# IS813: Gen AI - Implementation

Mohannad Elhamod

# Neural Nets in Language Modeling

**Continued…**

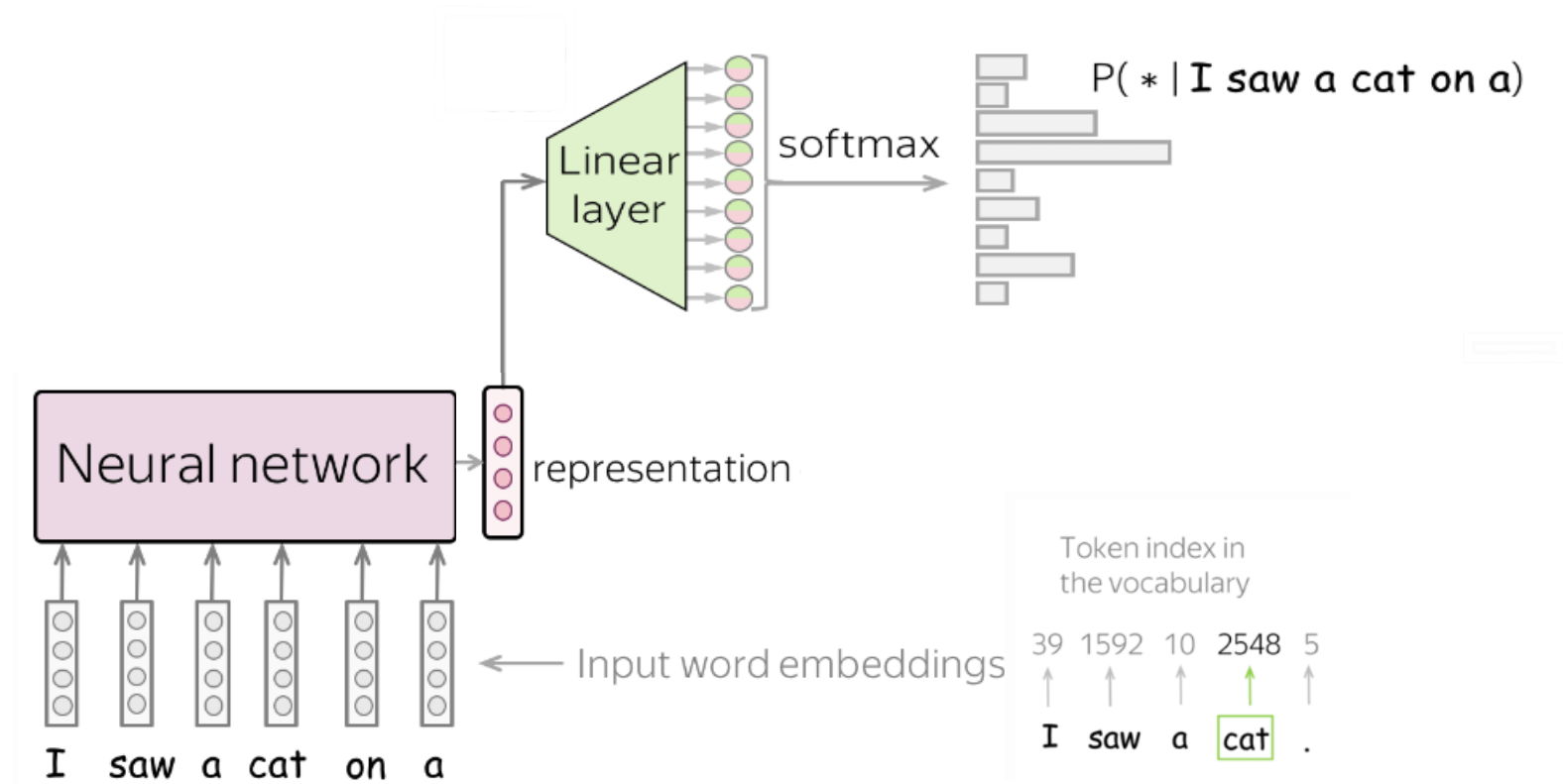**Boston University** Questrom School of Business

# Fast Forward...

As neural networks arrived at the scene, they were utilized for language modeling.

