

# NYU FRE 7773 - Week 7

---

*Machine Learning in Financial Engineering*

Ethan Rosenthal

# Intro to Deep Learning

---

*Machine Learning in Financial Engineering*

Ethan Rosenthal

# Artificial Neural Networks

---

# Feed Forward Network / Multilayer Perceptron

Data Input  
(x's)

Square  
Footage



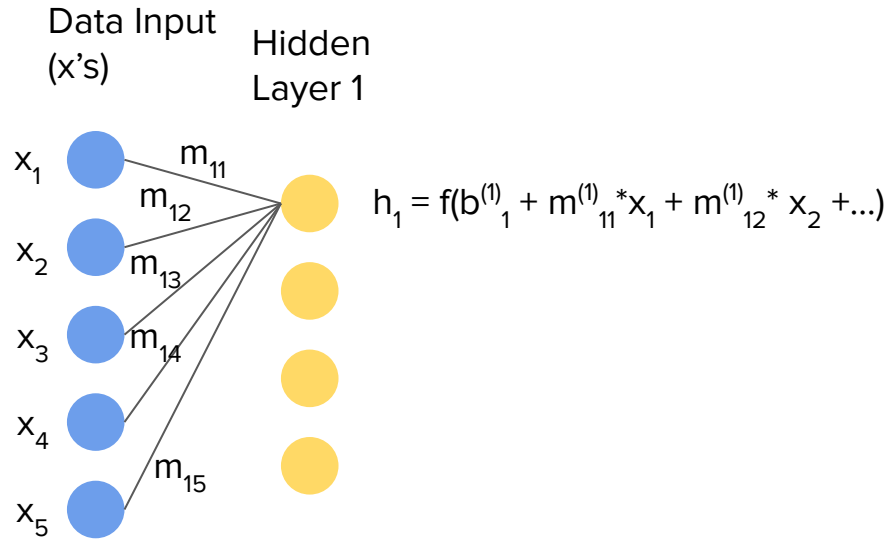
Distance



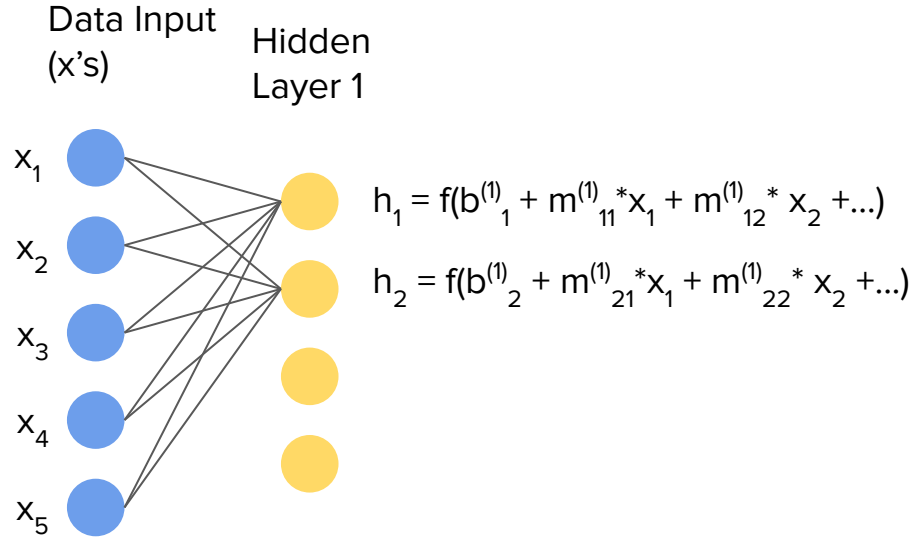
Bedrooms



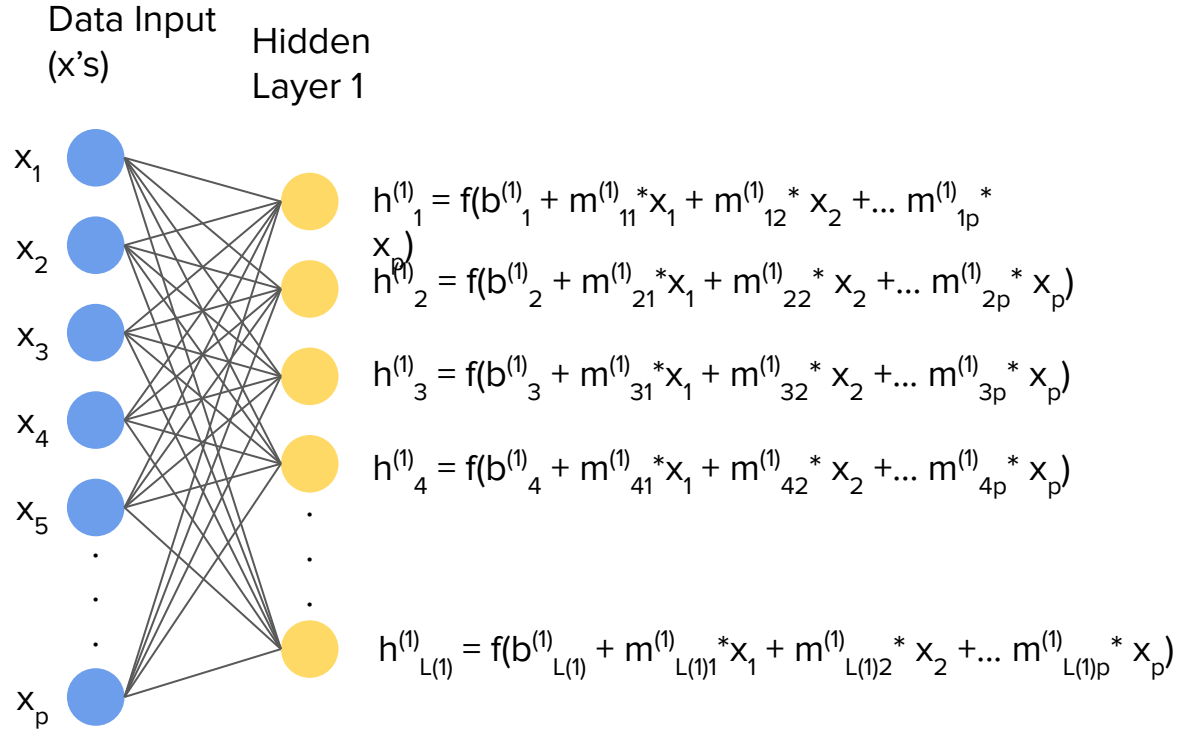
# Feed Forward Network / Multilayer Perceptron



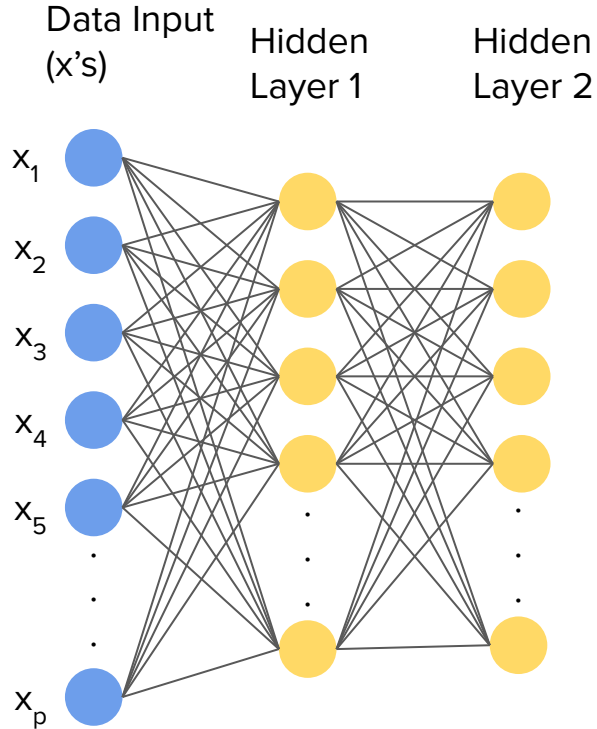
# Feed Forward Network / Multilayer Perceptron



# Feed Forward Network / Multilayer Perceptron



# Feed Forward Network / Multilayer Perceptron



$$h_1^{(2)} = f(b_1^{(2)} + m_{11}^{(2)} * h_1^{(1)} + m_{12}^{(2)} * h_2^{(1)} + \dots + m_{1L(1)}^{(2)} * h_{L(1)}^{(1)})$$

$$h_2^{(2)} = f(b_2^{(2)} + m_{21}^{(2)} * h_1^{(1)} + m_{22}^{(2)} * h_2^{(1)} + \dots + m_{2L(1)}^{(2)} * h_{L(1)}^{(1)})$$

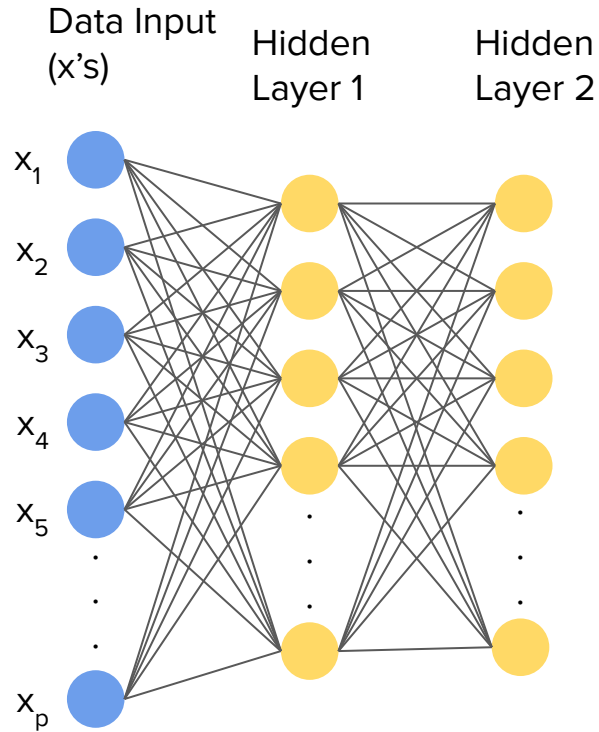
$$h_3^{(2)} = f(b_3^{(2)} + m_{31}^{(2)} * h_1^{(1)} + m_{32}^{(2)} * h_2^{(1)} + \dots + m_{3L(1)}^{(2)} * h_{L(1)}^{(1)})$$

$$h_4^{(2)} = f(b_4^{(2)} + m_{41}^{(2)} * h_1^{(1)} + m_{42}^{(2)} * h_2^{(1)} + \dots + m_{4L(1)}^{(2)} * h_{L(1)}^{(1)})$$

$$h_{L(2)}^{(2)} = f(b_{L(2)}^{(2)} + m_{L(2)1}^{(2)} * h_1^{(1)} + m_{L(2)2}^{(2)} * h_2^{(1)} + \dots + m_{L(2)L(1)}^{(2)} * h_{L(1)}^{(1)})$$

