Multi-Exit Kolmogorov–Arnold Networks: enhancing accuracy and parsimony

James Bagrow^{1,2,*} and Josh Bongard^{3,2}

¹Mathematics & Statistics, University of Vermont, Burlington, VT, United States

²Vermont Complex Systems Center, University of Vermont, Burlington, VT, United States

³Computer Science, University of Vermont, Burlington, VT, United States

*Corresponding author. Email: james.bagrow@uvm.edu, Homepage: bagrow.com

June 4, 2025

Kolmogorov-Arnold Networks (KANs) uniquely combine high accuracy with interpretability, making them valuable for scientific modeling. However, it is unclear a priori how deep a network needs to be for any given task, and deeper KANs can be difficult to optimize and interpret. Here we introduce multi-exit KANs, where each layer includes its own prediction branch, enabling the network to make accurate predictions at multiple depths simultaneously. This architecture provides deep supervision that improves training while discovering the right level of model complexity for each task. Multi-exit KANs consistently outperform standard, singleexit versions on synthetic functions, dynamical systems, and real-world datasets. Remarkably, the best predictions often come from earlier, simpler exits, revealing that these networks naturally identify smaller, more parsimonious and interpretable models without sacrificing accuracy. To automate this discovery, we develop a differentiable "learning to exit" algorithm that balances contributions from exits during training. Our approach offers scientists a practical way to achieve both high performance and interpretability, addressing a fundamental challenge in machine learning for scientific discovery.

Keywords— multi-exit and early-exit neural networks, scientific machine learning, interpretability and explainability, deep supervision, data-driven models, dynamical systems

1 Introduction

Machine learning has become indispensable for scientific discovery, enabling researchers to uncover complex patterns in data that traditional analytical methods cannot readily capture [1, 2, 3, 4]. Neural networks and related techniques excel at nonlinear regression and data-driven modeling of dynamical systems—fundamental challenges spanning physics, biology, chemistry, and engineering [1, 5, 4]. From climate modeling [6, 7] to protein folding prediction [8], these data-driven approaches can learn sophisticated functional relationships directly from observations, opening new pathways to understanding natural phenomena where first-principles models are unavailable or computationally intractable [9, 2, 10, 4].

Despite these advances, scientific applications face a fundamental challenge that is less pressing in many other machine learning domains: the need to achieve simultaneously both high predictive accuracy and model interpretability [2, 11]. While many commercial applications prioritize accuracy above all else, scientific modeling demands that researchers understand not just what the model predicts, but how and why it makes those predictions [2, 4]. Interpretability is essential for determining whether models capture genuine physical relationships rather than spurious correlations, for gaining scientific insight into the underlying phenomena, and for building the trust necessary to guide experimental design or inform policy decisions [12, 11, 13]. Unfortunately, accuracy and interpretability are in tension: the most accurate machine learning modelstypically deep neural networks with millions or billions of parameters-are often the least interpretable, functioning as "black boxes" that provide little insight into the mechanisms driving their predictions [14, 4]. This accuracy-interpretability trade-off remains a critical bottleneck for the adoption of machine learning in scientific discovery, where understanding the underlying relationships is often as important as making accurate predictions [2, 4].

Kolmogorov-Arnold Networks (KANs) have recently emerged as a promising solution to this accuracyinterpretability dilemma, representing one of the rare neural architectures that can achieve both high predictive performance and meaningful interpretability [15, 16, 17]. Motivated by the Kolmogorov-Arnold Representation Theorem, which shows that multivariate functions can be represented as compositions of univariate functions, KANs provide a divide-and-conquer approach to high-dimensional problems by breaking them down into manageable univariate components that can be learned directly from data [15]. This enables KANs to discover and represent complex functional relationships while maintaining the ability to visualize and interpret each learned univariate function individually. Results have demonstrated that KANs can achieve competitive accuracy with traditional deep networks on regression tasks and dynamical systems modeling while providing insights into the learned functional forms [15, 16, 18, 19]. This combination of accuracy and interpretability makes KANs well-suited for scientific applications where understanding the underlying mathematical relationships is as crucial as predictive performance.

Despite these promising qualities, KANs face challenges that can limit their effectiveness in scientific applications. Training KANs presents optimization difficulties, as the learnable univariate functions must be carefully parameterized and refined, with deeper networks often proving particularly challenging to optimize [15]. The architecture search problem—determining the appropriate number of layers and widths-also remains significant, as practitioners must balance expressiveness against parsimony while avoiding overfitting [20]. Seeking smaller, more parsimonious models without sacrificing accuracy is crucial because overly deep KANs lose interpretability as compositions of many univariate functions become difficult to understand, while smaller KANs better retain the interpretability that makes them valuable for scientific modeling [15]. These challenges suggest that architectural innovations are needed to better realize KANs' potential for scientific discovery.

In this paper, we introduce *multi-exit architectures* [21, 22] into KANs as a novel approach to address these challenges while preserving KANs' interpretability advantages. Multiexit networks, originally developed for deep networks to enable adaptive inference and computational efficiency, augment networks with additional prediction branches at intermediate layers, allowing models to make predictions at multiple depths [22, 23, 24]. When applied to KANs, this approach offers a novel way to tackle the architecture search challenge by enabling a single network to effectively explore multiple levels of complexity simultaneously [24]. Multi-exit KANs can help identify appropriate levels of model complexity for a given task: if early exits perform well, the network has found a parsimonious and more interpretable model that maintains accuracy, while deeper exits remain available when additional expressiveness is needed (Fig. 1). Furthermore, the multiexit approach enables deep supervision [25] during training, where gradients flow directly to earlier layers through the exit branches, potentially improving the optimization of deeper KANs.

Experiments demonstrate the effectiveness of multi-exit KANs across diverse scientific modeling tasks, showing consistent improvements in both accuracy and parsimony compared to traditional single-exit KANs. Our key contributions include: (1) the first application of multi-exit architectures to KANs, with a joint training framework that enables deep supervision; (2) empirical validation across regression problems, dynamical systems modeling, continual learning scenarios, and real-world datasets; (3) evidence that multi-exit KANs often achieve better performance at earlier exits, indicating more parsimonious models without sacrificing accuracy; and (4) a "learning to exit" algorithm that addresses the choice of extra hyperparameters when using multi-exits. Additionally, we provide insights into why multi-exit architectures are particularly well-suited for KANs and discuss the mechanisms underlying their improved performance. These results suggest that multiexit architectures represent a valuable enhancement to KANs for scientific applications, offering a principled approach to balancing accuracy and interpretability.

The rest of this paper is organized as follows. Section 2 provides background on Kolmogorov–Arnold networks and

multi-exit and early-exit neural architectures. Section 3 describes our approach for incorporating multiple exits into the KAN architecture. Section 4 presents results comparing singleexit and multi-exit KANs across multiple domains. Section 5 introduces the learning to exit method. Finally, we conclude in Sec 6 by discussing the implications of our findings, including why multi-exit architectures improve KAN performance, and directions for future work.

2 Background

We consider the problems of learning unknown functions from data: nonlinear regression,

$$\mathbf{y} = F(\mathbf{x}) \tag{1}$$

for $\mathbf{y} \in \mathbb{R}^{n \times m}$ and $\mathbf{x} \in \mathbb{R}^{n \times d}$, and data-driven modeling of continuous or discrete dynamical systems of the form

$$\frac{d\mathbf{x}}{dt} = \mathbf{F}(\mathbf{x}) \tag{2}$$

or

$$\mathbf{x}_{n+1} = \mathbf{F}(\mathbf{x}_n),\tag{3}$$

respectively, and subject to problem-relevant initial and/or boundary conditions. KANs have proven effective for both problems [15, 19, 18].

