

# Longitudinal data analysis in R: Deep learning & transformers

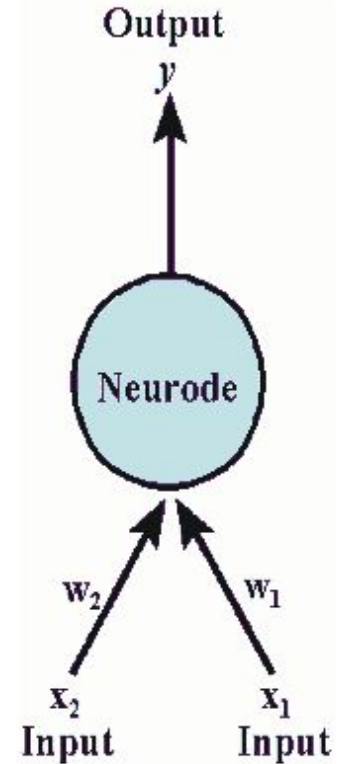
Nelson Nazzicari

*(Bioinformatician & data scientist)* CREA, Lodi (Italy)

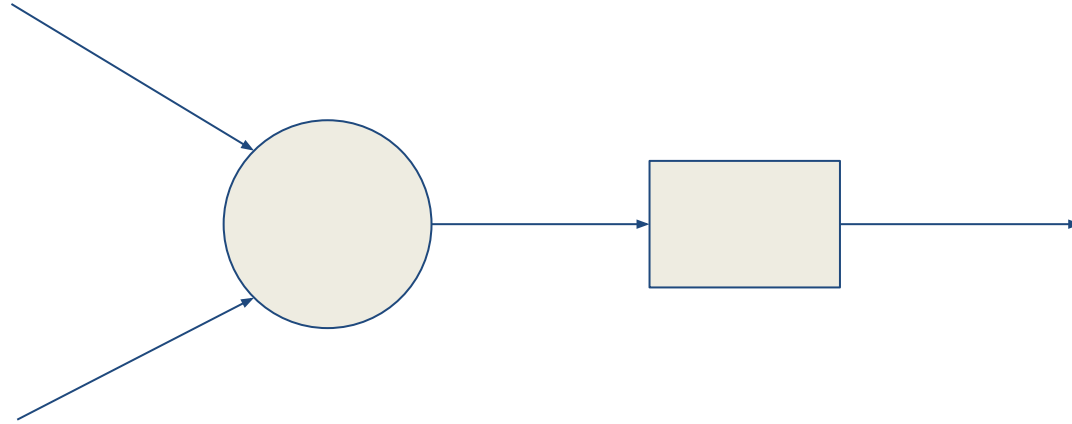


# How did it begin?

- **1943, McCulloch & Pitts:**
  - first mathematical model of a neural network
  - only binary input and output (0/1),
  - only used the threshold step activation function
  - did not incorporate weighting the different inputs
- **1950s:** pioneers started asking whether computers could be made to “think” → (A.I.)
- **1957, Rosenblatt:** introducing the **perceptron**:
  - weighting inputs
  - additional activation functions (e.g. sigmoid)

