

Got It AI
Tr[^]ansforming Conversational AI

Got It
Tr[^]



PhotoStudy

Transformers, Large Scale Pre-trained Language Models & Applications in Conversational AI

Chandra Khatri, Chief Scientist and Head of AI Research



Agenda

1. Background:

- a. Neural Networks
- b. Recurrent Neural Networks
- c. Language Models (LMs)
- d. Word Embeddings

2. Transformers:

- a. Building blocks of Transformers
- b. Large Scale Pre-trained LMs
- c. From Millions to Trillions of Parameters
- d. Different kinds of Transformers: Encoders, Decoders, Seq2Seq

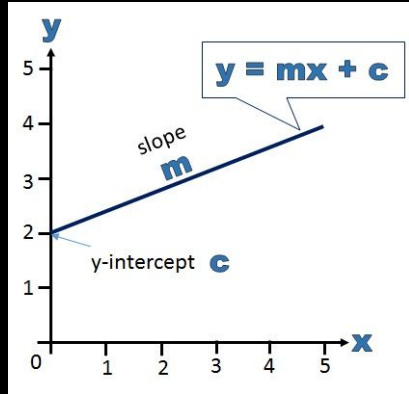
3. Applications:

- a. Transformers for NL Technologies: State of the Art and Recommendations
- b. Evolution of Conversational AI and how Transformers are democratizing the space
- c. Transformers at Got It AI

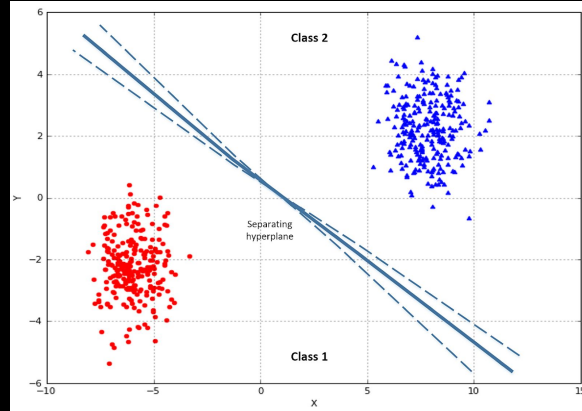
Neural Networks Background

Fundamental Unit of Neural Networks and Deep Learning:

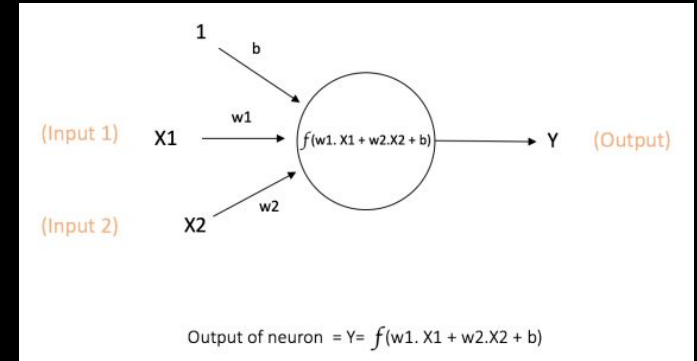
Linear Unit



Linear Equation



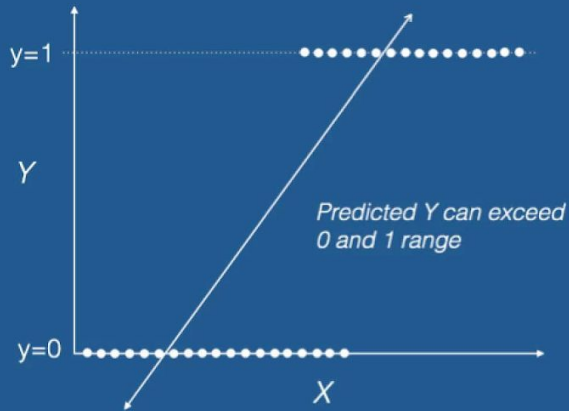
Linear Classification



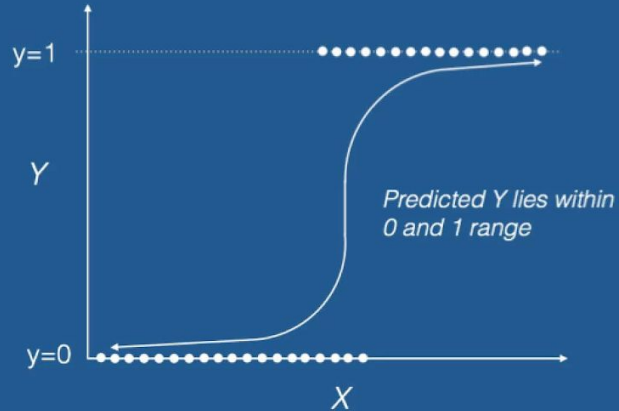
Mathematical representation of a
Linear Unit / Perceptron

