

# INF721

2024/2



# Deep Learning

## L18: Transformers (Part II)

# Logistics

## Announcements

- ▶ PA4 is due tonight! 11:59pm
- ▶ Midterm Exam II has been pushed back to Nov 25th

## Last Lecture

- ▶ Problems with RNNs
- ▶ Transformers
  - ▶ Self-attention
  - ▶ Multi-head Attention
  - ▶ Positional Embedding
  - ▶ Masked Multi-head

# Lecture Outline

- ▶ Contextual Word Embeddings
- ▶ General Pre-Training (GPT)
  - ▶ Fine-tuning
  - ▶ Generating Text
  - ▶ ChatGPT
- ▶ Bidirectional Encoder Representations from Transformers (BERT)
  - ▶ Masked Language Model
  - ▶ Word and Sentence Embedding

# Contextual Word Embeddings

Models like Word2Vec and Glove learn a static Embedding Matrix  $E$ , so they can't consider the context of a word to produce it's embedded vector.

$E =$

Man	Woman	King	Queen	Apple	Orange
-1	1	-0.95	0.97	0.00	0.01
0.01	0.02	0.93	0.95	-0.01	0.00
0.03	0.02	0.7	0.69	0.03	-0.02
0.04	0.01	0.02	0.01	0.95	0.97