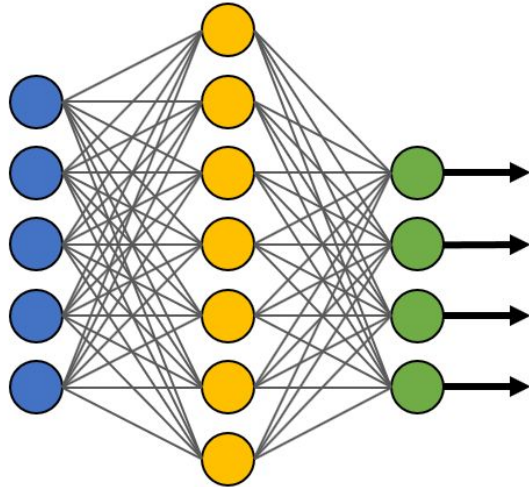


# Deep learning an introduction

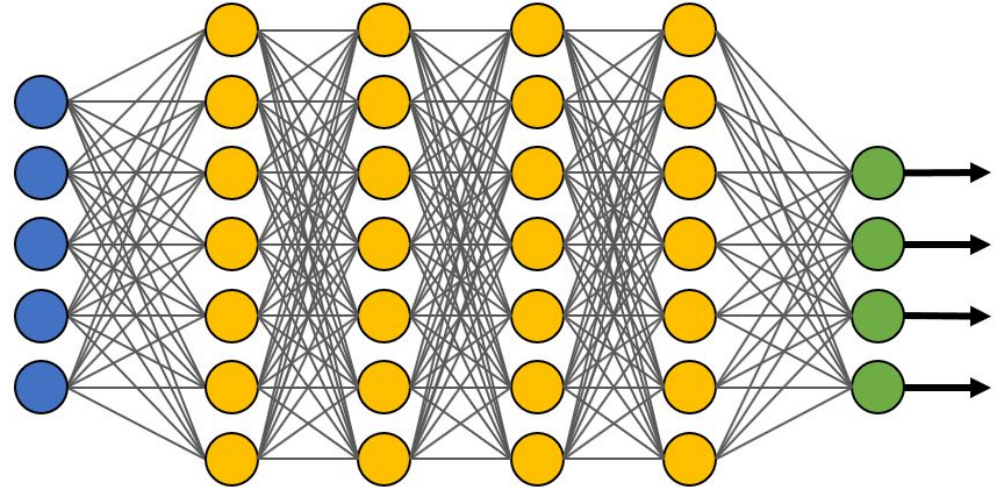


# Deep learning an introduction

## Simple Neural Network



## Deep Learning Neural Network

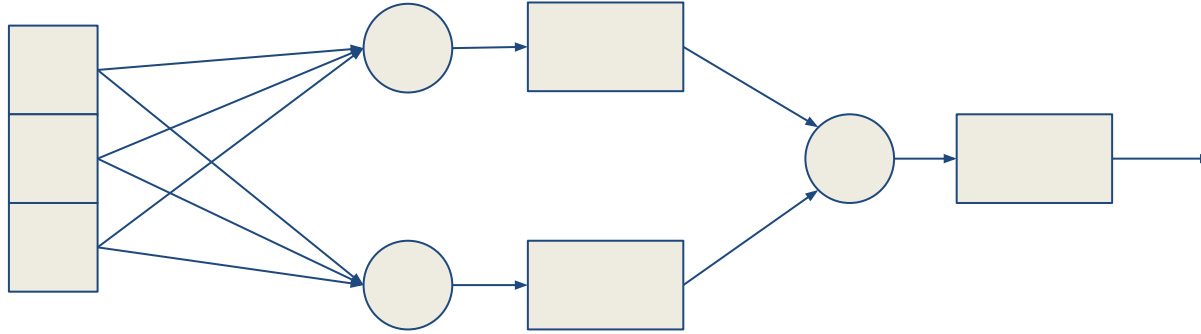


● Input Layer

● Hidden Layer

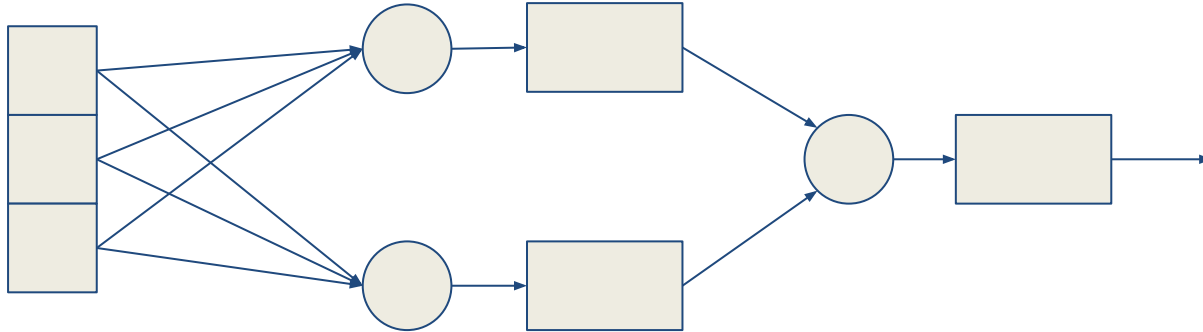
● Output Layer

# How to learn anything?



# Loss function

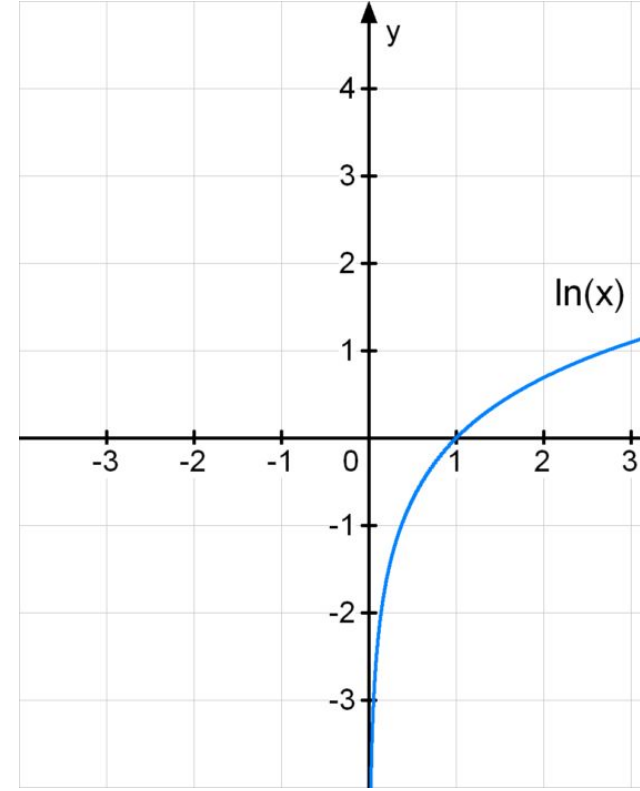
- Measure the distance from true value and prediction
  - Like a residual!
- **Theta**
- $\text{Loss}(\text{Theta}, X_i, y_i)$
- For regression: RMSE is your friend (differentiable)



# Loss function for binary classification

Binary Cross Entropy

$$\text{BCE}(y, \hat{y}) = - (y \ln(\hat{y}) + (1-y) \ln(1-\hat{y}))$$

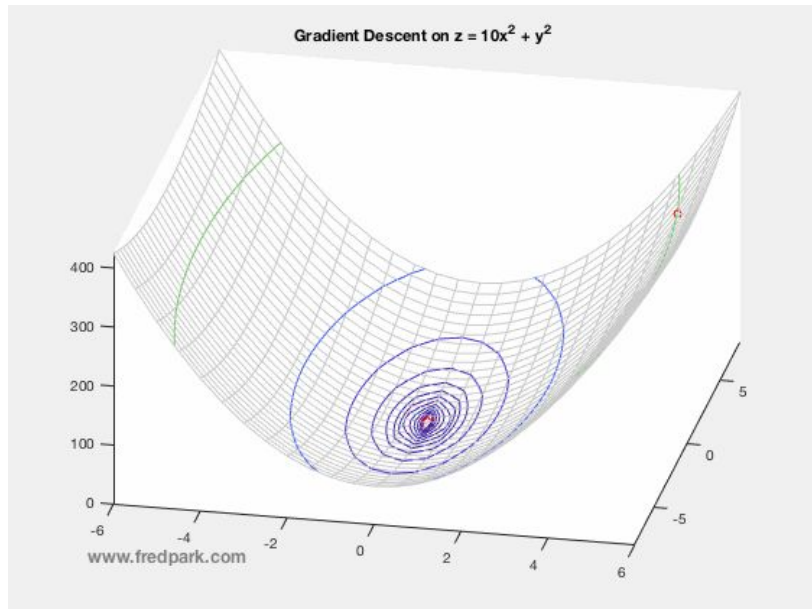


# Gradient descent

- TARGET: minimize the cost function  $J(\text{Theta})$
- ITERATIVELY:
  - Evaluate the gradient of  $J(\text{Theta})$
  - $\text{Delta} = - \text{learning\_rate} * \text{gradient}$
  - $\text{Theta} = \text{Theta} + \text{Delta}$

