



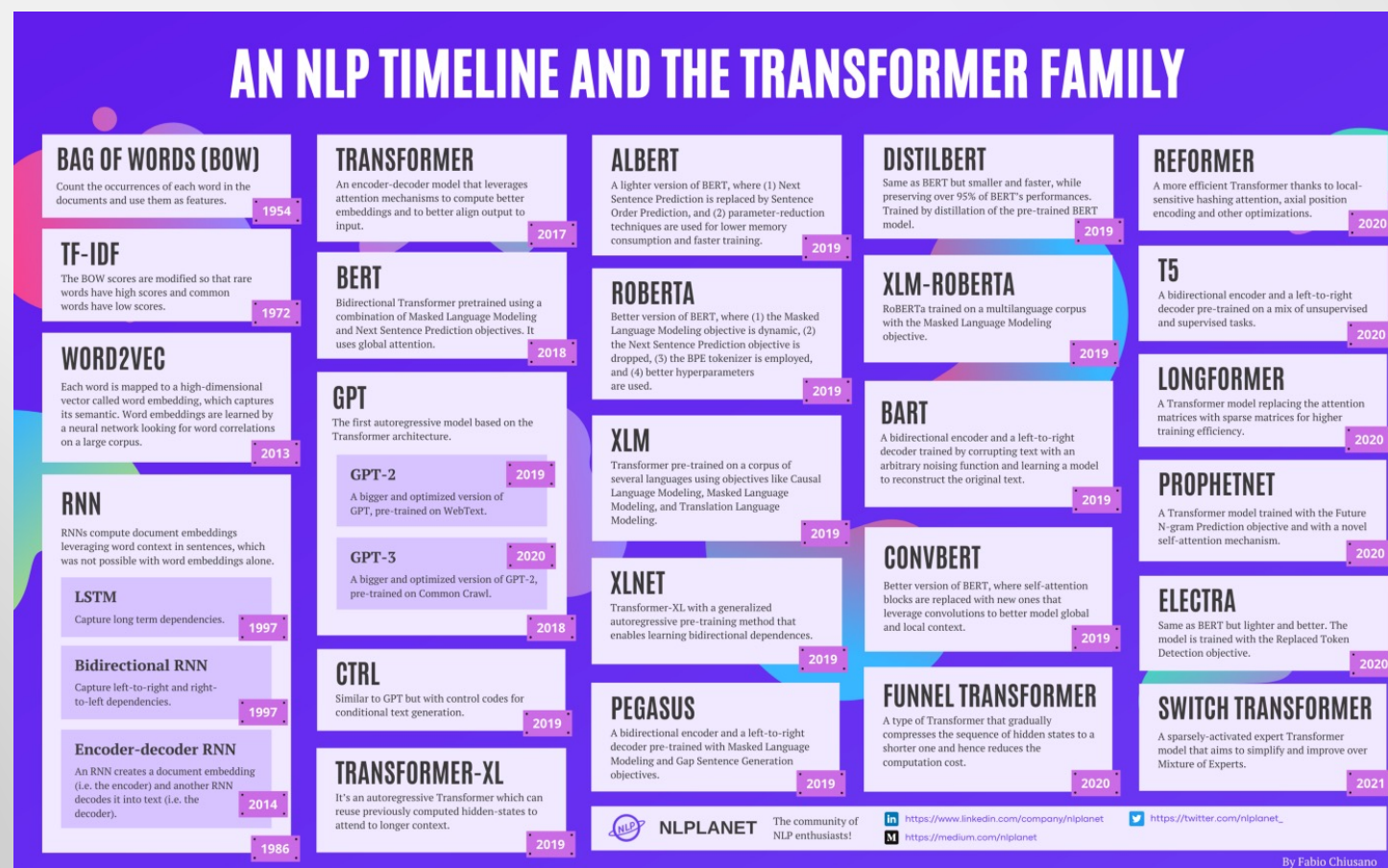
# C EURO<sup>2</sup>

Mastering Transformers: From Building Blocks to Real-World Applications

NLP Applications by Prof.Dr. Tuğba Taşkaya Temizel

