



# KAN: Kolmogorov-Arnold Networks

Ziming Liu<sup>1</sup> Yixuan Wang<sup>2</sup> Sachin Vaidya<sup>1</sup> Fabian Ruehle<sup>3</sup>  
Jim Halverson<sup>3</sup> Marin Soljacic<sup>1</sup> Thomas Y. Hou<sup>2</sup> Max Tegmark<sup>1</sup>  
<sup>1</sup>MIT <sup>2</sup>Caltech <sup>3</sup>Northeastern ICLR April 24-28, 2025



## Kolmogorov-Arnold Representation Theorem

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).
- $\Phi_q$  and  $\phi_{q,p}$  not necessarily smooth. In practice we may need more than two layers.

## Kolmogorov-Arnold Networks (KANs)

| Model             | Multi-Layer Perceptron (MLP)  | Kolmogorov-Arnold Network (KAN)  |
|-------------------|---|--|
| Theorem           | Universal Approximation Theorem   | Kolmogorov-Arnold Representation Theorem   |
| Formula (Shallow) | $f(\mathbf{x}) \approx \sum_{i=1}^{N(\epsilon)} a_i \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i)$                          | $f(\mathbf{x}) = \sum_{q=0}^{2n} \Phi_q \left( \sum_{p=1}^n \phi_{q,p}(x_p) \right)$ |
| 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)  |

We parametrize the learnable activation functions by B-splines.

## Approximation Theory

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

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

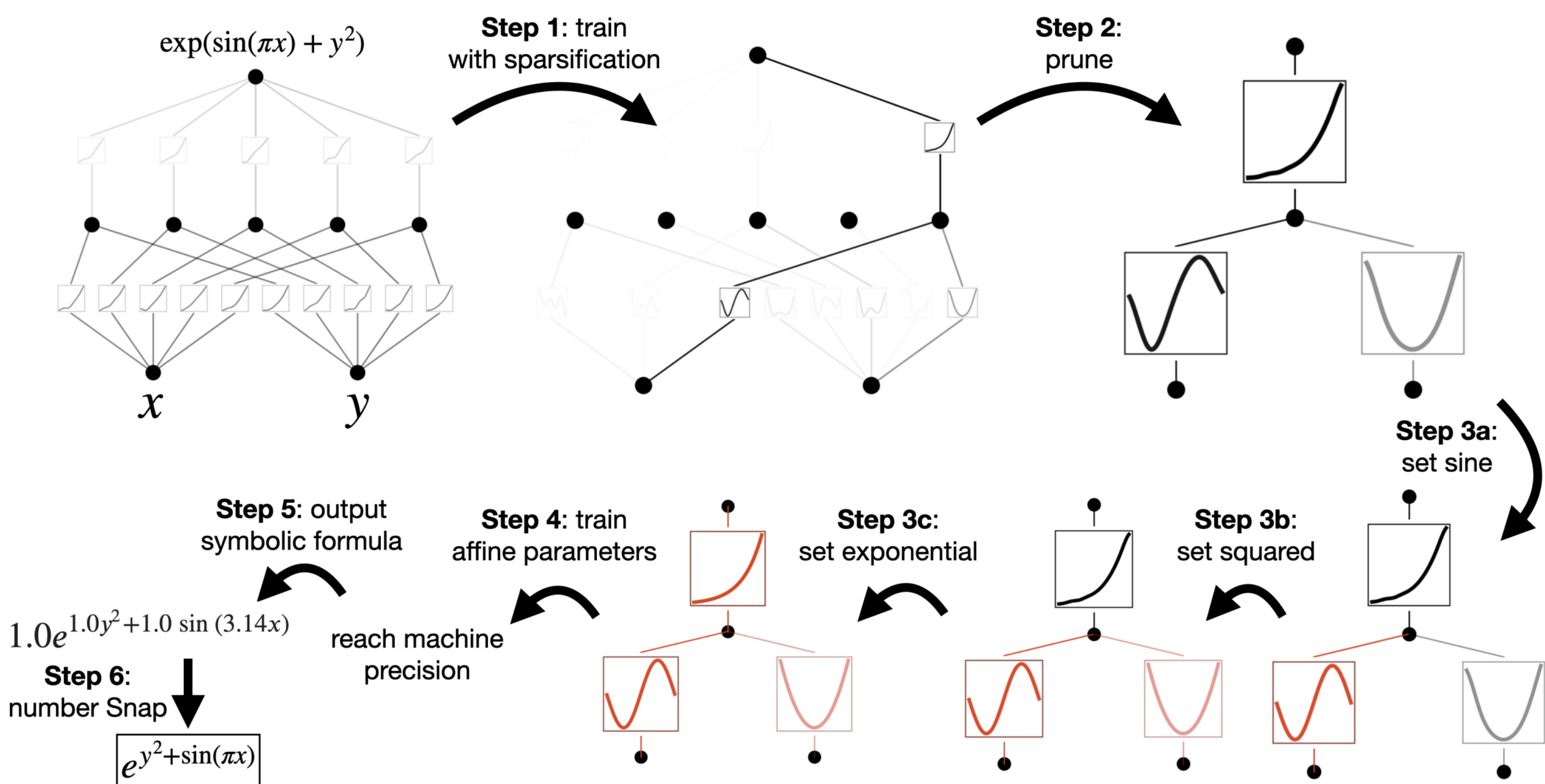
where  $\Phi_{l,i,j}$  are smooth with derivatives uniformly bounded up to  $k+1$ -th order. Then using  $k$ -th order B-splines with  $G+1$  grid points as activation functions, there exist  $\Phi_{l,i,j}^G$  such that for any  $0 \leq m \leq k$ , we have the bound