Linear Regression and Logistic Regression

Linear Regression

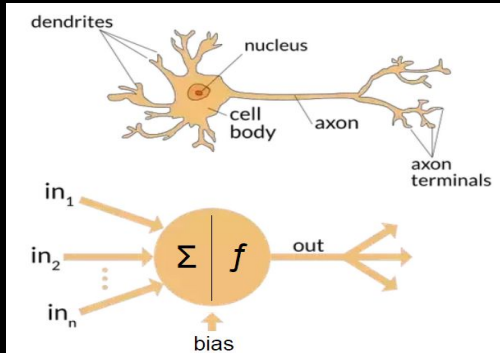


Logistic Regression

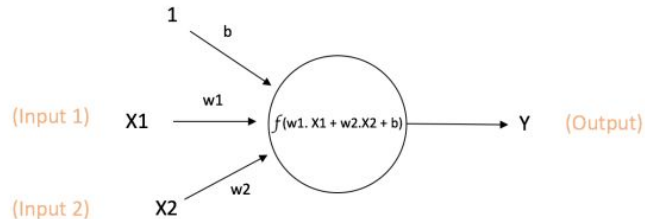


Neural Network

Neuron

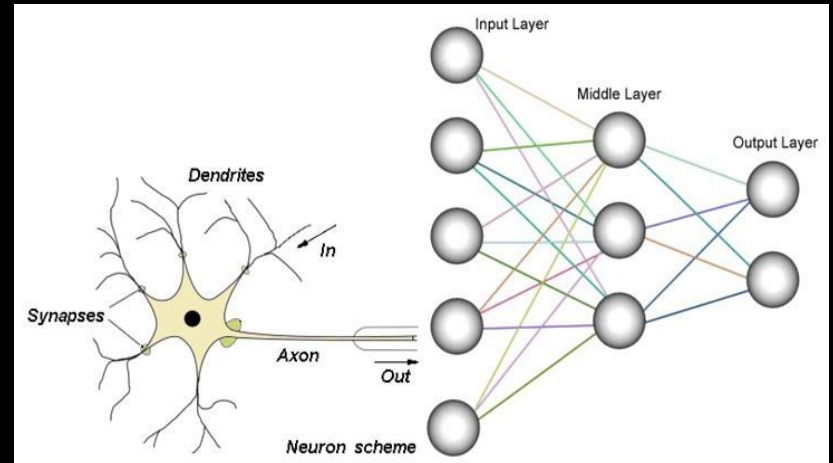


Linear Unit/Perceptron



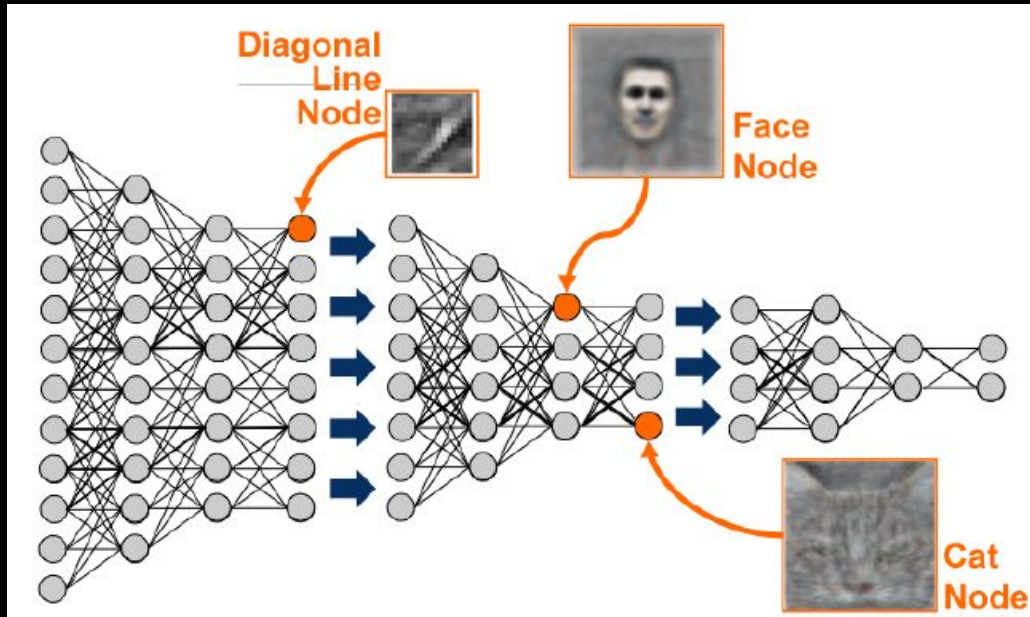
$$\text{Output of neuron} = Y = f(w1 \cdot X1 + w2 \cdot X2 + b)$$