With

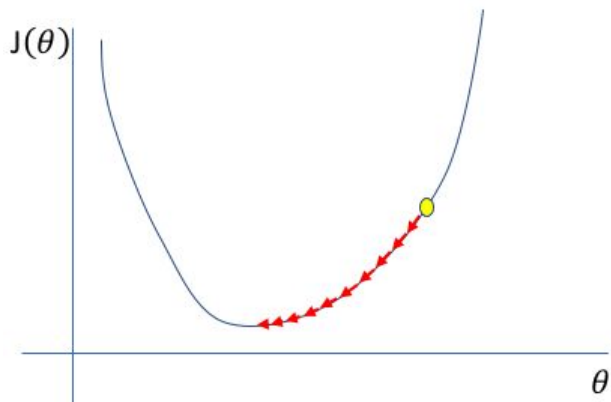
- $\text{Theta}$  : current solution
- $\text{Delta}$  : step in solution space
- $J(\text{Theta}) : \text{mean}(\text{Loss}(\text{Theta}, X_1, y_1) + \text{Loss}(\text{Theta}, X_2, y_2) + \dots)$
- Learning rate...





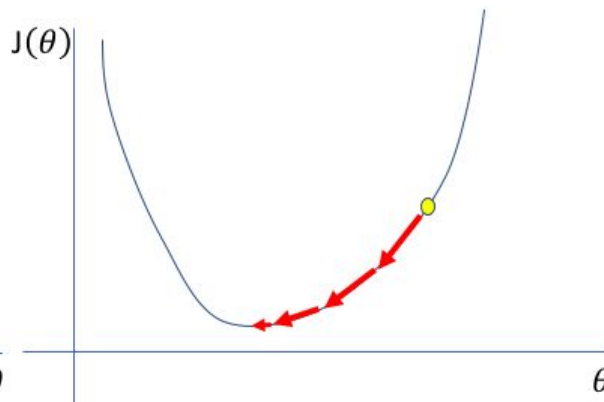
# Learning rate

Too low



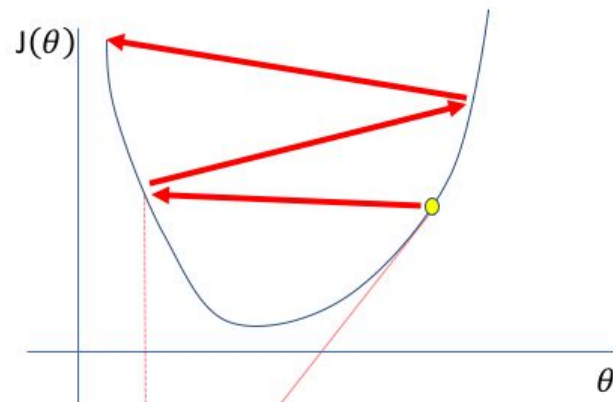
A small learning rate requires many updates before reaching the minimum point

Just right



The optimal learning rate swiftly reaches the minimum point

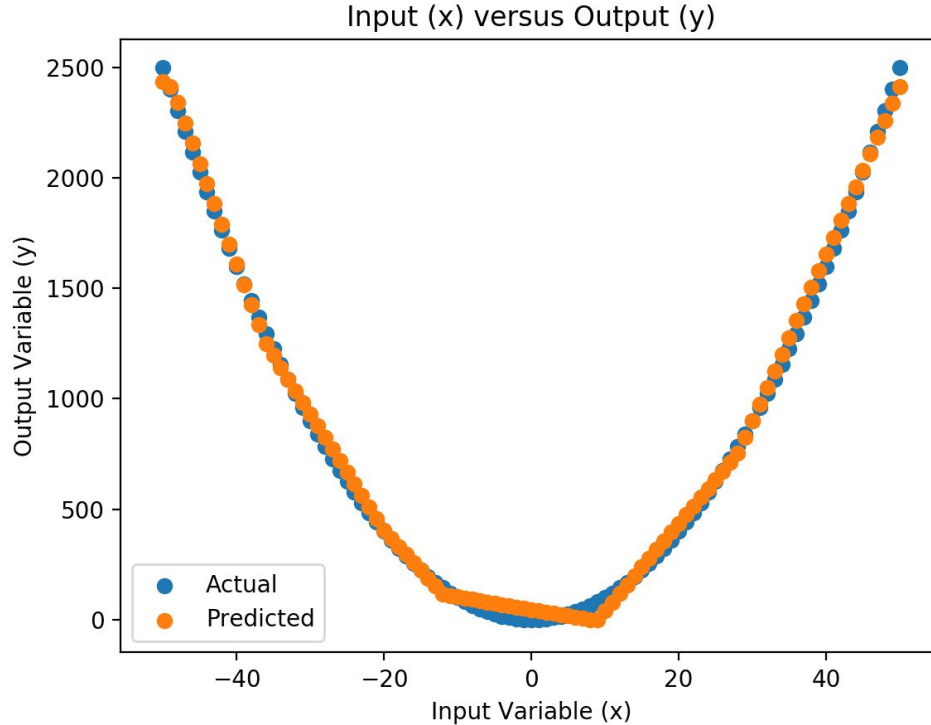
Too high



Too large of a learning rate causes drastic updates which lead to divergent behaviors



