{\frac{(k-2)!}{2\pi}} \leqslant \frac{2^{k+2}k}{R^{k-2}\partial^2 \mathrm{GELU}(0)} \sqrt{\frac{(k-2)!}{2\pi}}.$$

Hence, setting

$$R = \frac{10|\mathrm{GELU}|_{W^{3,\infty}(\mathbb{R})}}{\partial^2 \mathrm{GELU}(0) \cdot \varepsilon} \vee \max_{3 \leqslant k \leqslant m} \left( \frac{2^{k+2}k}{\partial^2 \mathrm{GELU}(0)} \sqrt{\frac{(k-2)!}{2\pi}} \right)^{1/(k-2)},$$

we ensure that for any $C \geqslant 1$

$$\max_{0 \leqslant k \leqslant 2} |\varphi_{sq} - f_{sq}|_{W^{k,\infty}([-C,C])} \leqslant C^{3-k}\varepsilon, \qquad \max_{3 \leqslant k \leqslant (m \vee 3)} |\varphi_{sq} - f_{sq}|_{W^{k,\infty}(\mathbb{R})} \leqslant \varepsilon.$$

This observation yields

$$\|\varphi_{sq} - f_{sq}\|_{W^{m,\infty}([-C,C])} \leqslant C^3\varepsilon, \quad \text{for all } C \geqslant 1.$$

The definition of $\varphi_{sq}$ given in (11) suggests that $L = \|W\|_\infty = 2$, $S = 6$ and

$$\log B \lesssim \log(R^2 \vee R^{-1} \vee 2) \lesssim \log(1/\varepsilon) + \log m,$$

where the last inequality uses Stirling's approximation. The proof is finished.

$\square$

Comparing our result from Lemma 3.6 to that presented in [Gühring and Raslan, 2021, Proposition 4.7], we observe that we provide approximation guarantees for the entire real line, rather than limiting our results to a specific segment. This is a key advantage for approximating functions on unbounded domains. Subsequently, we derive a straightforward corollary that provides an approximation error bound for the multiplication of two numbers.

**Corollary 3.7** (approximation of two number multiplication). *Define* $\mathrm{prod}_2 : (x, y) \mapsto x \cdot y$, *and let* $m \in \mathbb{N}$ *be arbitrary. Then, for any* $\varepsilon \in (0, 1)$, *there exists* $\varphi_{mul} \in \mathsf{NN}(L, W, S, B)$ *satisfying*

$$\|\varphi_{mul} - \mathrm{prod}_2\|_{W^{m,\infty}([-C,C]^2)} \leqslant C^3\varepsilon, \quad \text{for all } C \geqslant 1.$$

*In addition,* $L = 2$, $\|W\|_\infty \leqslant 4$, $S \leqslant 12$ *and* $\log B \lesssim \log(1/\varepsilon) + \log m$.

The proof of Corollary 3.7 can be found in Appendix A.2. Having derived the approximation guarantees for multiplication, we now turn to the approximation of multiple number multiplications, as outlined in the following lemma.

12

**Lemma 3.8** (approximating the multiplication of $d$ numbers). *Let $d, m \in \mathbb{N}$ with $d \geqslant 2$ be arbitrary, and define the function $\mathrm{prod}_d : (x_1, \ldots, x_d) \mapsto \prod_{i=1}^{d} x_i$. Then, for every $\varepsilon \in (0,1)$ and every $K \geqslant 1$, there exists $\varphi_{mul,d} \in \mathsf{NN}(L, W, S, B)$ such that*

$$(i) \quad \|\varphi_{mul,d} - \mathrm{prod}_d\|_{W^{m,\infty}([-K,K]^d)} \leqslant \varepsilon,$$

$$(ii) \quad \|\varphi_{mul,d}\|_{W^{m,\infty}(\mathbb{R}^d)} \leqslant \exp\{\mathcal{O}((m^2 + d)\log(mdK\log(1/\varepsilon)))\}.$$

*In addition,*

$$L \lesssim \log d, \quad \|W\|_\infty \vee S \lesssim d^2, \quad \log B \lesssim (\log(1/\varepsilon) + (d + m)\log K + m^2 d^2)\log d.$$

*Proof.* To improve readability, the proof is divided into several steps.

**Step 1: approximation error analysis.** We first prove the statement for $K = 1$ and then generalize it to any arbitrary $K \geqslant 1$. Overall, the resulting neural network is structured as a binary tree, in accordance with the methodology described in Schwab and Zech [2019]. We build an approximation of multiplication of $2^J$ numbers with $J = \lceil \log_2 d \rceil$. If $d < 2^J$, then a minor modification of the input layer implements a concatenation of the input vector with the vector of ones of length at most $d$. Now let

$$\varphi_j(x_{1:2^j}) = \varphi_{mul,j}\left(\varphi_{j-1,1}(x_{1:2^{j-1}}), \varphi_{j-1,2}(x_{2^{j-1}+1:2^j})\right), \quad 1 \leqslant j \leqslant J, \tag{12}$$

where $\varphi_{mul,j}$ is the neural network from Corollary 3.7 with accuracy parameter $\varepsilon_{mul}^{(j)}$ and the smoothness parameter $m$, $\varphi_{j-1,1}$ and $\varphi_{j-1,2}$ are identical copies of $\varphi_{j-1}$, and $\varphi_0$ represents the identity mapping. From Corollary 3.7 we deduce that

$$\|\varphi_{mul,j} - \mathrm{prod}_2\|_{W^{m,\infty}([-C,C]^2)} \leqslant C^3 \varepsilon_{mul}^{(j)} \quad \text{for all } C \geqslant 1 \text{ and } 1 \leqslant j \leqslant J. \tag{13}$$

To simplify the notation, we let $\varphi_j(x_{1:2^j}) = (\varphi_{mul,j} \circ (\varphi_{j-1,1}, \varphi_{j-1,2}))(x_{1:2^j})$. Now assume that for all $0 \leqslant j \leqslant J$ and $C \geqslant 1$, it holds that

$$\|\varphi_j - \mathrm{prod}_{2^j}\|_{W^{m,\infty}(\Omega_j)} = \varepsilon_j, \tag{14}$$

where $\Omega_j = [-1,1]^{2^j}$. Hence, for any $1 \leqslant j \leqslant J$, the triangle inequality suggests that

$$\|\varphi_j - \mathrm{prod}_{2^j}\|_{W^{m,\infty}(\Omega_j)}$$
$$= \|\varphi_{mul,j} \circ (\varphi_{j-1,1}, \varphi_{j-1,2}) - \mathrm{prod}_{2^j}\|_{W^{m,\infty}(\Omega_j)}$$
$$\leqslant \|(\varphi_{mul,j} - \mathrm{prod}_2) \circ (\varphi_{j-1,1}, \varphi_{j-1,2})\|_{W^{m,\infty}(\Omega_j)} + \|\varphi_{j-1,1} \cdot \varphi_{j-1,2} - \mathrm{prod}_{2^j}\|_{W^{m,\infty}(\Omega_j)}. \tag{15}$$

As for the first term of (15), we apply Lemma B.4 and arrive at

$$\|(\varphi_{mul,j} - \mathrm{prod}_2) \circ (\varphi_{j-1,1}, \varphi_{j-1,2})\|_{W^{m,\infty}(\Omega_j)}$$
$$\leqslant 16(e^2 m^4 \cdot 2 \cdot 4^{j-1})^m \|\varphi_{mul,j} - \mathrm{prod}_2\|_{W^{m,\infty}([-1-\varepsilon_{j-1}, 1+\varepsilon_{j-1}]^2)} (1 \vee \|\varphi_{j-1,1}\|_{W^{m,\infty}(\Omega_{j-1})}^m)$$
$$\leqslant 16(e^2 m^4 \cdot 2 \cdot 4^{j-1})^m (1 + \varepsilon_{j-1})^{m+3} \varepsilon_{mul}^{(j)}, \tag{16}$$

where the last inequality uses (13). As for the second term of (15), we apply Lemma B.3 and obtain that

$$\|\varphi_{j-1,1} \cdot \varphi_{j-1,2} - \mathrm{prod}_{2^j}\|_{W^{m,\infty}(\Omega_j)} \leqslant 2^m \|\varphi_{j-1,1} - \mathrm{prod}_{2^{j-1}}\|_{W^{m,\infty}(\Omega_{j-1})} \|\varphi_{j-1,2}\|_{W^{m,\infty}(\Omega_{j-1})}$$
$$+ 2^m \|\mathrm{prod}_{2^{j-1}}\|_{W^{m,\infty}(\Omega_{j-1})} \|\varphi_{j-1,2} - \mathrm{prod}_{2^{j-1}}\|_{W^{m,\infty}(\Omega_{j-1})}.$$

From (14) we deduce that

$$\|\varphi_{j-1,1} \cdot \varphi_{j-1,2} - \mathrm{prod}_{2^j}\|_{W^{m,\infty}(\Omega_j)} \leqslant 2^{m+1}\varepsilon_{j-1}(1+\varepsilon_{j-1}). \tag{17}$$

Therefore, combining (15), (16) and (17), we arrive at

$$\varepsilon_j \leqslant 16(e^2 m^4 \cdot 2 \cdot 4^{j-1})^m (1+\varepsilon_{j-1})^{m+3}\varepsilon_{mul}^{(j)} + 2^{m+1}\varepsilon_{j-1}(1+\varepsilon_{j-1}).$$

We find $\varepsilon_j$ in the form of $\varepsilon_j = 2^{\gamma_j}\varepsilon_1$ for each $1 \leqslant j \leqslant J$ with $\gamma_1 = 0$. Hence,

$$\varepsilon_j \leqslant 16(e^2 m^4 \cdot 2 \cdot 4^{j-1})^m 2^{(m+1)(\gamma_{j-1}+1)}\varepsilon_{mul}^{(j)} + 2^{m+1}2^{2\gamma_{j-1}+1}\varepsilon_1.$$

Setting

$$\varepsilon_{mul}^{(j)} = \varepsilon_1 \left(16(e^2 m^4 \cdot 2 \cdot 4^{j-1})^m 2^{(m+1)(\gamma_{j-1}+1)}\right)^{-1}, \quad 2 \leqslant j \leqslant J, \tag{18}$$

we have that

$$\gamma_j \leqslant m + 3 + 2\gamma_{j-1}, \quad 2 \leqslant j \leqslant J,$$

which yields that

$$\gamma_j \leqslant (m+3)4^j, \quad 1 \leqslant j \leqslant J.$$

Therefore,

$$\varepsilon_J \leqslant (m+3)2^{2\lceil \log_2 d \rceil}\varepsilon_1 \leqslant 4(m+3)d^2\varepsilon_1.$$

Choosing $\varepsilon_{mul}^{(1)} = \varepsilon(4(m+3)d^2)^{-1} \in (0,1)$ ensures that

$$\|\varphi_J - \mathrm{prod}_{2^J}\|_{W^{m,\infty}([-C,C])} \leqslant \varepsilon.$$

**Step 2: deriving the configuration of $\varphi_J$.** Due to the observation that $\gamma_j \lesssim md^2$, we deduce from Corollary 3.7 and (18) that, for all $1 \leqslant j \leqslant J$, we have $\varphi_{mul,j} \in \mathsf{NN}(L_{mul}, W_{mul}, S_{mul}, B_{mul})$ with

$$L_{mul} = 2, \quad \|W_{mul}\|_\infty \vee S_{mul} \lesssim 1, \quad \log B_{mul} \lesssim \log(1/\varepsilon) + m^2 d^2.$$

Let $\varphi_j \in \mathsf{NN}(L_j, W_j, S_j, B_j)$ for all $1 \leqslant j \leqslant J$. Then, from Lemma B.5 and Lemma B.6 we find that

$$L_J \leqslant J + 1, \quad \|W_J\|_\infty \vee S_J \lesssim 2^{J-1},$$
$$\log B_J \leqslant \log B_{J-1} + \log B_{mul} + \log \|W_{J-1}\|_\infty \lesssim (\log(1/\varepsilon) + m^2 d^2)\log d. \tag{19}$$

We now generalize the approximation result to the case when $K \geqslant 1$. Let

$$\varphi_{J,K}(x_1, \ldots, x_d) = K^d \varphi_J(x_1/K, \ldots, x_d/K).$$

Then, from the chain rule we obtain that

$$\|\varphi_{J,K} - \mathrm{prod}_d\|_{W^{m,\infty}([-K,K]^d)} \leqslant K^d \|\varphi_J - \mathrm{prod}_d\|_{W^{m,\infty}([-1,1]^d)} \leqslant K^d \varepsilon.$$

14

Therefore, taking the accuracy parameter $\varepsilon/K^d$ in $\varphi_{J,K}$, we deduce that for any $\varepsilon \in (0,1)$ there exists $\widetilde{\varphi}_{mul,d,K} \in \mathsf{NN}(\widetilde{L}, \widetilde{W}, \widetilde{S}, \widetilde{B})$ satisfying

$$\|\widetilde{\varphi}_{mul,d,K} - \mathrm{prod}_d\|_{W^{m,\infty}([-K,K]^d)} \leqslant \varepsilon.$$

Furthermore, (19) we find that

$$\widetilde{L} \lesssim \log d, \quad \|\widetilde{W}\|_\infty \vee \widetilde{S} \lesssim d, \quad \log \widetilde{B} \lesssim (\log(1/\varepsilon) + d \log K + m^2 d^2) \log d. \tag{20}$$

**Step 3: clipping the input.** Now let $\varphi_{clip}$ be a clipping operation approximation from Lemma 3.5 formulated with accuracy parameter $\varepsilon_{clip} \in (0,1)$ and clipping parameter $K$. Then, it holds that $\|\varphi_{clip}\|_{W^{0,\infty}(\mathbb{R})} \leqslant K + 5/2 \leqslant 4K$. Let $\varphi_{clip,d}$ be a parallel stacking of $d$ identical copies of $\varphi_{clip}$ that approximates a component-wise clipping. Let also $\widetilde{\varphi}_{mul,d,4K}$ has accuracy parameter $\varepsilon_{mul,d}$ and smoothness parameter $m + 1$. Then, it holds that

$$\|\widetilde{\varphi}_{mul,d,4K} \circ \varphi_{clip,d} - \mathrm{prod}_d \circ \mathrm{id}\|_{W^{m,\infty}([-K,K]^d)}$$
$$\leqslant \|(\widetilde{\varphi}_{mul,d,4K} - \mathrm{prod}_d) \circ \varphi_{clip,d}\|_{W^{m,\infty}([-K,K]^d)} + \|\mathrm{prod}_d \circ \varphi_{clip,d} - \mathrm{prod}_d \circ \mathrm{id}\|_{W^{m,\infty}([-K,K]^d)}.$$

Lemma B.4 suggests that

$$\|(\widetilde{\varphi}_{mul,d,4K} - \mathrm{prod}_d) \circ \varphi_{clip,d}\|_{W^{m,\infty}([-K,K]^d)} \leqslant \exp\{\mathcal{O}(m\log(md))\}\varepsilon_{mul,d}(\varepsilon_{clip} + K)^m$$
$$\leqslant \exp\{\mathcal{O}(m\log(mdK))\}\varepsilon_{mul,d}$$

and also

$$\|\mathrm{prod}_d \circ \varphi_{clip,d} - \mathrm{prod}_d \circ \mathrm{id}\|_{W^{m,\infty}([-K,K]^d)} \leqslant \exp\{\mathcal{O}(m\log(md))\}(K + 5/2)^d \varepsilon_{clip}(\varepsilon_{clip} + K)^{2m}$$
$$\leqslant \exp\{\mathcal{O}(m\log(mdK) + d\log K)\}\varepsilon_{clip}.$$

Therefore, setting

$$\log(1/\varepsilon_{mul,d}) \asymp \log(1/\varepsilon) + m\log(mdK), \quad \log(1/\varepsilon_{clip}) \asymp \log(1/\varepsilon) + m\log(mdK) + d\log K \tag{21}$$

for some $\varepsilon \in (0,1)$ ensures that

$$\|\widetilde{\varphi}_{mul,d,4K} \circ \varphi_{clip,d} - \mathrm{prod}_d \circ \mathrm{id}\|_{W^{m,\infty}([-K,K]^d)} \leqslant \varepsilon.$$

Moreover, Lemma 3.5 and Lemma B.4 imply that

$$\|\widetilde{\varphi}_{mul,d,4K} \circ \varphi_{clip,d}\|_{W^{m,\infty}(\mathbb{R}^d)} \leqslant \exp\{\mathcal{O}(m\log(md))\}(\varepsilon_{mul,d} + (4K)^d)(1 \vee \|\varphi_{clip,d}\|^m_{W^{m,\infty}(\mathbb{R}^d)})$$
$$\leqslant \exp\{\mathcal{O}((m^2 + d)\log(mdK\log(1/\varepsilon)))\},$$

where the last inequality uses (21). Recall that due to Lemma 3.5, Lemma B.6 and (21), we have that $\varphi_{clip,d} \in \mathsf{NN}(L_{clip}, W_{clip}, S_{clip}, B_{clip})$ with

$$L_{clip} \lesssim 1, \quad \|W_{clip}\|_\infty \vee S_{clip} \lesssim d, \quad \log B_{clip} \lesssim \log(1/\varepsilon) + m\log(mdK) + d\log K.$$

Finally, from Lemma B.5, (20) and (21) we deduce that $\varphi_{mul,d} = \widetilde{\varphi}_{mul,d,4K} \circ \varphi_{clip,d}$ has

$$L \lesssim \log d, \quad \|W\|_\infty \vee S \lesssim d^2, \quad \log B \lesssim (\log(1/\varepsilon) + (d+m)\log K + m^2 d^2)\log d.$$

The proof is complete.