2.1 Kolmogorov–Arnold networks

A KAN network is a multilayer feedforward neural network but unlike a traditional multilayer perceptron (MLP), the nonlinearities come from learnable, univariate activation functions (Fig. 1) associated with the connections between layers, and fixed summations (or multiplications) propagate signals between layers. In contrast, MLPs use fixed nonlinear activation functions on the units and learnable weights for summations associated with the connections between layers.

MLPs are motivated by the universal approximation theorem [26] while KANs are motivated by the KART, or *Kolmogorov–Arnold Representation Theorem* [27, 28, 29]: every multivariate continuous function on a finite domain can be expressed as a finite superposition of univariate continuous functions. More specifically, for $F : [0, 1]^d \rightarrow \mathbb{R}$,

$$F(\mathbf{x}) = F(x_1, x_2, \dots, x_d) = \sum_{q=1}^{2d+1} \Phi_q \left(\sum_{p=1}^d \phi_{q,p}(x_p) \right), \quad (4)$$

where $\phi_{q,p} : [0,1] \to \mathbb{R}$ and $\Phi_q : \mathbb{R} \to \mathbb{R}$. The KART shows that it is possible to represent any continuous function of multiple variables as a composition of one-dimensional functions. The innovation of KANs is to operationalize this by learning the univariate "activation" functions and stacking them in deep layers. As shown by Liu *et al.* [15], the stacking, which extends beyond the KART, often allows for smooth and interpretable activation functions which are not expected in the two-layer KART form given by Eq. (4).



Figure 1: Enhancing accuracy and interpretability of Kolmogorov–Arnold networks (KANs) with multiple exits, here illustrated on the toy problem $y = \sin(x_1) + \cos(x_2)$.

In a KAN, the activation functions are typically parameterized using B-splines, local polynomial approximations of the one-dimensional functions. Although many other functionfitting techniques have been considered, including radial basis functions [30], Fourier series [31], sinusoidal functions [32], Chebyshev polynomials [33], and wavelet-based representations [34], all of which have many advantages, particularly in terms of computational efficiency, for our purposes here we focus on the original B-spline approach.

Specifically, each activation function $\phi(x)$ is parameterized as a combination of a base function and a B-spline component:

$$\phi(x) = b(x) + \sum_{j} c_{j} B_{j}(x)$$
(5)

where b(x) is the base function, c_j are the learnable coefficients, and $B_j(x)$ are the B-spline basis functions. The base function serves as a residual connection similar to those in ResNets [25], facilitating gradient flow during training while allowing the B-spline component to learn nonlinear deviations. Common choices for the base function include the identity function b(x) = x, the SiLU function $b(x) = x/(1 + e^{-x})$, or the zero function b(x) = 0 when residual connections are omitted. Additionally, this formulation helps maintain the effectiveness of low-order B-splines even in deep networks by preventing their nested composition from creating numerically unstable high-order polynomials.

To form a deep network by stacking layers of learned activation functions, summation units (and, optionally, multiplication [16]) aggregate the outputs from the previous layer's

activation functions. The number of units across *L* layers is $[d = N_0, N_1, \ldots, N_L = m]$. The shape of the KAN is the vector of widths $[N_0, N_1, \ldots, N_L]$. (See Fig. 1 for an example of a KAN network with shape [2, 3, 2, 1].) Between layers *i* and *i* + 1 there are $N_i N_{i+1}$ activation functions. Following [15], $\phi_{\ell,j,i}(x)$ denotes the activation function connecting unit *i* in layer ℓ to unit *j* in layer $\ell + 1$ ($\ell = 0, \ldots, L - 1$; $i = 1, \ldots, N_\ell$; $j = 1, \ldots, N_{\ell+1}$). The summation units then combine the activation functions to propagate information through the network: for the signal $x_{\ell+1,j}$ into unit *j* in layer $\ell + 1$, giving

$$x_{\ell+1,j} = \sum_{i=1}^{N_{\ell}} \phi_{\ell,j,i}(x_{\ell,i}), \quad j = 1, \dots, n_{\ell+1}, \tag{6}$$

or, in matrix (broadcast) form, $\mathbf{x}_{\ell+1} = \Phi_{\ell}(\mathbf{x}_{\ell})$, where Φ_{ℓ} is the functional matrix containing the ϕ 's connecting layer ℓ and $\ell + 1$. Finally, the full network is represented by composing each Φ in sequence,

$$\operatorname{KAN}(\mathbf{x}) = (\Phi_{L-1} \circ \Phi_{L-2} \circ \cdots \circ \Phi_1 \circ \Phi_0) (\mathbf{x}).$$
(7)

It is this combination of superpositions of learned activation functions and layer-wise composition that gives KANs both their expressiveness and makes them distinct from MLPs. Beyond this architectural difference, KANs are considered more interpretable than MLPs because each activation function ϕ can be individually examined, as shown in Fig. 1.

KANs are trained with a loss function $\mathcal{L} = \mathcal{L}_{data} + \lambda \mathcal{L}_{reg}$ comprising a data loss and a regularization loss, with the hyperparameter λ determining regularization strength. Usually the data loss is mean squared error (MSE):

$$\mathcal{L}_{\text{data}} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2, \qquad (8)$$

where y_i is the true value for observation *i* and \hat{y}_i is the KAN's prediction for that observation, while the regularization loss is a combination of an L1 norm and an entropy both defined on the activation functions (for details, see Liu *et al.* [15]). Activation function parameters are learned via gradient-based optimization to minimize \mathcal{L} on training data. Training commonly includes a refinement process where the resolution of the B-spline grids is increased; gradually increasing grid resolution during training has been shown to improve KAN accuracy better than using a finer grid from the start [15]. KANs are typically optimized with quasi-Newton methods, particularly L-BFGS [35], also employed in this work, though first-order methods such as SGD or Adam [36] can be used as well.

2.2 Multi-exit and early exit networks

Multi-exit networks are a deep learning architecture where some or all hidden layers in the network have an additional branch point called an exit [21, 22, 23]. These exits are small subnetworks, often only a single layer, whose output can be used alongside or instead of the main trunk network's output. The most common application is for deep classifiers [21, 22, 37, 38], often for computer vision or language models. In a classification task, the network is typically given inputs of varying difficulties, For easy inputs, i.e., those far from the decision boundary, the network may be able to classify accurately using only basic features built by the earlier layers. But for difficult inputs, the network may need to utilize all the layers to build more complex features in order to make a successful classification [21]. By equipping the network with multiple exits, an easy input can reduce inference-time compute by using the output of an early exit only, saving both time and energy when making predictions [24]. The promise of efficiency gains motivates the study of early exit networks and *learning to exit* algorithms. This approach goes by various names in the literature, including *deeply supervised* networks, cascaded learning, conditional deep learning, and adaptive inference. For more on multi-exit and early-exit networks, see Scardapane et al. [24], Laskaridis et al. [39] or Rahmath et *al.* [40].

While multi-exits offer energy-efficient and fast inference, for our purposes, they provide a more important benefit: they allow *deep supervision* [21] of the network during training. By introducing a loss function that combines the outputs of all exits, training gradients enter directly into the earlier hidden layers. With appropriate losses, this can allow for a network that is trained more accurately or more efficiently, or both [24].