- N-grams look for exact prefixes, which is limiting…
- However, neural networks can learn more interesting relationships between the words.

Example: All humans are mortal. Socrates is a human. Therefore,

Socrates is mortal.

# General Model Architecture

Can you see any issue with inputting words in an NN?

**Boston University** Questrom School of Business
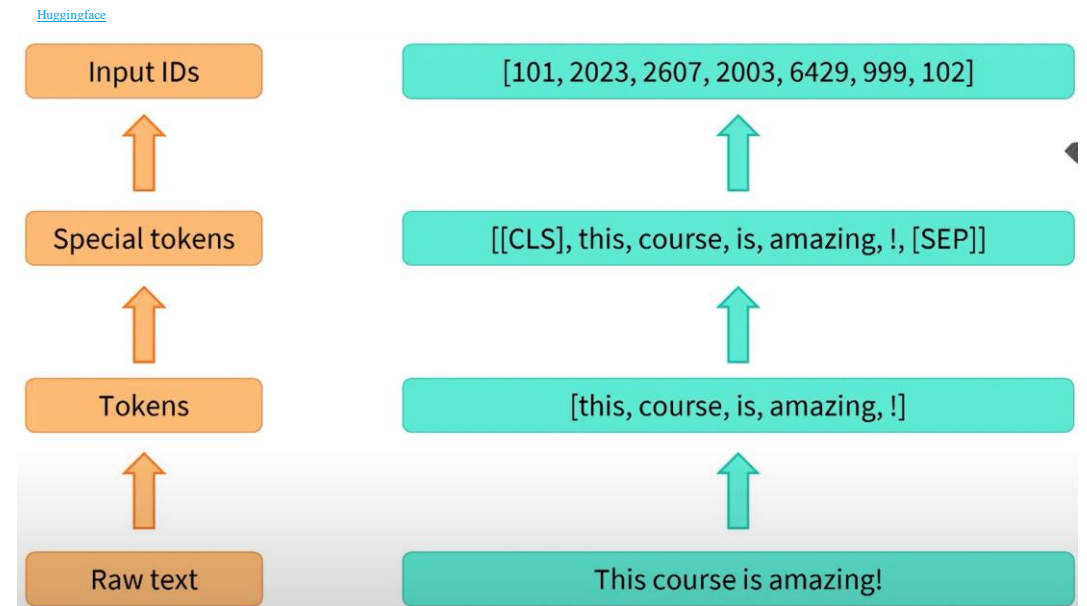
# What is an embedding?

- It is the numeric representation of data.

- [Example for images.](#)

# Word Embeddings

- We ideally want related words (i.e., *similar* meanings) to have smaller distances.

