

SUPN: SHALLOW UNIVERSAL POLYNOMIAL NETWORKS*

ZACHARY MORROW^{†‡}, MICHAEL PENWARDEN^{†‡}, BRIAN CHEN[‡], AURYA JAVEED[‡],
AKIL NARAYAN[§], AND JOHN D. JAKEMAN[‡]

Abstract. Deep neural networks (DNNs) and Kolmogorov–Arnold networks (KANs) are popular methods for function approximation due to their flexibility and expressivity. However, they typically require a large number of trainable parameters to produce a suitable approximation. Beyond making the resulting network less transparent, overparameterization creates a large optimization space, likely producing local minima in training that have quite different generalization errors. In this case, network initialization can have an outsized impact on the model’s out-of-sample accuracy. For these reasons, we propose shallow universal polynomial networks (SUPNs). These networks replace all but the last hidden layer with a single layer of polynomials with learnable coefficients, leveraging the strengths of DNNs and polynomials to achieve sufficient expressivity with far fewer parameters. We prove that SUPNs converge at the same rate as the best polynomial approximation of the same degree, and we derive explicit formulas for quasi-optimal SUPN parameters. We complement theory with an extensive suite of numerical experiments involving SUPNs, DNNs, KANs, and polynomial projection in one, two, and ten dimensions, consisting of over 13,000 trained models. On the target functions we numerically studied, for a given number of trainable parameters, the approximation error and variability are often lower for SUPNs than for DNNs and KANs by an order of magnitude. In our examples, SUPNs even outperform polynomial projection on non-smooth functions.

Key words. Machine learning, approximation theory, neural networks

AMS subject classifications. 41A46, 41A63, 65D15, 65D40, 68T07

1. Introduction. Many scientific applications involve computationally expensive input–output maps, denoted f , resulting from physical systems, such as state-space evolution [60] or the relationship between a model’s parameters and certain quantities of interest [54]. Accordingly, much work has been devoted toward constructing less costly approximations equipped with error estimates or convergence guarantees. One approach is focused on accelerating a numerical solver, such as a finite-element method, by simplifying a model’s governing equations [35] or constructing a reduced-order model [7] that diminishes the number of unknowns. Though this approach is focused on the physical system and not directly on f , such improvements will nonetheless reduce the expense of constructing f . Another approach is to approximate f directly. In practice, this requires sampling its inputs and outputs, and an approximation constructed this way is called a *data-driven surrogate model*. The efficiency of this approach depends on the dimensionality and nonlinearity of f . Many

[†]Z. M. and M. P. contributed equally to this work.

*Submitted to the editors on November 25, 2025.

Funding: The authors gratefully acknowledge funding from Sandia National Laboratories’ Laboratory Directed Research and Development program. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-NA0003525. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government. The publisher acknowledges that the U.S. Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this written work or allow others to do so, for U.S. Government purposes. The DOE will provide public access to results of federally sponsored research in accordance with the DOE Public Access Plan. SAND2025-14696O

[‡]Sandia National Laboratories, 1515 Eubank Blvd SE, Albuquerque, NM 87123 ({zbmorrow, mspenwa, brichen, asjavee, jdjakem}@sandia.gov).

[§]Scientific Computing and Imaging Institute and Department of Mathematics, University of Utah, 155 S 1400 E, Room 233, Salt Lake City, UT 84112 (akil@sci.utah.edu).

methods in this category are understandably focused largely on the high-dimensional case because a system may have a very large number of parameters mapped to some quantity of interest. In this paper, we present a highly expressive method for approximating nonlinear and low-regularity functions, which outperforms state-of-the-art alternatives, including deep neural networks and polynomial approximation.

Several different varieties of data-driven surrogate models exist, and we will purposefully save the technical formulation for a later section. Although an exhaustive list is beyond the scope of this paper, we offer a brief description. A data-driven surrogate of an input–output map can be broadly categorized as linear or nonlinear, depending on how the parameters appear within the surrogate’s ansatz. Linear methods express the resultant surrogate as a linear combination of chosen basis functions. This category includes polynomial chaos expansions [24, 75], sparse grids [57, 5], operator inference [60, 6], spectral methods [11, 27], and radial basis functions [10]. In contrast, the parameters of nonlinear methods are embedded within the surrogate’s constituent functions, which are generally nonlinear. The parameters of a nonlinear surrogate are determined with an iterative optimization procedure. Functional tensor trains [26] and deep neural networks (DNNs) are examples of this category.

Although linear surrogates are well-established methods with a rich body of theory, DNNs have attracted significant contemporary interest due to their flexibility and expressivity, resulting in a proliferation of architectures [16, 29, 64, 49, 72] and bespoke adaptations tailored to specific problems, such as AlphaFold [37]. Nonetheless, the search for a sufficiently accurate DNN is often hampered by vanishing or exploding gradients, as well as an unwieldy number of trainable parameters. Theoretical results establish that feedforward networks, also known as multi-layer perceptrons (MLPs), can approximate continuous functions arbitrarily well [32, 31, 16]. An alternative to the MLP is the Kolmogorov–Arnold network (KAN), based on the Kolmogorov–Arnold representation theorem [3, 42, 9]. In contrast to MLPs, KANs have learnable parameters within their activation functions, which are parametrized as, e.g., splines [19] or Chebyshev polynomials [4, p. 211]. In fact, KANs with a linear spline basis are equivalent to MLPs with a ReLU activation function [2]. Neural networks (NNs) have also been extended to learn discretizations of operators between function spaces. Kernel-based methods such as the Fourier neural operator (FNO) [46] and kernel neural operator (KNO) [50] seek to learn a sequence of kernel integral operators. Deep operator networks (DeepONets) [51] are an alternative approach consisting of two NNs, one learning a set of basis functions and the other learning a set of coefficients.

We avoid the previously discussed issues in deep learning of functions by adopting a fundamentally shallow approach. We introduce a single-layer tanh-activated network with an initial lift consisting of Chebyshev polynomials, which we call a *shallow universal polynomial network* (SUPN). This function-approximation method leverages the power of both polynomial approximation and NNs (via a nonlinear activation function) to produce an approximation that is both parsimonious and robust with respect to network parameter initialization. Our numerical experiments indicate that SUPNs are competitive against DNNs and KANs on smooth target functions, and they are competitive against polynomial projection on nonsmooth targets. Moreover, we observed that SUPNs perform *at least as well* as DNNs or KANs when the target function has tensor-product structure.

Figure 1 shows a visual comparison of SUPNs with other methods. SUPNs differ from MLPs through the Chebyshev lift, and they differ from KANs through the use of fixed activation functions and a single layer. In particular, previous discussions of Chebyshev KANs [66, 65] adopted a deep-learning approach that composed tanh-

activated Chebyshev layers together. In addition, [77] proposed a Chebyshev feature map where the learnable parameters are not basis coefficients but instead generalized “degrees” of Chebyshev polynomials; those networks still use a deep MLP after the Chebyshev lift.

The term *polynomial networks* has historically referred to methods that utilize polynomial activation functions. Despite the similarity in naming conventions, our method has little in common with these prior models, such as polynomial neural networks (PNNs) [18, 38, 30], deep PNNs [40, 13, 43], and sum-product networks [63]. What unites these methods is that they are reformulations or extensions of the so-called Group Method of Data Handling [33], which builds hierarchical polynomial regressions and is equivalent to an MLP with polynomial activation functions. However, it is well-established that MLPs with polynomial activation functions are not universal approximators [45, 67, 62]. In contrast, our proposed method, based on linear combinations of tanh-activated polynomials, *is* a universal approximator.

To motivate the derivation of SUPNs, enable their analysis, and demonstrate their approximation capabilities, the rest of the paper is organized as follows. [Section 2](#) provides a general framework for data-driven surrogate modeling and error estimation, as well as basic formulations of polynomial projection, MLPs, and KANs. [Section 3](#) presents the mathematical formulation of SUPNs and proves two universal approximation theorems. [Section 4](#) discusses the second-order optimizer utilized for training. [Section 5](#) studies several examples comparing SUPNs, polynomial projection, DNNs, and KANs in one, two, and ten dimensions. Existing literature currently lacks a systematic comparison of machine learning methods against polynomial approximations. [Section 6](#) summarizes our findings and discusses possible extensions of this work. For convenience, [Appendix A](#) summarizes the notation used throughout this paper.

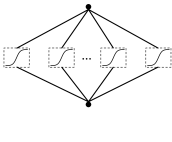
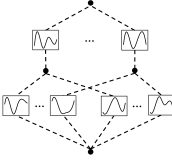
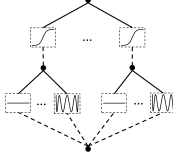
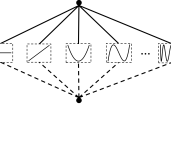


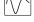

	Multi-layer Perceptron (MLP)	Kolmogorov-Arnold Network (KAN)	Shallow Universal Poly-Net (SUPN)	Polynomial Approximation
Formula	$f(x) = \sum_{i=1}^{L_1} c_i \sigma(w_i \cdot x + b_i)$	$f(x) = \sum_{i=1}^{L_1+1} \phi_i \left(\sum_{j=1}^{L_1} \phi_{i,j}(x) \right)$	$f(x) = \sum_{i=1}^N c_i \cdot \sigma \left(\sum_{j=0}^M a_{i,j} T_j(x) \right)$	$f(x) = \sum_{i=0}^K a_i P_i(x)$
Learnable Terms	$\{c_i, w_i, b_i\}$	$\{\phi_i, \phi_{i,j}\}$	$\{c_i, a_{i,j}\}$	$\{a_i\}$
Model				
<div style="display: flex; justify-content: space-around; align-items: center;"> <div>Learnable Coefficient </div> <div>No Coefficient </div> <div>Learnable Function </div> <div>Fixed Function </div> </div>				

FIG. 1. Illustration of function approximators considered in this work. Note that MLPs and KANs are shown in their single layer form for comparison, but are most often used in a deep (many-layer) configuration compared to SUPNs, which are always a single layer.

2. Surrogate models. Consider a model f mapping D bounded inputs to Q outputs. We can denote this model as $f : \mathbb{R}^D \supset \Omega \rightarrow \mathbb{R}^Q$.¹ The first step of surrogate modeling is to define a space \mathcal{F} of candidate functions, parameterized by an ansatz

¹Without loss of generality, we may take $\Omega = [-1, 1]^D$ since a compact interval $[a, b]$ can be mapped to $[-1, 1]$ with a straightforward affine transformation. We also often take $Q = 1$ for convenience, since the multi-output case constructs a scalar-valued result for each output dimension.

$\bar{g} = \bar{g}(\cdot; \boldsymbol{\theta})$ and $\boldsymbol{\theta} \in \mathbb{R}^P$ such that $\mathcal{F}_{\boldsymbol{\theta}} = \{\bar{g}(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathbb{R}^P\}$. At the highest level, the goal of surrogate modeling is to find a suitable approximation $\hat{f} \in \mathcal{F}_{\boldsymbol{\theta}}$. With $p \geq 1$ and μ a positive measure on Ω , the best approximation of f by $\mathcal{F}_{\boldsymbol{\theta}}$ in the μ -weighted $L^p(\Omega)$ norm is

$$(2.1) \quad \hat{f} := \arg \min_{g \in \mathcal{F}_{\boldsymbol{\theta}}} \|f - g\|_{L^p_{\mu}}.$$

An alternative form of (2.1), which will become useful later, is

$$(2.2) \quad \hat{f} = \bar{g}(\cdot; \hat{\boldsymbol{\theta}}), \quad \hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^P} \|f - \bar{g}(\cdot; \boldsymbol{\theta})\|_{L^p_{\mu}}.$$

The choice of norm and \bar{g} in (2.2) are crucial to the accuracy of \hat{f} . Common choices of \bar{g} are degree- n polynomials (denoted \mathbb{P}_n), sines and cosines with maximum frequency n , or neural networks of a given width and depth. Well-established theory exists for these ansatzes, and we will review some of this theory later in this section.

In practice, surrogates are constructed instead from a finite amount of data because the explicit functional form of f is rarely available. In the supervised learning context, one collects $K \in \mathbb{N}$ input–output pairs $\{(\mathbf{x}_k, f(\mathbf{x}_k))\}_{k \in [K]}$, where $[K] := \{1, \dots, K\}$. Data-driven surrogates approximate the continuous L^p norms with a corresponding discrete ℓ^p norm. The discrete minimization problem becomes

$$(2.3) \quad \begin{aligned} \tilde{f} &= \bar{g}(\cdot; \tilde{\boldsymbol{\theta}}), & \tilde{\boldsymbol{\theta}} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^P} \|\bar{g}(\boldsymbol{\theta}) - \mathbf{f}\|_{\ell^p_{\mathbf{w}}(\mathbb{R}^K)}, \\ \bar{g} &= (\bar{g}(\mathbf{x}_k; \boldsymbol{\theta}))_{k \in [K]}, & \mathbf{f} &= (f(\mathbf{x}_k))_{k \in [K]}, \end{aligned}$$

which is a direct adaptation of the best-approximation problem (2.2). The weights $\mathbf{w} \in \mathbb{R}^K$ are chosen to approximate the integral in the L^p norm using a quadrature rule, for example $w_k = 1/K$.

The final step is to minimize (2.3). Popular optimization methods include grid search and gradient-based optimizers. Grid searches discretize the feasible set for $\boldsymbol{\theta}$ and take the discrete minimum. However, grid searches often do not provide optimality guarantees, and P is typically quite large, so the number of samples required to cover the sample space (e.g., with low discrepancy) is practically prohibitive. Alternatively, when using gradient-based optimizers, the ℓ^p should be differentiable. As a result, practical implementation uses (2.3) with the ℓ^2 norm, yielding the familiar loss-minimization problem with mean-squared error (MSE):

$$(2.4) \quad \tilde{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^P} \mathcal{L}(\boldsymbol{\theta}), \quad \mathcal{L}(\boldsymbol{\theta}) = \sum_{k \in [K]} w_k (\bar{g}(\mathbf{x}_k; \boldsymbol{\theta}) - f(\mathbf{x}_k))^2.$$

In the special case where \bar{g} depends linearly on $\boldsymbol{\theta}$, the minimizer $\tilde{\boldsymbol{\theta}}$ is the solution to a linear system, which is significantly cheaper to solve than the nonlinear case.

2.1. Ansatzes. We now present three prominent choices of the ansatz \bar{g} . These ansatzes have associated theory that will be crucial in the theoretical analysis of section 3, and they will be used for comparison in section 5.

2.1.1. Polynomials. Consider how to approximate a function $f : \Omega \rightarrow \mathbb{R}$ with polynomials in some norm. For this paper, there are two important cases to consider: the L^∞ norm and the L^2 norm. The first case is addressed in the following theorems, which will be needed in the proofs of section 3.

THEOREM 2.1 (Stone–Weierstrass [15, 69, 74]). *Let $f : \Omega \rightarrow \mathbb{R}$ be continuous and $\Omega \subset \mathbb{R}^D$ be compact. For any $\epsilon > 0$, there exists a polynomial $q : \Omega \rightarrow \mathbb{R}$ such that $\|f - q\|_{L^\infty} < \epsilon$.*

THEOREM 2.2 (Jackson’s inequality [34]). *Let $f \in C^k([-1, 1], \mathbb{R})$ with $f^{(k)}$ Lipschitz continuous. There exists a constant $\beta_k > 0$ such that, for every $n \geq 1$, there is a polynomial $q_n \in \mathbb{P}_n$ for which*

$$\|f - q_n\|_{L^\infty} \leq \beta_k n^{-(k+1)}.$$

Now consider a polynomial approximation of degree at most P in the unweighted L^2 norm. For simplicity, we address $D = 1$. This corresponds to (2.2) with the L^2 norm and

$$(2.5) \quad \bar{g}_P(x; \boldsymbol{\theta}) = \sum_{p \in [P-1]_0} \theta_p L_p(x),$$

where $[P-1]_0 := \{0, \dots, P-1\}$ and L_p is the Legendre polynomial of degree p . The Legendre polynomials satisfy

$$\langle L_n, L_m \rangle = \delta_{nm} \frac{2}{2n+1}.$$

where

$$\langle f, g \rangle_{L^2_\mu(\Omega)} := \int_{\Omega} f(x)g(x) d\mu(x), \quad \|f\|_{L^2_\mu(\Omega)}^2 := \langle f, f \rangle_{L^2_\mu(\Omega)}.$$

In this case, the coefficients of the best approximation are

$$\hat{\theta}_p = \frac{\langle L_p, f \rangle}{\langle L_p, L_p \rangle} = \langle f, L_p \rangle \frac{2p+1}{2},$$

and we have

$$\lim_{P \rightarrow \infty} \|f - \bar{g}_P\|_{L^2} = 0.$$

The data-driven setting approximates $\langle f, \phi_p \rangle$ with a quadrature rule

$$(2.6) \quad \tilde{\theta}_p = \sum_{k \in [K]} w_k f(x_k) L_p(x_k).$$

Projection in higher dimensions. It will later become necessary to consider $D \geq 2$. In that case, the notion of “degree” must be generalized to allow for differing amounts of cross-interaction terms among the input space. This is accomplished by replacing the sum over $p \in [P-1]_0$ with a sum over $\mathbf{p} \in \Lambda \subset \mathbb{N}_0^D$, where Λ is a lower set and the multi-dimensional Legendre polynomials are products of the univariate ones:

$$L_{\mathbf{p}}(\mathbf{x}) = \prod_{d=1}^D L_{p_d}(x_d).$$

DEFINITION 2.3 (Lower set [12]). *A set $\Lambda \subset \mathbb{N}_0^D$ is called lower (or downward-closed) if and only if, for each $\mathbf{i} \in \Lambda$, $\{\mathbf{j} : \mathbf{j} \leq \mathbf{i}\} \subset \Lambda$. Here, $\mathbf{j} \leq \mathbf{i}$ if and only if $j_d \leq i_d$ for all $d \in [D]$.*

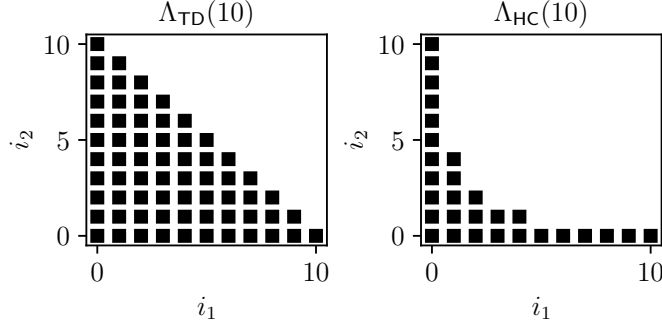


FIG. 2. Two examples of lower sets, the total-degree space (left) and hyperbolic cross-section (right).

