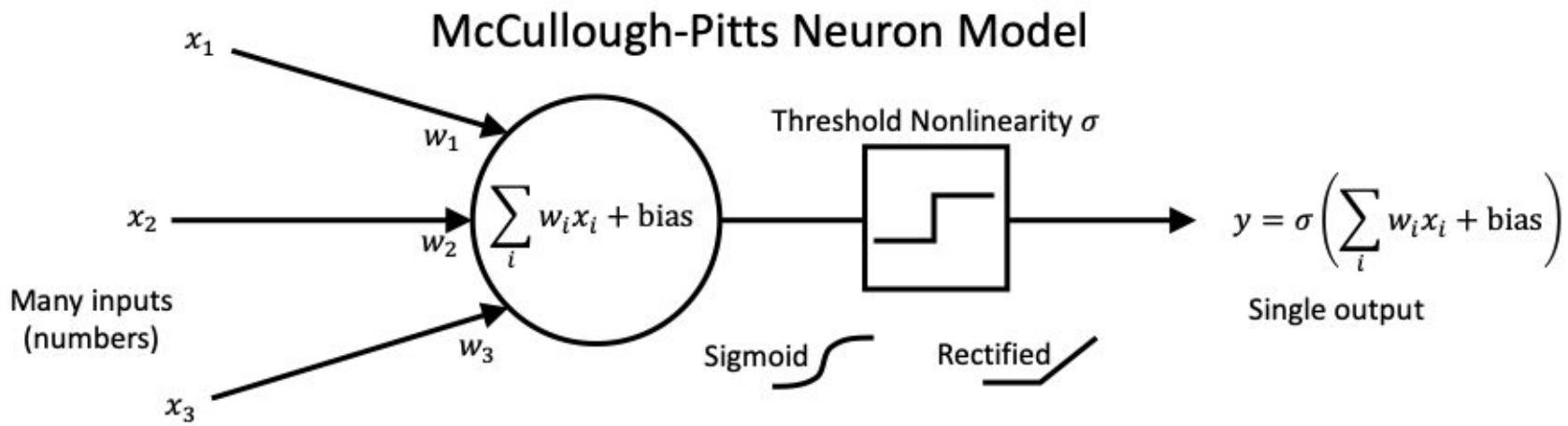
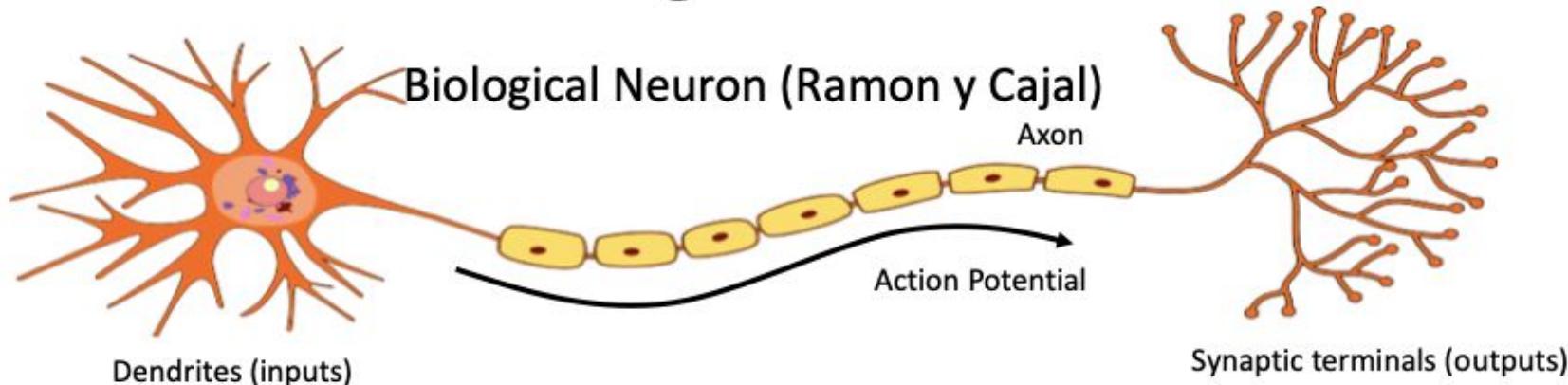


# Kolmogorov-Arnold Networks and Implicit Neural Representations (INRs)

Ben Wyant

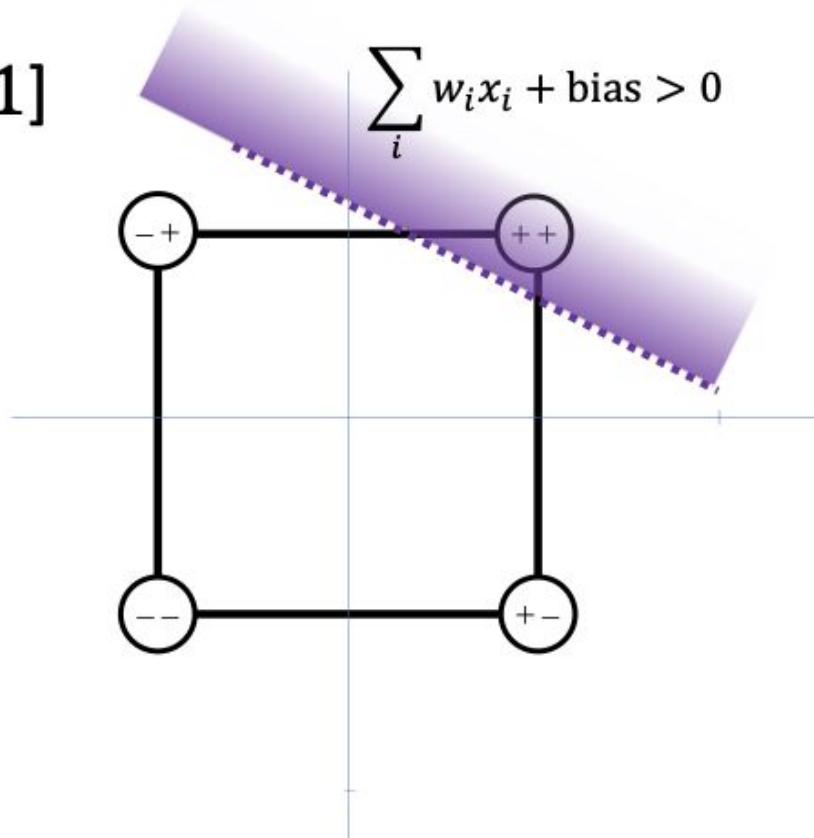
# McCullough-Pitts neuron



# What logic can be learned by one neuron?

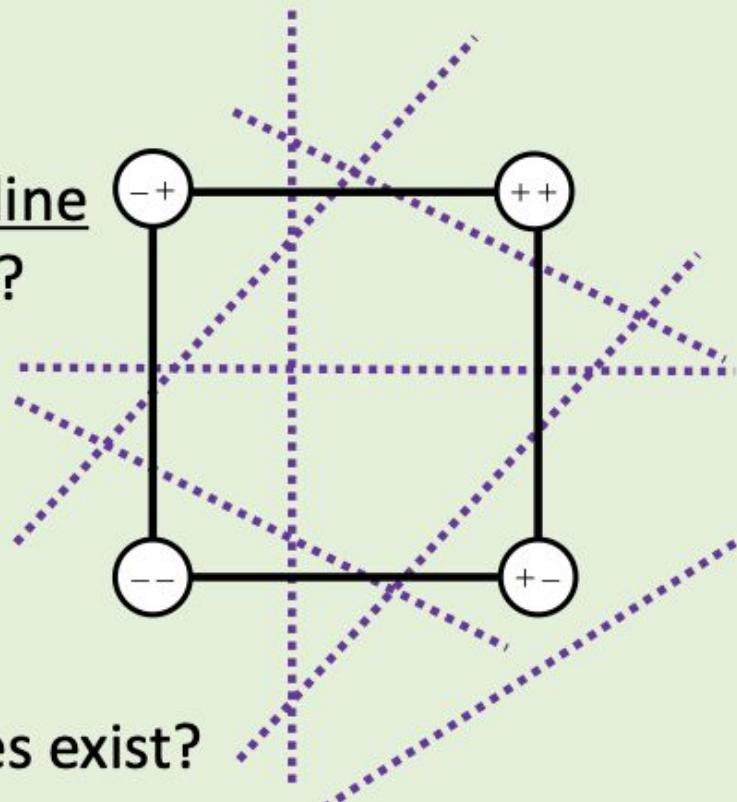
$$x_a, x_b = [\pm 1, \pm 1]$$

a	b	out
-1	-1	-1
-1	+1	-1
+1	-1	-1
+1	+1	+1



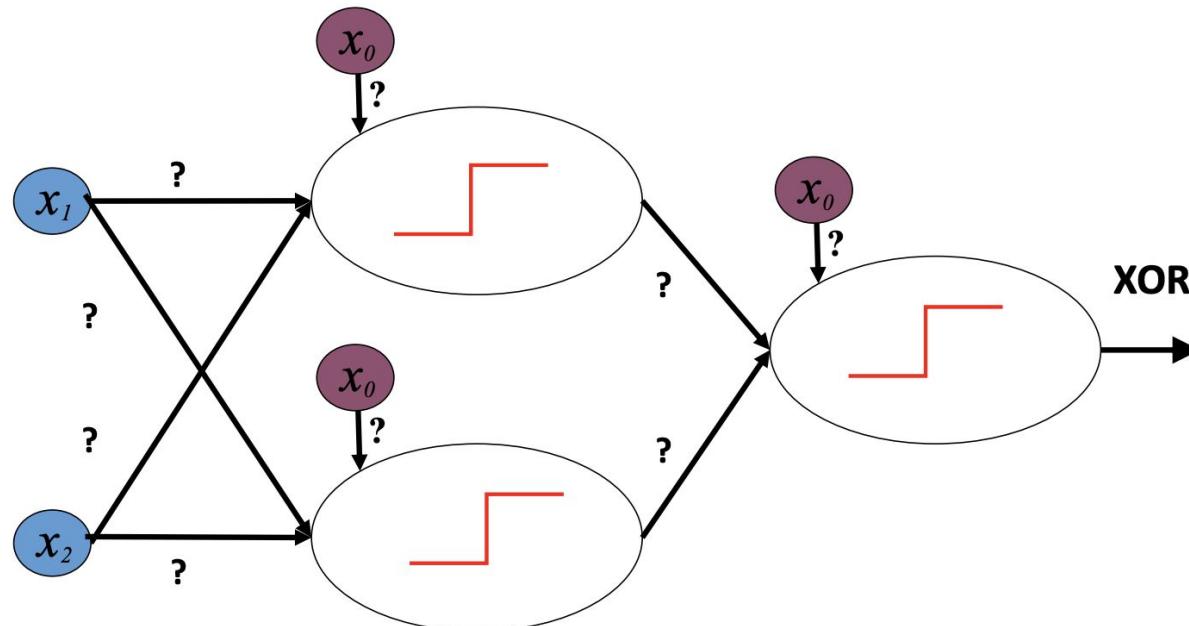
# What logic can be learned by one neuron?

How many ways can a straight line partition corners of a square?



How many d-variable truth-tables exist?

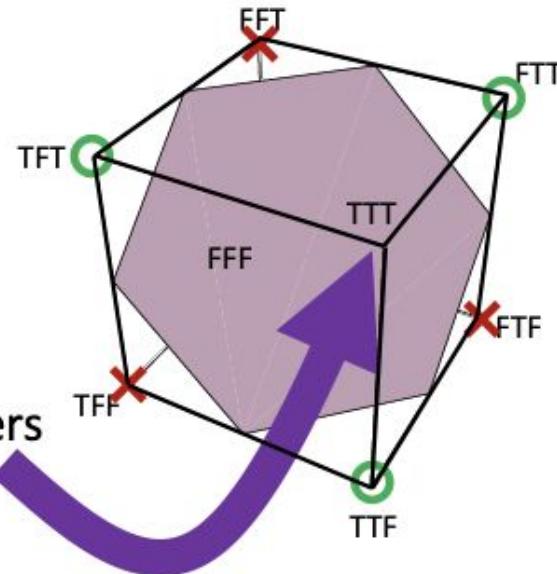
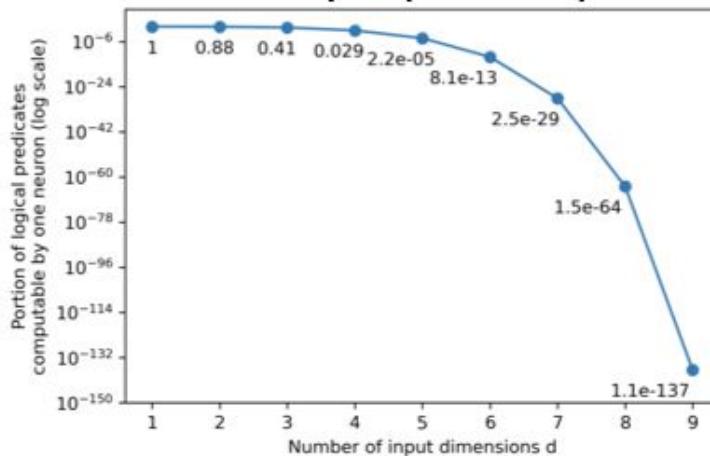
Combining perceptrons can make any Boolean function  
...if you can set the weights & connections right



How would you set the weights and connections to make XOR?

# Not all d-variable truth tables are linearly separable

Portion of linearly separable predicates



In a linear model, some corners determine other corners

In  $d=2$ , 88% of predicates are linearly separable.

In  $d=3$ , 41% of predicates are linearly separable.

In  $d=4$  it is only 2.9%.

[counting for  $d>10$  remains unsolved!]

# Puzzle

How many layers of neurons are needed to  
be able to express all truth tables on d  
inputs?