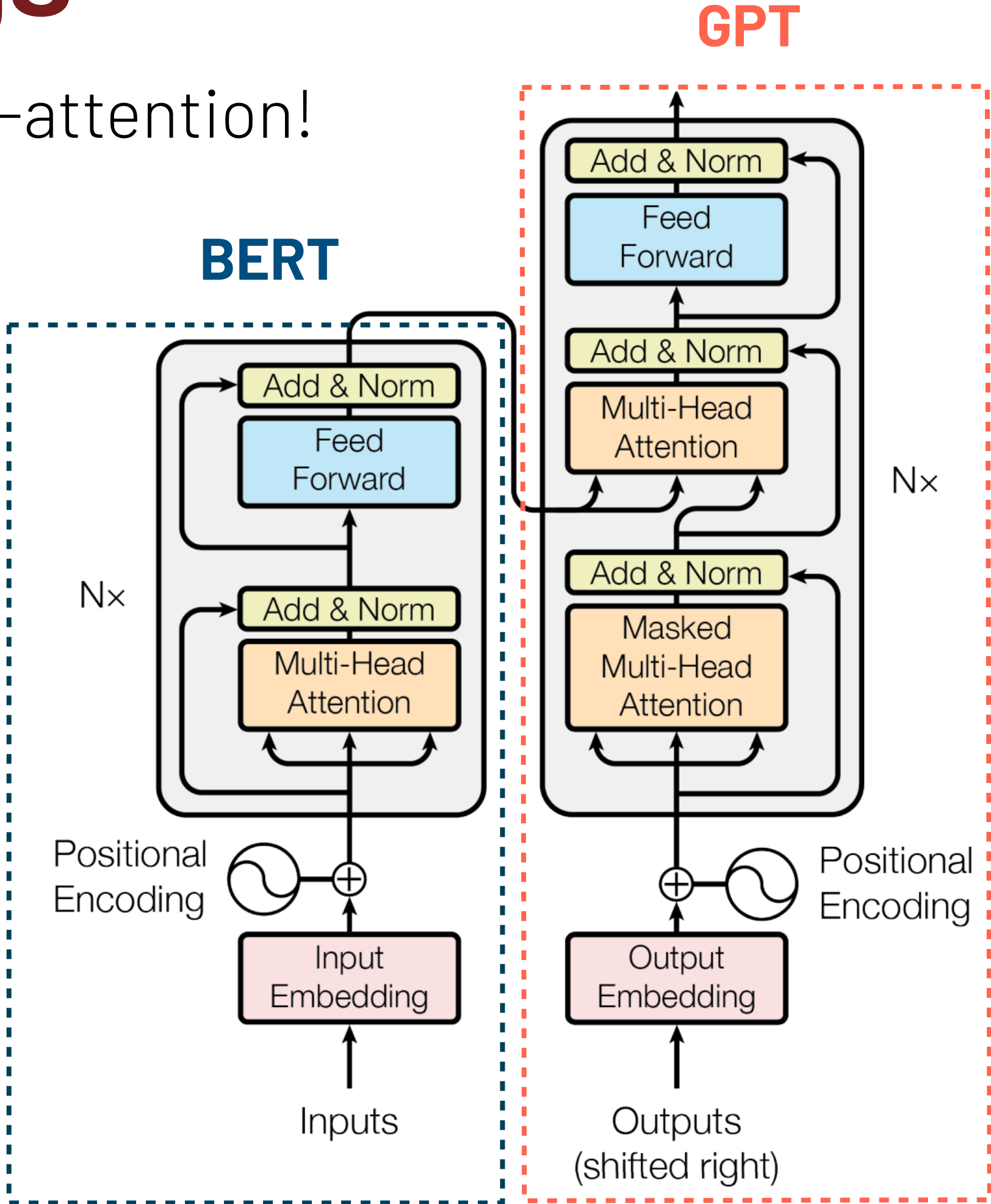
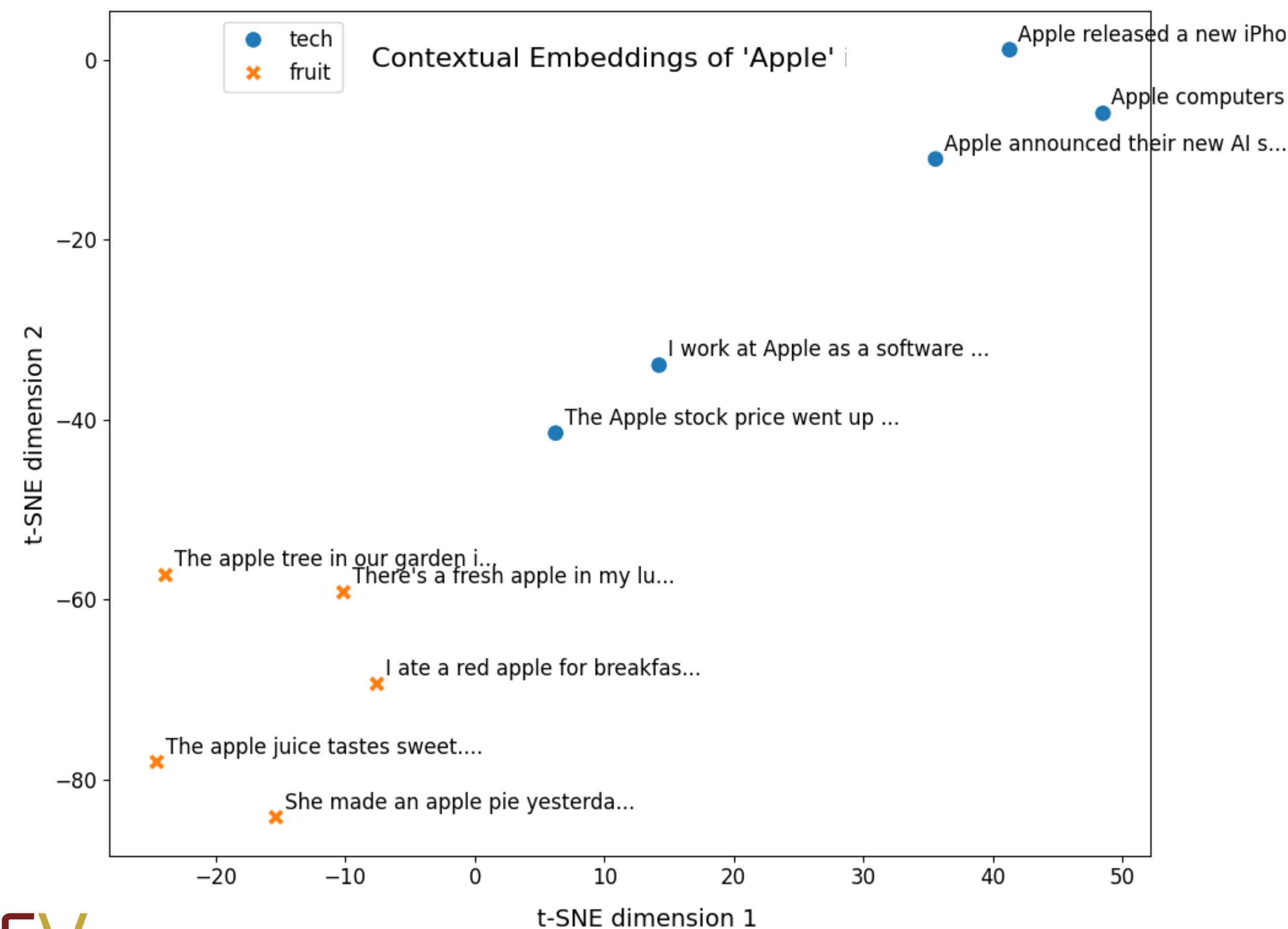
$$E \cdot o_{6257} = \begin{bmatrix} 0.01 \\ 0.00 \\ -0.02 \\ 0.97 \\ \dots \end{bmatrix}$$