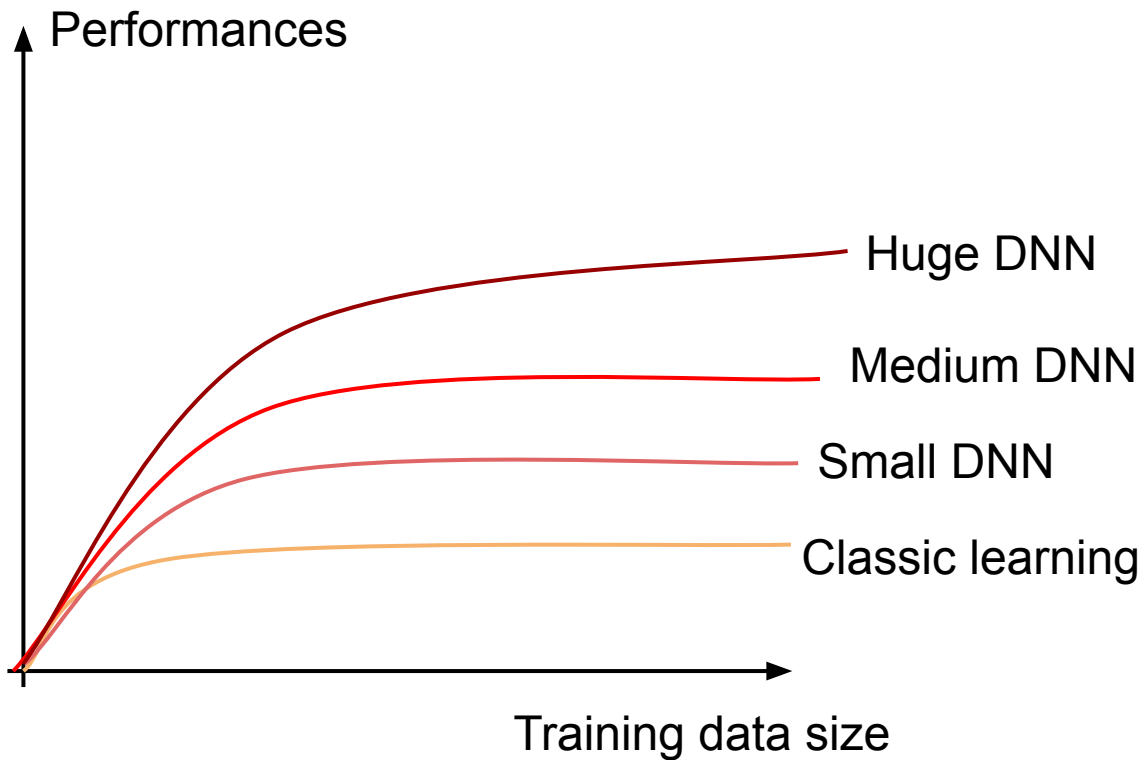
# Function approximators



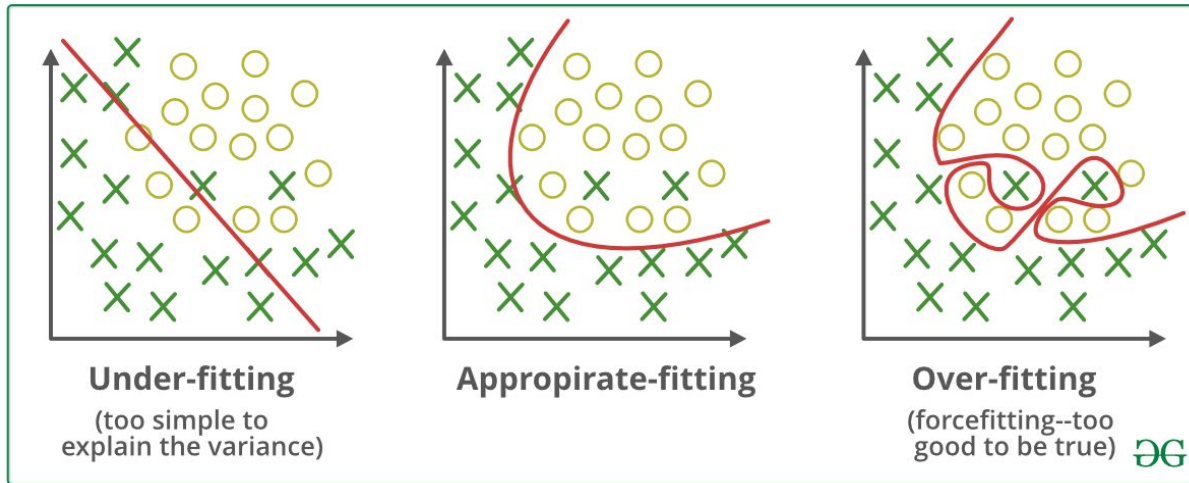
- (known) **quadratic function** (blue line)
- approximated with a **neural network model** [2 hidden layers with 10 nodes each] (orange line)



# Size matters



# The peril of overfitting



# Key concepts

- Weights
- Forward and back propagation
- Loss function, cost function
- Gradient Descent/Optimizers
- Universal function approximation
- Data hungry/overfitting



# Notebook on neural networks



# Network optimization: EfficientNet

<https://arxiv.org/abs/1905.11946>

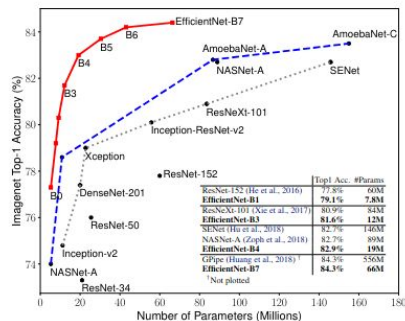
## EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

Mingxing Tan<sup>1</sup> Quoc V. Le<sup>1</sup>

### Abstract

Convolutional Neural Networks (ConvNets) are commonly developed at a fixed resource budget, and then scaled up for better accuracy if more resources are available. In this paper, we systematically study model scaling and identify that carefully balancing network depth, width, and resolution can lead to better performance. Based on this observation, we propose a new scaling method that uniformly scales all dimensions of depth/width/resolution using a simple yet highly effective *compound coefficient*. We demonstrate the effectiveness of this method on scaling up MobileNets and ResNet.