Example 2.4. Two common examples of lower sets are the hyperbolic-cross and total-degree spaces, shown in Figure 2 and given by

$$(2.7) \quad \Lambda_{HC}(M) = \left\{ \mathbf{i} \in \mathbb{N}_0^D : \prod_{d=1}^D (i_d + 1) \leq M + 1 \right\},$$

$$(2.8) \quad \Lambda_{TD}(M) = \left\{ \mathbf{i} \in \mathbb{N}_0^D : \sum_{d=1}^D i_d \leq M \right\},$$

respectively.

2.1.2. Multi-layer perceptrons. A multi-layer perceptron (MLP) is given by

$$(2.9) \quad \begin{cases} \tilde{f}(\mathbf{x}) &= \mathbf{W}_L \mathbf{y}_L(\mathbf{x}) \\ \mathbf{y}_{k+1}(\mathbf{x}) &= \sigma(\mathbf{W}_k \mathbf{y}_k(\mathbf{x}) + \mathbf{b}_k), \quad k \in [L-1] \\ \mathbf{y}_1(\mathbf{x}) &= \sigma(\mathbf{W}_0 \mathbf{x} + \mathbf{b}_0) \end{cases}.$$

Here, L is the network depth, $N = \dim(\mathbf{b}_k)$ is the uniform network width, and $\sigma(x)$ is a non-polynomial activation function applied component-wise. The network parameters $\boldsymbol{\theta}$ are the concatenation of the entries of each \mathbf{W}_k and \mathbf{b}_k . By counting the size of each \mathbf{W}_k and \mathbf{b}_k , the total number of parameters is $P = N(D+2) + (L-1)(N^2 + N)$.

We now state two universal approximation theorems for MLPs that will be foundational for the proofs in section 3.

THEOREM 2.5 (Shallow, wide MLPs [62]). *Let both $f : \Omega \rightarrow \mathbb{R}^Q$ and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be continuous. Let the network depth be fixed at $L = 1$. Then σ is non-polynomial if and only if, for any $\epsilon > 0$, there exists $N \in \mathbb{N}$, $\mathbf{W}_0 \in \mathbb{R}^{N \times D}$, $\mathbf{b}_0 \in \mathbb{R}^N$, and $\mathbf{W}_1 \in \mathbb{R}^{Q \times N}$ such that \tilde{f} in (2.9) satisfies*

$$\sup_{\mathbf{x} \in \Omega} \|f(\mathbf{x}) - \tilde{f}(\mathbf{x})\|_{\ell^\infty(\mathbb{R}^Q)} < \epsilon$$

where $\sigma(\cdot)$ is applied componentwise.

THEOREM 2.6 (Deep, narrow MLPs [39]). *Let $f : \Omega \rightarrow \mathbb{R}^Q$ be continuous, and let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be non-affine. Let the network width N be fixed such that $N \geq D + Q + 2$.*

Then there exists $L \in \mathbb{N}$, weights $\{\mathbf{W}_k\}_{k \in [L]_0}$, and biases $\{\mathbf{b}_k\}_{k \in [L-1]_0}$ such that $\tilde{f}(\mathbf{x})$ in (2.9) satisfies

$$\sup_{\mathbf{x} \in \Omega} \|f(\mathbf{x}) - \tilde{f}(\mathbf{x})\|_{\ell^\infty(\mathbb{R}^Q)} < \epsilon.$$

2.1.3. Kolmogorov–Arnold networks. A deep KAN [49, Eq. 2.10–2.12] has the form

$$(2.10) \quad \tilde{f}(\mathbf{x}) = \Phi_L \circ \Phi_{L-1} \circ \cdots \circ \Phi_0(\mathbf{x})$$

where the KAN layer $\Phi_\ell : \mathbb{R}^{N_\ell} \rightarrow \mathbb{R}^{N_{\ell+1}}$ is defined as

$$(2.11) \quad (\Phi_\ell(\mathbf{z}))_k = \mathbf{W}_\ell \text{SiLU}(\mathbf{z}) + \sum_{j \in [N_\ell]} \phi_{\ell,k,j}^{(p)}(z_j),$$

where $\ell \in [L]_0$, $k \in [N_{\ell+1}]$, and $\mathbf{W}_\ell \in \mathbb{R}^{N_{\ell+1} \times N_\ell}$. For a given spline order p , the functions ϕ are univariate B-splines of the form

$$\phi^{(p)}(x) = \sum_{n \in [G]} \alpha_n B_{n,p}(x),$$

where G is the size of the grid partition $\{x_n\}_{n \in [G]}$. The splines themselves are defined recursively via

$$B_{n,0} := \begin{cases} 1, & x_n \leq x < x_{n+1} \\ 0, & \text{otherwise} \end{cases}$$

and, for $p > 0$,

$$B_{n,p}(x) := \frac{x - x_n}{x_{n+p} - x_n} B_{n,p-1}(x) + \frac{x_{n+p+1} - x}{x_{n+p+1} - x_{n+1}} B_{n+1,p-1}(x).$$

The network parameters $\boldsymbol{\theta}$ are the concatenation of \mathbf{W}_ℓ and all B-spline coefficients over all KAN layers.

2.2. Empirical accuracy estimation. To evaluate the effectiveness a surrogate model, we identify some useful metrics such as *best approximation error*, *sampling error*, and *training robustness*. We will utilize these metrics in section 5. For clarity of presentation, we take $D = 1$ here.

2.2.1. Best approximation error. This error is determined purely by the choice of $\boldsymbol{\theta} \in \mathbb{R}^P$. The error under consideration is

$$(2.12) \quad \epsilon(\hat{\boldsymbol{\theta}}) := \inf_{g \in \mathcal{F}} \|f - g\|_{L_\mu^p} = \|f - \hat{f}\|_{L_\mu^p},$$

where \hat{f} is defined in (2.2).

2.2.2. Sampling error. This error is due to the replacement of the continuous L^2 norm in (2.2) with an empirical norm. Using this norm, we seek to minimize the empirical training loss (2.4). The corresponding generalization error is given by

$$\epsilon(\tilde{\boldsymbol{\theta}}) = \|f - \tilde{f}\|_{L_\mu^2}$$

where \tilde{f} is given by (2.3). The meaningful indicator about the effect of sampling would be a comparison between $\epsilon(\hat{\theta})$ and $\epsilon(\tilde{\theta})$. Fixing K , we examine three potential choices of grids on $[-1, 1]$, where $k \in [K]$:

- *Uniform sampling*:

$$x_k \stackrel{\text{iid}}{\sim} \mathcal{U}([-1, 1]), \quad w_k = 2/K$$

- *Equidistant sampling*:

$$x_k = -1 + \frac{2(k-1)}{K-1}, \quad w_k = 2/K$$

- *Gauss-Legendre quadrature* [4, p. 276]:

$$\{x_k\}_{k \in [K]} = L_K^{-1}(\{0\}), \quad w_k = \frac{-2}{(K+1)L_{K+1}(x_k)L'_K(x_k)}$$

2.2.3. Training robustness. When numerically optimizing the parameters of a nonlinear surrogate, typically we can only discover various local minima of (2.4) instead of a global minimum. However, different local minima of the training loss can result in dramatically different generalization error. Much attention has been given to the proper initialization of network weights for desirable generalization error, e.g. through explicit probability distributions [17, 22, 44, 53, 52, 25, 28] or *meta-learned* initializations [21, 56, 48, 61]. We examine training robustness by randomly initializing the network many times and reporting the variance of the test error:

$$\text{Var}_{\theta_{\text{init}} \sim \mu} [\epsilon(\tilde{\theta})]$$

where μ is the Kaiming uniform distribution [28].

3. Shallow universal polynomial networks (SUPNs). For one input dimension, we define a *shallow universal polynomial network* as any function of the form

$$(3.1) \quad f_{N,M}(x) := \sum_{n \in [N]} c_n \tanh \left(\sum_{m \in [M-1]_0} a_{n,m} T_m(x) \right), \quad x \in [-1, 1],$$

where c_n , and $a_{n,m}$ are trainable parameters and $T_m(x)$ is the degree- m Chebyshev polynomial. Intuitively, a SUPN replaces the first $L-1$ layers of a deep MLP with a Chebyshev polynomial, yielding $P = N(M+1)$. This substantially decreases the number of trainable network parameters, which accelerates training and improves interpretability of the neural network. For our purposes, T_m may be computed either trigonometrically or via the three-term recurrence relation [4, p. 211]. In principle, any polynomial basis on Ω could be used in (3.1), but we observed that the Chebyshev basis performs well in practice. We hypothesize that the special approximation properties of Chebyshev polynomials may be responsible for this performance. For instance, among all degree- n polynomials with unit L^∞ norm, T_n has largest leading coefficient, and each local extremum is ± 1 , which provides stability when learning $a_{n,m}$. We emphasize that SUPNs differ from Chebyshev KANs through the absence of repeated compositions.

In higher dimensions, we define a multi-dimensional SUPN as

$$(3.2) \quad f_{N,\Lambda}(\mathbf{x}) := \sum_{n \in [N]} c_n \tanh \left(\sum_{\mathbf{m} \in \Lambda} a_{n,\mathbf{m}} T_{\mathbf{m}}(\mathbf{x}) \right)$$

where $\Lambda \subset \mathbb{N}_0^D$ is a lower set and $T_{\mathbf{m}}(\mathbf{x}) = \prod_{d=1}^D T_{m_d}(x_d)$, as in [subsection 2.1.1](#). In this multi-dimensional setting, the number of trainable parameters for a SUPN is $P = N|\Lambda|$. Multi-dimensional SUPNs are able to model cross-interactions between the inputs with a single layer. In contrast, KANs require compositions to model interaction terms as each layer [\(2.11\)](#) only has *sums* of univariate functions. In this paper, we only use isotropic Λ , but we plan to incorporate anisotropy in Λ as future work, which would further increase the scalability of SUPNs to higher dimensions.

3.1. Universal approximation of SUPNs. We present two theorems for the universal approximation properties of [\(3.1\)](#), which rely on the fixed-width and fixed-depth universal approximation theorems for MLPs. We emphasize that these theorems do not construct the network attaining a given ϵ -accuracy, but they are necessary to establish the suitability of [\(3.1\)](#) as an approximation basis.

THEOREM 3.1 (fixed M , variable N). *Let $f : \Omega \rightarrow \mathbb{R}^Q$ be continuous with $\Omega \subset \mathbb{R}^D$ compact. Fix $M \geq D + 1$. For any $\epsilon > 0$, there exists $N \in \mathbb{N}$ such that $\|f - f_{N,M}\|_{L^\infty} < \epsilon$.*

PROOF 3.1. *Let $\epsilon > 0$ be arbitrary. By [Theorem 2.5](#), there exists $N \in \mathbb{N}$, $\mathbf{W}_0 \in \mathbb{R}^{N \times d}$, $\mathbf{b}_0 \in \mathbb{R}^N$, and $\mathbf{W}_1 \in \mathbb{R}^{D \times N}$ such that $\|f - \tilde{f}_N\|_{L^\infty} < \epsilon$, where*

$$\tilde{f}_N(\mathbf{x}) = \mathbf{W}_1 \tanh(\mathbf{W}_0 \mathbf{x} + \mathbf{b}),$$

with $\tanh(\cdot)$ applied componentwise. Each entry of $\mathbf{W}_0 \mathbf{x} + \mathbf{b}$ is a linear polynomial and can be exactly represented by the inner sum of [\(3.2\)](#) with $\{1, x_1, \dots, x_d\}$. \square

THEOREM 3.2 (fixed N , variable M). *Let $f : \Omega \rightarrow \mathbb{R}^Q$ be continuous, with $\Omega \subset \mathbb{R}^D$ compact. Fix $N \geq D + Q + 2$. For any $\epsilon > 0$, there exists $M \in \mathbb{N}$ such that $\|f - f_{N,M}\|_{L^\infty} < \epsilon$.*

PROOF 3.2. *Fix $N \geq D + Q + 2$, and let $\epsilon > 0$. By [Theorem 2.6](#), there exists $L \in \mathbb{N}$, $\{\mathbf{W}_k\}_{k \in [L]_0}$, and $\{\mathbf{b}_k\}_{k \in [L-1]_0}$ such that*

$$\sup_{\mathbf{x} \in \Omega} \|f(\mathbf{x}) - \tilde{f}_L(\mathbf{x})\|_{\ell^\infty(\mathbb{R}^Q)} < \epsilon/2,$$

where \tilde{f}_L is given by [\(2.9\)](#). Let $\tilde{f}_{L-1}(\mathbf{x})$ denote the output of layer, i.e. $\tilde{f}_L(\mathbf{x}) = \mathbf{W}_L \tanh(\tilde{f}_{L-1}(\mathbf{x}))$. By [Theorem 2.1](#) and the uniform continuity of \tanh and \tilde{f}_{L-1} , there exists a polynomial $p \in \mathbb{P}_\Lambda$ such that, for all $\mathbf{x} \in \Omega$,

$$\|f_{N,\Lambda}(\mathbf{x}) - \tilde{f}_L(\mathbf{x})\|_{\ell^\infty(\mathbb{R}^Q)} < \epsilon/2,$$

where $f_{N,\Lambda}(\mathbf{x}) = \mathbf{W}_L \tanh(p(\mathbf{x}))$. The triangle inequality completes the proof. \square

Remark 3.3. In principle, any non-polynomial activation function $\sigma(x)$ would suffice in [Theorem 3.1](#) and [Theorem 3.2](#). However, we observed numerical instabilities when $\sigma(x)$ is not bounded, such as the ReLU activation. Furthermore, in [\[20\]](#) the authors show networks with two \tanh hidden layers are as expressive as deeper ReLU networks, providing motivation for their use in a shallow formulation here.

3.2. Approximation properties inherited from polynomials. In this section, we present several theorems that bound the best-approximation error of SUPNs by the best-approximation error of polynomials. These theorems relax the assumptions of Subsection 3.1 and have considerably stronger accuracy guarantees. Proposition 3.4 shows that the best polynomial approximation of f in L_μ^2 can be approximated uniformly by SUPNs. Its proof motivates the use of $\sigma(x) = \tanh(x)$ as our activation function. Corollary 3.5 asserts that if a P -term Chebyshev series achieves a given accuracy, then there is a P -term SUPN achieves essentially the same accuracy. Hence, the approximation power of SUPNs can at least match polynomials. Theorem 3.6 and Theorem 3.7 derive explicit expressions for a quasi-optimal SUPN in the L_μ^2 and L^∞ norms. Finally, Theorem 3.8 obtains a best-approximation convergence rate that is polynomially fast in the smoothness of f .