- Demo

- Examples:
  1. Word2Vec (Google)
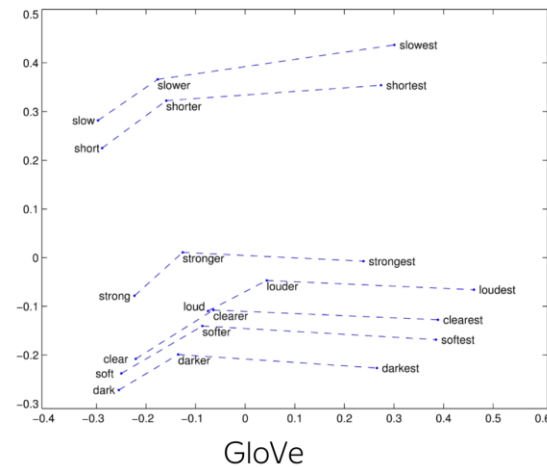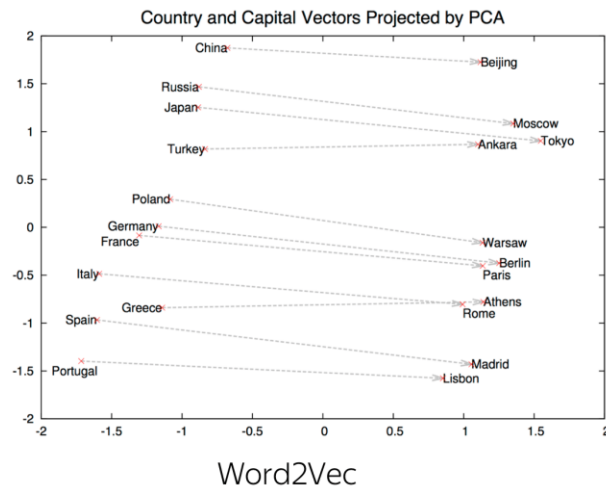  2. GloVe (Stanford)
  3. Train your own!

# Tokenization

- Computers only understand numbers.
- We need to convert the text into tokens (e.g., words).
- Each token can then be represented as a number.

Huggingface

| | |
|---|---|
| Input IDs | [101, 2023, 2607, 2003, 6429, 999, 102] |
| ↑ | ↑ |
| Special tokens | [[CLS], this, course, is, amazing, !, [SEP]] |
| ↑ | ↑ |
| Tokens | [this, course, is, amazing, !] |
| ↑ | ↑ |
| Raw text | This course is amazing! |

**Boston University** Questrom School of Business

# Word Embeddings
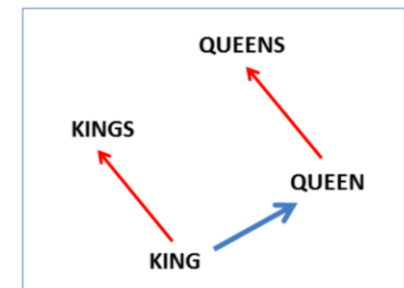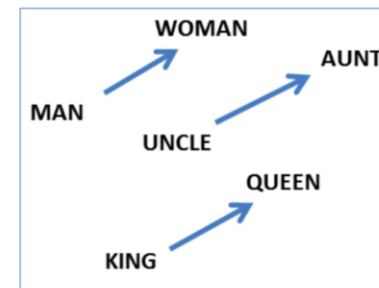
Since word embeddings carry *meaning*, certain directions in their space carry certain significance:
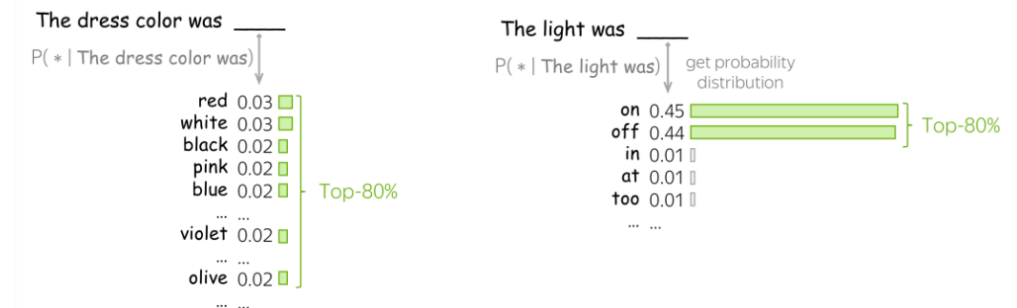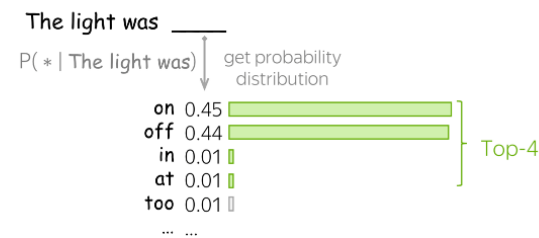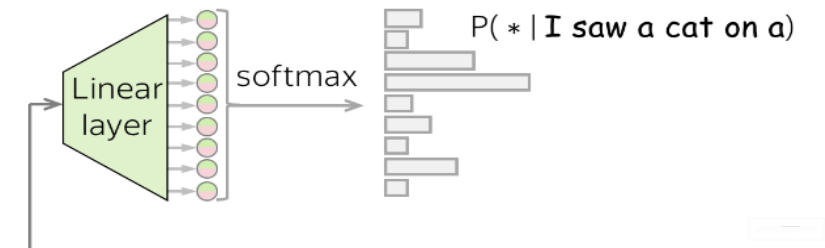
- Demo (dimensionality)