# Universal Approximation Theorem

Given:

any continuous function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,

a compact input domain  $X \subset \mathbb{R}^n$

and an error interval  $\varepsilon > 0$

There is guaranteed to be:

a finite network of neurons  $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$

that approximates  $|g(x) - f(x)| < \varepsilon$  for all  $x \in X$

Cybenko, 1989 proves this for sigmoid/step neurons

Hornik, 1991 proves almost any nonlinearity will work

# Visual Proof (Figures from Andrea Lörke)

Two neurons (step-up, step-down) can be summed to make a bump

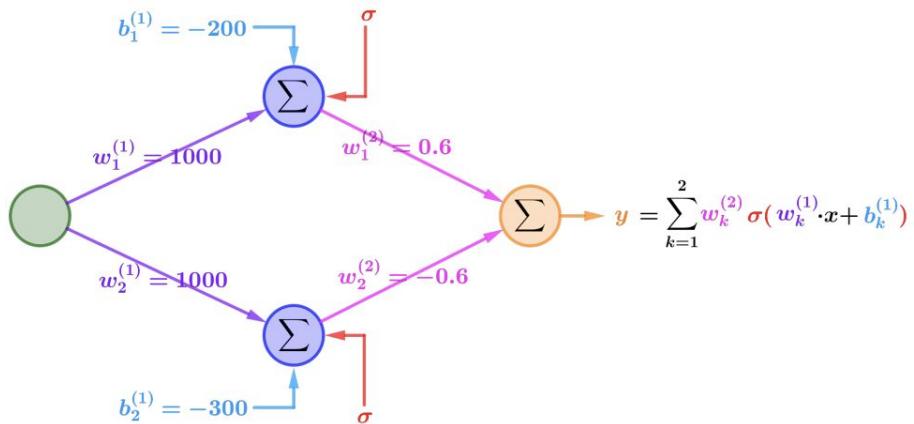


Figure 15: Neural network with 2 hidden nodes

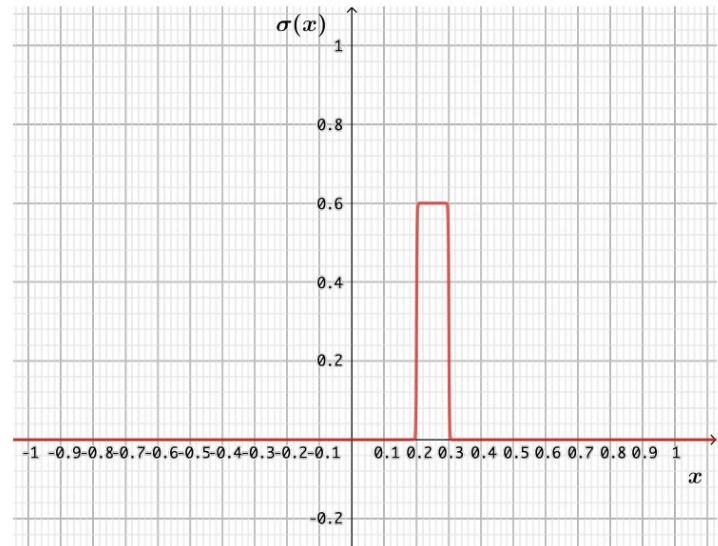


Figure 16: Coupled nodes

# Visual Proof (Figures from Andrea Lörke)

Many pairs of neurons can be summed up to sum up bumps

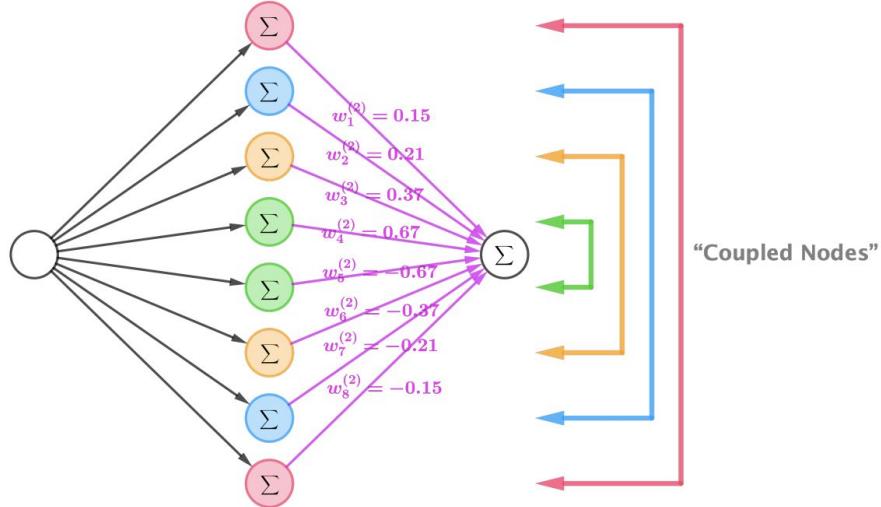


Figure 17: neural network with coupled nodes



Figure 18: Approximation (N=4)

# Visual Proof (Figures from Andrea Lörke)

By increasing the number of neuron pairs, error can be made  $< \varepsilon$



Figure 18: Approximation ( $N=4$ )

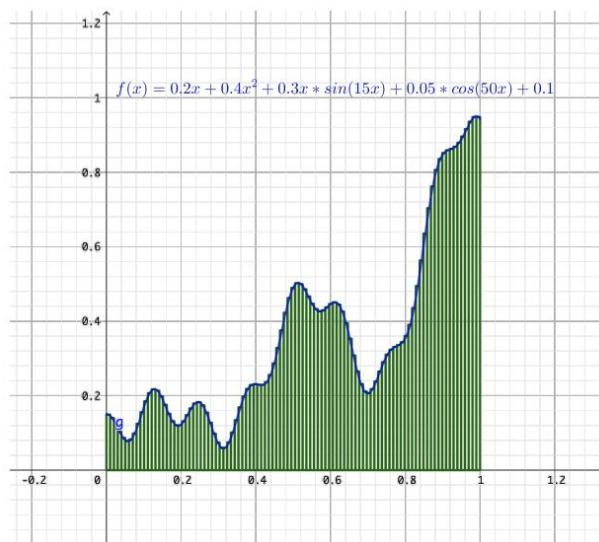


Figure 19: Approximation ( $N=100$ )

# Visual Proof (Figures from Andrea Lörke)

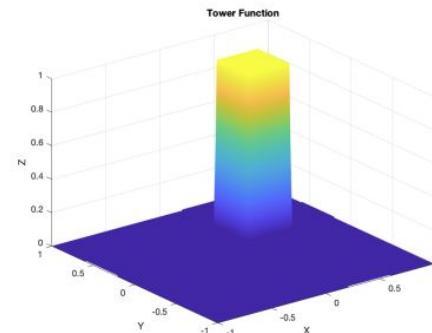
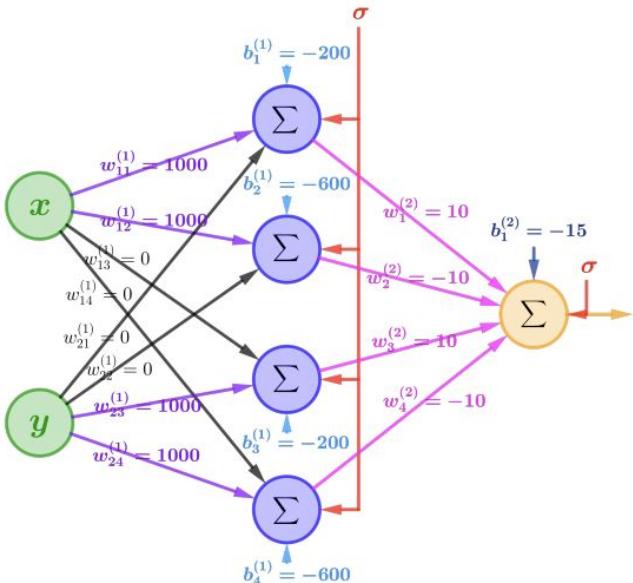


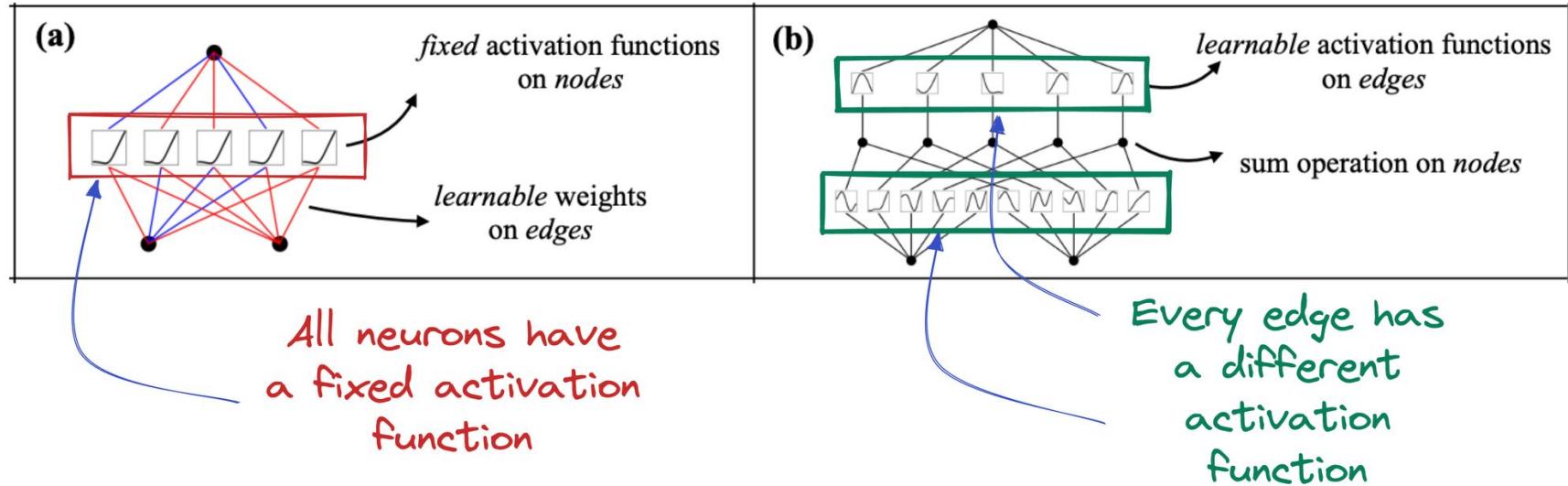
Figure 33:  
"Tower"-function

Figure 32: Nodes to compute a Tower-function  
at any Position

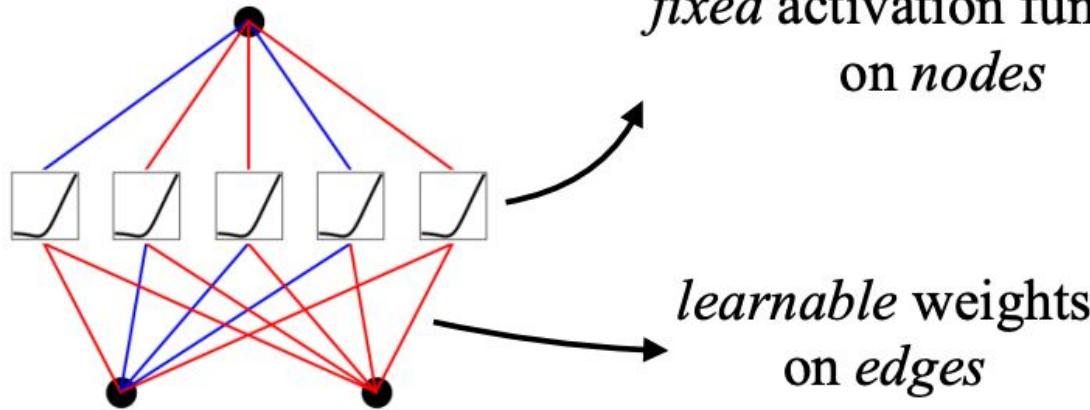
# Kolmogorov Arnold Network (KAN)

## Neural Network

## KAN



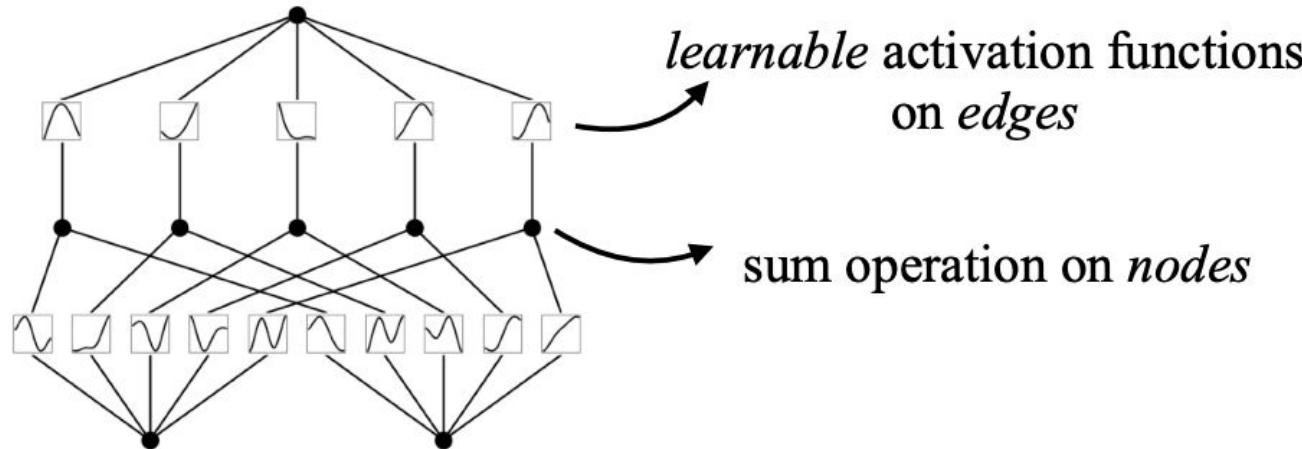
**(a)**



*fixed activation functions  
on nodes*

*learnable weights  
on edges*

**(b)**



*learnable activation functions  
on edges*

*sum operation on nodes*

# Kolmogorov-Arnold Representation Theorem

- Any multivariate continuous function can be represented as the composition of a **finite** number of continuous functions of a single variable

$$\sum_{j=1}^m \psi_j \left( \sum_{i=1}^n \phi_{ij}(x_i) \right)$$

- “Multivariate continuous function” is a function that accepts multiple parameters:

$$y = F(x_1, x_2, x_3, \dots, x_n)$$

- “A finite number of continuous functions of a single variable” means:

**Univariate functions**

$$\phi_1(x_1) + \phi_2(x_2) + \phi_3(x_3) + \dots + \phi_n(x_n)$$

The diagram illustrates the decomposition of a multivariate function into a sum of univariate functions. At the top, the text "Univariate functions" is written in a stylized, handwritten font. Three arrows originate from this text and point to the individual terms of the expression below:  $\phi_1(x_1)$ ,  $\phi_2(x_2)$ , and  $\phi_n(x_n)$ .