$\square$

Comparing our result from Lemma 3.8 to that presented in [De Ryck et al., 2021, Corollary 3.8], we observe a difference in the number of parameters: $\mathcal{O}(d^2)$ versus $\mathcal{O}(d \log d)$. We emphasize that, as a byproduct, we derived a neural network with the number of parameters $\mathcal{O}(d)$, but we employed clipping and concatenation to satisfy condition $(ii)$, which ultimately increased the parameter count. However, by adding clipping, we ensure that the approximation and its derivatives are bounded across the entire real line. Now, we turn to the approximation of monomials, as formulated in the following lemma.

**Corollary 3.9** (approximation of monomials). *Let $\mathbf{k} \in \mathbb{Z}_+^I$ for some $I \in \mathbb{N}$ such that $|\mathbf{k}| = d$, where $d \in \mathbb{N}$ with $d \geqslant 2$ is arbitrary. Define $\mathrm{prod}_{\mathbf{k}} : (x_1, \ldots, x_I) \mapsto \prod_{i=1}^I x_i^{k_i}$. Then, for every $\varepsilon \in (0,1)$, every $m \in \mathbb{N}$, and every $K \geqslant 1$, there exists a GELU network $\varphi_{mul,\mathbf{k}} \in \mathsf{NN}(L, W, S, B)$ such that*

$$
\begin{aligned}
(i) \quad & \|\varphi_{mul,\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \varepsilon, \\
(ii) \quad & \|\varphi_{mul,\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant \exp\{\mathcal{O}((m^2 + d) \log(mdK \log(1/\varepsilon)))\}.
\end{aligned}
$$

*In addition, $\varphi_{mul,\mathbf{k}}$ has*

$$
L \lesssim \log d, \quad \|W\|_\infty \vee S \lesssim (d \vee I)^3, \quad \log B \lesssim (\log(1/\varepsilon) + (d+m)\log K + m^2 d^2)\log d + \log I.
$$

We move the proof of Corollary 3.9 to Appendix A.3. The following lemma provides an approximation result for multivariate polynomials.

**Lemma 3.10** (approximation of multivariate polynomials). *Define $f_{\mathcal{A}} : x \mapsto \sum_{\mathbf{k} \in \mathcal{A}} a_{\mathbf{k}} x^{\mathbf{k}}$, where $x \in \mathbb{R}^I$ and $\mathcal{A} = \{\mathbf{k} \in \mathbb{Z}_+^I : |\mathbf{k}| \leqslant d\}$ for some $I, d \in \mathbb{N}$ with $d \geqslant 2$. Also assume that $|a_{\mathbf{k}}| \leqslant 1$ for all $\mathbf{k} \in \mathcal{A}$. Then, for every $\varepsilon \in (0,1)$, every natural $m \geqslant 3$ and every $K \geqslant 1$, there exists a neural network $\varphi_{\mathcal{A}} \in \mathsf{NN}(L, W, S, B)$ such that*

$$
\begin{aligned}
(i) \quad & \|f_{\mathcal{A}} - \varphi_{\mathcal{A}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \varepsilon, \\
(ii) \quad & \|\varphi_{\mathcal{A}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant \exp\{\mathcal{O}((m^2 + md + I) \log(mdKI \log(1/\varepsilon)))\}.
\end{aligned}
$$

*In addition, $\varphi_{\mathcal{A}}$ has*

$$
\begin{aligned}
& L \lesssim \log d, \quad \|W\|_\infty \vee S \lesssim (d+I)^{3+d \wedge I} \\
& \log B \lesssim (\log(1/\varepsilon) + m^2(d+I)\log(mdKI) + m^2 d^2)\log(d+I).
\end{aligned}
$$

*Proof.* Corollary 3.9 implies that for each $\mathbf{k} \in \mathcal{A}$ with $|\mathbf{k}| \geqslant 2$ there exists $\varphi_{\mathbf{k}}$ satisfying

$$
\|\varphi_{\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \varepsilon_{\mathbf{k}}, \tag{22}
$$

where $\varepsilon_{\mathbf{k}} \in (0,1)$ is accuracy parameter. Moreover,

$$
\|\varphi_{\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant \exp\{\mathcal{O}((m^2 + d) \log(mdK \log(1/\varepsilon_{\mathbf{k}})))\} \tag{23}
$$

and $\varphi_{\mathbf{k}} \in \mathsf{NN}(L_{\mathbf{k}}, W_{\mathbf{k}}, S_{\mathbf{k}}, B_{\mathbf{k}})$ with

$$
\begin{aligned}
& L_{\mathbf{k}} \lesssim \log d, \quad \|W_{\mathbf{k}}\|_\infty \vee S_{\mathbf{k}} \lesssim (d \vee I)^3, \\
& \log B_{\mathbf{k}} \lesssim (\log(1/\varepsilon_{\mathbf{k}}) + (d+m)\log K + m^2 d^2)\log d + \log I. \tag{24}
\end{aligned}
$$

As for $|\mathbf{k}| \in \{0, 1\}$, the approximation is exact, since it is implemented with a single linear layer. In order to build the final approximation, we have to implement a summation of GELU networks with different depth.

For this purpose, we add auxiliary identity layers. Let $\varphi_{id,\mathbf{k}}$ be an approximation of identity function from Lemma 3.2 formulated with the accuracy parameter $\varepsilon_{\mathbf{k}}$, $L_{id,\mathbf{k}} = 1 + \max_{\widetilde{\mathbf{k}} \in \mathcal{A}} L_{\widetilde{\mathbf{k}}} - L_{\mathbf{k}}$ number of layers and the scale parameter $\|\varphi_{\mathbf{k}}\|_{W^{0,\infty}([-K,K]^I)}$. Hence, the triangle inequality implies that

$$\|\varphi_{id,\mathbf{k}} \circ \varphi_{\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)}$$
$$\leqslant \|\varphi_{\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} + \|(\varphi_{id,\mathbf{k}} - \mathrm{id}) \circ \varphi_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)}. \tag{25}$$

Therefore, Lemma B.4 suggest that

$$\|(\varphi_{id,\mathbf{k}} - \mathrm{id}) \circ \varphi_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)}$$
$$\leqslant \exp\{\mathcal{O}(m \log m)\}\|\varphi_{id,\mathbf{k}} - \mathrm{id}\|_{W^{m,\infty}(\Omega_2)}(\|\varphi_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)}^m \vee 1),$$

where $\Omega_2 = [-\|\varphi_{\mathbf{k}}\|_{W^{0,\infty}([-K,K]^I)}, \|\varphi_{\mathbf{k}}\|_{W^{0,\infty}([-K,K]^I)}]$. As suggested by (22), we have that

$$\|\varphi_{\mathbf{k}}\|_{W^{0,\infty}([-K,K]^I)} \leqslant \|\varphi_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \|\mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} + \varepsilon_{\mathbf{k}} \leqslant 2d^m K^d. \tag{26}$$

Hence, it holds that

$$\|(\varphi_{id,\mathbf{k}} - \mathrm{id}) \circ \varphi_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \exp\{\mathcal{O}(m^2 d \log(mdK))\}\varepsilon_{\mathbf{k}}.$$

Then, from (22) and (25) we deduce that

$$\|\varphi_{id,\mathbf{k}} \circ \varphi_{\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \exp\{\mathcal{O}(m^2 d \log(mdK))\}\varepsilon_{\mathbf{k}}, \quad \text{for all } \mathbf{k} \in \mathcal{A}. \tag{27}$$

We now specify the configuration of each $\varphi_{id,\mathbf{k}}$, as suggested by (26) and Lemma 3.2:

$$L(\varphi_{id,\mathbf{k}}) \lesssim \log d, \quad \|W\|_\infty(\varphi_{id,\mathbf{k}}) \lesssim 1,$$
$$S(\varphi_{id,\mathbf{k}}) \lesssim \log d, \quad \log B(\varphi_{id,\mathbf{k}}) \lesssim (m + \log d) \log m + \log(1/\varepsilon_{\mathbf{k}}) + m^2 \log d + dm \log K.$$

Now (24) and Lemma B.5 imply that the composition $\varphi_{\mathbf{k}} \circ \varphi_{id,\mathbf{k}}$ has

$$L(\varphi_{\mathbf{k}} \circ \varphi_{id,\mathbf{k}}) \lesssim \log d, \quad \|W\|_\infty(\varphi_{\mathbf{k}} \circ \varphi_{id,\mathbf{k}}) \vee S(\varphi_{\mathbf{k}} \circ \varphi_{id,\mathbf{k}}) \lesssim (d \vee I)^3,$$
$$\log B(\varphi_{\mathbf{k}} \circ \varphi_{id,\mathbf{k}}) \lesssim (\log(1/\varepsilon_{\mathbf{k}}) + dm \log K + m^2 d^2) \log d + \log I.$$

Now setting $\varepsilon_{\mathbf{k}} = \varepsilon/|\mathcal{A}|$ and applying parallelization argument from Lemma B.6, for

$$\varphi_{\mathcal{A}}(x) = \sum_{\mathbf{k} \in \mathcal{A}} a_{\mathbf{k}} \cdot (\varphi_{id,\mathbf{k}} \circ \varphi_{\mathbf{k}})(x), \quad x \in \mathbb{R}^I,$$

we obtain that $\varphi_{\mathcal{A}}$ has

$$L(\varphi_{\mathcal{A}}) \lesssim \log d, \quad \|W\|_\infty(\varphi_{\mathcal{A}}) \vee S(\varphi_{\mathcal{A}}) \lesssim |\mathcal{A}|(d \vee I)^3$$
$$\log B(\varphi_{\mathcal{A}}) \lesssim \log(|\mathcal{A}|) + (\log(1/\varepsilon_{\mathbf{k}}) + dm \log K + m^2 d^2) \log d + \log I. \tag{28}$$

Since

$$|\mathcal{A}| \leqslant \binom{d+I}{d} \leqslant (d+I)^{d \wedge I} = \exp\{(d \wedge I) \log(d+I)\}, \tag{29}$$

17

then (27) yields

$$\|f_{\mathcal{A}} - \varphi_{\mathcal{A}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \exp\{\mathcal{O}(m^2(d+I)\log(mdKI))\}\varepsilon_{\mathbf{k}}.$$

Thus, setting

$$\log(1/\varepsilon_{\mathbf{k}}) = \log(1/\varepsilon) + m^2(d+I)\log(mdKI) \tag{30}$$

ensures that

$$\|f_{\mathcal{A}} - \varphi_{\mathcal{A}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \varepsilon.$$

Moreover, the configuration described in (28) is now

$$L(\varphi_{\mathcal{A}}) \lesssim \log d, \quad \|W\|_{\infty}(\varphi_{\mathcal{A}}) \vee S(\varphi_{\mathcal{A}}) \lesssim (d+I)^{3+d\wedge I}$$
$$\log B(\varphi_{\mathcal{A}}) \lesssim (\log(1/\varepsilon) + m^2(d+I)\log(mdKI) + m^2d^2)\log(d+I).$$

From (23), (29) and Lemma B.4 we deduce that

$$\|\varphi_{\mathcal{A}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant |\mathcal{A}| \max_{\mathbf{k}\in\mathcal{A}} \|\varphi_{id,\mathbf{k}} \circ \varphi_{\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)}$$
$$\leqslant \exp\{\mathcal{O}((d\wedge I)\log(d+I) + m\log m)\} \max_{\mathbf{k}\in\mathcal{A}} \|\varphi_{id,\mathbf{k}}\|_{W^{m,\infty}(\varphi_{\mathbf{k}}(\mathbb{R}^I))}(1 \vee \|\varphi_{\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)}^m).$$

As suggested by (23) and (30), we have that

$$\|\varphi_{\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant \exp\{\mathcal{O}(m^2+d)\log(mdKI\log(1/\varepsilon))\}.$$

From Lemma 3.2 and (26) we find that

$$\|\varphi_{id,\mathbf{k}}\|_{W^{m,\infty}(\varphi_{\mathbf{k}}(\mathbb{R}^I))} \leqslant \|\varphi_{id,\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m\log(dm\log(1/\varepsilon_{\mathbf{k}})) + d\log K)\}.$$

Therefore, due to (30), it holds that

$$\|\varphi_{\mathcal{A}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant \exp\{\mathcal{O}((m^2+md+I)\log(mdKI\log(1/\varepsilon)))\}.$$

The proof is complete.

$\square$

## 3.3 Approximation of the exponent and the division

Now, we address the approximation of nonlinear operations, including exponentiation and division. The following lemma provides quantitative bounds for the exponential function approximation.