Neural Network: Multi-Layer Perceptrons



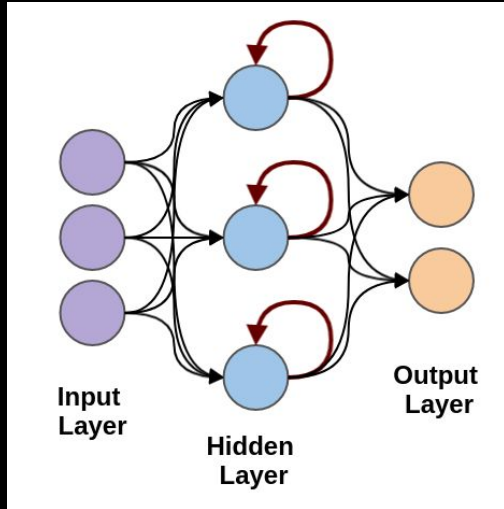
Deep Neural Networks

Fancy new term for Multi-layer Neural Networks with efficient ways of training



Recurrent Neural Network

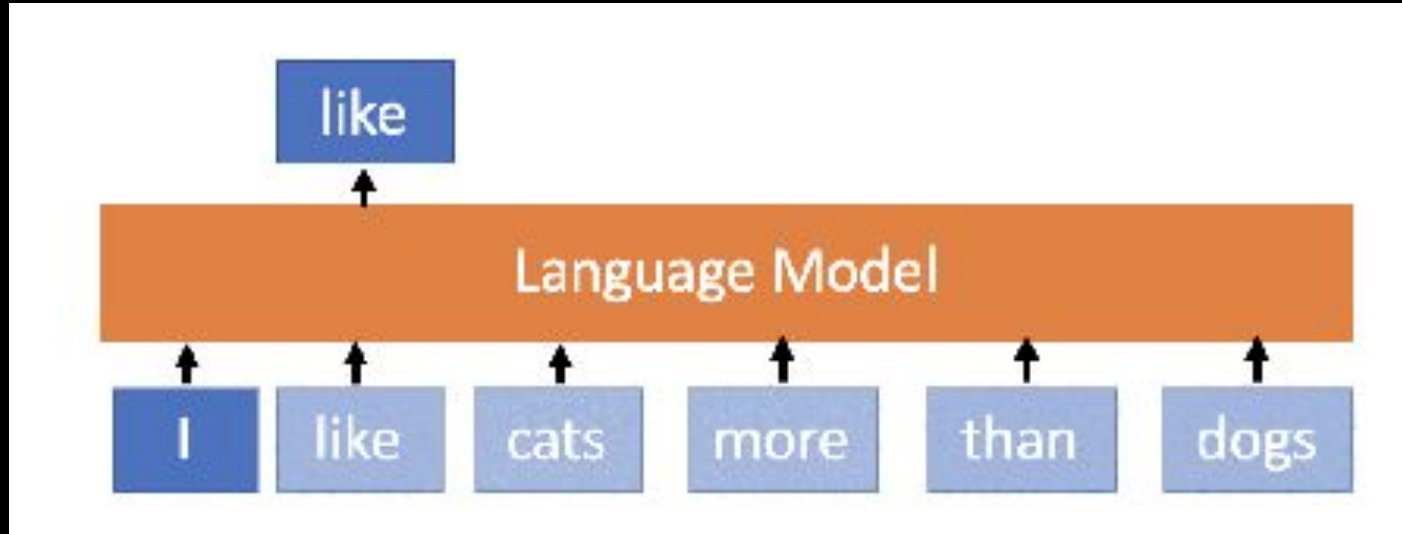
Fancy RNN Architecture



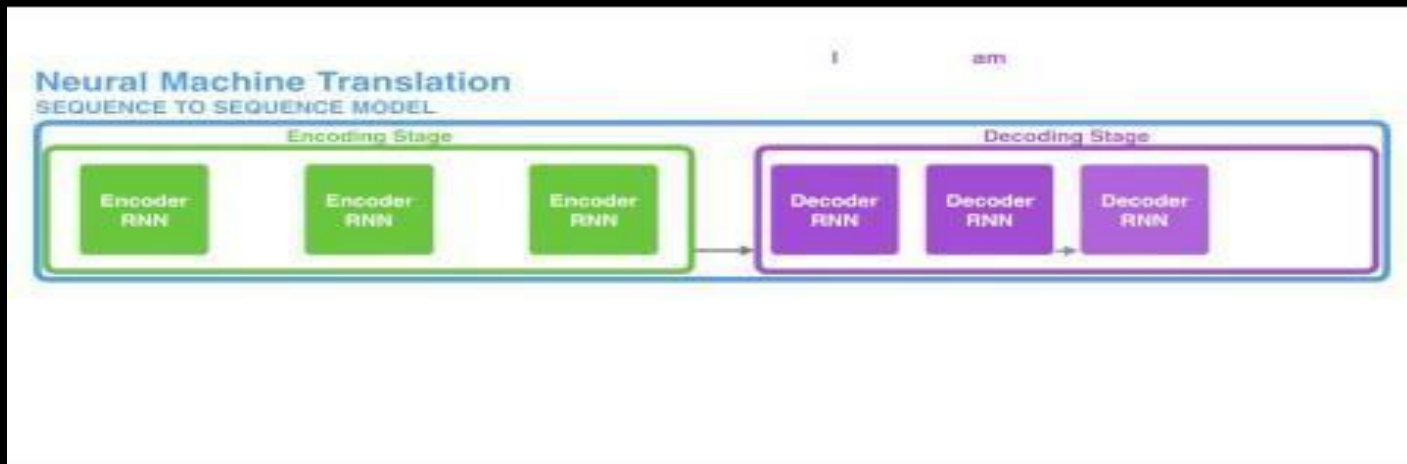
Reality: Memory, i.e. keep track of previous states for future prediction



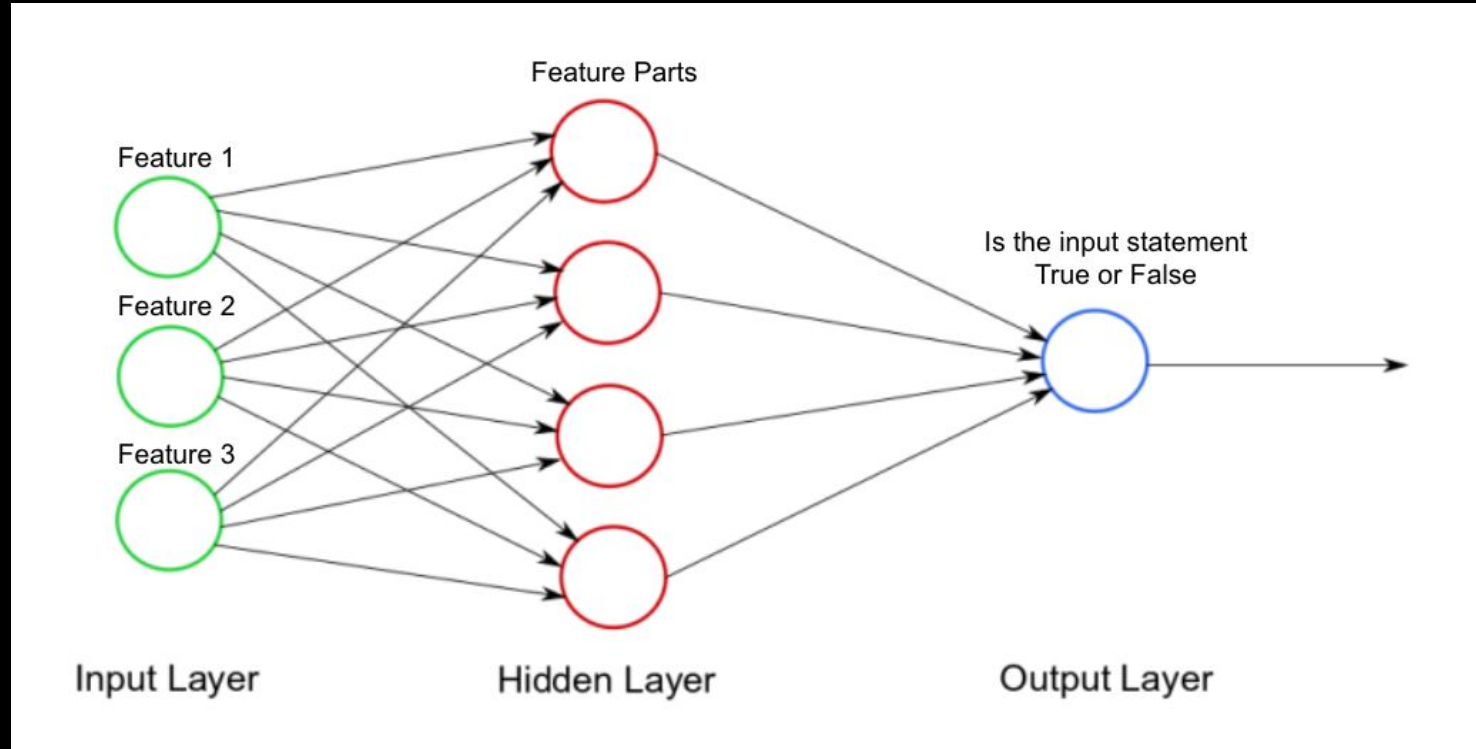
Language Models



Seq2Seq Illustration



Neural Network Features and Labels (Annotations)



Traditional NLP Featurization

Vocabulary Lookup Table

Sentence	the	is	visited	...	president	great	...	obama
Obama is the president of the us	1	1	0	...	1	0	...	1
Obama visited the great wall of china	1	0	1	...	0	1	..	1

One-hot encoding

Input/Sentence:

Obama is the president of the us

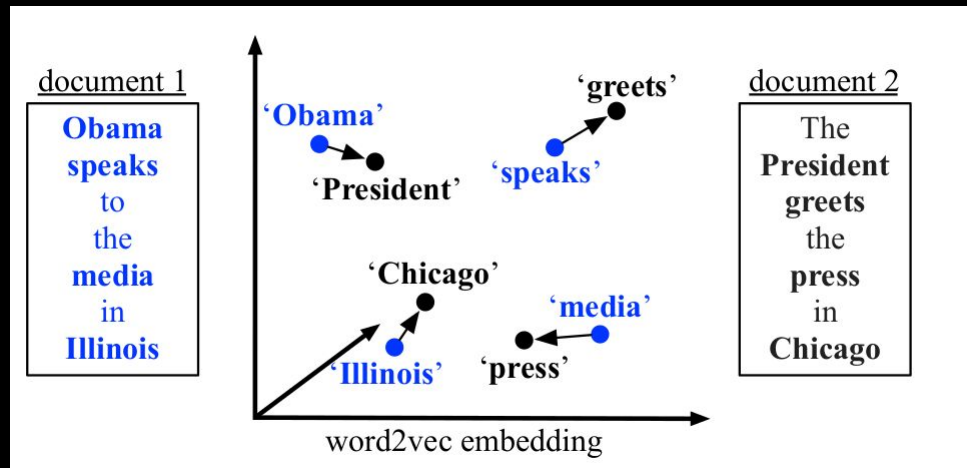
Features: Index of words appearing in the sentence

Vectorization: Initialize a vector of dimension V with all zeros except the ones appearing in the sentence
[1, 1, 0, 0, 1, 0, 1, 0, 0, 1,, 1]

This vector is the **feature representation** of the given sentence

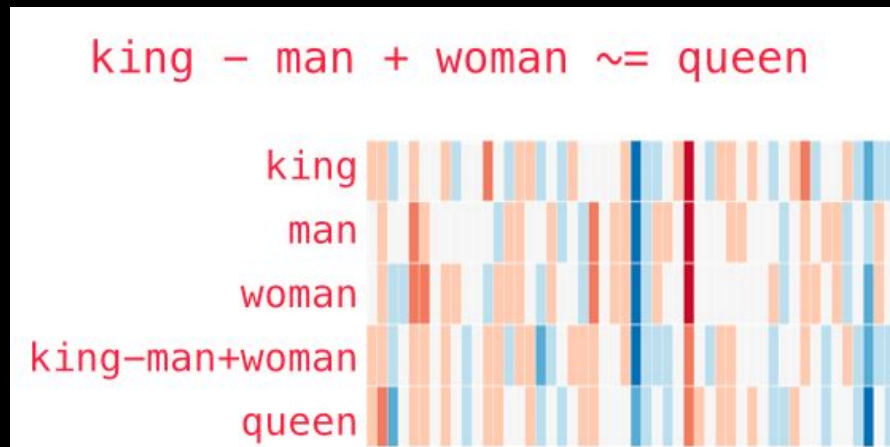
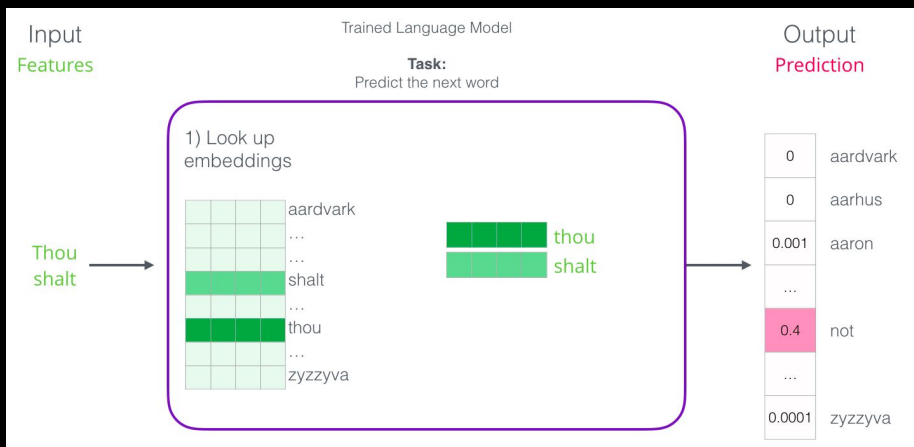
Word Embeddings: Semantic Repres. of Language Units

Word Embeddings



word	Dim 1	Dim 2	Dim ...	Dim 300
obama	0.1	0.5	...	0.9
president	0.15	0.4	...	0.85
chicago	0.3	0.7	...	0.6
illinois	0.25	0.72	...	0.58
the	0.01	0.02	...	0.07
...				
press

Word2Vec Illustration



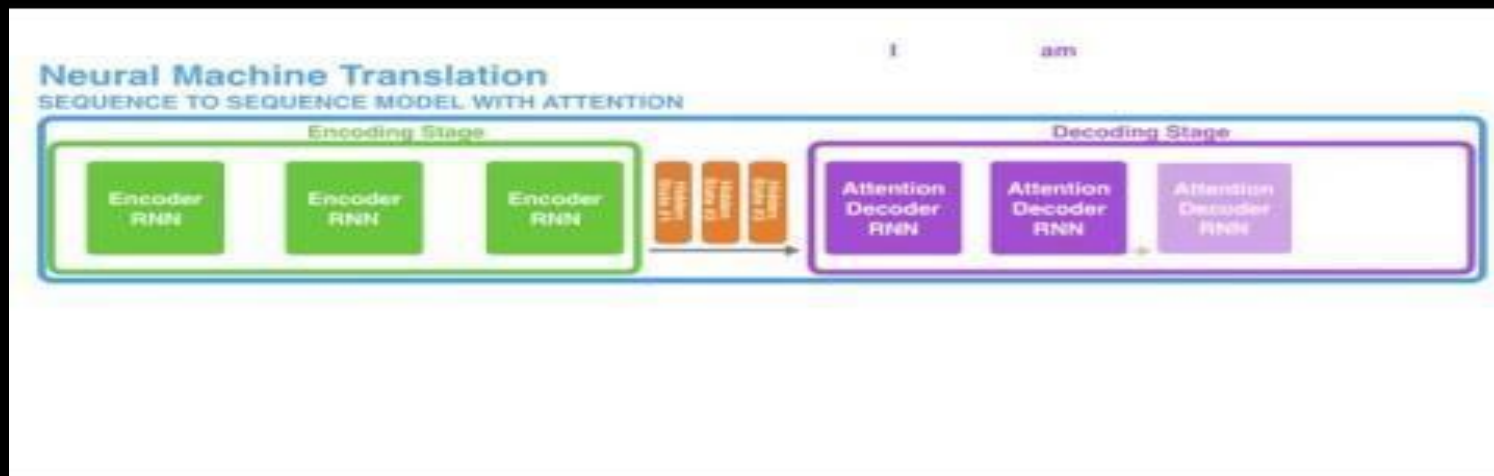
Are Word Embeddings contextual?

