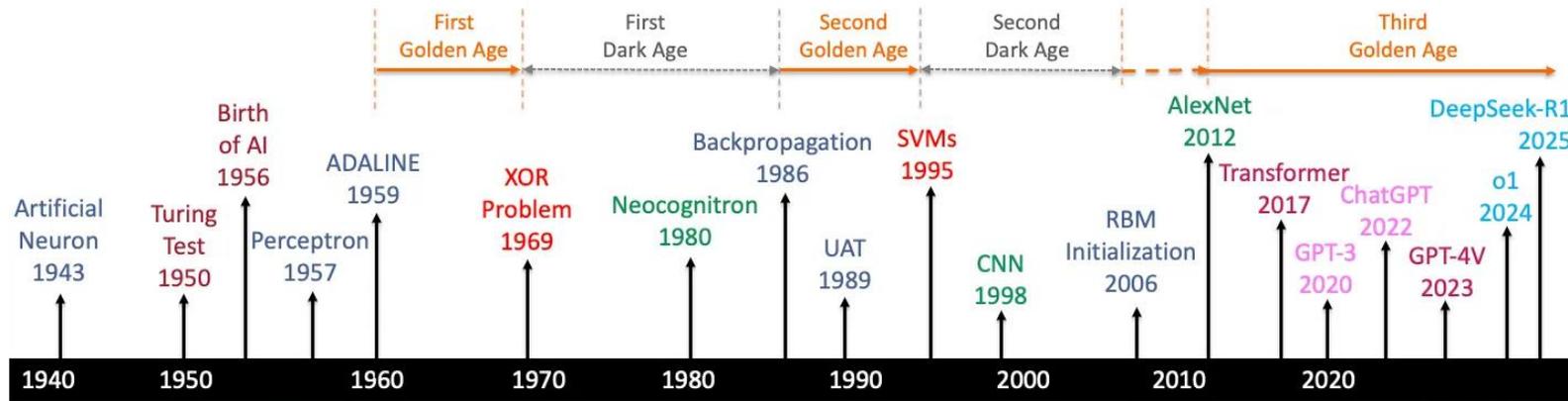


# Foundation Models for Healthcare

Introduction – Part II

# Historically

## A Brief History of AI with Deep Learning



McCulloch-Pitts

Rosenblatt

Widrow-Hoff

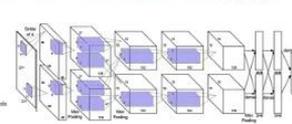
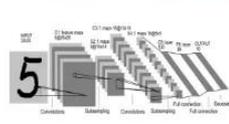
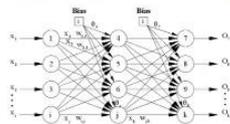
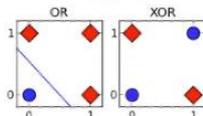
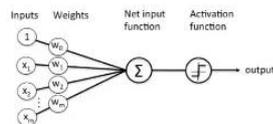
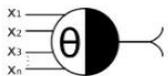
Minsky-Papert

Rumelhart, Hinton et al.

LeCun

Hinton-Ruslan Krizhevsky et al.

Vaswani



# ML Models



Input	Output	Function	Tasks	Application
<ul style="list-style-type: none"><li>• Structured</li><li>• Time series</li><li>• Text</li><li>• Image</li><li>• Audio</li><li>• Video</li><li>• Table</li><li>• Domain-specific</li><li>• Multi-D</li><li>• Multimodal</li></ul>	<ul style="list-style-type: none"><li>• Label</li><li>• Regression</li><li>• Text</li><li>• Image</li><li>• Audio</li><li>• Video</li><li>• Table</li><li>• Domain-specific</li></ul>	<ul style="list-style-type: none"><li>• Clustering</li><li>• Classification</li><li>• Prediction</li><li>• Regression</li><li>• Synthesis</li></ul>	<ul style="list-style-type: none"><li>• Recognition</li><li>• Detection</li><li>• Segmentation</li><li>• Captioning</li><li>• Image/text/synthesis</li><li>• Audio/Video/multimodal synthesis</li><li>• Question answering</li><li>• Text generation</li><li>• Autocompletion</li><li>• Translation</li><li>• Summarization</li><li>• Sentiment analysis</li><li>• Navigation</li><li>• Search/IR</li><li>• Recommendation</li><li>•</li></ul>	<ul style="list-style-type: none"><li>• Customer service</li><li>• Social Media</li><li>• Marketing</li><li>• Code generation</li><li>• Automated reporting</li><li>• Domain-specific</li></ul>

