



Transformer-based
language models

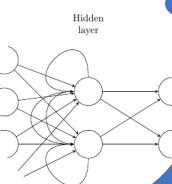
Where are we?



Bag of word, TF-IDF



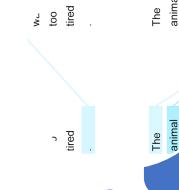
Word Embedding



Recurrent network (RNN, LSTM, GRU)

Sequence 2

Sequence

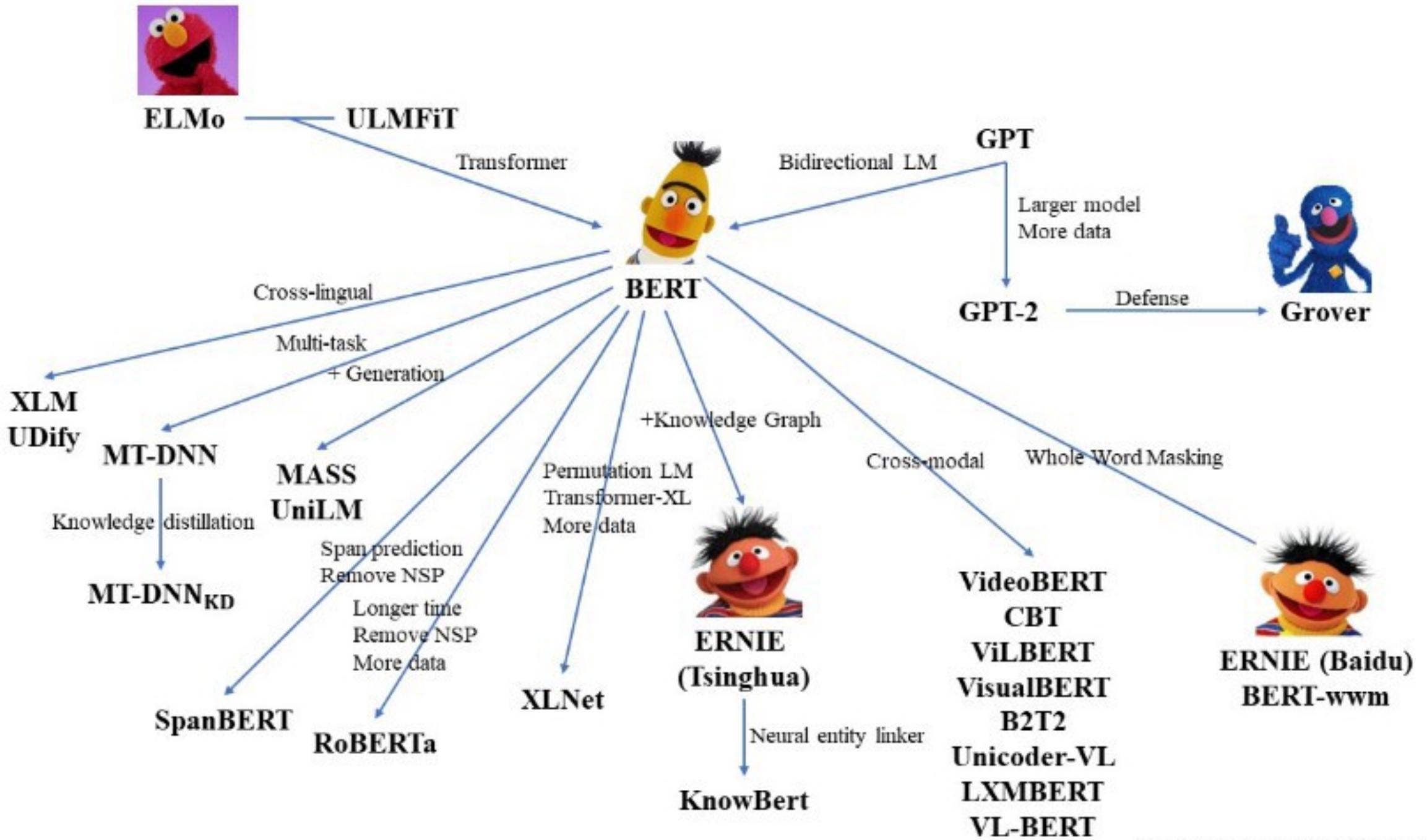


Pretrained Language Model (ULMFiT, BERT, XLNet, GPT)



What is Language Modeling?