Transformers: Attention is all you need!

What are Transformers? Key Ideas

- Attention based Neural Networks
- Contextual Embeddings using:
 - Word position
 - Dynamic attention
- Auto-Regressive:
 - Use output generated at time “t” to generate at “t+1” time
- More Scalable:
 - Parallelizable: Vertical Architecture [RNN - Horizontal]
 - Faster, Deeper,

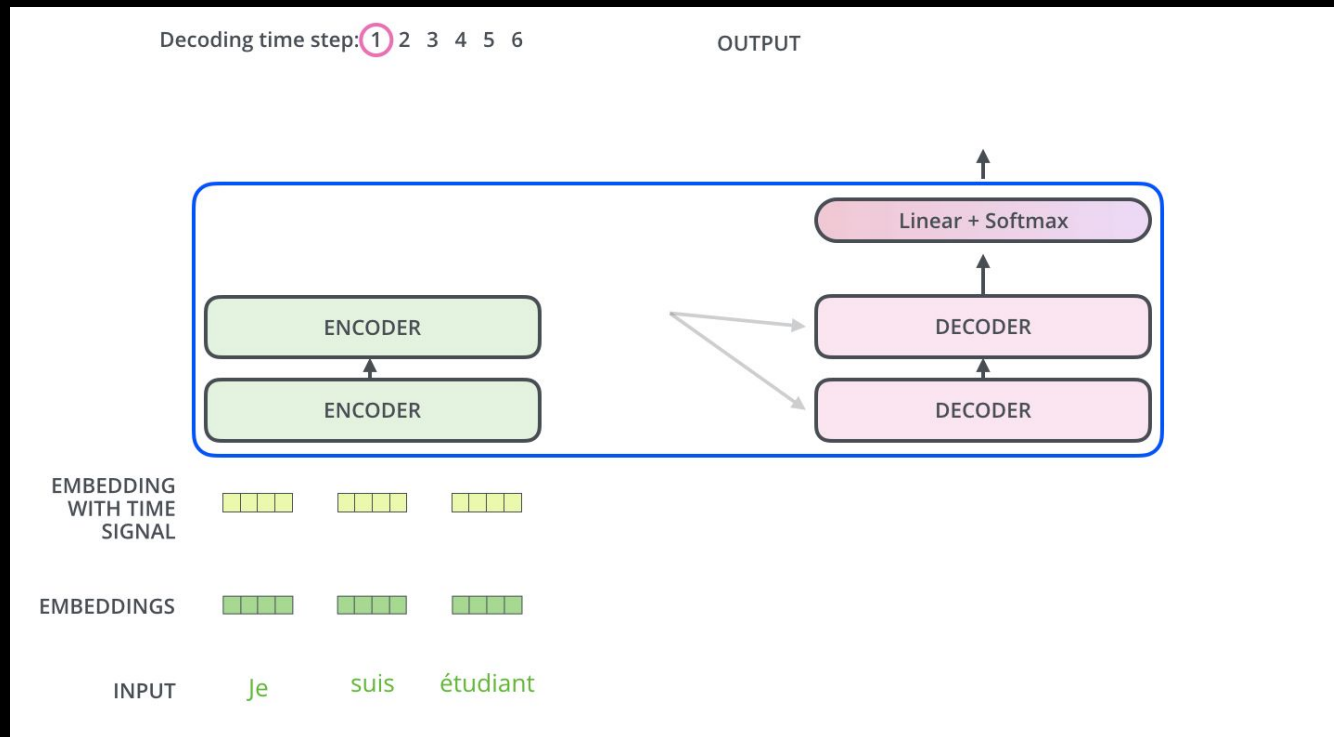
Seq2Seq with Attention



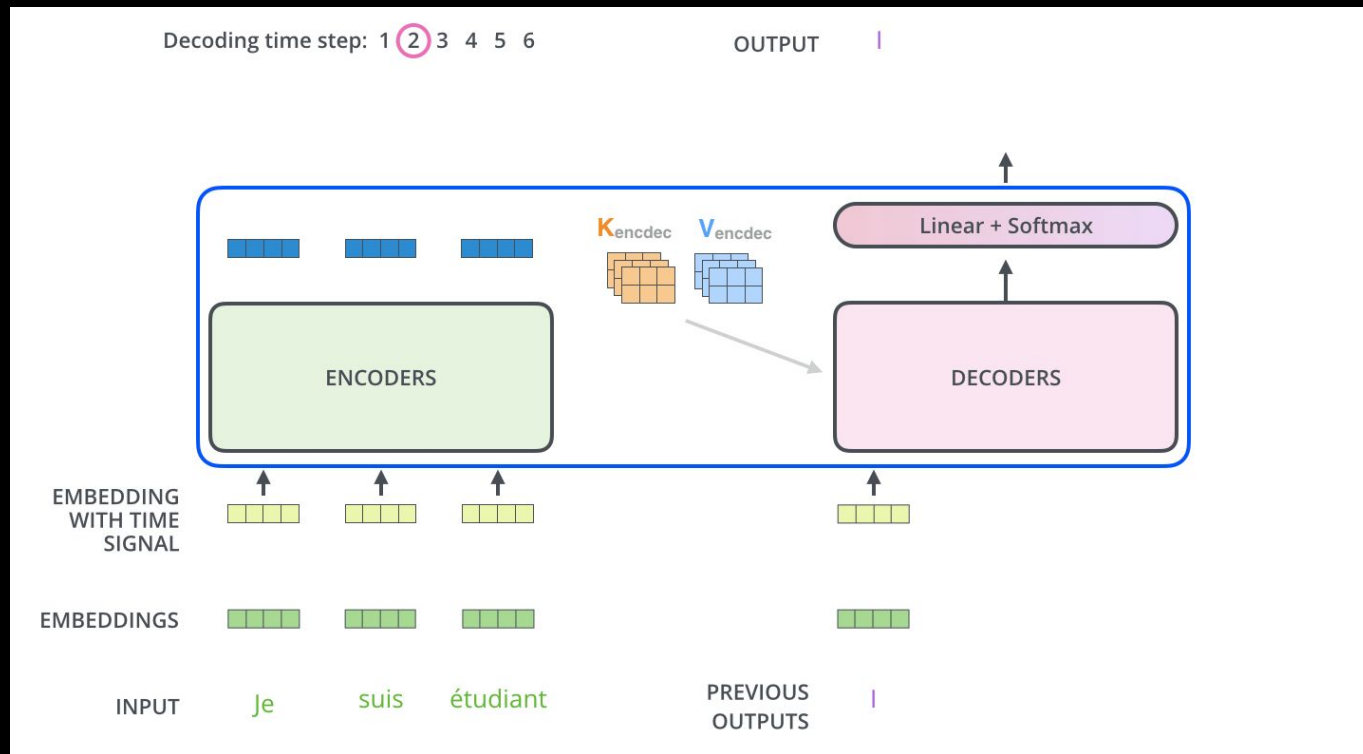
Attention in Neural Networks

. la maison de Léa <end> .