# Foundation Model



## Input

- Trained on a variety of data input

## Output

- Different types of output based on downstream use

## Function

- Defined based on the scope of ability, range of uses, breadth of tasks or types of output and input
- Encoder
- Decoder
- Encoder-Decoder

## Tasks

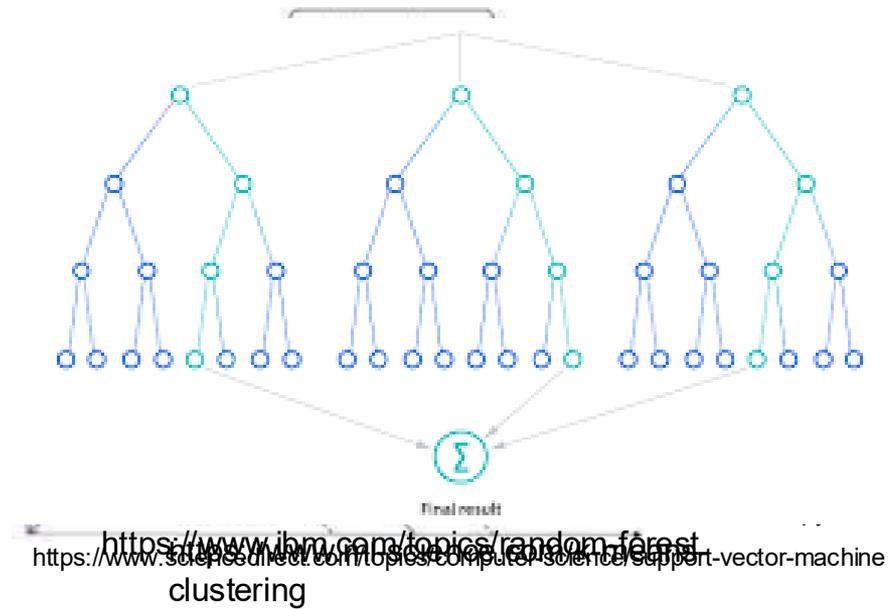
- A single model can accomplish multiple tasks:
  - Sentence completion
  - Sentiment classification
  - Summarization
  - ....
- Works for arbitrary input
- Can be fine-tuned

## Application

- Same model can serve multiple applications
  - Customer service
  - Code generation
  - ...

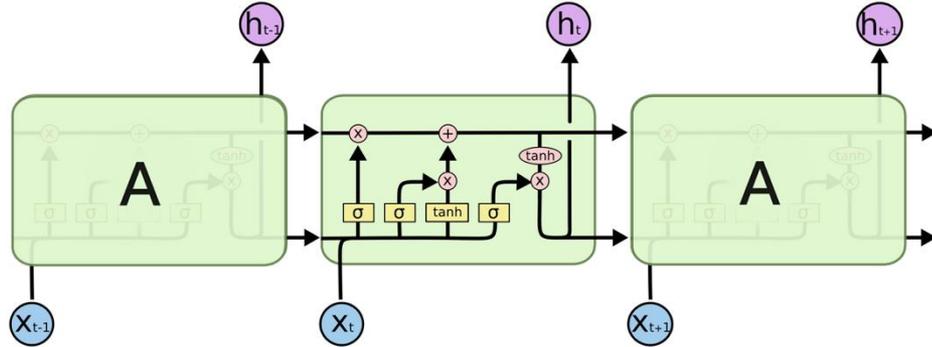
# Are these Foundational Models?

- Decision Trees
- Clustering
- Support Vector Machines
- Random Forests



# Are these Foundational Models?

- Multi-layer perceptron
- RNN
- LSTM



The repeating module in an LSTM contains four interacting layers.