We argue in this paper that multi-exits are a natural extension of KANs and, unlike alternative forms of deep supervision such as DenseNet-style forward connections [41] (see Discussion), this form of deep supervision is especially appropriate to the interpretability advantages of KANs.

3 Adding exits to KANs

A multi-exit KAN augments a standard (single-exit) KAN of shape $[d = N_0, N_1, ..., N_L = m]$ with additional exits as follows. For each layer ℓ of the KAN, beginning from the input and continuing until the second-to-last hidden layer¹, add an exit layer, another KAN network, with shape $[N_{\ell}, m]$. The *m*-dimensional output of these exits can then be used to predict the same output as the main trunk exit. We illustrate a multiexit KAN alongside the corresponding single-exit version in Fig. 1.

To train a multi-exit KAN requires a data loss and a regularization loss. The regularization loss can remain the same as in standard KANs, but applied to all activation functions across the network including those in the exits. The data loss, on the other hand, must now accommodate multiple predictions. A multi-exit KAN with *K* exits will now emit *K* outputs, denoted $\hat{y}^{(0)}, \hat{y}^{(1)}, \dots, \hat{y}^{(K-1)}$, where $\hat{y}^{(0)}$ is the output of the exit connected directly to the KAN input layer and $\hat{y}^{(K-1)}$ is the output of the main trunk. To train the entire KAN across all exits, enabling deep supervision [21] of all layers, requires a joint loss function that combines all these outputs. A

¹The last hidden layer already has an exit, the main trunk output.

straightforward option that we focus on is a weighted average of the individual exit data losses:

$$\mathcal{L}_{\text{multi}} = \sum_{k=0}^{K-1} w_k \mathcal{L}_k, \qquad (9)$$

where $\mathcal{L}_k = \sum_i \left(y_i - \hat{y}_i^{(k)}\right)^2 / n$ is the MSE for exit *k* and w_k is the weight for exit *k*. The exit weights satisfy $\sum_k w_k = 1$. These weights become a hyperparameter to be tuned by the researcher using validation data and this tuning process in practice has been straightforward. Equation (9) is quite flexible as the exit weights allow us to prioritize certain exits by weighting them more heavily than others, and individual exits can even be disabled by zeroing their weights. However, for a network with many exits, these K - 1 degrees of freedom become a large search space. Therefore, after our main results, in Sec. 5 we propose and apply a "learning to exit" method to automatically learn *w* alongside the KAN parameters.

Number of parameters How many more parameters are added to a KAN by adding exits? The number of parameters in a KAN depends on its architecture, which dictates the number of activation functions, and the number of parameters per activation function. The parameters per activation function depends on the number and order of the spline bases (Eq. (5)), assuming B-splines are used to model the activation functions. This is the same for activation functions in the main trunk and in the exits, so we only need to consider the number of activation functions to determine the overhead added to a KAN by adding exits.

The number of activation functions between layers in a standard KAN is the product of the layer widths, so a KAN with shape $[d = N_0, N_1, \dots, N_L = m]$ will have

$$N_{\rm act} = dN_1 + N_1 N_2 + \dots + N_{L-1} m \tag{10}$$

activation functions. A multi-exit KAN of the same shape will have those activation functions plus an additional *m* activation functions for each additional exit:

$$N_{\text{act}} = dN_1 + N_1 N_2 + \dots + N_{L-2} N_{L-1} + (d + N_1 + \dots + N_{L-1}) m.$$
(11)

Notice that the sum of layer widths will be smaller than the sum of products of adjacent layer widths, (unless the layers are all one unit), so, unless *m* is large, the main trunk dominates the number of parameters in the KAN. Indeed, for the case of uniform layer width, $N_{\ell} = d$ for all ℓ , the main trunk will have $(L-1)d^2 + dm$ activation functions, including the original exit, and the newly added exits will introduce (L-1)dm activation functions in total. In this case, the new exits will contribute fewer activation functions than the main branch when m < d, typical of regression problems (Eq (1)) and a nearly equal number (fewer by dm) of activation functions when m = d, typical of dynamical systems modeling (Eqs. (2) or (3)).

Also, note that no one prediction made by the model will use all the activation functions, even if they were all used during training. Thus, multi-exits provide their benefits with reasonable, often modest, parameter overhead.

4 Results

Experiments compare single-exit and multi-exit KANs on various regression tasks of known functional forms (Sec. 4.1; Figs. 2 and 3; Table 1), on multi-step forecasting of dynamical systems (Sec. 4.2; Figs. 4 and 5), on a model of continual learning (Sec. 4.3, Fig. 6), and on three real-world datasets ((Sec. 4.4, Table 2). For each task, a manual architecture search identifies good KAN shapes and exit weights, while holding all other hyperparameters fixed. Performance is assessed with the root mean squared error (*RMSE*) of its predictions on test data (as well as R^2 values for the real-world data). Regarding the exit weights, Sec. 5 explores learning the weights automatically with a "learning to exit" approach.

4.1 Regression tasks

Our first experiment uses the sinc function,

$$f(x) = \frac{\sin(\pi x)}{\pi x}.$$
 (12)

This one-dimensional function works well as a test because it features both local oscillations and global decay, and approximation methods often struggle due to its sharp spectral cutoffs, making it a simple but challenging benchmark for function approximation.

For the sinc function a KAN of shape [1, 2, 2, 2, 1] performed well. See Appendix A for full details on training settings and hyperparameters and data generation. This KAN may seem deep for such a function. KANs can support multiplication units [16], although they are not strictly necessary due to the KART, so we decided to forgo them for simplicity and therefore expected a deeper KAN to perform better for this task, hence the aforementioned shape. Parameterizing the multi-exit weighted loss with a simple linear ramp, w = [1, 2, 3, 4] (unnormalized) worked well: performance on test data was good, with the final exit showing an order of magnitude lower error than the equivalent single-exit network (Fig. 2). In fact, three of the four exits, all but Exit 0, outperformed the single-exit network, indicating robust and parsimonious (parameter-efficient) capture of the data.

Next was the function

$$f(x_1, x_2) = \sin\left(2\pi x_1^2\right) \sin\left(4\pi x_2^2\right).$$
 (13)

This nonlinear function tests multivariate approximation through spatially-varying frequencies that challenge learning. Models used a KAN shape of [2, 3, 2, 1] and, for the multi-exit KAN, exit weights w = [1, 2, 1]. As shown in Fig. 3, we found good results for the single-exit KAN but even better for the multi-exit KAN.

In the multi-exit KAN, unsurprisingly, the initial exit, which lacks any composability, is unable to represent this function.



Figure 2: Multi-exits reduce error in 1D nonlinear regression.

The next exit, however, does well, achieving error one fourth that of the larger, single-exit model (RMSE = 0.0045 vs. 0.0232). The final exit at RMSE = 0.007 also outperforms the same-size single-exit model. This result demonstrates that multi-exit KANs can identify the right level of complexity, with the middle exit outperforming both simpler and more complex alternatives.

Our final regression experiment uses a sample of equations from the *Feynman Equation* dataset, a standard function approximation benchmark [42, 43]. These equations cover a range of functional forms and complexity levels, while representing practically relevant physical relationships.

Each equation was converted to the dimensionless form indicated in the table and used to generate data (see Appendix A). KANs with shape [d, 5, 5, 5, 5, 1]) were fitted to each dataset. The multi-exit KAN found good results with w = [0, 0, 1, 1, 3/2] (only Exits 2–4 are active) on Eq. I.27.6, and we proceeded to use this w across all equations. For each multi-exit KAN, Table 1 reports the smallest *RMSE* across its exits.