12 Sept 2023

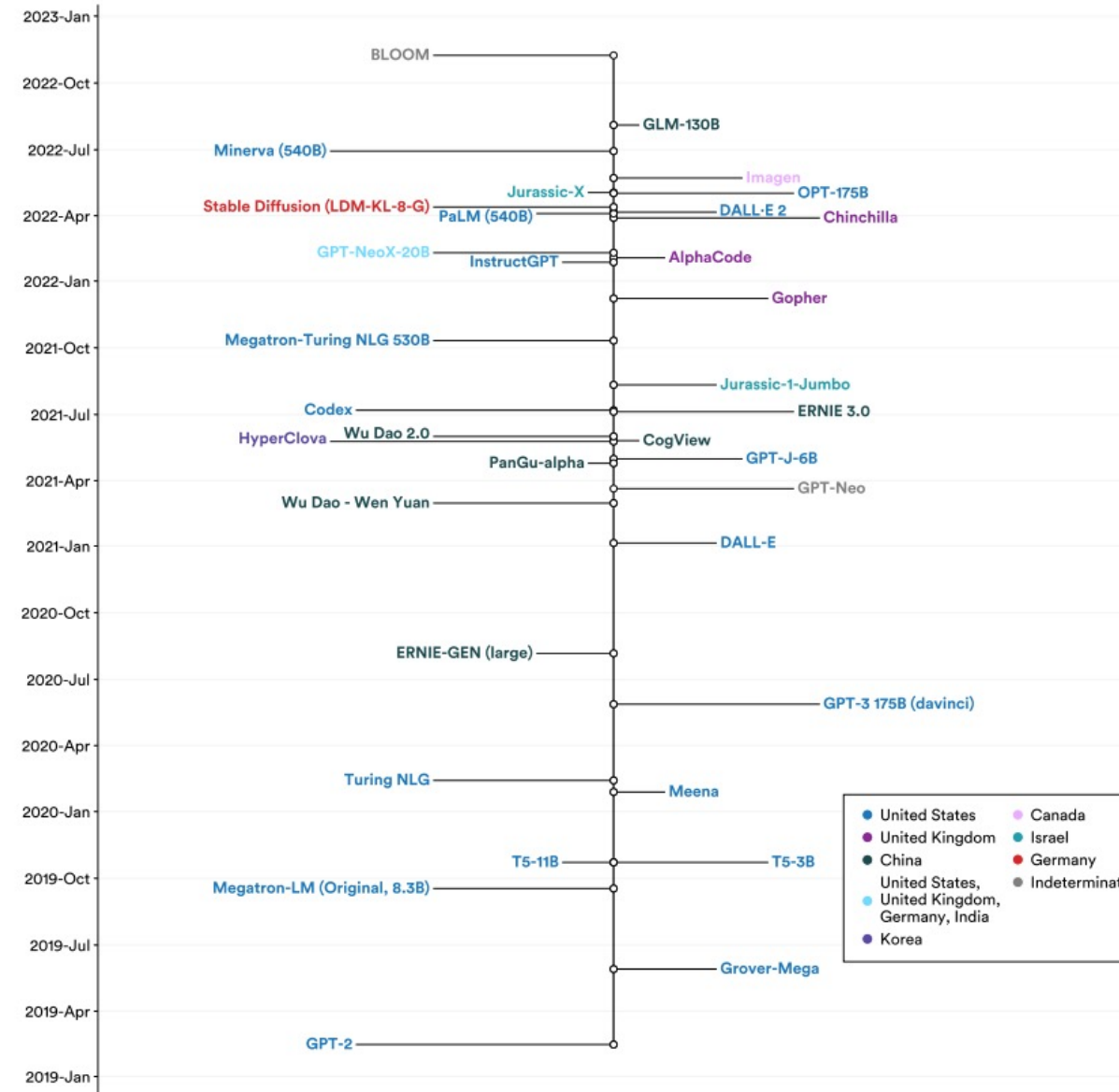
# A Brief History in Transformers



<https://medium.com/nlplanet/a-brief-timeline-of-nlp-from-bag-of-words-to-the-transformer-family-7caad8bbba56>