semantic: $v(king) - v(man) + v(woman) \approx v(queen)$

syntactic: $v(kings) - v(king) + v(queen) \approx v(queens)$

Word2Vec

GloVe

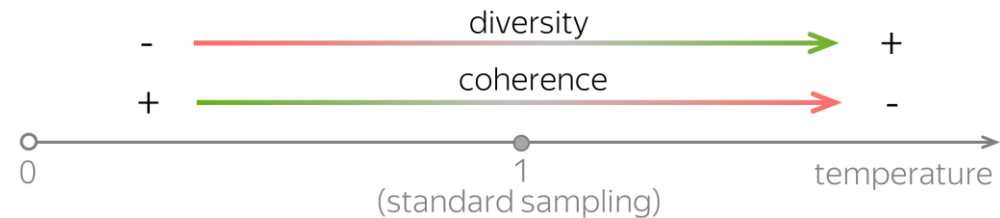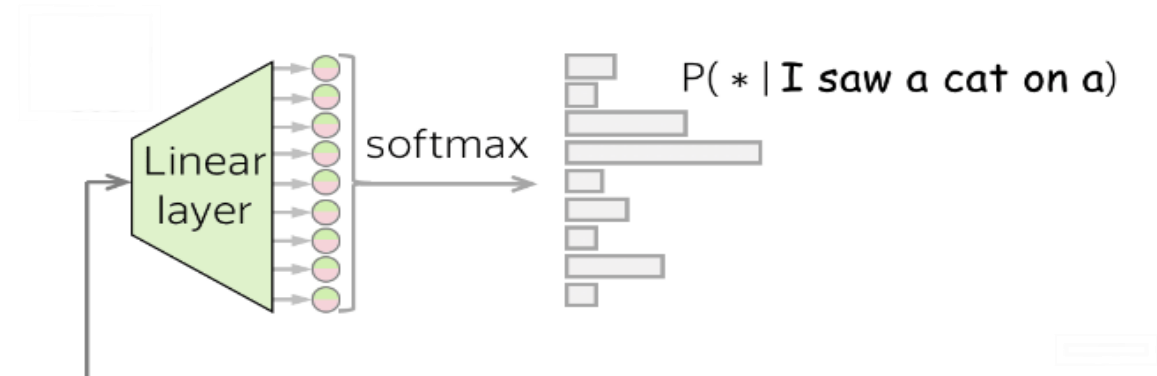# Sampling The Distribution



- ## Always take top probability?
  - ### That makes the model deterministic (no creativity).

- ## Alternative?
  - ### Top-k or top-p.