As shown in Table 1, multi-exit networks achieved lower test *RMSE* than the single-exit model in nine of ten cases. Interestingly, and in line with our observations from the 2D regression shown in Fig. 3, for half of the problems, the best performing exit is not the final exit, indicating we find more parsimonious (smaller) models with multi-exits that are also more accurate than larger, single-exit models.

4.2 Data-driven modeling of dynamical systems

Beyond regression problems, experiments study how multi-exit architectures perform as models of dynamical systems, which present challenging time-series forecasting problems due to their chaotic attractors, evaluating both one-step and multi-step (closed-loop) prediction tasks.

Two dynamical systems are considered. The first is the *Ikeda* map [44, 45], a famous example of a practically motivated

discrete-time chaotic dynamical system that does not admit an accurate sparse representation:

$$x_{n+1} = 1 + \mu (x_n \cos (\phi_n) - y_n \sin (\phi_n)),$$

$$y_{n+1} = \mu (x_n \sin(\phi_n) + y_n \cos(\phi_n)),$$
(14)

where $\phi_n = 0.4 - 6(1 + x_n^2 + y_n^2)^{-1}$ and bifurcation parameter $\mu = 0.9$. KANs, unlike sparse regression methods, have been show to model the Ikeda map well [19].

KANs found good results modeling the Ikeda map with a [2, 4, 4, 2] shape and, for the multi-exit KAN, exit weights w = [0, 0, 1, 2] (the first two exits were disabled). (Note that Panahi *et al.* [19] used a fixed grid G = 10 while we used grid refinement (see Appendix) for both single- and multi-exit KANs.) For one-step prediction the multi-exit KAN achieved *RMSE* = 4.560×10^{-3} compared to the single-exit's 5.484×10^{-3} , an improvement of 16.8%. For multi-step prediction, as shown in Fig. 4 the multi-exit KAN tracks the dynamics for approximately twice as many timesteps as the single-exit KAN before inevitably diverging due to chaos.

The second dynamical system we consider is a continuoustime model of a three-population *ecosystem*:

$$\frac{dN}{dt} = N\left(1 - \frac{N}{K}\right) - x_p y_p \frac{NP}{N + N_0},$$

$$\frac{dP}{dt} = x_p P\left(y_p \frac{N}{N + N_0} - 1\right) - x_q y_q \frac{PQ}{P + P_0}, \quad (15)$$

$$\frac{dQ}{dt} = x_q Q\left(y_q \frac{P}{P + P_0} - 1\right),$$

where *N*, *P*, and *Q* are the primary producer, herbivore, and carnivore populations, respectively, and the carrying capacity *K* acts as bifurcation parameter. To model a chaotic system, we set K = 0.98, $x_p = 0.4$, $y_p = 2.009$, $x_q = 0.08$, $y_q = 2.876$, $N_0 = 0.16129$, and $P_0 = 0.5$, ensuring the system exhibits a chaotic attractor [46]. As with the Ikeda map, KANs are known to model this system well [19].

KANs performed well modeling the ecosystem with shape [3, 3, 3, 3] KANs and w = [2, 1, 1/2] for the multi-exit KAN, achieving very good multi-step prediction (Fig. 5) but slightly worse one-step prediction (*RMSE* = 3.774×10^{-4} for the multi-exit compared to 3.171×10^{-4} for the single-exit KAN). The single-exit KAN, while still performing well, appears to be overfit to the local trajectory while the multi-exit KAN better captured the underlying attractor structure.

4.3 Continual learning

As a further demonstration of the usefulness of augmenting KANs with multi-exits, we consider a toy model of continual learning [15, 47] (Fig. 6, top row). Here the function to represent is a one-dimensional line of five peaks, a multi-modal mixture of Gaussian functions:

$$f(x) = \sum_{i=1}^{5} \exp\left(-300(x-c_i)^2\right),$$
 (16)



Figure 3: Multi-exits reduce error in 2D nonlinear regression.

Feynman Eq.	Original Formula	Dimensionless formula	# vars	RMSE (single)	RMSE (multi) (Exit)
I.6.20	$\frac{e^{-\frac{\theta^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$	$\frac{e^{-\frac{\theta^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$	2	1.90×10^{-4}	1.03 × 10⁻⁵ (3)
I.6.20b	$\frac{e^{-\frac{(\theta-\theta_1)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$	$\frac{e^{-\frac{(\theta-\theta_1)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}}$	3	1.12×10^{-3}	1.16×10^{-3} (2)
I.9.18	$\frac{Gm_1m_2}{(x_2-x_1)^2+(x_2-x_1)^2+(z_2-z_1)^2}$	$\frac{a}{(b-1)^2+(c-d)^2+(e-f)^2}$	6	5.03×10^{-2}	2.09×10^{-2} (3)
I.12.11	$q\left(E_f + Bv\sin\theta\right)$	$1 + a \sin \theta$	2	4.94	2.41 × 10⁻² (4)
I.13.12	$Gm_1m_2(1/r_2-1/r_1)$	a(1/b-1)	2	1.55×10^{-1}	1.31 × 10 ⁻² (4)
I.15.3x	$\frac{x-ut}{\sqrt{1-\frac{u^2}{2}}}$	$\frac{1-a}{\sqrt{1-b^2}}$	2	1.32×10^{-2}	2.86×10^{-3} (2)
I.16.6	$\frac{u+v}{1+\frac{uv}{2}}$	$\frac{a+b}{1+ab}$	2	2.05×10^{-3}	3.04 × 10⁻⁴ (2)
I.18.4	$\frac{m_1r_1+m_2r_2}{m_1+m_2}$	$\frac{1+ab}{1+a}$	2	2.53×10^{-4}	1.90 × 10⁻⁴ (4)
I.26.2	$asin(n sin \theta_2)$	$asin(n sin \theta_2)$	2	1.22×10^{-3}	7.62 × 10 ⁻⁴ (4)
I.27.6	$\frac{1}{n/d_2+1/d_1}$	$\frac{1}{1+ab}$	2	1.77×10^{-5}	1.54×10^{-5} (4)

Table 1: Performance on Feynman Equation dataset.

where the centers c_i of peaks *i* were evenly spaced. For this experiment, we generated 100 equally spaced points around each peak, for a total of 500 samples. KANs can easily fit such data but there is a wrinkle: the model is not trained on all the data at once. Instead, it sees the data in phases, one peak at a time (Fig. 6, top row).

The question becomes whether a KAN can learn a new peak without losing its representation of previous peaks. As argued by Liu *et al.* [15], shallow KANs are well adapted to this task due to the local nature of their spline-based activation functions: updating one region of a spline will not affect the fit of other regions, enabling the KAN to retain the form of a previous peak when incorporating the next peak. However, they also note that KANs lose this ability as they get deeper, since the composition of splines across multiple layers reduces their locality, opening the door for catastrophic forgetting. Do multi-exit KANs with their deep supervision retain more locality and exhibit less forgetting?

As seen in Fig. 6, we can answer in the affirmative. While both architectures display some forgetting, with previously learned peaks changing with subsequent data, the effect is worse for the single-exit KAN (shape [1, 5, 5, 1]), especially when learning the second peak. Compared to the single-exit KAN of the same architecture (Fig 6, middle row), the multiexit KAN (bottom row) with exit weights w = [1, 1, 2e3] better tracks the underlying function across training phases, and when finished displays less than half the error of the single-exit KAN (*RMSE* = 0.086 vs. 0.19).