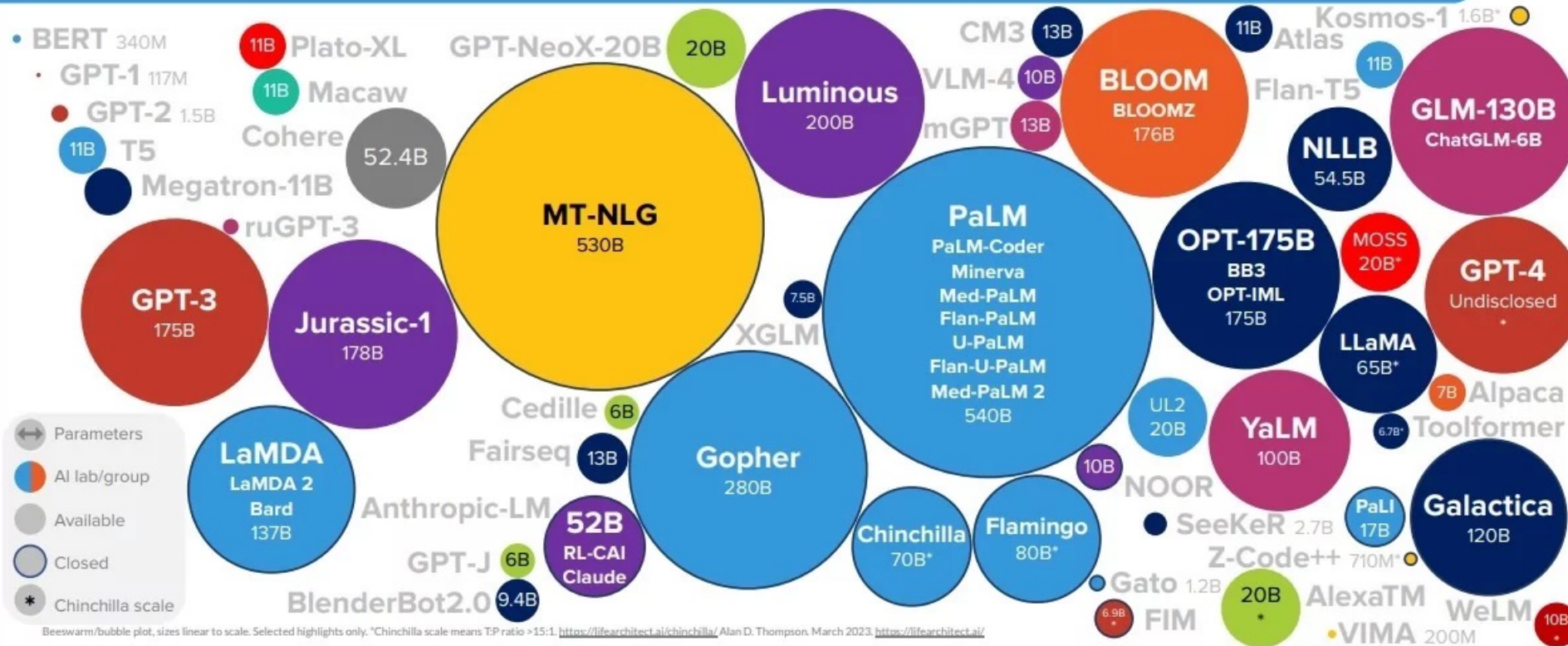
## Timeline and National Affiliation of Select Large Language and Multimodal Model Releases

Source: AI Index, 2022 | Chart: 2023 AI Index Report





# LANGUAGE MODEL SIZES TO MAR/2023



LifeArchitect.ai/models

A fast growing area:  
 LLaMa, LLaMa2  
 Il's Falcon 140B  
 GPT4  
 Claude  
 Falcon  
 Cohere  
 ...