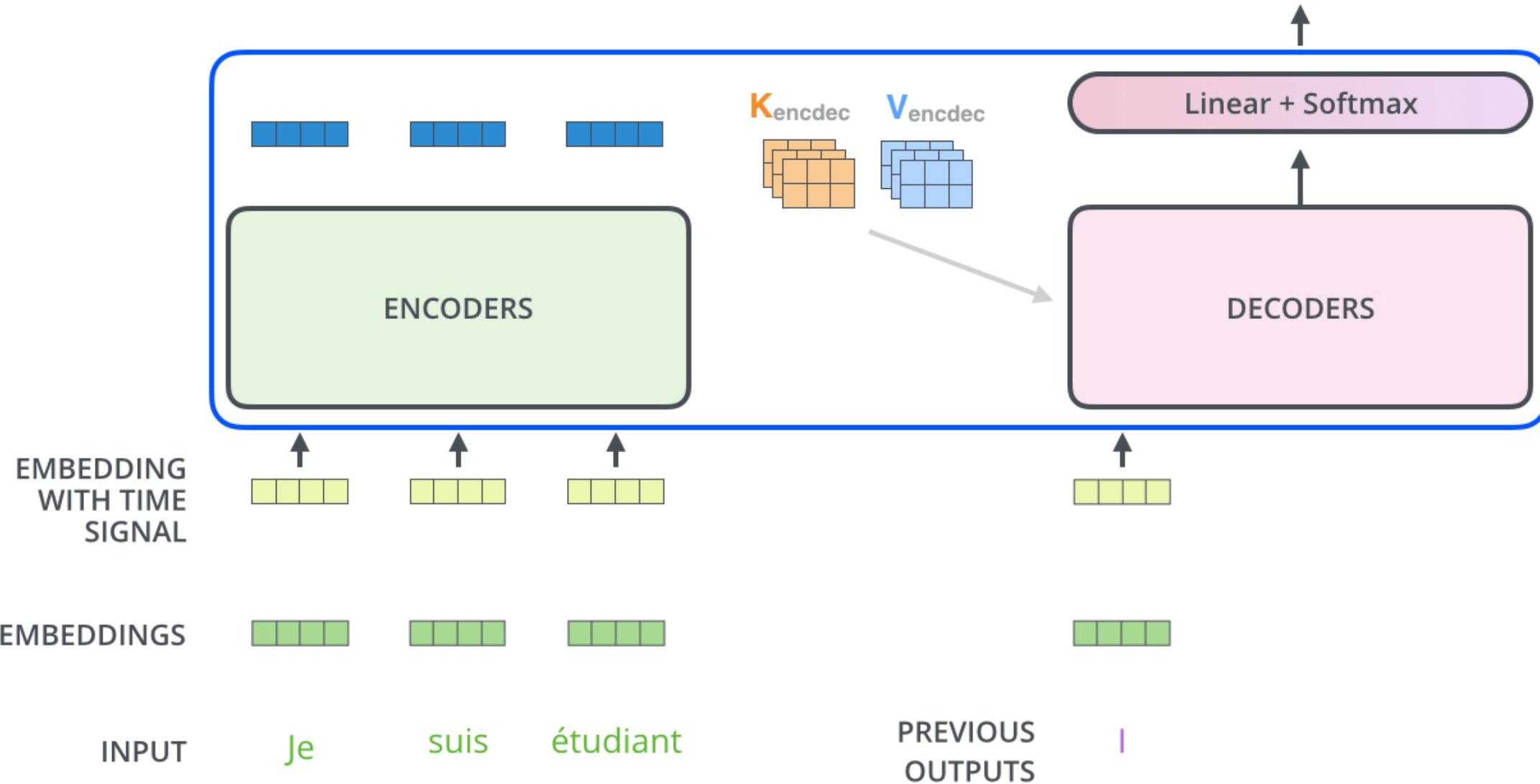
- Probability distribution over strings or text.
- $S = \text{"Tôi đi học"}$
- $P(S) = P(\text{"Tôi"}) \times P(\text{"đi" | "Tôi"}) \times P(\text{"học" | "Tôi đi"})$
- $P(S) > P(\text{"đi học Tôi"})$



Decoding time step: 1 2 3 4 5 6

OUTPUT

|



BERT

- Bidirectional Encoder Representations from Transformers.
- Use the Transformer Encoder architecture.
- Introduced in 2018 by Google AI.

~~un~~supervised

1 - ~~Semi-supervised~~ training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

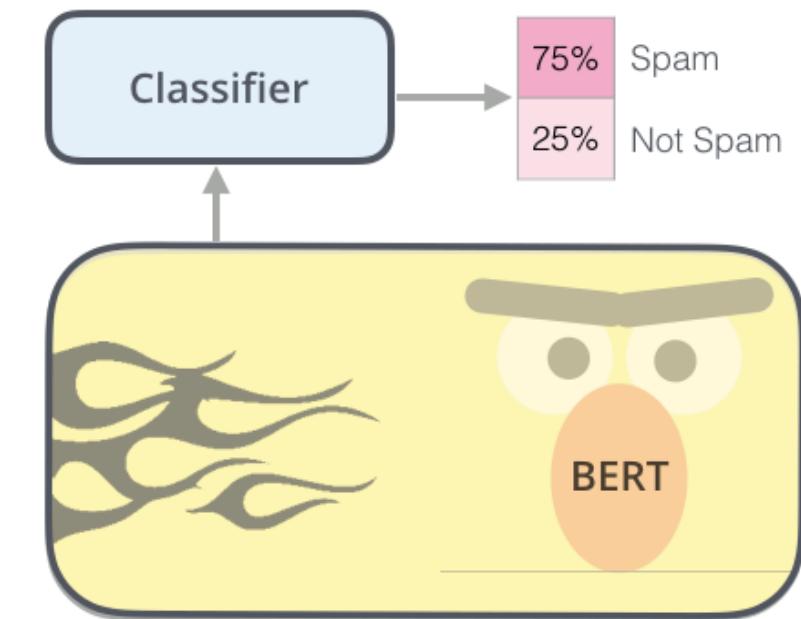
Semi-supervised Learning Step



Predict the masked word
(language modeling)

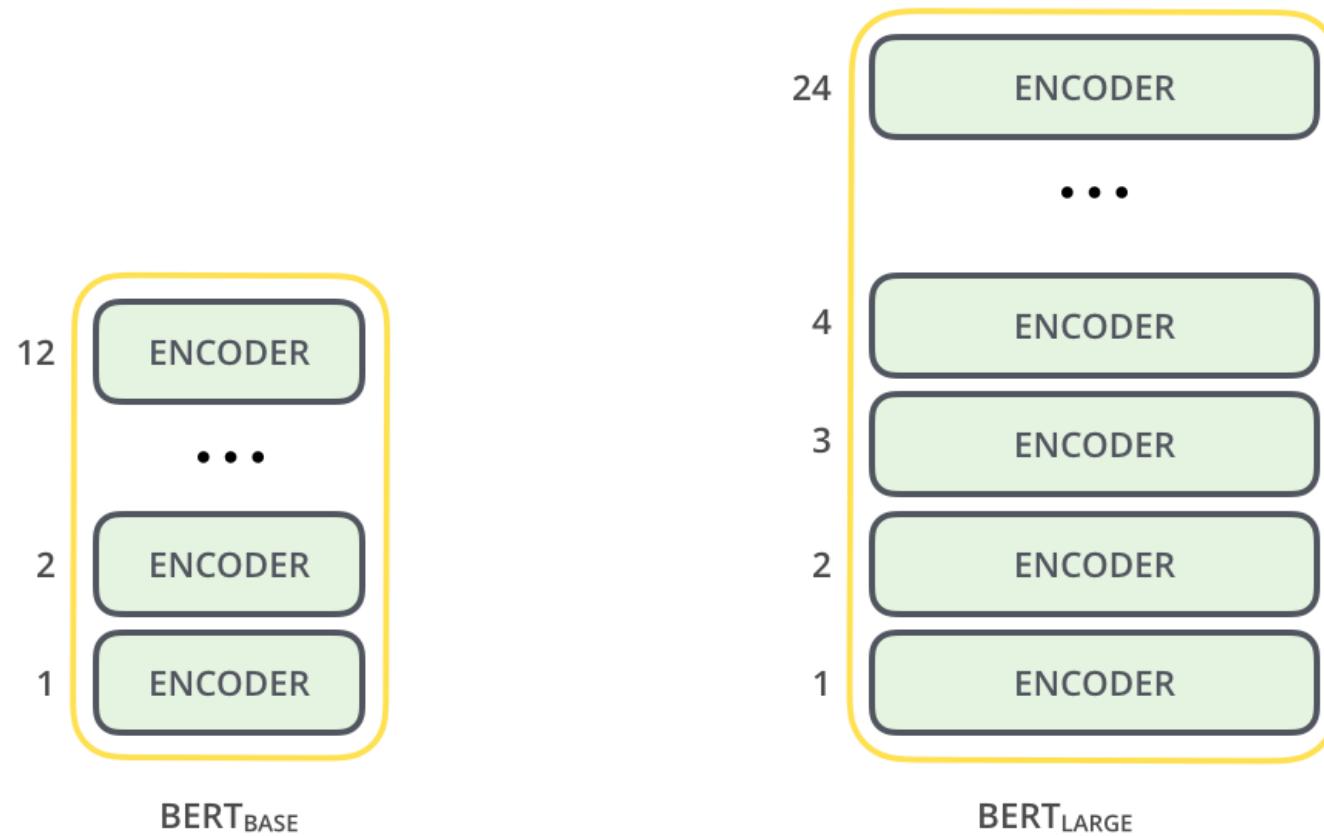
2 - Supervised training on a specific task with a labeled dataset.