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

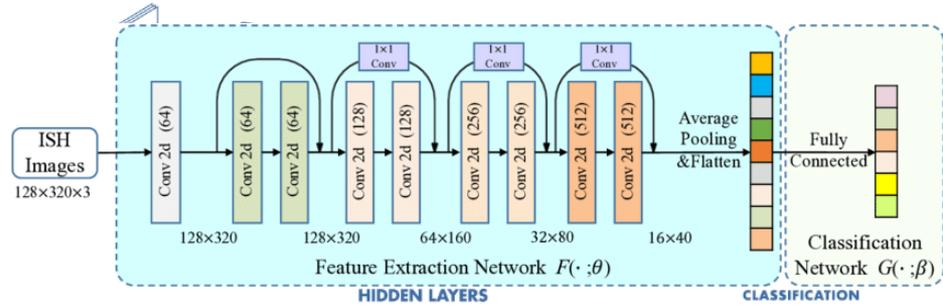
Layer  $L_1$

Layer  $L_2$

<http://deeplearning.stanford.edu/tutorial/supervised/MultiLayerNeuralNetworks/>

# Are these Foundational Models?

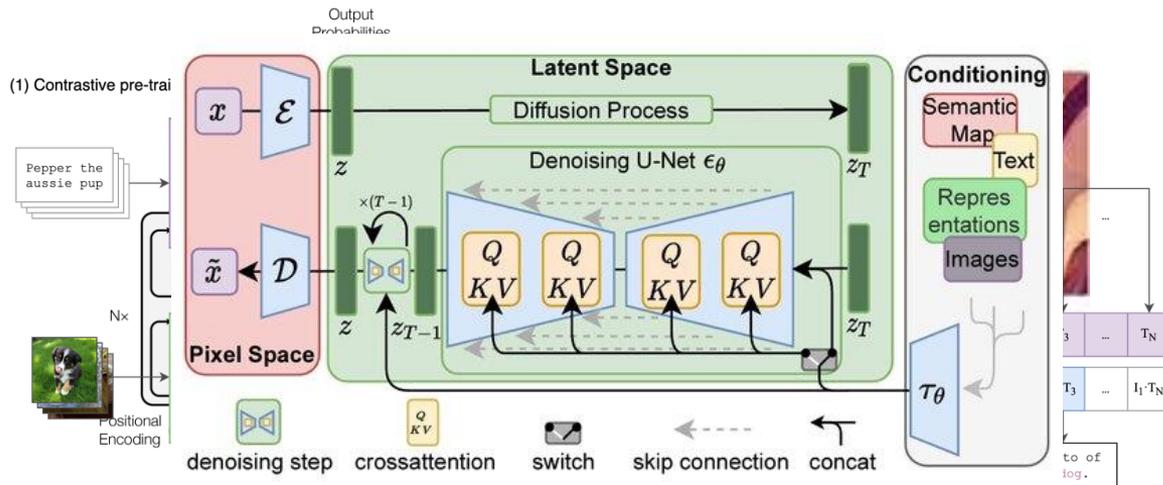
- Convolutional neural networks
- ResNet
- DNN



<https://www.linkedin.com/pulse/deep-residual-networks-resnet-ayoub-kirouane>  
<https://towardsdatascience.com/convolutional-neural-networks-explained-9cc5188c4939>