<https://pureinsights.com/blog/2023/what-are-large-language-models-llms-search-and-ai-perspectives/>  
<https://sungkim11.medium.com/list-of-open-sourced-fine-tuned-large-language-models-llm-8d95a2e0dc76>  
<https://datasciencedojo.com/blog/best-large-language-models/>

# Language Models

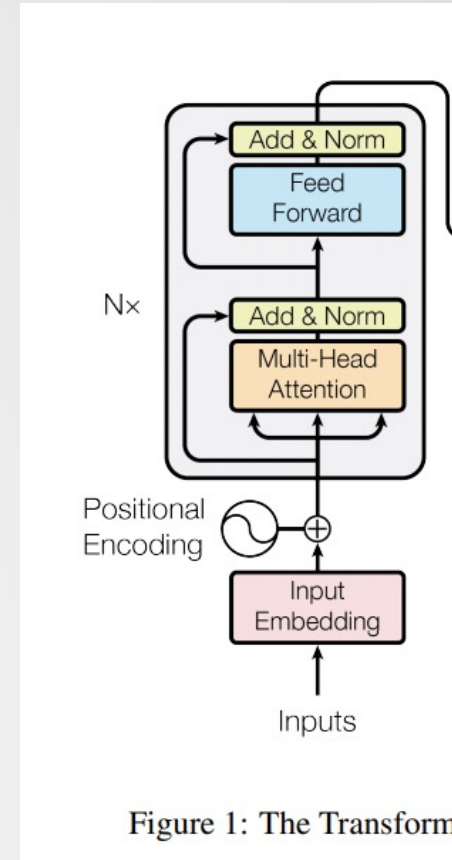
- **Language modeling:** predicting upcoming words from prior word context.
- **Neural net-based language models** turn out to have many advantages over the n-gram language models:
  - They can handle much longer histories, and they can generalize over contexts of similar words.
  - For a training set of a given size, a neural language model has much higher predictive accuracy than an n-gram language model.
  - They underlie many of the models for tasks like machine translation, dialog, and language generation.

# NLP Applications

- Language models, once created, can be adapted and used for numerous NLP tasks.
- Motivation: Unlabeled text corpora are abundant, but labeled data for learning these specific tasks is scarce.
- Large gains on these tasks can be realized by generative pre-training of a language model on a diverse corpus of unlabeled text, followed by discriminative fine-tuning on each specific task.
- It falls under the category of **semi-supervised learning** for natural language.
  - Semi-supervised learning combines many unlabeled data with small labeled data for training.

# BERT (Bidirectional Encoder Representations from Transformers)

- BERT's model architecture is a multi-layer bidirectional Transformer **encoder based** on the original implementation described in Vaswani et al. (2017)
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems, pages 6000–6010
- Language model where k% of the input words were masked, and then were predicted





# BERT (Bidirectional Encoder Representations from Transformers)

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

**Model:**



**Dataset:**



**Objective:**

Predict the masked word (language modeling)

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

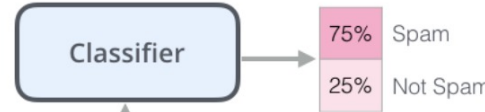
## Supervised Learning Step

**Model:**  
(pre-trained in step #1)



**Dataset:**

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



Widely used for  
classification tasks  
and sentiment  
prediction

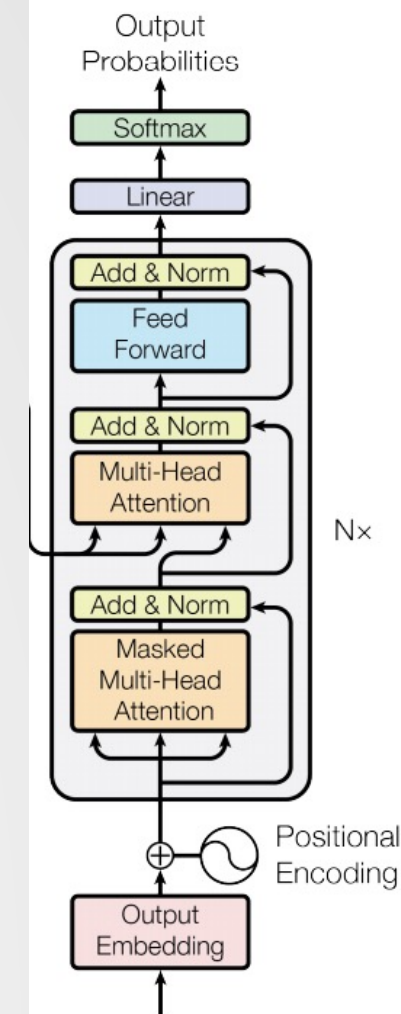
The two steps of how BERT is developed. You can download the model pre-trained in step 1 (trained on un-annotated data), and only worry about fine-tuning it for step 2. [\[Source for book icon\]](https://alammar.github.io/illustrated-bert/).