Dynamic Attention



Transformers Language Models



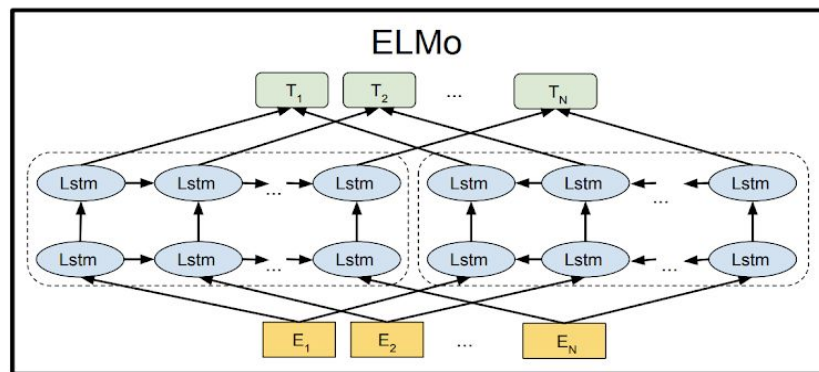
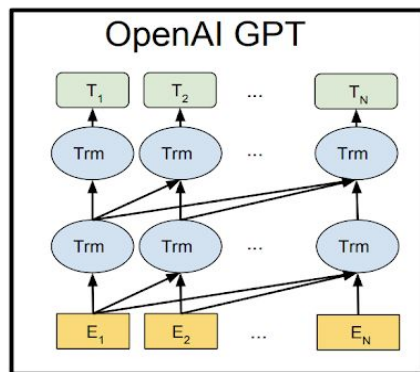
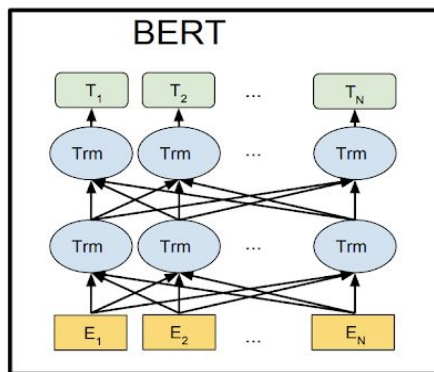
Autoregressive Language Models

There are a million ways

Contextual Word and Sentence Embeddings: For better NLU/NLP

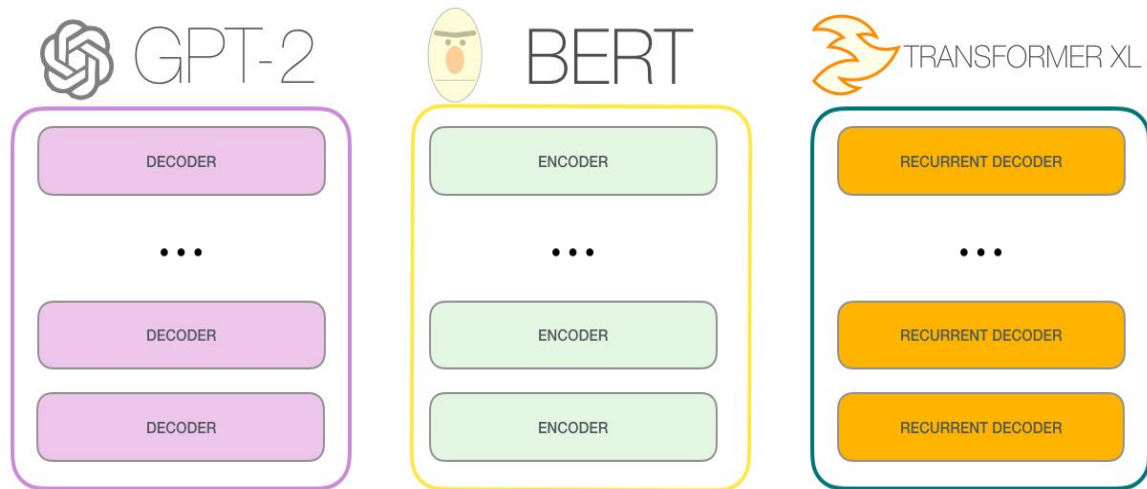
Transformers

LSTMs/RNNs



Transformer, a Neural Architecture, which learns **contextual word and sentence embeddings**. Trained with 100s of millions, billions of parameters