To go even further, we use neural architecture



### EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

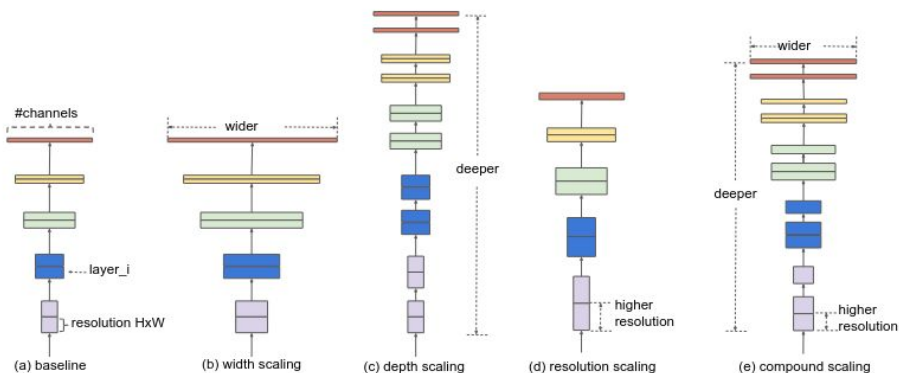
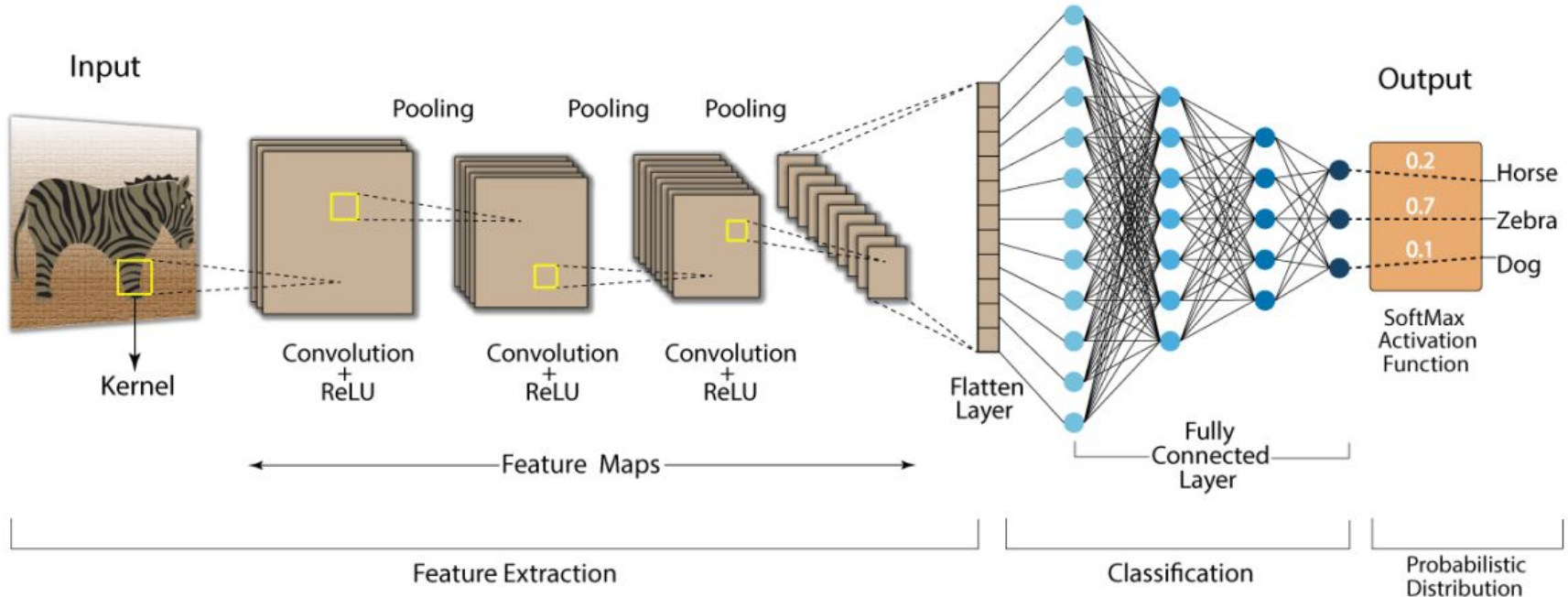


Figure 2. **Model Scaling.** (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

# Convolutional neural networks

## Convolution Neural Network (CNN)

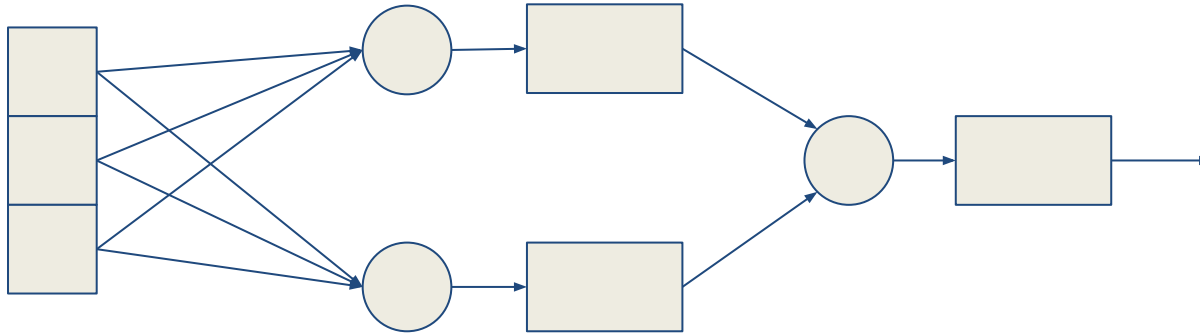




# “Traditional” approach to longitudinal data

## RNNs (Recurrent Neural Networks)

- LSTMs/GRUs (Long Short-Term Memory / Gated Recurrent Units)
- Limitations: vanishing gradients, slow training, difficulty modeling long-term dependencies.