# Generative Pre-Training (GPT-1)

## Some Model Details

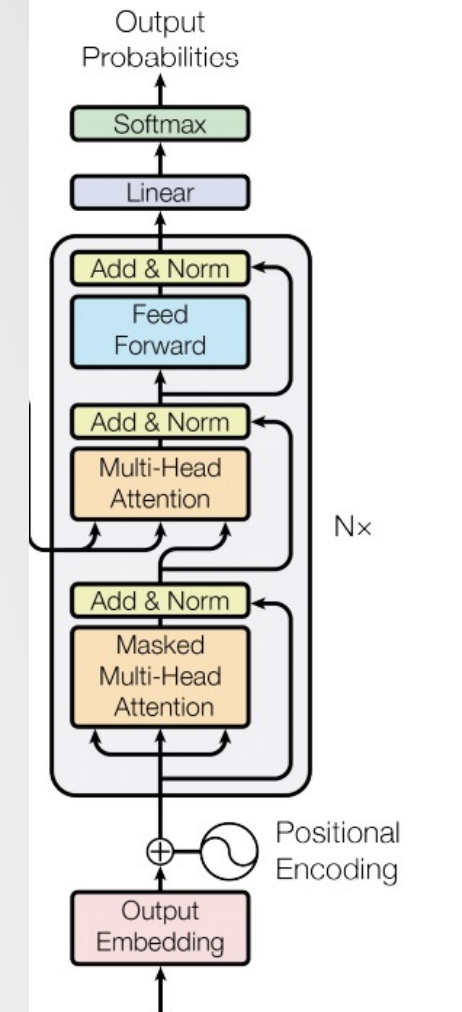
- They aim to predict the subsequent token as a part of a language modelling task.
  - Note that the original transformer was used for machine translation.
- Trained 12-layer decoder-only transformer with masked self-attention heads (768 dimensional states and 12 attention heads).
  - Total number of parameters is 117M parameters.
- Generating text that is both coherent and contextually consistent is a significant challenge in many NLP tasks, such as text generation, machine translation, and content summarization.
  - Decoder-only transformer architectures have been designed to address this problem. E.g.: GPT (Generative Pre-trained Transformer) models