**Lemma 3.11** (approximation of the exponential function). *Define $f_{exp} : x \mapsto e^{-x}$, and let $m \in \mathbb{N}$ be arbitrary. Then, for any $\varepsilon \in (0,1)$ and $0 \leqslant A \leqslant 1$, there exists a neural network $\varphi_{exp} \in \mathsf{NN}(L,W,S,B)$ such that*

$$(i) \quad \|\varphi_{exp} - f_{exp}\|_{W^{m,\infty}([-A,+\infty))} \leqslant \varepsilon,$$
$$(ii) \quad \|\varphi_{exp}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m^2\log(m\log(1/\varepsilon)))\}.$$

*Furthermore,*

$$L \lesssim \log m + \log\log(1/\varepsilon), \quad \|W\|_{\infty} \vee S \lesssim m^{12}\log^4(1/\varepsilon), \quad \log B \lesssim m^{11}\log^3(1/\varepsilon).$$

18

*Proof.* For some $r \in \mathbb{N}$ with $r \geqslant 3$ and $K \geqslant 2$, which will be determined later, consider the approximation accuracy of a Tailor expansion $f_r(x) = \sum_{i=0}^{r-1} \frac{(-1)^i x^i}{i!}$

$$\|f_{exp} - f_r\|_{W^{m,\infty}([-4A,4K])} \leqslant \max_{0 \leqslant m' \leqslant m} \left( \frac{e^{4A}(4K)^{r-m'}}{(r-m')!} \right) \leqslant e^{4A} \max_{0 \leqslant m' \leqslant m} \left( \frac{4e(K \vee A)}{r-m'} \right)^{r-m'},$$

where in the last inequality we used Stirling's approximation for the factorial. Thus, setting

$$r = \lceil m + 4Ke^2 + 4A + \log(2/\varepsilon_0) \rceil \geqslant 3,$$

where $\varepsilon_0 \in (0,1)$ and will be optimized further in the proof. Next, we obtain that

$$\|f_{exp} - f_r\|_{W^{m,\infty}([-4A,4K])} \leqslant \varepsilon_0/2. \tag{31}$$

Now Lemma 3.10 suggests that there exists a GELU network $\widetilde{\varphi}_{exp}$ formulated with the accuracy parameter $\varepsilon_0/2$, the scaling parameter $4K$, the smoothness parameter $m+1$, and the maximum power of the monomial $r-1$ such that

$$\|f_r - \widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-4A,4K])} \leqslant \|f_r - \widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-4K,4K])} \leqslant \varepsilon_0/2,$$

which together with (31) immediately implies that

$$\|f_{exp} - \widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-4A,4K])} \leqslant \|f_{exp} - f_r\|_{W^{m,\infty}([-4A,4K])} + \|f_r - \widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-4A,4K])} \leqslant \varepsilon_0. \tag{32}$$

In addition, (32) suggests that

$$\|\widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-4A,4K])} \leqslant \|\widetilde{\varphi}_{exp} - f_{exp}\|_{W^{m,\infty}([-4A,4K])} + \|f_{exp}\|_{W^{m,\infty}([-4A,4K])} \leqslant 2e^{4A}. \tag{33}$$

Next, substituting the choice of $r$ into the configuration of $\widetilde{\varphi}_{exp} \in \mathsf{NN}(\widetilde{L}, \widetilde{W}, \widetilde{S}, \widetilde{B})$ outlined in Lemma 3.10 yields

$$\widetilde{L} \lesssim \log r \lesssim \log(m + K + \log(1/\varepsilon_0)), \quad \|\widetilde{W}\|_\infty \vee \widetilde{S} \lesssim r^4 \lesssim m^4 + K^4 + \log^4(1/\varepsilon_0),$$

$$\log \widetilde{B} \lesssim (\log(1/\varepsilon_0) + m^2 r \log(mrK) + m^2 r^2) \log r \lesssim m^2(m^3 + K^3 + \log^3(1/\varepsilon_0)). \tag{34}$$

Let $\varphi_{clip}$ be an approximation of clipping operation from Lemma 3.5 with the accuracy parameter $\varepsilon_{clip} \in (0,1)$ and the scale parameter $(A+K)/2 \geqslant 1$. Then we have that

$$\varphi_{clip,-A,K}(x) = \varphi_{clip}(x + (A-K)/2) + (K-A)/2$$

has the following properties:

$$(i) \quad \|\varphi_{clip,-A,K} - \mathrm{id}\|_{W^{m,\infty}([-A,K])} \leqslant \|\varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-(A+K)/2,(A+K)/2])} \leqslant \varepsilon_{clip},$$

$$(ii) \quad -5/2 - A \leqslant \varphi_{clip,-A,K}(x) \leqslant K + 5/2, \quad \text{for all } x \in \mathbb{R}$$

$$(iii) \quad \|\varphi_{clip,-A,K} - K - 1/2\|_{W^{0,\infty}([K,+\infty))} \leqslant \varepsilon_{clip} + 1,$$

$$(iv) \quad \|\varphi_{clip,-A,K}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m \log(m \log(1/\varepsilon_{clip})) + \log(A+K))\}.$$

Hence, from properties $(i)$, $(ii)$ and Lemma B.4 we obtain for $\varphi_{exp} = \widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K}$ that

$$\|\widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K} - \widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-A,K])} \leqslant \exp\{\mathcal{O}(m \log(mK))\} \|\widetilde{\varphi}_{exp}\|_{W^{m+1,\infty}([-4A,4K])} \varepsilon_{clip}.$$

From (33) we find that

$$\|\widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K} - \widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-A,K])} \leqslant \exp\{\mathcal{O}(m\log(mK))\}\varepsilon_{clip}.$$

Thus, (32) implies that for

$$\log(1/\varepsilon_{clip}) \asymp \log(1/\varepsilon_0) + m\log mK \tag{35}$$

we have

$$\|\varphi_{exp} - f_{exp}\|_{W^{m,\infty}([-A,K])} \leqslant 2\varepsilon_0. \tag{36}$$

Property $(iii)$ together with Lemma B.4 suggests that

$$\|\widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K}\|_{W^{m,\infty}([K,+\infty))}$$
$$\leqslant \exp\{\mathcal{O}(m\log m)\}\|\widetilde{\varphi}_{exp}\|_{W^{m,\infty}([K-3/2,K+5/2])}(1 \vee \|\varphi_{clip,-A,K}\|_{W^{m,\infty}(\mathbb{R})}^m).$$

From (32), (35) and property $(iv)$ we find that

$$\|\widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K}\|_{W^{m,\infty}([K,+\infty))} \leqslant \exp\{\mathcal{O}(m^2\log(mK\log(1/\varepsilon_0)))\}(\varepsilon_0 + e^{-K}).$$

Thus, for $K = 2 \vee \log(1/\varepsilon_0)$ we have that

$$\|\widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K}\|_{W^{m,\infty}([K,+\infty))} \leqslant \exp\{\mathcal{O}(m^2\log(m\log(1/\varepsilon_0)))\}\varepsilon_0.$$

This and the triangle inequality imply that

$$\|\widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K} - f_{exp}\|_{W^{m,\infty}([K,+\infty))} \leqslant \exp\{\mathcal{O}(m^2\log(m\log(1/\varepsilon_0)))\}\varepsilon_0.$$

Therefore, from (36) we deduce that

$$\|\varphi_{exp} - f_{exp}\|_{W^{m,\infty}([-A,+\infty))} \leqslant \exp\{\mathcal{O}(m^2\log(m\log(1/\varepsilon_0)))\}\varepsilon_0.$$

Hence, setting

$$\log(1/\varepsilon_0) \asymp m^2\log m + m^3\log(1/\varepsilon) \tag{37}$$

ensures that

$$\|\varphi_{exp} - f_{exp}\|_{W^{m,\infty}([-A,+\infty))} \leqslant \varepsilon.$$

Moreover, properties $(ii)$, $(iv)$ together with Lemma B.4, (33), (35) and (37) yields

$$\|\widetilde{\varphi}_{exp} \circ \varphi_{clip,-A,K}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \|\widetilde{\varphi}_{exp}\|_{W^{m,\infty}([-4A,4K])}\exp\{\mathcal{O}(m^2\log(m\log(1/\varepsilon_{clip})) + m\log(2K))\}$$
$$\leqslant \exp\{\mathcal{O}(m^2\log(m\log(1/\varepsilon)))\}.$$

From (35), (37) and Lemma 3.5 we find that the configuration of $\varphi_{clip} \in \mathsf{NN}(L_{clip}, W_{clip}, S_{clip}, B_{clip})$ is

$$L_{clip} \vee \|W_{clip}\|_\infty \vee S_{clip} \lesssim 1, \quad \log B_{clip} \lesssim \log(mK/\varepsilon_{clip}) \lesssim m^2\log m + m^3\log(1/\varepsilon).$$

Therefore, (34), (37) and Lemma B.5 suggest that the configuration of $\varphi_{exp}$ is

$$L \lesssim \log m + \log\log(1/\varepsilon), \quad \|W\|_\infty \vee S \lesssim m^{12}\log^4(1/\varepsilon), \quad \log B \lesssim m^{11}\log^3(1/\varepsilon).$$

This completes the proof.

$\square$

Comparing our result presented in Lemma 3.11 with [Yakovlev and Puchkin, 2025, Corollary F.3], we observe a less favorable configuration scaling. Specifically, the number of parameters in Lemma 3.11 scales as $\mathcal{O}(\log^4(1/\varepsilon))$, compared to $\mathcal{O}(\log^2(1/\varepsilon))$. Nevertheless, we extend the approximation guarantees to high-order Sobolev norms.

Now, we focus on the approximation of the division operation, beginning by approximating the reciprocal function in a straightforward manner, as suggested by the following lemma.

**Lemma 3.12** (naive approximation of the reciprocal function). *Let $0 < a \leqslant b \leqslant 2$ such that $b/a \geqslant 5/4$ and $a < 1$. Define $f_{rec} : x \mapsto 1/x$, where $x > 0$. Then, for every $\varepsilon \in (0,1)$ and every $m \in \mathbb{N}$ such that $m \geqslant 3$, there exists $\varphi_{rec} \in \mathsf{NN}(L, W, S, B)$ satisfying*

$$(i) \quad \|\varphi_{rec} - f_{rec}\|_{W^{m,\infty}([a,b])} \leqslant \varepsilon,$$

$$(ii) \quad \|\varphi_{rec}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m^2 \log(mN/a) + m^2 \log \log(1/\varepsilon) + mN)\}.$$

*Moreover, $\varphi_{rec}$ has*

$$L \lesssim \log((mb/a)\log(1/\varepsilon)), \quad \|W\|_\infty \vee S \lesssim (mb/a)^4 \log^4(m/\varepsilon a),$$
$$\log B \lesssim (m^4 b^2/a^2) \log^2(1/\varepsilon a) \log^2((mb/a)\log(1/\varepsilon a)).$$

The proof of Lemma 3.12 is deferred to Appendix A.4. Overall, the proof is similar to that of [Yakovlev and Puchkin, 2025, Lemma A.8], but extends it to Sobolev norms. The following result constructs a strong approximator by leveraging the weak approximators derived in Lemma 3.12, drawing inspiration from [Yakovlev and Puchkin, 2025, Lemma A.4].

**Lemma 3.13** (reciprocal function approximation). *Define $f_{rec} : x \mapsto 1/x$ for any $x > 0$. Let also $a_0 = 2^{-N}$ for some $N \in \mathbb{N}$ such that $N \geqslant 3$. Then, for every $\varepsilon \in (0,1)$ and every $m \in \mathbb{N}$ with $m \geqslant 3$, there exists a GELU network $\varphi_{rec} \in \mathsf{NN}(L, W, S, B)$ such that*

$$(i) \quad \|\varphi_{rec} - f_{rec}\|_{W^{m,\infty}([a_0,1])} \leqslant \varepsilon,$$

$$(ii) \quad \|\varphi_{rec}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m^3 N + m^3 \log(m \log(1/\varepsilon)))\}.$$

*In addition, the network has*

$$L \lesssim \log(mN \log(1/\varepsilon)), \quad \|W\|_\infty \vee S \lesssim m^8 N(N^4 + m^4 \log^4(1/\varepsilon)),$$
$$\log B \lesssim m^8(N^4 + m^4 \log^4(1/\varepsilon)).$$

*Proof.* The proof proceeds in multiple steps.

**Step 1: introducing basic approximators.** Let $N \in \mathbb{N}$ with $N \geqslant 3$ and let $a_i = 2^{-N+i}$ for each $i \in \{-1, 1, \ldots, N+1\}$. Let also $\varphi_{rec}$ be in the following form:

$$\varphi_{rec}(x) = \sum_{i=1}^{N} q(\varphi_i, \psi_i), \quad \varphi_i = \varphi_{id,i} \circ \varphi_{rec,i}, \; \psi_i = \psi_{id,i} \circ \psi_{pou,i}, \tag{38}$$

where $q$ is a GELU network from Corollary 3.7, which approximates multiplication with accuracy parameter $\varepsilon_{mul}$. Networks $\{\psi_{pou,i}\}_{i=1}^{N}$ form a partition of unity according to Lemma 3.4 with accuracy $\varepsilon_{pou}$. In

addition, the networks $\{\varphi_{rec,i}\}_{i=1}^N$ serve as local approximators of the reciprocal function from Lemma 3.12 with the accuracy $\varepsilon_{rec}/2$ and the parameters $a = a_{i-2}$ and $b = a_{i+1} \wedge 1$. Hence, we have that

$$\max_{1 \leqslant i \leqslant N} \|\varphi_{rec,i} - f_{rec}\|_{W^{m,\infty}([a_{i-2}, a_{i+1} \wedge 1])} \leqslant \varepsilon_{rec}/2. \tag{39}$$

This implies that

$$\|\varphi_{rec,i}\|_{W^{m,\infty}([a_{i-2}, a_{i+1} \wedge 1])} \leqslant \varepsilon_{rec} + \exp\{\mathcal{O}(m \log(m/a_0))\} \leqslant \exp\{\mathcal{O}(mN + m \log m)\} \tag{40}$$

and also

$$\|\varphi_{rec,i}\|_{W^{0,\infty}([a_{i-2}, a_{i+1} \wedge 1])} \leqslant \varepsilon_{rec} + 1/a_{-1} \leqslant 4/a_0. \tag{41}$$

The networks $\{\varphi_{id,i}\}_{i=1}^N$ approximate the identity operation (see Lemma 3.2) with the accuracy parameter $\varepsilon_{id,\varphi} \in (0,1)$, the scale parameter $4/a_0$. Similarly, the networks $\{\psi_{id,i}\}_{i=1}^N$ aim to approximate the identity operation with the accuracy parameter $\varepsilon_{id,\psi} \in (0,1)$ and the scale parameter $\|\psi_{pou,i}\|_{W^{m,\infty}(\mathbb{R})}$. The number of layers of the identity networks will be specified later in the proof. We also put the smoothness parameter $m \in \mathbb{N}$ with $m \geqslant 3$ for all the networks. First, we derive the approximation accuracy of $\varphi_i$. Note that according to Lemma 3.12, for all $i \in \{1, \ldots, N\}$, we have $\varphi_{rec,i} \in \mathsf{NN}(L_{rec}, W_{rec}, S_{rec}, B_{rec})$ with

$$L_{rec} \lesssim \log(m \log(1/\varepsilon_{rec})), \quad \|W_{rec}\| \vee S_{rec} \lesssim m^4 \log^4(m/\varepsilon_{rec}a_0) \lesssim m^4 N^4 + m^4 \log^4(m/\varepsilon_{rec}),$$
$$\log B_{rec} \lesssim m^4 \log^4(m/\varepsilon_{rec}a_0) \lesssim m^4 N^4 + m^4 \log^4(m/\varepsilon_{rec}). \tag{42}$$

Therefore, each $\varphi_{id,i}$ has at most $L_{id,\varphi} \lesssim \log(m \log(1/\varepsilon_{rec}))$ the number of layers. Now Lemma 3.2 in conjunction with (40), (41) and Lemma B.4 imply that

$$\max_{1 \leqslant i \leqslant N} \|\varphi_{id,i} \circ \varphi_{rec,i} - \varphi_{rec,i}\|_{W^{m,\infty}([a_{i-2}, a_{i+1} \wedge 1])} \leqslant \exp\{\mathcal{O}(m^2 N + m^2 \log m)\}\varepsilon_{id,\varphi}.$$

Hence setting

$$\log(1/\varepsilon_{id,\varphi}) \asymp \log(1/\varepsilon_{rec}) + m^2 N + m^2 \log m \tag{43}$$

guarantees that

$$\max_{1 \leqslant i \leqslant N} \|\varphi_{id,i} \circ \varphi_{rec,i} - \varphi_{rec,i}\|_{W^{m,\infty}([a_{i-2}, a_{i+1} \wedge 1])} \leqslant \varepsilon_{rec}/2.$$

This and (39) imply that

$$\max_{1 \leqslant i \leqslant N} \|\varphi_i - f_{rec}\|_{W^{m,\infty}([a_{i-2}, a_{i+1} \wedge 1])} \leqslant \varepsilon_{rec}. \tag{44}$$

Furthermore, for $\varphi_{id,i} \in \mathsf{NN}(L_{id,\varphi}, W_{id,\varphi}, S_{id,\varphi}, B_{id,\varphi})$ we have from (42), (43) and Lemma 3.2 that

$$L_{id,\varphi} \lesssim \log(m \log(1/\varepsilon_{rec})), \quad \|W_{id,\varphi}\|_\infty \lesssim 1,$$
$$S_{id,\varphi} \lesssim \log(m \log(1/\varepsilon_{rec})), \quad \log B_{id,\varphi} \lesssim m^2 N + m^2 \log m + \log(1/\varepsilon_{rec}) \log m. \tag{45}$$