Supervised Learning Step



Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

Architecture



Pretraining

- Two unsupervised tasks:
 1. Masked Language Model
 2. Next Sentence Prediction

Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

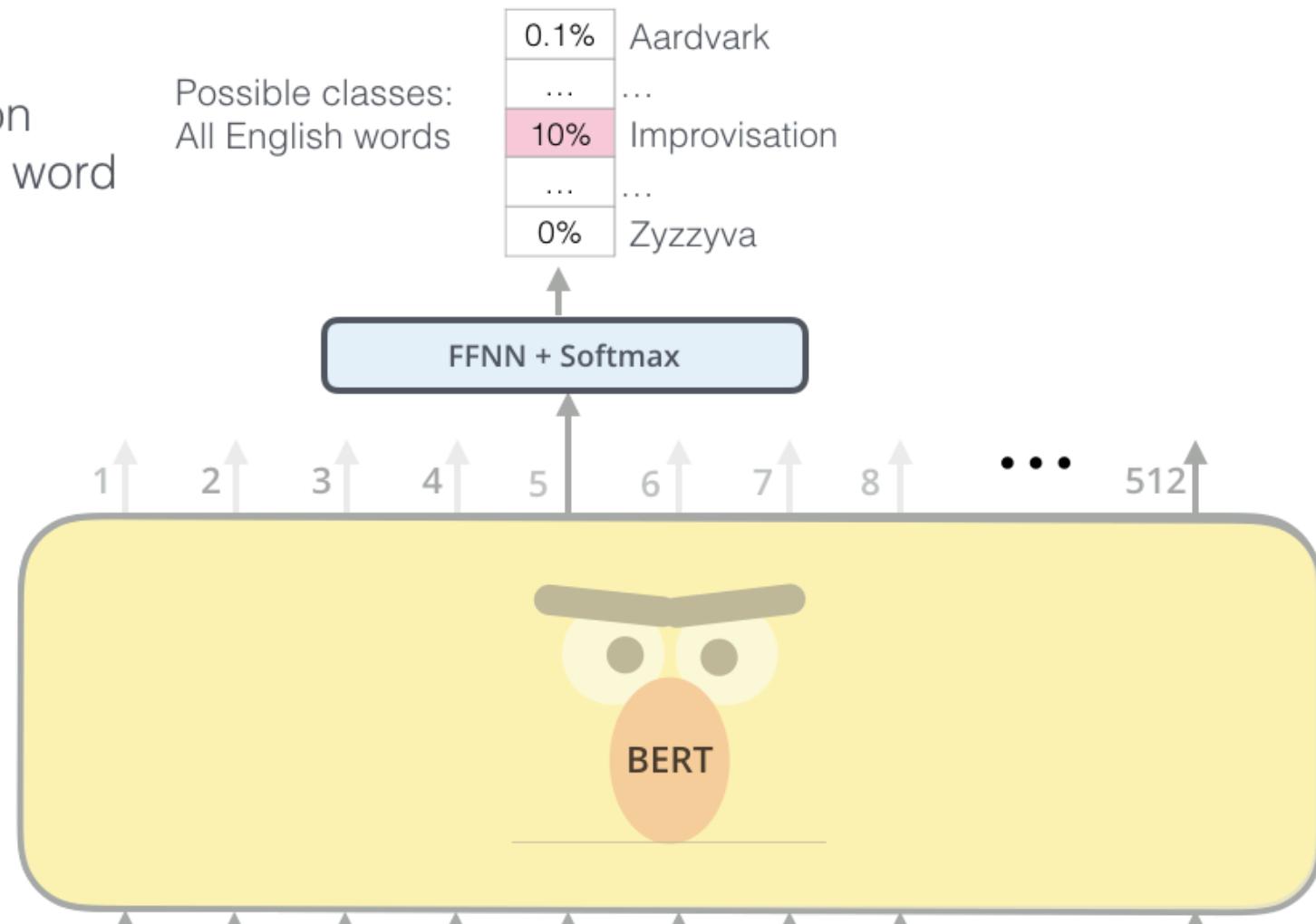
0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zzyzyva