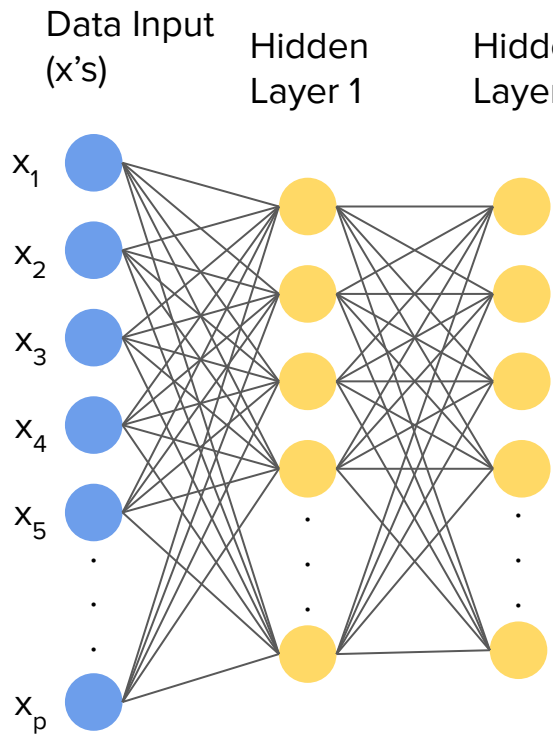


# Feed Forward Network / Multilayer Perceptron



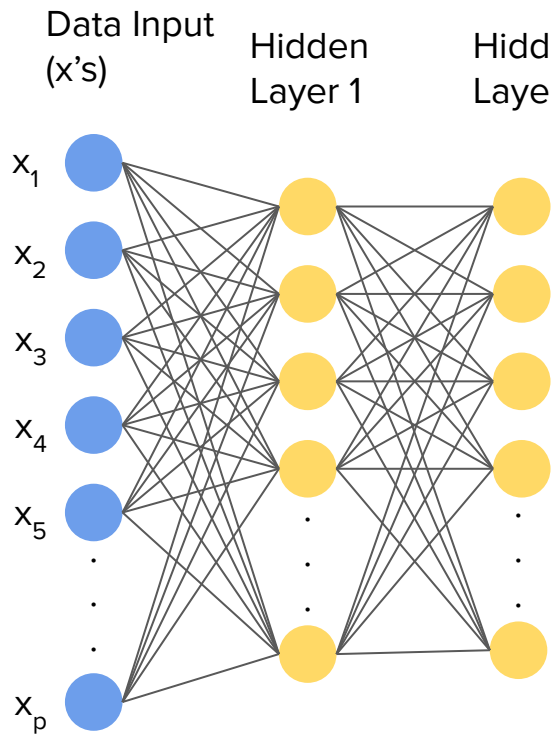
$$h^{(2)}_1 = f(b^{(2)}_1 + m^{(2)}_{11} * h^{(1)}_1 + m^{(2)}_{12} * h^{(1)}_2 + \dots + m^{(2)}_{1L(1)} * h^{(1)}_{L(1)})$$

# Feed Forward Network / Multilayer Perceptron



$$\begin{aligned}
 h^{(2)}_1 &= f(b^{(2)}_1 + m^{(2)}_{11} * h^{(1)}_1 + m^{(2)}_{12} * h^{(1)}_2 + \dots m^{(2)}_{1L(1)} * h^{(1)}_{L(1)}) \\
 &= f(b^{(2)}_1 + \\
 &\quad m^{(2)}_{11} * f(b^{(1)}_1 + m^{(1)}_{11} * x_1 + m^{(1)}_{12} * x_2 + \dots m^{(1)}_{1p} * x_p) \\
 &\quad + m^{(2)}_{12} * f(b^{(1)}_2 + m^{(1)}_{21} * x_1 + m^{(1)}_{22} * x_2 + \dots m^{(1)}_{2p} * x_p) \\
 &\quad + \dots \\
 &\quad m^{(2)}_{1L(1)} * f(b^{(1)}_{L(1)} + m^{(1)}_{L(1)1} * x_1 + m^{(1)}_{L(1)2} * x_2 + \dots m^{(1)}_{L(1)p} * x_p) \\
 &\quad )
 \end{aligned}$$

# Feed Forward Network / Multilayer Perceptron

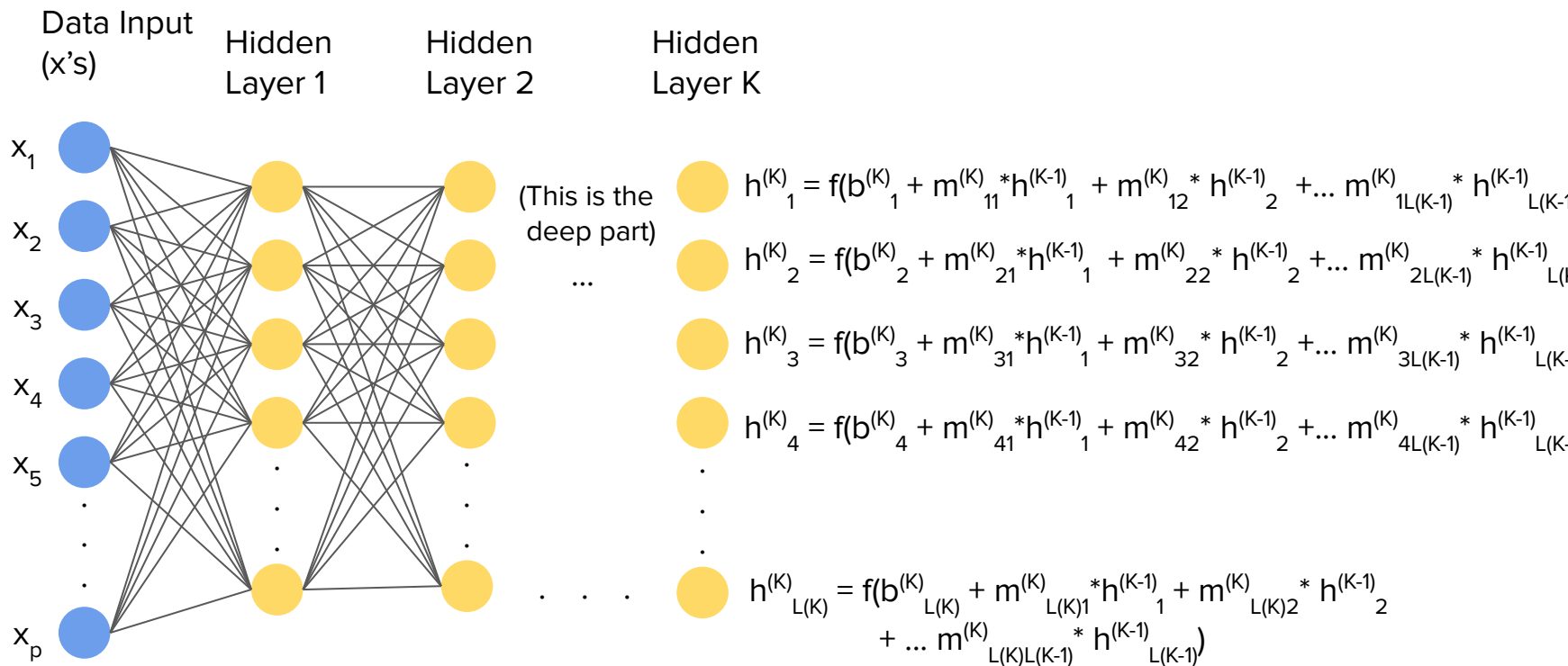


$$\begin{aligned}
 h^{(2)}_1 &= f(b^{(2)}_1 + m^{(2)}_{11} * h^{(1)}_1 + m^{(2)}_{12} * h^{(1)}_2 + \dots m^{(2)}_{1L(1)} * h^{(1)}_{L(1)}) \\
 &= f(b^{(2)}_1 + \\
 &\quad m^{(2)}_{11} * f(b^{(1)}_1 + m^{(1)}_{11} * x_1 + m^{(1)}_{12} * x_2 + \dots m^{(1)}_{1p} * x_p) \\
 &\quad + m^{(2)}_{12} * f(b^{(1)}_2 + m^{(1)}_{21} * x_1 + m^{(1)}_{22} * x_2 + \dots m^{(1)}_{2p} * x_p) \\
 &\quad + \dots \\
 &\quad m^{(2)}_{1L(1)} * f(b^{(1)}_{L(1)} + m^{(1)}_{L(1)1} * x_1 + m^{(1)}_{L(1)2} * x_2 + \dots m^{(1)}_{L(1)p} * x_p) \\
 &\quad )
 \end{aligned}$$