Therefore, from Lemma B.5 we find that $\varphi_i = \varphi_{rec,i} \circ \varphi_{id,i} \in \mathsf{NN}(L_{rec}, W_{rec}, S_{rec}, B_{rec})$, where the parameters of the neural network class are presented in (42). Next, we deduce from (43), Lemma 3.2 and Lemma 3.12 that

$$\|\varphi_{rec,i} \circ \varphi_{id,i}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m \log(m + \|\varphi_{id,i}\|_{W^{m,\infty}(\mathbb{R})}))\}\|\varphi_{rec,i}\|_{W^{m,\infty}(\mathbb{R})}$$
$$\leqslant \exp\{\mathcal{O}(m^2 N + m^2 \log(mN \log(1/\varepsilon_{rec})))\}. \tag{46}$$

22

Following this, consider the approximation properties of $\psi_{id,i} \circ \psi_{pou,i}$. From Lemma 3.4 and Lemma B.4 we find that

$$\max_{1 \leqslant i \leqslant N} \|\psi_{id,i} \circ \psi_{pou,i}\|_{W^{m,\infty}([a_0,1]\setminus[a_{i-2},a_{i+1}\wedge 1])} \leqslant \exp\{\mathcal{O}(m \log m)\}\|\psi_{id,i}\|_{W^{m,\infty}([-\varepsilon_{pou},\varepsilon_{pou}])}(1 \vee \varepsilon_{pou}^m)$$

$$\leqslant \exp\{\mathcal{O}(m \log m)\}(\varepsilon_{pou} + \varepsilon_{id,\psi})$$

and also

$$\|\psi_{id,i} \circ \psi_{pou,i} - \psi_{pou,i}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m \log m)\}\varepsilon_{id,\psi}(1 \vee \|\psi_{pou,i}\|_{W^{m,\infty}(\mathbb{R})}^m)$$

$$\leqslant \exp\{\mathcal{O}(m^2 N + m^2 \log(m \log(1/\varepsilon_{pou})))\}\varepsilon_{id,\psi}.$$

Therefore, setting

$$\log(1/\varepsilon_{id,\psi}) \asymp m^2 \log(1/\varepsilon_{pou}) + m^2 N + m^2 \log m \tag{47}$$

ensures that

$$\max_{1 \leqslant i \leqslant N} \|\psi_{id,i} \circ \psi_{pou,i} - \psi_{pou,i}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \varepsilon_{pou}. \tag{48}$$

and

$$\max_{1 \leqslant i \leqslant N} \|\psi_{id,i} \circ \psi_{pou,i}\|_{W^{m,\infty}([a_0,1]\setminus[a_{i-2},a_{i+1}\wedge 1])} \leqslant \exp\{\mathcal{O}(m \log m)\}\varepsilon_{pou}. \tag{49}$$

From (48) and Lemma 3.4 we also deduce that

$$\max_{1 \leqslant i \leqslant N} \|\psi_i\|_{W^{m,\infty}(\mathbb{R})} \leqslant \varepsilon_{pou} + \max_{1 \leqslant i \leqslant N} \|\psi_{pou,i}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(mN + m \log(m \log(1/\varepsilon_{pou})))\}. \tag{50}$$

Now derive the approximation error bound for $\varphi_{rec}$ defined in (38). The triangle inequality suggests that

$$\|\varphi_{rec} - f_{rec}\|_{W^{m,\infty}([a_0,1])} \leqslant \underbrace{\left\|f_{rec}\left(1 - \sum_{i=1}^{N} \psi_i\right)\right\|_{W^{m,\infty}([a_0,1])}}_{(A)}$$

$$+ \underbrace{\left\|\sum_{i=1}^{N}(f_{rec} - \varphi_i)\psi_i\right\|_{W^{m,\infty}([a_0,1])}}_{(B)} + \underbrace{\left\|\sum_{i=1}^{N} \varphi_i \cdot \psi_i - q(\varphi_i, \psi_i)\right\|_{W^{m,\infty}([a_0,1])}}_{(C)}. \tag{51}$$

**Step 2: bounding term** $(A)$. From Lemma B.3 we deduce that

$$\left\|f_{rec}\left(1 - \sum_{i=1}^{N} \psi_i\right)\right\|_{W^{m,\infty}([a_0,1])} \leqslant 2^m \|f_{rec}\|_{W^{m,\infty}([a_0,1])} \left\|1 - \sum_{i=1}^{N} \psi_i\right\|_{W^{m,\infty}([a_0,1])}.$$

Since $\sum_{i=1}^{N} \psi_{pou,i}(x) = 1$ for all $x \in \mathbb{R}$ due to Lemma 3.4 and $\|f_{rec}\|_{W^{m,\infty}([a_0,1])} \leqslant \exp\{\mathcal{O}(m \log(m/a_0))\}$, we obtain from (48) that

$$\left\|f_{rec}\left(1 - \sum_{i=1}^{N} \psi_i\right)\right\|_{W^{m,\infty}([a_0,1])} \leqslant \exp\{\mathcal{O}(m \log(m/a_0))\}N\varepsilon_{pou}$$

$$\leqslant \exp\{\mathcal{O}(mN + m \log m)\}\varepsilon_{pou}. \tag{52}$$

**Step 3: bounding term** $(B)$. The triangle inequality suggests that

$$\left\|\sum_{i=1}^{N}(f_{rec}-\varphi_i)\psi_i\right\|_{W^{m,\infty}([a_0,1])} \leqslant N \max_{1\leqslant i\leqslant N}\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_0,1])}. \tag{53}$$

Moreover, for each $1 \leqslant i \leqslant N$ we have that

$$\begin{aligned}&\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_0,1])}\\&\quad = \|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_{i-2},a_{i+1}\wedge 1])} \vee \|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_0,1]\setminus[a_{i-2},a_{i+1}\wedge 1])}. \end{aligned} \tag{54}$$

Now we analyze each term separately. As for the first term, we obtain from (44), (50) and Lemma B.3 that

$$\begin{aligned}\max_{1\leqslant i\leqslant N}\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_{i-2},a_{i+1}\wedge 1])} &\leqslant 2^m\varepsilon_{rec}\|\psi_i\|_{W^{m,\infty}(\mathbb{R})}\\&\leqslant \exp\{\mathcal{O}(mN+m\log(m\log(1/\varepsilon_{pou})))\}\varepsilon_{rec}.\end{aligned} \tag{55}$$

As for the second term, (46) together with (49), the fact that $\|f_{rec}\|_{W^{m,\infty}([a_0,1])} \leqslant \exp\{\mathcal{O}(m\log m+mN)\}$ and Lemma B.3 yield that for all $1 \leqslant i \leqslant N$

$$\begin{aligned}\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_0,1]\setminus[a_{i-2},a_{i+1}\wedge 1])} &\leqslant \exp\{\mathcal{O}(m\log m)\}(\|\varphi_i\|_{W^{m,\infty}(\mathbb{R})} + \|f_{rec}\|_{W^{m,\infty}([a_0,1])})\varepsilon_{pou}\\&\leqslant \exp\{\mathcal{O}(m^2N+m^2\log(mN\log(1/\varepsilon_{rec})))\}\varepsilon_{pou}.\end{aligned}$$

Therefore, setting

$$\log(1/\varepsilon_{pou}) \asymp \log(1/\varepsilon_{rec}) + m^2N + m^2\log(mN\log(1/\varepsilon_{rec})) \tag{56}$$

guarantees

$$\max_{1\leqslant i\leqslant N}\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_0,1]\setminus[a_{i-2},a_{i+1}\wedge 1])} \leqslant \frac{\varepsilon_{rec}}{3N}. \tag{57}$$

Moreover, from (55) we deduce that

$$\max_{1\leqslant i\leqslant N}\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_{i-2},a_{i+1}\wedge 1])} \leqslant \exp\{\mathcal{O}(mN+m\log(m\log(1/\varepsilon_{rec})))\}\varepsilon_{rec}.$$

Thus, setting

$$\log(1/\varepsilon_{rec}) \asymp mN + m^2\log(1/\varepsilon) \tag{58}$$

ensures that

$$\max_{1\leqslant i\leqslant N}\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_{i-2},a_{i+1}\wedge 1])} \leqslant \frac{\varepsilon}{3N}. \tag{59}$$

The combination of (54), (57) and (59) imply that

$$\max_{1\leqslant i\leqslant N}\|(f_{rec}-\varphi_i)\psi_i\|_{W^{m,\infty}([a_0,1])} \leqslant \frac{\varepsilon}{3N}.$$

Therefore, we deduce from (53) that

$$\left\|\sum_{i=1}^{N}(f_{rec}-\varphi_i)\psi_i\right\|_{W^{m,\infty}([a_0,1])} \leqslant \varepsilon/3. \tag{60}$$

**Step 4: bounding term $(C)$.**  We first note from (46) and (58) that

$$\max_{1\leqslant i\leqslant N}\|\varphi_i\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m^2N+m^2\log(m\log(1/\varepsilon)))\}. \tag{61}$$

In addition, from (50), (56) and (58) we find that

$$\max_{1\leqslant i\leqslant N}\|\psi_i\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(mN+m\log(m\log(1/\varepsilon_{pou})))\}$$
$$\leqslant \exp\{\mathcal{O}(mN+m\log(m\log(1/\varepsilon)))\}. \tag{62}$$

These observations together with Corollary 3.7 and Lemma B.4 imply that

$$\left\|\sum_{i=1}^{N}\varphi_i\cdot\psi_i - q(\varphi_i,\psi_i)\right\|_{W^{m,\infty}([a_0,1])} \leqslant \sum_{i=1}^{N}\exp\{\mathcal{O}(m\log(m+\|\varphi_i\|_{W^{m,\infty}(\mathbb{R})}+\|\psi_i\|_{W^{m,\infty}(\mathbb{R})}))\}\varepsilon_{mul}$$
$$\leqslant \exp\{\mathcal{O}(m^3N+m^3\log(m\log(1/\varepsilon)))\}\varepsilon_{mul}.$$

Thus, setting

$$\log(1/\varepsilon_{mul}) \asymp m^3N+m^3\log(m/\varepsilon) \tag{63}$$

guarantees that

$$\left\|\sum_{i=1}^{N}\varphi_i\cdot\psi_i - q(\varphi_i,\psi_i)\right\|_{W^{m,\infty}([a_0,1])} \leqslant \varepsilon/3. \tag{64}$$

**Step 5: combining $(A)$, $(B)$ and $(C)$ together.**  From (52), (56) and (58) we deduce that the term $(A)$ is evaluated as

$$\left\|f_{rec}\left(1-\sum_{i=1}^{N}\psi_i\right)\right\|_{W^{m,\infty}([a_0,1])} \leqslant \varepsilon/3.$$

Therefore, combining this bound with (51), (60) and (64) yields

$$\|\varphi_{rec}-f_{rec}\|_{W^{m,\infty}([a_0,1])} \leqslant \varepsilon.$$

In addition, from Lemma B.4 we deduce that

$$\|\varphi_{rec}\|_{W^{m,\infty}(\mathbb{R})} \leqslant N\max_{1\leqslant i\leqslant N}\|q(\phi_i,\psi_i)\|_{W^{m,\infty}(\mathbb{R})}$$
$$\leqslant N\exp\{\mathcal{O}(m\log m)\}\max_{1\leqslant i\leqslant N}\|q\|_{W^{m,\infty}(\varphi_i(\mathbb{R})\times\psi_i(\mathbb{R}))}(1\vee\|\varphi_i\|_{W^{m,\infty}(\mathbb{R})}^m\vee\|\psi_i\|_{W^{m,\infty}(\mathbb{R})}^m).$$

Corollary 3.7 together with (61) and (62)

$$\|\varphi_{rec}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \max_{1 \leqslant i \leqslant N} N \exp\{\mathcal{O}(m \log(m + \|\varphi_i\|_{W^{m,\infty}(\mathbb{R})} + \|\psi_i\|_{W^{m,\infty}(\mathbb{R})}))\}$$

$$\leqslant \exp\{\mathcal{O}(m^3 N + m^3 \log(m \log(1/\varepsilon)))\}.$$

**Step 6: deriving the configuration of $\varphi_{rec}$.** First, from (42) and (58) it follows that for each $i \in \{1, \ldots, N\}$, we have $\varphi_i \in \mathsf{NN}(L_{rec}, W_{rec}, S_{rec}, B_{rec})$ with

$$L_{rec} \lesssim \log(mN) + \log\log(1/\varepsilon), \quad \|W_{rec}\|_\infty \vee S_{rec} \lesssim m^8(N^4 + m^4 \log^4(1/\varepsilon)),$$

$$\log B_{rec} \lesssim m^8(N^4 + m^4 \log^4(1/\varepsilon)).$$

Second, from (47), (56) and (58) we deduce that

$$\log(1/\varepsilon_{pou}) \lesssim m^2(N + \log(m/\varepsilon)), \quad \log(1/\varepsilon_{id,\psi}) \lesssim m^4(N + \log(m/\varepsilon)). \tag{65}$$

Hence, (45), (58) and Lemma 3.2 imply that for each $1 \leqslant i \leqslant N$ we have $\psi_{id,i} \in \mathsf{NN}(L_{id,\psi}, W_{id,\psi}, S_{id,\psi}, B_{id,\psi})$ with

$$L_{id,\psi} \vee S_{id,\psi} \lesssim L_{id,\varphi} \lesssim \log(mN \log(1/\varepsilon)), \quad \|W_{id,\psi}\|_\infty \lesssim 1,$$

$$\log B_{id,\psi} \lesssim (m + L_{id,\psi}) \log m + \log(1/\varepsilon_{id,\psi}) + m \log\left(\max_{1 \leqslant i \leqslant N} \|\psi_{pou,i}\|_{W^{m,\infty}(\mathbb{R})}\right).$$

Form (65) and Lemma 3.4 it follows that

$$\log B_{id,\psi} \lesssim m^4(N + \log(m/\varepsilon)).$$

Lemma 3.4 together with (65) imply that for all $1 \leqslant i \leqslant N$ it holds that $\psi_{pou,i} \in \mathsf{NN}(L_{pou}, W_{pou}, S_{pou}, B_{pou})$ with

$$L_{pou} \vee \|W_{pou}\| \vee S_{pou} \lesssim 1, \quad \log B_{pou} \lesssim \log(1/\varepsilon_{pou}) + mN + m \log m \lesssim m^2(N + \log(m/\varepsilon)).$$

Therefore, Lemma B.5 yields that $\psi_i = \psi_{id,i} \circ \psi_{pou,i} \in \mathsf{NN}(L_{id,\psi}, W_{id,\psi}, S_{id,\psi}, B_{id,\psi})$. In addition, from (63) and Corollary 3.7 we find that $q \in \mathsf{NN}(L_{mul}, W_{mul}, S_{mul}, B_{mul})$ with

$$L_{mul} \vee \|W_{mul}\| \vee S_{mul} \lesssim 1, \quad \log B_{mul} \lesssim \log m + \log(1/\varepsilon_{mul}) \lesssim m^3(N + \log(m/\varepsilon)).$$

Thus, applying Lemma B.5 and Lemma B.6, we obtain that

$$L \lesssim L_{mul} + L_{id,\psi} \lesssim \log(mN \log(1/\varepsilon)), \quad \|W\|_\infty \vee S \lesssim N\|W_{rec}\|_\infty \lesssim m^8 N(N^4 + m^4 \log^4(1/\varepsilon)),$$

$$\log B \lesssim \log B_{rec} + \log B_{id,\psi} + \log \|W_{rec}\|_\infty + \log N \lesssim m^8(N^4 + m^4 \log^4(1/\varepsilon)).$$

The proof is complete.

$\square$

By comparing our Lemma 3.13 to [Yakovlev and Puchkin, 2025, Lemma A.4], we see that the parameter count for high-order Sobolev approximation is $\mathcal{O}(N^5 + N \log^4(1/\varepsilon))$, slightly exceeding their bound of $\mathcal{O}(N^4 + N \log^3(1/\varepsilon))$. Nevertheless, we generalize the approximation capabilities to high-order Sobolev norms.

Finally, we present a result on the division approximation, utilizing the reciprocal-based approach we have developed.