BERT (RoBERTa) vs GPT (1, 2, 3) vs Transformer XL



Encoders

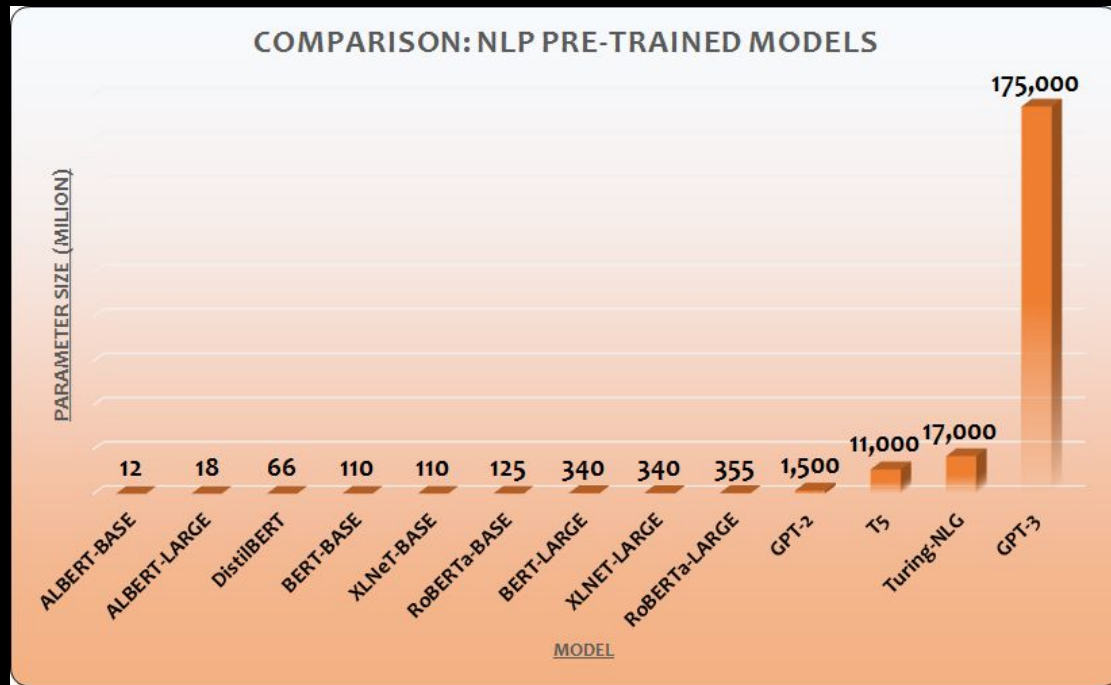
- Usually bi-directional
- Great for getting good input representations
- Classification, Tagging, Extractive Summarization

Decoders:

- Uni-directional (forward)
- Great for generative tasks
- NLG, Abstractive Summarization, Translation, QA

Lots of Transformers!!

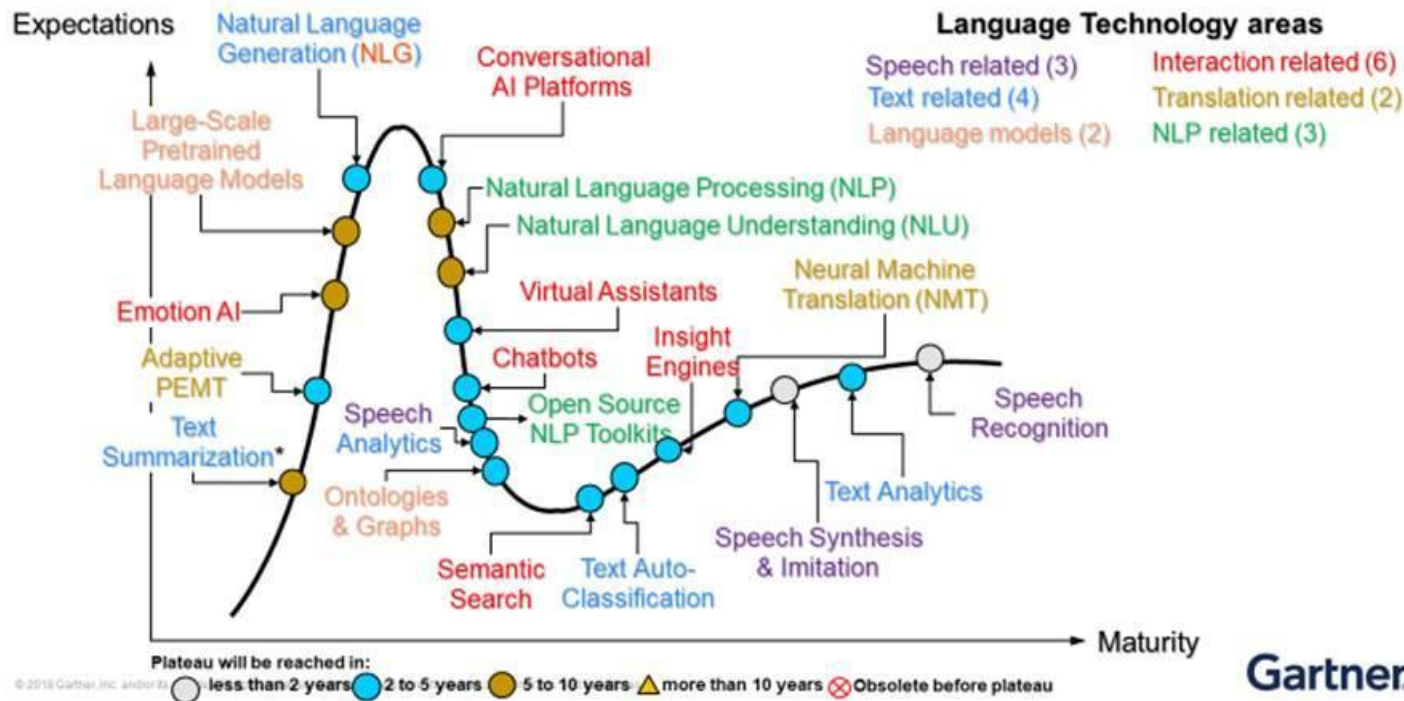
Which one to use?



