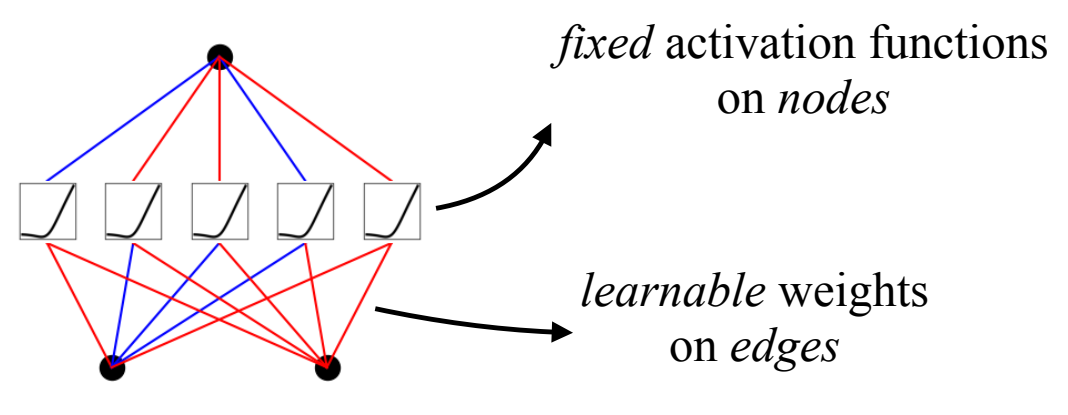
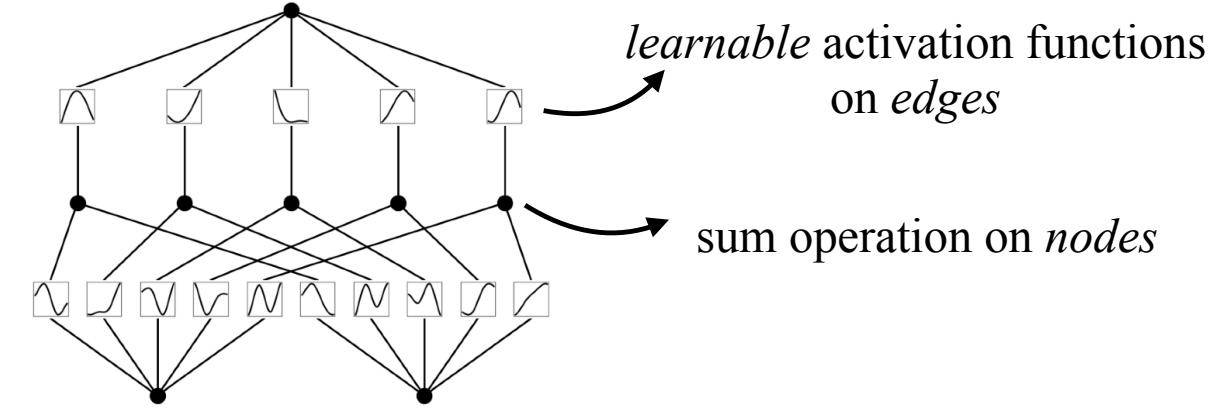
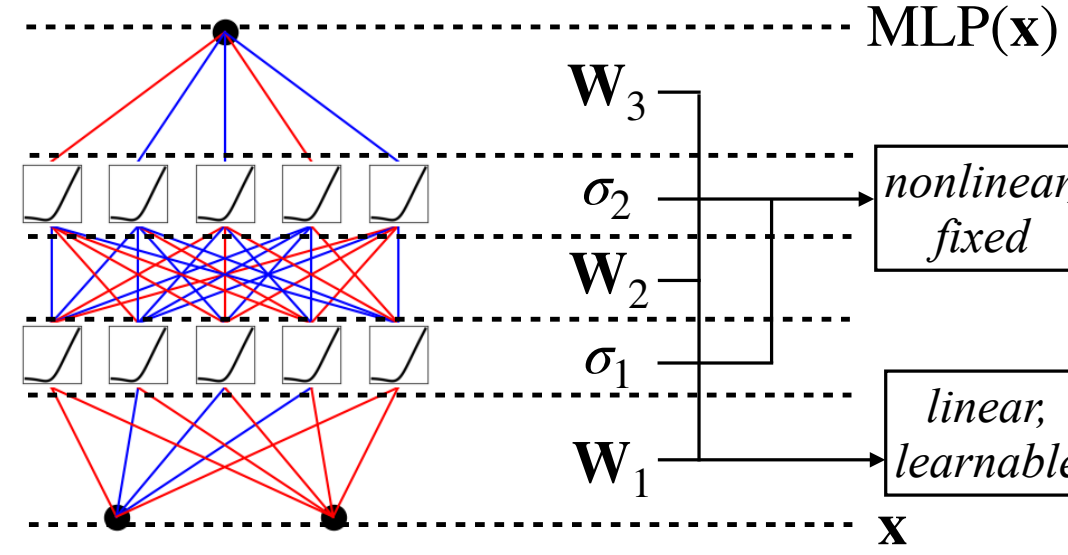
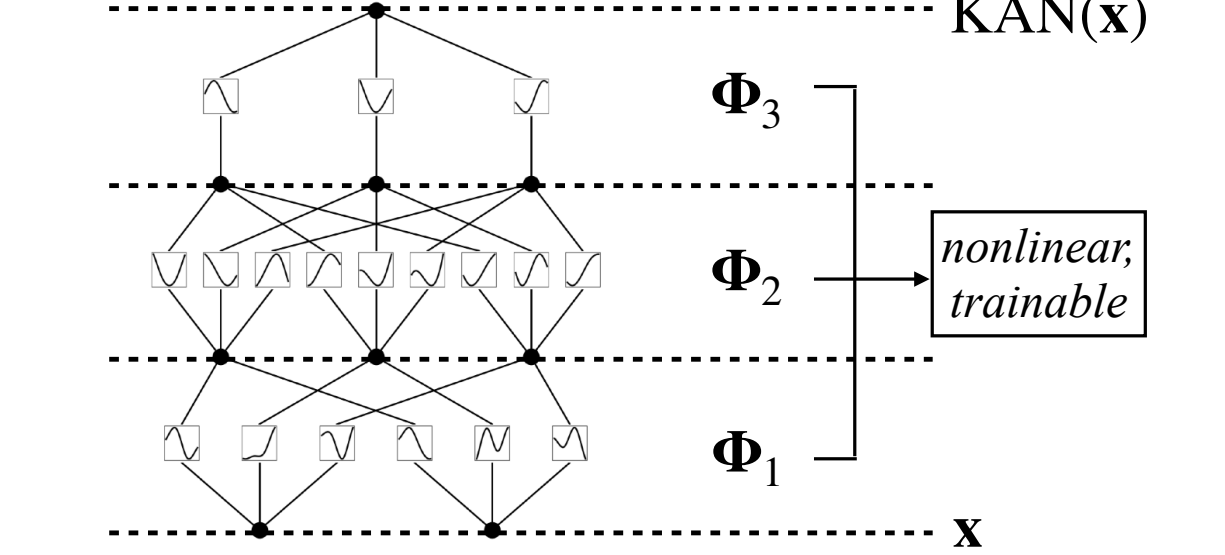