**Lemma 3.14** (division operation approximation). *Define* $\mathrm{div} : (x, y) \mapsto x/y$ *for any* $x \in \mathbb{R}$ *and* $y > 0$. *Let also* $a_0 = 2^{-N}$ *for* $N \in \mathbb{N}$ *with* $N \geqslant 3$. *Then, for every* $\varepsilon \in (0, 1)$ *and every* $m \in \mathbb{N}$ *such that* $m \geqslant 3$, *there exists a GELU network* $\varphi_{div} \in \mathsf{NN}(L, W, S, B)$ *satisfying*

$$
\begin{aligned}
(i) &\quad \|\varphi_{div} - \mathrm{div}\|_{W^{m,\infty}([-1,1] \times [a_0, 1])} \leqslant \varepsilon, \\
(ii) &\quad \|\varphi_{div}\|_{W^{m,\infty}(\mathbb{R}^2)} \leqslant \exp\{\mathcal{O}(m^4 N + m^4 \log(m \log(1/\varepsilon)))\}.
\end{aligned}
$$

*Furthermore, the network* $\varphi_{div}$ *has*

$$
L \lesssim \log(mN \log(1/\varepsilon)), \quad \|W\|_\infty \vee S \lesssim m^{21} N^5 \log^4(1/\varepsilon), \quad \log B \lesssim m^{24} N^4 \log^4(1/\varepsilon).
$$

The proof of Lemma 3.14 can be found in Appendix A.5.

# References

A. Abdeljawad and T. Dittrich. Weighted Sobolev approximation rates for neural networks on unbounded domains. Preprint. ArXiv:2411.04108, 2024.

E. Abdo, L. Chai, R. Hu, and X. Yang. Error estimates of physics-informed neural networks for approximating Boltzmann equation. Preprint. ArXiv:2407.08383, 2024.

M. Á. Alejo, L. Cossetti, L. Fanelli, C. Muñoz, and N. Valenzuela. Error bounds for physics informed neural networks in nonlinear Schrödinger equations placed on unbounded domains. Preprint. ArXiv:2409.17938, 2024.

I. Azangulov, G. Deligiannidis, and J. Rousseau. Convergence of diffusion models under the manifold hypothesis in high-dimensions. Preprint. ArXiv:2409.18804, 2024.

D. Belomestny, A. Naumov, N. Puchkin, and S. Samsonov. Simultaneous approximation of a smooth function and its derivatives by deep neural networks with piecewise-polynomial activations. *Neural Networks*, 161:242–253, 2023.

G. M. Constantine and T. H. Savits. A multivariate Faa di Bruno formula with applications. *Transactions of the American Mathematical Society*, 348(2):503–520, 1996.

T. De Ryck, S. Lanthaler, and S. Mishra. On the approximation of functions by tanh neural networks. *Neural Networks*, 143:732–750, 2021.

J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019.

K. Fukumizu, T. Suzuki, N. Isobe, K. Oko, and M. Koyama. Flow matching achieves almost minimax optimal convergence. In *The Thirteenth International Conference on Learning Representations*, 2025.

I. Gühring and M. Raslan. Approximation rates for neural networks with encodable weights in smoothness spaces. *Neural Networks*, 134:107–130, 2021.

D. Hendrycks and K. Gimpel. Gaussian Error Linear Units (GELUs). Preprint. ArXiv:1606.08415, 2016.

R. Nakada and M. Imaizumi. Adaptive approximation and generalization of deep neural network with intrinsic dimensionality. *Journal of Machine Learning Research*, 21(174):1–38, 2020a.

R. Nakada and M. Imaizumi. Adaptive approximation and generalization of deep neural network with intrinsic dimensionality. *Journal of Machine Learning Research*, 21(174):1–38, 2020b.

K. Oko, S. Akiyama, and T. Suzuki. Diffusion models are minimax optimal distribution estimators. In *International Conference on Machine Learning*, pages 26517–26582. PMLR, 2023.

N. Puchkin, E. Gorbunov, N. Kutuzov, and A. Gasnikov. Breaking the heavy-tailed noise barrier in stochastic optimization problems. In *International Conference on Artificial Intelligence and Statistics*, pages 856–864. PMLR, 2024.

C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.

F. Scarselli and A. C. Tsoi. Universal approximation using feedforward neural networks: A survey of some existing methods, and some new results. *Neural networks*, 11(1):15–37, 1998.

C. Schwab and J. Zech. Deep learning in high dimension: Neural network expression rates for generalized polynomial chaos expansions in UQ. *Analysis and Applications*, 17(01):19–55, 2019.

C. Schwab and J. Zech. Deep learning in high dimension: Neural network approximation of analytic functions in $(L^2(\mathbb{R}^d, \gamma_d))$. Preprint. ArXiv:2111.07080, 2021.

M. Shoeybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro. Megatron-LM: Training multi-billion parameter language models using model parallelism. Preprint. ArXiv:1909.08053, 2019.

R. Tang and Y. Yang. Adaptivity of diffusion models to manifold structures. In S. Dasgupta, S. Mandt, and Y. Li, editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 1648–1656. PMLR, 2024.

T. D. van Nuland. Noncompact uniform universal approximation. *Neural Networks*, 173:106181, 2024.

K. Yakovlev and N. Puchkin. Generalization error bound for denoising score matching under relaxed manifold assumption. In *Proceedings of Thirty Eighth Conference on Learning Theory*, volume 291 of *Proceedings of Machine Learning Research*, pages 5824–5891. PMLR, 2025.

D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:103–114, 2017.

# A  Deferred proofs

## A.1  Proof of Lemma 3.2

We first prove the statement for $K = 1$ and then generalize to $K \geqslant 1$. Let

$$\varphi_j = \varphi_{id,j} \circ \varphi_{j-1}, \quad 2 \leqslant j \leqslant L,$$

where $\varphi_1 = \text{id}$ and $\varphi_{id,j}$ is a GELU network from Lemma 3.1 that approximates the identity operation with the accuracy parameter $\varepsilon_{id}^{(j)}$. Formally, for each $2 \leqslant j \leqslant L$ we have

$$\|\varphi_{id,j} - \text{id}\|_{W^{m,\infty}([-C,C])} \leqslant C^2 \varepsilon_{id}^{(j)} \quad \text{for all } C \geqslant 1. \tag{66}$$

For every $1 \leqslant j \leqslant L$ we introduce $\varepsilon_j = \|\varphi_j - \mathrm{id}\|_{W^{m,\infty}([-1,1])}$. Therefore, the triangle inequality for every $2 \leqslant j \leqslant L$ implies that

$$\|\varphi_{id,j} \circ \varphi_{j-1} - \mathrm{id}\|_{W^{m,\infty}([-1,1])} \leqslant \|\varphi_{j-1} - \mathrm{id}\|_{W^{m,\infty}([-1,1])} + \|(\mathrm{id} - \varphi_{id,j}) \circ \varphi_{j-1}\|_{W^{m,\infty}([-1,1])}.$$

Next, applying Lemma B.4, we obtain that

$$\varepsilon_j \leqslant \varepsilon_{j-1} + 16(e^2 m^4)^m \|\varphi_{id,j} - \mathrm{id}\|_{W^{m,\infty}([-1-\varepsilon_{j-1},1+\varepsilon_{j-1}])}(1 \vee \|\varphi_{j-1}\|_{W^{m,\infty}([-1,1])}^m),$$

Therefore, (66) suggests that

$$\varepsilon_j \leqslant \varepsilon_{j-1} + 16(e^2 m^4)^m (1 + \varepsilon_{j-1})^{m+2} \varepsilon_{id}^{(j)}.$$

Now choosing

$$\varepsilon_{id}^{(j)} = 16(e^2 m^4)^{-m} \varepsilon_2 \in (0,1), \quad 2 \leqslant j \leqslant L, \tag{67}$$

we find that

$$\varepsilon_j \leqslant \varepsilon_{j-1} + (1 + \varepsilon_{j-1})^{m+2} \varepsilon_{j-1}, \quad 2 \leqslant j \leqslant L. \tag{68}$$

Suppose that for each $2 \leqslant j \leqslant L$, the approximation error is given by $\varepsilon_j = 2^{\gamma_j} \varepsilon_2$ with $\gamma_2 = 0$. We also set a helper $\gamma_1 = 0$. Hence, considering (68), we conclude that

$$2^{\gamma_j} \varepsilon_2 \leqslant 2^{\gamma_{j-1}} \varepsilon_2 + 2^{(m+2)(\gamma_{j-1}+1)} \varepsilon_2 \leqslant 2^{(m+3)\gamma_{j-1}+m+3} \varepsilon_2.$$

Therefore, $\gamma_j \leqslant (2(m+3))^j$ for each $2 \leqslant j \leqslant L$. Setting $\varepsilon_{id}^{(2)} = \varepsilon'(2(m+3))^{-L}$ for some $\varepsilon' \in (0,1)$, we deduce from (67) that for any $2 \leqslant j \leqslant L$

$$\log(1/\varepsilon_{id}^{(j)}) \lesssim m \log m + \log(1/\varepsilon_{mul}^{(2)}) \lesssim (m + L) \log m + \log(1/\varepsilon').$$

Moreover,

$$\|\varphi_L - \mathrm{id}\|_{W^{m,\infty}([-1,1])} \leqslant \varepsilon'.$$

Next, using Lemma B.5 and Lemma 3.1 we find that $\varphi_L \in \mathsf{NN}(L, W_{id}, S_{id}, B_{id})$ with

$$\|W_{id}\|_\infty \lesssim 1, \quad S_{id} \lesssim L, \quad \log B_{id} \lesssim \log m + \max_{2 \leqslant j \leqslant L} \log(1/\varepsilon_{id}^{(j)}) \lesssim (m + L) \log m + \log(1/\varepsilon'). \tag{69}$$

The generalization to the case when $K \geqslant 1$ is trivial. Let $\varphi_{L,K}(x) = K\varphi_L(x/K)$ for any $x \in \mathbb{R}$. Then we have that

$$\|\varphi_{L,K} - \mathrm{id}\|_{W^{m,\infty}([-K,K])} \leqslant K\|\varphi_L - \mathrm{id}\|_{W^{m,\infty}([-1,1])} \leqslant K\varepsilon'. \tag{70}$$

As a final step, we add a clipping operation to ensure that the resulting function has finite norm on a real line. Let $\varphi_{clip}$ be a clipping operation approximation from Lemma 3.5 with the accuracy parameter $\varepsilon'$ and the scale parameter $K$. Thus, we have that

$$\|\varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-K,K])} \leqslant \varepsilon', \quad \|\varphi_{clip}\|_{W^{0,\infty}(\mathbb{R})} \leqslant 4K,$$
$$\|\varphi_{clip}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m \log(m \log(1/\varepsilon')) + \log(2K))\}. \tag{71}$$

In addition, $\varphi_{clip} \in \mathsf{NN}(L_{clip}, W_{clip}, S_{clip}, B_{clip})$ with

$$L_{clip} \vee \|W_{clip}\|_{\infty} \vee S_{clip} \lesssim 1, \quad \log B_{clip} \lesssim \log(Km/\varepsilon'). \tag{72}$$

Then it holds due to the triangle inequality and (70) that

$$\|\varphi_{L,4K} \circ \varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-K,K])} \leqslant 4K\varepsilon' + \|\varphi_{L,4K} \circ \varphi_{clip} - \varphi_{L,4K}\|_{W^{m,\infty}([-K,K])}. \tag{73}$$

Hence, Lemma B.4 together with (71) imply that

$$\|\varphi_{L,4K} \circ \varphi_{clip} - \varphi_{L,4K}\|_{W^{m,\infty}([-K,K])} \leqslant \exp\{\mathcal{O}(m\log m)\}\|\varphi_{L,4K}\|_{W^{m+1,\infty}([-4K,4K])}\varepsilon'(\varepsilon' + K)^{2m}.$$

Next we note that (69) and (70) are true if the smoothness parameter is $m+1$ instead of $m$. Then we have from (70) that

$$\|\varphi_{L,4K} \circ \varphi_{clip} - \varphi_{L,4K}\|_{W^{m,\infty}([-K,K])} \leqslant \exp\{\mathcal{O}(m\log(mK))\}\varepsilon'.$$

Therefore, (73) is evaluated as

$$\|\varphi_{L,4K} \circ \varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-K,K])} \leqslant \exp\{\mathcal{O}(m\log(mK))\}\varepsilon'.$$

Thus, setting

$$\log(1/\varepsilon') \asymp \log(1/\varepsilon) + m\log(mK) \tag{74}$$

ensures that

$$\|\varphi_{L,4K} \circ \varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-K,K])} \leqslant \varepsilon.$$

In addition, from (70), (71), (74) and Lemma B.4 we obtain that

$$\|\varphi_{L,4K} \circ \varphi_{clip}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m\log m)\}4K(1+\varepsilon')\exp\{\mathcal{O}(m\log(m\log(1/\varepsilon'))) + \log(2K))\}$$
$$\leqslant \exp\{\mathcal{O}(m\log(m\log(1/\varepsilon))) + \log(2K))\}.$$

Finally, Lemma B.5 combined with (69), (72) and (74) yields that $\varphi_{id} = \varphi_{L,4K} \circ \varphi_{clip}$ has

$$\|W\|_{\infty} \lesssim 1, \quad S \lesssim L, \quad \log B \lesssim (m+L)\log m + \log(1/\varepsilon) + m\log(K).$$

This finishes the proof.

$\square$

## A.2  Proof of Corollary 3.7

We are going to reduce the multiplication to the case of square operation by letting

$$\varphi_{mul}(x,y) := \frac{1}{4}\left(\varphi_{sq}(x+y) - \varphi_{sq}(x-y)\right), \tag{75}$$

where $\varphi_{sq}$ is a GELU network from Lemma 3.6 with the accuracy parameter $\varepsilon/4$. Therefore, using the observation that for any $\alpha = (\alpha_1, \alpha_2)^{\top} \in \mathbb{Z}_+^2$, it holds that

$$D^{\alpha}\varphi_{mul}(x,y) = \frac{1}{4}(D^{|\alpha|}\varphi_{sq}(x+y) - (-1)^{\alpha_2}\varphi_{sq}(x-y)), \quad \text{for all } x,y \in \mathbb{R},$$

leads to

$$\|\varphi_{mul} - \mathrm{prod}_2\|_{W^{m,\infty}([-C,C]^2)} \leqslant \frac{1}{2}\|\varphi_{sq} - f_{sq}\|_{W^{m,\infty}([-2C,2C])} \leqslant C^3\varepsilon,$$

where $C \geqslant 1$ is arbitrary, and the last inequality uses Lemma 3.6. Finally, in view of (75), we deduce that the summation argument outlined in Lemma B.6 yields the configuration in the statement. This completes the proof.

$\square$

## A.3  Proof of Corollary 3.9

We first introduce a flatten operation as follows:

$$\mathrm{flat}_{\mathbf{k}}(x_1, \ldots, x_I) = (\underbrace{x_1, \ldots, x_1}_{k_1 \text{ times}}, \ldots, \underbrace{x_I, \ldots, x_I}_{k_I \text{ times}})^\top.$$

Now let $\varphi_{mul,d}$ be a neural network from Lemma 3.8 with the accuracy parameter $\widetilde{\varepsilon} \in (0,1)$, which will be specified a bit later in the proof, and scale parameter $K$. Therefore, using Lemma B.4, we derive an approximation accuracy for $\varphi_{mul,\mathbf{k}} = \varphi_{mul,d} \circ \mathrm{flat}_{\mathbf{k}}$:

$$\|\varphi_{mul,\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} = \|(\varphi_{mul,d} - \mathrm{prod}_d) \circ \mathrm{flat}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)}$$
$$\leqslant \exp\{\mathcal{O}(m\log(md))\}\|\varphi_{mul,d} - \mathrm{prod}_d\|_{W^{m,\infty}([-K,K]^d)}K^m.$$

Next, the approximation property of $\varphi_{mul,d}$ implies that

$$\|\varphi_{mul,\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \exp\{\mathcal{O}(m\log(mdK))\}\widetilde{\varepsilon}.$$

To continue, we set

$$\log(1/\widetilde{\varepsilon}) \asymp \log(1/\varepsilon) + m\log(mdK) \tag{76}$$

and arrive at

$$\|\varphi_{mul,\mathbf{k}} - \mathrm{prod}_{\mathbf{k}}\|_{W^{m,\infty}([-K,K]^I)} \leqslant \varepsilon.$$

In addition, Lemma 3.8 together with Lemma B.4 suggest that

$$\|\varphi_{mul,\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)} = \|\varphi_{mul,d} \circ \mathrm{flat}_{\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant \exp\{\mathcal{O}((m^2+d)\log(mdK\log(1/\widetilde{\varepsilon})))\}.$$

From (76) we obtain that

$$\|\varphi_{mul,\mathbf{k}}\|_{W^{m,\infty}(\mathbb{R}^I)} \leqslant \exp\{\mathcal{O}((m^2+d)\log(mdK\log(1/\varepsilon)))\}.$$

To finalize the proof, we formulate the configuration of $\varphi_{mul,\mathbf{k}}$. Note that $\mathrm{flat}_{\mathbf{k}}$ is implemented using a single linear layer without a bias term, and its weight matrix contains binary values. Consequently, the choice of $\widetilde{\varepsilon}$ given in (76) combined with Lemma 3.8 and the concatenation result outlined in Lemma B.5 ensures that $\varphi_{mul,\mathbf{k}}$ has $L \lesssim \log d$, $S \vee \|W\|_\infty \lesssim (d \vee I)^3$ and

$$\log B \lesssim (\log(1/\varepsilon) + (d+m)\log K + m^2d^2)\log d + \log I.$$

This completes the proof.