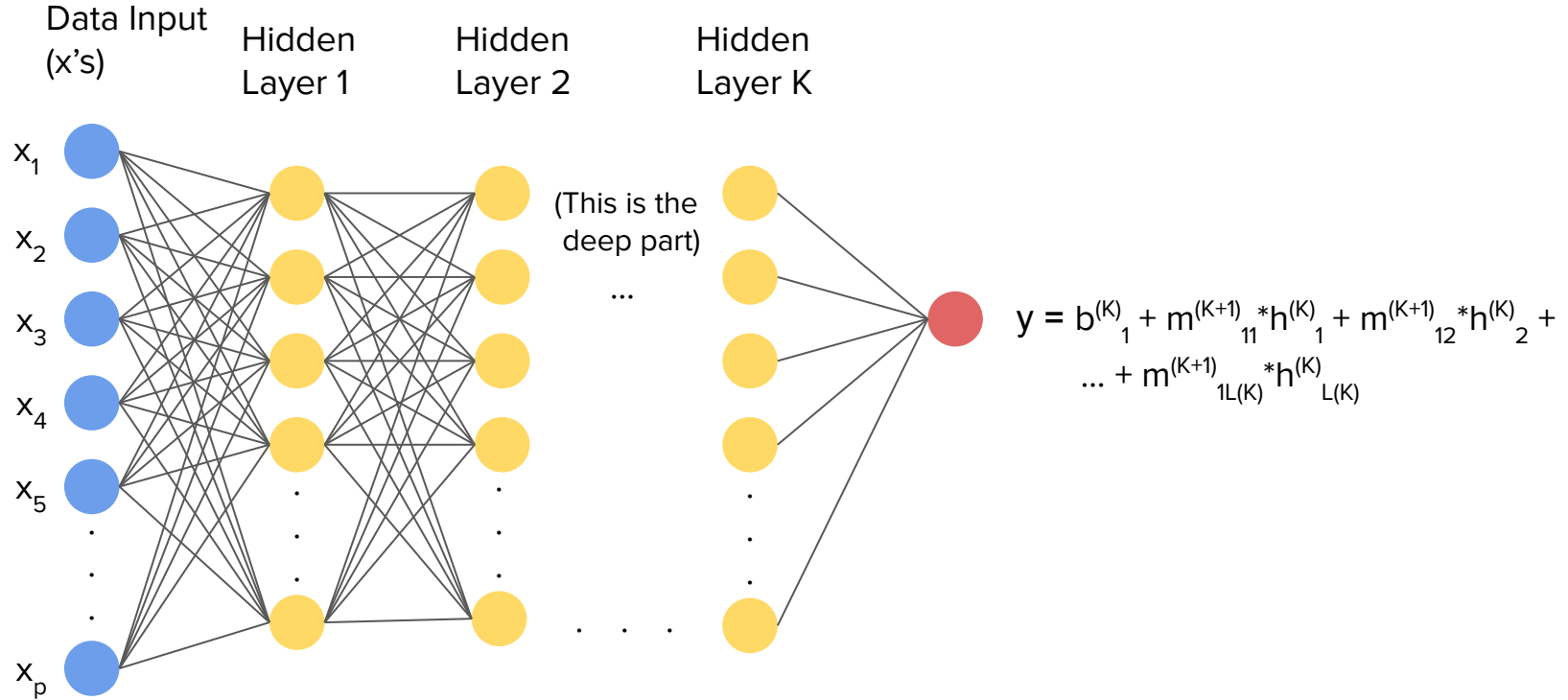
$$h^{(2)}_1 = f(b^{(2)}_1; m^{(2)}_1; h^{(1)}) \quad , \quad h^{(1)}_1 = f(b^{(1)}_1; m^{(1)}_1; x)$$

$$h^{(2)}_1 = f(b^{(2)}_1; m^{(2)}_1; f(b^{(1)}_1; m^{(1)}_1; x))$$

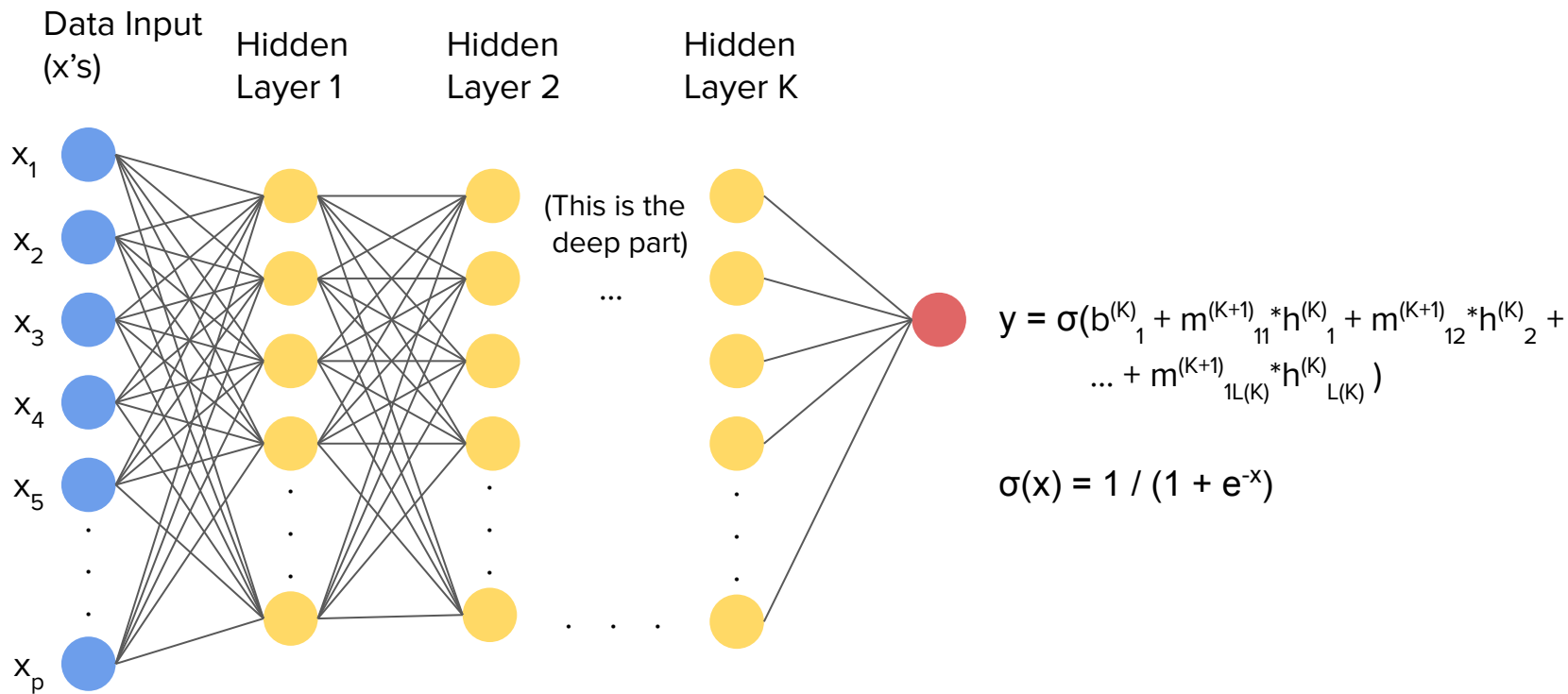
# Feed Forward Network / Multilayer Perceptron



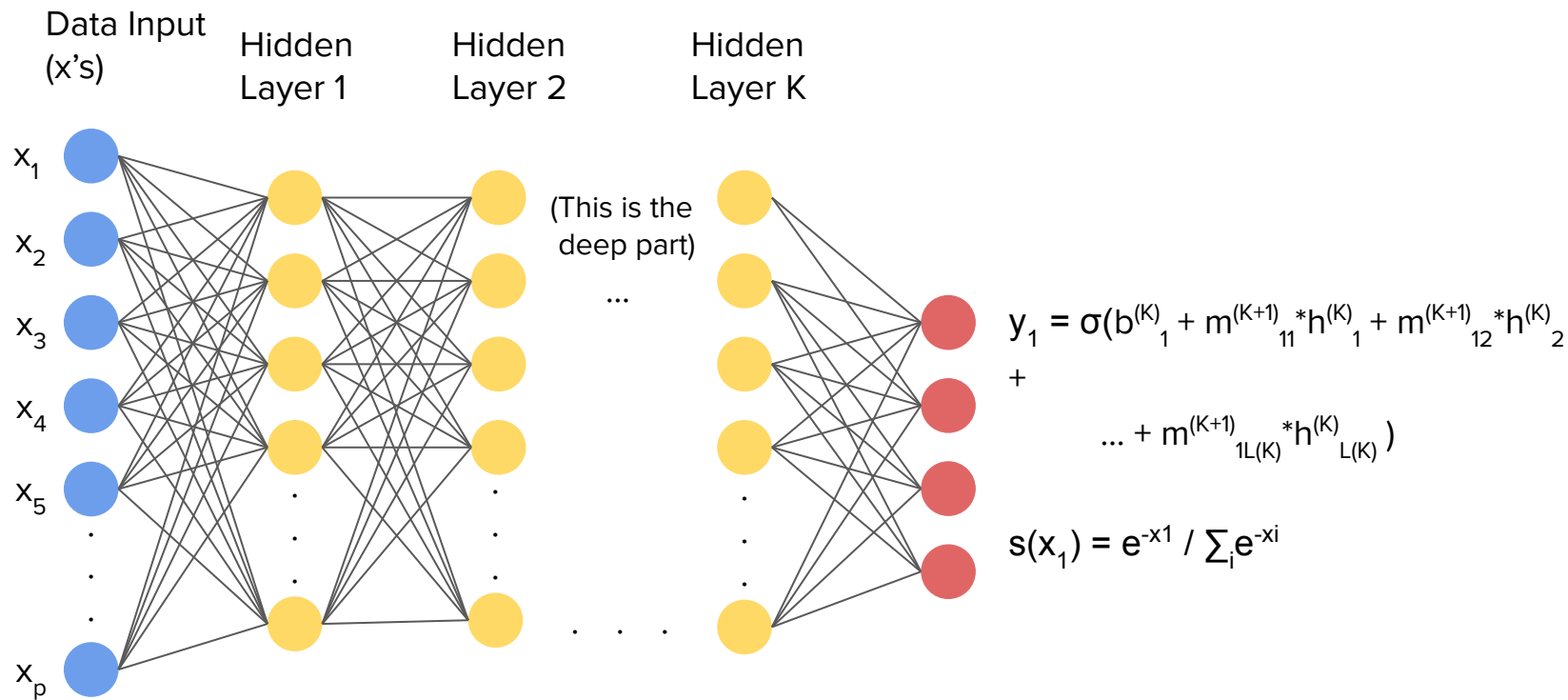
# Feed Forward Network / Multilayer Perceptron