# The rise of transformers

- Transformers don't process sequentially, they use self-attention to look at the entire sequence at once.
- Originally developed for NLP (Natural Language Processing), but now used for time series, longitudinal data, etc.
- "Attention is All You Need" <https://arxiv.org/abs/1706.03762>
- Strengths: parallelizable, better modeling of long-range dependencies
- Always data-hungry
- Can be computationally expensive
- Beware of the hyperparameters flood



# Key concept: tokens

- Mapping input to a dictionary of known tokens (100k+)
- <https://platform.openai.com/tokenizer>



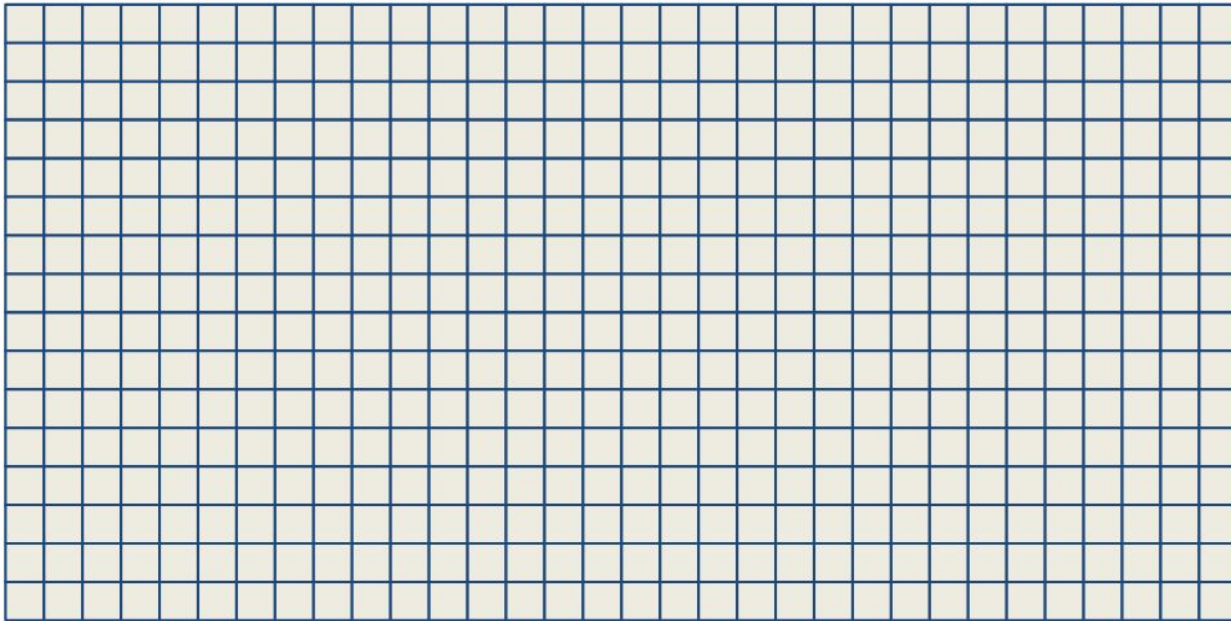
# Key concept: embedding

- Mapping tokens to high dimensional space
  - chatGPT 3.0  $\rightarrow$  12,288 dimensions
  - More recent models: undisclosed