PROPOSITION 3.4 (Polynomial approximation proximity). *Let $f \in W_\mu^{k,p}(\Omega; \mathbb{R})$, with μ a probability measure, $\Omega \subset \mathbb{R}^D$ compact, $k \in \mathbb{N}_0$, and $1 \leq p \leq \infty$. For $\Lambda \in \mathbb{N}_0^D$ and with $\mathbb{P}_\Lambda := \text{span}_{\mathbf{m} \in \Lambda} T_{\mathbf{m}}$, define*

$$\epsilon_\Lambda(f) := \inf_{q \in \mathbb{P}_\Lambda} \|f - q\|_{W_\mu^{k,p}}.$$

For any $\delta > 0$, there is a polynomial $\tilde{q} \in \mathbb{P}_\Lambda$ and a SUPN of the form (3.2) with $P = |\Lambda| + 1$ parameters such that

$$(3.3) \quad \|f - \tilde{q}\|_{W_\mu^{k,p}} < \begin{cases} (1 + \delta)\epsilon_\Lambda(f), & \text{if } \epsilon_\Lambda(f) > 0 \\ \delta, & \text{if } \epsilon_\Lambda(f) = 0 \end{cases} \quad \text{and}$$

$$(3.4) \quad \|f_{N,\Lambda} - \tilde{q}\|_{W^{k,\infty}} < \begin{cases} \delta \epsilon_\Lambda(f), & \text{if } \epsilon_\Lambda(f) > 0 \\ \delta, & \text{if } \epsilon_\Lambda(f) = 0 \end{cases}.$$

PROOF 3.4. *We deal explicitly with $k = 0$, $\epsilon_\Lambda(f) > 0$. An analogous proof holds for $\epsilon_\Lambda(f) = 0$. We give the proof for $k > 0$ in the Supplementary Materials. By definition of $\epsilon_\Lambda(f)$, there is some polynomial $\tilde{q} \in \mathbb{P}_\Lambda$ satisfying (3.3). Let $\tilde{q} = \sum_{\mathbf{m} \in \Lambda} \alpha_{\mathbf{m}} T_{\mathbf{m}}$, and define $R := \sum_{\mathbf{m} \in \Lambda} |\alpha_{\mathbf{m}}|$. If $R = 0$, set $c_1 = 0$ in (3.2) with $N = 1$ to achieve $f_{N,\Lambda} = q = 0$, achieving (3.4). Otherwise, consider $R > 0$. Note that,*

$$\sup_{\mathbf{x} \in \Omega} |\tilde{q}(\mathbf{x})| \leq \sup_{\mathbf{x} \in \Omega} \sum_{\mathbf{m} \in \Lambda} |\alpha_{\mathbf{m}}| |T_{\mathbf{m}}(\mathbf{x})| \leq R.$$

We will construct a SUPN with $(N, M) = (1, |\Lambda|)$ that well-approximates \tilde{q} . Let $\sigma(y) = \tanh(y)$. Around $y = 0$, $\sigma(y) \approx y$, and using $|\sigma''(y)| < 2|y|$, Taylor's Theorem with the integral form of the remainder gives, for any $r > 0$,

$$|\sigma(y) - y| < r^3, \quad |y| < r.$$

Now consider the $(N, M) = (1, |\Lambda|)$ SUPN in (3.2), with the parameter assignment

$$a_{1,\mathbf{m}} = \frac{\alpha_{\mathbf{m}}}{S}, \quad c_1 = S, \quad S := \sqrt{\frac{R^3}{\delta \epsilon_\Lambda(f)}}.$$

Then, for any $\mathbf{x} \in \Omega$,

$$|f_{N,\Lambda}(\mathbf{x}) - \tilde{q}(\mathbf{x})| = S \left| \sigma\left(\frac{\tilde{q}(\mathbf{x})}{S}\right) - \frac{\tilde{q}(\mathbf{x})}{S} \right| \leq S \left(\frac{R}{S}\right)^3 = \delta \epsilon_\Lambda(f).$$

The $k > 0$ proof introduces no significantly new ideas but is more technical; see *Supplementary Materials*. \square

The next result uses the L^∞ (or $W^{k,\infty}$) approximability by SUPNs of the best-approximating polynomial to establish explicit error bounds of $f_{N,\Lambda}$ in the $W_\mu^{k,p}$ norm.

COROLLARY 3.5. *Under assumptions of [Proposition 3.4](#), for any $\delta > 0$, there is a SUPN of the form [\(3.2\)](#) with $P = |\Lambda| + 1$ parameters such that*

$$(3.5) \quad \|f - f_{N,\Lambda}\|_{W_\mu^{k,p}} < \begin{cases} (1 + \delta)\epsilon_\Lambda(f), & \text{if } \epsilon_\Lambda(f) > 0 \\ \delta, & \text{if } \epsilon_\Lambda(f) = 0 \end{cases}.$$

In contrast to the generality of approximability in Sobolev spaces we have shown above, if we specialize to $k = 0$ and $p = 2$, we can construct the near-optimal SUPN weights under the L_μ^2 norm. This can be done by explicitly using the argumentation of [Proof 3.4](#).

THEOREM 3.6 (Constructive SUPN in L_μ^2 norm). *The SUPN guaranteed under [Proposition 3.4](#) with $k = 0$ and $p = 2$ is given by*

$$a_{1,\mathbf{m}} = \frac{\tilde{\alpha}_{\mathbf{m}}}{S}, \quad c_1 = S, \quad S := \begin{cases} \sqrt{\frac{R^3}{\delta}} & \text{if } \epsilon_\Lambda(f) = 0 \\ \sqrt{\frac{R^3}{\delta \epsilon_\Lambda(f)}} & \text{if } \epsilon_\Lambda(f) > 0 \end{cases}$$

where

$$\alpha_{\mathbf{m}} = \frac{\langle \phi_{\mathbf{m}}, f \rangle_{L_\mu^2}}{\langle \phi_{\mathbf{m}}, \phi_{\mathbf{m}} \rangle_{L_\mu^2}}, \quad R = \sum_{\mathbf{m} \in \Lambda} |\alpha_{\mathbf{m}}|,$$

$$\langle \phi_{\mathbf{m}}, \phi_{\mathbf{n}} \rangle_{L_\mu^2} = \delta_{\mathbf{m},\mathbf{n}} \|\phi_{\mathbf{n}}\|_{L_\mu^2}^2, \quad \phi_{\mathbf{m}} \in \mathbb{P}_\Lambda,$$

and $\{\tilde{\alpha}_{\mathbf{m}}\}_{\mathbf{m} \in \Lambda}$ are the transformed coefficients of the polynomial $\phi_{\mathbf{m}}$ to the Chebyshev polynomial $T_{\mathbf{m}}$. Furthermore, the error $\epsilon_\Lambda(f)$ can be computed via

$$(3.6) \quad \epsilon_\Lambda(f) = \sqrt{\|f\|_{L_\mu^2}^2 - \sum_{\mathbf{m} \in \Lambda} |\alpha_{\mathbf{m}}|^2}$$

PROOF 3.6. Apply [Proposition 3.4](#) with the optimal L_μ^2 projection coefficients. \square

Now, the optimal L_μ^2 projection onto Chebyshev polynomials with the measure $d\mu = dx/\sqrt{1-x^2}$ is nearly optimal under the L^∞ norm [\[23\]](#). Hence, we can explicitly construct the network weights guaranteed by the universal approximation theorem ([Theorem 3.2](#)), and we can even relax its assumptions on N .

THEOREM 3.7 (Constructive SUPN in L^∞ norm). *Let $f \in C(\Omega)$ with $\Omega \subset \mathbb{R}$ compact. For any $\delta > 0$, we can construct a SUPN [\(3.1\)](#) with $P = M + 1$ parameters satisfying*

$$\|f - f_{1,M}\|_{L^\infty} < (2 \ln M + 3) \inf_{p \in \mathbb{P}_M} \|f - p\|_{L^\infty} + \delta$$

using weights from [Theorem 3.6](#) with $S := \sqrt{R^3/\delta}$ and μ being the Chebyshev measure.

PROOF 3.7. As shown in [\[23\]](#), the Chebyshev series

$$\tilde{f} = \sum_{m \in [M]_0} \alpha_m T_m(x), \quad \alpha_m = \frac{\langle T_m, f \rangle_\mu}{\langle T_m, T_m \rangle_\mu}, \quad d\mu(x) = \frac{dx}{\sqrt{1-x^2}}$$

is nearly optimal as a degree- M minimax polynomial, with a multiplicative term arising from the Lebesgue constant λ_M , yielding

$$\|f - \tilde{f}\|_{L^\infty} < (1 + \lambda_M) \inf_{p \in \mathbb{P}_M} \|f - p\|_{L^\infty}.$$

Furthermore, [23] demonstrates the bound $\lambda_M \leq 2 \ln M + 2$. Using $S := \sqrt{R^3/\delta}$ and the weights from [Theorem 3.6](#) within [Proof 3.4](#) yields

$$\|\tilde{f} - f_{1,M}\|_{L^\infty} < \delta,$$

which completes the proof. \square

Finally, we can use well-established bounds on polynomial approximation error in terms of the target function’s smoothness class to get an explicit convergence rate of SUPNs.

THEOREM 3.8 (SUPN L^∞ convergence rate). *Let $f \in C^k(\Omega)$ with $\Omega \subset \mathbb{R}$ compact, $f^{(k)}$ Lipschitz continuous, and $\epsilon_{M,\infty}(f) > 0$. Then there exists $\beta > 0$ such that, for any $M \in \mathbb{N}$ and $\delta > 0$, there is a SUPN of the form [\(3.1\)](#) satisfying*

$$(3.7) \quad \|f - f_{1,M}\|_{L^\infty} < (1 + \delta)\beta M^{-(k+1)}.$$

PROOF 3.8. *The result is immediate from [Proposition 3.4](#) and [Theorem 2.2](#). \square*

4. Second-order optimizer. The use of second-order optimization methods for training SUPNs is inspired by the similar mathematical structure of that problem and classical inverse problems where such methods have worked well. Second-order optimization methods approximate a Newton step on the optimality condition $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = 0$, as opposed to approximating the direction of steepest descent, as is done in first-order methods. Empirical advantages gained by the use of second-order methods over first-order methods include robustness (i.e., the performance of these methods requires less hyperparameter tuning), the ability to navigate past saddle points and flat regions, and more rapid rates of convergence. Second-order methods have also been explored in the context of training deep neural networks; an empirical study of these methods can be found in, e.g., [76]. In our experiments, we use a fixed batch of data, so there is no stochasticity in the calculation of the gradients and Hessians (e.g., as a result of minibatching). However, trust region methods can also be applied in the case of stochastic gradients or Hessians; the complexity and convergence of trust region methods in this setting has been explored in [8, 59].

In the case of training SUPNs, DNNs, and KANs, which are generally nonconvex, the Hessians of these networks with respect to their parameters are generally not positive-definite. Trust-region methods are able to handle indefinite Hessian estimates. At every iteration n of a trust-region method, we approximate the loss function $\mathcal{L}(\boldsymbol{\theta})$ locally near the current parameters $\boldsymbol{\theta}_n$ using a model, which is often quadratic:

$$(4.1) \quad \mathcal{L}(\boldsymbol{\theta}_n + \mathbf{s}) \approx \mathcal{L}(\boldsymbol{\theta}_n) + (\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_n))^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H} \mathbf{s},$$

where \mathbf{H} is either the Hessian $\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta}_n)$ or an approximation thereof. This model is (approximately) minimized with the constraint that $\|\mathbf{s}\|$ is less than some trust radius Δ_n :

$$(4.2) \quad \left. \begin{array}{ll} \min_{\mathbf{s} \in \mathbb{R}^P} & (\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_n))^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H} \mathbf{s} \\ \text{s.t.} & \|\mathbf{s}\| < \Delta_n \end{array} \right\}.$$

Note that we have dropped $\mathcal{L}(\theta_n)$ as it does not depend on the step \mathbf{s} . Letting \mathbf{s}_n be a solution (or approximate solution) of the above optimization problem (called the “subproblem”), the proposed iterate $\theta_n + \mathbf{s}_n$ is either accepted or rejected. Note that an exact subproblem solution is not required; finding an \mathbf{s} that satisfies a fraction of Cauchy decrease is enough (i.e., a fraction of the decrease attained by the minimizer in the direction of the negative gradient). The trust-region radius is adjusted from one iteration to the next by assessing the accuracy of the model. Additional details and specific implementations can be found in, e.g., [58, 14].

For our experiments, we use a trust region method, where the subproblem solver is the Steihaug–Toint conjugate gradient (CG) method [70, 68]. We accelerate this optimization method using a limited-memory BFGS Hessian approximation [47] as a preconditioner. This optimization algorithm is implemented using the Rapid Optimization Library [36]. This method does not require fully forming a Hessian or computing its inverse; it only requires Hessian–vector products, which can be efficiently calculated with automatic differentiation.

Example 4.1. In Figure 3, we show training losses for the Adam [41], BFGS [47], and second-order trust-region methods after an initial 1000 epochs of burn-in with Adam (not shown). Adam and BFGS use a learning rate of 10^{-2} and are implemented in PyTorch. The target function to learn is $f(x) = |x|$. Both DNNs and SUPNs converge attain a lower error in far fewer epochs with Newton–CG than with Adam or BFGS. This demonstrates the utility of using a second-order method to determine best-approximation error.

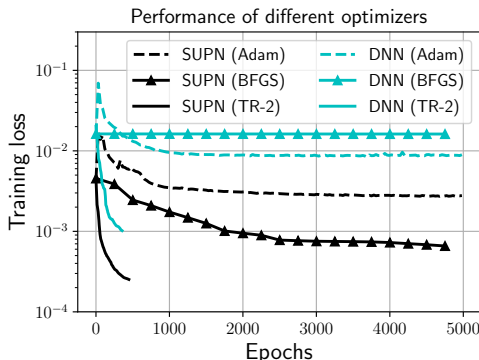


FIG. 3. Training losses with Adam, BFGS, and a second-order trust-region method.

5. Numerical experiments. We present several experiments in one, two, and ten dimensions to demonstrate the utility of SUPNs. These experiments are designed to show how SUPNs, DNNs, KANs, and polynomial projection perform on functions of varying regularity.

We briefly summarize the commonalities among all examples. Every DNN uses tanh activations in each layer except the output layer. KANs use a grid size of 5, a spline order of 3, and the SiLU activation function, in keeping with [49]. The Supplementary Materials contain specific values for N and M in each example. All experiments were performed on Nvidia A100 GPUs, except the 10D case, which used H100 GPUs due to memory requirements. The loss functions for neural networks are highly nonconvex in general, so we use a large number of allowable optimization

iterations. Training consisted of a burn-in period of 5000 epochs with Adam (learning rate 10^{-3}), followed by trust-region Newton–CG. The trust-region algorithm has a maximum of 1000 Newton steps, a gradient tolerance of 10^{-6} and a step tolerance of 5×10^{-5} . At each optimization step, CG has an absolute tolerance of 10^{-4} , a relative tolerance of 10^{-2} , and a maximum of 500 iterations. Finally, five independent realizations of initial network weights are performed for each experiment. Accuracy on a validation set is monitored during training, and the reported accuracy is the test-set accuracy at the minimum validation-set error.

5.1. One dimension. We first compare best-approximation error and training robustness between SUPNs, DNNs, KANs, and polynomial projection on a suite of test problems. Then, we establish sampling requirements of SUPNs at models corresponding to low, medium, and high best-approximation accuracy. Next, we demonstrate that the order of convergence of SUPNs on the Runge function is independent of its bandwidth. Finally, we show that SUPNs can fit a highly oscillatory sine function from [1, Fig. 5] with an extended domain.

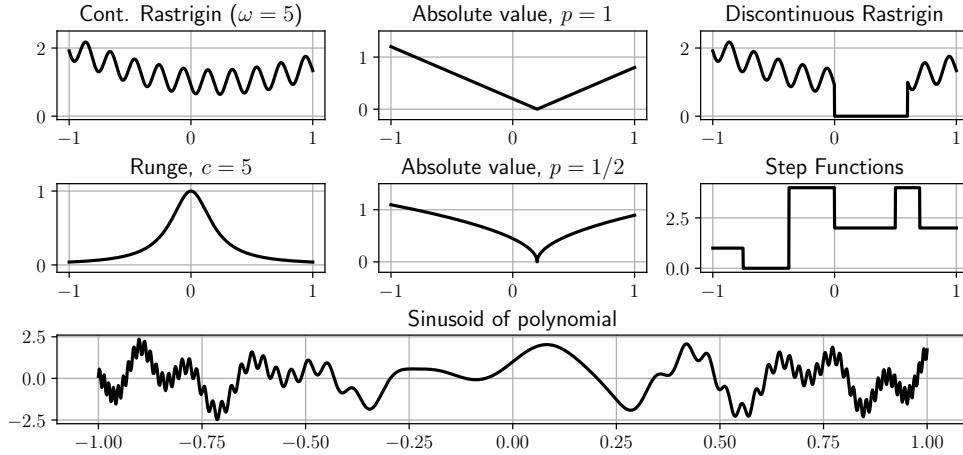


FIG. 4. 1D target functions for approximation.

5.1.1. Target functions. We consider seven 1D functions $f : [-1, 1] \rightarrow \mathbb{R}$, shown in Figure 4. These functions are defined below.

1. Continuous Rastrigin:

$$(5.1) \quad f_1(x; \omega) = 2(x - 0.2)^2 - \frac{1}{2.77} \cos(2\pi\omega x - 1.22) + 1$$

2. Discontinuous Rastrigin:

$$f_2(x) = \begin{cases} 0, & 0 \leq x \leq 0.6 \\ f_1(x; \omega = 5), & \text{otherwise} \end{cases}$$

3. Absolute value of order $p \in (0, 1]$:

$$f_3(x; p) = |x - 0.2|^p$$

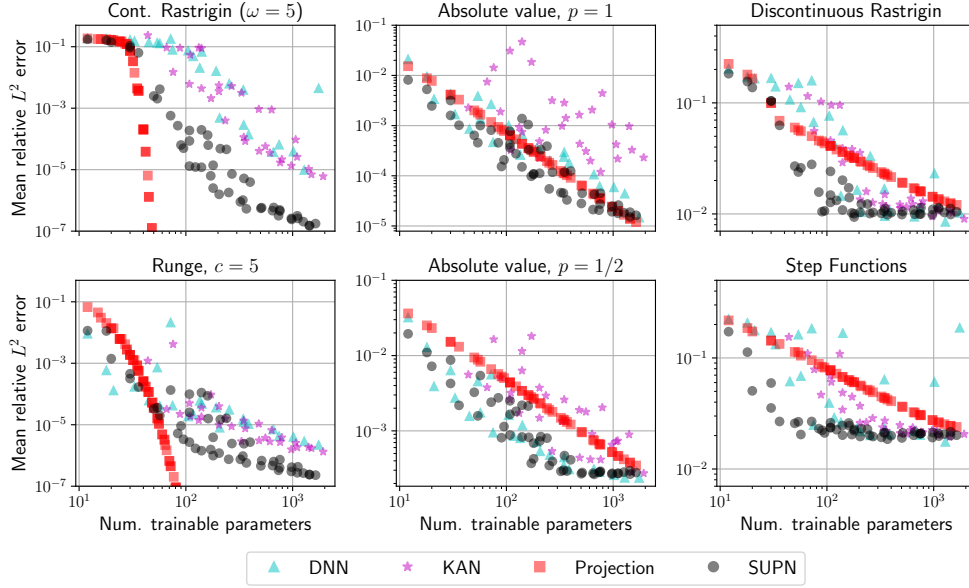


FIG. 5. Relative L^2 errors vs. trainable parameters using SUPN, DNN, KAN, and polynomial projection for f_k , $k \in [5]$.