## Univariate functions

$$\phi_1(x_1) + \phi_2(x_2) + \phi_3(x_3) + \cdots + \phi_n(x_n)$$

- The above sum is passed through one more function  $\psi$  and then a composition is applied

$$\psi\left(\sum_{i=1}^n \phi_i(x_i)\right)$$

$$\sum_{j=1}^m \psi_j\left(\sum_{i=1}^n \phi_{ij}(x_i)\right)$$

Multivariate  
continuous  
function

Composition of  
univariate functions

$$F(x_1, x_2, x_3, \dots, x_n) = \sum_{j=1}^m \psi_j \left( \sum_{i=1}^n \phi_{ij}(x_i) \right)$$

Multivariate continuous  
function

Composition of  
univariate functions

$$\begin{aligned} F(x_1, x_2, x_3, \dots, x_n) &= \psi_1(\phi_{11}(x_1) + \phi_{21}(x_2) + \dots + \phi_{n1}(x_n)) \\ &\quad + \psi_2(\phi_{12}(x_1) + \phi_{22}(x_2) + \dots + \phi_{n2}(x_n)) \\ &\quad \vdots \\ &\quad + \psi_m(\phi_{1m}(x_1) + \phi_{2m}(x_2) + \dots + \phi_{nm}(x_n)) \end{aligned}$$

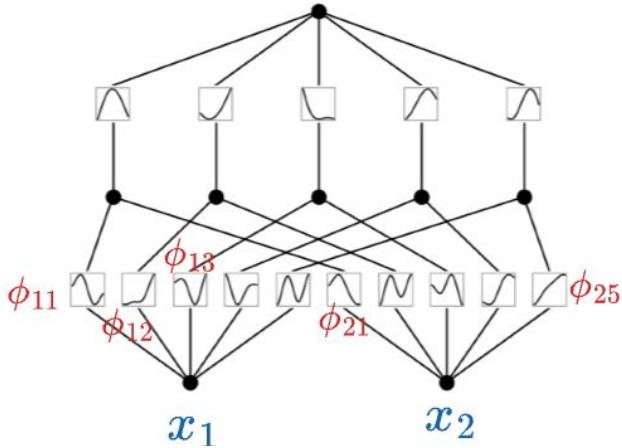
Multivariate continuous  
function

Composition of  
univariate functions

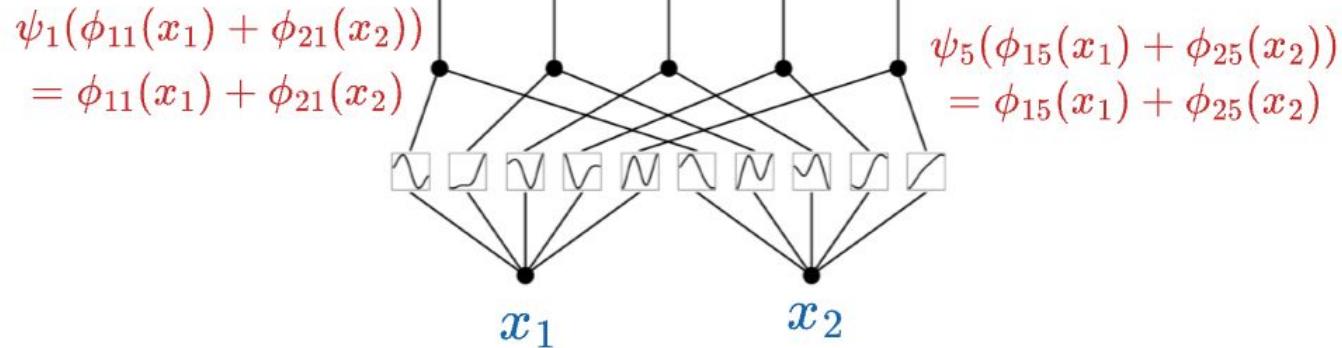
$$F(x, y) = xy$$

$$F(x, y) = \boxed{\exp} (\boxed{\log(x)} + \boxed{\log(y)})$$

$\psi \quad \phi_x \quad \phi_y$



In this case  $\psi(z) = z$



# Key Differences

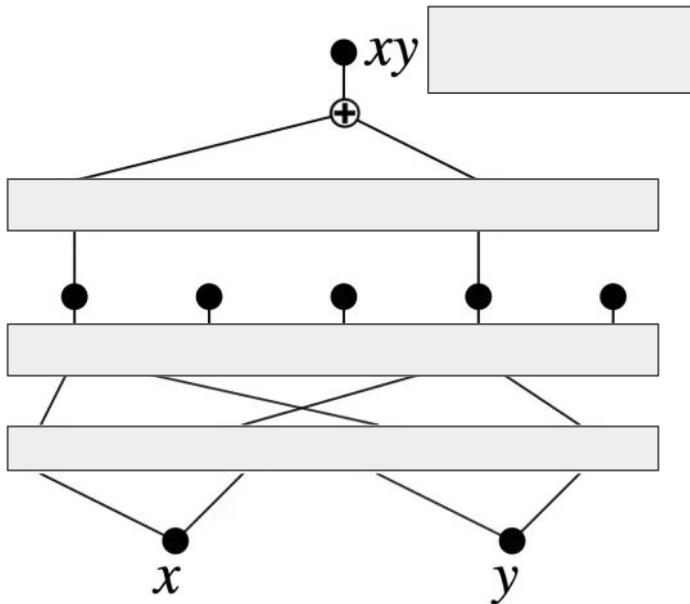
- Universal Approximation Theorem deals with **approximating functions** using neural networks
- Kolmogorov-Arnold theorem provides a way to **exactly represent** continuous functions through sums of univariate functions



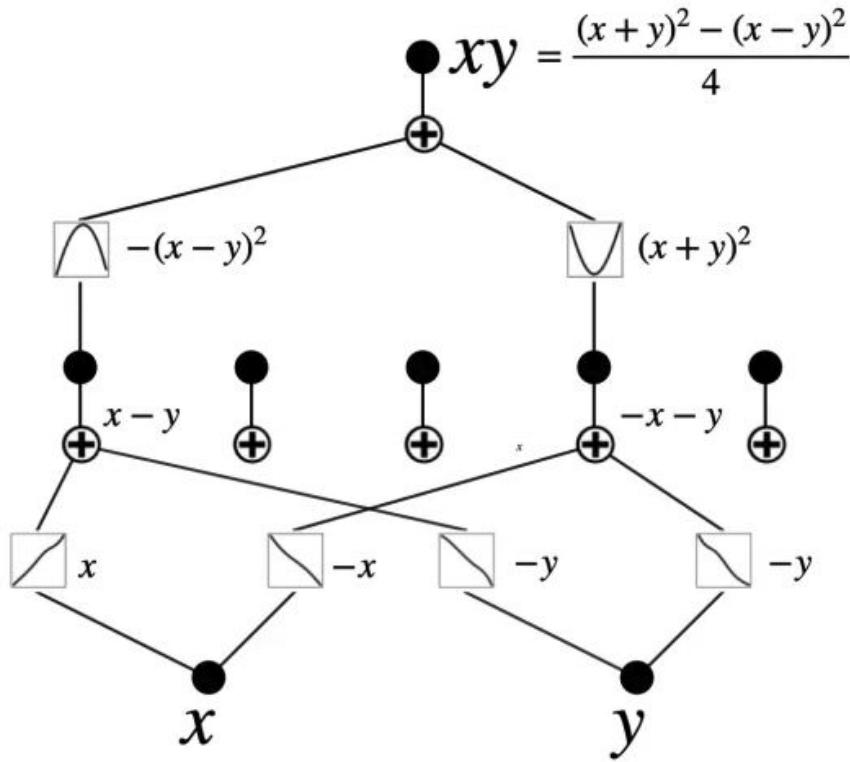
Can MLPs do multiplication?

# Puzzle

Can you come up with another way to represent  $xy$ ?

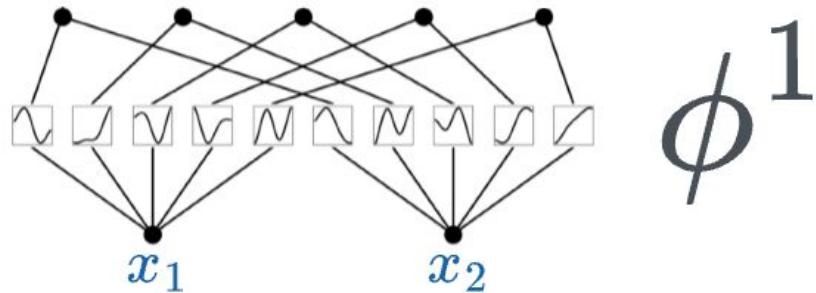


# Answer



# KAN matrix representation

- A single KAN layer is taking the input  $(x_1, x_2, \dots, x_n)$  and applying a transformation  $\phi$
- $\phi^1$  can be represented with the matrix below where  $n$  is the number of inputs and  $m$  is the number of outputs to layer  $l_1$



$$\phi^1 = \begin{bmatrix} \phi_{11} & \phi_{12} & \dots & \phi_{1n} \\ \phi_{21} & \phi_{12} & \dots & \phi_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{m1} & \phi_{m2} & \dots & \phi_{mn} \end{bmatrix}$$

Transformation matrix  $\phi^1$       Input  $X$

$$z^1 = \begin{bmatrix} \phi_{11} & \phi_{12} & \dots & \phi_{1n} \\ \phi_{21} & \phi_{12} & \dots & \phi_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{m1} & \phi_{m2} & \dots & \phi_{mn} \end{bmatrix} * \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Output of  
first layer

Function  
applied to  $x_1$

Function  
applied to  $x_n$

$$z^1 = \begin{bmatrix} \boxed{\phi_{11}(x_1)} + \phi_{12}(x_2) + \dots + \boxed{\phi_{1n}(x_n)} \\ \phi_{21}(x_1) + \phi_{22}(x_2) + \dots + \phi_{2n}(x_n) \\ \vdots \\ \phi_{m1}(x_1) + \phi_{m2}(x_2) + \dots + \phi_{mn}(x_n) \end{bmatrix}$$

# KAN Formula Representation

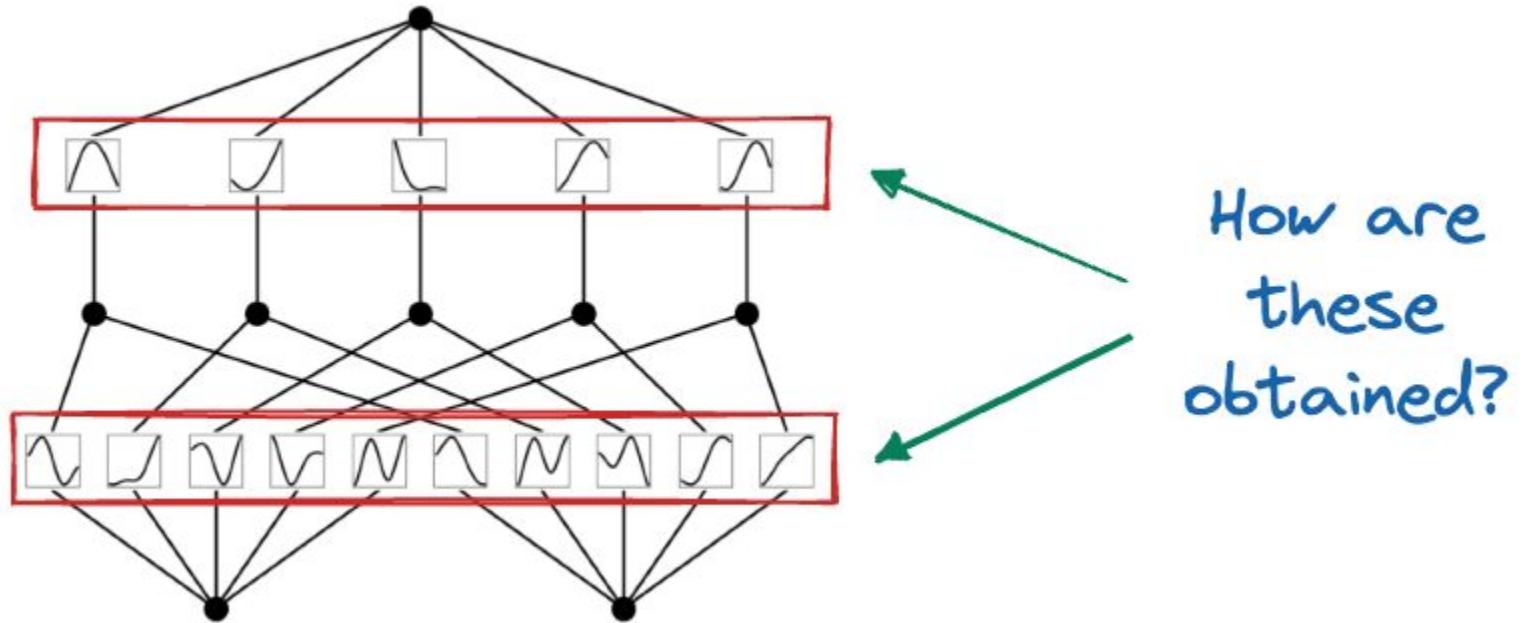
- The entire KAN network can be condensed into one formula:

$$KAN(x) = \phi^L ( \phi^{L-1} ( \dots ( \phi^2 ( \phi^1(x) ) ) ) )$$

Where:

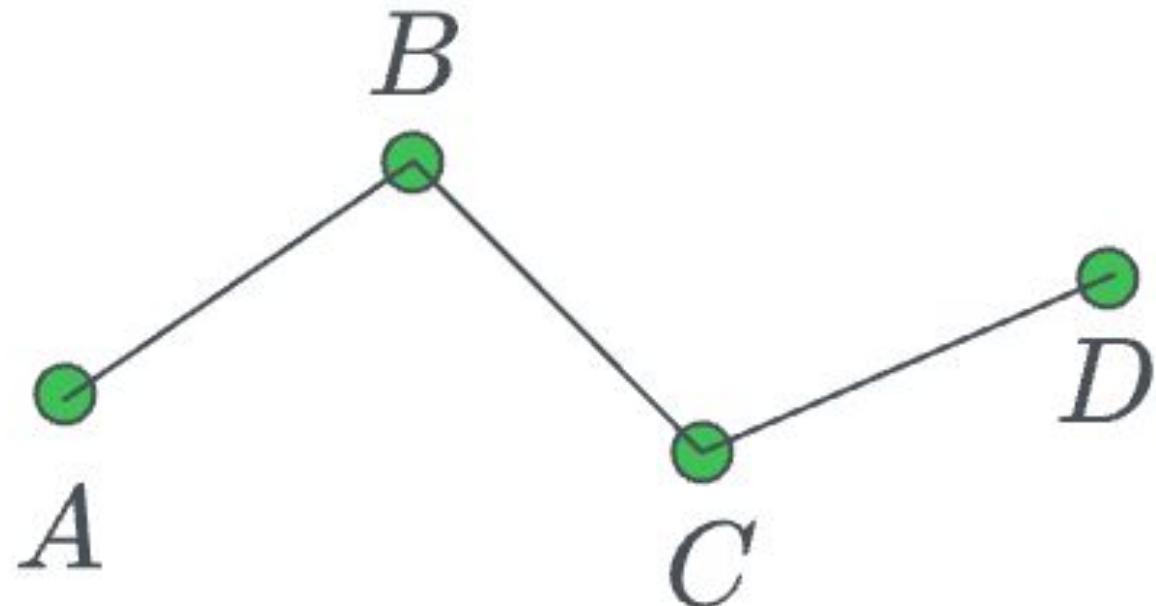
- $x$  is the input vector
- $\phi^k$  is the function transformation matrix of layer  $k$
- $KAN(x)$  is the output of the KAN

# How are KANs trained?



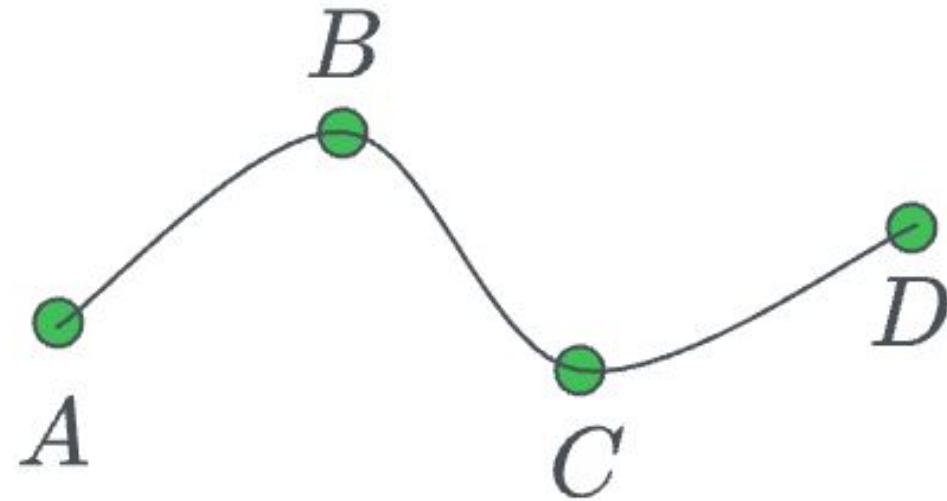
# Let's talk about splines

Character



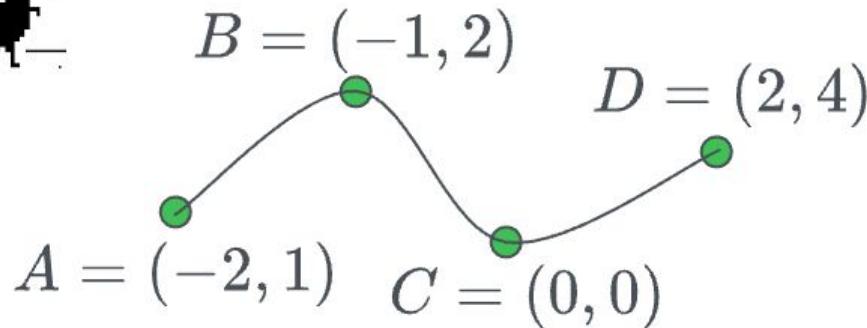
# Bezier Curves

Character



$$f(x) = ax^3 + bx^2 + cx + d$$

## Character



Solve these equations

$$A \rightarrow 1 = a(-2)^3 + b(-2)^2 + c(-2) + d$$

$$B \rightarrow 2 = a(-1)^3 + b(-1)^2 + c(-1) + d$$

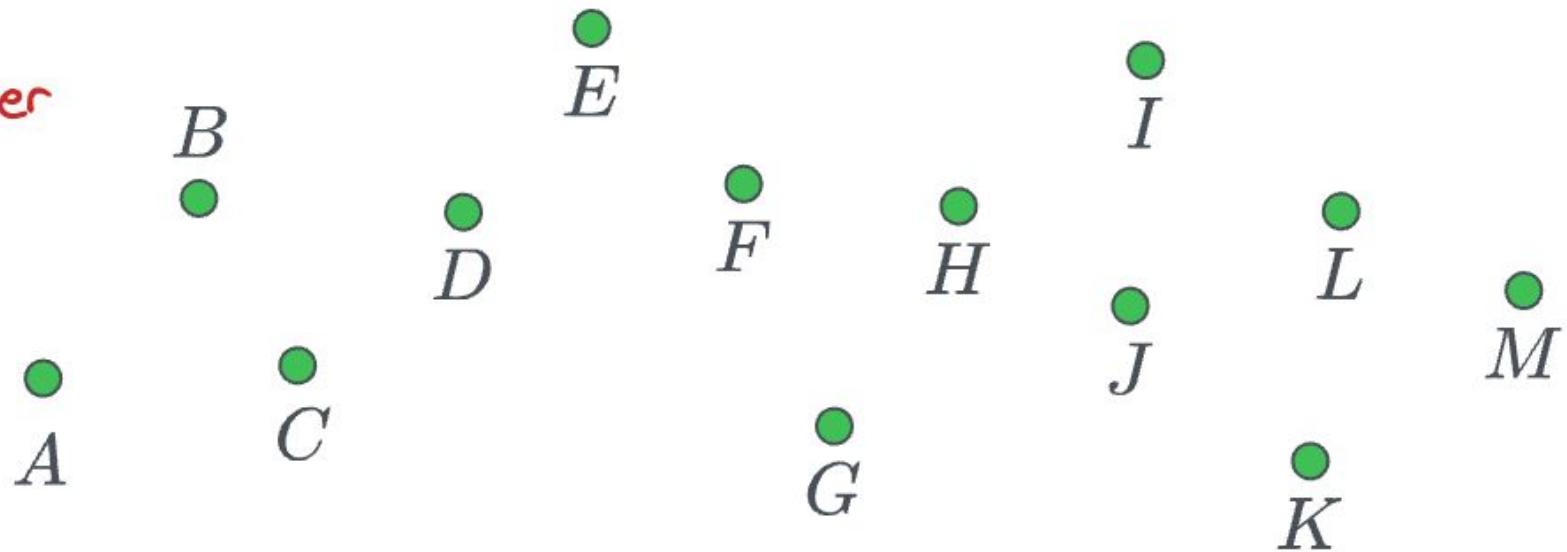
$$C \rightarrow 0 = a(0)^3 + b(0)^2 + c(0) + d$$

$$D \rightarrow 4 = a(2)^3 + b(2)^2 + c(2) + d$$

- We need 4 points to solve for an equation with 4 variables

# What do we do if we have hundreds of data points?

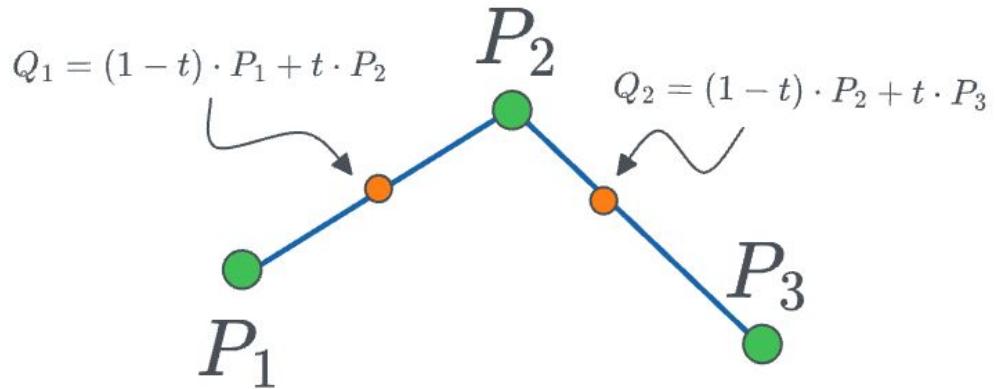
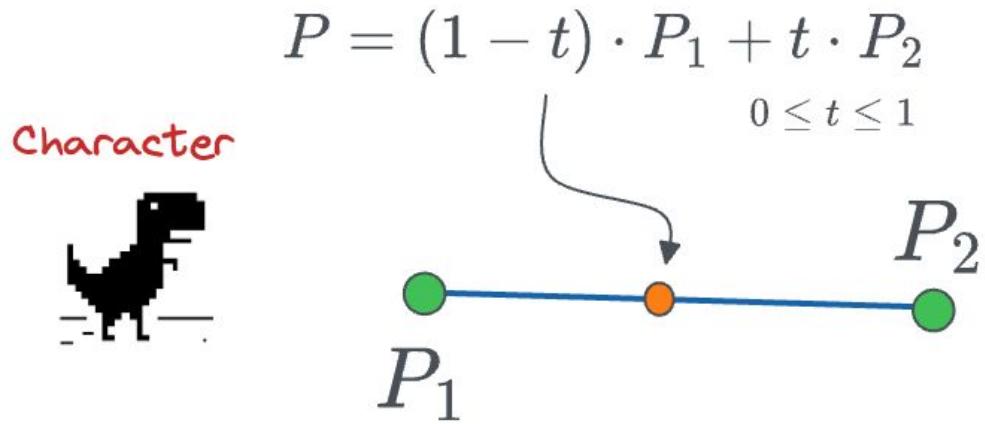
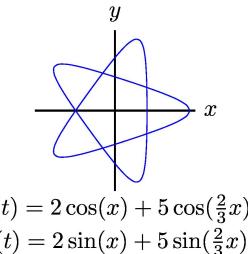
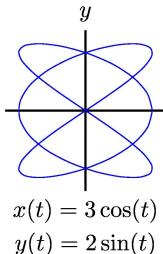
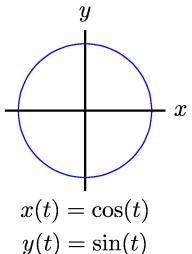
Character



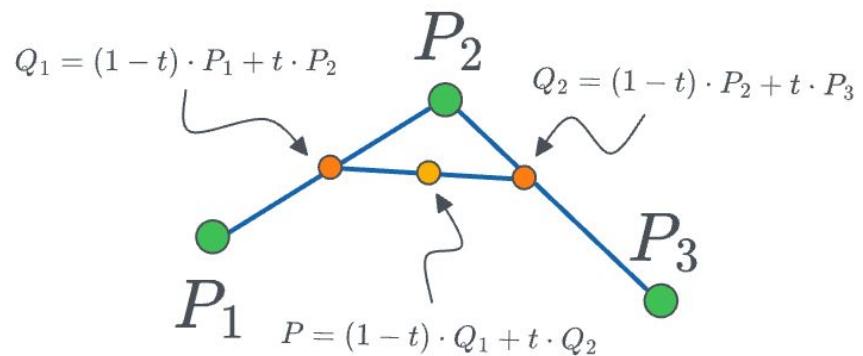
# Bezier Curves

- When  $t= 0$ , P will be  $P_1$
- When  $t= 1$ , P will be  $P_2$
- $t \in [0,1]$

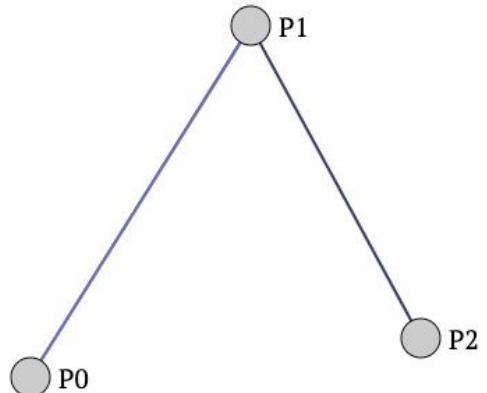
Some other parametric equations:



Between points  $Q_1$  and  $Q_2$  we can create a linear interpolation resulting in the **output** curve



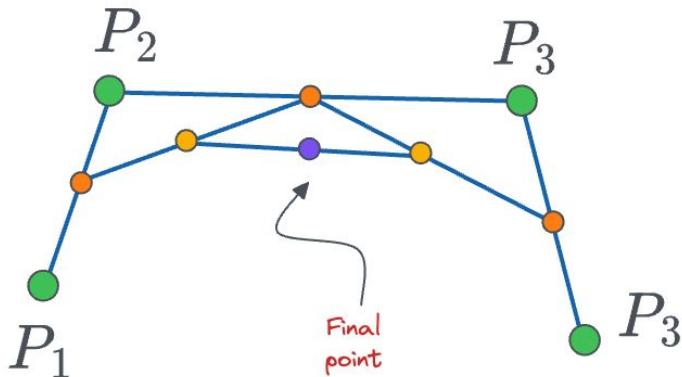
$$B(t) = (1 - t)^2 P_1 + 2(1 - t)t P_2 + t^2 P_3, \quad t \in [0, 1]$$



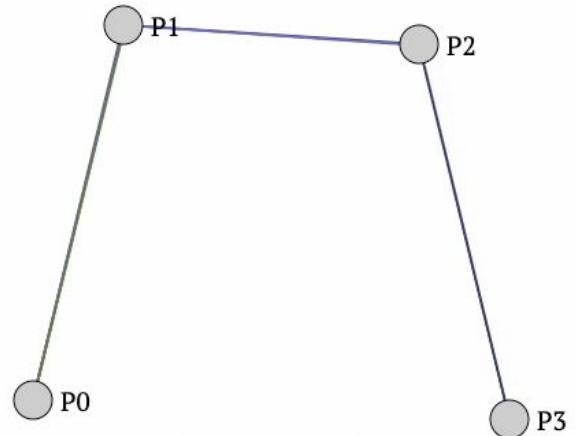
**t=0.00**

# Cubic Bezier Curve

We can extend this same idea to 4 points which results in this new **output** curve



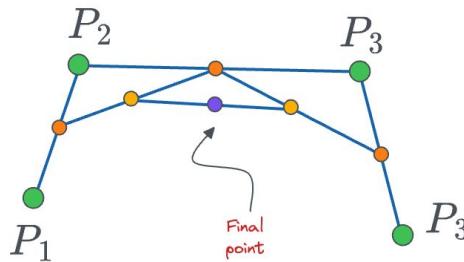
$$B(t) = (1 - t)^3 P_1 + 3(1 - t)^2 t P_2 + 3(1 - t)t^2 P_3 + t^3 P_4, \quad t \in [0, 1]$$