Kolmogorov–Arnold Networks (KANs) [2]

Model	Multi-Layer Perceptron (MLP)	Kolmogorov-Arnold Network (KAN)
Theorem	Universal Approximation Theorem	Kolmogorov-Arnold Representation Theorem
Formula (Shallow)	$f(\mathbf{x}) \approx \sum_{i=1}^{N(e)} a_i \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i)$	$f(\mathbf{x}) = \sum_{q=0}^{2n} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right)$
Model (Shallow)	(a) 	(b) 
Formula (Deep)	$\text{MLP}(\mathbf{x}) = (\mathbf{W}_3 \circ \sigma_2 \circ \mathbf{W}_2 \circ \sigma_1 \circ \mathbf{W}_1)(\mathbf{x})$	$\text{KAN}(\mathbf{x}) = (\Phi_3 \circ \Phi_2 \circ \Phi_1)(\mathbf{x})$
Model (Deep)	(c) 	(d) 

For a continuous $f : [0, 1]^n \rightarrow \mathbb{R}$

$$f(\mathbf{x}) = f(x_1, \dots, x_n) = \sum_{q=0}^{2n} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right). \quad (1)$$

where $\phi_{q,p} : [0, 1] \rightarrow \mathbb{R}$ and $\Phi_q : \mathbb{R} \rightarrow \mathbb{R}$ are continuous.