4. Linear combination of step functions:

$$f_4(x) = \begin{cases} 1, & -1 \leq x < -0.75 \\ 0, & -0.75 \leq x < -0.375 \\ 4, & -0.375 \leq x < 0 \text{ or } 0.5 \leq x < 0.7 \\ 2, & 0 \leq x < 0.5 \text{ or } 0.7 \leq x \leq 1 \end{cases}$$

5. Runge function ($c \geq 1$):

$$(5.2) \quad f_5(x; c) = \frac{1}{1 + (cx)^2}$$

6. Sinusoid of polynomial:

$$f_6(x) = \sin(2\pi^2 x) + \cos(\pi^3 x^2) + \cos(\pi^4 x^3) \sin(\pi^4 x^3)$$

5.1.2. Best-approximation error. We first examine the error of a local minimizer of (2.2). We use a 2000-point Gauss–Legendre quadrature rule in the loss function (2.4). The validation set consists of 3001 equidistant points on $[-1, 1]$. The test set consists of 17001 equidistant points on $[-1, 1]$, and we report the test error at the minimum validation error encountered during training. We vary N , M , width and depth to create networks with $P \leq 2000$; exact details are in the Supplementary Materials. We perform five initializations of network weights from the Kaiming uniform distribution [28].

In Figure 5, we demonstrate that SUPNs are able to attain a given error tolerance with significantly fewer parameters than DNNs or KANs, often by an order of magnitude. For the continuously differentiable targets, i.e. continuous Rastrigin and Runge,

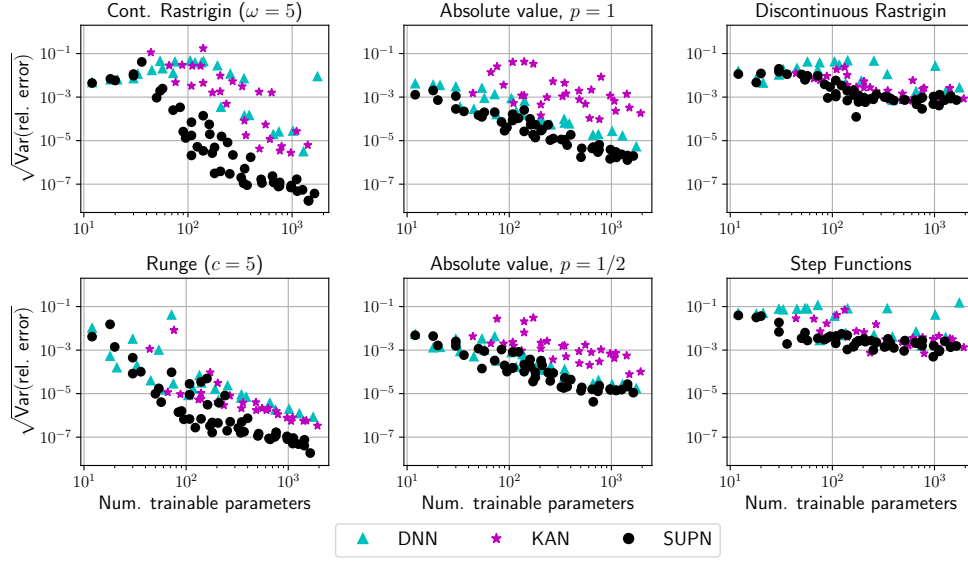


FIG. 6. Standard deviation of relative error over 5 realizations of initial network weights.

projection outperforms all the other methods given sufficiently many trainable parameters, but SUPNs outperform DNNs and KANs. For $f(x) = |x|$, SUPNs have a very slight advantage over DNNs and projection when $10^2 \leq P \leq 10^3$. For $f(x) = |x|^{1/2}$, SUPNs and DNNs track each other very closely, outperforming polynomial projection; the infinite derivative at $x = 0$ causes stagnation around $P \approx 500$. For discontinuous targets, SUPNs have the fewest parameters at the lowest error threshold. Similar trends hold for the L^∞ norm (except for discontinuous targets, where all methods fail); for brevity, further details may be found in the Supplementary Materials.

Figure 6 shows the standard deviation of relative error over five realizations of initial network weights. The standard deviation decreases as the mean relative error decreases, except in a handful of very small networks, which give reliably high error. Comparing Figure 5 and Figure 6 at a fixed number of trainable parameters shows that SUPNs are almost always more robust than DNNs and KANs for a given parameter count. The caveat is $f(x) = |x|^{1/2}$, where DNNs and SUPNs are *equally* robust.

5.1.3. Effect of hyperparameters. Figure 7 shows a heat map of SUPN, DNN, KAN, and polynomial projection best-approximation error on the continuous Rastrigin target, for a subset of the (N, M) pairs used to make Figure 5. The bottom row shows how P varies with different choices of N and M . This plot indicates that overly narrow networks (small N) can saturate the approximation capability of a degree- M polynomial in practice (e.g., initializing differently than Theorem 3.6). Likewise, networks with a low polynomial degree M will require a large number N of basis functions. SUPNs have the desirable feature of continuously improving in accuracy in (N, M) , instead of abruptly achieving high accuracy at some arbitrary (based on the target function) hyperparameter condition like DNNs or KANs. These observations are consistent with Theorem 2.5 and Theorem 2.6.

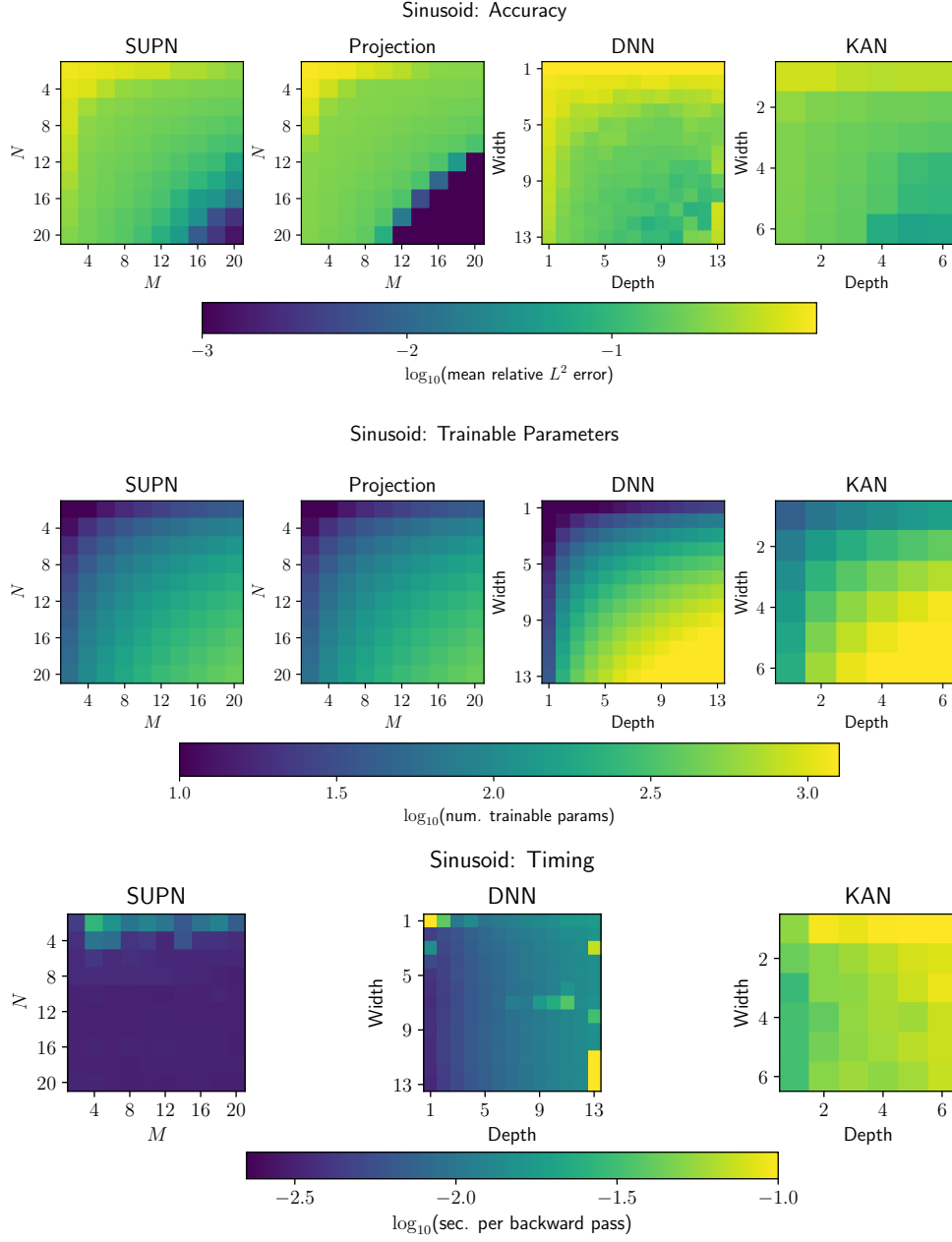


FIG. 7. Top: Heat map of test errors for $f_7(x)$ over hyperparameter sweep. Middle: Corresponding parameter counts. Bottom: Time required to compute one backward pass; projection is omitted because it requires no backward passes. Direct numerical comparisons are in Figure 5. Compared to DNNs and KANs, SUPNs display continuous accuracy improvement in (N, M) along with far lower parameter counts and relative errors.

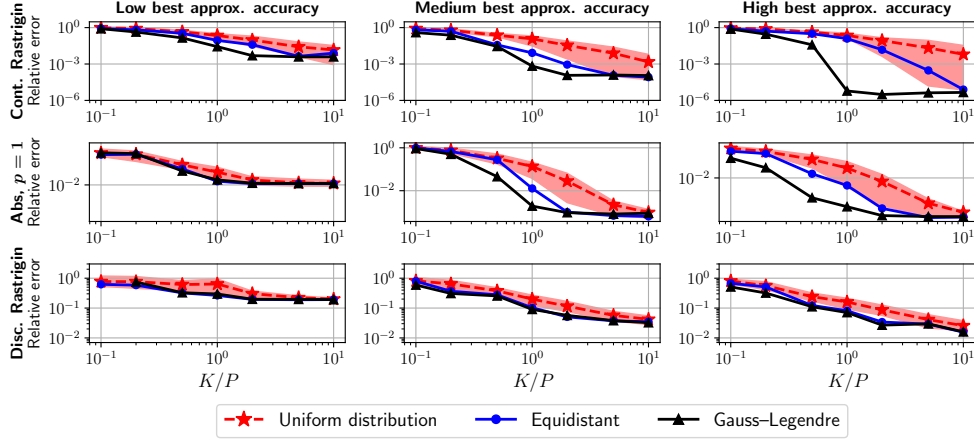


FIG. 8. Mean L^2 error (over network initialization) of finite-sampling experiments for different combinations of target function, best approximation accuracy, and sampling strategy. For the “uniform” results, the line is the mean over 10 realizations of training data and 5 weight initializations, while the shaded region is the 10th to 90th percentile.

5.1.4. Sampling error. Next, we empirically investigate sampling requirements for SUPNs. We vary the ratio between the number of training data K and the number of SUPN parameters, $P = N(M + 1)$. We examine models corresponding to qualitatively low, medium, and high best-approximation accuracy from Figure 5, using the sampling techniques discussed in subsection 2.2. Figure 8 indicates that Gauss-Legendre requires the fewest number of samples to achieve best-approximation error. The specific ratio appears to depend on the smoothness of the target function. We see up to $K = P$ for continuous Rastrigin, and $K = 2P$ for absolute value and discontinuous Rastrigin.

5.1.5. Convergence on Runge function. The Runge function (5.2) is a famous example of a function that can be difficult to approximate with naive methods [4]. Consider approximating the parameterized Runge function with polynomial projection in the L^2 norm. It can be shown that the best approximation error (2.12) using the P -term Legendre series (2.5) is bounded by

$$\epsilon(\hat{\theta}) \lesssim \exp(-\beta P/c)$$

for some universal constant $\beta > 0$. See, e.g., [73, Theorem 2.1] or [71, Theorem 8.2]. This bound gives spectral convergence, but the convergence rate decays with c . Figure 9 shows spectral convergence for polynomial projection that degrades as c increases. In contrast, SUPNs display a finite order of convergence, but that order appears to be independent of c .

5.1.6. Sinusoid of polynomial. This function uses an extended domain of the problem in [1, Sec. 3.1], in which the authors observe MLPs struggling to fit high-frequency components of the function. We display these results separately from the other functions because larger networks are required than for the experiments in Figure 5. Figure 10 demonstrates that SUPNs and polynomial projection have low error with a small number of parameters, while KANs and especially MLPs struggle

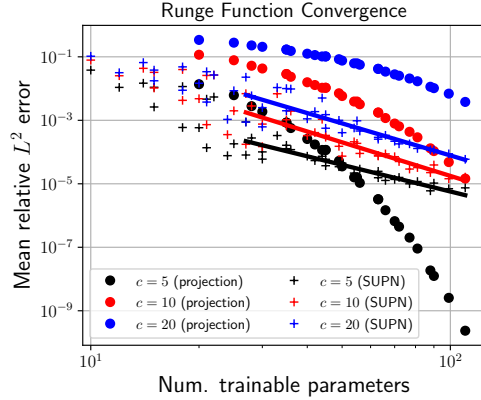


FIG. 9. Convergence of SUPN and polynomial projection on Runge function for $c = 5, 10, 20$. Lines of best fit are shown for SUPN.

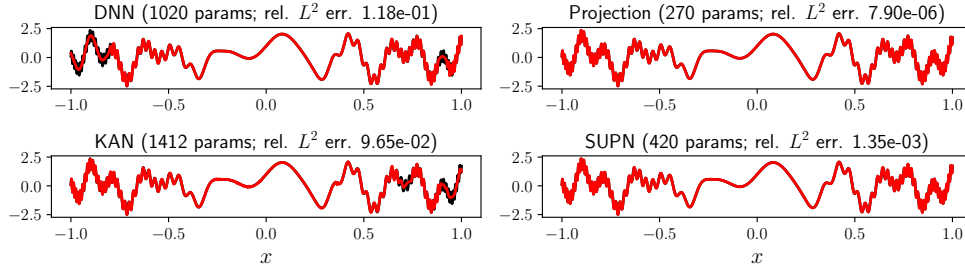


FIG. 10. Surrogate predictions on sinusoid example corresponding to median-error model.

at higher frequencies.

5.2. Two dimensions. We now present best-approximation results for functions with two input dimensions. The results of a finite-sampling study are in the Supplementary Materials and qualitatively similar to the one-dimensional case. For both SUPN and projection, we consider the isotropic hyperbolic cross-section Λ_{HC} (2.7) and total-degree space Λ_{TD} (2.8) to describe admissible polynomial degrees.

5.2.1. Target functions. We consider three 2D functions $f : [-1, 1]^2 \rightarrow \mathbb{R}$, based on the one-dimensional Rastrigin function (5.1) with $\omega = 5$. These target functions are shown in Figure 11 and are selected to have different amounts of regularity.

1. Continuous Rastrigin (sum):

$$f_7(x, y) = f_1(x; \omega = 5) + f_1(y; \omega = 5)$$

2. Continuous Rastrigin (radial):

$$f_8(x, y) = f_1(r; \omega = 5), \quad r = \sqrt{(x - 0.2)^2 + (y - 0.2)^2}$$

3. Discontinuous Rastrigin ($\omega = 5$):

$$f_9(x, y) = \begin{cases} 0, & r \in [0.3, 0.5] \\ f_1(|x - 0.2|)f_1(|y - 0.2|), & \text{otherwise} \end{cases}$$

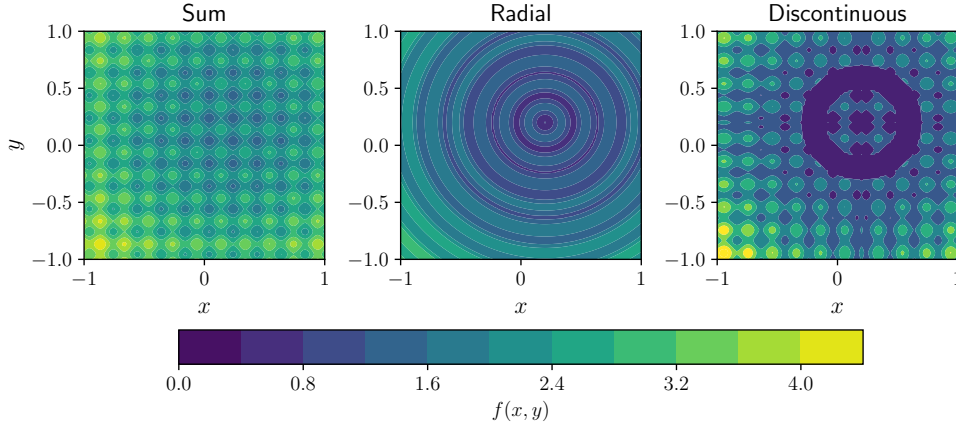


FIG. 11. 2D target functions for approximation.