**t=0.00**

# Bezier Curve Recap

These coefficients match with the binomial expansion of the polynomial  $(1 + x)^n$



$$B(t) = (1-t)^3 P_1 + 3(1-t)^2 t P_2 + 3(1-t)t^2 P_3 + t^3 P_4, \quad t \in [0, 1]$$

2 points

$$B(t) = (1-t)P_1 + tP_2, \quad t \in [0, 1]$$

3 points

$$B(t) = (1-t)^2 P_1 + 2(1-t)t P_2 + t^2 P_3, \quad t \in [0, 1]$$

4 points

$$B(t) = (1-t)^3 P_1 + 3(1-t)^2 t P_2 + 3(1-t)t^2 P_3 + t^3 P_4, \quad t \in [0, 1]$$

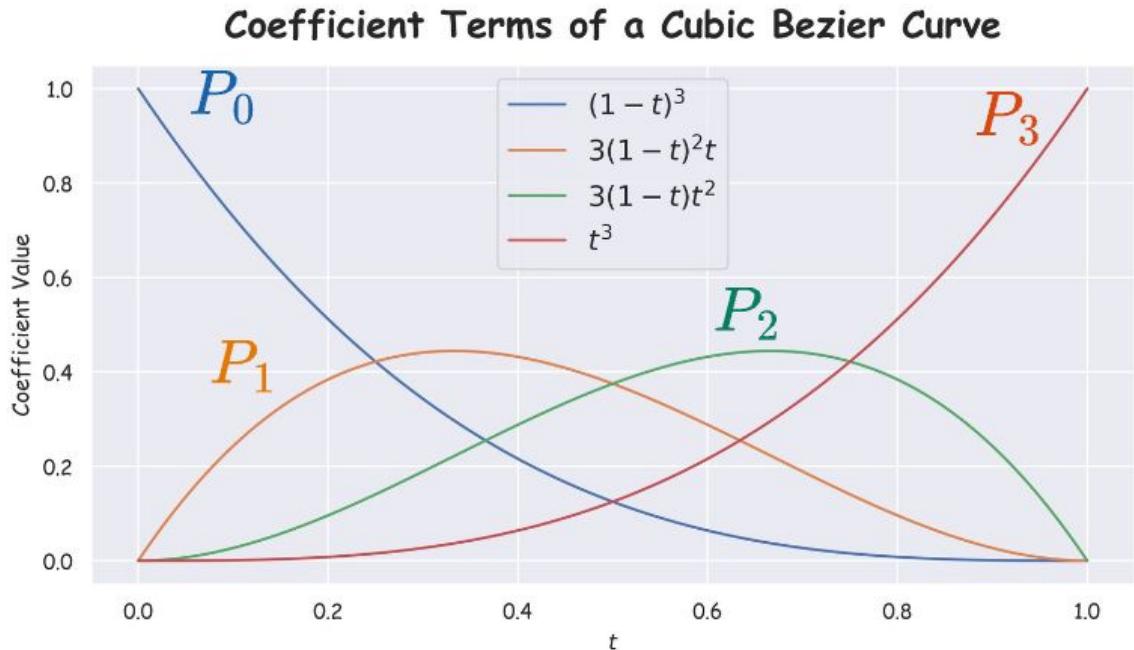
$$B(t) = \sum_{i=0}^n \binom{n}{i} (1-t)^{n-i} t^i P_i ; \quad \binom{n}{i} = \frac{n!}{i!(n-i)!}$$

$$= \sum_{i=0}^n c_{i,n}(t) P_i$$

Binomial coefficient

# Plotting the Bezier Curve

- At  $t= 0$ , except  $P_0$  the coefs of all points are 0 so the curve starts from  $P_0$
- From  $t \in (0.25, 0.5)$ ,  $P_1$  is max, the curve is closest to  $P_1$  in that duration
- From  $t \in (0.5, 0.75)$ ,  $P_2$  is max, the curve is closest to  $P_2$  in that duration
- At  $t= 1$ , except  $P_3$ , coefs of all points are 0 so the curve ends at  $P_3$



# Problem with Bezier Curves

- Having 100 data points will result in a polynomial of degree 99, which will be computationally expensive
- Factorial terms are going to be expensive as well

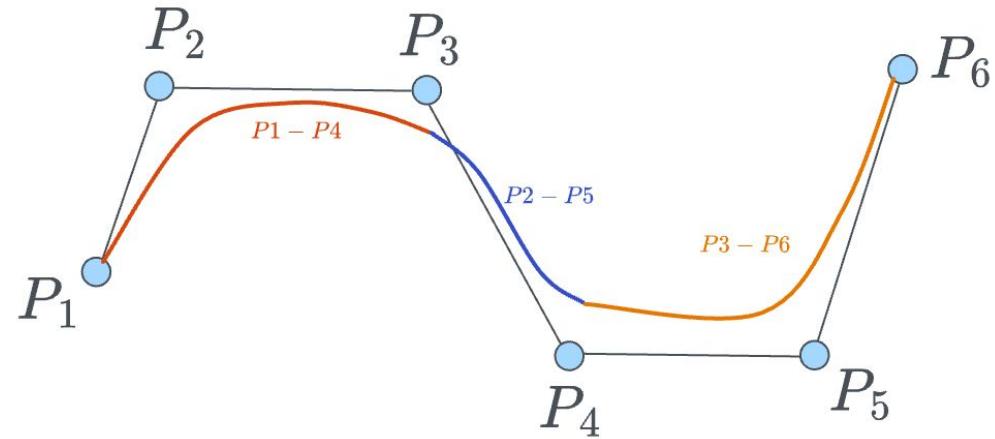
$$B(t) = \sum_{i=0}^n \binom{n}{i} (1-t)^{n-i} t^i P_i ; \quad \binom{n}{i} = \frac{n!}{i!(n-i)!}$$

$$= \sum_{i=0}^n c_{i,n}(t) P_i$$

Binomial coefficient

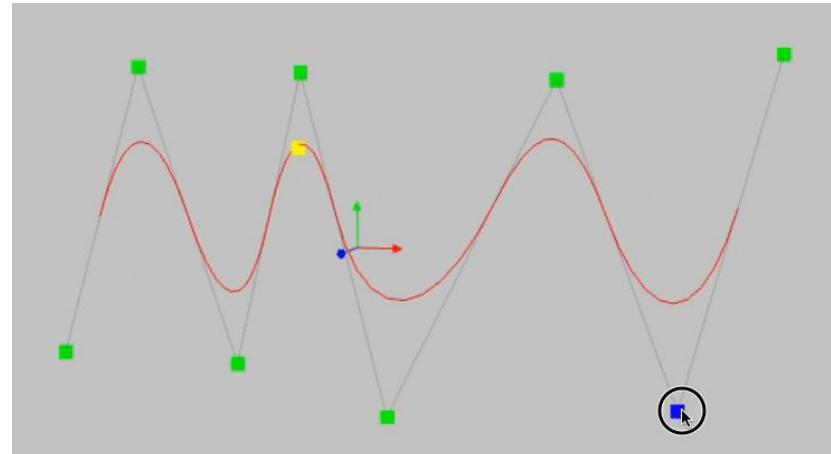
# B-Splines

With 6 points we could generate a degree 5 Bezier curve - this is still expensive



Instead, we can create curves of smaller degrees, say 3, and connect them

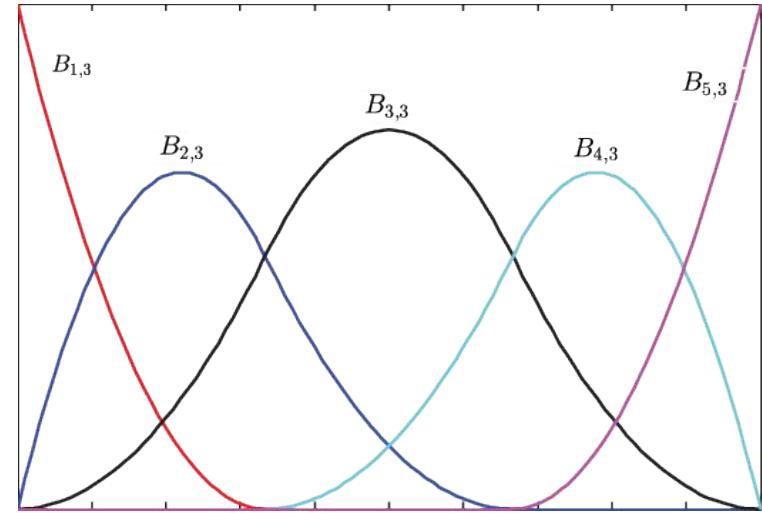
When we have  $n$  control points and create  $k$  degree polynomial Bezier curves we get  $(n-k)$  Bezier curves in the final B-spline



# How do we connect these curves?

We must ensure  $C^2$  continuity

- The curvature is smooth across the joins (also called knots)
- $P_i$ : Control points that define the shape of the curve
- $N_{i,k}(t)$ : B-spline basis functions of degree  $k$  associated with each control point  $P_i$



$$S(t) = \sum_{i=0}^n P_i \cdot N_{i,k}$$

control points                      Basis function

<https://www.dailydoseofds.com/a-beginner-friendly-introduction-to-kolmogorov-arnold-networks-kan/>

[https://www.researchgate.net/figure/Uniform-B-spline-basis-functions-for-N-5-k-3-and-X-0001323111\\_fiq2\\_282803553](https://www.researchgate.net/figure/Uniform-B-spline-basis-functions-for-N-5-k-3-and-X-0001323111_fiq2_282803553)

# Building a KAN layer

By adjusting the positions of the control points during training, the KAN model can dynamically shape the activation functions that best fit the data

- During training, update the positions of these control points through backprop
- Use gradient descent to adjust the control points

$$\phi(x) = w( b(x) + \text{spline}(x) )$$

- $b(x)$  is a basis function which acts as residual connection
- $b(x)$  helps mitigate the vanishing gradient problem
- spline ( $x$ ) is learnable, specifically, the points  $c_i$
- In principle,  $w$  is redundant, but it is included because it helps to better control magnitudes of the activation function

$$\phi(x) = w \left( b(x) + \text{spline}(x) \right)$$

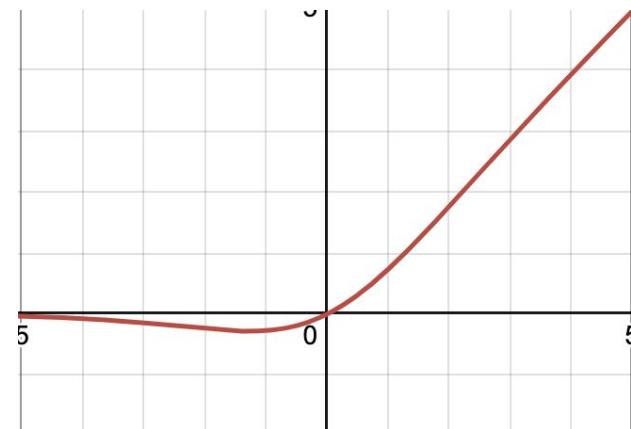
learnable      fixed

$$b(x) = \frac{x}{1 + e^{-x}}$$

$$\text{spline}(x) = \sum_i c_i B_i(x)$$

learnable

$$b(x) = \text{silu}(x)$$



# Initializing the weights

$$\phi(x) = \textcolor{red}{w_b} b(x) + \textcolor{green}{w_s} \text{spline}(x)$$

$$\text{spline}(x) \approx 0 \qquad \qquad \qquad \textcolor{green}{w_s} = 1$$

$$c_i \sim \mathcal{N}(0, \sigma^2)$$

small value

$$\textcolor{red}{w_b} \sim N\left(0, \frac{2}{n_{\text{in}} + n_{\text{out}}}\right)$$

Now we can train a network

- Initialize each of the  $\phi$  matrices:

$$\phi^1 = \begin{bmatrix} w(b(x) + \sum_i c_i B_i(x)) \\ \phi_{11}(\cdot) & \phi_{12}(\cdot) & \dots & \phi_{1n}(\cdot) \\ \phi_{21}(\cdot) & \phi_{22}(\cdot) & \dots & \phi_{2n}(\cdot) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{m1}(\cdot) & \phi_{m2}(\cdot) & \dots & \phi_{mn}(\cdot) \end{bmatrix}$$

- Run the forward pass

$$KAN(x) = \phi^L (\phi^{L-1} (\dots (\phi^2 (\phi^1(x)))) )$$

- Calculate loss and run backpropagation

# KAN Simplification techniques - Sparsification

For MLPs, L1 regularization of linear weights is used to favor sparsity. KANs can adapt this high-level idea, but need two modifications:

1. There is no linear “weight” in KANs. Linear weights are replaced by learnable activation functions, so we should define the L1 norm of these activation functions.
2. L1 regularization is insufficient for sparsification of KANs; instead an additional entropy regularization is necessary

# L1 norm for KAN

- L1 norm of an activation function  $\phi$  to be its average magnitude over its  $N_p$  inputs

$$|\phi|_1 = \frac{1}{N_p} \sum_{s=1}^{N_p} |\phi(x^{(s)})|$$

- For a KAN layer  $\Phi$  with  $n_{\text{in}}$  inputs and  $n_{\text{out}}$  outputs, the L1 norm of  $\Phi$  is the sum of L1 norms of all the activation functions:

$$|\Phi|_1 = \sum_{i=1}^{n_{\text{in}}} \sum_{j=1}^{n_{\text{out}}} |\phi_{i,j}|_1$$

# KAN entropy

- We define the entropy of  $\Phi$  to be:

$$S(\Phi) \equiv - \sum_{i=1}^{n_{\text{in}}} \sum_{j=1}^{n_{\text{out}}} \frac{|\phi_{i,j}|_1}{|\Phi|_1} \log \left( \frac{|\phi_{i,j}|_1}{|\Phi|_1} \right)$$

- The total training objective is the prediction loss plus L1 entropy regularization of all the KAN layers:

$$\ell_{\text{total}} = \ell_{\text{pred}} + \lambda \left( \mu_1 \sum_{l=0}^{L-1} |\Phi_l|_1 + \mu_2 \sum_{l=0}^{L-1} S(\Phi_l) \right)$$