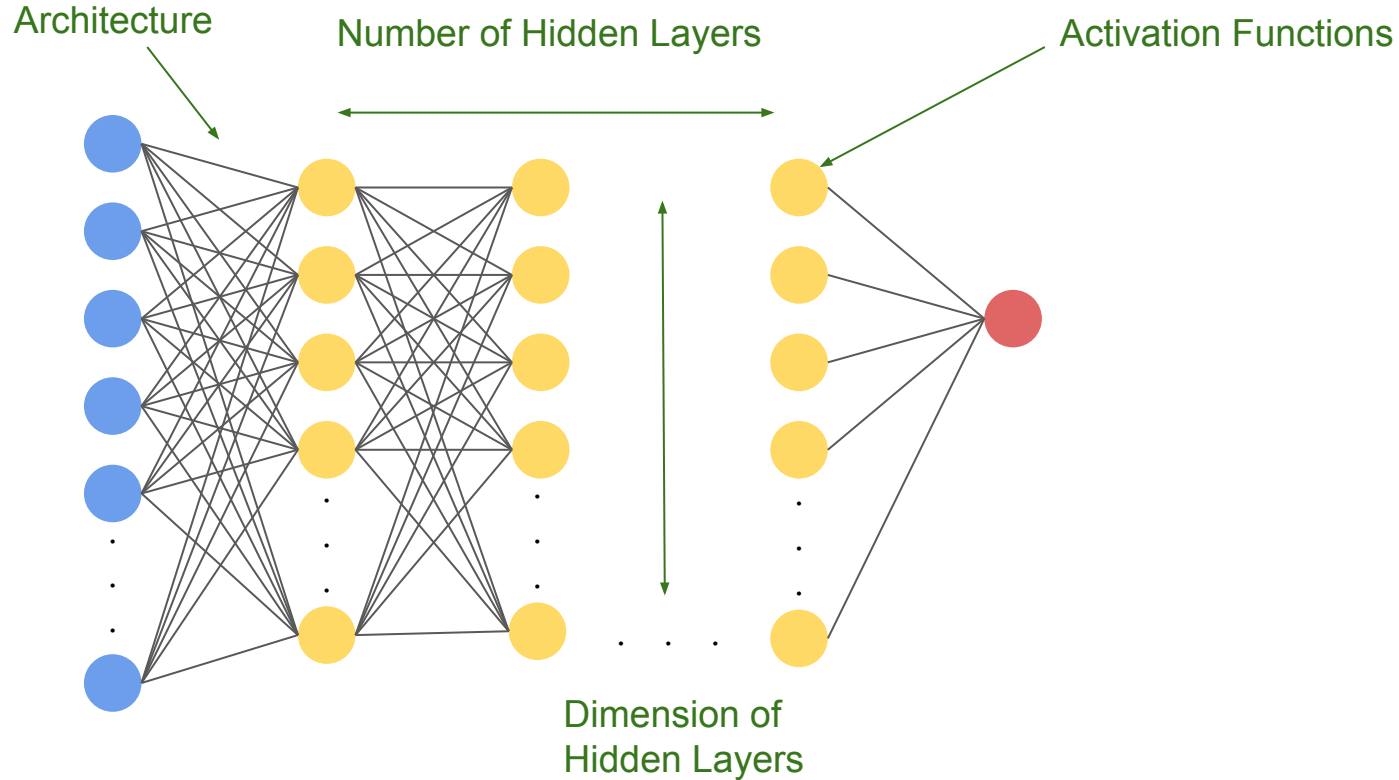
# Feed Forward Network – Binary Classification