Randomly mask
15% of tokens

1 [CLS] 2 Let's 3 stick 4 to 5 [MASK] 6 in 7 this 8 skit ... 512

Input

[CLS] Let's stick to improvisation in this skit



Predict likelihood
that sentence B
belongs after
sentence A

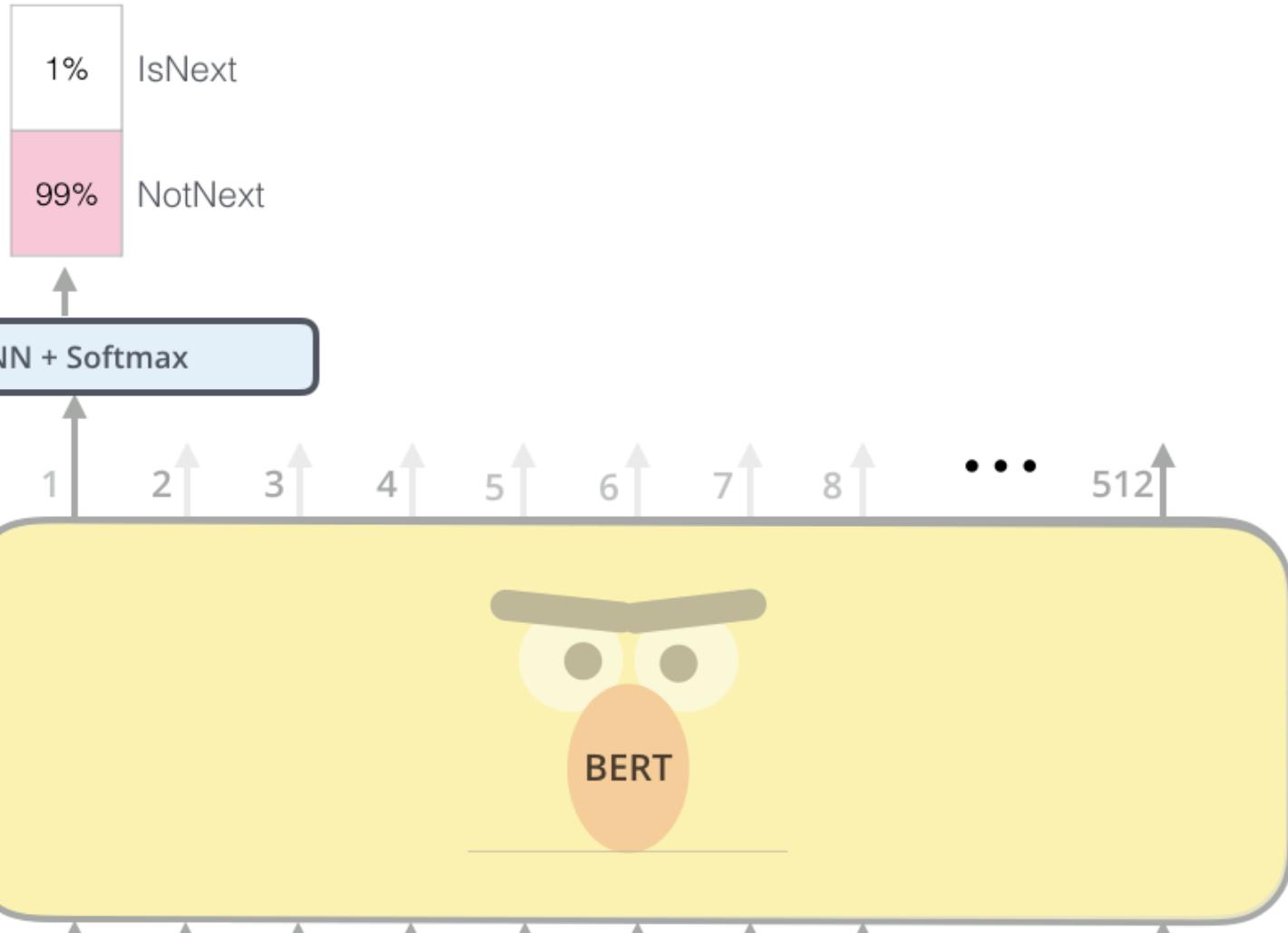
Tokenized
Input

Input

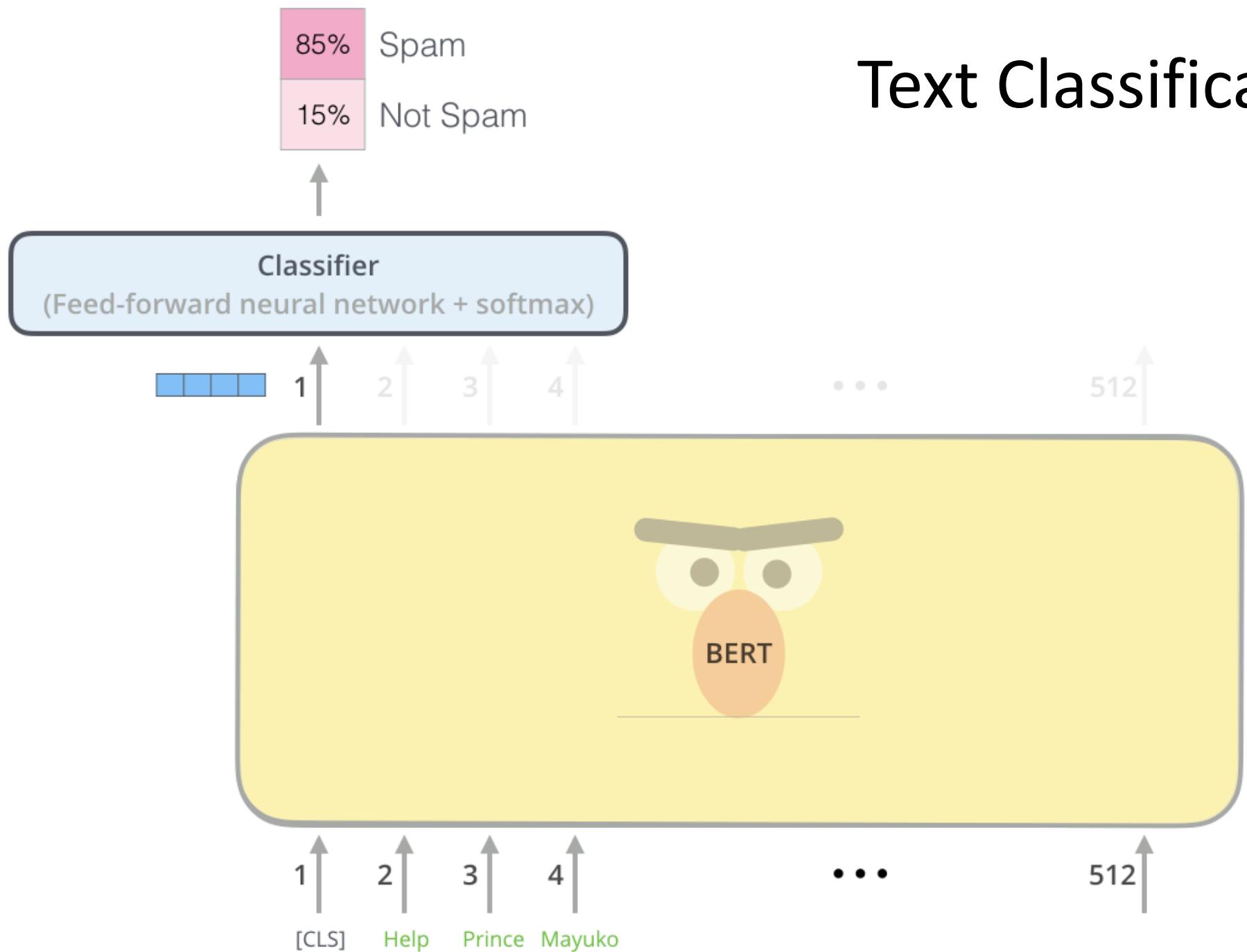
[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Sentence A