4.4 Real-world data

Now we consider how multi-exit KANs perform on three real-world datasets:

• *Airfoil Noise*. Predict the self-noise (scaled sound pressure, in dB) of a NACA 0012 airfoil for different angles



Figure 4: Multi-step prediction of the Ikeda map (Eq. (14). The multi-exit KAN tracks the dynamics well for about twice as many steps into the future (shaded regions) as the single-exit KAN. (Blue: ground truth; orange: KAN prediction.)

of attack, free stream velocity, and other features. Data originated from anechoic wind tunnel experiments [48].

- Power Plant Energy. Predict the electrical power output (in MW) for different ambient atmospheric conditions, temperature, pressure, relative humidity, and the steam turbine pressure (or vacuum). Data originated from a 480 MW combined cycle power plant with two gas turbines, one steam turbine, and two heat recovery steam generators, collected over a six-year period (2006-2011) [49, 50].
- Superconductor Critical Temperature. Predict critical temperature (in K) of superconductors based on material properties. The dataset contains many features and statistical variants, so for ease of experimentation, we selected five representative features capturing composition complexity, electronic structure, and chemical bonding properties: number of elements, weighted mean valence, valence entropy, weighted mean first ionization energy, and mean electron affinity. These data originated from the Superconducting Material Database maintained by Japan's National Institute for Materials Science [51, 52].

All data were retrieved from the UCI Machine Learning Repository [53] (accessed: 23 May 2025). From each dataset 1k observations for training and 1k for testing were randomly sampled, except for the smaller Airfoil dataset, which contains only 1503 observations in total so 750 observations for training and 750 for testing were randomly sampled. Besides sampling for training/testing and selecting features for the superconductivity data, no other filtering or preprocessing was performed.



Figure 5: Multi-step prediction of the ecosystem (Eq (15)). The multiexit KAN encodes the dynamics well, not diverging significantly until t > 1000.

Single-exit KANs generally performed well with one hidden layer (Table 2), but we also considered zero- and two-layer KANs, and used the same shapes for the corresponding multiexit KANs (a KAN with shape [d, m] can only have one exit). Experimentation led to exit weights w that performed well, although for both shape and weight, there is likely room for improvement. All other hyperparameters and training settings (grid refinement, etc.) were unchanged, leaving even more room to improve. For multi-exit KANs, in all cases the best performing exit was either Exit 0 or 1.

As seen in Table 2, multi-exits improved predictive performance on all three datasets. Interestingly, for all datasets, the two- and three-exit KANs outperformed the single-exit KANs of the same size and smaller, at earlier exits. For instance, on Airfoil, the three-exit KAN achieved RMSE = 5.129 at *Exit 0*, improving on the single-exit KANs of the same size (RMSE = 5.338) by 3.9%. Multi-exit KANs achieved accurate fits while simultaneously being more parsimonious. While additional training of the one-layer KAN could potentially close this gap, this result nevertheless suggests that multi-exit KANs improve performance simultaneously across shallower and deeper architectures.

5 Learning to exit

The exit weight hyperparameter complicates the architecture search problem, already a potential challenge when estimating



Figure 6: Mitigating catastrophic forgetting with multi-exits in a toy model (Eq. (16)) of continual learning.

Dataset	Obs.	Features	Shape	Exit weight	Single		Multi		
					RMSE	R^2	RMSE	R^2	Exit
Airfoil	1500	5	[5,1]	_	5.338	0.421			_
			[5, 5, 1]	[1, 2e3]	5.897	0.293	5.635	0.355	1
			[5, 5, 5, 1]	[1e3, 100, 1]	7.130	-0.0332	5.129	0.465	0
Power plant	2000	4	[4,1]	_	4.462	0.934	_		_
			[4, 4, 1]	[5,2]	4.496	0.932	4.451	0.934	0
			[4, 4, 4, 1]	[100, 1e3, 1]	4.601	0.929	4.406	0.935	1
Superconductivity	2000	5 (of 81)	[5,1]	_	19.524	0.658	_		_
			[5, 5, 1]	[10, 8]	18.556	0.691	18.353	0.698	1
			[5, 5, 5, 1]	[1, 10, 1]	19.278	0.666	18.445	0.695	1

Table 2: Performance of single- and multi-exit KANs on three datasets.

KAN models. Indeed, from our experiments, it is not always obvious *a priori* what exit weights will best optimize the KAN. Sometimes uniform weights work well, or an increasing or decreasing ramp, or sometimes even a heavy weight on the last exit(s) can be beneficial. While KANs train quite quickly on these smaller scientific problems, allowing rapid iteration to explore the weight simplex, nevertheless, it would be useful, and data-efficient, to avoid this trial-and-error process. To this end, here we introduce and apply a basic "learning to exit" algorithm which treats the exit weights as learnable parameters, eliminating them entirely from the architecture search.

5.1 Learnable exit weights

Consider a multi-exit KAN with *L* exits producing outputs $\{\hat{y}^{(0)}, \hat{y}^{(1)}, \dots, \hat{y}^{(L-1)}\}$ for a given input. Until now multi-exit KANs were trained using a fixed weighted loss (Eq. (9)):

$$\mathcal{L}_{\text{multi}} = \sum_{i} w_i \mathcal{L}_i \left(\hat{y}^{(i)}, y \right), \qquad (17)$$

where w_i (i = 0, ..., L - 1) are predetermined constants and \mathcal{L}_i is the loss function for exit *i*. In our learning to exit framework, we introduce *exit logits* $\boldsymbol{\theta}_w = \{\theta_0, \theta_1, ..., \theta_{L-1}\}$ to be optimized. A softmax transformation connects these to Eq. 17,

$$w_i(\theta_i) = \frac{\exp(\theta_i)}{\sum_j \exp(\theta_j)},$$
(18)

guaranteeing positive, normalized exit weights. The loss function becomes:

$$\mathcal{L}_{\text{joint}}(\boldsymbol{\theta}_{\text{KAN}}, \boldsymbol{\theta}_{w}) = \sum_{i} w_{i}\left(\theta_{i}\right) \mathcal{L}_{i}\left(\hat{y}^{(i)}\left(\boldsymbol{\theta}_{\text{KAN}}^{(i)}\right), y\right), \quad (19)$$

where θ_{KAN} represents the KAN parameters and $\theta_{\text{KAN}}^{(i)}$ denote the KAN parameters for layers up to and including exit *i*.

5.2 Optimization procedure

The optimization problem is formulated as:

$$\boldsymbol{\theta}_{\text{KAN}}^*, \boldsymbol{\theta}_w^* = \arg\min_{\boldsymbol{\theta}_{\text{KAN}}, \boldsymbol{\theta}_w} \mathcal{L}_{\text{joint}}(\boldsymbol{\theta}_{\text{KAN}}, \boldsymbol{\theta}_w).$$
 (20)

Both sets of parameters are updated simultaneously using gradient-based optimization. The gradients with respect to the exit logits are

$$\frac{\partial \mathcal{L}_{\text{joint}}}{\partial \theta_i} = \sum_j \mathcal{L}_j \left(\hat{y}^{(j)}, y \right) \frac{\partial w_j}{\partial \theta_i}, \qquad (21)$$

where, from Eq. (18),

$$\frac{\partial w_j}{\partial \theta_i} = \begin{cases} w_i(1-w_i) & \text{if } i=j, \\ -w_iw_j & \text{if } i\neq j. \end{cases}$$
(22)

As before, L-BFGS was used for the joint optimization, but other methods, such as Adam, could be used instead. Weight logits were initialized with uniform values ($\theta_w = 0$), but other initializations may be worth exploring. All other training settings and hyperparameters were unchanged.