# Pruning

- After training with sparsification penalty, we may want to prune the network
- KANs are pruned on the node level (rather than edge level)

For each node (say the  $i^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer, the incoming and outgoing score are defined as:

$$I_{l,i} = \max_k (|\phi_{l-1,i,k}|_1), \quad O_{l,i} = \max_j (|\phi_{l+1,j,i}|_1)$$

A node is considered to be important if both incoming and outgoing scores are greater than a hyperparameter  $\theta = 10^{-2}$  by default

# Symbolic Regression Process

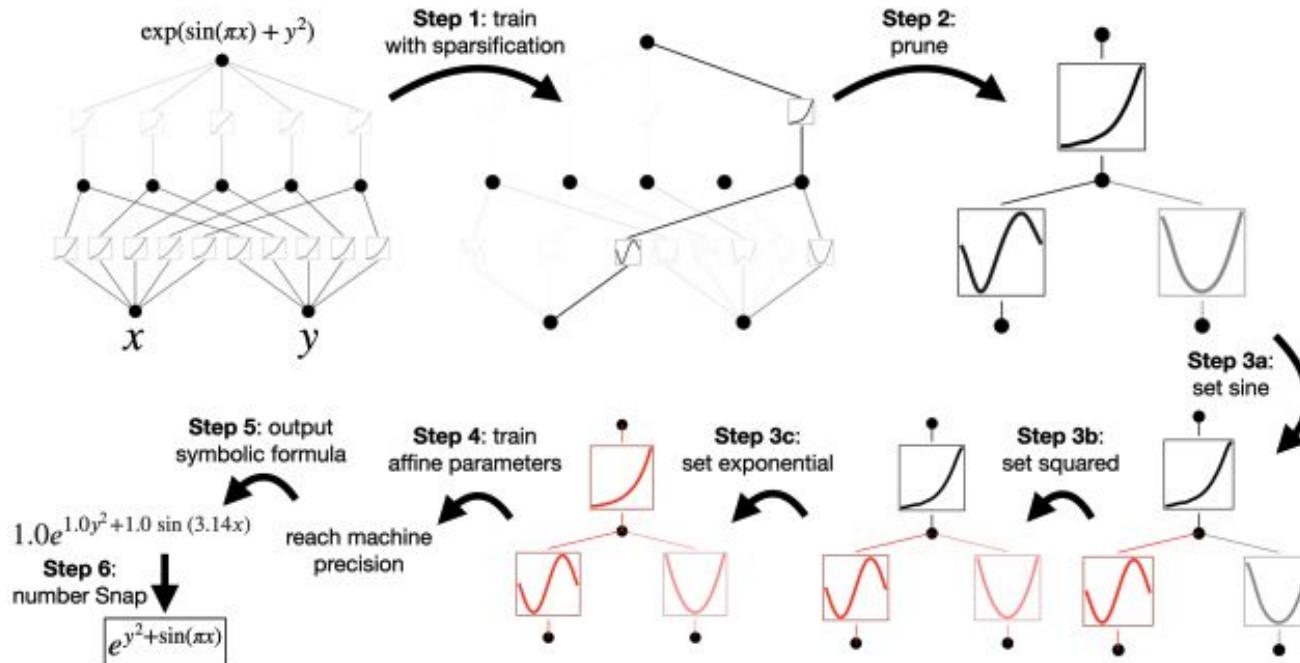
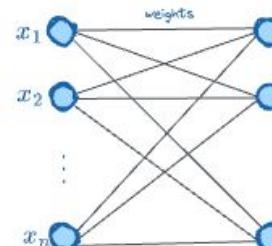
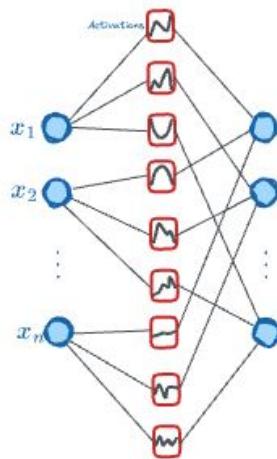


Figure 2.4: An example of how to do symbolic regression with KAN.

# KAN vs MLP parameter count

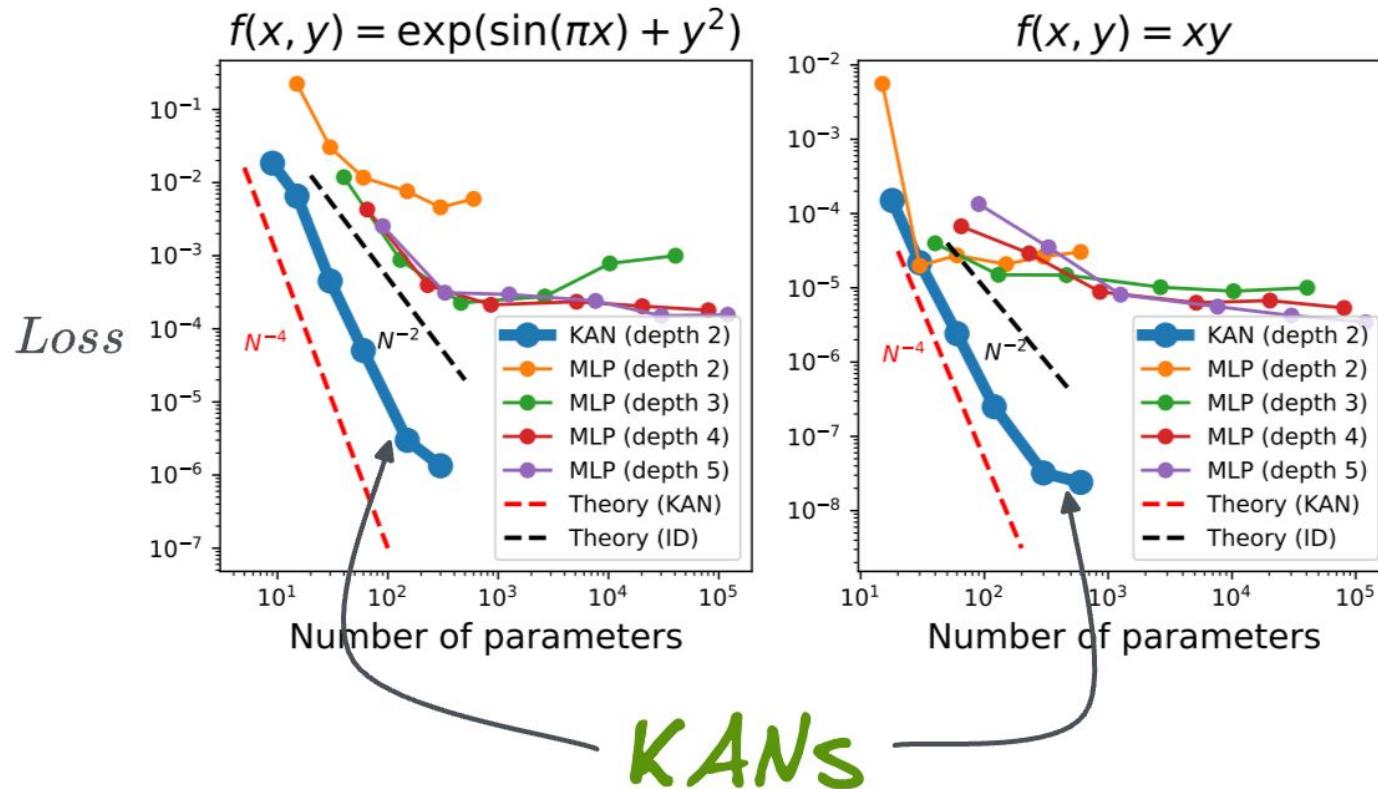
*KAN Layer*      *MLP Layer*



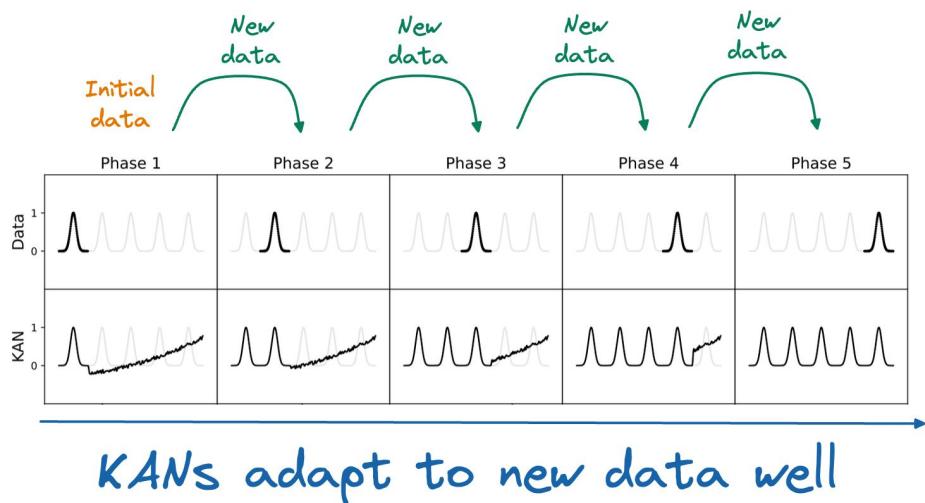
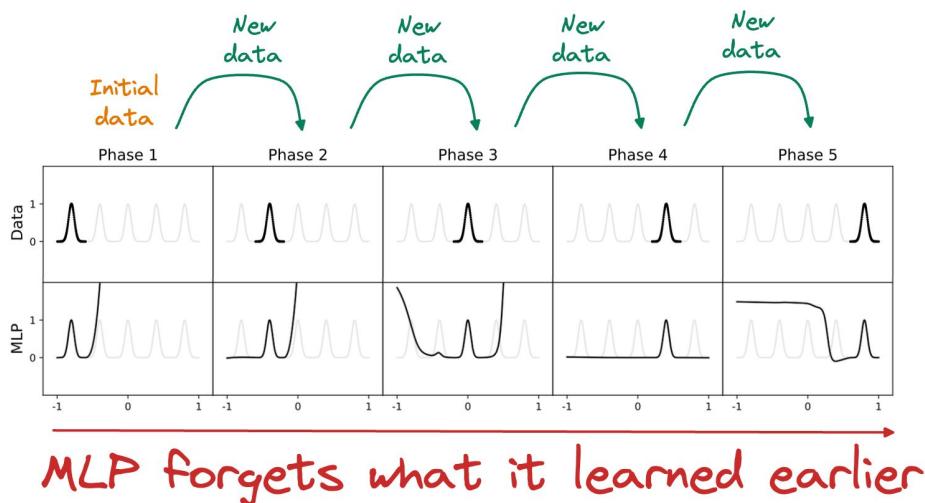
$$\sim N^2 L(G + k)$$

$$N^2 L$$

# KANs are accurate with fewer parameters



# Continual Learning



# Interpretability

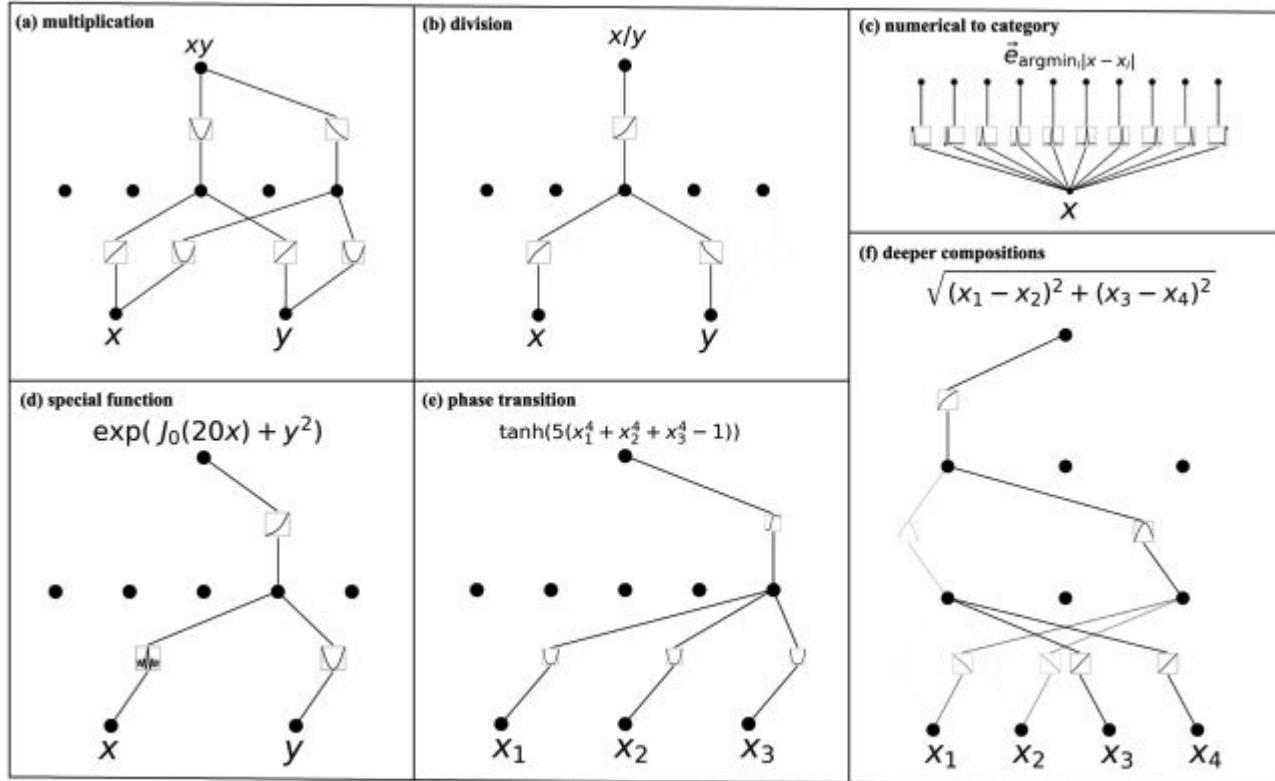
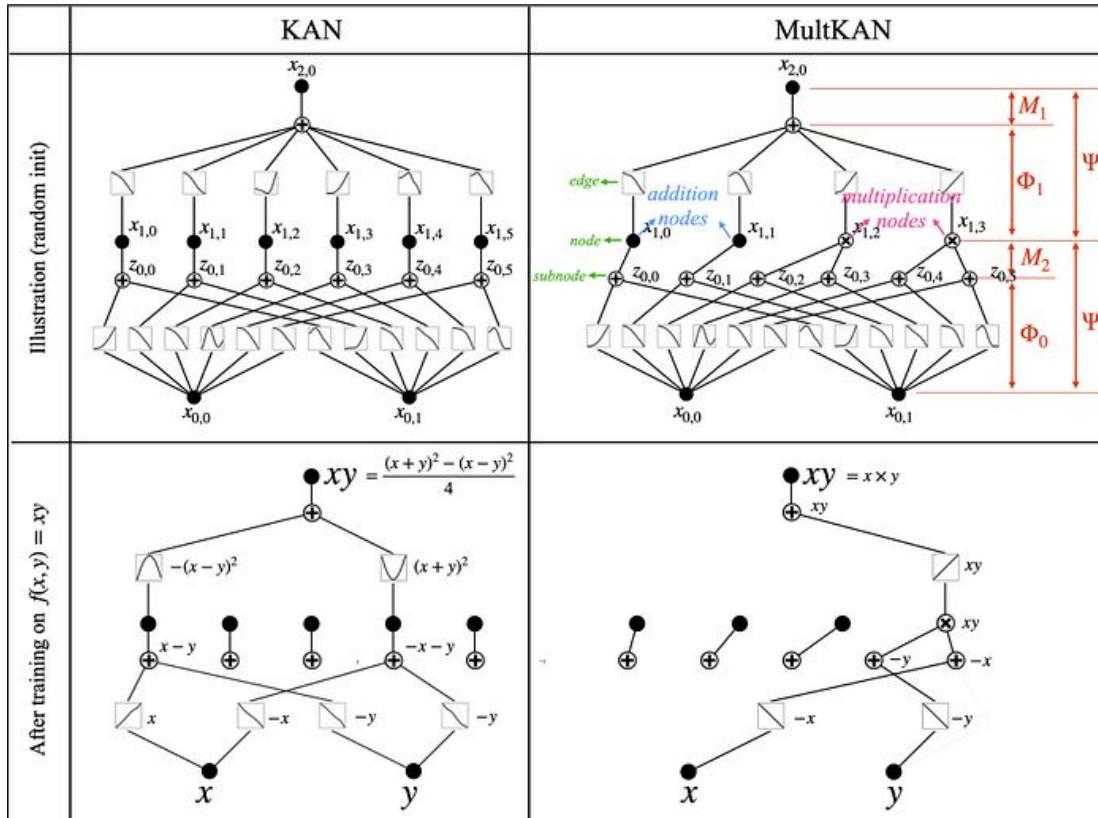