Sentence B



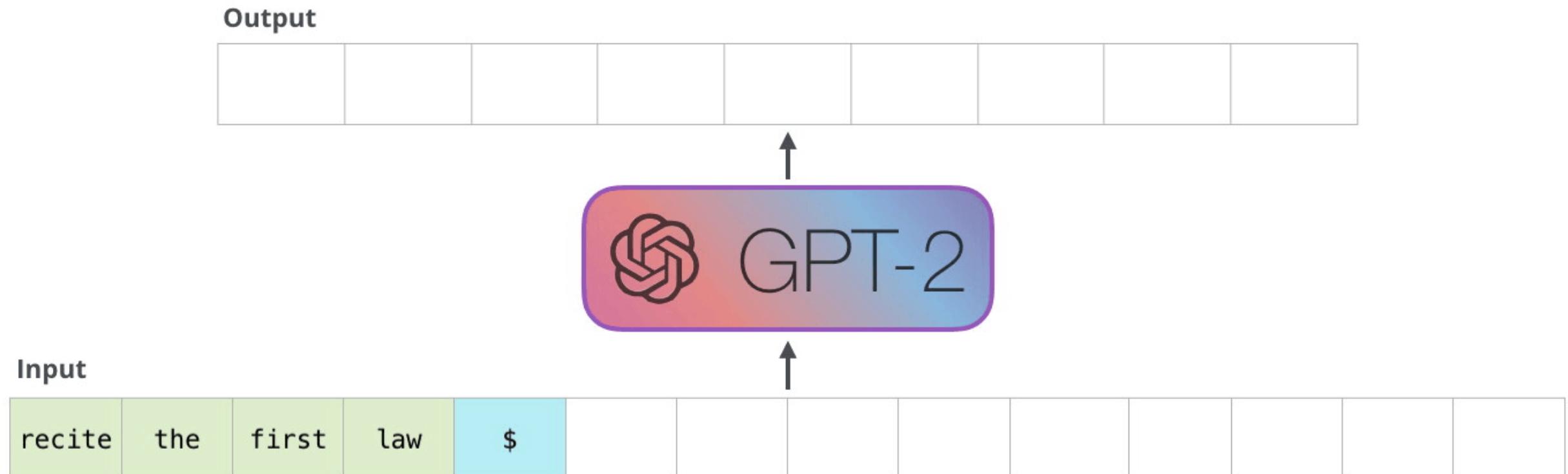
Text Classification



GPT

- Generative Pre-trained Transformer
- Use the Transformer Decoder architecture.
- Introduced in 2018 by OpenAI.

How it works?



Evolution of GPT

Model	Number of parameters	Training data size	Year
GPT	110M	4GB	2018
GPT-2	1.5B	40GB	2019
GPT-3	175B	≈2TB	2020

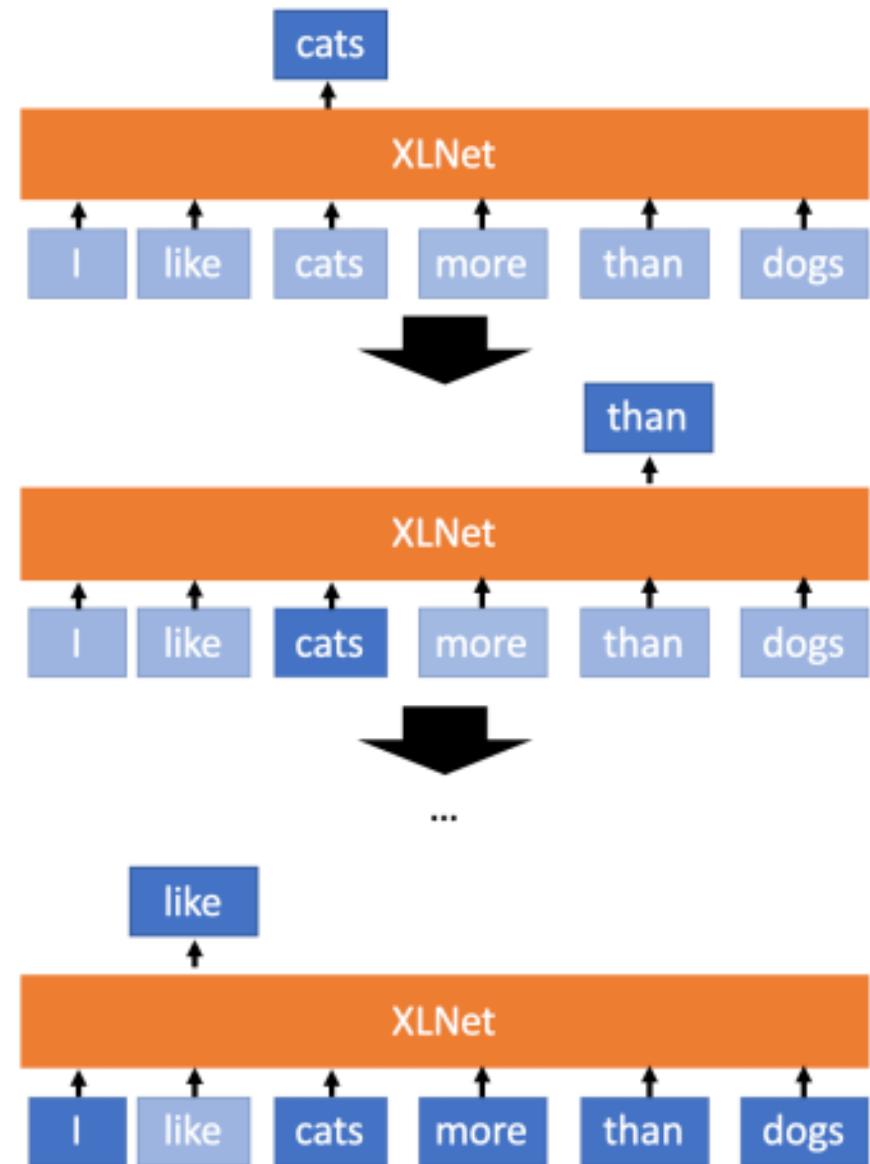
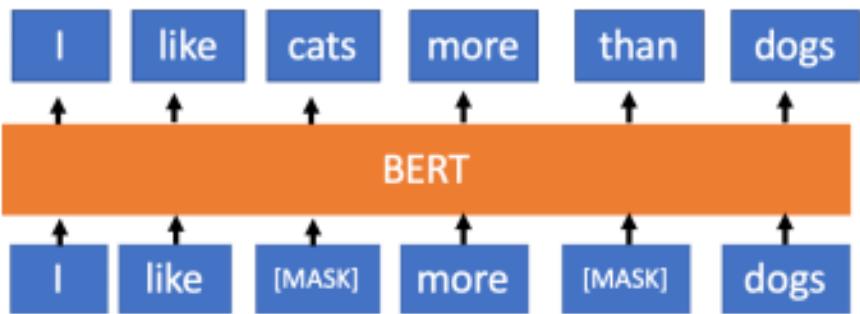
XLNet