# Are these Foundational Models?

- Auto-encoders
- GANs
- Transformers
- GPT
- CLIP
- Diffusion models



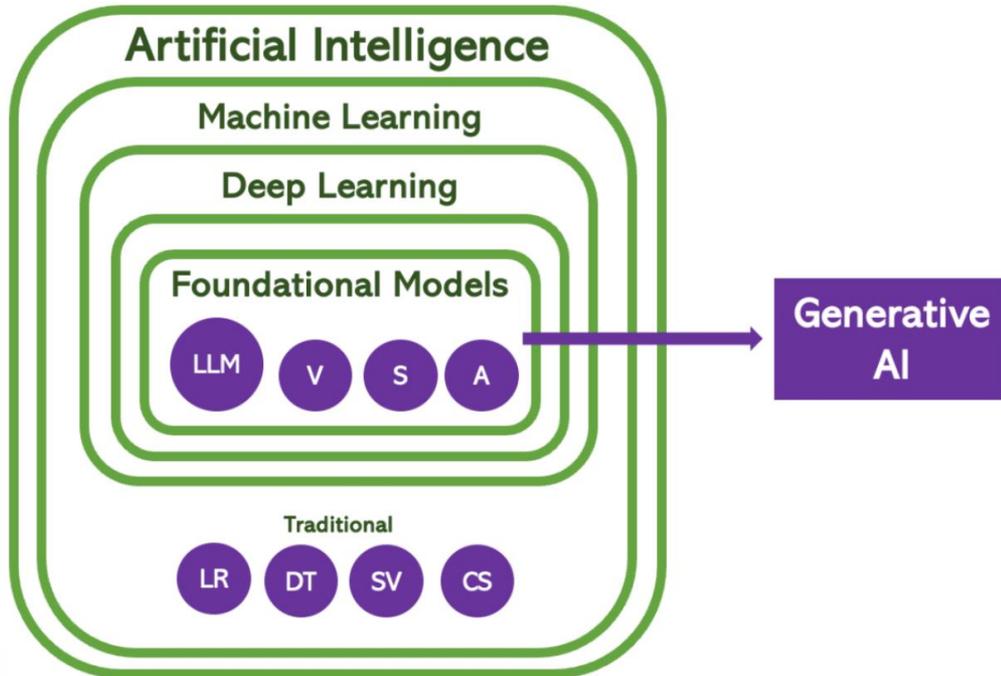
[https://openaccess.thcf.org/content/2022/papers/Rombach\\_High-Resolution\\_Image\\_Synthesis\\_using\\_Latent\\_Diffusion\\_Models\\_CVPR\\_2022\\_paper.pdf](https://openaccess.thcf.org/content/2022/papers/Rombach_High-Resolution_Image_Synthesis_using_Latent_Diffusion_Models_CVPR_2022_paper.pdf)

<https://arxiv.org/abs/1706.03762>

# How has the landscape changed?

Traditional ML	Foundational models
Engineered features	Auto-derived features
Novel architectures but are they useful now?	Mostly transformer variants.
Multiple independent instances of the product	Single instance adapted for use cases through fine-tuning
Trained with smaller, hand-selected and supervised datasets	Trained with self-supervised on very large datasets
Trained for a specific purpose, often classification.	Trained once and applied to different use cases.
Less compute resources	Power hungry GPU-dependent
Deployment was simple as library calls	Full ML Ops behind training, fine-tuning and deployment.
Humans focused on design and inventions of architectures, less focus on data	More focus on data since architectures are standardized. Most improvements are practical with less radical changes.
Less focus on data collection	A whole industry of annotators, alignment, and LLM training data collection
Custom effort for every instance	Less ongoing maintenance after development
Fewer parameters	Large number of parameters
More domain specificity, less generalizability	Lack of domain specificity
Risk of overfitting	Dataset bias

# Evolution of ML models



<https://www.linkedin.com/pulse/ai-alphabet-soup-from-buzzwords-brilliance-brian-gruttadauria-fxtie/>

# What is the general field headed?

- Foundational models still have limitations
  - Not truly domain independent despite the large knowledge given
  - Reasoning powers are still limited.
  - Largely unsuccessful in healthcare world
- New thinking is needed for AI to make a larger impact.
  - Domain understanding
    - New datasets interdisciplinary emerging
  - Domain reasoning
  - New architectures
    - From just scaling models to architectural efficiency, better compute/memory, and flexible designs for innovation
      - Best Paper at NeurIPs - gating mechanisms in transformer-based LLMs to improve efficiency and stability (bringing back LSTM ideas into transformers)
- More from larger to smaller models
- Ethical concerns still looming
- Real cognition still far away

# Deep dive into some of the basics

- Gauge understanding of the fundamentals
- Demos

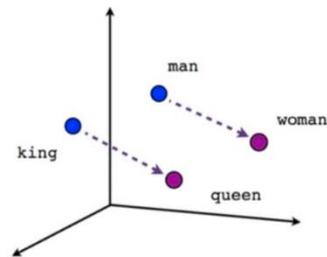
# Masked language prediction – A key element of foundational modeling

I would like to say a few more ADJECTIVE words about the most important invention of the twentieth century. I am not referring to science AN INVENTION or even to the discovery of science A FOOD. The most important ADJECTIVE invention, in my opinion, is the sneaker. If it were not for sneakers, our lives PART OF BODY (PLURAL) would be dirty, cold, and miserable ADJECTIVE. Sneakers keep me from skidding if the rocks PLURAL NOUN are slippery, and when I run, they keep me from stubbing my foot PLURAL NOUN.