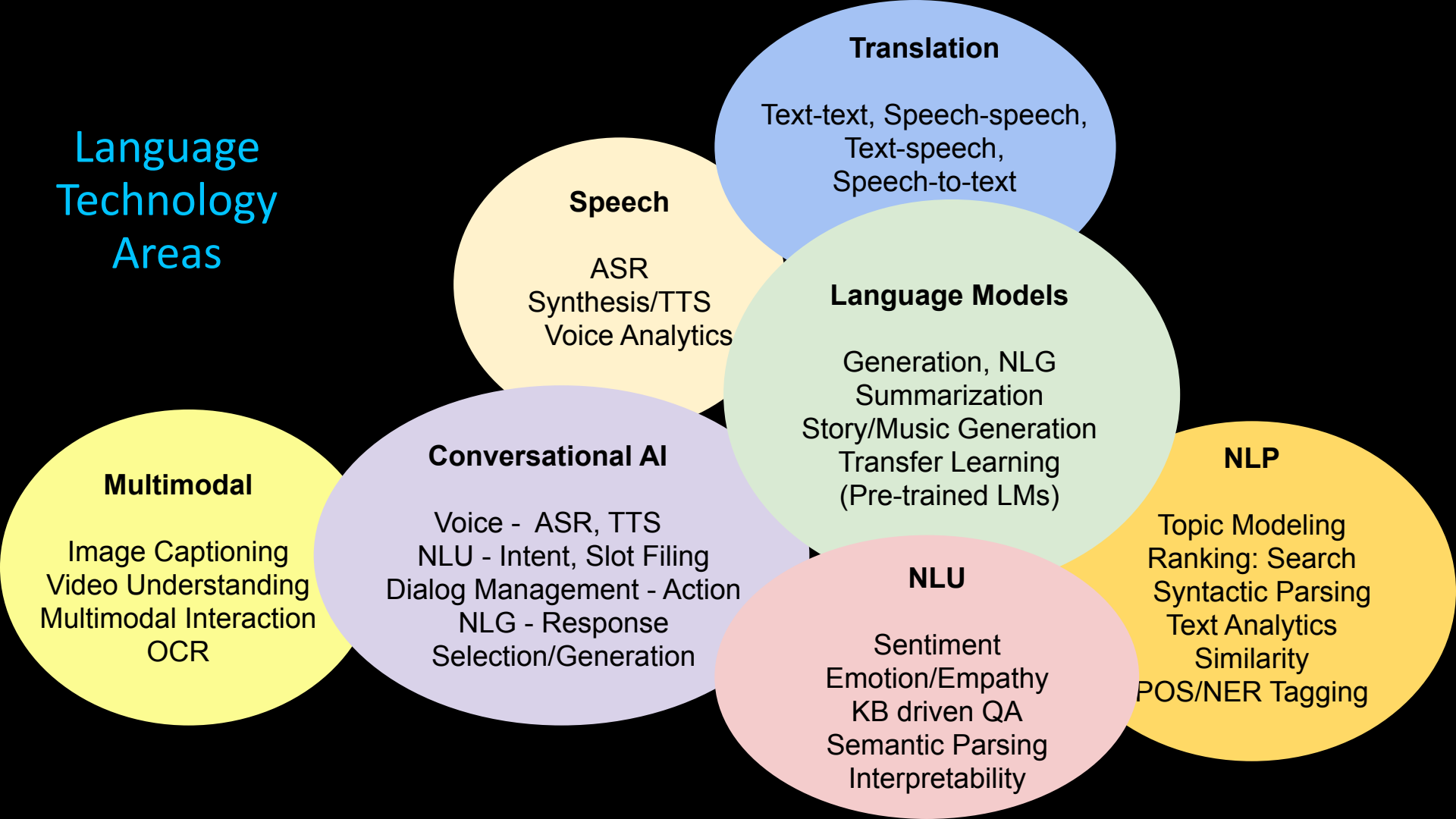
T5XXL: 1,700,000

Applications

Hype Cycle for Natural Language Technologies, 2020

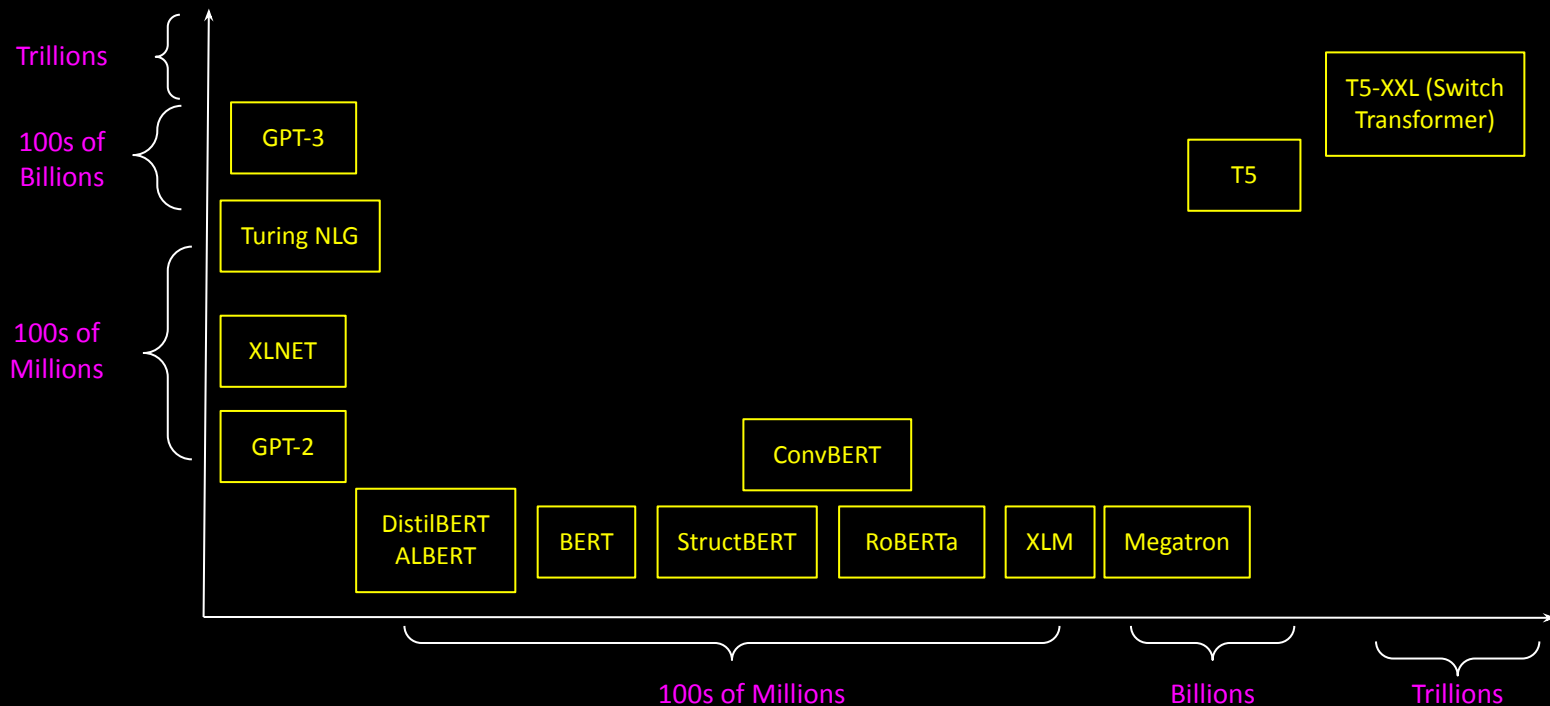


Language Technology Areas



Popular Transformer Architectures & their Applications

Decoders: Better for Generation Applications



Encoders: Better for Classification, Ranking Applications