5.3 Results

To evaluate the learning-to-exit model, we apply it to the three datasets studied in Sec. 4.4, focusing on the three-exit architectures. Comparing to Table 2, the results were promising. On every dataset, the learned model outperformed the single-exit KAN of the same shape. Further, in two of the three cases, the new model outperformed *every* model in Table 2. Specifically, for Airfoil the new model achieved *RMSE* = 4.947 and for Superconductivity *RMSE* = 18.126, beating the previous best *RMSE* = 5.129 and 18.353, respectively. On the other hand, for Power plant, it achieved *RMSE* = 4.558, better than the single-exit result of *RMSE* = 4.601 but worse than the multi-exit's *RMSE* = 4.406.

Interestingly, the learned exit weights varied for all three datasets, despite all being initialized to w = [1, 1, 1]/3. For Airfoil, the final weights were, to machine precision, w = [1, 0, 0] (focus on first exit), for Superconductivity, w = [0, 0, 1] (focus on last exit), and for Power plant, w = [0.002, 0.848, 0.150] (mixed focus). Power's exit logits also converged more slowly than the other two, which may relate to the weaker relative performance.

The dynamics of exit selection are worth further study for one, we expected a regularization term on the exit logits would be necessary, but these results imply otherwise—yet our findings already show that the learning-to-exit model has promise.

6 Discussion

Augmenting KANs with multiple exits improves their performance and often their parsimony. When a multi-exit KAN performs well at an early exit, it suggests that the full network is deeper than necessary and functions with fewer levels of composition are sufficient to model the given data. While in principle a deeper single-exit KAN can be encouraged to simplify, either through regularization or through linearizing the later activation functions—in fact, both single-exit and multi-exit methods are likely to benefit in general from finetuning hyperparameters and training settings—nevertheless the multi-exit approach is a promising alternative to achieving this parsimony.

What is the mechanism of action behind the success of multi-exits? The first and more obvious mechanism is that of deep supervision. The compositional nature of learned activation functions makes gradient flow through many layers a challenge during training. By connecting the loss function directly to the earlier layers, training will allow for better conditioning of the activation functions and weights within those layers. In this sense, the multi-exits fulfill a role similar to that of DenseNet [41]-style forward connections. In DenseNet architectures, the forward connections link input and hidden layers directly to the final output layer, and backpropagation can then reach deeper into the network for training. (The other common form of deep supervision, residual connections (ResNet) [25] is already commonly used in KANs; see Eq. (5).) In fact, forward connections in KANs could be viewed as another useful form of deep supervision. However, they create a large number of outputs at the final layer, which may necessitate a more complex final functional form, hindering KAN's goal of interpretability. Multi-exits, in contrast, may be a better alternative thanks to their potential for parsimony.

A second and less obvious mechanism of action is implicit regularization through the optimization method. In quasi-Newton methods such as BFGS and L-BFGS, a Hessian approximation captures curvature information across all parameters simultaneously, enabling the optimizer to find parameter configurations that balance competing objectives from different exits. In other words, the curvature approximation regularizes the parameters during training. In our experiments using L-BFGS, the presence of Exit 0 only, the exit connected directly to the input, still led to improvements in KAN performance. When there are no other intermediary exits, deep supervision can't condition the earlier layers—Exit 0 is not on the computational path of the last exit. Yet, through the optimization process, the network is still able to find better solutions. For instance, in the experiments with real world data (Sec. 4.4), multi-exit KANs with only one hidden layer still improved, albeit quite modestly, over single-exit KANs. Implicit regularization was absent when training with first-order methods such as SGD or Adam [36], which update parameters based on individual gradient moments rather than capturing the joint parameter dependencies, though such methods still provide the benefits of deep supervision. This second mechanism, while weaker than deep supervision, offers another reason why L-BFGS is well-suited for KAN optimization.

Despite the promising results presented in this work, some limitations should be acknowledged. Perhaps the most serious concern is the extra need to set the exit weights when fitting a multi-exit KAN. In all our experiments, setting *w* required only brief coarse-grained tuning (and it is addressed by our learning-to-exit algorithm) but in settings where data are scarce, it may be difficult to optimize the exit weights without overfitting. Second, we focused our comparisons on single-exit versus multi-exit KANs to isolate the effect of the architectural change, leaving broader comparisons for future work. Likewise, future work should consider the effects of different hyperparameter values and training settings, including other sources of regularization such as pruning [15], since it is important in practice to ascertain the optimal settings for different applications and whether multi-exits benefit from or remain useful with such fine-tuning. Third, while multi-exit architectures often identify more parsimonious models, the interpretability gains are indirect-the exits themselves do not enhance the interpretability of individual activation functions, but rather help identify simpler network configurations. Fourth, our learning-to-exit algorithm should be studied further and, while effective, can surely be improved. Finally, the additional computational overhead during training, though modest, may become more significant for very deep networks with many exits.

Several promising directions are worth pursuing. First, receiving multiple predictions from a KAN immediately brings to mind the idea of ensemble learning [54]. However, ensembles benefit from uncorrelated or de-correlated models, but the different exits in a multi-exit KANs are not independent. Investigating methods to encourage heterogeneity among exits while maintaining their collaborative training could enable ensemble learning. This may also lead to uncertainty quantification capabilities [55, 56]—a key goal for KANs [57, 58]—through the natural variation of predictions across exits. Second, our learning-to-exit framework leaves room for improvement, and incorporating ideas from differentiable architecture search [59] more generally could benefit KANs. Third, exploring exit architectures beyond simple single-layer KAN exits may yield better performance, although doing so without harming parsimony may be difficult. Fourth, extending multi-exit architectures from B-splines to other KAN variants (Fourier-KANs, Wavelet-KANs, Physics-Informed KANs, etc.) could reveal whether the benefits generalize across different basis functions. Finally, applying multi-exit KANs to larger-scale scientific problems, particularly in domains like climate modeling or molecular dynamics where both accuracy and interpretability are crucial, would provide valuable insights into their practical utility.

Conclusion We have introduced multi-exit architectures for Kolmogorov-Arnold Networks, demonstrating that augmenting KANs with additional outputs consistently improves their performance across diverse scientific modeling tasks. Our experiments revealed that multi-exit KANs often achieve their best performance at earlier exits, indicating that they successfully identify more parsimonious models without sacrificingand often improving-accuracy. While multi-exit architectures introduce additional hyperparameters, our results show the benefits substantially outweigh this added complexity. This finding is particularly valuable for scientific applications where interpretability is paramount, as simpler models with fewer compositional layers are inherently more interpretable. These results suggest that multi-exit architectures represent a natural and effective enhancement to KANs, offering researchers a principled approach to finding the right balance between model complexity and performance.

A Materials and Methods

Implementation and training KANs were implemented with the PyKAN library v0.2.8 (https://github.com/ KindXiaoming/pykan), based on PyTorch v2.6.0 [60]. To fit KAN models using training data, the Limited-memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS) algorithm [35] was used with history size 10, strong Wolfe conditions for line search, and convergence tolerances of 10^{-32} for gradient norm, parameter changes, and curvature conditions. Unless otherwise noted, fitting employed a progressive grid refinement strategy, iteratively refining the spline basis functions through a sequence of increasing grid sizes G = 3, 5, 10, 20, with 30 optimization steps performed at each grid resolution. This approach allows the model to first capture coarse-grained patterns before learning finer details, leading to more stable convergence and better generalization performance [15]. All other training settings and hyperparameters were kept at the default values of PyKAN, including the learning rate of 1.0 and the default spline order K = 3 and grid update schedule. Multiplication units were not used. The same training procedure and settings was used throughout, for both single-exit and multi-exit KANs. Fine-tuning these settings would likely further improve performance, but both types of KANs would be expected to benefit. Our source code will be available at https://github.com/bagrow/multi-exit-KAN upon acceptance of this manuscript.