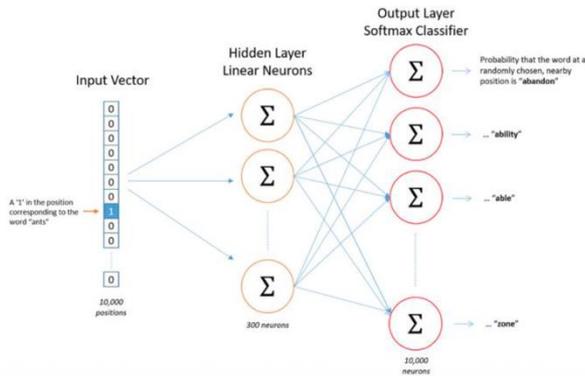
(source: Medium article 2018)

# Word2Vec model

- From Words to numeric vectors
- 2-layer NN
- Input : one-hot vector of vocabulary
- Hidden layer : A vector per vocabulary word serving as the embedding
- Output: predicted probabilities of seeing another word within distance k from the input word
- Captures both syntactic and semantic similarity

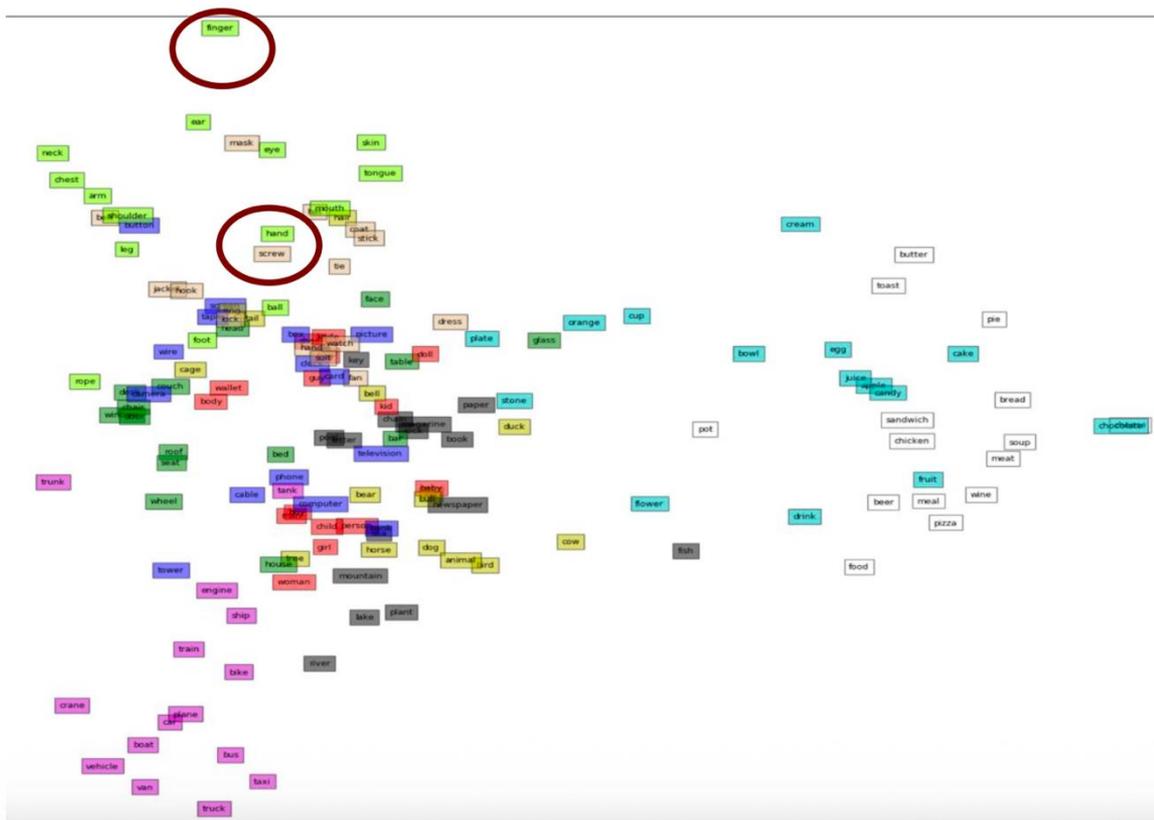


Demo



Understanding the neural network training of Word2Vec model

# How good are these representations?



# Transformers

- Your kitchen sink of good ideas all rolled into one architecture
  - Tokenizing the input text
  - Positional encoding of tokens to preserve order
  - Attention – the Q,K,V paradigm
  - Multi-head attention – for more combinations to explore
  - Multi-layer encoders for increased abstractions
  - Linear layers for aggregating it all together
  - Explores this in parallel for all tokens unlike RNN
  - Representations created more powerful than CNNs!