# Feed Forward Network – Multiclass Classification

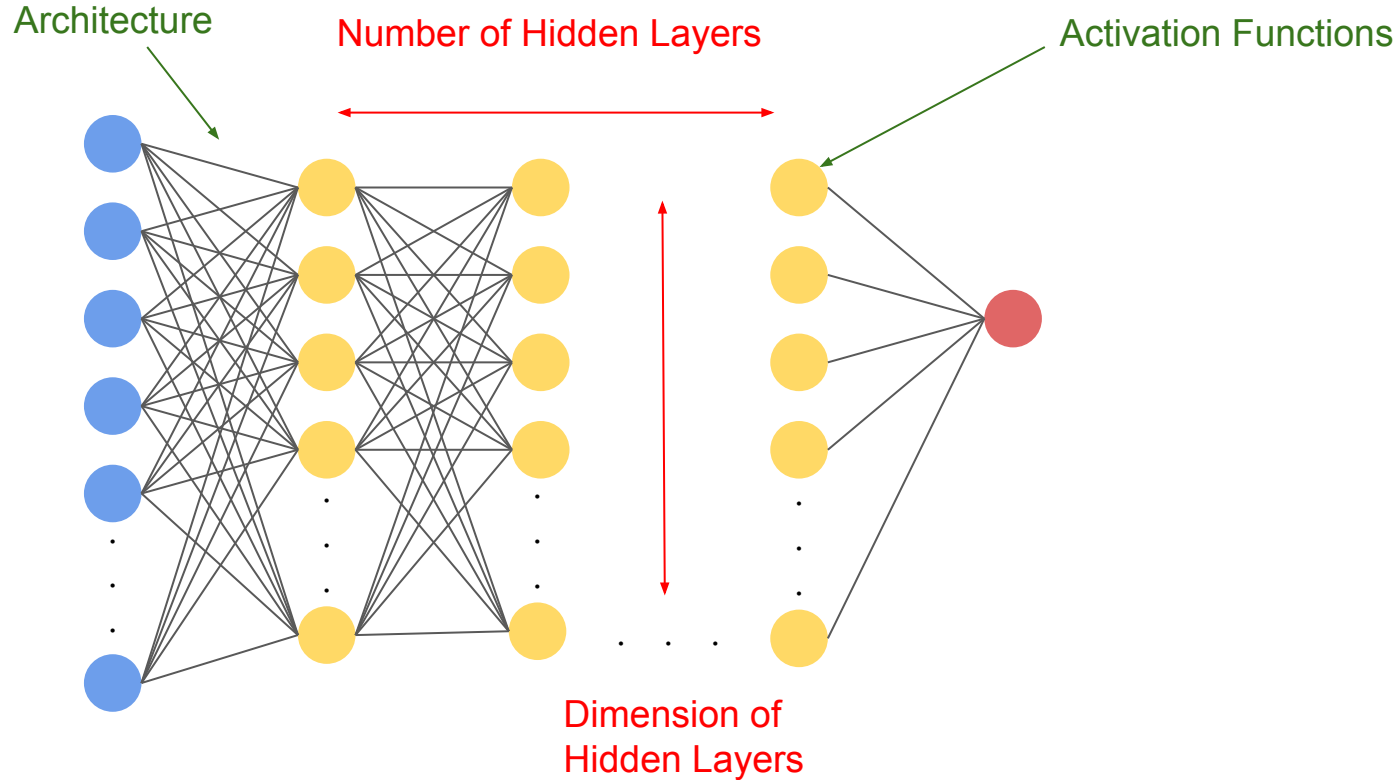


# Knobs at our disposal

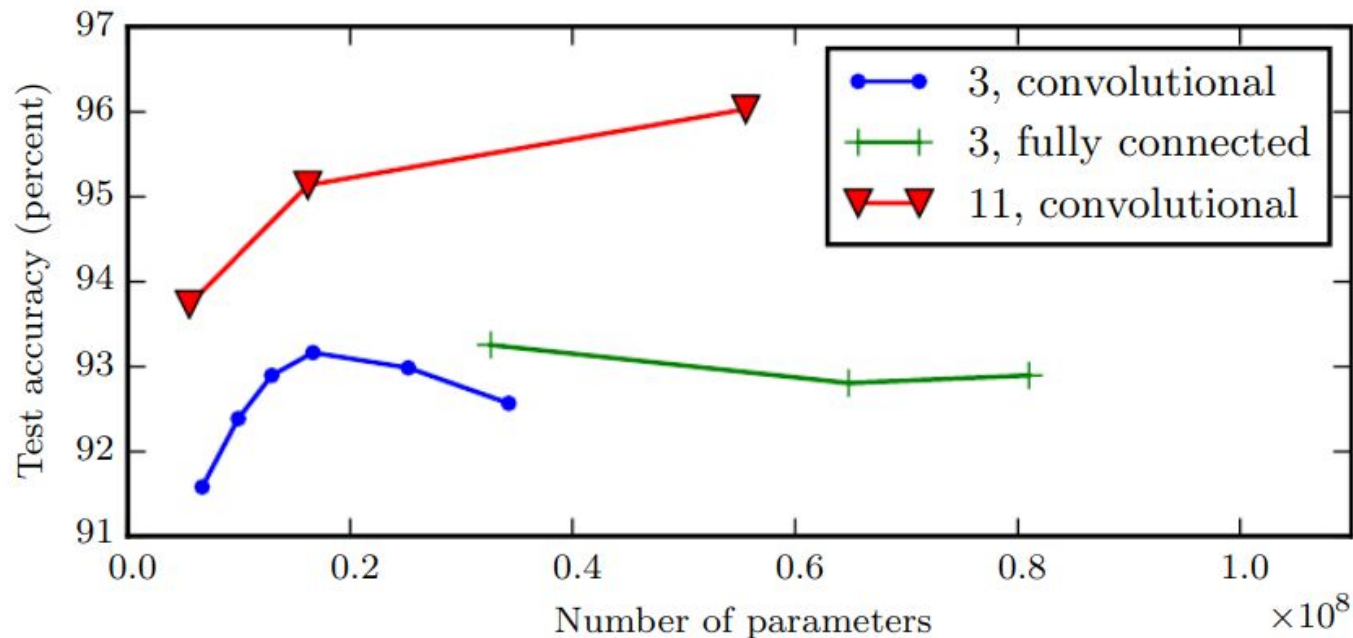




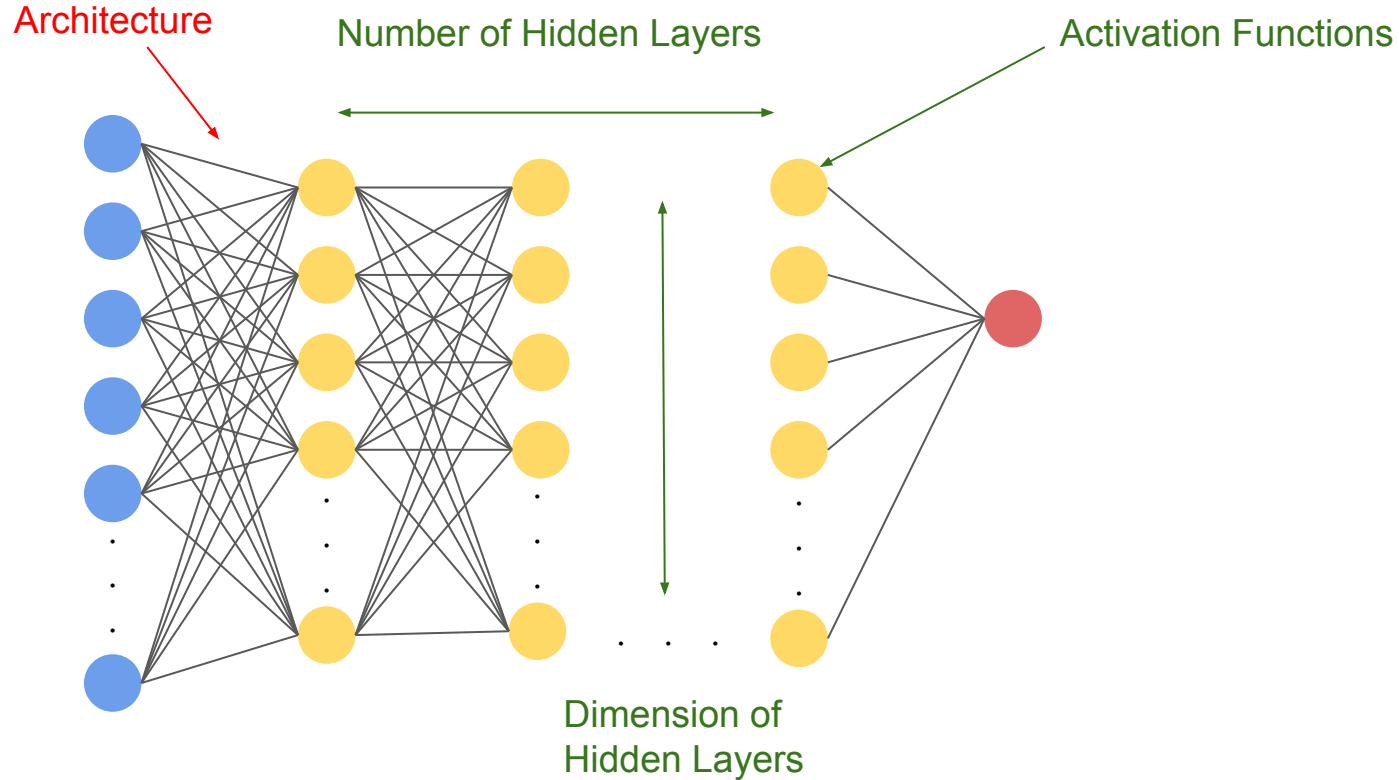
# Knobs at our disposal



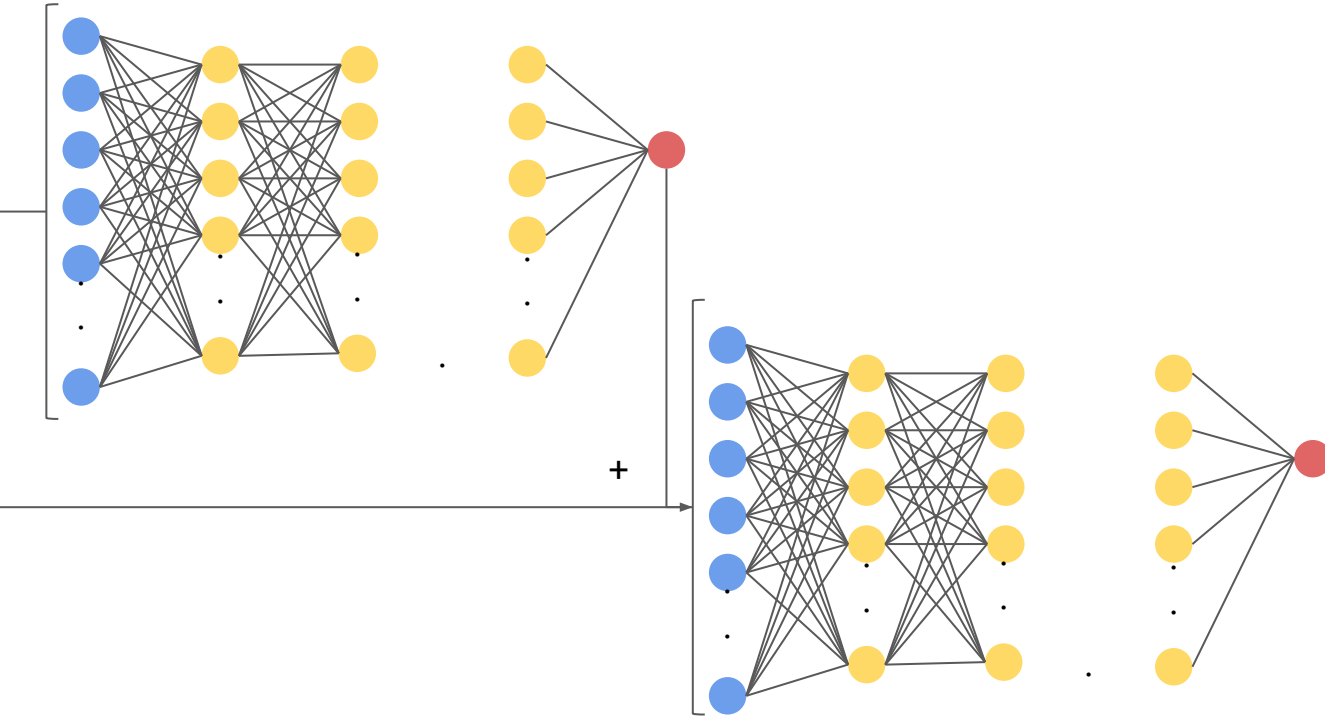
# Knobs at our disposal - Depth vs. Width