$\square$

## A.4 Proof of Lemma 3.12

For some $r \in \mathbb{N}$ with $r \geqslant m$, which will be optimized later, consider $f_r(x) = \frac{1}{b} \sum_{i=0}^{r-1} (1 - x/b)^i$. We next note that for all $x \in [a, b]$

$$\frac{1}{x} - f_r(x) = \frac{1}{b} \sum_{i=0}^{\infty} (1 - x/b)^i - \frac{1}{b} \sum_{i=0}^{r-1} (1 - x/b)^i = \frac{(1 - x/b)^r}{x}$$

Then, by lemma B.3, the approximation accuracy of $f_r$ is

$$\begin{aligned}
\|f_{rec} - f_r\|_{W^{m,\infty}([a,b])} &\leqslant 2^m \|f_{rec}\|_{W^{m,\infty}([a,b])} (b \wedge 1)^{-m} r^m (1 - a/b)^{r-m} \\
&\leqslant \left( \frac{2r}{b \wedge 1} \right)^m a^{-m-1} (1 - a/b)^{r-m} m! \\
&\leqslant \frac{m!}{a} \left( \frac{2r}{a(b \wedge 1)} \right)^m \exp\left( -\frac{(r-m)a}{b} \right),
\end{aligned}$$

where the last inequality uses $1 + x \leqslant e^x$ for any $x \in \mathbb{R}$. Therefore, setting $r = \lceil m + \frac{b}{a} \log(1/\varepsilon') \rceil$ for some $\varepsilon' \in (0, 1)$ guarantees that

$$\|f_{rec} - f_r\|_{W^{m,\infty}([a,b])} \leqslant \frac{m!}{a} \left( \frac{2r}{a(b \wedge 1)} \right)^m \varepsilon'.$$

We now set $\varepsilon' = \frac{a}{4m!} \left( \frac{a(b \wedge 1)}{2r} \right)^m \varepsilon$ with $\varepsilon \in (0, 1)$, which leads to

$$\|f_{rec} - f_r\|_{W^{m,\infty}([a,b])} \leqslant \varepsilon/4. \tag{77}$$

We also deduce from Stirling's approximation that

$$\begin{aligned}
\log(1/\varepsilon') &\lesssim \log(1/\varepsilon) + m \log(m/ab) + m \log\left( m + \frac{b}{a} \log(1/\varepsilon') \right) \\
&\lesssim \log(1/\varepsilon) + m \log(m/a) + m \log \log(1/\varepsilon').
\end{aligned}$$

The last inequality suggests that $\log(1/\varepsilon') \lesssim \log(1/\varepsilon) + m \log(m/a)$, since the inequality $x \lesssim a + b \log x$ yields $x \lesssim a + b \log b$ for any positive $a$, $b$ and $x$. Let $\varphi_{part}$ be a GELU network from Lemma 3.10 with the accuracy parameter $\varepsilon_{part}$, the scale parameter $K = 1 + 1/b$, the parameter $I = 1$ and $d = r$. Then, for $f_{part}(x) = \sum_{i=0}^{r-1} x^i$ and $\widetilde{\varphi}_{rec} = (1/b) \varphi_{part} \circ (1 - \mathrm{id}/b)$ it holds that

$$\|\widetilde{\varphi}_{rec} - f_r\|_{W^{m,\infty}([a,b])} \leqslant b^{-1} \|(\varphi_{part} - f_{part}) \circ (1 - \mathrm{id}/b)\|_{W^{m,\infty}([a,b])} \leqslant (b \wedge 1)^{-m-1} \varepsilon_{part},$$

where the last inequality follows from the chain rule. Thus, setting

$$\log(1/\varepsilon_{part}) \asymp \log(1/\varepsilon) + m \log(1/a), \tag{78}$$

we obtain from (77) that

$$\|\widetilde{\varphi}_{rec} - f_{rec}\|_{W^{m,\infty}([a,b])} \leqslant \varepsilon/2. \tag{79}$$

32

In addition, Lemma 3.10 together with Lemma B.5 imply that $\widetilde{\varphi}_{rec} \in \mathsf{NN}(L_{rec}, W_{rec}, S_{rec}, B_{rec})$ with

$$L_{rec} \lesssim \log r \lesssim \log((mb/a)\log(1/\varepsilon)), \quad \|W_{rec}\|_\infty \vee S_{rec} \lesssim r^4 \lesssim (mb/a)^4 \log^4(m/\varepsilon a),$$
$$\log B_{rec} \lesssim (\log(1/\varepsilon_{part}) + m^2 r \log(mr) + m^2 r^2) \log r$$
$$\lesssim (m^4 b^2/a^2) \log^2(1/\varepsilon a) \log^2((mb/a)\log(1/\varepsilon a)). \tag{80}$$

We further observe that through rescaling of the parameters $a$ and $b$ and increasing the smoothness parameter $m$, the bound (79) remains valid for $a/2$ and $2b$, while preserving the configuration specified (80) remains the same. Fromally, we have that

$$\|\widetilde{\varphi}_{rec} - f_{rec}\|_{W^{m+1,\infty}([a/2,2b])} \leqslant \varepsilon/2. \tag{81}$$

Furthermore, the derived bound together with the fact that $\|f_{rec}\|_{W^{m,\infty}([a/2,2b])} \leqslant \exp\{\mathcal{O}(m\log m + mN)\}$ imply that

$$\|\widetilde{\varphi}_{rec}\|_{W^{m,\infty}([a/2,2b])} \leqslant \exp\{\mathcal{O}(m\log m + mN)\}. \tag{82}$$

Now let $\varphi_{clip}$ be the clipping operation approximation from Lemma 3.5 with the precision parameter $\varepsilon_{clip}$ and the scale parameter $(4b/a - 4) \geqslant 1$. Let also

$$\breve{\varphi}_{clip}(x) = \varphi_{clip}(x - 4 - 4b/a) + 4b/a + 4, \quad x \in \mathbb{R}.$$

Therefore, from Lemma 3.5 we find that $\breve{\varphi}_{clip}$ satisfies

$$\begin{aligned}
(i) & \quad \|\breve{\varphi}_{clip} - \mathrm{id}\|_{W^{m,\infty}([8,8b/a])} \leqslant \|\varphi_{clip} - \mathrm{id}\|_{W^{m,\infty}([-(4b/a-4),4b/a-4])} \leqslant \varepsilon_{clip}, \\
(ii) & \quad 11/2 \leqslant \breve{\varphi}_{clip}(x) \leqslant 8b/a + 5/2, \quad \text{for all } x \in \mathbb{R}, \\
(iii) & \quad \|\breve{\varphi}_{clip}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m\log m + m\log\log(1/\varepsilon_{clip}) + \log(b/a))\}.
\end{aligned}$$

Moreover, for $\widetilde{\varphi}_{clip}(x) = (a/8)\breve{\varphi}_{clip}(8x/a)$ we deduce from the chain rule and property $(i)$ that

$$\|\widetilde{\varphi}_{clip} - \mathrm{id}\|_{W^{m,\infty}([a,b])} \leqslant \exp\{\mathcal{O}(m\log(1/a))\}\|\breve{\varphi}_{clip} - \mathrm{id}\|_{W^{m,\infty}([8,8b/a])}$$
$$\leqslant \exp\{\mathcal{O}(m\log(1/a))\}\varepsilon_{clip}. \tag{83}$$

From property $(ii)$ it follows that

$$a/2 \leqslant \widetilde{\varphi}_{clip}(x) \leqslant 2b, \quad \text{for all } x \in \mathbb{R}. \tag{84}$$

In addition, property $(iii)$ yields

$$\|\widetilde{\varphi}_{clip}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m\log(m/a) + m\log\log(1/\varepsilon_{clip}))\}. \tag{85}$$

Combining (83), (84) and Lemma B.4, we obtain for $\varphi_{rec} = \widetilde{\varphi}_{rec} \circ \widetilde{\varphi}_{clip}$ that

$$\|\widetilde{\varphi}_{rec} \circ \widetilde{\varphi}_{clip} - \widetilde{\varphi}_{rec}\|_{W^{m,\infty}([a,b])} \leqslant \exp\{\mathcal{O}(m\log(m/a))\}\|\widetilde{\varphi}_{rec}\|_{W^{m+1,\infty}([a/2,2b])}\varepsilon_{clip}.$$

From (82) we find that

$$\|\widetilde{\varphi}_{rec} \circ \widetilde{\varphi}_{clip} - \widetilde{\varphi}_{rec}\|_{W^{m,\infty}([a,b])} \leqslant \exp\{\mathcal{O}(mN + m\log(m/a))\}\varepsilon_{clip}.$$

Choosing

$$\log(1/\varepsilon_{clip}) \asymp \log(1/\varepsilon) + mN + m\log(m/a) \tag{86}$$

and combining the derived bound with (81), it follows that

$$\|\varphi_{rec} - f_{rec}\|_{W^{m,\infty}([a,b])} \leqslant \varepsilon.$$

We also find from (82), (84), (85) and Lemma B.4 that

$$\|\varphi_{rec}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m\log(m + \|\widetilde{\varphi}_{clip}\|_{W^{m,\infty}(\mathbb{R})}))\}\|\widetilde{\varphi}_{rec}\|_{W^{m,\infty}([a/2,2b])}$$
$$\leqslant \exp\{\mathcal{O}(m^2\log(m/a) + m^2\log\log(1/\varepsilon_{clip}) + mN)\}.$$

The choice of $\varepsilon_{clip}$ given in (86) suggests that

$$\|\varphi_{rec}\|_{W^{m,\infty}(\mathbb{R})} \leqslant \exp\{\mathcal{O}(m^2\log(mN/a) + m^2\log\log(1/\varepsilon) + mN)\}.$$

Lemma 3.5 together with (86) imply that $\widetilde{\varphi}_{clip} \in \mathsf{NN}(L_{clip}, W_{clip}, S_{clip}, B_{clip})$ with

$$L_{clip} \vee \|W_{clip}\|_{\infty} \vee S_{clip} \lesssim 1, \quad \log B_{clip} \lesssim \log(1/\varepsilon) + mN + m\log(m/a).$$

Therefore, applying Lemma B.5, we obtain that $\varphi_{rec} \in \mathsf{NN}(L_{rec}, W_{rec}, S_{rec}, B_{rec})$ with the parameters specified in (80). The proof is complete.

$\square$

## A.5 Proof of Lemma 3.14

**Step 1: approximation error decomposition.** Let $\varphi_{rec}$ be a reciprocal function approximation from Lemma 3.13 with the accuracy parameter $\varepsilon_0 \in (0,1)$ and $\varphi_{id}$ be the identity approximation from Lemma 3.2 with the accuracy parameter $\varepsilon_0$, the scale parameter 1, and the number of layers of $\varphi_{rec}$. Let also $\varphi_{mul}$ be a multiplication network from Corollary 3.7 with the precision parameter $\varepsilon_0$. We put the smoothness parameter $m+1$ for all the networks. We also refer to $f_{rec}(x) = 1/x$ for any $x > 0$ as a reciprocal function. Therefore, for $\varphi_{div} = \varphi_{mul}(\varphi_{id}, \varphi_{rec})$ we have due to the triangle inequality that

$$\|\varphi_{mul}(\varphi_{id}, \varphi_{rec}) - \mathrm{id} \cdot f_{rec}\|_{W^{m,\infty}([-1,1]\times[a_0,1])} \tag{87}$$
$$\leqslant \underbrace{\|\varphi_{mul}(\varphi_{id}, \varphi_{rec}) - \varphi_{mul}(\mathrm{id}, f_{rec})\|_{W^{m,\infty}([0,1]\times[a_0,1])}}_{(A)} + \underbrace{\|\varphi_{mul}(\mathrm{id}, f_{rec}) - \mathrm{id} \cdot f_{rec}\|_{W^{m,\infty}([-1,1]\times[a_0,1])}}_{(B)}.$$

Next, we evaluate the terms $(A)$ and $(B)$ individually.

**Step 2: bounding term $(A)$.** Lemma B.4 together with Corollary 3.7 suggests that

$$\|\varphi_{mul}(\varphi_{id}, \varphi_{rec}) - \varphi_{mul}(\mathrm{id}, f_{rec})\|_{W^{m,\infty}([-1,1]\times[a_0,1])}$$
$$\leqslant \exp\{\mathcal{O}(m\log(m + \|\mathrm{id}\|_{W^{m,\infty}([-1,1])} + \|f_{rec}\|_{W^{m,\infty}([a_0,1])}))\}\varepsilon_0.$$

Note that Stirling's approximation yields

$$\|f_{rec}\|_{W^{m,\infty}([a_0,1])} \leqslant a_0^{-(m+1)}m! = \exp\{\mathcal{O}(mN + m\log m)\}. \tag{88}$$

This observation implies that

$$\|\varphi_{mul}(\varphi_{id}, \varphi_{rec}) - \varphi_{mul}(\mathrm{id}, f_{rec})\|_{W^{m,\infty}([-1,1]\times[a_0,1])} \leqslant \exp\{\mathcal{O}(m^2 N + m^2 \log m)\}\varepsilon_0. \qquad (89)$$

**Step 3: bounding term** $(B)$. The bound is obtained in a similar way. Formally, Lemma B.4 implies that

$$\|\varphi_{mul}(\mathrm{id}, f_{rec}) - \mathrm{id} \cdot f_{rec}\|_{W^{m,\infty}([-1,1]\times[a_0,1])}$$
$$\leqslant \exp\{\mathcal{O}(m \log(m + \|\mathrm{id}\|_{W^{m,\infty}([-1,1])} + \|f_{rec}\|_{W^{m,\infty}([a_0,1])}))\}\varepsilon_0.$$

Using (88), we arrive at

$$\|\varphi_{mul}(\mathrm{id}, f_{rec}) - \mathrm{id} \cdot f_{rec}\|_{W^{m,\infty}([-1,1]\times[a_0,1])} \leqslant \exp\{\mathcal{O}(m^2 N + m^2 \log m)\}\varepsilon_0. \qquad (90)$$

**Step 4: combining** $(A)$ **and** $(B)$ **together.** From (87), (89) and (90) we deduce that setting

$$\log(1/\varepsilon_0) \asymp m^2 N + m^2 \log m + \log(1/\varepsilon), \qquad (91)$$

ensures that

$$\|\varphi_{mul}(\varphi_{id}, \varphi_{rec}) - \mathrm{div}\|_{W^{m,\infty}([-1,1]\times[a_0,1])} \leqslant \varepsilon.$$

Moreover, using Corollary 3.7 together with Lemmata 3.2, 3.13, and B.4, we deduce that

$$\|\varphi_{div}\|_{W^{m,\infty}(\mathbb{R}^2)} \leqslant \exp\{\mathcal{O}(m \log(m + \|\varphi_{id}\|_{W^{m,\infty}(\mathbb{R})} + \|\varphi_{rec}\|_{W^{m,\infty}(\mathbb{R})}))\}$$
$$\leqslant \exp\{\mathcal{O}(m^4 N + m^4 \log(m \log(1/\varepsilon)))\}.$$

**Step 5: deriving the configuration of** $\varphi_{div}$**.** We find from (91) and Lemma 3.13 that $\varphi_{rec}$ belongs to the neural network class $\mathsf{NN}(L_{rec}, W_{rec}, S_{rec}, B_{rec})$ with

$$L_{rec} \lesssim \log(mN \log(1/\varepsilon)), \quad \|W_{rec}\|_\infty \vee S_{rec} \lesssim m^8 N(N^4 + m^4 \log^4(1/\varepsilon_0)) \lesssim m^{21} N^5 \log^4(1/\varepsilon),$$
$$\log B_{rec} \lesssim m^{24} N^4 \log^4(1/\varepsilon).$$

Moreover, the bound for $L_{rec}$ together with Lemma 3.2 suggest that $\varphi_{id} \in \mathsf{NN}(L_{id}, W_{id}, S_{id}, B_{id})$ with

$$L_{id} \vee S_{id} \lesssim L_{rec} \lesssim \log(mN \log(1/\varepsilon)), \quad \|W_{id}\|_\infty \lesssim 1,$$
$$\log B_{id} \lesssim (m + L_{rec}) \log m + \log(1/\varepsilon_0) \lesssim m^2(N + \log m) + \log m \cdot \log(1/\varepsilon).$$

In addition, due to Corollary 3.7, it holds that $\varphi_{mul} \in \mathsf{NN}(L_{mul}, W_{mul}, S_{mul}, B_{mul})$ with

$$L_{mul} \vee \|W_{mul}\| \vee S_{mul} \lesssim 1, \quad \log B_{mul} \lesssim m^2 N + m^2 \log m + \log(1/\varepsilon).$$

Therefore, Lemma B.6 and Lemma B.5 imply that $\varphi_{div}$ has

$$L \lesssim \log(mN \log(1/\varepsilon)), \quad \|W\|_\infty \vee S \lesssim m^{21} N^5 \log^4(1/\varepsilon), \quad \log B \lesssim m^{24} N^4 \log^4(1/\varepsilon).$$

The proof is complete.