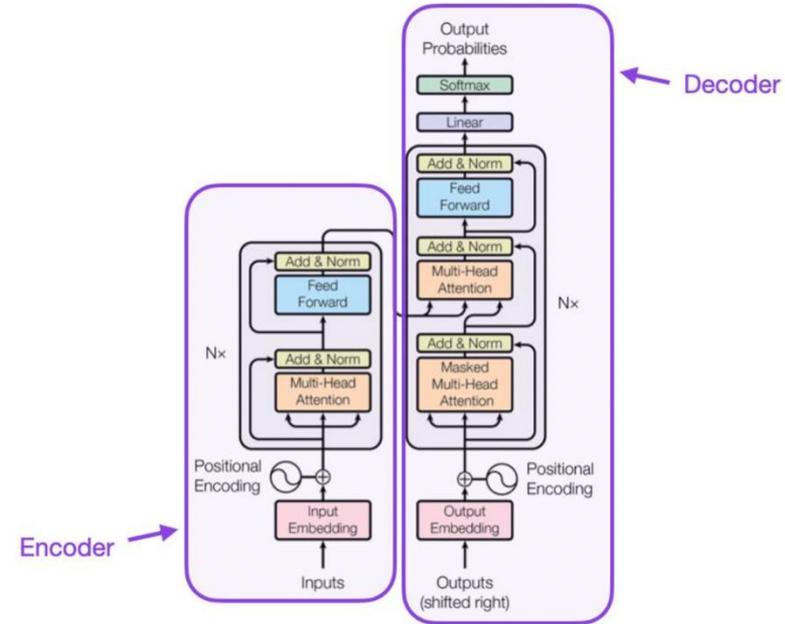


Figure 1: The Transformer - model architecture.

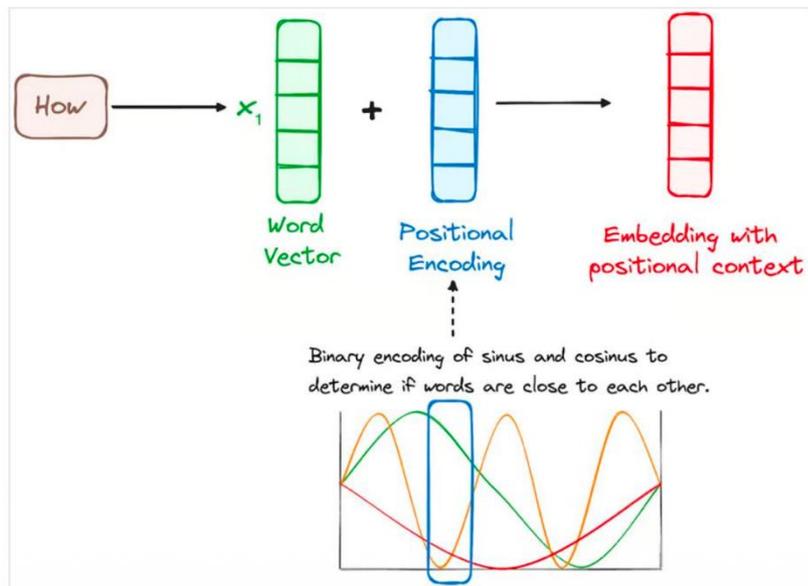
<https://www.linkedin.com/pulse/lm-transformer-architecture-shivasish-mahapatra-kj9qf/>

# Tokenizing the input

- Word2Vec encoding of base vocabulary
  - 30,000+ words
- Word-piece tokenization for handling out-of-vocabulary
- E.g. A sequence of numbers or letters that opens a combination lock
  - ['<s>', 'a', 'Gsequence', 'Gof', 'Gnumbers', 'Gor', 'Gletters', 'Gthat', 'Gopens', 'Ga', 'Gcombination', 'Glock', '</s>']
- This combination of terms makes sense to me to describe this beautifying phenomenon .
  - ['<s>', 'This', 'Gcombination', 'Gof', 'Gterms', 'Gmakes', 'Gsense', 'Gto', 'Gme', 'Gto', 'Gdescribe', 'Gthis', 'Gbeaut', 'ifying', 'Gphenomenon', 'G.', '</s>']
- Tokens may be non-sensical from language perspective

# Capturing sequence information

- Allows for sequence and nearness of words.
- But is additive operation the only possibility?