- Summing and composition of univariate functions. Potentially address the **curse of dimensionality** (COD).
- Φ_q and $\phi_{q,p}$ not necessarily smooth. We may need more than two layers.

We parametrize the learnable activation functions by B-splines.

Approximation Theory

Suppose that a function $f(\mathbf{x})$ admits a smooth representation

$$f = (\Phi_{L-1} \circ \Phi_{L-2} \circ \dots \circ \Phi_1 \circ \Phi_0)\mathbf{x}, \quad (2)$$

where $\Phi_{l,i,j}$ are smooth with derivatives uniformly bounded up to $k+1$ -th order. Then using k -th order B-splines with $G+1$ grid points as activation functions, there exist $\Phi_{l,i,j}^G$ such that for any $0 \leq m \leq k$, we have the bound

$$\|f - (\Phi_{L-1}^G \circ \Phi_{L-2}^G \circ \dots \circ \Phi_1^G \circ \Phi_0^G)\mathbf{x}\|_{C^m} \leq CG^{-k-1+m}. \quad (3)$$

In particular for L^2 or RMSE, we have the scaling law $k+1$. Informally, such functions are dense in the class of continuous functions, by [1].

Leveraging the 1D structure to get better scaling laws

Expressiveness and Spectral bias: KANs and MLPs

- ReLU- k MLP with width W , depth L can be represented by Spline- k -KAN with width W , depth $2L$, grid size 2.
- k -KAN with width W , depth L , grid size G can be represented by ReLU- k MLP with width $(G+2k+1)W^2$, depth $2L$.
- $O(G^2W^4L)$ parameter count for MLP, $O(GW^2L)$ for KAN in this formulation.
- Sharp if we restrict the depth of MLP.
- MLPs have difficulty learning high frequencies, i.e. have spectral bias.
- Splines, as building blocks of KANs, have no spectral bias.

Function Fitting

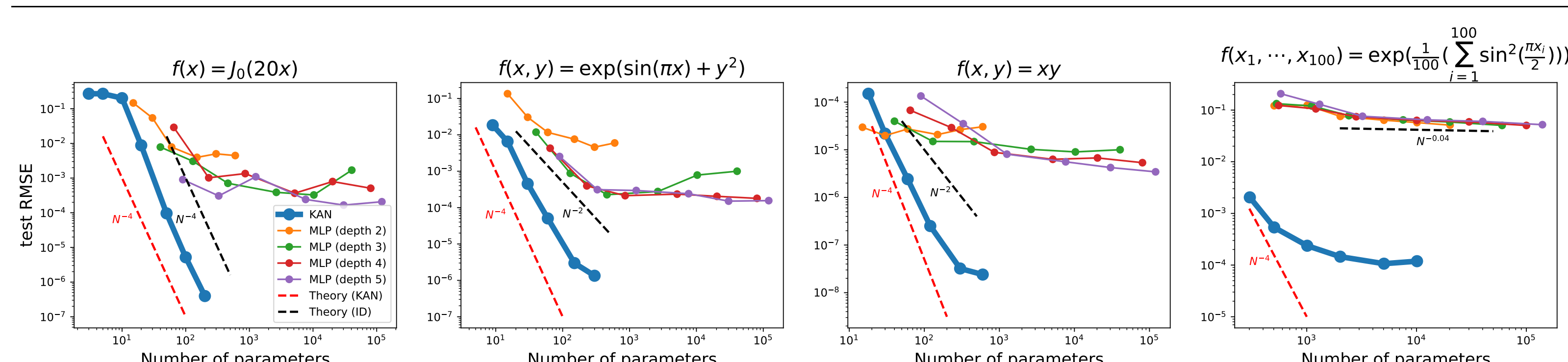


Figure 1. KANs almost saturate the fastest scaling law by theory ($\alpha = 4$), while MLPs scale slowly and plateau.