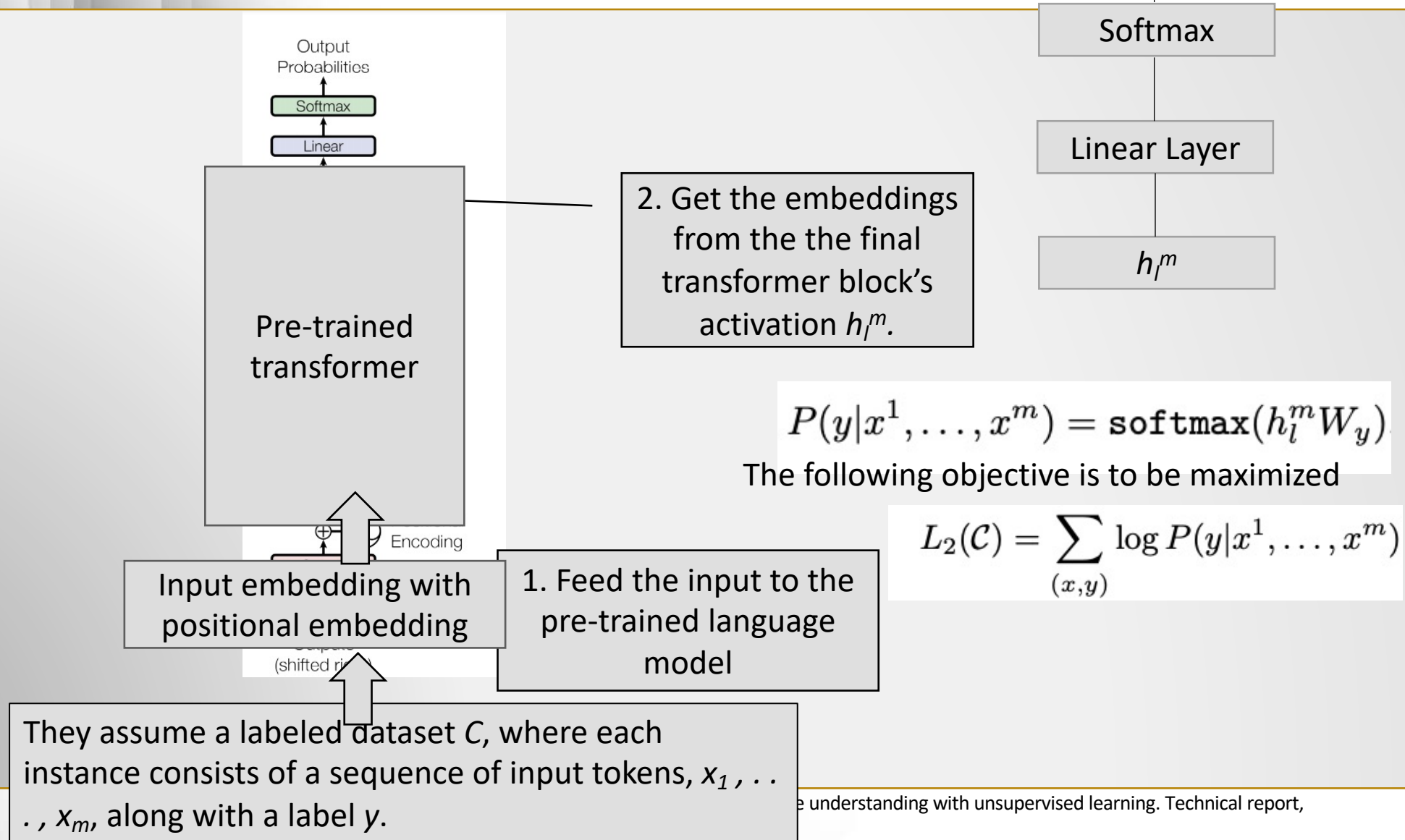
# Generative Pre-Training (GPT-1)

## Some Model Details

- For the position-wise feed-forward networks, they used 3072 dimensional inner states.
- They used the Adam optimization scheme with a max learning rate of  $2.5e-4$ .
- They used a byte pair encoding (BPE) vocabulary with 40,000 merges [53] and residual, embedding, and attention dropouts with a rate of 0.1 for regularization.
- They used the Gaussian Error Linear Unit (GELU).
- Their language model did not have access to subsequent words to the right of current word.



# Generative Pre-Training (GPT-1) Architecture



# Generative Pre-Training (GPT-1)

## Architecture: Auxiliary Training Objectives

- Adding auxiliary unsupervised training objectives is an alternative form of semi-supervised learning.
- They found that including language modeling as an auxiliary objective to the fine-tuning helped learning by
  - (a) improving generalization of the supervised model, and
  - (b) accelerating convergence.
- Specifically, they optimize the following objective (with weight  $\lambda$ ):

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

$\lambda$  was set to 0.5

Learning objective of  
the classification task.

Learning objective of  
the language model  
task.



