# KAN 2.0: Kolmogorov-Arnold Networks Meet Science

Deep Learning for the Natural Sciences

Jonathan Wheeler

### **Table of contents**



1

# **Motivation**

**Current Landscape** 

### The AI + Science Landscape

# Ů

#### **Current State**

- Al is transformational
- Deep learning dominates current approaches

#### **Major Achievements**

- AlphaFold: Near-perfect protein structure prediction
- Al Weather Forecasting: Beyond traditional accuracy



#### **Remaining Challenges**

- Limited scientific understanding from AI systems
- Need for more transparent Al approaches

#### **Application-driven Science (MLP)**

- What (is the output)?
- Works for well

formulated problems with clear objectives

Less interpritible MLP
 models



#### Curiosity-driven Science (KAN)

- Why (is this the result)?
- Scientific discoveries require mathematical understanding and clear principles
- KAN's offer more
  interpretable structure
  through edge functions



2

# Understanding KANs

Theoretical introductions

#### The Kolmogorov-Arnold Theorem

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

- States that any continuous function in a high-dimensional space can be represented as a finite composition of univariate continuous functions and additions.
- Example:

xy = exp(logx + logy)

#### **Two-layer Network**

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



Key problem is has to be smooth functions

### **Key Breakthrough**

$$\mathbf{x}_{l+1} = \underbrace{\begin{pmatrix} \phi_{l,1,1}(\cdot) & \phi_{l,2,1}(\cdot) & \cdots & \phi_{l,n_{l},1}(\cdot) \\ \phi_{l,1,2}(\cdot) & \phi_{l,2,2}(\cdot) & \cdots & \phi_{l,n_{l},2}(\cdot) \\ \vdots & \vdots & & \vdots \\ \phi_{l,1,n_{l+1}}(\cdot) & \phi_{l,2,n_{l+1}}(\cdot) & \cdots & \phi_{l,n_{l},n_{l+1}}(\cdot) \end{pmatrix}}_{\mathbf{\Phi}_{l}} \mathbf{x}_{l}$$

Generalizing the original Kolmogorov-Arnold representation, extending it to arbitrary depths

$$f(\mathbf{x}) = \sum_{i_{L-1}=1}^{n_{L-1}} \phi_{L-1,i_{L},i_{L-1}} \left( \sum_{i_{L-2}=1}^{n_{L-2}} \cdots \left( \sum_{i_{2}=1}^{n_{2}} \phi_{2,i_{3},i_{2}} \left( \sum_{i_{1}=1}^{n_{1}} \phi_{1,i_{2},i_{1}} \left( \sum_{i_{0}=1}^{n_{0}} \phi_{0,i_{1},i_{0}}(x_{i_{0}}) \right) \right) \right) \cdots \right)$$

# **Basic KAN Architecture**

- nodes
  (summation)
- edges (learnable activation functions)
  B-Splines
- Symbolic regression gives output formula



 $\phi(x)$ 

k = 3

## Comparison

MLPs:

- Universal representation theorem based
- fixed nonlinear activation
- Each neuron processes weighted sum of inputs

#### KANs:

- Kolmogorov-Arnold representation based
- Edge-based univariate transformations
- Node-wise signal summation





 $\mathrm{KAN}(\mathbf{x}) = (\mathbf{\Phi}_{L-1} \circ \mathbf{\Phi}_{L-2} \circ \cdots \circ \mathbf{\Phi}_1 \circ \mathbf{\Phi}_0)\mathbf{x}.$ 

# **MultKAN**

- Nodes(input) and subnodes(output)
- Standard KAN layer plus multiplication layer
- No additional trainable parameters
- Learns multiplication



 $\operatorname{MultKAN}(\mathbf{x}) = (\Psi_L \circ \Psi_{L-1} \circ \cdots \circ \Psi_1 \circ \Psi_0) \mathbf{x}.$ 

**Problem**: symbolic formulas are not always possible

#### Solution:

- 1. Build in scientific knowledge to KANs (section 3)
- 2. Extract out scientific knowledge from KANs (section 4)