<https://roywangyx.github.io/>

Function Fitting: 1D Waves of Different Frequencies

$$f(x) = \sum A_i \sin(2\pi k_i x + \varphi_i), \quad k = (5, 10, \dots, 45, 50).$$

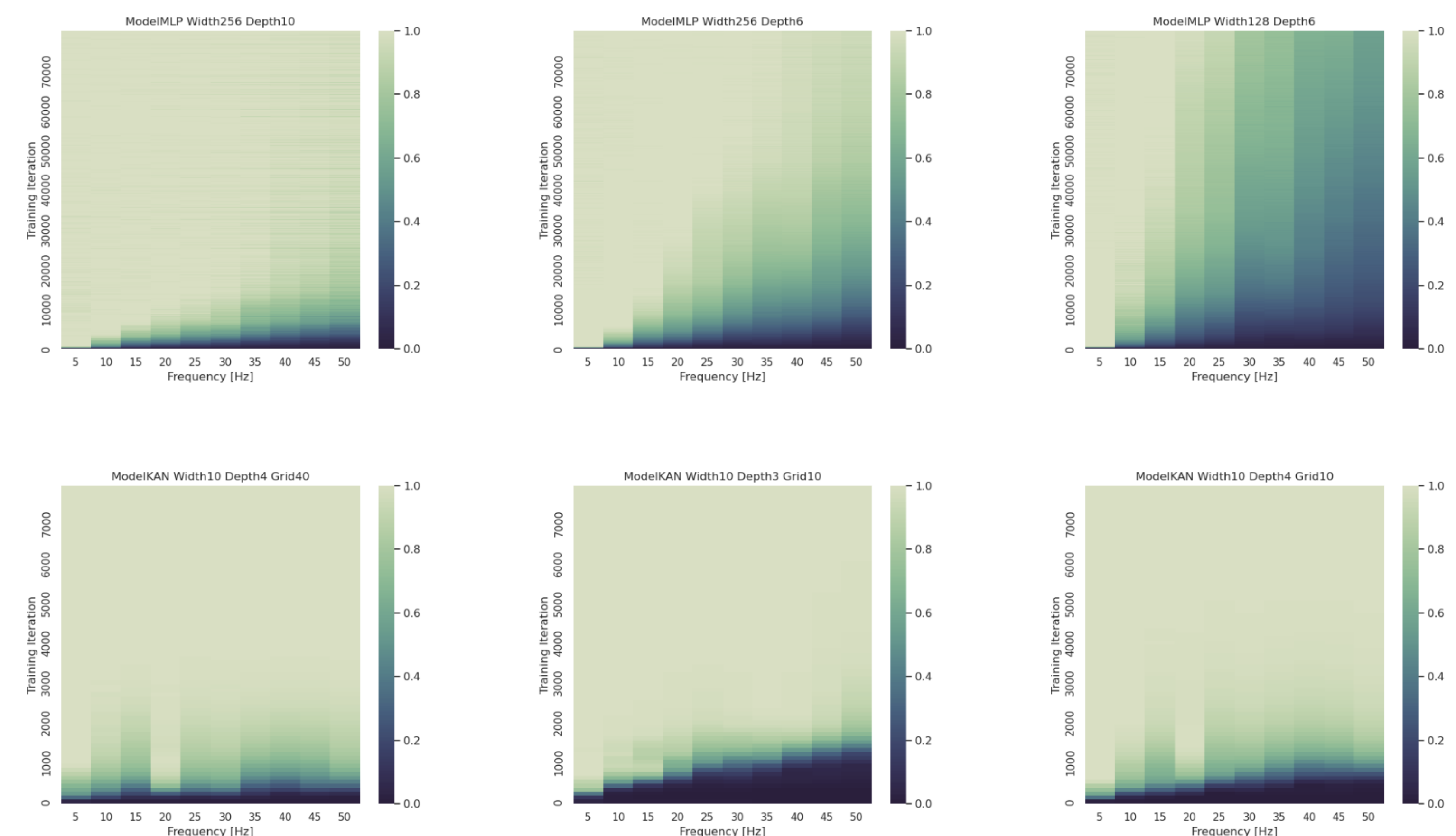


Figure 2. KANs learn high frequencies much faster, even with 10x fewer epochs.

PDE Solving: 2D Poisson of Different Frequencies

MLP of width 256, depth 6; KAN of width 10, depth 2, grid size 20.

$$-\Delta u = f, \quad u = \sin(\pi x) \sin(\pi y) + \frac{1}{k} \sin(k\pi x) \sin(k\pi y).$$