Experimental details For the 1D and 2D regression tasks, n = 1000 training and n = 1000 testing points were generated, with $\mathbf{x} \sim U([x_{\min}, x_{\max}]^d)$ and y generated from \mathbf{x} , without additional noise, according to the given equation. For the fits illustrated in Fig. 1, n = 1000 training points and n = 200testing points were used, as well as a single G = 5 grid, 30 optimization steps, KAN shape [2, 3, 2, 1] and w = [1, 1, 1]. The Feynman Equations dataset used the same data generating process as the 1D and 2D regression tasks, but each exogenous variable's range was given by the range of values in the released AI Feynman dataset [42]. For the dynamical systems experiments, data for each system was generated and split into training and testing folds following the procedure and parameters of Panahi et al. [19]. The continual learning experiment used 100 samples per peak, a grid size of 20 without refinement or updates, 10 L-BFGS steps per phase, and all other settings at PyKAN defaults. Details for the three real-world datasets were covered in Sec. 4.4.

References

- G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, "Machine learning and the physical sciences," *Reviews of Modern Physics*, vol. 91, no. 4, p. 045002, 2019. 1
- [2] R. Roscher, B. Bohn, M. F. Duarte, and J. Garcke, "Explainable machine learning for scientific insights and discoveries," *Ieee Access*, vol. 8, pp. 42200–42216, 2020. 1

- [3] Y. Xu, X. Liu, X. Cao, C. Huang, E. Liu, S. Qian, X. Liu, Y. Wu, F. Dong, C.-W. Qiu, *et al.*, "Artificial intelligence: A powerful paradigm for scientific research," *The Innovation*, vol. 2, no. 4, 2021. 1
- [4] H. Wang, T. Fu, Y. Du, W. Gao, K. Huang, Z. Liu, P. Chandak, S. Liu, P. Van Katwyk, A. Deac, *et al.*, "Scientific discovery in the age of artificial intelligence," *Nature*, vol. 620, no. 7972, pp. 47–60, 2023. 1
- [5] G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and L. Yang, "Physics-informed machine learning," *Nature Reviews Physics*, vol. 3, no. 6, pp. 422–440, 2021. 1
- [6] K. Kashinath, M. Mustafa, A. Albert, J. Wu, C. Jiang, S. Esmaeilzadeh, K. Azizzadenesheli, R. Wang, A. Chattopadhyay, A. Singh, *et al.*, "Physics-informed machine learning: case studies for weather and climate modelling," *Philosophical Transactions of the Royal Society A*, vol. 379, no. 2194, p. 20200093, 2021. 1
- [7] R. Lam, A. Sanchez-Gonzalez, M. Willson, P. Wirnsberger, M. Fortunato, F. Alet, S. Ravuri, T. Ewalds, Z. Eaton-Rosen, W. Hu, A. Merose, S. Hoyer, G. Holland, O. Vinyals, J. Stott, A. Pritzel, S. Mohamed, and P. Battaglia, "Learning skillful medium-range global weather forecasting," *Science*, vol. 382, no. 6677, pp. 1416–1421, 2023. 1
- [8] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, *et al.*, "Highly accurate protein structure prediction with alphafold," *nature*, vol. 596, no. 7873, pp. 583–589, 2021. 1
- [9] F. J. Montáns, F. Chinesta, R. Gómez-Bombarelli, and J. N. Kutz, "Data-driven modeling and learning in science and engineering," *Comptes Rendus Mécanique*, vol. 347, no. 11, pp. 845–855, 2019. Data-Based Engineering Science and Technology. 1
- [10] W. Bradley, J. Kim, Z. Kilwein, L. Blakely, M. Eydenberg, J. Jalvin, C. Laird, and F. Boukouvala, "Perspectives on the integration between first-principles and data-driven modeling," *Computers & Chemical Engineering*, vol. 166, p. 107898, 2022.
- [11] A. Bell, I. Solano-Kamaiko, O. Nov, and J. Stoyanovich, "It's just not that simple: an empirical study of the accuracy-explainability trade-off in machine learning for public policy," in *Proceedings* of the 2022 ACM conference on fairness, accountability, and transparency, pp. 248–266, 2022. 1
- [12] A. Ferrario and M. Loi, "How explainability contributes to trust in ai," in *Proceedings of the 2022 ACM conference on fairness*, accountability, and transparency, pp. 1457–1466, 2022. 1
- [13] R. Van Noorden and J. M. Perkel, "Ai and science: what 1,600 researchers think," *Nature*, vol. 621, no. 7980, pp. 672–675, 2023. 1
- [14] D. Castelvecchi, "Can we open the black box of ai?," *Nature News*, vol. 538, no. 7623, p. 20, 2016. 1
- [15] Z. Liu, Y. Wang, S. Vaidya, F. Ruehle, J. Halverson, M. Soljacic, T. Y. Hou, and M. Tegmark, "KAN: Kolmogorov–Arnold networks," in *The Thirteenth International Conference on Learning Representations*, 2025. 1, 2, 4, 6, 7, 11
- [16] Z. Liu, P. Ma, Y. Wang, W. Matusik, and M. Tegmark, "Kan 2.0: Kolmogorov-arnold networks meet science," *arXiv preprint arXiv:2408.10205*, 2024. 1, 3, 5

- [17] J. D. Toscano, V. Oommen, A. J. Varghese, Z. Zou, N. Ahmadi Daryakenari, C. Wu, and G. E. Karniadakis, "From PINNs to PIKANs: Recent advances in physics-informed machine learning," *Machine Learning for Computational Science and Engineering*, vol. 1, no. 1, pp. 1–43, 2025. 1
- [18] B. C. Koenig, S. Kim, and S. Deng, "KAN-ODEs: Kolmogorov– Arnold network ordinary differential equations for learning dynamical systems and hidden physics," *Computer Methods in Applied Mechanics and Engineering*, vol. 432, p. 117397, 2024. 1, 2
- [19] S. Panahi, M. Moradi, E. M. Bollt, and Y.-C. Lai, "Data-driven model discovery with Kolmogorov–Arnold networks," *Phys. Rev. Res.*, vol. 7, p. 023037, Apr 2025. 1, 2, 6, 11
- [20] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *Journal of Machine Learning Research*, vol. 20, no. 55, pp. 1–21, 2019. 2
- [21] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu, "Deeplysupervised nets," in *Artificial intelligence and statistics*, pp. 562– 570, Pmlr, 2015. 2, 4
- [22] S. Teerapittayanon, B. McDanel, and H. Kung, "BranchyNet: Fast inference via early exiting from deep neural networks," in 2016 23rd International Conference on Pattern Recognition (ICPR), pp. 2464–2469, 2016. 2, 4
- [23] P. Panda, A. Sengupta, and K. Roy, "Conditional deep learning for energy-efficient and enhanced pattern recognition," in *Proceedings of the 2016 Conference on Design, Automation & Test in Europe*, DATE '16, (San Jose, CA, USA), p. 475–480, EDA Consortium, 2016. 2, 4
- [24] S. Scardapane, M. Scarpiniti, E. Baccarelli, and A. Uncini, "Why should we add early exits to neural networks?," *Cognitive Computation*, vol. 12, no. 5, pp. 954–966, 2020. 2, 4
- [25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference* on computer vision and pattern recognition, pp. 770–778, 2016. 2, 3, 10
- [26] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural Networks*, vol. 2, no. 5, pp. 359–366, 1989. 2
- [27] A. N. Kolmogorov, On the representation of continuous functions of several variables by superpositions of continuous functions of a smaller number of variables. American Mathematical Society, 1961. 2
- [28] V. I. Arnold, "On functions of three variables," *Collected Works: Representations of Functions, Celestial Mechanics and KAM Theory*, 1957–1965, pp. 5–8, 2009. 2
- [29] 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, pp. 953–956, 1957. 2
- [30] Z. Li, "Kolmogorov–Arnold networks are Radial Basis Function networks," arXiv preprint arXiv:2405.06721, 2024. 3
- [31] J. Xu, Z. Chen, J. Li, S. Yang, W. Wang, X. Hu, and E. C.-H. Ngai, "FourierKAN-GCF: Fourier Kolmogorov–Arnold network–an effective and efficient feature transformation for graph collaborative filtering," *arXiv preprint arXiv:2406.01034*, 2024. 3