# Textual Entailment

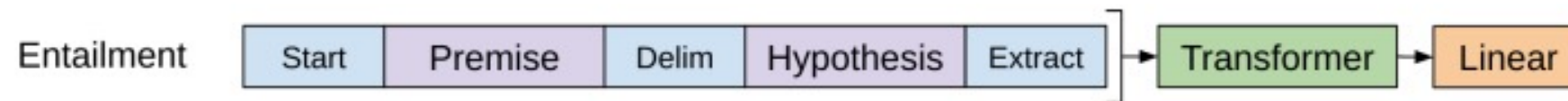
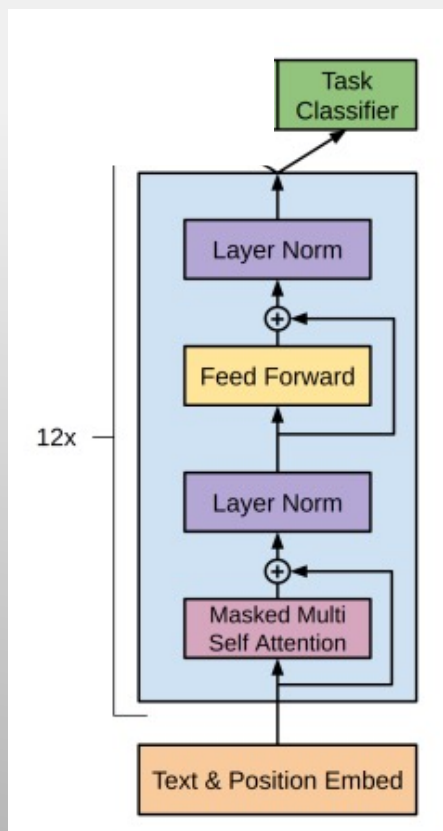
- **Textual entailment recognition** is the task of deciding, given two text fragments, whether the meaning of one text is entailed (can be inferred) from another text.
- It involves reading a pair of sentences and judging the relationship between them from one of entailment, contradiction or neutral.

Textual entailment example from Stanford Natural Language Inference (SNLI) corpus:

SNLI corpus (version 1.0) is a collection of 570k human-written English sentence pairs manually labeled for balanced classification with the labels *entailment*, *contradiction*, and *neutral*.

Text	Judgments	Hypothesis
A man inspects the uniform of a figure in some East Asian country.	contradiction C C C C C	The man is sleeping
An older and younger man smiling.	neutral N N E N N	Two men are smiling and laughing at the cats playing on the floor.
A black race car starts up in front of a crowd of people.	contradiction C C C C C	A man is driving down a lonely road.
A soccer game with multiple males playing.	entailment E E E E E	Some men are playing a sport.
A smiling costumed woman is holding an umbrella.	neutral N N E C N	A happy woman in a fairy costume holds an umbrella.

# Textual Entailment



The premise  $p$  and hypothesis  $h$  token sequences are concatenated, with a delimiter token (\$) in between.

All transformations include adding randomly initialized start and end tokens(<s>,<e>).

E.g.

Input: <s>A man inspects the uniform of a figure in some East Asian country\$The man is sleeping <e>).

Output: contradict, neutral, entailed (class variable)

# Text Similarity

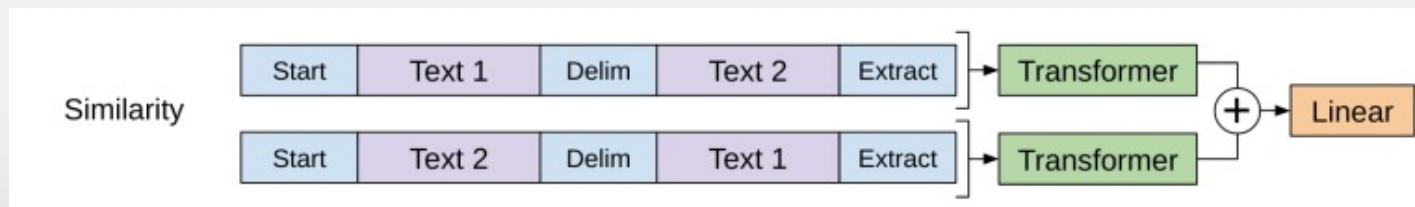
- **Semantic similarity** (or paraphrase detection) tasks involve predicting whether two sentences are semantically equivalent or not.
- The challenges lie in recognizing the rephrasing of concepts, understanding negation, and handling syntactic ambiguity.

Sematically similar example:

The sun is shining brightly in the clear blue sky.

A radiant sun is illuminating the cloudless sky

# Text Similarity



- There is no inherent ordering of the two sentences being compared.
  - Output: true/false
- They modify the input sequence to contain both possible sentence orderings (with a delimiter in between) and process each independently to produce two sequence representations added element-wise before being fed into the linear output layer.