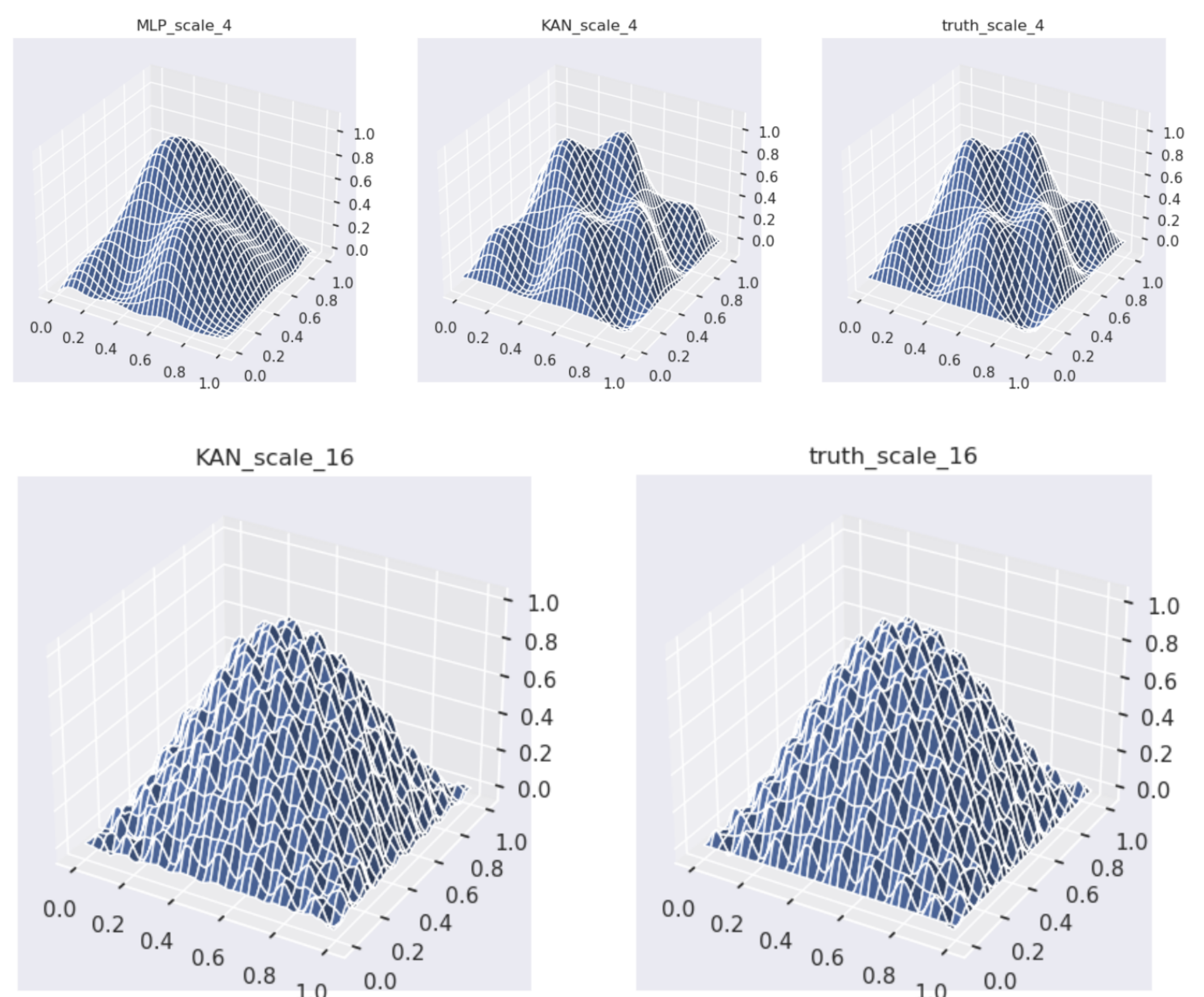


Figure 3. MLPs struggle to learn high frequency information without frequency encoding.

References

- Ming-Jun Lai and Zhaiming Shen. The kolmogorov superposition theorem can break the curse of dimensionality when approximating high dimensional functions. *arXiv preprint arXiv:2112.09963*, 2021.
- Ziming Liu, Yixuan Wang, Sachin Vaidya, Fabian Ruehle, James Halverson, Marin Soljačić, Thomas Y Hou, and Max Tegmark. Kan: Kolmogorov-arnold networks. *arXiv preprint arXiv:2404.19756, ICLR 2025 oral*, 2024.
- Yixuan Wang, Jonathan W Siegel, Ziming Liu, and Thomas Y Hou. On the expressiveness and spectral bias of kans. *arXiv preprint arXiv:2410.01803, ICLR 2025*, 2024.



ICLR