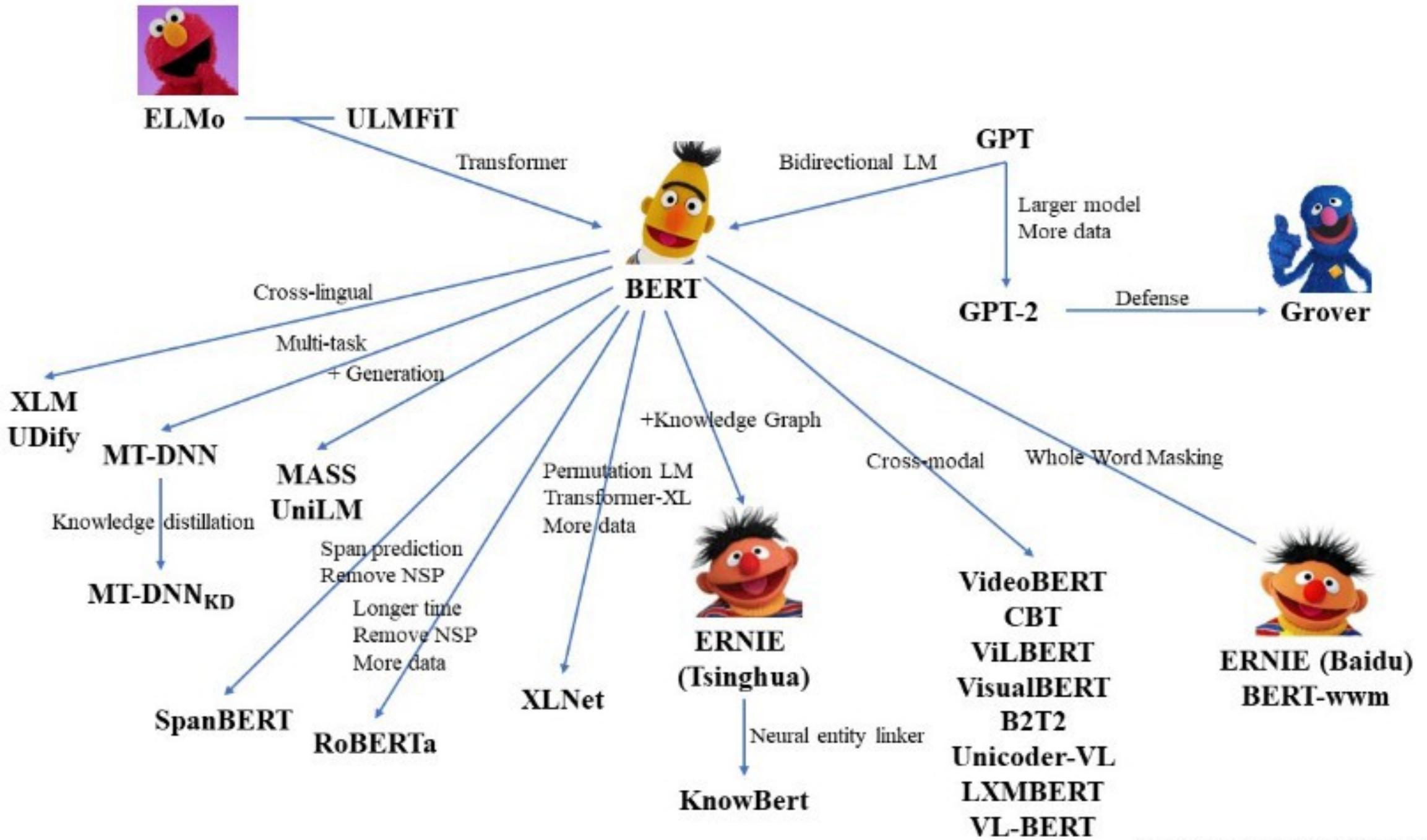
- Autoencoding (BERT):
 - [MASK] tokens do not appear during finetuning \Rightarrow pretrain-finetuning discrepancy.
 - Assume the predicted tokens are independent of each other given the unmasked tokens. Example: “New York is a city” \Rightarrow “[MASK] [MASK] is a city”
- Autoregressive (GPT):
 - Only trained to encode a unidirectional context (forward or backward).