5.2.2. Best-approximation error. In the loss function (2.4), we use a tensor-product Gauss–Legendre rule with 200 nodes in each dimension. The validation set consists of a 130×130 grid of equispaced points, and the test set consists of a 450×450 grid of equispaced points. Due to the computational cost of a tensor-product grid, the validation and test sets are commensurate in size with the training set, unlike in subsection 5.1.2. As before, we report the test error at the minimum validation error seen during training.

Figure 12 shows that SUPNs are a competitive method in 2D. For a sum of one-dimensional Rastrigin functions, SUPNs dramatically outperform DNNs and KANs. SUPNs with a properly chosen index set Λ even outperform polynomial projection with a misspecified Λ . Interestingly, DNNs outperform all methods on the radial target, which we explain by the lack of tensor-product structure in f_8 causing difficulty for polynomial-based methods. However, with a transform into polar coordinates (r, θ) , the radial results would become qualitatively similar to the top left panel of Figure 5. Moreover, the horizontal spread of KANs on the radial target is much wider than for SUPNs, so the performance of KANs is less predictable with respect to its width and depth. On the discontinuous target, SUPNs have the lowest error among all methods at a given parameter count. Robustness on these 2D examples is qualitatively similar to Figure 6 and may be found in the Supplementary Materials.

5.3. Ten dimensions. We now demonstrate that SUPNs can scale into moderately large dimensions. We consider the anisotropic and purposefully nonlinear target function

$$f_{\text{aniso}}(\mathbf{x}) = \exp(x_1 - 0.7) \sin(1.3x_2) + 0.2 \cos(2\pi x_3) + 0.01|x_4 - 0.27|x_5 \\ + 0.1|x_6|x_7 + 0.05 \exp(-(x_8 - 0.3)^2/16) + 0.1x_9x_{10}$$

where $\mathbf{x} \in [-1, 1]^{10}$. We use 10^5 samples from the Halton sequence for training. For validation, we use the next 2×10^5 Halton elements, and for testing, the next 2×10^5 Halton elements. As in subsection 5.2, all index sets are isotropic.

Figure 13 displays the results. Asymptotically, SUPNs are an order of magnitude more accurate than both KANs and DNNs. To ensure we show best-approximation

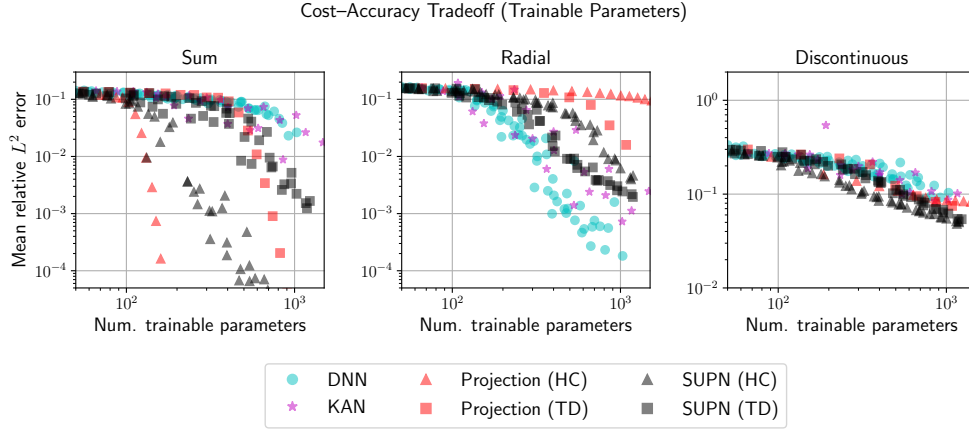


FIG. 12. Mean L^2 error versus number of trainable parameters P for 2D examples (f_k , $k = 7, 8, 9$).

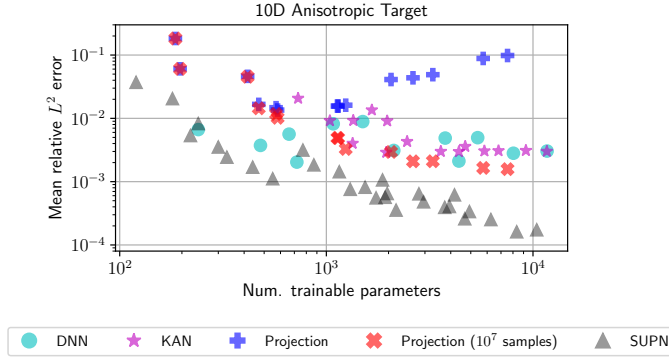


FIG. 13. Mean relative L^2 error on anisotropic target function with $D = 10$ input dimensions.

error rather than sampling error, we also show projection using 10^7 Halton elements to compute the L^2 inner products. With only 10^5 samples, there is significant quadrature error in the higher-order Legendre expansion coefficients ($P \gtrsim 10^3$).

6. Conclusion and future work. We have presented shallow universal polynomial networks, a parsimonious and provably convergent surrogate model that avoids common issues arising from highly over-parameterized models. Intuitively, SUPNs replace the learned bases of a DNN or KAN with a learned polynomial, which dramatically reduces the parameter count of the network. The early layers of an NN learn the same simplistic building blocks, creating unnecessarily complex loss landscapes which can inhibit the learning of high frequency features shown in Section 5. An extensive set of numerical experiments validates that SUPNs outperform DNNs and KANs with on one-, two-, and ten-dimensional target functions that exhibit tensor-product structure. Additionally, the performance of SUPNs is significantly more robust with respect to network initialization.

In future work, we see SUPNs as potential building-blocks for operator learning in

one, two, and three spatial dimensions, which are the predominant use-case of neural networks in physics-informed cases. Although MLPs work well in high dimensions, they cannot match traditional approximation methods in lower dimensions. Physics-informed SUPNs may overcome issues fitting higher frequencies in PDEs such as the Helmholtz equation without the need for domain decomposition such as in [55]. Additionally, we plan to investigate adaptive methods for populating Λ so that SUPNs can exploit anisotropy to scale to even higher dimensions.

Appendix A. Symbols and Notation.

D	Input dimension	$\mathcal{L}(\cdot)$	Loss function
Q	Output dimension	θ	Learnable parameters
P	Parameter dimension	$\hat{\theta}$	Best-approximation parameters
Ω	Compact subset of \mathbb{R}^D	$\tilde{\theta}$	Empirically trained parameters
T_i	Chebyshev polynomial of deg. i	σ	Activation function
L_i	Legendre polynomial of deg. i	$B_{j,i}$	B-spline of deg. i on partition j
K	Training set size	ϵ	Generalization error
$[N]$	$\{1, 2, \dots, N\}$	$[N]_0$	$\{0, 1, \dots, N\}$
\mathbb{P}_n	Polynomials of degree n		

TABLE 1
Symbols and notation.

REFERENCES

- [1] D. W. ABUEIDDA, P. PANTIDIS, AND M. E. MOBASHER, *DeepOKAN: Deep operator network based on Kolmogorov Arnold networks for mechanics problems*, Comput. Meth. Appl. Mech. Eng., 436 (2025), p. 117699, <https://doi.org/10.1016/j.cma.2024.117699>.
- [2] J. A. ACTOR, G. HARPER, B. SOUTHWORTH, AND E. C. CYR, *Leveraging KANs for expedient training of multichannel MLPs via preconditioning and geometric refinement*, 2025, <https://arxiv.org/abs/2505.18131>.
- [3] V. I. ARNOLD, *On the representation of functions of several variables as a superposition of functions of a smaller number of variables*, Springer Berlin Heidelberg, 2009, pp. 25–46.
- [4] K. E. ATKINSON, *An Introduction to Numerical Analysis*, Wiley, 1989.
- [5] V. BARTHELMANN, E. NOVAK, AND K. RITTER, *High dimensional polynomial interpolation on sparse grids*, Adv. Comput. Math., 12 (2000), pp. 273–288, <https://doi.org/10.1023/A:1018977404843>.
- [6] P. BENNER, S. GUGERCIN, AND K. WILLCOX, *A survey of projection-based model reduction methods for parametric dynamical systems*, SIAM Rev., 57 (2015), pp. 483–531, <https://doi.org/10.1137/130932715>.
- [7] P. BENNER, M. OHLBERGER, A. COHEN, AND K. WILLCOX, *Model Reduction and Approximation: Theory and Algorithms*, SIAM, 2017.
- [8] R. BOLLAPRAGADA, R. H. BYRD, AND J. NOCEDAL, *Exact and inexact subsampled Newton methods for optimization*, IMA J. Numer. Anal., 39 (2018), pp. 545–578, <https://doi.org/10.1093/imanum/dry009>.
- [9] J. BRAUN AND M. GRIEBEL, *On a constructive proof of Kolmogorov’s superposition theorem*, Constr. Approx., 30 (2009), pp. 653–675, <https://doi.org/10.1007/s00365-009-9054-2>.
- [10] M. D. BUHMANN, *Radial Basis Functions: Theory and Implementations*, Cambridge University Press, 2003.
- [11] C. CANUTO, M. Y. HUSSAINI, A. QUARTERONI, AND T. A. ZANG JR, *Spectral Methods: Fundamentals in Single Domains*, Springer, 2006.
- [12] A. CHKIFA, N. DEXTER, H. TRAN, AND C. G. WEBSTER, *Polynomial approximation via compressed sensing of high-dimensional functions on lower sets*, Math. Comput., 87 (2018), pp. 1415–1450, <https://doi.org/10.1090/mcom/3272>.
- [13] G. G. CHRYSOS, S. MOSCHOLOU, G. BOURITSAS, J. DENG, Y. PANAGAKIS, AND S. ZAFEIRIOU, *Deep polynomial neural networks*, IEEE Trans. Pattern Anal. Mach. Intell., 44 (2021),

- pp. 4021–4034.
- [14] A. R. CONN, N. I. M. GOULD, AND P. L. TOINT, *Trust Region Methods*, SIAM, 2000, <https://doi.org/10.1137/1.9780898719857>.
 - [15] N. COTTER, *The Stone–Weierstrass theorem and its application to neural networks*, IEEE Trans. Neural Netw., 1 (1990), pp. 290–295, <https://doi.org/10.1109/72.80265>.
 - [16] G. CYBENKO, *Approximation by superpositions of a sigmoidal function*, Math. Control Signals Syst., 2 (1989), pp. 303–314, <https://doi.org/10.1007/BF02551274>.
 - [17] E. C. CYR, M. A. GULIAN, R. G. PATEL, M. PEREGO, AND N. A. TRASK, *Robust training and initialization of deep neural networks: An adaptive basis viewpoint*, in Proceedings of the First Mathematical and Scientific Machine Learning Conference, J. Lu and R. Ward, eds., vol. 107 of Proceedings of Machine Learning Research, 2020, pp. 512–536, <https://proceedings.mlr.press/v107/cyr20a.html>.
 - [18] S. DAS, *The polynomial neural network*, Inf. Sci., 87 (1995), pp. 231–246, [https://doi.org/https://doi.org/10.1016/0020-0255\(95\)00133-6](https://doi.org/https://doi.org/10.1016/0020-0255(95)00133-6).
 - [19] C. DE BOOR, *A Practical Guide to Splines*, Springer, 1978, <https://doi.org/10.2307/2006241>.
 - [20] T. DE RYCK, S. LANTHALER, AND S. MISHRA, *On the approximation of functions by tanh neural networks*, Neural Netw., 143 (2021), pp. 732–750, <https://doi.org/https://doi.org/10.1016/j.neunet.2021.08.015>.
 - [21] C. FINN, P. ABBEEL, AND S. LEVINE, *Model-agnostic meta-learning for fast adaptation of deep networks*, in International Conference on Machine Learning, 2017, pp. 1126–1135.
 - [22] D. FOKINA AND I. OSELEDETS, *Growing axons: Greedy learning of neural networks with application to function approximation*, 2020, <https://arxiv.org/abs/1910.12686>.
 - [23] K. O. GEDDES, *Near-minimax polynomial approximation in an elliptical region*, SIAM J. Numer. Anal., 15 (1978), pp. 1225–1233, <https://doi.org/10.1137/0715083>.
 - [24] R. G. GHANEM AND P. D. SPANOS, *Stochastic Finite Element Method: Response Statistics*, Springer, 1991, pp. 101–119.
 - [25] X. GLOROT AND Y. BENGIO, *Understanding the difficulty of training deep feedforward neural networks*, in Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, Y. W. Teh and M. Titterton, eds., vol. 9 of Proceedings of Machine Learning Research, 2010, pp. 249–256, <https://proceedings.mlr.press/v9/glorot10a.html>.
 - [26] A. GORODETSKY, S. KARAMAN, AND Y. MARZOUK, *A continuous analogue of the tensor-train decomposition*, Comput. Meth. Appl. Mech. Eng., 347 (2019), pp. 59–84, <https://doi.org/10.1016/j.cma.2018.12.015>.
 - [27] D. GOTTLIEB AND S. A. ORSZAG, *Numerical Analysis of Spectral Methods*, SIAM, 1977, <https://doi.org/10.1137/1.9781611970425>.
 - [28] K. HE, X. ZHANG, S. REN, AND J. SUN, *Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification*, in 2015 IEEE International Conference on Computer Vision, 2015, pp. 1026–1034, <https://doi.org/10.1109/ICCV.2015.123>.
 - [29] S. HOCHREITER AND J. SCHMIDHUBER, *Long short-term memory*, Neural Comput., 9 (1997), pp. 1735–1780, <https://doi.org/10.1162/neco.1997.9.8.1735>.
 - [30] J. J. HOPFIELD, *Neural networks and physical systems with emergent collective computational abilities*, Proc. Nat. Acad. Sci. USA, 79 (1982), pp. 2554–2558, <https://doi.org/10.1073/pnas.79.8.2554>.
 - [31] K. HORNIK, *Approximation capabilities of multilayer feedforward networks*, Neural Netw., 4 (1991), pp. 251–257, [https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T).
 - [32] K. HORNIK, M. STINCHCOMBE, AND H. WHITE, *Multilayer feedforward networks are universal approximators*, Neural Netw., 2 (1989), pp. 359–366, [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8).
 - [33] A. G. IVAKHNENKO, *Polynomial theory of complex systems*, IEEE Trans. Syst. Man Cybern., SMC-1 (1971), pp. 364–378, <https://doi.org/10.1109/TSMC.1971.4308320>.
 - [34] D. JACKSON, *The Theory of Approximation*, American Mathematical Society, 1930, ch. 1. Theorem VII.
 - [35] J. D. JAKEMAN, M. PEREGO, D. T. SEIDL, T. A. HARTLAND, T. R. HILLEBRAND, M. J. HOFFMAN, AND S. F. PRICE, *An evaluation of multi-fidelity methods for quantifying uncertainty in projections of ice-sheet mass change*, Earth Syst. Dyn., 16 (2025), pp. 513–544, <https://doi.org/10.5194/esd-16-513-2025>, <https://esd.copernicus.org/articles/16/513/2025/>.
 - [36] A. JAVED, D. KOURI, D. RIDZAL, AND G. VON WINCKEL, *Get ROL-ing: An introduction to Sandia’s Rapid Optimization Library*, in 7th International Conference on Continuous Optimization, 2022.
 - [37] J. JUMPER, R. EVANS, A. PRITZEL, T. GREEN, M. FIGURNOV, O. RONNEBERGER, K. TUNYASUVUNAKOOL, R. BATES, A. ZIDEK, A. POTAPENKO, A. BRIDGLAND, C. MEYER, S. A. A. KOHL, A. J. BALLARD, A. COWIE, B. ROMERA-PAREDES, S. NIKOLOV, R. JAIN,