$\square$

# B  Auxiliary results

**Lemma B.1** (evaluation of Hermite polynomials, Puchkin et al. [2024], Appendix D). *For any $n \in \mathbb{N}$ we define a "probabilist's" Hermite polynomial*

$$\mathcal{H}_n(x) = (-1)^n e^{x^2/2} \frac{\mathrm{d}^n}{\mathrm{d}x^n} e^{-x^2/2}, \quad x \in \mathbb{R}.$$

*Then it holds that*

$$\max_{x \in \mathbb{R}} \left| \mathcal{H}_n(x) e^{-x^2/4} \right| \leqslant \sqrt{n!} \quad \text{for all } n \in \mathbb{N}.$$

**Lemma B.2** (properties of GELU acitvation function). *For any $k \in \mathbb{N}$ we have the following bounds for the Sobolev seminorms:*

$$|\mathrm{GELU}|_{W^{k,\infty}(\mathbb{R})} \leqslant \begin{cases} 1 + 1/\sqrt{2\pi}, & k = 1, \\ (k+1)\sqrt{\frac{(k-2)!}{2\pi}}, & k \geqslant 2 \end{cases}$$

*For $k = 0$ we have that*

$$\|\mathrm{GELU}\|_{W^{0,\infty}([-C,C])} \leqslant C, \quad \text{for all } C > 0.$$

*In addition, for any $A \geqslant 0$ and $m \in \mathbb{N}$, the tails behave as follows:*

$$\|\mathrm{GELU} - \mathrm{id}\|_{W^{m,\infty}([A,+\infty))} \vee \|\mathrm{GELU}\|_{W^{m,\infty}((-\infty,-A])} \leqslant 2e^{-A^2/4}\sqrt{m!}.$$

*Proof.* We first recall that

$$\mathrm{GELU}(x) = x \cdot \Phi(x), \quad \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-t^2/2} \, \mathrm{d}t,$$

which immediately implies that $\|\mathrm{GELU}\|_{W^{0,\infty}([-C,C])} \leqslant C$ for any $C > 0$ and

$$\partial^1 \mathrm{GELU}(x) = \Phi(x) - \frac{1}{\sqrt{2\pi}} \frac{\mathrm{d}}{\mathrm{d}x} e^{-x^2/2}. \tag{92}$$

Hence, using Lemma B.1 together with the observation that $\partial^k e^{-x^2/2} = (-1)^k e^{-x^2/2} \mathcal{H}_k(x)$ for any $k \in \mathbb{N}$, where

$$\mathcal{H}_k(x) = (-1)^k e^{x^2/2} \partial^k e^{-x^2/2}, \quad x \in \mathbb{R},$$

we obtain that $|\mathrm{GELU}|_{W^{1,\infty}(\mathbb{R})} \leqslant 1 + 1/\sqrt{2\pi}$. Subsequently, from (92) we deduce that for any $k \in \mathbb{N}$ with $k \geqslant 2$

$$\partial^k \mathrm{GELU}(x) = \frac{1}{\sqrt{2\pi}} \left( \partial^{k-2}(e^{-x^2/2}) - \partial^k(e^{-x^2/2}) \right). \tag{93}$$

Applying Lemma B.1, we have that

$$|\text{GELU}|_{W^{k,\infty}(\mathbb{R})} \leqslant \frac{1}{\sqrt{2\pi}} \left( \sqrt{(k-2)!} + \sqrt{k!} \right) \leqslant (k+1)\sqrt{\frac{(k-2)!}{2\pi}},$$

which validates the first claim of the statement. Now focus on the behavior of tails. First, consider

$$\|\text{GELU} - \text{id}\|_{W^{0,\infty}([A,+\infty))} = \sup_{x \geqslant A} x(1 - \Phi(x)) \leqslant \sup_{x \geqslant A} x e^{-x^2/2} \leqslant e^{-A^2/4}\sqrt{2}e^{-1/2},$$

where the penultimate inequality uses Gaussian tails and the last inequality follows from the observation that $xe^{-x^2/4} \leqslant \sqrt{2}e^{-1/2}$ for all $x \in \mathbb{R}$. Similarly,

$$\|\text{GELU}\|_{W^{0,\infty}((-\infty,A])} = \sup_{x \leqslant -A} |x\Phi(x)| = \sup_{x \geqslant A} x(1 - \Phi(x)) \leqslant e^{-A^2/4}\sqrt{2}e^{-1/2}. \tag{94}$$

As for the derivatives, we have

$$\begin{aligned}
|\text{GELU} - \text{id}|_{W^{1,\infty}([A,+\infty))} &\leqslant \sup_{x \geqslant A}(1 - \Phi(x)) + (\sqrt{2\pi})^{-1} \sup_{x \geqslant A} xe^{-x^2/2} \\
&\leqslant e^{-A^2/2} + (\sqrt{2\pi})^{-1}\sqrt{2}e^{-1/2}e^{-A^2/4} \\
&\leqslant 2e^{-A^2/4}
\end{aligned}$$

and also

$$\begin{aligned}
|\text{GELU}|_{W^{1,\infty}((-\infty,-A])} &\leqslant \sup_{x \leqslant -A} \Phi(x) + (\sqrt{2\pi})^{-1} \sup_{x \leqslant -A} |xe^{-x^2/2}| \\
&\leqslant e^{-A^2/2} + (\sqrt{2\pi})^{-1}\sqrt{2}e^{-1/2}e^{-A^2/4} \\
&\leqslant 2e^{-A^2/4}.
\end{aligned}$$

Now for any natural $k \geqslant 2$ we have from (93) that

$$|\text{GELU}|_{W^{k,\infty}((-\infty,-A]\cup[A,+\infty))} \leqslant (\sqrt{2\pi})^{-1} \left( \sup_{|x| \geqslant A} |e^{-x^2/2}\mathcal{H}_k(x)| + \sup_{|x| \geqslant A} |e^{-x^2/2}\mathcal{H}_{k-2}(x)| \right).$$

Hence, Lemma B.1 implies that

$$|\text{GELU}|_{W^{k,\infty}((-\infty,-A]\cup[A,+\infty))} \leqslant \sqrt{\frac{2}{\pi}}e^{-A^2/4}\sqrt{k!}.$$

Combining all together, we have that

$$\|\text{GELU} - \text{id}\|_{W^{m,\infty}([A,+\infty))} \vee \|\text{GELU}\|_{W^{m,\infty}((-\infty,-A])} \leqslant 2e^{-A^2/4}\sqrt{m!}.$$

The proof is now complete.

$\square$

**Lemma B.3** (De Ryck et al. [2021], Lemma A.6). *Let $d \in \mathbb{N}$, $k \in \mathbb{Z}_+$, $\Omega \subseteq \mathbb{R}^d$ and $f, g \in W^{k,\infty}(\Omega)$. Then it holds that*

$$\|f \cdot g\|_{W^{k,\infty}(\Omega)} \leqslant 2^k \|f\|_{W^{k,\infty}(\Omega)}\|g\|_{W^{k,\infty}(\Omega)}.$$

37

**Lemma B.4** ([De Ryck et al. [2021], Lemma A.7). *Let $d, m, n \in \mathbb{N}$ and let also $\Omega_1 \subseteq \mathbb{R}^d$, $\Omega_2 \subseteq \mathbb{R}^m$, $f \in C^n(\Omega_1, \Omega_2)$ and $g \in C^n(\Omega_2, \mathbb{R})$. Then it holds that*

$$\|g \circ f\|_{W^{n,\infty}(\Omega_1)} \leqslant 16(e^2 n^4 m d^2)^n \|g\|_{W^{n,\infty}(\Omega_2)} \max_{1 \leqslant i \leqslant m} (\|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n \vee 1).$$

*Moreover, if $g \in C^{n+1}(\Omega_2, \mathbb{R})$ and $\widetilde{f} \in C^n(\Omega_1, \Omega_2)$, then*

$$\|g \circ f - g \circ \widetilde{f}\|_{W^{n,\infty}(\Omega_1)}$$
$$\leqslant 32(e^2 n^5 m^2 d^2)^n \|g\|_{W^{n+1,\infty}(\Omega_2)} \max_{1 \leqslant i \leqslant m} \|(f)_i - (\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)} \left(1 \vee \|(f)_i\|_{W^{n,\infty}(\Omega_1)}^{2n} \vee \|(\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)}^{2n}\right).$$

*Proof.* We reprove Lemma A.7 from De Ryck et al. [2021], correcting a minor technical oversight in the original derivation. Specifically, their bound omits $\max(1, \cdot)$ term, which we include here for correctness. We begin with the multivariate Faà di Bruno formula [Constantine and Savits, 1996, Theorem 2.1]. For $\boldsymbol{\nu} \in \mathbb{Z}_+^d$ with $|\boldsymbol{\nu}| = q$ for some $q \in \mathbb{N}$ with $q \leqslant n$ it holds that

$$\partial^{\boldsymbol{\nu}}(g \circ f) = \sum_{1 \leqslant |\boldsymbol{\lambda}| \leqslant q} \partial^{\boldsymbol{\lambda}} g \sum_{p(\boldsymbol{\nu}, \boldsymbol{\lambda})} (\boldsymbol{\nu}!) \prod_{j=1}^q \frac{(f_{l_j})^{k_j}}{k_j!(l_j!)^{k_j}}, \tag{95}$$

where $(f_\mu)_i = \partial^{\boldsymbol{\mu}} f_i$ for $1 \leqslant i \leqslant m$. In addition,

$$p(\boldsymbol{\nu}, \boldsymbol{\lambda}) = \big\{(\mathbf{k}_1, \ldots, \mathbf{k}_q; \mathbf{l}_1, \ldots, \mathbf{l}_q) : \text{for some } 1 \leqslant s \leqslant q,$$
$$\mathbf{k}_i = 0_m \text{ and } \mathbf{l}_i = 0_d \text{ for all } 1 \leqslant i \leqslant q - s; \ |\mathbf{k}_i| > 0 \text{ for all } q - s + 1 \leqslant i \leqslant n;$$
$$\text{and } 0_d \prec \mathbf{l}_{q-s+1} \prec \cdots \prec \mathbf{l}_q \text{ are such that}$$
$$\sum_{i=1}^n \mathbf{k}_i = \boldsymbol{\lambda}, \ \sum_{i=1}^n |\mathbf{k}_i| \mathbf{l}_i = \boldsymbol{\nu}\big\},$$

where we write $\boldsymbol{a} \prec \boldsymbol{b}$ if either $|\boldsymbol{a}| \leqslant |\boldsymbol{b}|$ or $|\boldsymbol{a}| = |\boldsymbol{b}|$ and $a_1 < b_1$ or $|\boldsymbol{a}| = |\boldsymbol{b}|$ and for some $1 \leqslant k \leqslant d - 1$ we have $a_{k+1} < b_{k+1}$ with $a_1 = b_1, \ldots, a_k = b_k$. It is evident that in (95) we have $\sum_{i=1}^n |\mathbf{k}_i| \leqslant n$ and, hence, the number of $(\mathbf{k}_1, \ldots, \mathbf{k}_n)$ satisfying the definition of $p(\boldsymbol{\nu}, \boldsymbol{\lambda})$ is bounded by $|P_{n,nm+1}|$, which is then evaluated as $\sqrt{\pi} e^n (mn)^n$ according to Lemma 2.1 from De Ryck et al. [2021]. Similarly, the number of $(\mathbf{l}_1, \ldots, \mathbf{l}_n)$ satisfying the definition of $p(\boldsymbol{\nu}, \boldsymbol{\lambda})$ is bounded by $|P_{n,dn+1}|$, which in turn, is bounded by $\sqrt{\pi} e^n (dn)^n$. This results in

$$|p(\boldsymbol{\nu}, \boldsymbol{\lambda})| \leqslant \pi (e^2 n^2 m d)^n. \tag{96}$$

Finally, evaluate

$$|\{\boldsymbol{\lambda} \in \mathbb{Z}_+^d \ : \ 1 \leqslant |\boldsymbol{\lambda}| \leqslant q\}| \leqslant |P_{n,d+1}| \leqslant \sqrt{\pi} e^n d^n, \quad |\partial^{\boldsymbol{\lambda}} g| \leqslant \|g\|_{W^{n,\infty}(\Omega_2)}, \quad \boldsymbol{\nu}! \leqslant n! \tag{97}$$

and

$$\prod_{j=1}^n (f_{l_j})^{k_j} \leqslant 1 \vee \max_{1 \leqslant i \leqslant m} \|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n, \tag{98}$$

Therefore, Stirling's approximation implies that

$$\max_{\boldsymbol{\nu}\in\mathbb{Z}_+^d,\,1\leqslant|\boldsymbol{\nu}|\leqslant n}\|\partial^{\boldsymbol{\nu}}g\circ f\|_{W^{0,\infty}(\Omega_1)}\leqslant\sqrt{\pi}(ed)^n\|g\|_{W^{n,\infty}(\Omega_2)}\pi(e^2n^2md)^n n!(1\vee\max_{1\leqslant i\leqslant m}\|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n)$$

$$\leqslant 16(e^2n^4md^2)^n\|g\|_{W^{n,\infty}(\Omega_2)}(1\vee\max_{1\leqslant i\leqslant m}\|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n).$$

For $\boldsymbol{\nu}=0_d$ we have that

$$\|g\circ f\|_{W^{0,\infty}(\Omega_1)}\leqslant\|g\|_{W^{m,\infty}(\Omega_2)}\leqslant 16(e^2n^4md^2)^n\|g\|_{W^{n,\infty}(\Omega_2)}(1\vee\max_{1\leqslant i\leqslant m}\|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n).$$