XLNet

- XLNet combines pros from both while avoiding their cons.
- Techniques:
 - Permutation Language Modeling
 - Two-Stream Self-Attention for Target-Aware Representations
 - Incorporating Ideas from Transformer-XL
 - Modeling Multiple Segments

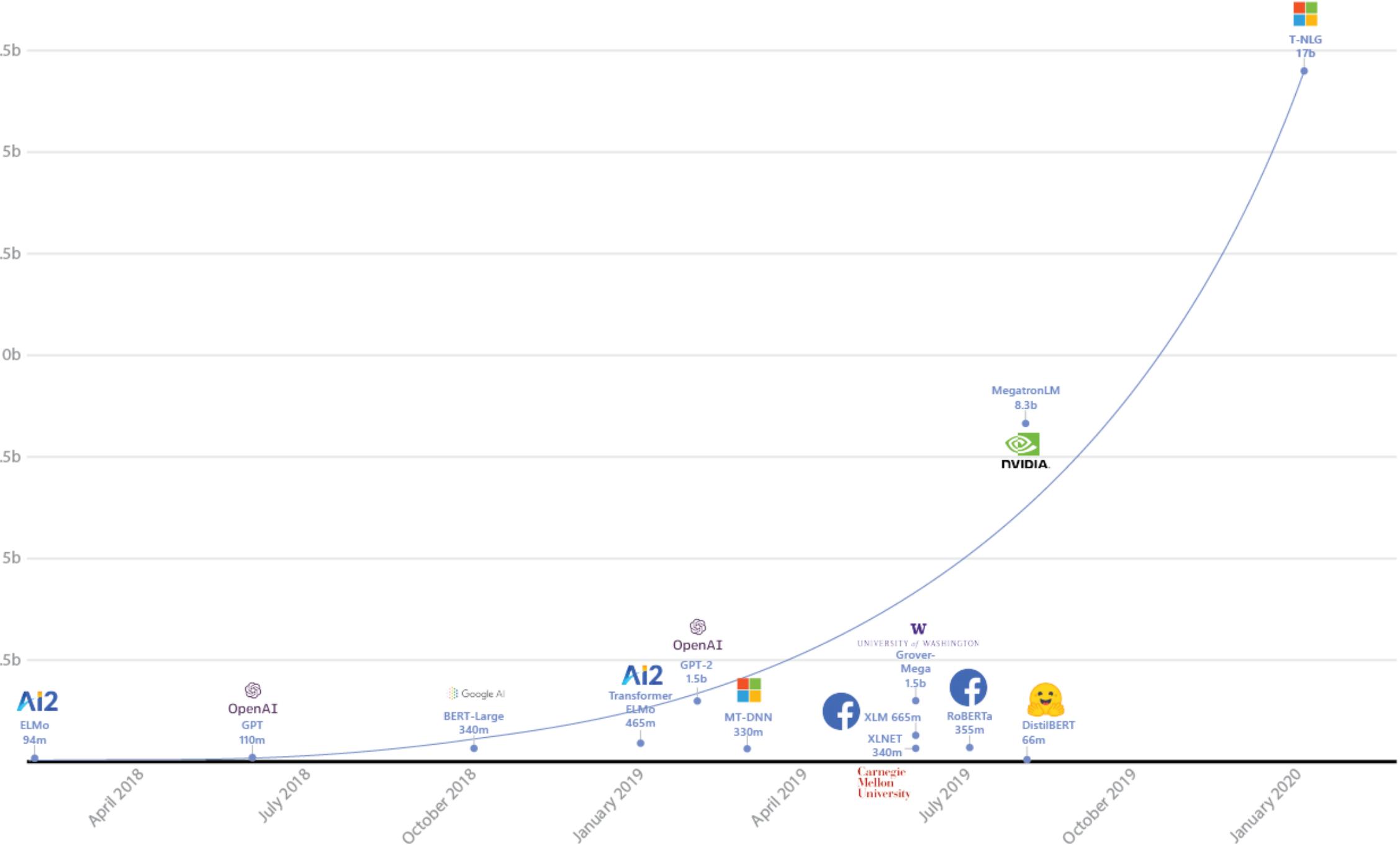
Permutation Language Modeling

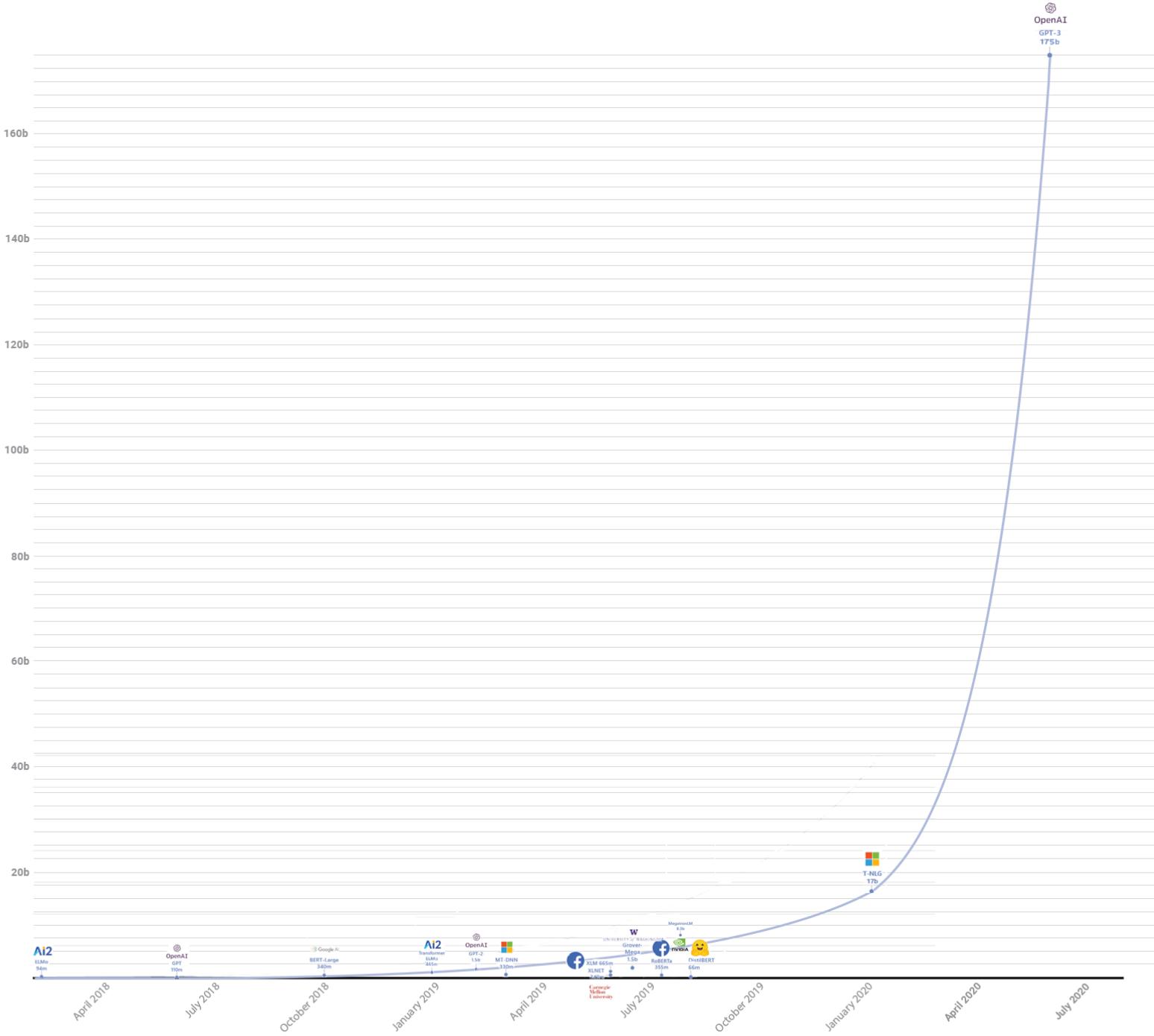




The GPT-3 Hype

- 175 billion parameters => 700 GB of memory!
- 10x larger than the 2nd largest model (Turing-NLG).

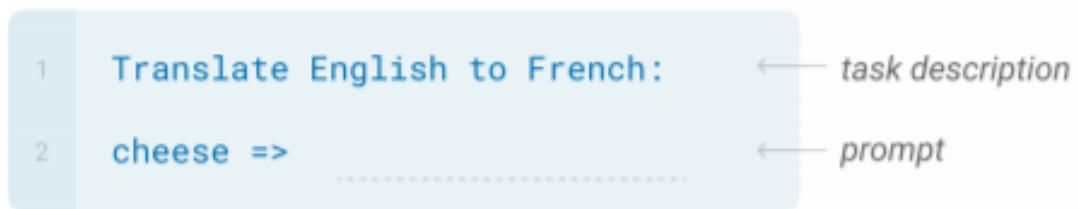




The GPT-3 Hype

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



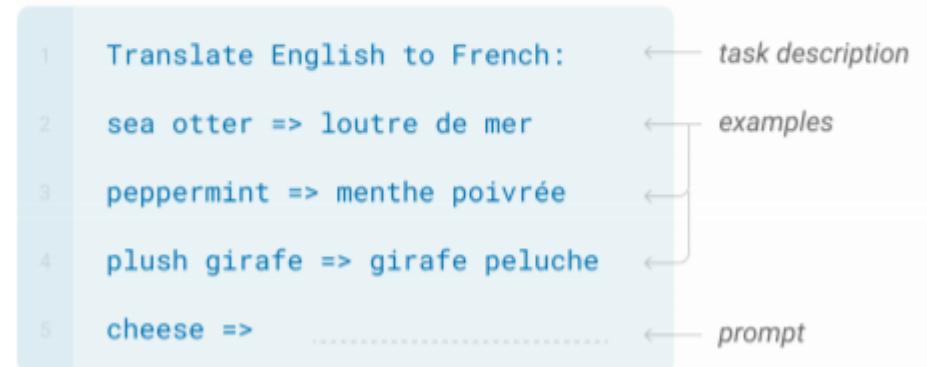
One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Applications: AI Dungeon ([link](#))

←

...

You are Duy, a survivor trying to survive in a post apocalyptic world by scavenging among the ruins of what is left. You have a backpack and a canteen. You haven't eaten in two days so you're desperately searching for food. You enter a rundown building and walk down a flight of stairs. At the bottom you find two human heads on a pole, skin still hanging off. You don't know what it means but you run out of the building.

> You find some food.

You walk back outside and see a store. You run inside and find that it is completely empty. You begin looking through the shelves and find some canned foods. You take them with you as you leave the building.

> You see some zombies.

You begin to see a few shamblers wandering around so you pick two that look slow and easy kill and begin leading them to where the others are.