3

# Building Scientific Knowledge

Knowledge Embedding

### **Overview of embedding knowledge**

incorporate available inductive biases into KANs



**KAN** 

### **1.** Adding important features to KANs

$$y = f(x_1, x_2, \cdots, x_n)$$

- Addition of auxiliary input variable
- Adds expressive power and interpretability to the network

$$a = a(x_1, x_2, \dots, x_n)$$
  
 $y = f(x_1, \cdots, x_n, a)$ 





### 2. Modular structures

- Allows for extracting knowledge from clusters of neurons
- KAN allows for explicit creation of modular structures.
- Two types of modularity:
  - 1. Separability
  - 2. Generalized Symmetry



start\_layer\_id, '[nodes\_id]->[subnodes\_id]->[nodes\_id]...'



#### **Generalized symmetry**

$$f(x_{1}, x_{2}, x_{3}, \cdots) = g(h(x_{1}, x_{2}), x_{3}, \cdots)$$

0, '[0,1]->[0,1]->[0,1]->[0]'



# **3. Compiling Symbolic Formulas**

Combine expressivity of neural networks with symbolic explainable formulas

#### Two steps:

- 1. Compile formula into a KAN (introducing domain knowledge)
- 2. Train KAN on data (learn "new physics")

### **KANPILER + Learning**

- Parse Formula → Tree, Nodes = expressions, Edges = operations/functions
- 2. Transform Tree  $\rightarrow$  KAN Structure Moves leaf nodes to input layer
- 3. Create Final Graph by combining variables in first layer





# Extracting Knowledge

#### **Overview of Knowledge Extraction**



## 1. Identifying Important Features

- Assign scores to important inputs
- Previously: L1 norm (only considers local information)
  - $E_{I,i,j}$  = standard deviation of the function on each edge
  - $N_{II}$  = standard deviation at each node
- Update: Iterative computation of scores
  - node (attribution) score:  $A_{II}$
  - edge (attribution) score:  $B_{l,i,j}$

$$B_{l-1,i,j} = A_{l,j} \frac{E_{l,j}}{N_{l+1,j}}, \quad A_{l-1,i} = \sum_{j=0}^{n_l} B_{l-1,i,j}, \quad l = L, L-1, \cdots, 1.$$

Example



### 2. Identifying modular structures





anatomical modularity

functional modularity

# **Anatomical Modularity**

- Spatial proximity of neurons
- Autoswap: neuron swapping methods (preserves functionality)
- Easy to visually identify modular structures





### **Functional Modularity: separability**

Additively separable ~ 

- Compute the hessian matrix
- If  $H_{ii} = 0$  for all  $1 \le i \le k$ and  $k + 1 \le j \le n$  then additively separable
- For multiplicative separability, convert to additive separability by taking the log

$$f(x_1, x_2, \cdots x_n) = g(x_1, \dots, x_k) + h(x_{k+1}, \dots, x_n)$$
  
 $f(\mathbf{H} \equiv \nabla^T \nabla f(\mathbf{H}_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j})$ 

$$\log |f(x_1, x_2, \cdots, x_n)| = \log |g(x_1, \dots, x_k)| + \log |h(x_{k+1}, \dots, x_n)|$$

$$\mathbf{H}_{ij} \equiv rac{\partial^2 \mathrm{log} |f|}{\partial x_i \partial x_j}$$



#### 3. Identifying symbolic formulas

Trick A: discover and leverage modular structures



#### 3. Identifying symbolic formulas

Trick B: Sparse initialization — a more natural fit



#### 3. Identifying symbolic formulas

Trick C: Hypothesis testing - easy to try different ideas.



5

# Real-World Applications Application areas

# 1. Discovering Lagrangians

- Core idea: Learn lagrangian from data
  - Phase space: position and velocity
- Predict accelerations using Euler-Lagrange equations

How KANs Help:

- 1. Numerical Stability:
  - Pre-encode kinetic energy
- 2. Interpretability:

$$\begin{split} \mathbf{L}(\mathbf{q}, \dot{q}) &= T - V from(q, \dot{q}, \ddot{q}) \\ \ddot{q} &= (\nabla_{\dot{q}} \nabla_{\dot{q}}^T L)^{-1} \left[ \nabla_q L - \nabla_q \nabla_{\dot{q}}^T \dot{q} \right] \end{split}$$

• Symbolic regression extracts explicit formulas from KAN edges

#### **KAN** solution

- Key Results:
  - **Stable Training**: Fewer numerical errors.
  - Explainable Models: Learns physics-aligned formulas.