$e_{6257}$

- For example, the word **apple** will have the same representation in both sentences:  
*"I want a glass of **apple** juice"*  
*"I work at **apple**"*
- But they have completely different meanings, because of their different context.

# Contextual Word Embeddings

Transformers learn **contextual representations** via self-attention!



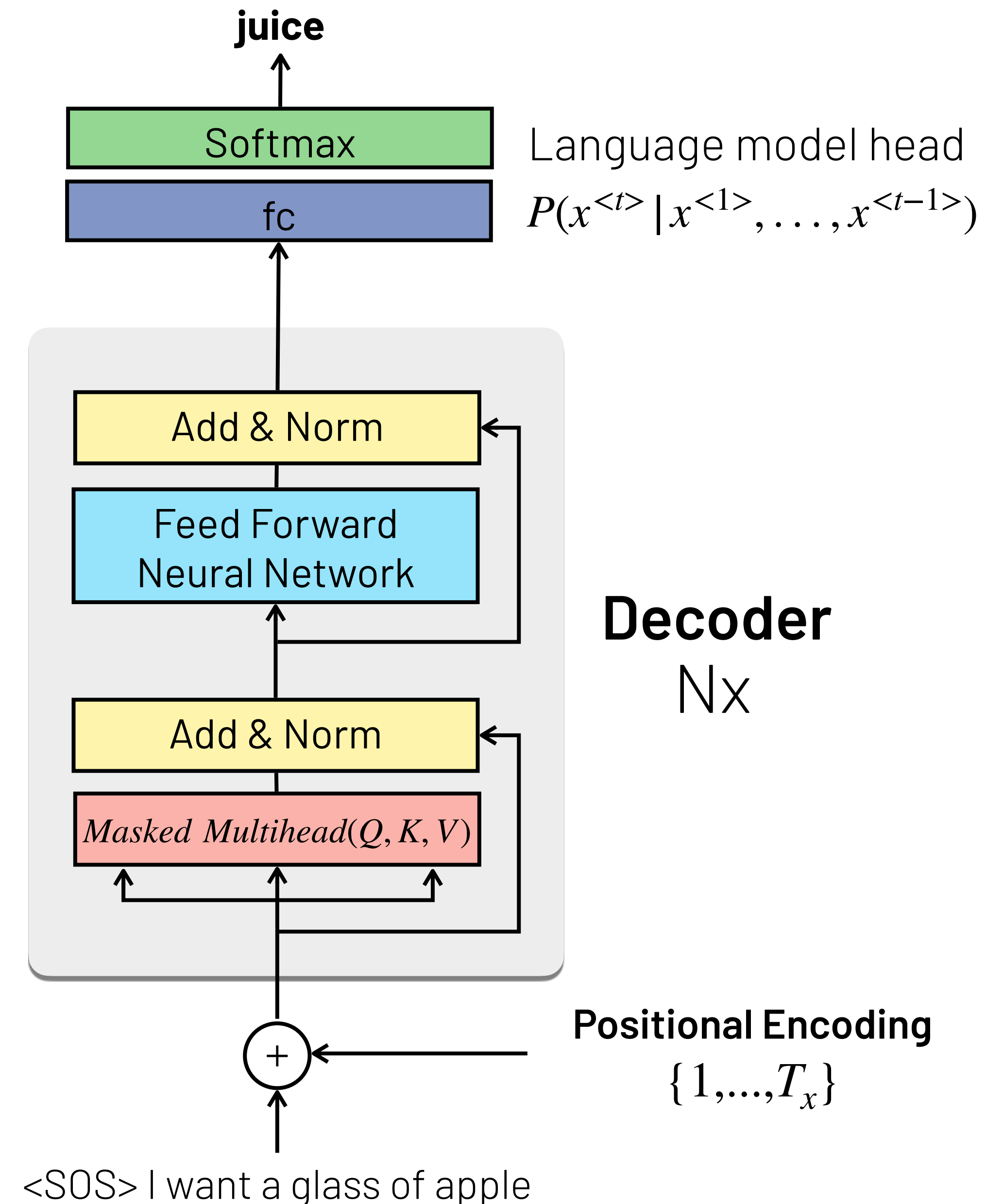
We can learn contextual embeddings using either the **Encoder** or the **Decoder** to train a Language Model!