Tip: Join the community! Check out the "Explore" section in the sidebar to view adventures and scenarios published by other players.



Do

What do you do?

G

>

Limitations

- Computation
- Size
- Transparency
- Can't capture the actual meaning of words. From GPT-3 Authors:

“GPT-3 samples still sometimes repeat themselves semantically at the document level, start to lose coherence over sufficiently long passages, contradict themselves, and occasionally contain non-sequitur sentences or paragraphs.”

Human: How many eyes does a giraffe have?

GPT-3: A giraffe has two eyes.

Human: How many eyes does my foot have?

GPT-3: Your foot has two eyes.

Human: How many eyes does a spider have?

GPT-3: A spider has eight eyes.

Human: How many eyes does the sun have?

GPT-3: The sun has one eye.



What's next?



- Human brain: **100 trillion*** synapses

- * Could be 1000+ trillion. ANN != BNN.

- GPT-3: **175 billion** parameters

- Cost: \$4.6m

- GPT-4?: **100 trillion** parameters

- Cost (2020): \$2.6 billion
 - Cost (2024): \$325 million
 - Cost (2028): \$40 million
 - Cost (2032): \$5 million

Reference

- Devlin, J. et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.” NAACL-HLT (2019).
- Yang, Z. et al. “XLNet: Generalized Autoregressive Pretraining for Language Understanding.” NeurIPS (2019).
- Brown, T. et al. “Language Models are Few-Shot Learners.” ArXiv abs/2005.14165 (2020): n. pag.
- Lex Fridman. GPT-3 vs Human Brain.
https://www.youtube.com/watch?v=kpiY_LemaTc
- Illustrations: <https://jalammar.github.io/>