- J. ADLER, T. BACK, S. PETERSEN, D. REIMAN, E. CLANCY, M. ZIELINSKI, M. STEINEGER, M. PACHOLSKA, T. BERGHAMMER, S. BODENSTEIN, D. SILVER, O. VINYALS, A. W. SENIOR, K. KAVUKCUOGLU, P. KOHLI, AND D. HASSABIS, *Highly accurate protein structure prediction with AlphaFold*, *Nature*, 596 (2021), pp. 583–589, <https://doi.org/10.1038/s41586-021-03819-2>.
- [38] S. C. KAK, *Feedback neural networks: New characteristics and a generalization*, *Circuits Syst. and Signal Process.*, 12 (1993), pp. 263–278, <https://doi.org/10.1007/BF01189877>.
- [39] P. KIDGER AND T. LYONS, *Universal approximation with deep narrow networks*, 2020, <https://arxiv.org/abs/1905.08539>.
- [40] J. KILEEL, M. TRAGER, AND J. BRUNA, *On the expressive power of deep polynomial neural networks*, *Adv. Neural Inf. Process. Syst.*, 32 (2019).
- [41] D. P. KINGMA AND J. BA, *Adam: A method for stochastic optimization*, 2017, <https://arxiv.org/abs/1412.6980>.
- [42] A. N. KOLMOGOROV, *On the representations of continuous functions of many variables by superposition of continuous functions of one variable and addition*, in *Dokl. Akad. Nauk USSR*, vol. 114, 1957, pp. 953–956, <http://mi.mathnet.ru/dan22050>.
- [43] K. KUBJAS, J. LI, AND M. WIESMANN, *Geometry of polynomial neural networks*, *Algebr. Stat.*, 15 (2024), pp. 295–328.
- [44] H. LEE, Y. KIM, S. Y. YANG, AND H. CHOI, *Improved weight initialization for deep and narrow feedforward neural network*, *Neural Netw.*, 176 (2024), p. 106362, <https://doi.org/10.1016/j.neunet.2024.106362>.
- [45] M. LESHNO, V. Y. LIN, A. PINKUS, AND S. SCHOCKEN, *Multilayer feedforward networks with a nonpolynomial activation function can approximate any function*, *Neural Netw.*, 6 (1993), pp. 861–867, [https://doi.org/https://doi.org/10.1016/S0893-6080\(05\)80131-5](https://doi.org/https://doi.org/10.1016/S0893-6080(05)80131-5).
- [46] Z. LI, N. KOVACHKI, K. AZIZZADENESHELI, B. LIU, K. BHATTACHARYA, A. STUART, AND A. ANANDKUMAR, *Fourier neural operator for parametric partial differential equations*, 2021, <https://arxiv.org/abs/2010.08895>.
- [47] D. C. LIU AND J. NOCEDAL, *On the limited memory BFGS method for large scale optimization*, *Math. Program.*, 45 (1989), pp. 503–528.
- [48] X. LIU, X. ZHANG, W. PENG, W. ZHOU, AND W. YAO, *A novel meta-learning initialization method for physics-informed neural networks*, *Neural Comput. Appl.*, 34 (2022), pp. 14511–14534.
- [49] Z. LIU, Y. WANG, S. VAIDYA, F. RUEHLE, J. HALVERSON, M. M. SOLJACIĆ, T. Y. HOU, AND M. TEGMARK, *KAN: Kolmogorov–Arnold networks*, 2025, <https://arxiv.org/abs/2404.19756>.
- [50] M. LOWERY, J. TURNAGE, Z. MORROW, J. D. JAKEMAN, A. NARAYAN, S. ZHE, AND V. SHANKAR, *Kernel neural operators (KNOs) for scalable, memory-efficient, geometrically-flexible operator learning*, 2024, <https://arxiv.org/abs/2407.00809>.
- [51] L. LU, P. JIN, G. PANG, Z. ZHANG, AND G. E. KARNIAKAKIS, *Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators*, *Nat. Mach. Intell.*, 3 (2021), pp. 218–229, <https://doi.org/10.1038/s42256-021-00302-5>.
- [52] L. LU, S. YANHUI, AND G. E. KARNIAKAKIS, *Collapse of deep and narrow neural nets*, 2018, <https://arxiv.org/abs/1808.04947>.
- [53] L. LU, S. YEONJONG, S. YANHUI, AND G. E. KARNIAKAKIS, *Dying ReLU and initialization: Theory and numerical examples*, *Commun. Comput. Phys.*, 28 (2020), pp. 1671–1706, <https://doi.org/10.4208/cicp.OA-2020-0165>.
- [54] Z. MORROW AND M. STOYANOV, *A method for dimensionally adaptive sparse trigonometric interpolation of periodic functions*, *SIAM J. Sci. Comput.*, 42 (2020), pp. A2436–A2460, <https://doi.org/10.1137/19M1283483>.
- [55] B. MOSELEY, A. MARKHAM, AND T. NISSEN-MEYER, *Finite basis physics-informed neural networks (FBPINNs): A scalable domain decomposition approach for solving differential equations*, *Adv. Comput. Math.*, 49 (2023), p. 62.
- [56] A. NICHOL, J. ACHIAM, AND J. SCHULMAN, *On first-order meta-learning algorithms*, 2018, <https://arxiv.org/abs/1803.02999>.
- [57] F. NOBILE, R. TEMPONE, AND C. G. WEBSTER, *A sparse grid stochastic collocation method for partial differential equations with random input data*, *SIAM J. Numer. Anal.*, 46 (2008), pp. 2309–2345, <https://doi.org/10.1137/060663660>.
- [58] J. NOCEDAL AND S. WRIGHT, *Numerical Optimization*, Springer, 2nd ed., 2006, <https://doi.org/10.1007/978-0-387-40065-5>.
- [59] T. O’LEARY-ROSEBERRY, N. ALGER, AND O. GHATTAS, *Inexact Newton methods for stochastic nonconvex optimization with applications to neural network training*, 2019, <https://arxiv.org/abs/1905.06738>.

- [60] B. PEHERSTORFER AND K. WILLCOX, *Data-driven operator inference for nonintrusive projection-based model reduction*, Comput. Meth. Appl. Mech. Eng., 306 (2016), pp. 196–215, <https://doi.org/10.1016/j.cma.2016.03.025>.
- [61] M. PENWARDEN, S. ZHE, A. NARAYAN, AND R. M. KIRBY, *A metalearning approach for physics-informed neural networks (PINNs): Application to parameterized PDEs*, J. Comput. Phys., 477 (2023), p. 111912.
- [62] A. PINKUS, *Approximation theory of the MLP model in neural networks*, Acta Numer., 8 (1999), p. 143–195, <https://doi.org/10.1017/S0962492900002919>.
- [63] H. POON AND P. DOMINGOS, *Sum-product networks: A new deep architecture*, in 2011 IEEE International Conference on Computer Vision Workshops, 2011, pp. 689–690.
- [64] O. RONNEBERGER, P. FISCHER, AND T. BROX, *U-net: Convolutional networks for biomedical image segmentation*, in International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2015, pp. 234–241, https://doi.org/10.1007/978-3-319-24574-4_28.
- [65] K. SHUKLA, J. D. TOSCANO, Z. WANG, Z. ZOU, AND G. E. KARNIADAKIS, *A comprehensive and FAIR comparison between MLP and KAN representations for differential equations and operator networks*, Comput. Meth. Appl. Mech. Eng., 431 (2024), p. 117290, <https://doi.org/10.1016/j.cma.2024.117290>.
- [66] S. S. SIDHARTH, A. R. KEERTHANA, R. GOKUL, AND K. P. ANAS, *Chebyshev polynomial-based Kolmogorov–Arnold networks: An efficient architecture for nonlinear function approximation*, 2024, <https://arxiv.org/abs/2405.07200>.
- [67] S. SONODA AND N. MURATA, *Neural network with unbounded activation functions is universal approximator*, Appl. Comput. Harmon. Anal., 43 (2017), pp. 233–268.
- [68] T. STEihaug, *The conjugate gradient method and trust regions in large scale optimization*, SIAM J. Numer. Anal., 20 (1983), pp. 626–637, <https://doi.org/10.1137/0720042>.
- [69] M. H. STONE, *Applications of the theory of Boolean rings to general topology*, Trans. Am. Math. Soc., 41 (1937), pp. 375–481, <https://doi.org/10.1090/S0002-9947-1937-1501905-7>.
- [70] P. L. TOINT, *Towards an efficient sparsity exploiting Newton method for minimization*, in Sparse Matrices and Their Uses, 1981, <https://api.semanticscholar.org/CorpusID:115391681>.
- [71] L. N. TREFETHEN, *Approximation Theory and Approximation Practice*, SIAM, 2012.
- [72] A. VASWANI, N. SHAZEER, N. PARMAR, J. USZKOREIT, L. JONES, A. N. GOMEZ, L. KAISER, AND I. POLOSUKHIN, *Attention is all you need*, Adv. Neural Inf. Process. Syst., 30 (2017), https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- [73] H. WANG AND S. XIANG, *On the convergence rates of Legendre approximation*, Math. Comput., 81 (2012), pp. 861–877, <https://doi.org/10.1090/S0025-5718-2011-02549-4>.
- [74] K. WEIERSTRASS, *Über die analytische Darstellbarkeit sogenannter willkürlicher Functionen reeller Argumente*, Königl. Akad. der Wiss., 9 (1885), pp. 1–37, <https://doi.org/10.1017/CBO9781139567886.002>.
- [75] D. XIU AND G. E. KARNIADAKIS, *The Wiener–Askey polynomial chaos for stochastic differential equations*, SIAM J. Sci. Comput., 24 (2002), pp. 619–644, <https://doi.org/10.1137/S1064827501387826>.
- [76] P. XU, F. ROOSTA, AND M. W. MAHONEY, *Second-order Optimization for Non-convex Machine Learning: An Empirical Study*, SIAM, 2020, pp. 199–207, <https://doi.org/10.1137/1.9781611976236.23>.
- [77] Z. XU, Y. CHEN, AND D. XIU, *Chebyshev feature neural network for accurate function approximation*, J. Mach. Learn. Model. Comput., 6 (2025), <https://doi.org/10.1615/JMachLearnModelComput.2025056536>.

SUPPLEMENTARY MATERIALS: SUPN: SHALLOW UNIVERSAL POLYNOMIAL NETWORKS*

ZACHARY MORROW^{†‡}, MICHAEL PENWARDEN^{†‡}, BRIAN CHEN[‡], AURYA JAVEED[‡],
AKIL NARAYAN[§], AND JOHN D. JAKEMAN[‡]

SM1. One dimension.

SM1.1. L^∞ best-approximation error. In Figure SM1.1, we observe similar convergence trends in learnable parameter count for L^∞ as we do in L^2 for continuous function approximation cases. The exceptions are that KANs perform noticeably worse on $f(x) = |x|^p$, and SUPNs are now evenly matched on the Runge function. For discontinuous cases, all methods suffer in L^∞ .

SM1.2. Architecture details. We specify the architectures of the 1D networks in Table SM1.1. We performed polynomial projection with the same number of degrees of freedom as SUPNs.

SM2. Two dimensions.

SM2.1. L^∞ best-approximation error. L^∞ convergence on the 2D examples (Figure SM2.1) is qualitatively similar to the L^2 case, with two exceptions. First, uniform approximation of a discontinuous function by polynomials is error-prone. Second, the L^∞ accuracy of projection on the radial example has seriously degraded from the L^2 setting, but it appears that (for Λ_{TD} , at least) accuracy begins to improve after $P \approx 1000$. This is consistent with the lack of tensor-product structure in the radial target function.

SM2.2. L^2 robustness. Figure SM2.2 shows that for $P \gtrsim 100$, SUPNs have between 5-100x less variability in their performance than DNNs and KANs at the same error tolerance. These results comport with those for the one-dimensional examples. Compared to the 1D robustness plots, we do observe an early *increase* in variability up to $P \approx 100$. In this regime, networks have sufficiently few parameters to yield reliable predictions that are nonetheless poor.

SM2.3. Finite sampling. Figure SM2.3 shows the results of a finite sampling experiment. The horizontal axis is the ratio of the univariate grid size to the square-root of the number of network parameters. The largest ratio needed to obtain best-

[†]Z. M. and M. P. contributed equally to this work.

*Submitted to the editors on November 25, 2025.

Funding: The authors gratefully acknowledge funding from Sandia National Laboratories' Laboratory Directed Research and Development program. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government. The publisher acknowledges that the U.S. Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this written work or allow others to do so, for U.S. Government purposes. The DOE will provide public access to results of federally sponsored research in accordance with the DOE Public Access Plan. SAND2025-14696O

[‡]Sandia National Laboratories, 1515 Eubank Blvd SE, Albuquerque, NM 87123 ({[zbmorrow](mailto:zbmorrow@sandia.gov), [mshpenwa](mailto:mshpenwa@sandia.gov), [brichen](mailto:brichen@sandia.gov), [asjavee](mailto:asjavee@sandia.gov), [jdjakem](mailto:jdjakem@sandia.gov)}@sandia.gov).

[§]Scientific Computing and Imaging Institute and Department of Mathematics, University of Utah, 155 S 1400 E, Room 233, Salt Lake City, UT 84112 (akil@sci.utah.edu).

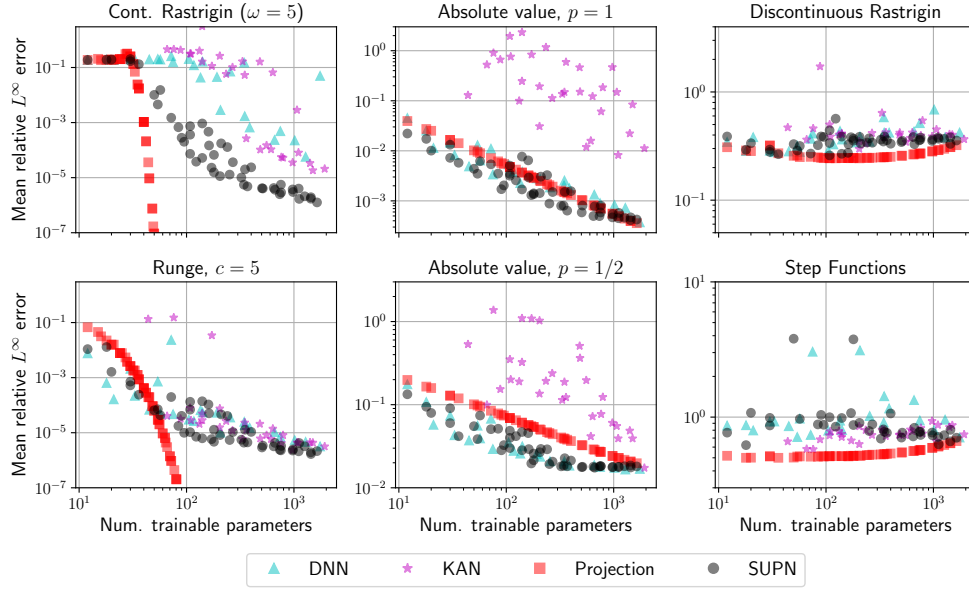


FIG. SM1.1. Relative L^∞ errors vs. trainable parameters using SUPN, DNN, KAN, and polynomial projection for f_k , $k \in [6]$.

	SUPN	DNN	KAN
N (width)	3, 5, 9, 18, 27, 35, 40	2, 3, 5, 9, 12	[6]
M (depth)	3, 5, 9, 18, 27, 35, 40	2, 3, 5, 9, 12	[5]

TABLE SM1.1
Network architectures for 1D examples.

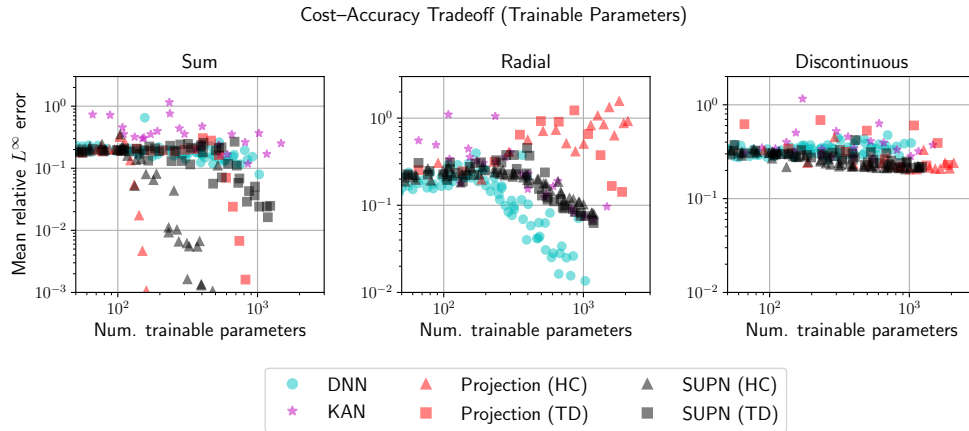


FIG. SM2.1. Relative L^∞ errors in 2D.

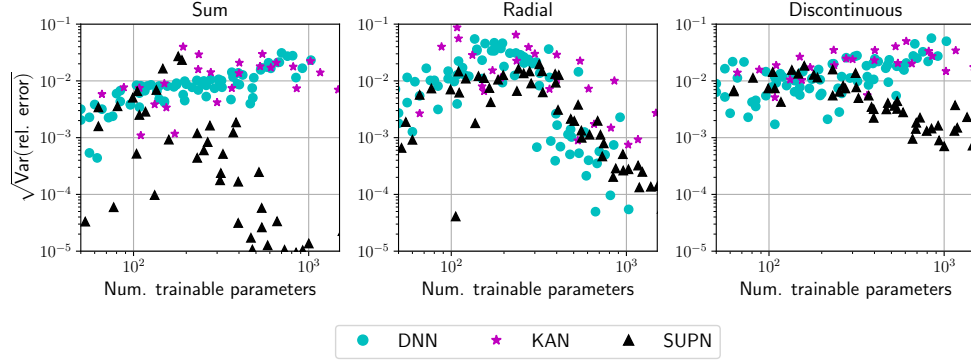


FIG. SM2.2. Standard deviation of relative error over 5 realizations of initial network weights for 2D examples. To reduce clutter, only the better-performing of Λ_{TD} and Λ_{HC} are shown for SUPN.