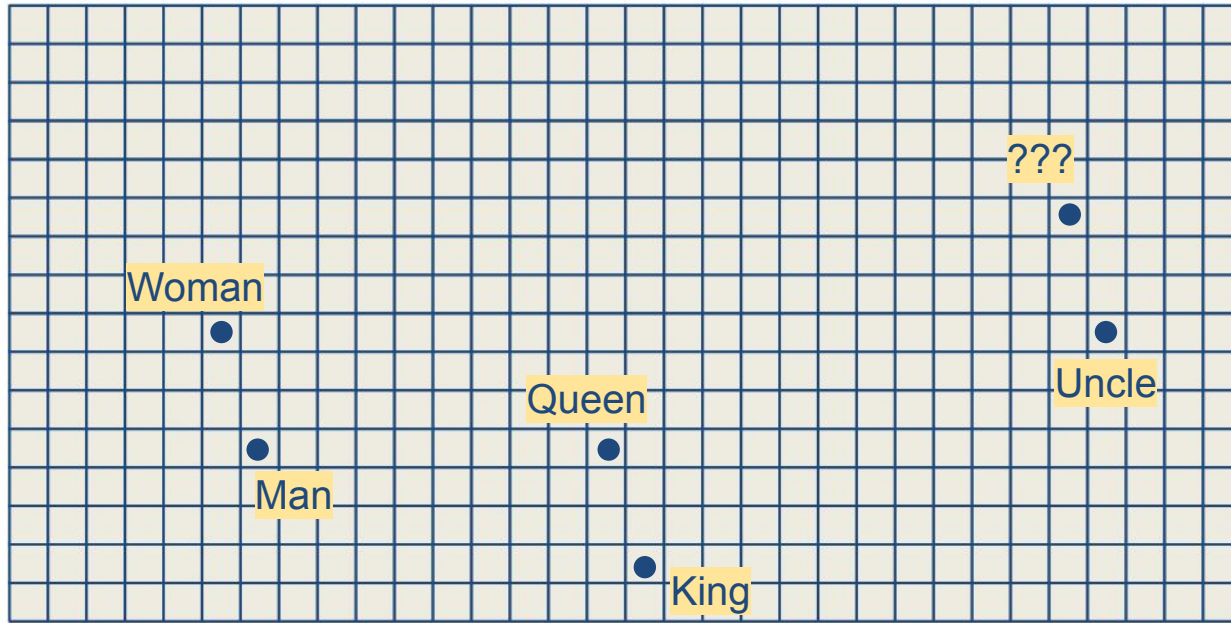
# Key concept: embedding

- Mapping tokens to high dimensional space
  - chatGPT 3.0  $\rightarrow$  12,288 dimensions
  - More recent models: undisclosed



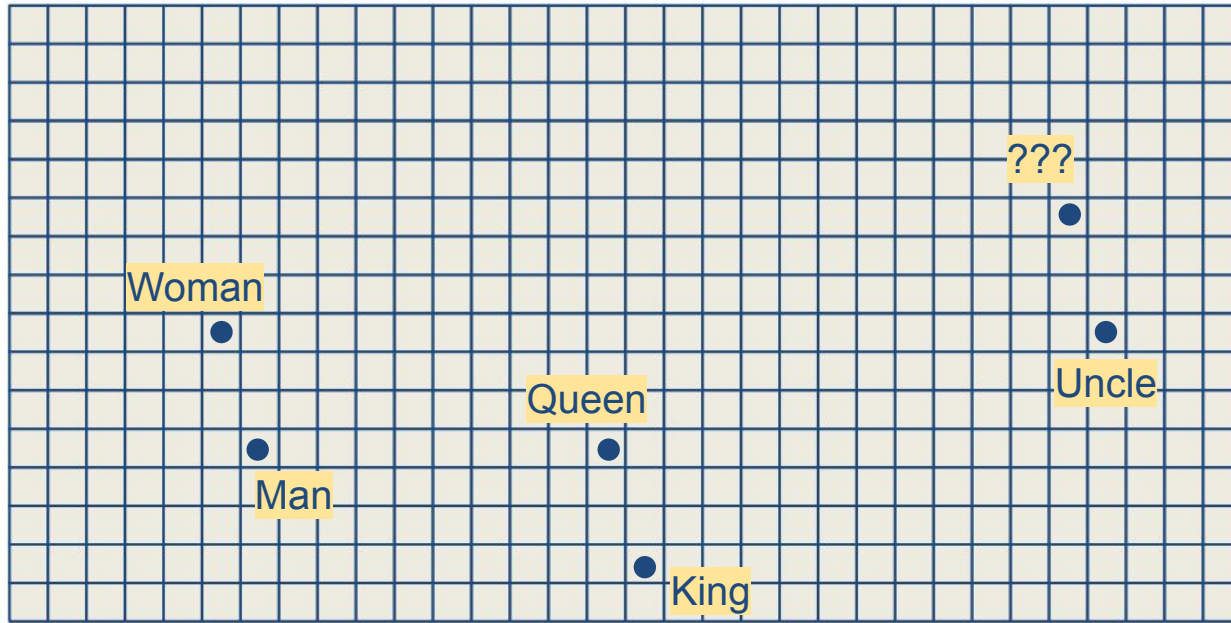
# Key concept: embedding

- Mapping tokens to high dimensional space
  - chatGPT 3.0 → 12,288 dimensions
  - More recent models: undisclosed



# Key concept: embedding

- Mapping tokens to high dimensional space
  - chatGPT 3.0 → 12,288 dimensions
  - More recent models: undisclosed