Lena-voita

# Sampling The Distribution

- Some words have way higher probability than others.
- This can be manually tuned through temperature.
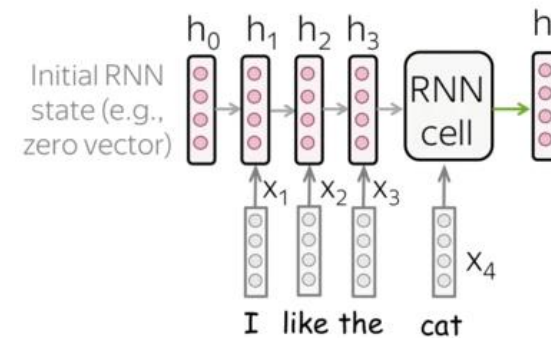- Demo



Lena-voita

# Measuring The Metric

- What are we looking for?
  - A model that is not surprised by the new text it seen.

- We use perplexity.
  - Takes values between 1 and number of possible tokens.
  - Smaller is better.
  - Demo

# Fast Forward...

- There exists many types of Neural Nets for language modeling:
  - CNNs
  - RNNs
  - LSTMs…

- Generally, Neural Nets learn an *embedding* that represents the _entire prefix_ to predict the next word.

Lena-voita

# Attention!

**Dzmitry Bahdanau**
Jacobs University Bremen, Germany

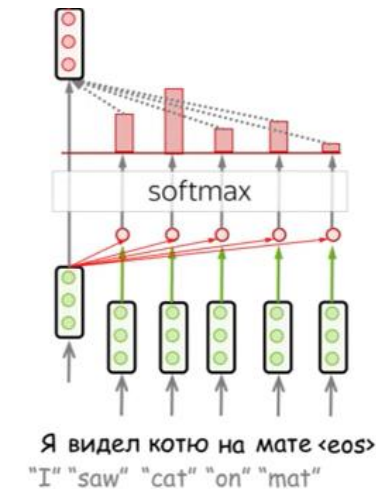**KyungHyun Cho**    **Yoshua Bengio**[*]
Université de Montréal

- These types of Neural Nets, however, suffered from various issues:
  - E.g., *catastrophic forgetting*, where earlier context in longer sentences tends to be forgotten.

- In 2015, *attention* in Neural Nets was invented:
  - It allowed models to attend to different parts of the sentence (instead of a single representation).



softmax

Я видел котю на мате <eos>
"I" "saw" "cat" "on" "mat"