	SUPN (HC)	SUPN (TD)	DNN	KAN
N (width)	1, 3, 5, 7, 9, 11, 13, 15	1, 3, 5, 7, 9, 11, 13, 15	[10]	[5]
L (depth)	0, 3, 7, 11, 16, 22, 28, 34	0, 2, 4, 6, 8, 10, 12, 14	[10]	[5]

TABLE SM2.1

Network architectures for 2D examples.

approximation accuracy is $K_{1D} \approx 5\sqrt{P}$. Since our networks have $P < 1600$, then our tensor-product Gauss–Legendre rule with 200 nodes in each dimension is sufficient.

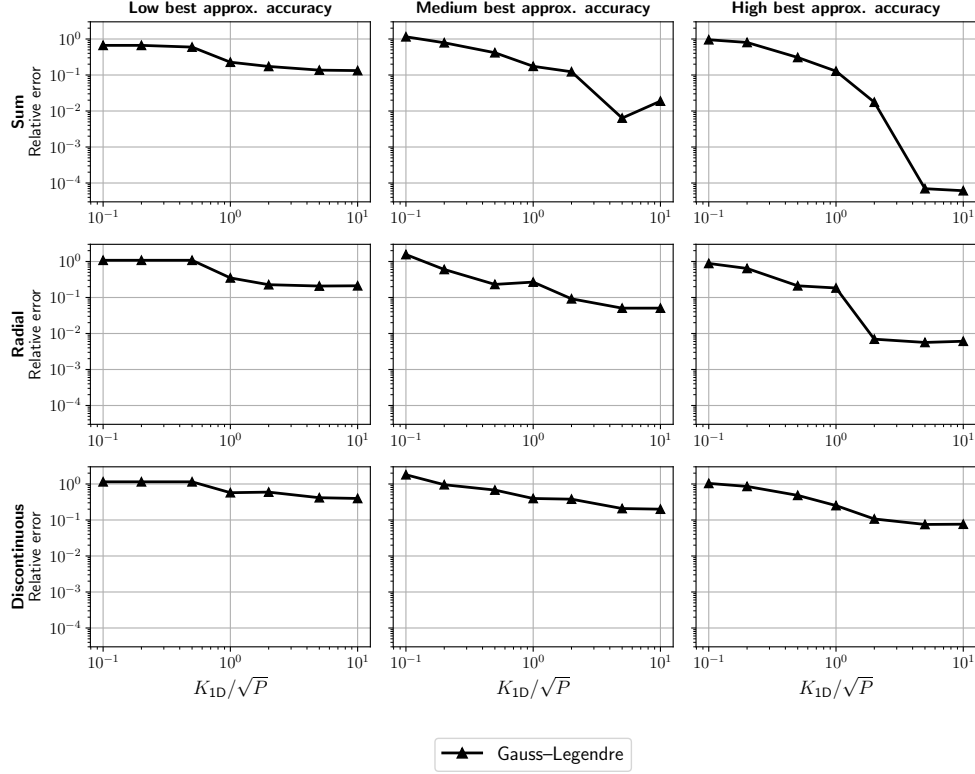
SM2.4. Architecture details. Table SM2.1 describes the architecture for the 2D targets. Here, L refers to the contour of either $\Lambda_{HC}(L)$ or $\Lambda_{TD}(L)$.

SM3. Ten dimensions.

SM3.1. Architecture details. Table SM3.1 describes the architecture for the 10D targets.

SM4. Best-approximation predictions. In this section, we plot the predictions corresponding to the 1D step function and continuous Rastrigin targets. These plots provide a qualitative complement to the best-approximation error plots in our results. We sweep over $(N, M) \in [10] \times [10]$. For KANs and DNNs, N denotes network width, and M denotes network depth.

SM4.1. Step functions. Figures SM4.1 to SM4.4 demonstrate that the assumption of continuity on the target function is truly necessary for L^∞ convergence. All methods can exhibit persistent spikes over very small sub-intervals of Ω at the discontinuities, analogous to Gibbs effects. Note that we distinguish between the training and testing grid. These spikes are not present on the training grid, since the error of the approximation is included in the loss and therefore minimized at that grid. However, at finer resolutions during testing, there is no constraint (it is necessary to not have a constraint to fit the discontinuities in the first place during training). To fit the discontinuities at the training grid, the basis functions become extreme which

FIG. SM2.3. *Finite-sampling error on the 2D examples.*

	SUPN (HC)	DNN	KAN
N (width)	5, 10, 15, 20	20, 40, 60	5, 10, 15
L (depth)	[7]	[4]	[5]

TABLE SM3.1

Network architectures for 10D example.

can randomly induce spikes in the underlying approximation which is seen in L^∞ but not L^2 since they contribute very little in that norm. This is a limitation of all methods, but is particularly bad for SUPNs, and could be a future research direction to mitigate the effect while still maintaining the ability to fit discontinuities at the training grid level.

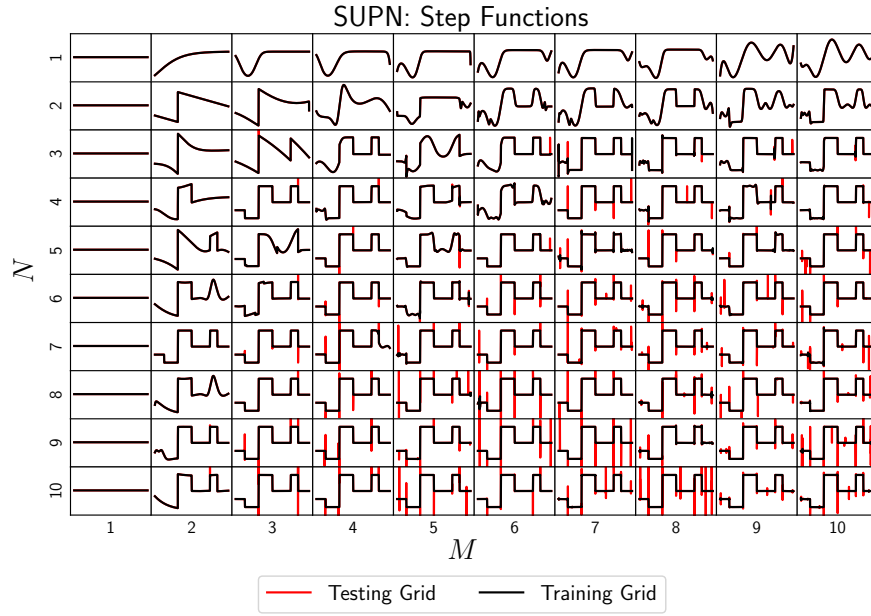


FIG. SM4.1. Plot of best approximation with SUPN on the step functions over (N, M) hyperparameter sweep.

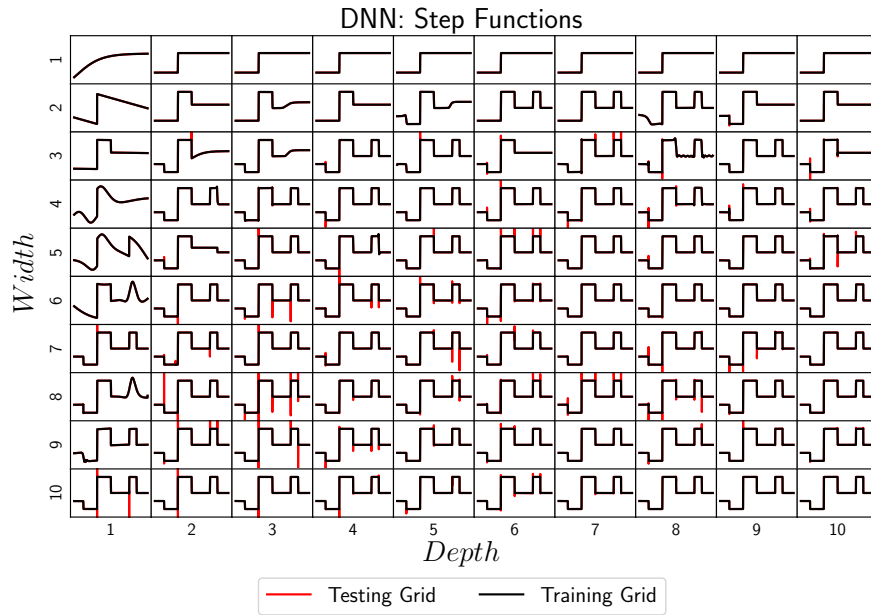


FIG. SM4.2. Plot of best approximation with DNN on the step functions over $(Width, Depth)$ hyperparameter sweep.

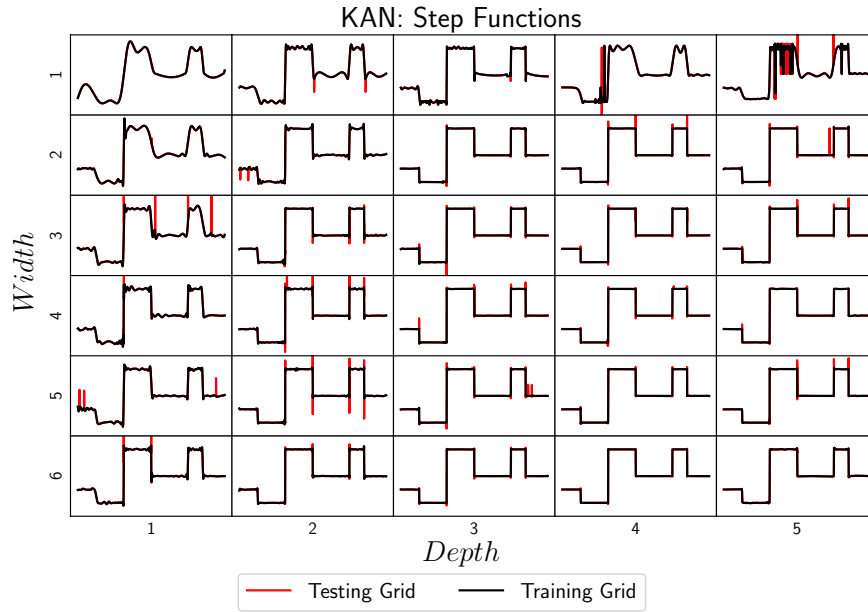


FIG. SM4.3. Plot of best approximation with KAN on the step functions over (N, M) hyperparameter sweep.

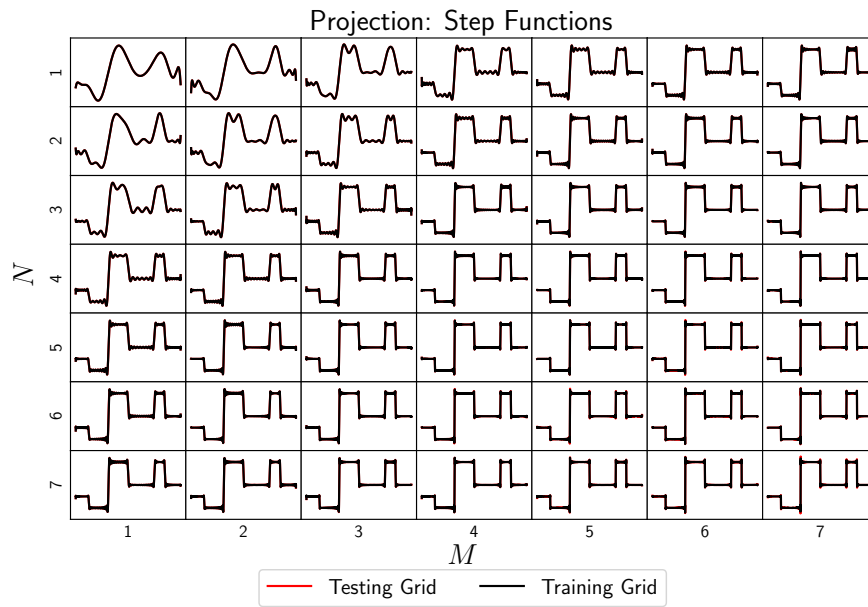


FIG. SM4.4. Plot of best approximation with polynomial projection on the step functions over (N, M) hyperparameter sweep.

SM4.1.1. Zoomed-in Discontinuities. Figures SM4.5 to SM4.7 provide a more nuanced perspective to the discontinuity approximation and how it drives the lack of L_∞ convergence in Figure SM1.1. Consider the initialization 0 results in Figure SM4.5 for the DNN on the test discretization. Despite the approximation perfectly fitting the training points, the method has no knowledge of the analytical location of the discontinuity, so it is misaligned with the target which is observable at higher sampling. Because of this, the L_∞ error for the DNN in this case is 4. Regardless of the spiking behavior observed in all NN-based models (which can be more clearly seen in Figures SM4.5 to SM4.7), they have no chance of ever converging in L_∞ due to simply lacking the exact location of the discontinuity at higher discretizations. Other noteworthy observations from these plots is how smooth polynomial projection is compared to the NN-based methods when zoomed in, as well as the fact KAN sometimes fails to pass through the training points as seen in Figure SM4.7 initialization 0.

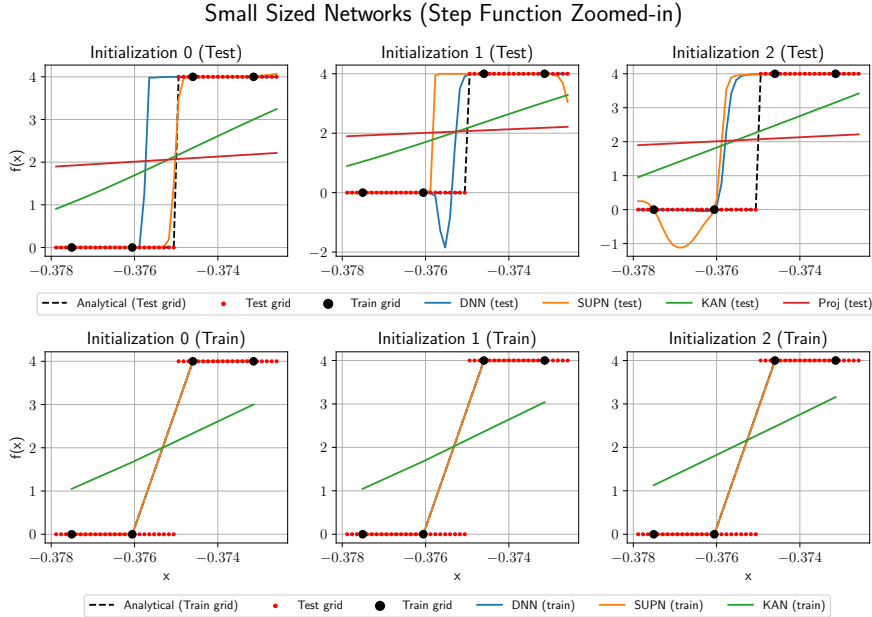


FIG. SM4.5. Plot of zoomed in discontinuity at $x = -0.375$ for the step function target. The network approximations shown are drawn from the **5th largest** width/depth or M/N models shown in Figures SM4.1 to SM4.3, i.e., along the diagonal, for three different model initializations.

SM4.2. Continuous Rastrigin. Figures SM4.8 to SM4.11 visually show the approximation methods as a function of architecture size at both the training and testing discretizations. Unlike Subsection SM4.1, we see very little generalization error between the training and testing discretizations, and spiking is not induced at higher discretizations. This gives credence to the idea that the NN-based models learn extreme basis functions to fit the discontinuity more accurately than polynomial projection. For continuous functions, this phenomena is not observed and models converge with respect to their expressability under proper optimization.

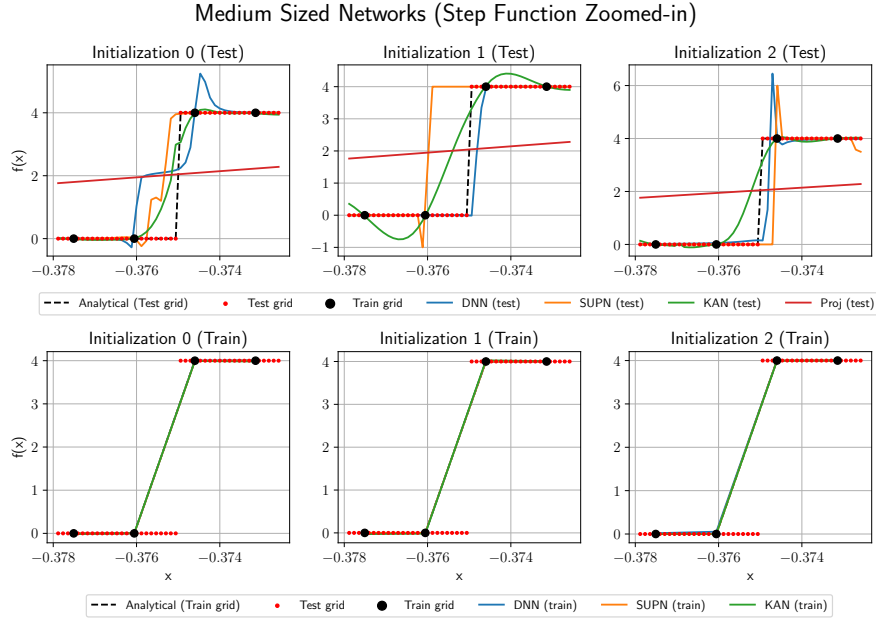


FIG. SM4.6. Plot of zoomed in discontinuity at $x = -0.375$ for the step function target. The network approximations shown are drawn from the **3rd largest** width/depth or M/N models shown in [Figures SM4.1 to SM4.3](#), i.e., along the diagonal, for three different model initializations.