Hence, the first claim holds true. Now using (95), we deduce that

$$|\partial^{\boldsymbol{\nu}}(g\circ f)-\partial^{\boldsymbol{\nu}}(g\circ\widetilde{f})|$$

$$\leqslant\sum_{1\leqslant|\boldsymbol{\lambda}|\leqslant q}|\partial^{\boldsymbol{\lambda}}[g]\circ f-\partial^{\boldsymbol{\lambda}}[g]\circ\widetilde{f}|\sum_{p(\boldsymbol{\nu},\boldsymbol{\lambda})}(\boldsymbol{\nu}!)\prod_{j=1}^q\frac{(f_{l_j})^{k_j}}{k_j!(l_j!)^{k_j}}$$

$$+\sum_{1\leqslant|\boldsymbol{\lambda}|\leqslant q}|\partial^{\boldsymbol{\lambda}}[g]\circ\widetilde{f}|\sum_{p(\boldsymbol{\nu},\boldsymbol{\lambda})}(\boldsymbol{\nu}!)\frac{|\prod_{j=1}^q(f_{l_j})^{k_j}-\prod_{j=1}^q(\widetilde{f}_{l_j})^{k_j}|}{\prod_{j=1}^q k_j!(l_j!)^{k_j}}.$$

First, bound the first term. From (97) and (98) we find that

$$\sum_{1\leqslant|\boldsymbol{\lambda}|\leqslant q}|\partial^{\boldsymbol{\lambda}}[g]\circ f-\partial^{\boldsymbol{\lambda}}[g]\circ\widetilde{f}|\sum_{p(\boldsymbol{\nu},\boldsymbol{\lambda})}(\boldsymbol{\nu}!)\prod_{j=1}^q\frac{(f_{l_j})^{k_j}}{k_j!(l_j!)^{k_j}}$$

$$\leqslant|\{\boldsymbol{\lambda}\in\mathbb{Z}_+^d\ :\ 1\leqslant|\boldsymbol{\lambda}|\leqslant q\}|\cdot|p(\boldsymbol{\nu},\boldsymbol{\lambda})|\cdot n!\cdot(1\vee\max_{1\leqslant i\leqslant m}\|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n)|\partial^{\boldsymbol{\lambda}}[g]\circ f-\partial^{\boldsymbol{\lambda}}[g]\circ\widetilde{f}|.$$

Now mean value theorem together with (96) and (97) suggests that

$$\sum_{1\leqslant|\boldsymbol{\lambda}|\leqslant q}|\partial^{\boldsymbol{\lambda}}[g]\circ f-\partial^{\boldsymbol{\lambda}}[g]\circ\widetilde{f}|\sum_{p(\boldsymbol{\nu},\boldsymbol{\lambda})}(\boldsymbol{\nu}!)\prod_{j=1}^q\frac{(f_{l_j})^{k_j}}{k_j!(l_j!)^{k_j}}$$

$$\leqslant 16m(e^2n^4md^2)^n(1\vee\max_{1\leqslant i\leqslant m}\|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n)\|g\|_{W^{n+1,\infty}(\Omega_2)}\max_{1\leqslant i\leqslant m}\|(f)_i-(\widetilde{f})_i\|_{W^{0,\infty}(\Omega_1)}.\quad(99)$$

Second, evaluate the second term, using (96) and (97):

$$\sum_{1\leqslant|\boldsymbol{\lambda}|\leqslant q}|\partial^{\boldsymbol{\lambda}}[g]\circ\widetilde{f}|\sum_{p(\boldsymbol{\nu},\boldsymbol{\lambda})}(\boldsymbol{\nu}!)\frac{|\prod_{j=1}^q(f_{l_j})^{k_j}-\prod_{j=1}^q(\widetilde{f}_{l_j})^{k_j}|}{\prod_{j=1}^q k_j!(l_j!)^{k_j}}$$

$$\leqslant\sqrt{\pi}e^nd^n\cdot\pi(e^2n^2md)^n\cdot\|g\|_{W^{n,\infty}(\Omega_2)}n!\cdot\sum_{j=1}^q|(f_{l_j})^{k_j}-(\widetilde{f}_{l_j})^{k_j}|\prod_{u<j}|(f_{l_u})^{k_u}|\prod_{u>j}|(f_{l_u})^{k_u}|.$$

Therefore, Stirling's approximation suggests that

$$\sum_{1\leqslant|\boldsymbol{\lambda}|\leqslant q}|\partial^{\boldsymbol{\lambda}}[g]\circ\widetilde{f}|\sum_{p(\boldsymbol{\nu},\boldsymbol{\lambda})}(\boldsymbol{\nu}!)\frac{|\prod_{j=1}^q(f_{l_j})^{k_j}-\prod_{j=1}^q(\widetilde{f}_{l_j})^{k_j}|}{\prod_{j=1}^q k_j!(l_j!)^{k_j}}$$

$$\leqslant 16(e^2n^4md^2)^n\|g\|_{W^{n,\infty}(\Omega_2)}\sum_{j=1}^q|(f_{l_j})^{k_j}-(\widetilde{f}_{l_j})^{k_j}|\max_{1\leqslant i\leqslant m}(1\vee\|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n\vee\|(\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)}^n).$$

39

We next note that

$$\sum_{j=1}^{q} |(f_{l_j})^{k_j} - (\widetilde{f}_{l_j})^{k_j}| = \sum_{j=1}^{q} \sum_{i=1}^{m} |(f_{l_j})_i^{(k_j)_i} - (\widetilde{f}_{l_j})_i^{(k_j)_i}| \prod_{u<i} |(f_{l_j})_u^{(k_j)_u}| \prod_{u>i} |(\widetilde{f}_{l_j})_u^{(k_j)_u}|$$

$$\leqslant \sum_{j=1}^{q} |k_j| \cdot \max_{1 \leqslant i \leqslant m} \|(f)_i - (\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)} (1 \vee \|(f)_i\|_{W^{n,\infty}(\Omega_1)}^n \vee \|(\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)}^n).$$

Therefore, we obtain that

$$\sum_{1 \leqslant |\boldsymbol{\lambda}| \leqslant q} |\partial^{\boldsymbol{\lambda}}[g] \circ \widetilde{f}| \sum_{p(\boldsymbol{\nu},\boldsymbol{\lambda})} (\boldsymbol{\nu}!) \frac{|\prod_{j=1}^{q}(f_{l_j})^{k_j} - \prod_{j=1}^{q}(\widetilde{f}_{l_j})^{k_j}|}{\prod_{j=1}^{q} k_j!(l_j!)^{k_j}}$$

$$\leqslant 16n(e^2 n^4 m d^2)^n \|g\|_{W^{n,\infty}(\Omega_2)} \max_{1 \leqslant i \leqslant m} \|(f)_i - (\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)} (1 \vee \|f_i\|_{W^{n,\infty}(\Omega_1)}^{2n} \vee \|(\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)}^{2n}).$$

$$\tag{100}$$

For $\boldsymbol{\nu} = 0_d$ it holds that

$$\|g \circ f - g \circ \widetilde{f}\|_{W^{0,\infty}(\Omega_1)} \leqslant m \|g\|_{W^{1,\infty}(\Omega_2)} \max_{1 \leqslant i \leqslant m} \|(f)_i - (\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)}.$$

Thus, from (99) and (100) we conclude that

$$\|g \circ f - g \circ \widetilde{f}\|_{W^{n,\infty}(\Omega_1)}$$
$$\leqslant 32(e^2 n^5 m^2 d^2)^n \|g\|_{W^{n+1,\infty}(\Omega_2)} \max_{1 \leqslant i \leqslant m} \|(f)_i - (\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)} (1 \vee \|(f)_i\|_{W^{n,\infty}(\Omega_1)}^{2n} \vee \|(\widetilde{f})_i\|_{W^{n,\infty}(\Omega_1)}^{2n}).$$

The proof is complete.

$$\square$$

**Lemma B.5** (concatenation of neural networks)**.** *Given $K \in \mathbb{N}$ with $K \geqslant 2$. Then, for any neural networks $\varphi^{(k)} \in \mathsf{NN}(L_k, W_k, S_k, B_k)$ with $1 \leqslant k \leqslant K$ such that $\varphi^{(k)} : \mathbb{R}^{d_k} \to \mathbb{R}^{d_{k+1}}$, there exists a neural network $h = \varphi^{(K)} \circ \varphi^{(K-1)} \cdots \circ \varphi^{(1)} \in \mathsf{NN}(L, W, S, B)$ satisfying*

$$L \leqslant 1 + \sum_{k=1}^{K} (L_k - 1), \quad S \leqslant \sum_{k=1}^{K} S_k + 2 \sum_{k=1}^{K-1} \|W_k\|_\infty \cdot \|W_{k+1}\|_\infty,$$

$$\|W\|_\infty \leqslant \max_{1 \leqslant k \leqslant K} \|W_k\|_\infty, \quad B \leqslant 2 \max_{1 \leqslant k \leqslant K-1} \left[(B_k \vee 1)(B_{k+1} \vee 1)(\|W_k\|_\infty \vee \|W_{k+1}\|_\infty)\right].$$

*Proof.* It suffices to prove the statement for $K = 2$, since one can easily generalize it to $K \geqslant 3$ by induction. Recall that each $\varphi^{(j)}$ for $j \in \{1, 2\}$ admits the representation given in (1). Specifically,

$$\varphi^{(j)}(x) = -b_{L_j}^j + A_{L_j}^j \circ \mathrm{GELU}_{b_{L_j-1}^j} \circ A_{L_j-1}^j \circ \mathrm{GELU}_{b_{L_j-2}^j} \circ \cdots \circ A_2^j \circ \mathrm{GELU}_{b_1^j} \circ A_1^j \circ x.$$

Therefore, we deduce that

$$\varphi^{(2)} \circ \varphi^{(1)} \circ x = -b_{L_2}^2 + A_{L_2}^2 \circ \mathrm{GELU}_{b_{L_2-1}^2} \circ \cdots \circ \mathrm{GELU}_{A_1^2 b_{L_1}^1 + b_1^2} \circ A_1^2 A_{L_1}^1 \circ \cdots \circ \mathrm{GELU}_{b_1^1} \circ A_1^1 \circ x.$$

Consequently, it follows that $\varphi^{(2)} \circ \varphi^{(1)} \in \mathsf{NN}(L, W, S, B)$ with

$$L \leqslant L_1 + L_2 - 1, \quad \|W\|_\infty \leqslant \|W_1\|_\infty \vee \|W_2\|_\infty,$$
$$S \leqslant S_1 + S_2 + 2\|W_1\|_\infty \cdot \|W_2\|_\infty, \quad B \leqslant 2(B_1 \vee 1)(B_2 \vee 1)(\|W_1\|_\infty \vee \|W_2\|_\infty).$$

Hence, the base case holds. The result then follows by induction.

$\square$

**Lemma B.6** (parallelization of neural networks). *Let $K \in \mathbb{N}$ with $K \geqslant 2$ and let neural networks $\varphi^{(k)} \in \mathsf{NN}(L_k, W_k, S_k, B_k)$ for $1 \leqslant k \leqslant K$. Assume further that $L_k = L$ for all $1 \leqslant k \leqslant K$. Then, the following holds:*

*(i) if $\varphi^{(k)} : \mathbb{R}^{d_k} \to \mathbb{R}$ for each $1 \leqslant k \leqslant K$, then there exists a neural network $\varphi \in \mathsf{NN}(L, W, S, B)$ such that $\varphi(x) = (\varphi^{(1)}(x_1), \dots, \varphi^{(K)}(x_K))^\top$ for all $x = (x_1^\top, \dots, x_K^\top)^\top$, where $x_k \in \mathbb{R}^{d_k}$ for any $1 \leqslant k \leqslant K$. In addition, there exists $\varphi_{sum} \in \mathsf{NN}(L, W, S, B_{sum})$, which implements the summation, that is,*

$$\varphi_{sum}(x) = \sum_{k=1}^{K} \varphi^{(k)}(x_k), \quad \text{for all } x = (x_1^\top, \dots, x_K^\top)^\top.$$

*(ii) if $\varphi^{(k)} : \mathbb{R}^p \to \mathbb{R}$ for some $p \in \mathbb{N}$ for every $1 \leqslant k \leqslant K$, then there exists a neural network $\varphi \in \mathsf{NN}(L, W, S, B)$ satisfying $\varphi(x) = (\varphi^{(1)}(x), \dots, \varphi^{(K)}(x))^\top$ for all $x \in \mathbb{R}^p$. Moreover, there exists a summation network $\varphi_{sum} \in \mathsf{NN}(L, W, S, B_{sum})$ such that*

$$\varphi_{sum}(x) = \sum_{k=1}^{K} \varphi^{(k)}(x), \quad \text{for all } x \in \mathbb{R}^p.$$

*Furthermore, in both cases it holds that*

$$\|W\|_\infty \leqslant \sum_{k=1}^{K} \|W^{(k)}\|_\infty, \quad S \leqslant \sum_{k=1}^{K} S^{(k)}, \quad B \leqslant \max_{1 \leqslant k \leqslant K} B^{(k)}, \quad B_{sum} \leqslant K \max_{1 \leqslant k \leqslant K} B^{(k)}.$$

*Proof.* As for the case $(i)$, from (1) we find that $\varphi^{(j)}$ for each $1 \leqslant j \leqslant K$ has the following form:

$$\varphi^{(j)}(x_j) = -b_L^j + A_L^j \circ \mathrm{GELU}_{b_{L-1}^j} \circ A_{L-1}^j \circ \mathrm{GELU}_{b_{L-2}^j} \circ \cdots \circ A_2^j \circ \mathrm{GELU}_{b_1^j} \circ A_1^j \circ x_j.$$

Following Nakada and Imaizumi [2020b], we introduce

$$\widetilde{A}_l = \begin{pmatrix} A_l^1 & 0 & \dots & 0 \\ 0 & A_l^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & A_l^K \end{pmatrix}, \quad \widetilde{b}_l = \begin{pmatrix} b_l^1 \\ \vdots \\ b_l^K \end{pmatrix}, \quad \text{for all } 1 \leqslant l \leqslant L. \tag{101}$$

Hence, as suggested by (1), for

$$\varphi(x) = -\widetilde{b}_L + \widetilde{A}_L \circ \mathrm{GELU}_{\widetilde{b}_{L-1}} \circ \widetilde{A}_{L-1} \circ \mathrm{GELU}_{\widetilde{b}_{L-2}} \circ \cdots \circ \widetilde{A}_2 \circ \mathrm{GELU}_{\widetilde{b}_1} \circ \widetilde{A}_1 \circ x$$

we have that $\varphi(x) = (\varphi^{(1)}(x_1), \ldots, \varphi^{(K)}(x_K))^\top$ for all $x = (x_1^\top, \ldots, x_K^\top)^\top$. Furthermore, the configuration of the network $\varphi$ coincides with that from the statement of the lemma. As for the summation network, we let

$$\varphi_{sum}(x) = -\bar{b}_L + \bar{A}_L \circ \mathrm{GELU}_{\widetilde{b}_{L-1}} \circ \widetilde{A}_{L-1} \circ \mathrm{GELU}_{\widetilde{b}_{L-2}} \circ \cdots \circ \widetilde{A}_2 \circ \mathrm{GELU}_{\widetilde{b}_1} \circ \widetilde{A}_1 \circ x,$$

where

$$\bar{A}_L = \begin{pmatrix} A_L^1 & A_L^2 & \ldots & A_L^K \end{pmatrix}, \quad \bar{b}_L = \sum_{k=1}^K b_L^k. \tag{102}$$

Hence, it follows that

$$\varphi_{sum}(x) = \sum_{k=1}^K \varphi^{(k)}(x_k), \quad \text{for all } x = (x_1^\top, \ldots, x_K^\top)^\top.$$

The configuration of $\varphi_{sum}$ immediately follows from (102). The proof of the case $(ii)$ is identical to the considered one. The only difference is that in (101) for $l = 1$ we define

$$\widetilde{A}_1 = \begin{pmatrix} A_1^{1\top} & A_1^{2\top} & \ldots & A_1^{K\top} \end{pmatrix}^\top, \quad \widetilde{b}_1 = \begin{pmatrix} b_1^{1\top} & b_1^{2\top} & \ldots & b_1^{K\top} \end{pmatrix}^\top.$$

The proof is finished.

$\square$