# Knobs at our disposal

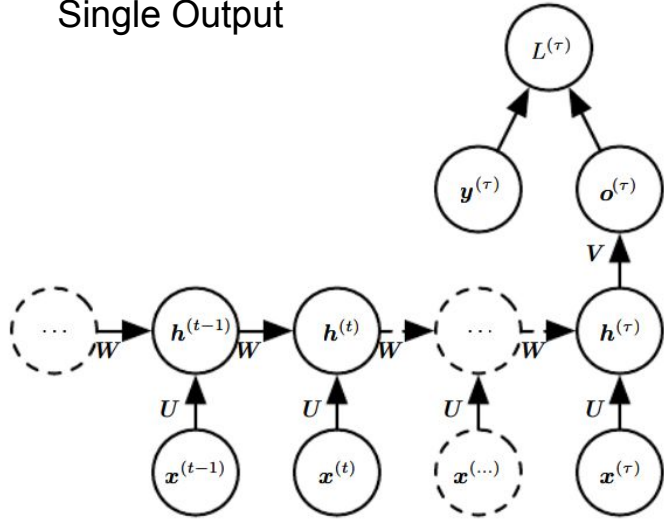


# Architecture - Residual Networks

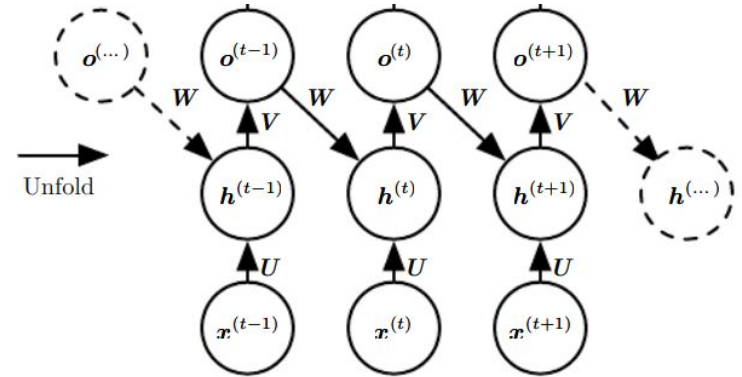


# Architecture - Recurrent Neural Networks

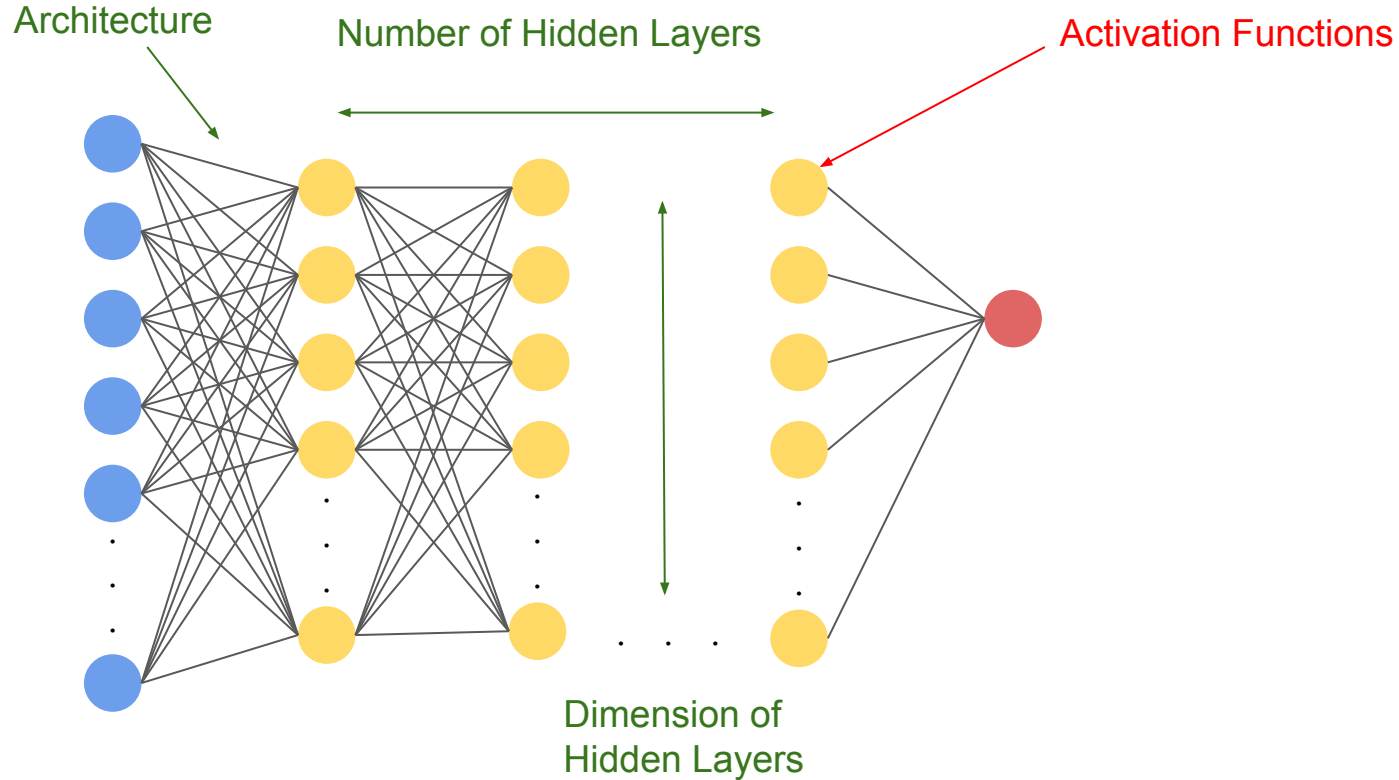
Sequence with  
Single Output



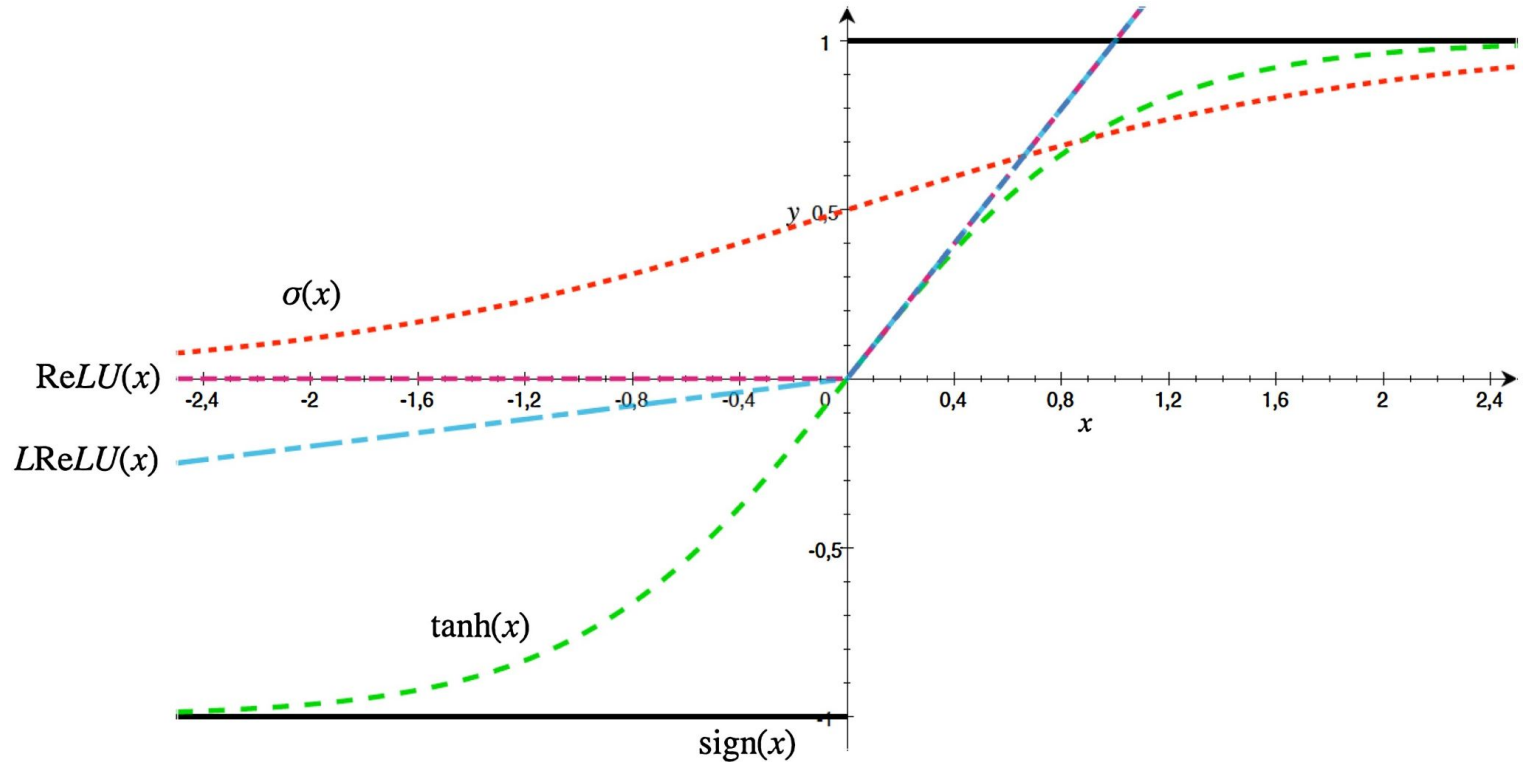
Sequence with  
Output at each  
Step



# Knobs at our disposal



# Activation Functions



# Why Neural Nets?

---



# Universal Approximation Theorem



Any feedforward network with a hidden layer and “squashing” activation can arbitrarily fit any function



# Universal Approximation Theorem

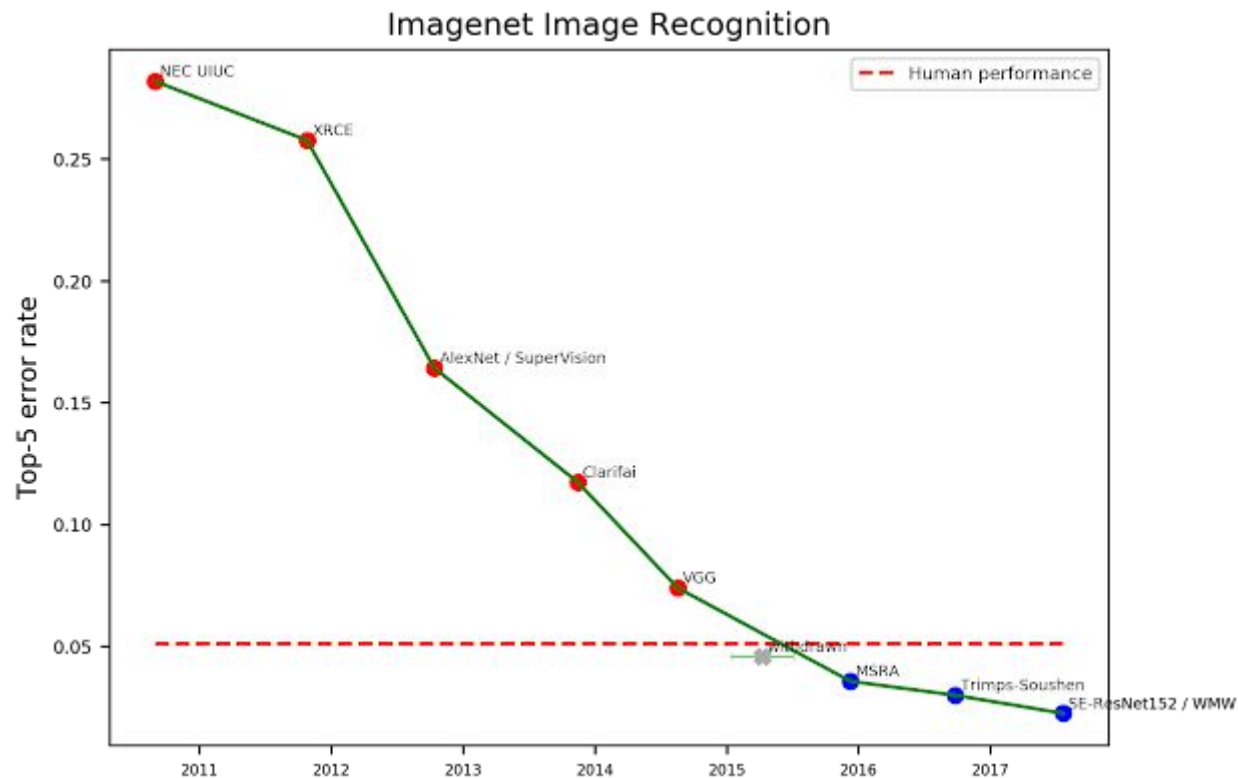


Any feedforward network with a hidden layer and “squashing” activation can arbitrarily fit any function



(provided there are enough hidden units)

Also it just works really well



# Image Classification

Google Photos

Search "Portugal"