Lena-voita

**Boston University** Questrom School of Business
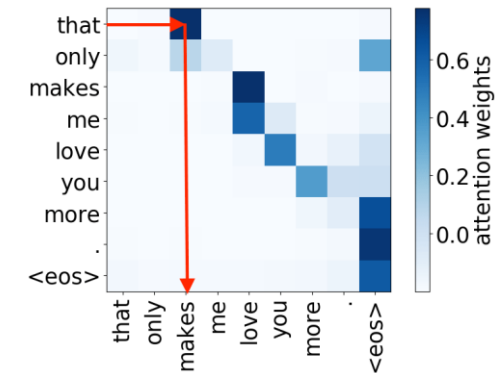
# Attention!

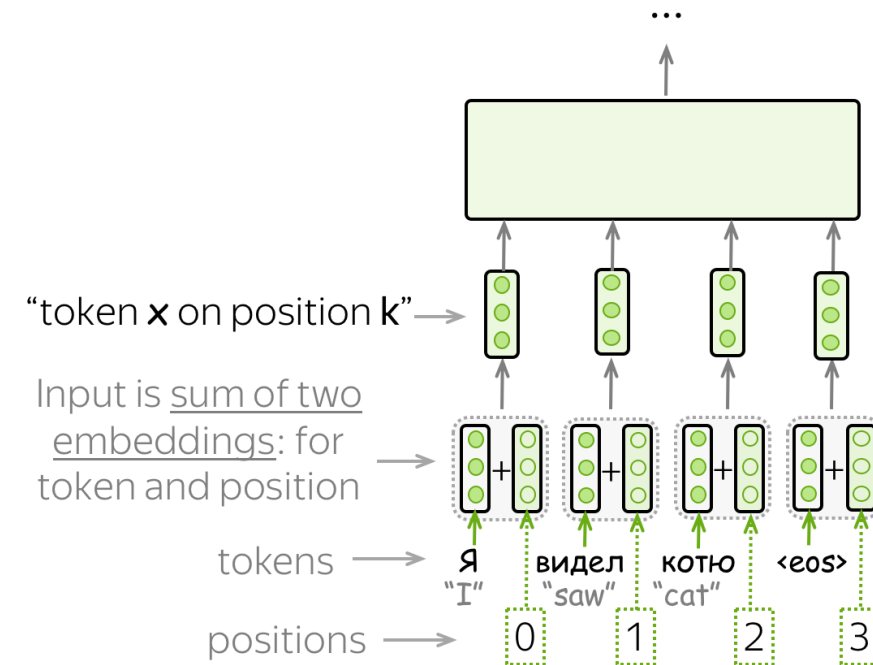- Once each part has its own embedding, different types of *relationships* can be learned!
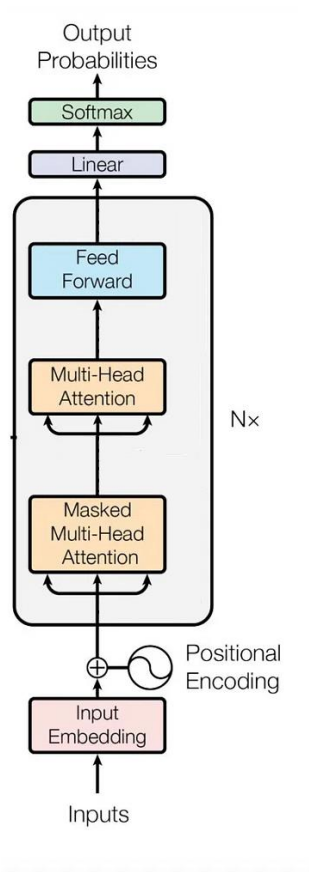


Subject -> verb

Lena-voita

# Order Matters: Positional Encoding!

- Since token embeddings do not contain information about the location of the word, they should be combined with a *positional encoding.*

"token **x** on position **k**" →

Input is <u>sum of two embeddings</u>: for token and position →

tokens →   Я        видел      котю      ‹eos›
           "I"      "saw"      "cat"

positions →    0      1       2      3

Lena-voita

# The Transformer is born!

# Models in the wild

# Model Types

We are not going to get into technical details, but certain models may be more fit for certain tasks:

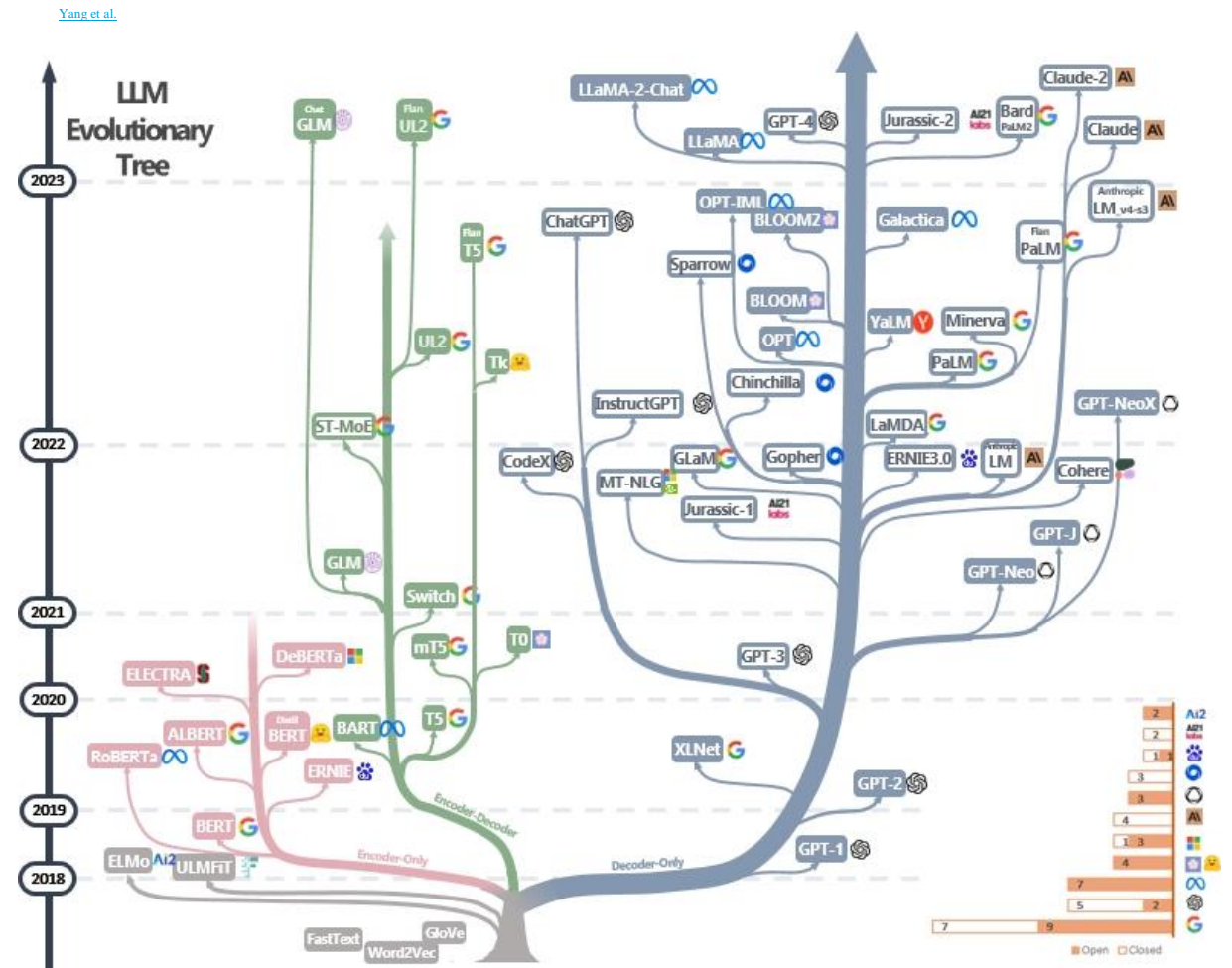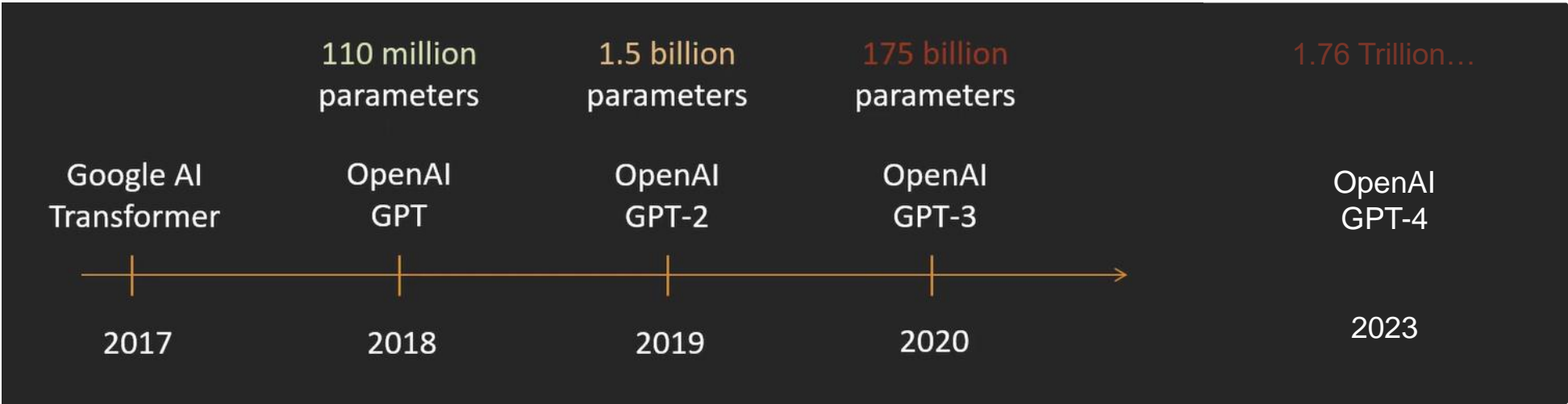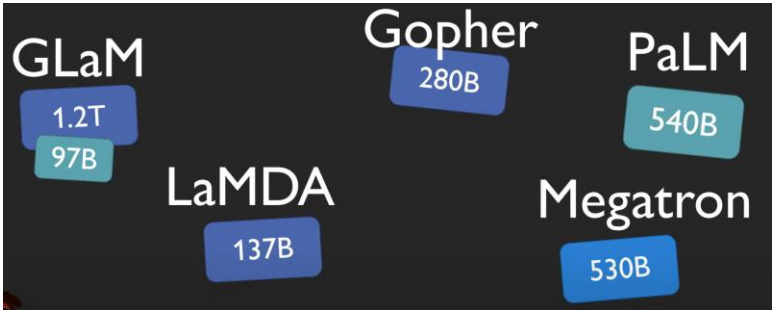| Model | Examples | Tasks |
|---|---|---|
| Encoder | ALBERT, BERT, DistilBERT, ELECTRA, RoBERTa | Sentence classification, named entity recognition, extractive question answering |
| Decoder | CTRL, GPT, GPT-2, Transformer XL | Text generation |
| Encoder-decoder | BART, T5, Marian, mBART | Summarization, translation, generative question answering |