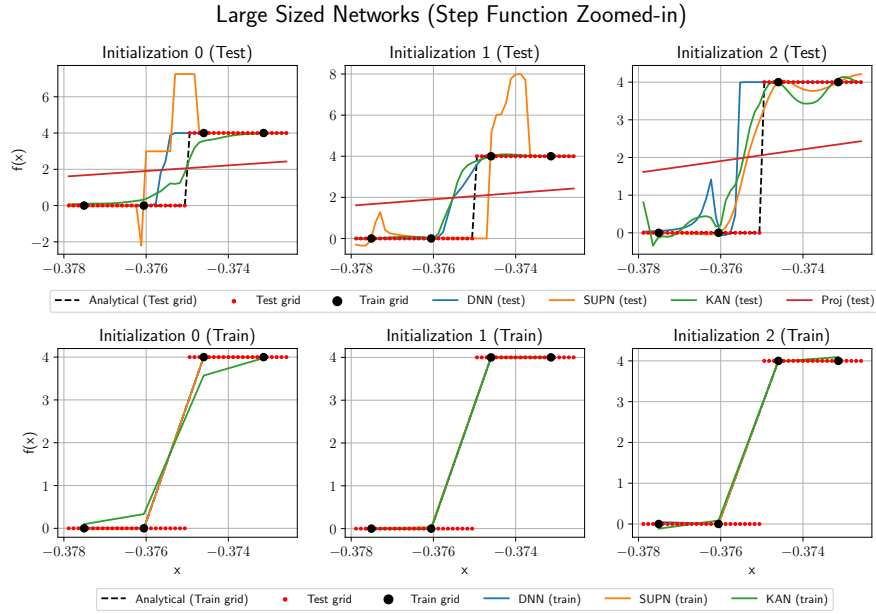


FIG. SM4.7. Plot of zoomed in discontinuity at $x = -0.375$ for the step function target. The network approximations shown are drawn from the **largest** width/depth or M/N models shown in [Figures SM4.1 to SM4.3](#), i.e., along the diagonal, for three different model initializations.

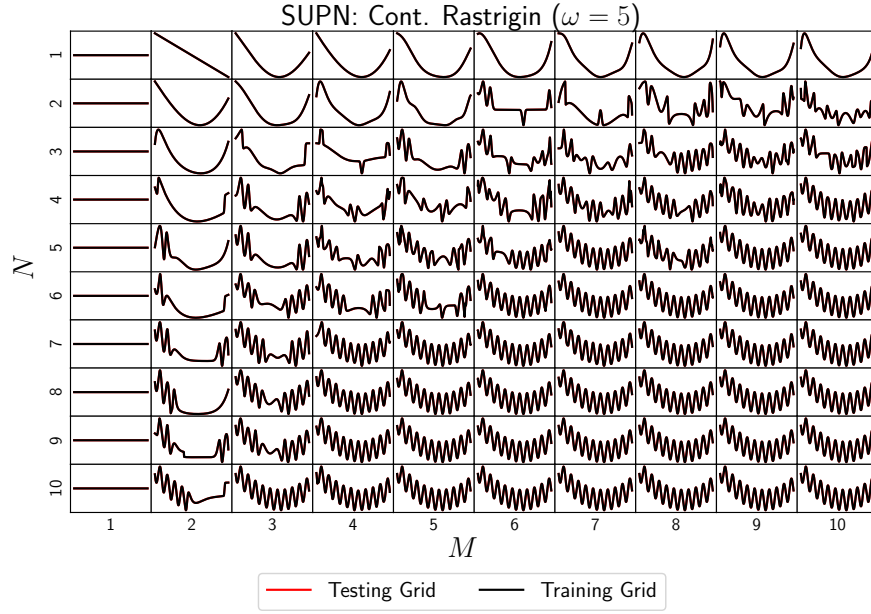


FIG. SM4.8. Plot of best approximation with SUPN on the continuous Rastrigin function over (N, M) hyperparameter sweep.

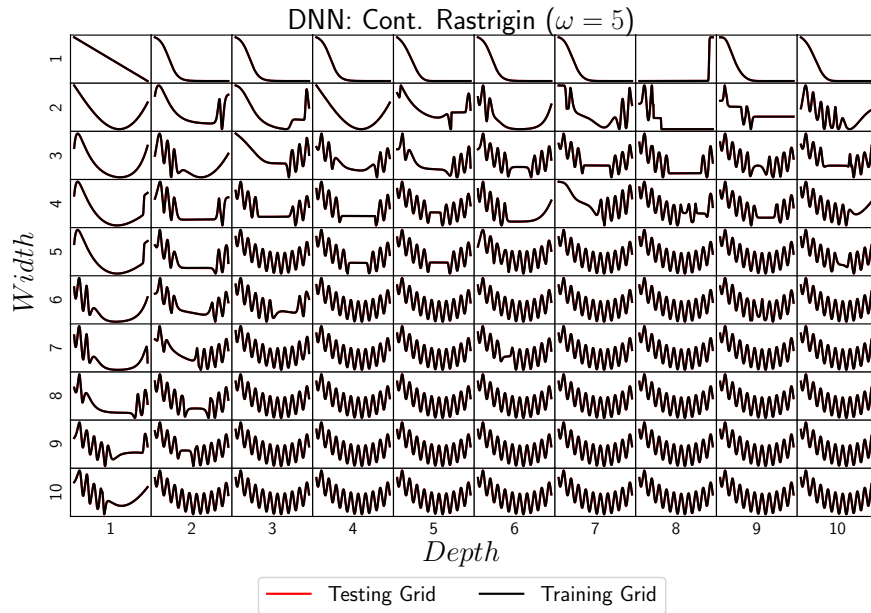


FIG. SM4.9. Plot of best approximation with DNN on the continuous Rastrigin function over $(Width, Depth)$ hyperparameter sweep.

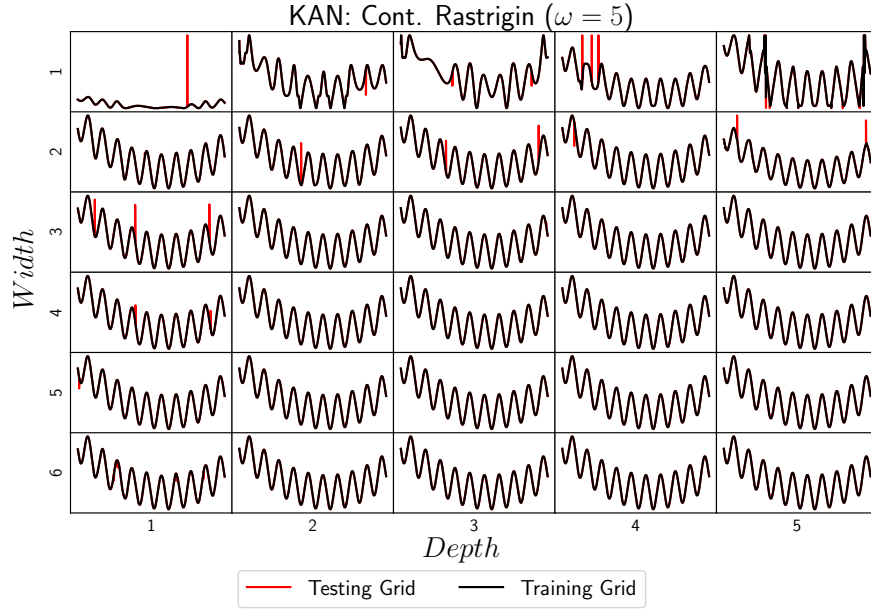


FIG. SM4.10. Plot of best approximation with KAN on the continuous Rastrigin function over $(\text{Width}, \text{Depth})$ hyperparameter sweep.

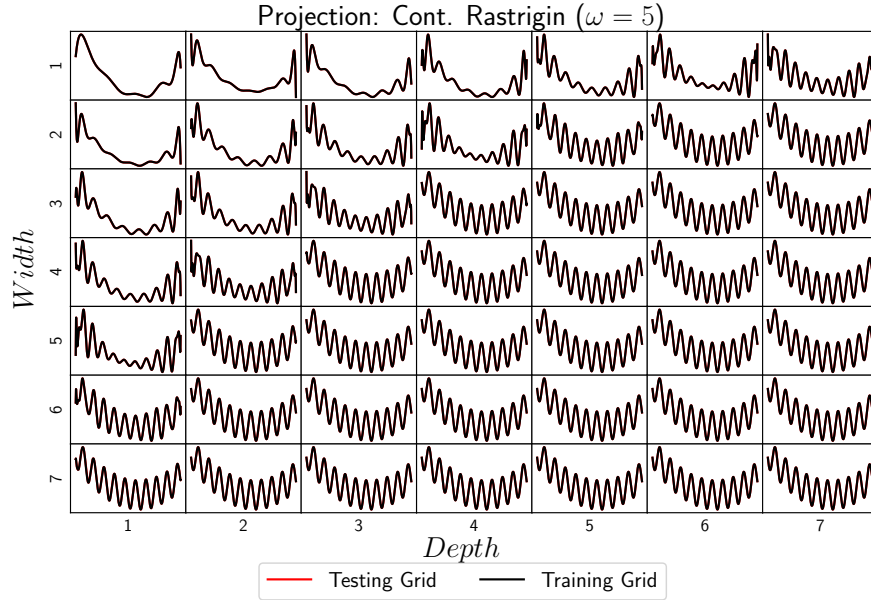


FIG. SM4.11. Plot of best approximation with polynomial projection on the continuous Rastrigin function over (N, M) hyperparameter sweep. Each grid cell uses $N(M + 1)$ terms, i.e. the same as an SUPN.

SM5. Sobolev space SUPN estimates. We provide the proof of Proposition 3.4 for $k > 0$ and $\epsilon_\Lambda(f) > 0$. This proof is conceptually the same as the $k = 0$ proof, but involves more opaque constants. By definition of $\epsilon_\Lambda(f)$, there is some polynomial $\tilde{q} \in \mathbb{P}_\Lambda$ satisfying (3.3). Let $\tilde{q} = \sum_{\mathbf{m} \in \Lambda} \alpha_{\mathbf{m}} T_{\mathbf{m}}$, and define

$$R := \|\tilde{q}\|_{W^{k,\infty}(\Omega)} = \max_{|\beta| \leq k} \sup_{x \in \Omega} \left| \frac{\partial^\beta}{\partial \mathbf{x}^\beta} \tilde{q}(x) \right| < \infty.$$

If $R = 0$, set $c_1 = 0$ in (3.2) with $N = 1$ to achieve $f_{N,\Lambda} = q = 0$, achieving (3.4). Otherwise, consider $R > 0$.

With $\sigma(y) = \tanh y$, we have that,

$$\sigma^{(n)}(0) = \begin{cases} 0, & n \text{ even} \\ \frac{2^{n+1}(2^{n+1}-1)B_{n+1}}{n+1}, & n \text{ odd} \end{cases}$$

where B_n is the n th Bernoulli number. An additional property of σ we require is that there are absolute numbers, (c_n, d_n) , such that,

$$\left| \sigma^{(n)}(y) - \sigma^{(n)}(0) \right| \leq c_n |y|, \quad |y| \leq d_n \leq 1.$$

For example, the above is true for any function with a bounded derivative of order $n + 2$, which includes the \tanh function. To make computations simpler, we use the looser but sufficient estimate,

$$\left| \sigma^{(n)}(y) \right| \leq a_n, \quad |y| \leq \tilde{d}_n,$$

where $a_n = \max\{1, \sigma^{(n)}(0)\}$, where $\tilde{d}_n \in (0, d_n)$.

To compute an order- k derivative of an SUPN, Faà di Bruno's formula states that, for a function q with sufficient differentiable smoothness,

$$\begin{aligned} \frac{d^k}{dx^k} \sigma(q(x)) &= \sum_{\ell=0}^k \sigma^{(\ell)}(q(x)) B_{k,\ell} \left(q^{(1)}, \dots, q^{(k-\ell+1)} \right) \\ &= \sigma'(q(x)) B_{k,1} \left(q^{(1)}, \dots, q^{(k)} \right) + \sum_{\substack{\ell=0 \\ \ell \neq 1}}^k \sigma^{(\ell)}(q(x)) B_{k,\ell} \left(q^{(1)}, \dots, q^{(k-\ell+1)} \right), \end{aligned}$$

where $B_{k,\ell}$ is the k -variate (exponential) incomplete Bell polynomial of degree ℓ . We note that $B_{k,1}(s_1, \dots, s_k) = s_k$ for all k , and $B_{k,0} = 0$ for any $k > 0$, so that our formula simplifies to,

$$\frac{d^k}{dx^k} \sigma(q(x)) = \sigma'(q(x)) q^{(k)}(x) + \sum_{\ell=2}^k \sigma^{(\ell)}(q(x)) B_{k,\ell} \left(q^{(1)}, \dots, q^{(k-\ell+1)} \right),$$

We will need suprema of the incomplete Bell polynomials over a compact isotropic hypercube of edglength R :

$$\overline{B}_k := \max_{\ell=0,\dots,k} \|B_{k,\ell}(z_1, \dots, z_k)\|_{L^\infty([-R,R]^k)}.$$

A final property of the Bell polynomials we need is that for $r \geq 2$, $B_{r,\ell}(z_1, \dots, z_r)$ is homogeneous of order ℓ , i.e., $B_{r,\ell}(cz_1, \dots, cz_r) = c^\ell B_{r,\ell}(z_1, \dots, z_r)$.

For simplicity, we take the multi-index β specifying the derivative to be a (multiple of a) cardinal unit vector, i.e., $\beta = r\mathbf{e}_s$, where $r \in [1, k]$, and $\mathbf{e}_s \in \mathbb{R}^d$ is the cardinal unit vector in direction s . Then $\partial_\beta f = \partial_s^r f$. For any $r \in [1, k]$,

$$\begin{aligned} \frac{\partial^r}{\partial x_s^r} \sigma(q(\mathbf{x})) - q^{(r,s)}(\mathbf{x}) &= \underbrace{(\sigma'(q(x)) - 1) q^{(r)}(x)}_{(A)} \\ &\quad + \underbrace{\sum_{\ell=2}^r \sigma^{(\ell)}(q(x)) B_{r,\ell} \left(q^{(1,s)}, \dots, q^{(r-\ell+1,s)} \right)}_{(B)} \end{aligned}$$

where for shorthand we've written $q^{(r,s)}$ to denote the r -th partial derivative of q in the s th variable. Now, define q as,

$$q(x) = \frac{1}{S} \tilde{q}(x), \quad S = \frac{1}{\sqrt{\delta\epsilon_\Lambda(f)}} \max \left\{ 1, \frac{R}{\min\{\tilde{d}_1, \dots, \tilde{d}_k\}}, R^{3/2}, \bar{B}_r \max\{a_2, \dots, a_r\} \right\}$$

Then:

$$\begin{aligned} |(A)| &\leq |\sigma'(q(x)) - 1| R/S \leq |q(x)|^2 R/S \leq R^3/S^3 \leq \delta\epsilon_\Lambda(f) \frac{1}{S} \\ |(B)| &\leq \sum_{\ell=2}^r |\sigma^{(\ell)}(q(x))| \left| B_{r,\ell} \left(q^{(1,s)}, \dots, q^{(r-\ell+1,s)} \right) \right| \\ &\leq \sum_{\ell=2}^r a_\ell \left| B_{r,\ell} \left(q^{(1,s)}, \dots, q^{(r-\ell+1,s)} \right) \right| \\ &\leq \sum_{\ell=2}^r a_\ell S^{-\ell} \left| B_{r,\ell} \left(\tilde{q}^{(1,s)}, \dots, \tilde{q}^{(r-\ell+1,s)} \right) \right| \\ &\leq \sum_{\ell=2}^r a_\ell S^{-\ell} \bar{B}_r \\ &\leq \frac{1}{S} \sum_{\ell=2}^r a_\ell S^{-1} \bar{B}_r \\ &\leq \delta\epsilon_\Lambda(f) \frac{1}{S} \end{aligned}$$

The same estimate holds for all β satisfying $|\beta| \leq k$ by taking iterated derivatives in each variable. Now consider the $(N, M) = (1, |\Lambda|)$ SUPN in (3.2), with the parameter assignment

$$a_{1,\mathbf{m}} = \frac{\alpha_{\mathbf{m}}}{S}, \quad c_1 = S,$$

so that $f_{N,\Lambda}(\mathbf{x}) = S\sigma(q(\mathbf{x}))$. We have shown, for any $\mathbf{x} \in \Omega$ and $|\beta| \leq k$,

$$\left| \frac{\partial^\beta}{\partial \mathbf{x}^\beta} f_{N,\Lambda}(\mathbf{x}) - \frac{\partial^\beta}{\partial \mathbf{x}^\beta} \tilde{q}(\mathbf{x}) \right| = S \left| \frac{\partial^\beta}{\partial \mathbf{x}^\beta} \sigma(q(\mathbf{x})) - \frac{\partial^\beta}{\partial \mathbf{x}^\beta} q(\mathbf{x}) \right| \leq \delta\epsilon_\Lambda(f),$$

which completes the proof.