#### 2. Constitutive laws

- Constitutive Laws:
  - Define material behavior under forces or deformation (e.g., Hooke's Law).
  - Linear for small deformations (Hooke's law) —
  - Nonlinear for larger deformations (e.g., Neo-Hookean).
- Why Use KANs?

$$-P_{12} = \mu(F_{11}F_{21} + F_{12}F_{22} + F_{13}F_{23})$$

 $^{-}P_{12} = \mu(F_{12} + F_{21})$ 

- predict the P (stress) tensor from the F (deformation) tensor
- Incorporates prior knowledge (e.g., linear laws).
- Extracts interpretable symbolic formulas via symbolic regression.

### **Results**:

#### With Prior Knowledge:

- The initial KAN is based on the linear constitutive law.
- Helps guide training but may converge to poor local minima
- Without Prior Knowledge:
  - Start with a randomly initialized KAN.
  - Requires more data to train but avoids being biased

 $P_{12} = \mu(F_{12} + F_{21})$ 



 $P_{12} = \mu(F_{11}F_{21} + F_{12}F_{22} + F_{13}F_{23})$ 



 $P_{11} = 0.42(F_{11}^2 + F_{12}^2 + F_{13}^2 - 1) + 0.28\log(|F|)$ 



# **Future**

#### Software



#### (b) "software version"-scale plane

# Outlook

#### • Scalability Challenge

- Limited interpretability at large scales
- Need better methods for managing complexity

#### • Key Research Areas

- Advanced interpretability methods
- Scaling to larger problems
- Expansion beyond physics

#### KAN VS MLP

MLP:

- Lower interpribility
- Well established and proven

KAN:

- 10x Slower training time
- Higher interpretability
- Continual learning
- More parameter efficient

#### **References**

1. Z. Liu, Y. Wang, S. Vaidya, F. Ruehle, J. Halverson, M. Soljačić, T. Y. Hou, and M. Tegmark, "KAN: Kolmogorov-Arnold Networks," *arXiv preprint*, arXiv:2404.19756, Jun. 2024. [Online]. Available: <u>https://doi.org/10.48550/arXiv.2404.19756</u>

2. Z. Liu, P. Ma, Y. Wang, W. Matusik, and M. Tegmark, "KAN 2.0: Kolmogorov-Arnold Networks meet science," *arXiv preprint*, arXiv:2408.10205, Aug. 2024. [Online]. Available: <u>https://doi.org/10.48550/arXiv.2408.10205</u>

$$\phi(x) = w_b b(x) + w_s \operatorname{spline}(x).$$
  
 $b(x) = \operatorname{silu}(x) = x/(1 + e^{-x})$   
 $\operatorname{spline}(x) = \sum_i c_i B_i(x)$ 

#### **1. Discovering conserved quantities**

- Conserved quantities as differential equation solving
- $z \in R^d$  governed by the equation dz/dt = f(z)
- H(z) is a conserved quantity if  $f(z) \cdot \nabla H(z) = 0$  for all Z

#### Harmonic Oscillator

- System Description:
  - Phase space: z = (x, p)
    - *x*: position
    - *p*: momentum
- Evolution:
  - $\circ \quad d(x,p)/dt = (p, -x)$

#### **Harmonic Oscillator**

#### **Energy Conservation:**

- Energy function:  $H = \frac{1}{2}(x^2 + p^2)$
- Phase space: z = (x, p)
- Conservation proof:
  - Force vector: f(z) = (p, -x)
  - Gradient of *H*:  $\nabla H(z) = (x, p)$
  - Dot product:  $f(z) \cdot \nabla H(z) = p \cdot x + (-x) \cdot p = 0$
  - Zero result proves H is conserved

#### KAN:

- Parametrize *H* into network
- Train with loss function