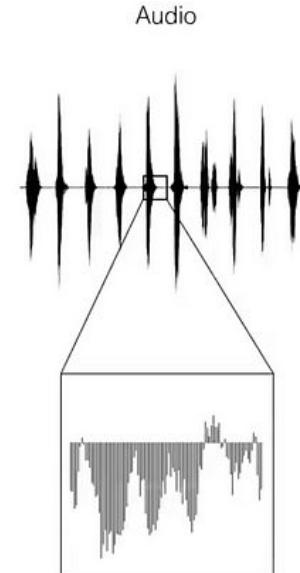
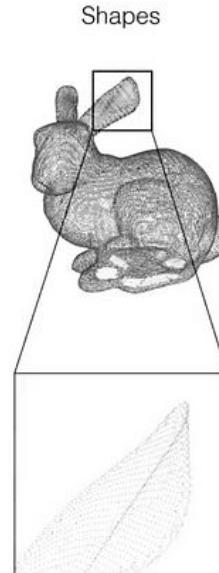
Figure 4.1: KANs are interpretable for simple symbolic tasks

# MultKAN



# Implicit Neural Representations (INRs)

- Also referred to as coordinate-based representations or neural fields, INRs are a way to parameterize signals of all kinds



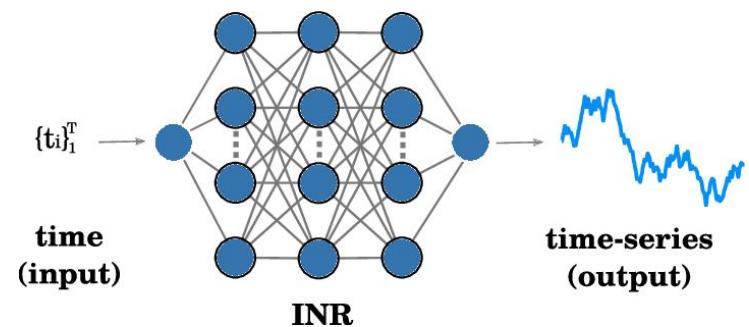
# Implicit Neural Representations

Coordinate-based networks model signals as continuous mappings:

$$f : \mathbb{R}^d \rightarrow \mathbb{R}^m$$

where coordinates  $x \in \mathbb{R}^d$

map to attributes  $y \in \mathbb{R}^m$



# Differentiable Signal Processing

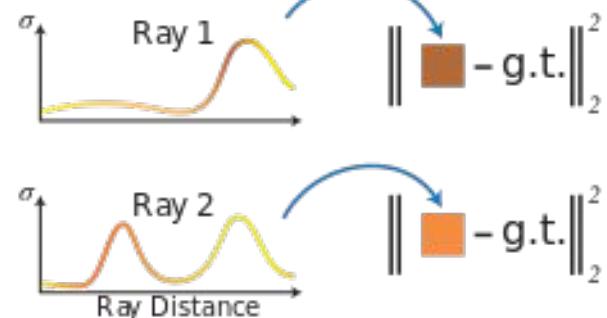
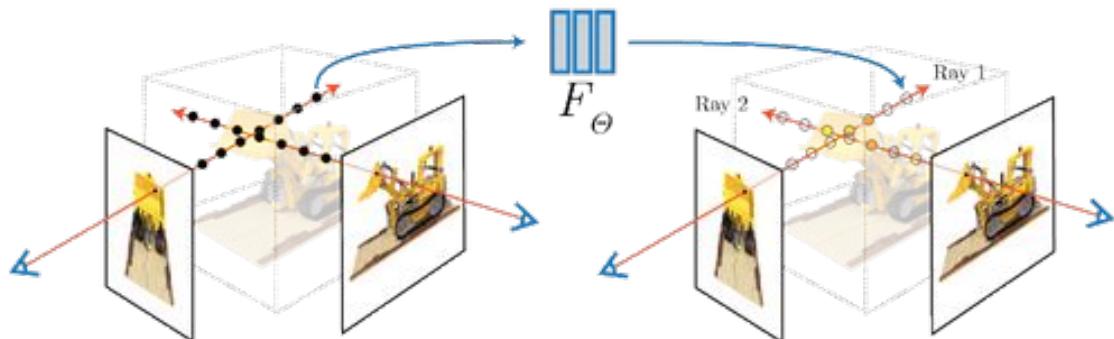
- Unlike discrete representations (grids or meshes) INRs learn a continuous mapping, this means every point in the domain has a derivative
- Many objective functions depend on derivatives - continuous differentiability means gradient-based optimizations can find optimal solutions not just approximations



Solving the wave equation

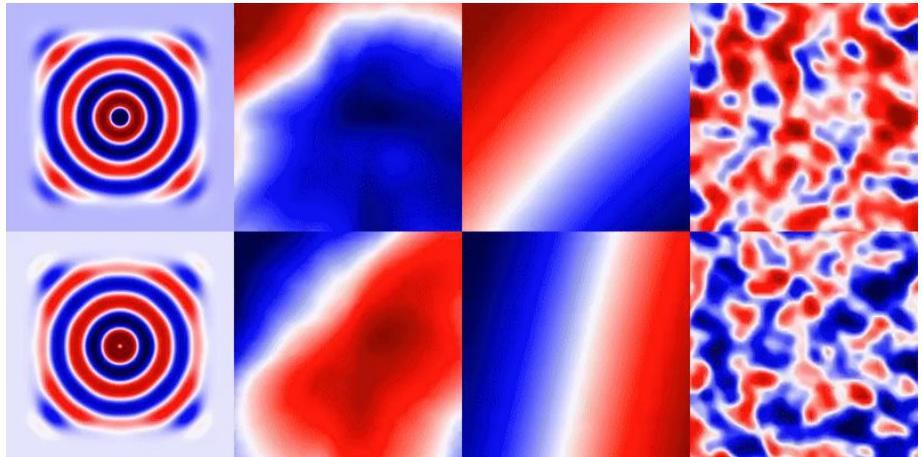
# Neural Radiance Fields (NeRF)

- NeRF uses a fully-connected deep network
  - Input is a single 5D coordinate: spatial location  $(x,y,z)$ , and viewing direction  $(\theta, \varphi)$
  - Output is the volume density and view-dependent emitted radiance at that spatial location



# Physics-informed applications

- Compact neural encoding allows rapid evaluation across different physical parameters without full resimulation
- Continuous mappings allow complex physical phenomena to be stored compactly without discretization
- The differentiable nature of INRs allows exact computation of derivatives, improving PDE solution accuracy



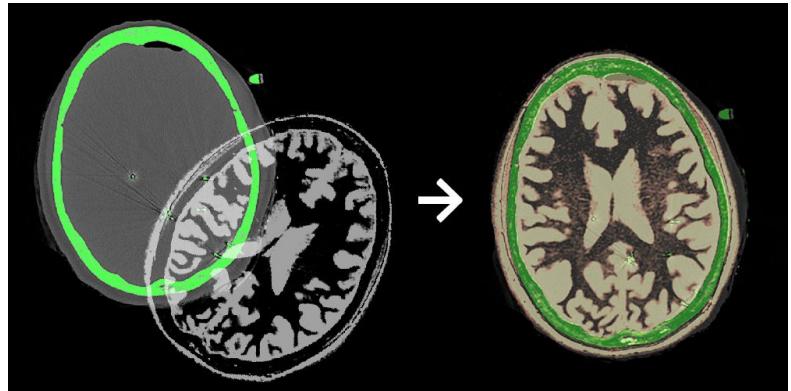
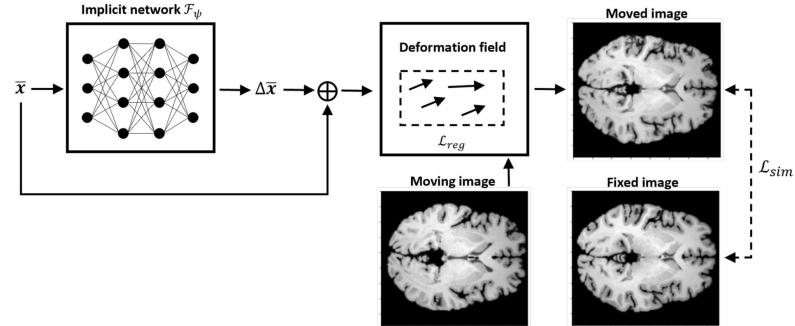
$$H(m)\Phi(x) = -f(x), \text{ with } H(m) = (\Delta + m(x)w^2)$$

# Medical Applications

If you have MRIs undersampled in the k-space (for an accelerated MRI scan) an INR can be used to reconstruct it with missing data

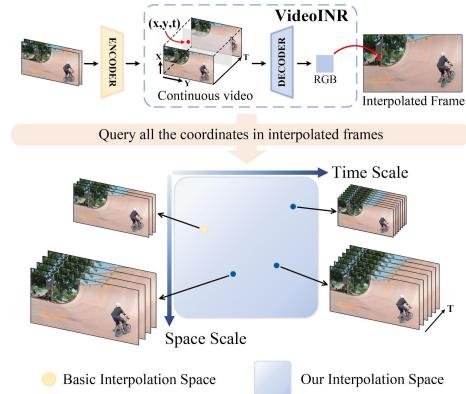
For image registration, INRs naturally model continuous deformation fields between images, producing mappings that preserve topology while allowing computation of exact spatial derivatives for optimization

INRs support arbitrary slice orientations and zoom levels after training, making it well-suited for interactive medical visualization and analysis



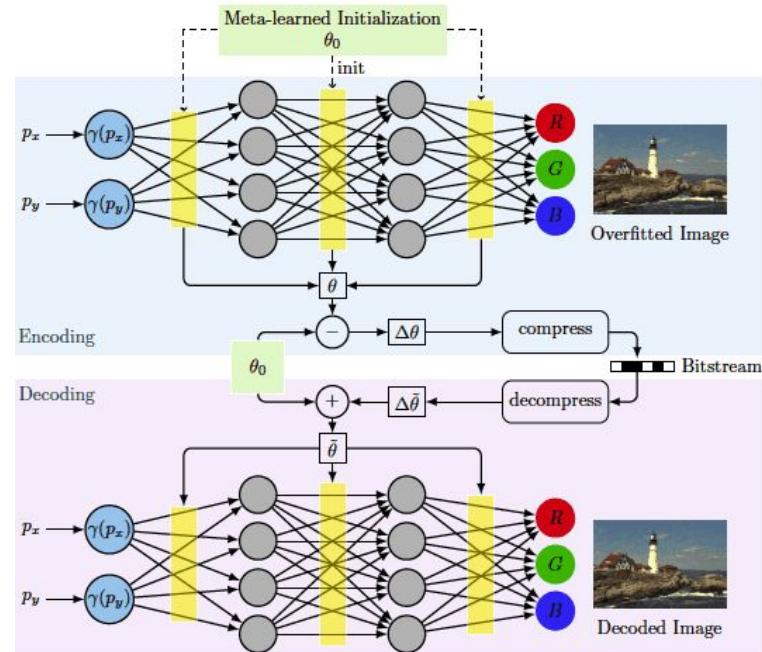
# Resolution Independence (VideoINR)

- $(x,y,t)$  are inputs to the model
- Model learns the resulting frame
- Allows for upsampling in resolution and framerate



# Memory efficiency

- INRs achieve high compression ratios by encoding complex signals through compact networks
- You can use the same network structure to learn signals in different domains
- INRs preserve sharp features and do not introduce blocking artifacts or ringing



# INR Challenges

**Spectral bias** - neural networks prioritize learning the low frequency modes first

- Even if a high frequency component has a larger amplitude in the target signal, the network will still learn to approximate lower frequency components that are less prominent first