# General Pre-Training (GPT)

**Language Model** based on the **Transformer Decoder**

1. **Train a language model** in unlabelled text

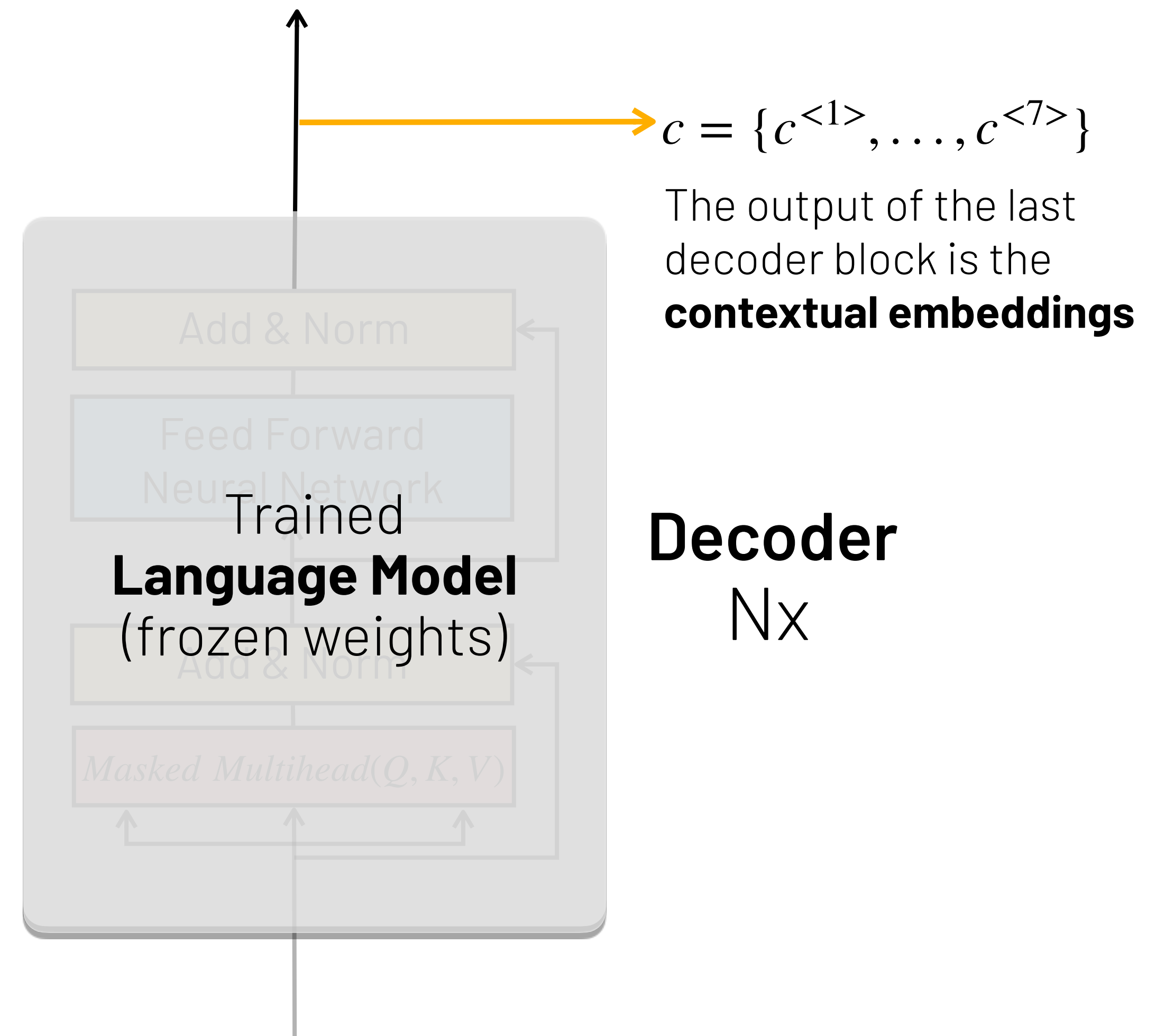
No need for another Multihead attention layer because we are not doing translation!