$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

$$\begin{aligned} f(a+b) &= \begin{bmatrix} \cos(a+b) \\ \sin(a+b) \end{bmatrix} \\ &= \begin{bmatrix} \cos a \cos b - \sin a \sin b \\ \cos a \sin b + \sin a \cos b \end{bmatrix} \\ &= \begin{bmatrix} \cos b & -\sin b \\ \sin b & \cos b \end{bmatrix} \begin{bmatrix} \cos a \\ \sin a \end{bmatrix} \\ &= R(b)f(a) \end{aligned}$$

Why sinusoidal?

where  $R(b)$  is the rotation matrix for angle  $b$ . A similar identity holds if  $f(a)$  is replaced with  $f(\omega a)$  for any  $\omega \in \mathbb{R}$ .

In contrast, suppose  $f$  is the binary encoding:

$$f(a) = \begin{bmatrix} a \bmod 2 \\ \lfloor a \div 2 \rfloor \bmod 2 \\ \lfloor a \div 4 \rfloor \bmod 2 \\ \lfloor a \div 8 \rfloor \bmod 2 \\ \dots \end{bmatrix}$$

Then  $f(a+1)$  is not a linear function of  $f(a)$ .

<https://www.datacamp.com/tutorial/how-transformers-work>

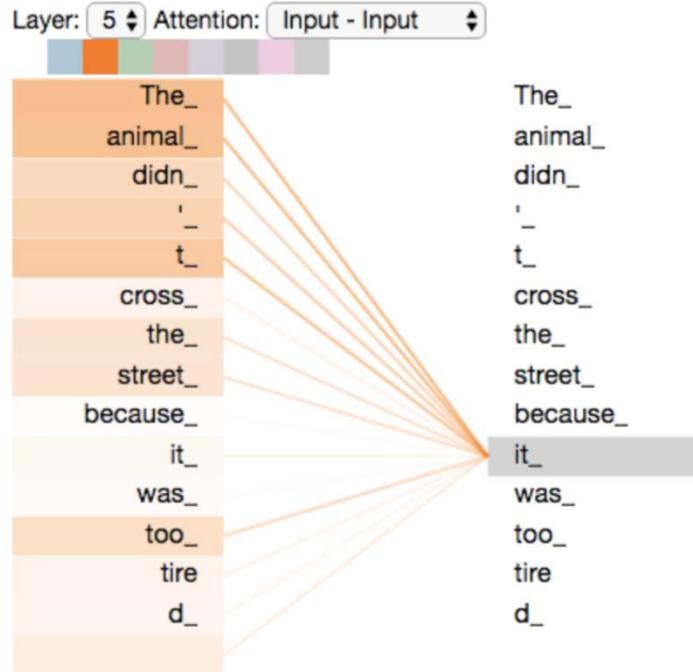
# Attention

- What is (Q,K,V)
  - From Information retrieval literature
    - Key is the index or pointer to the content, q=query, V=actual value stored
    - Also a model for human memory retrieval
  - In the context of sequence modeling:
    - Capturing relations between words.
      - Given a word Q, and its neighbors x-steps away K, can you predict the likelihood of a word V?
      - In Word2Vec we did a fixed K (e.g. immediate neighbor)
      - More general masked language modeling



# Attention

e.g. I am going to my home and play with toy house.



The output embedding for “house” is mainly attending to “toy” and “play”, while the output embedding for “home” is attending to “my” and “going”. Even though home and house are synonyms and more related.

$$\text{Attention}(q, k, v) = \text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

from to

Attention weights

vector dimensionality of K, V

Normalizes the score

<https://jalammar.github.io/illustrated-transformer/>

Does it resolve reference to context?

How is this different from cosine similarity?