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

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

## KAN for Interpretability: Symbolic Training [2]



## Function Fitting

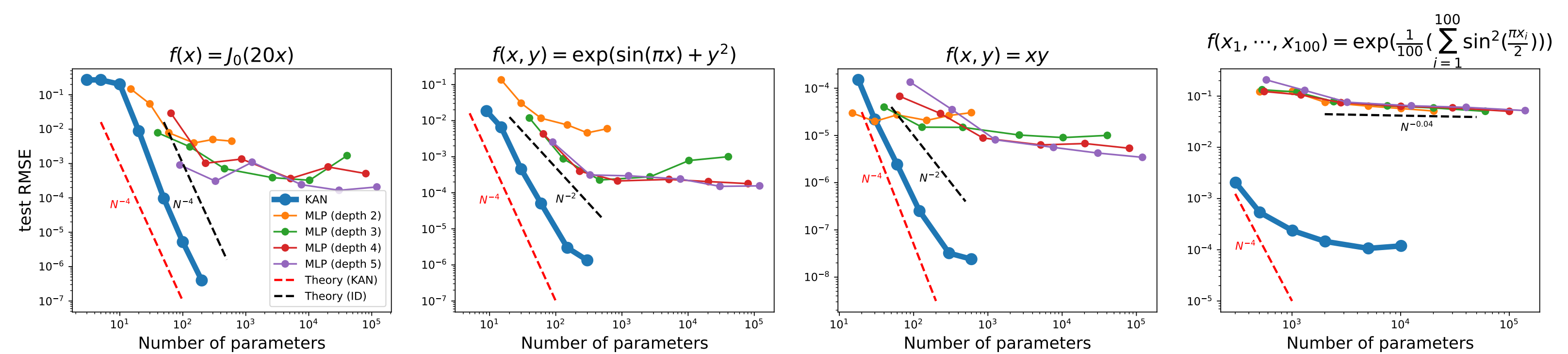


Figure 1. KANs almost saturate the fastest scaling law by theory ( $\alpha = 4$ ), while MLPs scale slowly and plateau.

## Image Fitting



Figure 2. KANs outperform MLP with frequency encoding tricks, due to the ability to capture high frequency [4].

## Scaling up KANs

| Problem       | Model  | PSNR / L2 Error      | Training Time (s) |
|---------------|--|----------------------|-------------------|
| Image Fitting | KAN [2,128,128,128,128,1], G=[100,10,10,10,10] | <b>45.76</b>         | 1809              |
| Image Fitting | MLP [2,404,404,404,404,1]                      | 22.09                | <b>182</b>        |
| Image Fitting | SIREN 1 [2,128,128,128,128,1]                  | 27.34                | 254               |
| Image Fitting | SIREN 2 [2,404,404,404,404,1]                  | 30.79                | 407               |
| Image Fitting | MLP_RFF [2,404,404,404,404,1]                  | 26.26                | 195               |
| Allen-Cahn    | KAN [2,5,5,1], G=5                             | $3.4 \times 10^{-3}$ | 2801              |
| Allen-Cahn    | MLP [2,128,128,128,1]                          | $1.5 \times 10^{-1}$ | <b>478</b>        |
| Allen-Cahn    | MLP [2,128,128,128,1] (10x training)           | $3.9 \times 10^{-4}$ | 4766              |
| Darcy Flow    | KAN [2,10,1], G=20                             | $3.9 \times 10^{-4}$ | 66                |
| Darcy Flow    | KAN [2,100,1], G=10                            | $4.3 \times 10^{-6}$ | 107               |
| Darcy Flow    | KAN [2,10,10,10,10,1], G=5                     | $8.5 \times 10^{-5}$ | 123               |
| Darcy Flow    | MLP [2,128,128,128,1]                          | $3.0 \times 10^{-5}$ | <b>30</b>         |
| Darcy Flow    | MLP [2,128,128,128,1] (10x training)           | $4.5 \times 10^{-6}$ | 277               |
| Darcy Flow    | MLP_RFF [2,128,128,128,1]                      | $5.9 \times 10^{-6}$ | 31                |

Figure 3. KANs can scale up on GPUs and are fast. Examples of image fitting and PDE solving using Adam.

## 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, Pingchuan Ma, Yixuan Wang, Wojciech Matusik, and Max Tegmark. Kan 2.0: Kolmogorov-arnold networks meet science. *arXiv preprint arXiv:2408.10205*, 2024.
- 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.