Upload



Photos



Explore



Sharing



Print store

LIBRARY



Favorites



Albums



Utilities



Archive

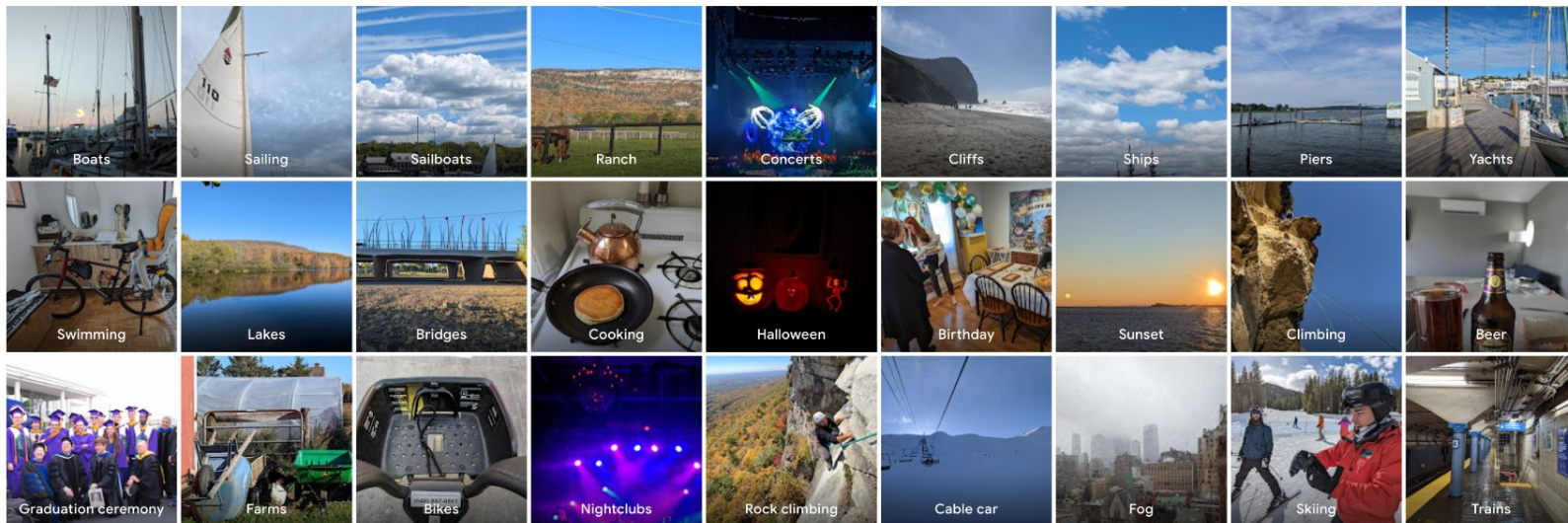


Locked Folder

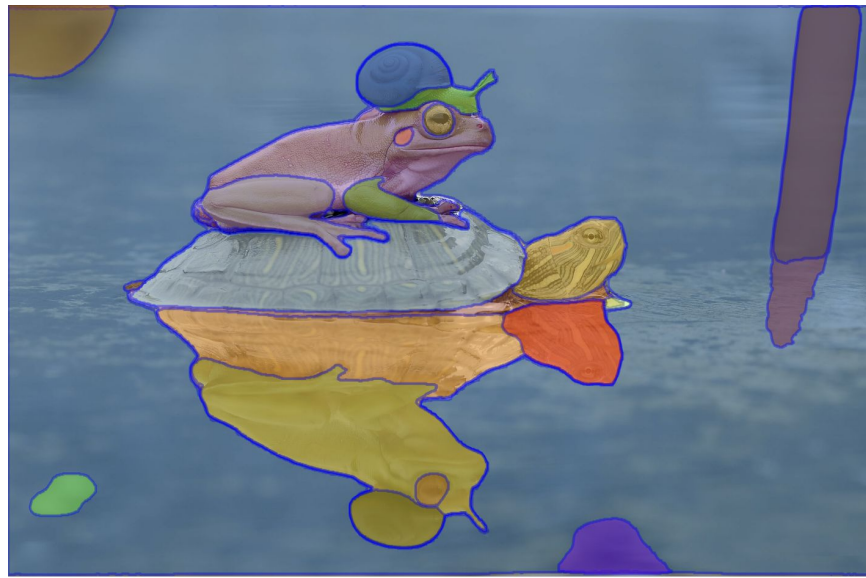


Trash

Things



# Image Segmentation



# Image Generation

"a photograph of an astronaut riding a horse"



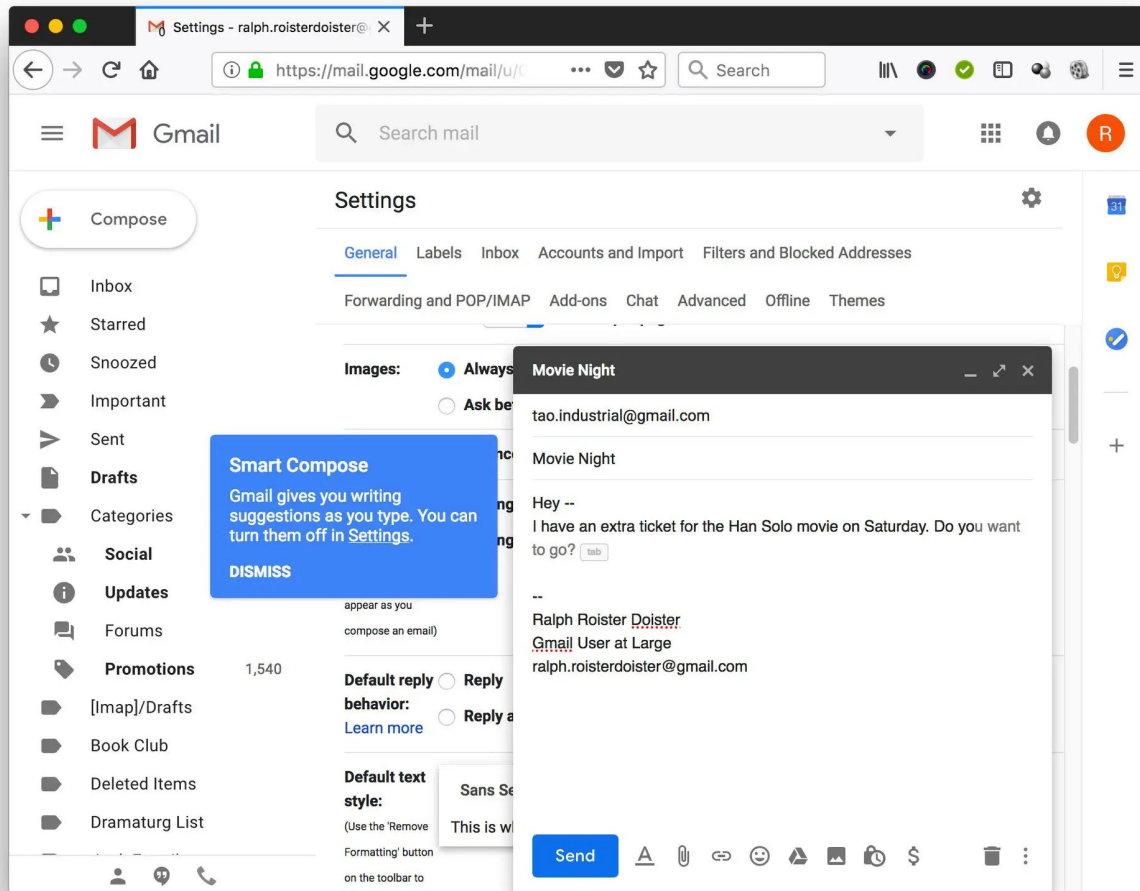
Reference

# Translation



[Link](#)

# Autocomplete





# ChatGPT

ER

What breakthroughs has deep learning been responsible for over the last 10 years?



Deep learning has been responsible for numerous breakthroughs in various fields over the last decade. Here are some of the notable breakthroughs and advancements:



1. Image and Video Recognition:
  - Image classification and object recognition: Deep learning techniques, especially Convolutional Neural Networks (CNNs), have achieved remarkable accuracy in tasks like ImageNet classification.
  - Object detection: Models like YOLO (You Only Look Once) and Faster R-CNN have made significant progress in real-time object detection in images and videos.
  - Image and video segmentation: Deep learning models have improved the accuracy and speed of pixel-wise image and video segmentation, with applications in medical imaging, autonomous vehicles, and more.
2. Natural Language Processing (NLP):
  - Language models: Transformer-based models like GPT-3 and BERT have revolutionized NLP by achieving state-of-the-art results in various language understanding tasks.
  - Language generation: Deep learning models have been used for natural language generation tasks, including text summarization, chatbots, and content creation.
  - Multilingual NLP: Deep learning models have facilitated multilingual NLP applications, allowing for cross-lingual understanding and translation.
3. Reinforcement Learning:
  - AlphaGo: Deep reinforcement learning was at the core of DeepMind's AlphaGo, which defeated the world champion Go player. It marked a significant milestone in AI's ability to tackle complex board games.
4. Healthcare:
  - Medical image analysis: Deep learning has enabled better diagnosis and analysis of medical images, such as X-rays, CT scans, and MRIs, helping in early disease detection.
  - Drug discovery: Deep learning models have accelerated drug discovery by predicting potential drug candidates and their interactions.
5. Autonomous Vehicles:
  - Self-driving cars: Deep learning algorithms have played a crucial role in the development of autonomous vehicles, allowing them to perceive and navigate the environment.
6. Speech Recognition:
  - Automatic speech recognition (ASR): Deep learning has improved the accuracy of ASR systems, making voice assistants and transcription services more effective.

# Speech to Text

<https://huggingface.co/spaces/openai/whisper>

## Whisper

Whisper is a general-purpose speech recognition model. It is trained on a large dataset of diverse audio and is also a multi-task model that can perform multilingual speech recognition as well as speech translation and language identification. This demo cuts audio after around 30 secs.

You can skip the queue by using google colab for the space:



 [Open in Colab](#)


Record from microphone


Transcribe


# Text to Speech


<https://huggingface.co/suno/bark>


 **suno/bark** 


 like 486


 Text-to-Speech


 Transformers


 PyTorch


 13 languages


 bark


 text-to-audio


 audio


 Inference Endpoints


 License: mit


 **Model card**


 Files

 Community 35

 Train

 Deploy

 Use in Transformers


 Edit model card



## Bark


Bark is a transformer-based text-to-audio model created by [Suno](#). Bark can generate highly realistic, multilingual speech as well as other audio - including music, background noise and simple sound effects. The model can also produce nonverbal communications like laughing, sighing and crying. To support the research community, we are providing access to pretrained model checkpoints ready for inference.

The original github repo and model card can be found [here](#).

Downloads last month  
**47,963**





 **Hosted inference API** 

 Text-to-Speech

Compute

This model can be loaded on the Inference API on-demand.