# Multi-layer encoders

- Why multi-layer when we are already doing multi-head attention?
  - With multiple attention heads, each head learns a separate aspect of the data
  - Multiple layers build progressively more complex and abstract representations of the input sequence. Each layer in the encoder builds upon the previous one, refining the representations and capturing richer contextual information.
    - Similar to CNNs going over larger and larger areas of the image
  - By stacking multiple layers, the model can capture context not just at the local level, but also at a more global level, understanding how different parts of the sequence relate to each other
    - The first layer gives context for the actual adjacent.
    - And the feedforward layer after that attention layer recognizes phrases.
    - The next attention layer gives context to phrases about other phrases, so that “threw the ball” and “over the wall” get paired together.
    - And then higher and higher features get recognized each layer and the next attention contextualizes those

# Why linear layers or feed-forward layers?

- Map input representations to output predictions
- Transforming data within the model's architecture.
- Project vectors into different dimensions
- Enables the model to learn complex relationships, transform data, and generate meaningful output

# Summarizing the transformer model

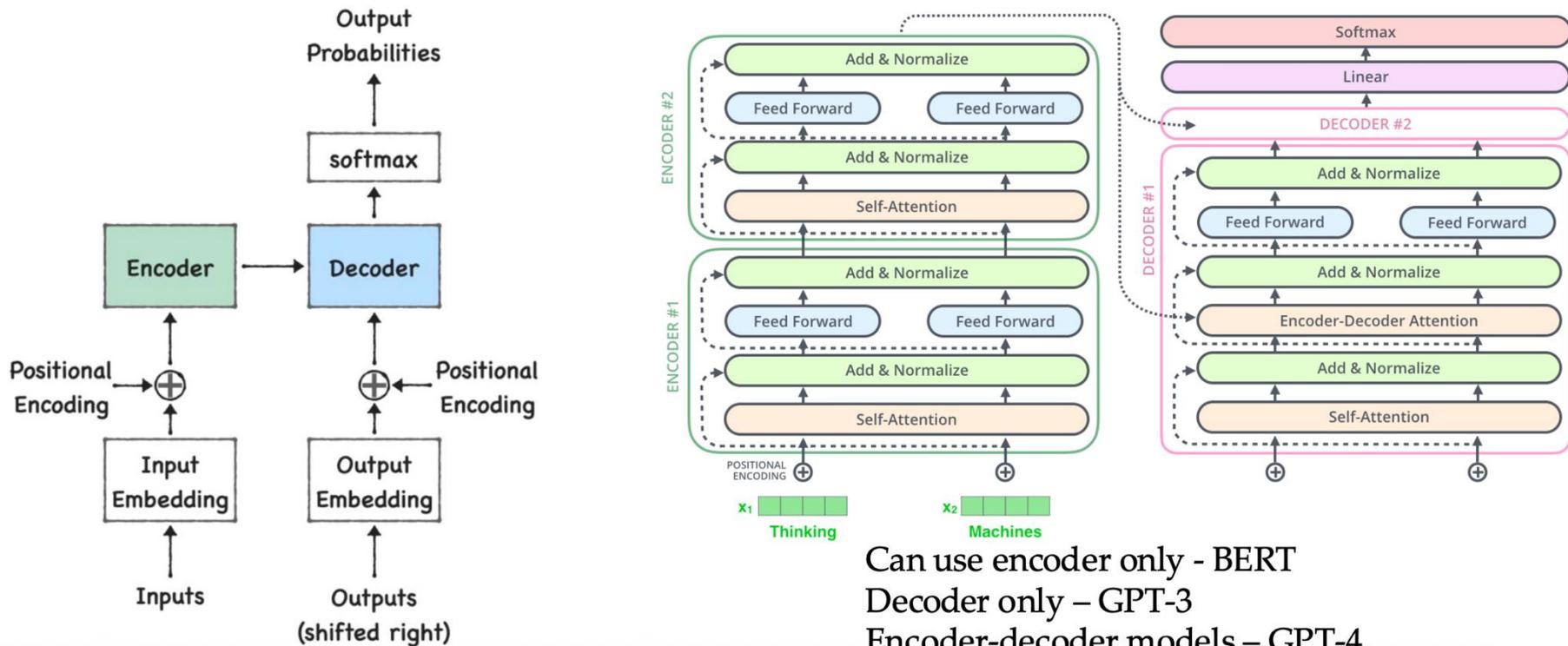


Figure 1: The original encoder-decoder architecture of a transformer model (adapted from Vaswani et al, 2017)