$$\text{Man} - \text{Woman} + \text{Uncle} = ???$$





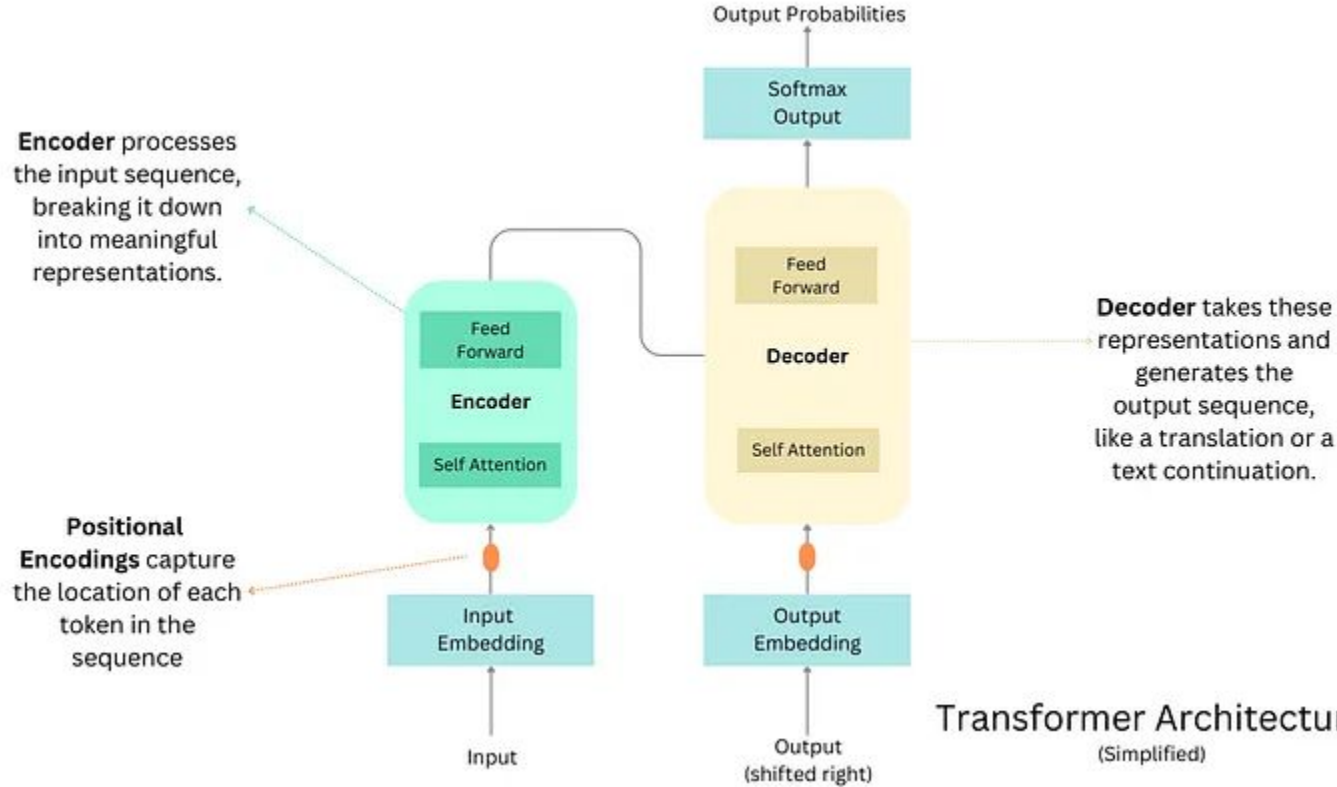
# Key concept: self attention

The cat run through the door because it was open

The cat run through the door because it was hungry

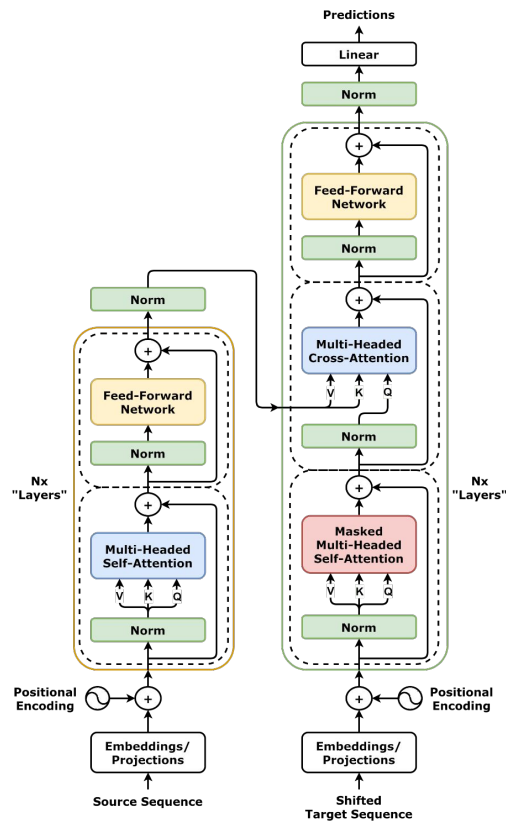


# Transformer architecture



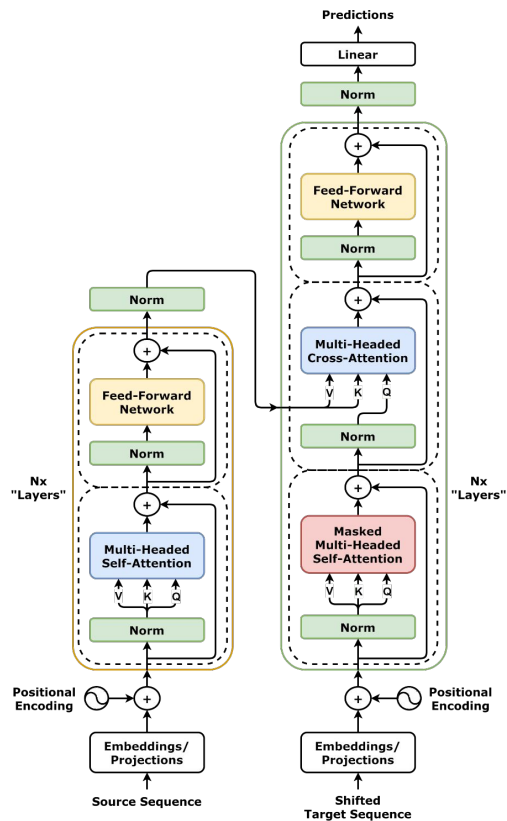
From <https://medium.com/@theaveragegal/transformer-architecture-simplified-3fb501d461c8>

# Transformer architecture



# Transformer architecture

- Input embeddings (including positional encodings — important for sequences!)
- Self-attention: each element attends to others with learned weights
- V (value), K (key), Q (query)
- Feed-forward networks applied to each element
- Layer normalization and residual connections for stability
- Context size is fixed



# Transformers for sequential/longitudinal data



Adaptations to handle time information:

- Positional encoding
- Continuous-time embeddings (e.g., in patient data where time between visits is variable)

Examples of applications:

- Time Series Forecasting (e.g., stock prices, weather)
- Healthcare (e.g., modeling patient trajectories over multiple visits — models like Transformer for EHR data)
- Sensor data analysis (e.g., predictive maintenance)



# Notebook on transformer architecture



# Transformers for sequential/longitudinal data



- Time-series transformers (Informer, TimesNet, PatchTST, etc.)
- Clinical transformers (e.g., BEHRT, Med-BERT)
- Specialized designs: sometimes replace vanilla attention with local/global attention to handle very long sequences
- Model: ProtBERT (meta) / AlphaFold (deepmind)
  - These models treat amino acid sequences as text, where each amino acid is like a word. They use transformers (like BERT or GPT-style architectures) to model the dependencies between residues.



# Notebook on prot BERT