 JSON Output  Maximize

# 4D View Synthesis

<https://zju3dv.github.io/4k4d/>

## 4K4D: Real-Time 4D View Synthesis at 4K Resolution

Zhen Xu<sup>1</sup> Sida Peng<sup>1</sup> Haotong Lin<sup>1</sup> Guangzhao He<sup>1</sup> Jiaming Sun<sup>2</sup> Yujun Shen<sup>3</sup> Hujun Bao<sup>1</sup>  
Xiaowei Zhou<sup>1</sup>

<sup>1</sup>Zhejiang University <sup>2</sup>Image Derivative Inc. <sup>3</sup>Ant Group

 Paper

 arXiv

 Video (Coming Soon)

 Code

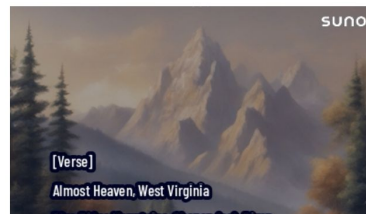
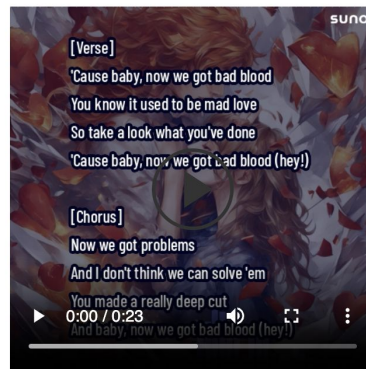
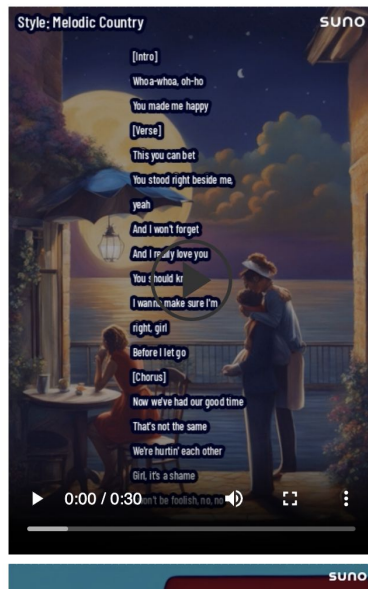
Real-time rendering demo on the DNA-Rendering, ENeRF-Outdoor and Mobile-Stage dataset. The videos might take a few moments to load.



# Text to Music

<https://suno-ai.notion.site/Chirp-v1-Examples-cc71e6c0c79f4e03acf39aa5d5a3dd0>

9



# The Bitter Lesson

- Compute Rules Everything Around Me

# Optimization

---

## Recall Optimizing Linear Regression

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \left( y_i - \vec{X}_i \cdot \vec{\beta} \right)^2$$

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=0}^p X_{ij} \beta_j \right)^2$$

$$\frac{\partial \mathcal{L}}{\partial \vec{\beta}} = 0$$

$$\vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$$



# Optimizing a Neural Network

$$\mathcal{L}(f(\mathbf{X}; \theta); \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{X}_i; \theta))^2$$

$$\frac{\partial \mathcal{L}}{\partial \theta_j} = 0$$

# Stochastic Gradient Descent

$$\mathcal{L}(f(\mathbf{X}; \theta); \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{X}_i; \theta))^2$$

$$\frac{\partial \mathcal{L}}{\partial \theta_j} = 0$$

1. Sample a Batch of data,  $m$ , and calculate the loss on this

$$\frac{1}{m} \sum_{i=1}^m \mathcal{L}(f(\mathbf{X}_i; \theta); \mathbf{y}_i)$$

# Stochastic Gradient Descent

$$\mathcal{L}(f(\mathbf{X}; \theta); \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{X}_i; \theta))^2$$

$$\frac{\partial \mathcal{L}}{\partial \theta_j} = 0$$

1. Sample a Batch of data,  $m$ , and calculate the loss on this

$$\frac{1}{m} \sum_{i=1}^m \mathcal{L}(f(\mathbf{X}_i; \theta); \mathbf{y}_i)$$

2. Calculate the gradient of the loss with respect to a parameter

$$\hat{g}_j = \frac{1}{m} \nabla_{\theta_j} \sum_{i=1}^m \mathcal{L}(f(\mathbf{X}_i; \theta); \mathbf{y}_i)$$

# Stochastic Gradient Descent

$$\mathcal{L}(f(\mathbf{X}; \theta); \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{X}_i; \theta))^2$$

$$\frac{\partial \mathcal{L}}{\partial \theta_j} = 0$$

1. Sample a Batch of data,  $m$ , and calculate the loss on this

$$\frac{1}{m} \sum_{i=1}^m \mathcal{L}(f(\mathbf{X}_i; \theta); \mathbf{y}_i)$$

2. Calculate the gradient of the loss with respect to a parameter

$$\hat{g}_j = \frac{1}{m} \nabla_{\theta_j} \sum_{i=1}^m \mathcal{L}(f(\mathbf{X}_i; \theta); \mathbf{y}_i)$$

3. Update the parameter proportional to the gradient.

$$\theta_j \leftarrow \theta - \epsilon \hat{g}_j$$

# Stochastic Gradient Descent (SGD)

- We don't have to hold all of the data in memory.
- We don't have to perform calculations on all of the data at the same time.
- We don't have to calculate all of the gradients at the same time.
- Surprisingly, this noisy approximation works!

# Stochastic Gradient Descent (SGD)

$$\theta_j \leftarrow \theta - \epsilon \hat{g}_j$$

- There are many modifications to this equation.
- The learning rate,  $\epsilon$ , is a crucial hyperparameter.
- The learning rate can change as optimization progresses (e.g. decay).
- We can scale parameter-specific learning rates based upon their gradients (e.g. Adagrad).
- We can accelerate learning (“momentum”) by using the gradient from the previous iteration.

# Backprop

---

# The Chain Rule of Calculus

$$y = g(x)$$

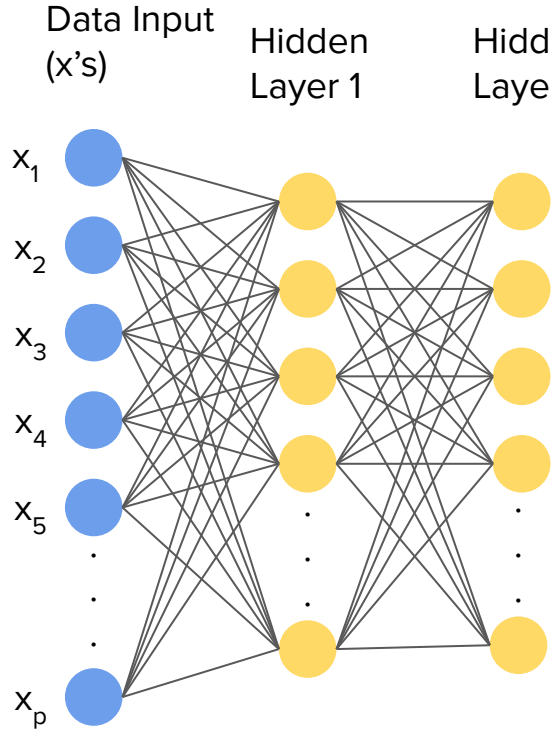
$$z = f(g(x)) = f(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$



# Feed Forward Network / Multilayer Perceptron



$$\begin{aligned}
 h^{(2)}_1 &= f(b^{(2)}_1 + m^{(2)}_{11} * h^{(1)}_1 + m^{(2)}_{12} * h^{(1)}_2 + \dots m^{(2)}_{1L(1)} * h^{(1)}_{L(1)}) \\
 &= f(b^{(2)}_1 + \\
 &\quad m^{(2)}_{11} * f(b^{(1)}_1 + m^{(1)}_{11} * x_1 + m^{(1)}_{12} * x_2 + \dots m^{(1)}_{1p} * x_p) \\
 &\quad + m^{(2)}_{12} * f(b^{(1)}_2 + m^{(1)}_{21} * x_1 + m^{(1)}_{22} * x_2 + \dots m^{(1)}_{2p} * x_p) \\
 &\quad + \dots \\
 &\quad m^{(2)}_{1L(1)} * f(b^{(1)}_{L(1)} + m^{(1)}_{L(1)1} * x_1 + m^{(1)}_{L(1)2} * x_2 + \dots m^{(1)}_{L(1)p} * x_p) \\
 &\quad )
 \end{aligned}$$

$$\begin{aligned}
 h^{(2)}_1 &= f(b^{(2)}_1; m^{(2)}_1; h^{(1)}) \quad , \quad h^{(1)}_1 = f(b^{(1)}_1; m^{(1)}_1; x) \\
 h^{(2)}_1 &= f(b^{(2)}_1; m^{(2)}_1; f(b^{(1)}_1; m^{(1)}_1; x))
 \end{aligned}$$

## Chain Rule → Backprop

$$y = g(x)$$

$$z = f(g(x)) = f(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

$$h_1^{(2)} = f(b_1^{(2)}; m_1^{(2)}; h_1^{(1)}) \quad , \quad h_1^{(1)} = f(b_1^{(1)}; m_1^{(1)}; x)$$

$$h_1^{(2)} = f(b_1^{(2)}; m_1^{(2)}; f(b_1^{(1)}; m_1^{(1)}; x))$$

$$\frac{\partial h_1^{(2)}}{\partial m_1^{(1)}} = \sum_j \frac{\partial h_1^{(2)}}{\partial h_j^{(1)}} \frac{\partial h_j^{(1)}}{\partial m_1^{(1)}}$$

## Chain Rule → Backprop

$$y = g(x)$$

$$z = f(g(x)) = f(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

$$h_1^{(2)} = f(b_1^{(2)}; m_1^{(2)}; h_1^{(1)}) \quad , \quad h_1^{(1)} = f(b_1^{(1)}; m_1^{(1)}; x)$$

$$h_1^{(2)} = f(b_1^{(2)}; m_1^{(2)}; f(b_1^{(1)}; m_1^{(1)}; x))$$

$$\frac{\partial h_1^{(2)}}{\partial m_1^{(1)}} = \sum_j \frac{\partial h_1^{(2)}}{\partial h_j^{(1)}} \frac{\partial h_j^{(1)}}{\partial m_1^{(1)}}$$

# Backpropagation

We can efficiently compute gradients for all parameters by:

1. Computing in reverse order.
2. Storing and updating gradient computations.