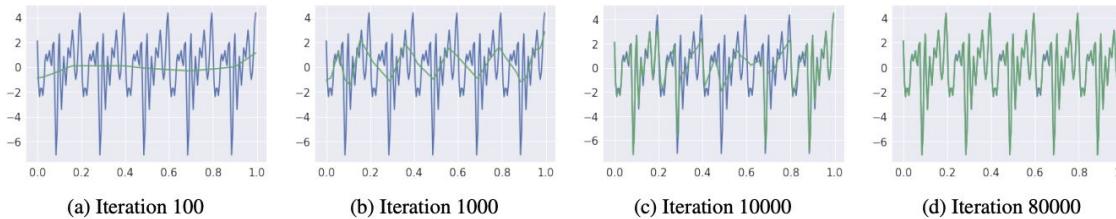
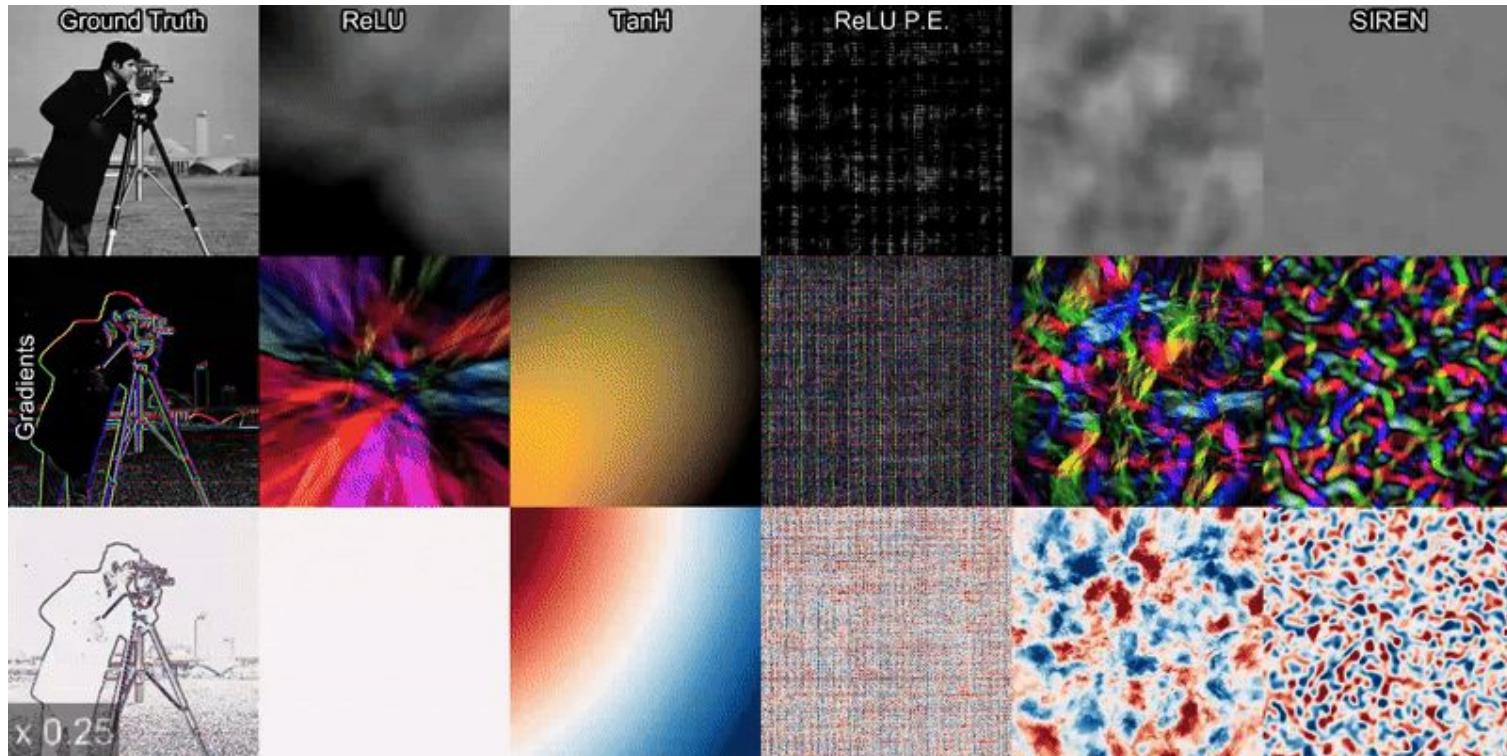


Figure 2. The learnt function (green) overlayed on the target function (blue) as the training progresses. The target function is a superposition of sinusoids of frequencies  $\kappa = (5, 10, \dots, 45, 50)$ , equal amplitudes and randomly sampled phases.

**Inference** requires a separate network evaluation for each point, making sampling of the domain expensive vs. grid based lookups

INRs also suffer from the **curse of dimensionality** as the required number of training samples grows exponentially with input dimensions

# Sinusoidal Representation Networks (SIREN)



# Periodic Activations for INRs

$$F(\mathbf{x}, \Phi, \nabla_{\mathbf{x}}\Phi, \nabla_{\mathbf{x}}^2\Phi, \dots) = 0, \quad \Phi : \mathbf{x} \mapsto \Phi(\mathbf{x})$$

This formulation takes as input the spatial or spatio-temporal coordinates  $x \in \mathbb{R}^m$  and, optionally, derivatives of  $\Phi$  with respect to these coordinates.

$$\Phi(\mathbf{x}) = \mathbf{W}_n (\phi_{n-1} \circ \phi_{n-2} \circ \dots \circ \phi_0)(\mathbf{x}) + \mathbf{b}_n,$$

$$\mathbf{x}_i \mapsto \phi_i(\mathbf{x}_i) = \sin(\mathbf{W}_i \mathbf{x}_i + \mathbf{b}_i).$$

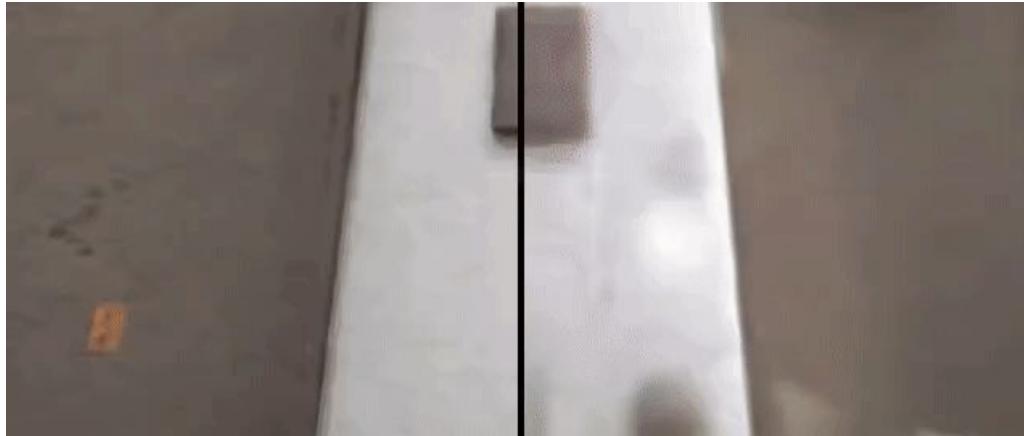
$\phi_i : \mathbb{R}^{M_i} \rightarrow \mathbb{R}^{N_i}$  is the  $i^{th}$  layer of the network

# Simple example: fitting an image

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3, x \rightarrow \Phi(x)$$

dataset  $D = \{(x_i, f(x_i))\}_i$  of pixel coordinates  $x_i = (x_i, y_i)$  with RGB colors  $f(x_i)$

$$\mathcal{L} = \sum_i \|\Phi(x_i) - f(x_i)\|^2$$



# SIREN initialization scheme

The key idea in the initialization scheme is to preserve the distribution of activations through the network so that the final output does not depend on the number of layers

- The authors propose to draw weights with:

$$c = 6 \text{ so that } w_i \sim \mathcal{U} \left( -\sqrt{6/n}, \sqrt{6/n} \right)$$

This is because when each component of weights a layer  $\mathbf{w}$  is uniformly distributed such as

$$w_i \sim \mathcal{U} \left( -c/\sqrt{n}, c/\sqrt{n} \right)$$

The dot product of converges to the normal distribution  $\mathbf{w}^T \mathbf{x} \sim \mathcal{N} (0, c^2/6)$  this ensures the input to each sine activation is normally distributed with a standard deviation of 1

Finally, the first sine layer needs to be initialized with weights so that the sin function  $(\omega_0 \cdot \mathbf{W}\mathbf{x} + b)$  spans multiple periods over  $[-1, 1]$

# WIRE: Wavelet Implicit Neural Representations

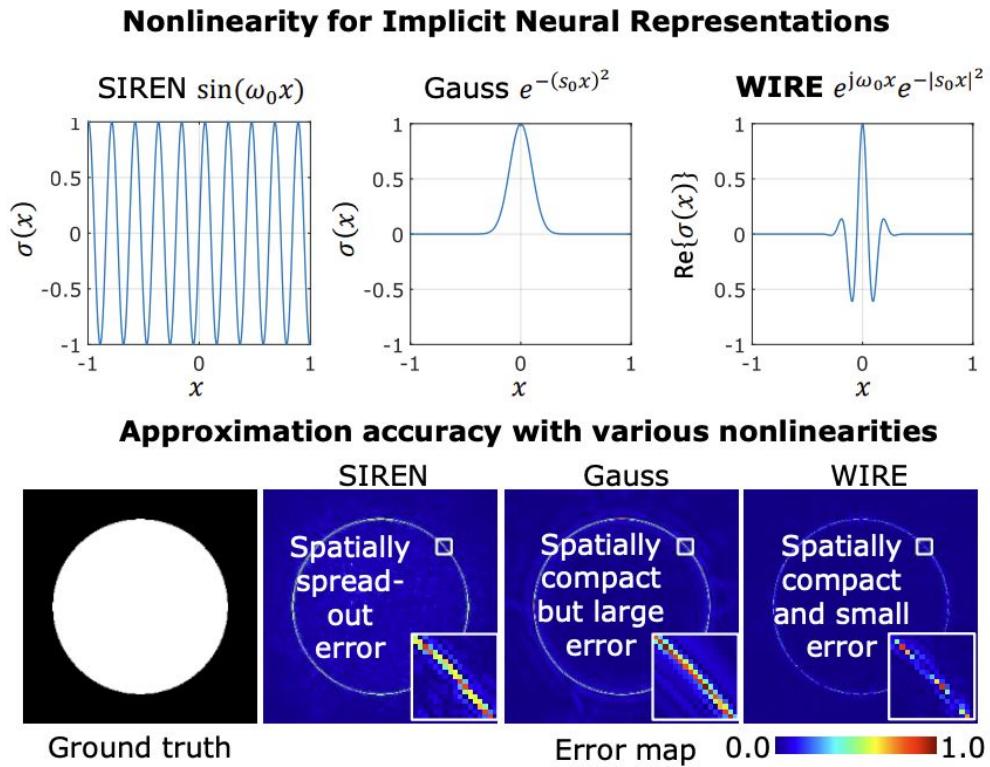
WIRE uses the continuous complex Gabor wavelet  $\psi$  for its activation:

$$\sigma(x) = \psi(x; \omega_0, s_0) = e^{j\omega_0 x} e^{-|s_0 x|^2}$$

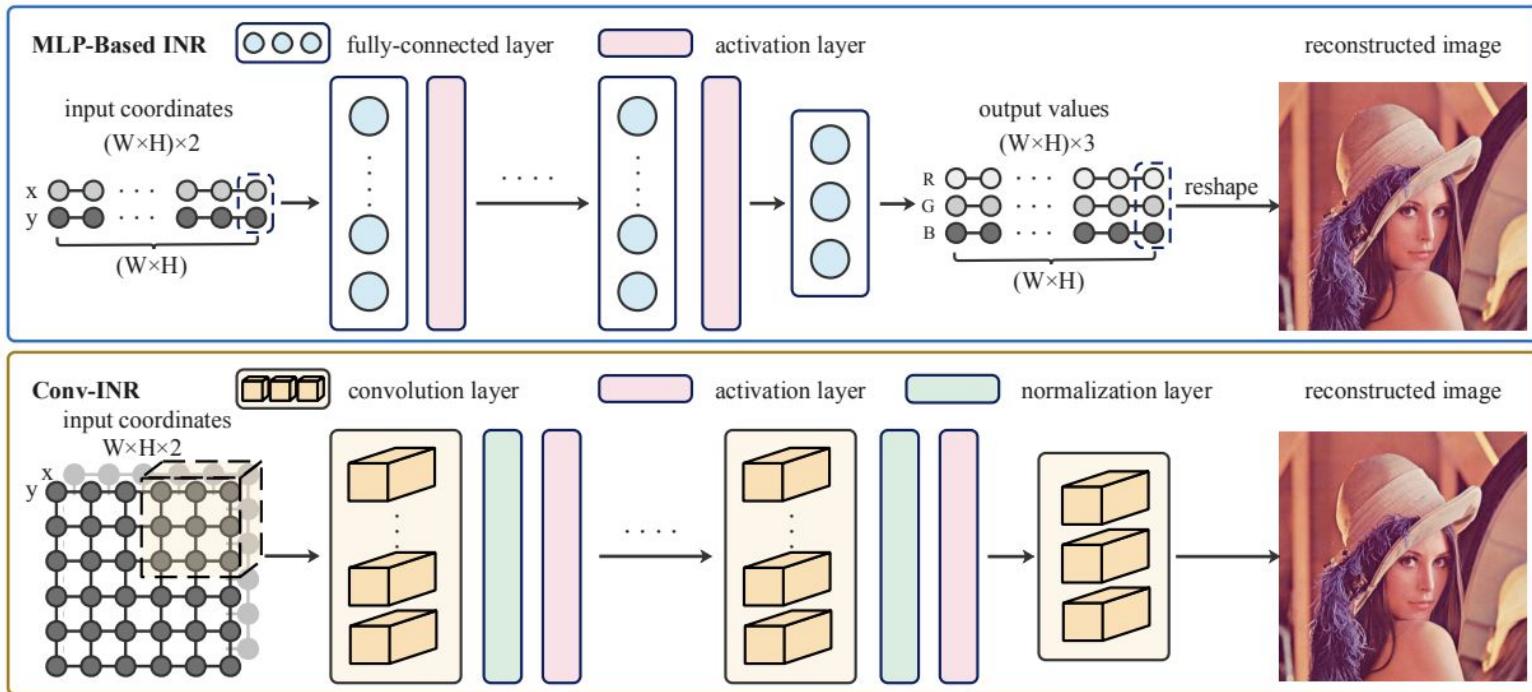
$\omega_0$  controls the frequency  
 $s_0$  controls the width

The first layer activations have the form:

$$y_1 = \psi(W_1 x + b_1; \omega_0, s_0)$$



# Conv-INR



**Fig. 1:** Pipelines of MLP-based INR and Conv-INR. Take the image fitting task as the example.