Output of first transformer

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{bmatrix} + \begin{bmatrix} b_{11} & b_{12} & b_{13} & \dots & b_{1n} \\ b_{21} & b_{22} & b_{23} & \dots & b_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & b_{m3} & \dots & b_{mn} \end{bmatrix} = \begin{bmatrix} a_{11} + b_{11} & a_{12} + b_{12} & a_{13} + b_{13} & \dots & a_{1n} + b_{1n} \\ a_{21} + b_{21} & a_{22} + b_{22} & a_{23} + b_{23} & \dots & a_{2n} + b_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{m1} + b_{m1} & a_{m2} + b_{m2} & a_{m3} + b_{m3} & \dots & a_{mn} + b_{mn} \end{bmatrix}$$

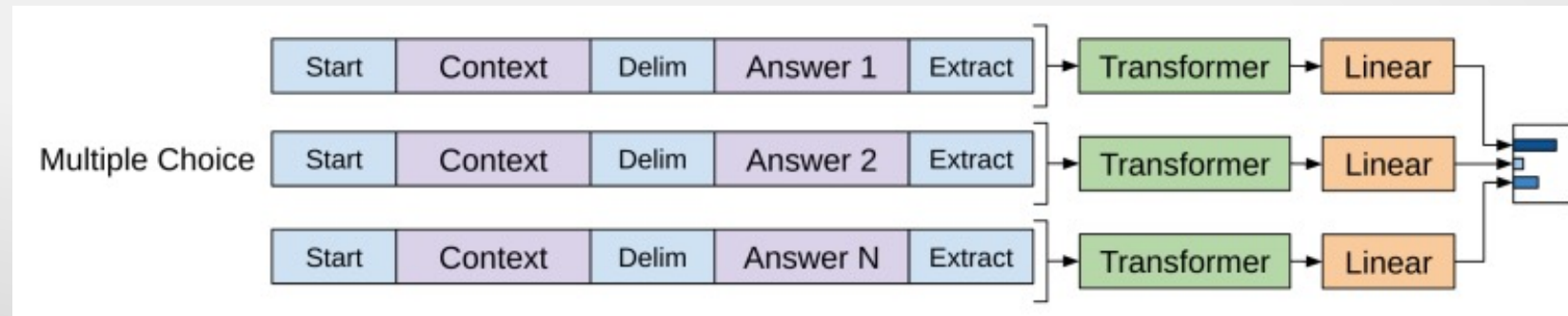
$[a_{ij}]_{m \times n}$ 
 $[b_{ij}]_{m \times n}$ 
 $[a_{ij} + b_{ij}]_{m \times n}$

Output of second transformer

Element-wise added



# Question Answering



We are given a context document  $z$ , a question  $q$ , and a set of possible answers  $\{a_k\}$ .

We concatenate the document context and question with each possible answer, adding a delimiter token in between to get  $[z;q;\$;a_k]$ .

Each of these sequences is processed independently with their model and then normalized via a softmax layer to produce an output distribution over possible answers.

# Summary

- **Encoder based models:** Encoder part is responsible for understanding and extracting the relevant information from the input text. It generates a representation of the input text.
  - E.g. BERT, RoBERTa (good for label prediction)
- **Decoder based models:** Its masked multi-head attention component facilitates text generation by letting the model not seeing the subsequent tokens (autoregressive decoding lets generating token one at a time conditioning on the previously generated tokens).
  - E.g. GPT1, GPT2
- **Encoder-decoder hybrids:** Encoder-decoder models are typically used for natural language processing tasks that involve understanding input sequences and generating output sequences, often with different lengths and structures.
  - E.g. BART, T5. (machine translation, summarization tasks...)

# Thanks!



**EuroHPC**  
Joint Undertaking

This project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 951732. The JU receives support from the European Union's Horizon 2020 research and innovation programme and Germany, Bulgaria, Austria, Croatia, Cyprus, Czech Republic, Denmark, Estonia, Finland, Greece, Hungary, Ireland, Italy, Lithuania, Latvia, Poland, Portugal, Romania, Slovenia, Spain, Sweden, United Kingdom, France, Netherlands, Belgium, Luxembourg, Slovakia, Norway, Switzerland, Turkey, Republic of North Macedonia, Iceland, Montenegro