Javinkarla

**Boston University** Questrom School of Business

# Why so many?

Where do the differences come from?

- Data.
- Model type and size.
- Hyperparameters (context size, embedding size,…).
- Training process (the cost function, fine-tuning, human feedback, etc.).



Yang et al.

# The GPT Evolution...

GLaM
1.2T
97B

Gopher
280B

PaLM
540B

LaMDA
137B

Megatron
530B

110 million parameters

1.5 billion parameters

175 billion parameters

1.76 Trillion...

Google AI Transformer

OpenAI GPT

OpenAI GPT-2

OpenAI GPT-3

OpenAI GPT-4

2017

2018

2019

2020

2023

AI Coffee Beak with Letitia

Book Corpus WebText

Common Crawl + ...

????

1,038 books (around 74M sentences and 1G words) of 16 different sub-genres (e.g., Romance, Historical, Adventure, etc.)

Over 240 billion pages. Petabytes of data.

**Boston University** Questrom School of Business

BOSTON UNIVERSITY

# The GPT Evolution...





AI Coffee Beak with Letitia

# Different model sizes



Jay Alammar

# Exploring Your Options

- [OpenAI model reference](#)


- [HuggingFace tasks](#)
- [HuggingFace models](#)

# How much training does it take?

# Pre-trained Models: Democratizing AI

- Most of us don't have the expertise, data, or resources to train anything close to these impressive large models.

- Instead:

  - *Zero-shot Learning:* We can use open-source models out-of-the-box, even though they have never seen our data before.

  - *Transfer learning/Fine-Tuning:* Can be used as a base for further training (e.g., if the training data is non-public legal documents).

# Example: Instruct LLMs



Coursera

# In-Class Work

**HuggingFace**

Boston University Questrom School of Business

# Resources

- [Meaning and calculation of perplexity.](#)

- [Video: LLMs vs The Brain](#)

- [Video: Deciding which pre-trained model to fine-tune](#)

**Boston University** Questrom School of Business