- [32] E. Reinhardt, D. Ramakrishnan, and S. Gleyzer, "SineKAN: Kolmogorov–Arnold networks using sinusoidal activation functions," *Frontiers in Artificial Intelligence*, vol. Volume 7 - 2024, 2025. 3
- [33] S. Sidharth, A. Keerthana, R. Gokul, and K. Anas, "Chebyshev polynomial-based Kolmogorov–Arnold networks: An efficient architecture for nonlinear function approximation," *arXiv preprint arXiv:2405.07200*, 2024. 3
- [34] Z. Bozorgasl and H. Chen, "Wav-kan: Wavelet Kolmogorov– Arnold networks," arXiv preprint arXiv:2405.12832, 2024. 3
- [35] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Mathematical programming*, vol. 45, no. 1, pp. 503–528, 1989. 4, 11
- [36] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2014. 4, 10
- [37] W. Zhou, C. Xu, T. Ge, J. McAuley, K. Xu, and F. Wei, "Bert loses patience: Fast and robust inference with early exit," in *Advances in Neural Information Processing Systems* (H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, eds.), vol. 33, pp. 18330–18341, Curran Associates, Inc., 2020. 4
- [38] J. Xin, R. Tang, Y. Yu, and J. Lin, "BERxiT: Early exiting for BERT with better fine-tuning and extension to regression," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume* (P. Merlo, J. Tiedemann, and R. Tsarfaty, eds.), (Online), pp. 91– 104, Association for Computational Linguistics, Apr. 2021. 4
- [39] S. Laskaridis, A. Kouris, and N. D. Lane, "Adaptive inference through early-exit networks: Design, challenges and directions," in *Proceedings of the 5th International Workshop on Embedded and Mobile Deep Learning*, EMDL'21, (New York, NY, USA), p. 1–6, Association for Computing Machinery, 2021. 4
- [40] H. Rahmath P, V. Srivastava, K. Chaurasia, R. G. Pacheco, and R. S. Couto, "Early-exit deep neural network - a comprehensive survey," ACM Comput. Surv., vol. 57, Nov. 2024. 4
- [41] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017. 4, 10
- [42] S.-M. Udrescu and M. Tegmark, "AI Feynman: A physicsinspired method for symbolic regression," *Science advances*, vol. 6, no. 16, p. eaay2631, 2020. 6, 11
- [43] S.-M. Udrescu, A. Tan, J. Feng, O. Neto, T. Wu, and M. Tegmark, "AI Feynman 2.0: Pareto-optimal symbolic regression exploiting graph modularity," *Advances in Neural Information Processing Systems*, vol. 33, pp. 4860–4871, 2020. 6
- [44] K. Ikeda, "Multiple-valued stationary state and its instability of the transmitted light by a ring cavity system," *Optics communications*, vol. 30, no. 2, pp. 257–261, 1979. 6
- [45] S. Hammel, C. Jones, and J. V. Moloney, "Global dynamical behavior of the optical field in a ring cavity," *Journal of the Optical Society of America B*, vol. 2, no. 4, pp. 552–564, 1985.
- [46] K. McCann and P. Yodzis, "Nonlinear dynamics and population disappearances," *The American Naturalist*, vol. 144, no. 5, pp. 873–879, 1994. 6

- [47] H. van Deventer and A. S. Bosman, "Distal interference: Exploring the limits of model-based continual learning," arXiv preprint arXiv:2402.08255, 2024. 6
- [48] T. F. Brooks, D. S. Pope, and M. A. Marcolini, "Airfoil self-noise and prediction," Tech. Rep. NASA-RP-1218, NASA Langley Research Center, July 1989. NASA Reference Publication 1218.
- [49] H. Kaya, P. Tüfekci, and F. S. Gürgen, "Local and global learning methods for predicting power of a combined gas & steam turbine," in *Proceedings of the international conference* on emerging trends in computer and electronics engineering ICETCEE, pp. 13–18, 2012. 8
- [50] P. Tüfekci, "Prediction of full load electrical power output of a base load operated combined cycle power plant using machine learning methods," *International Journal of Electrical Power & Energy Systems*, vol. 60, pp. 126–140, 2014. 8
- [51] K. Hamidieh, "A data-driven statistical model for predicting the critical temperature of a superconductor," *Computational Materials Science*, vol. 154, pp. 346–354, 2018. 8
- [52] C. for Basic Research on Materials, "MDR SuperCon datasheet ver.240322." 8
- [53] M. Kelly, R. Longjohn, and K. Nottingham, "The UCI machine learning repository," n.d. Accessed: 23 May 2025. 8
- [54] O. Sagi and L. Rokach, "Ensemble learning: A survey," Wiley interdisciplinary reviews: data mining and knowledge discovery, vol. 8, no. 4, p. e1249, 2018. 11
- [55] R. C. Smith, Uncertainty Quantification: Theory, Implementation, and Applications. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2013. 11
- [56] M. Abdar, F. Pourpanah, S. Hussain, D. Rezazadegan, L. Liu, M. Ghavamzadeh, P. Fieguth, X. Cao, A. Khosravi, U. R. Acharya, *et al.*, "A review of uncertainty quantification in deep learning: Techniques, applications and challenges," *Information fusion*, vol. 76, pp. 243–297, 2021. 11
- [57] M. M. Hassan, "Bayesian kolmogorov arnold networks (bayesian_kans): A probabilistic approach to enhance accuracy and interpretability," *arXiv preprint arXiv:2408.02706*, 2024. 11
- [58] A. Mollaali, C. B. Moya, A. A. Howard, A. Heinlein, P. Stinis, and G. Lin, "Conformalized-kans: Uncertainty quantification with coverage guarantees for kolmogorov-arnold networks (kans) in scientific machine learning," *arXiv preprint arXiv:2504.15240*, 2025. 11
- [59] H. Liu, K. Simonyan, and Y. Yang, "DARTS: Differentiable architecture search," in *International Conference on Learning Representations*, 2019. 11
- [60] A. Paszke, "PyTorch: An imperative style, high-performance deep learning library," *arXiv preprint arXiv:1912.01703*, 2019.
 11