# General Pre-Training (GPT)

## Language Model based on the Transformer Decoder

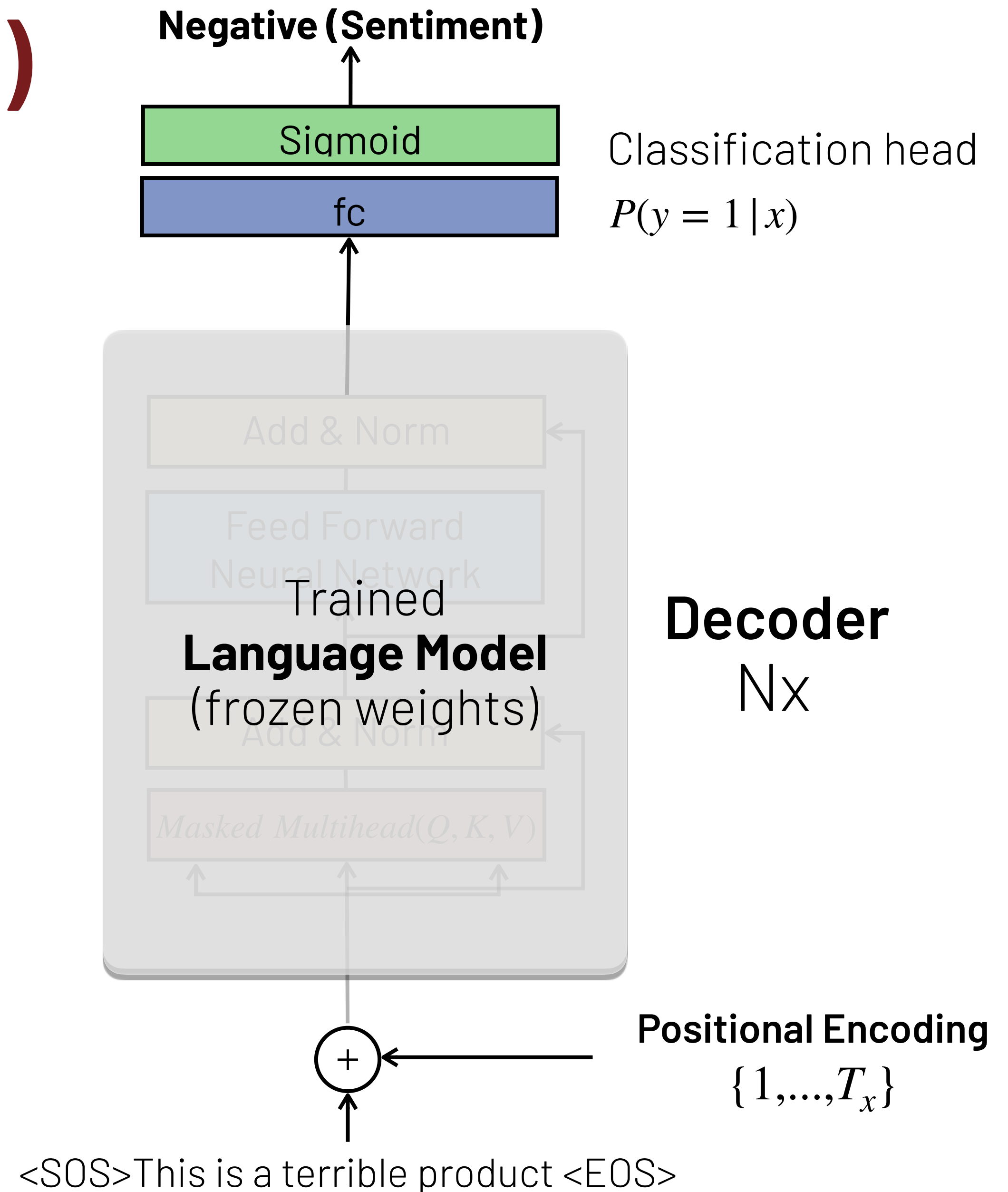
1. Train a language model in unlabelled text
2. **Freeze the model** up to the last transformer block and use it **to extract contextual embeddings**



# General Pre-Training (GPT)

## Language Model based on the Transformer Decoder

1. Train a language model in unlabelled text
2. Freeze the model up to the last transformer block and use it to extract contextual embeddings
3. **Substitute** the Language Model head by a **Classification** head (e.g., sentiment analysis)





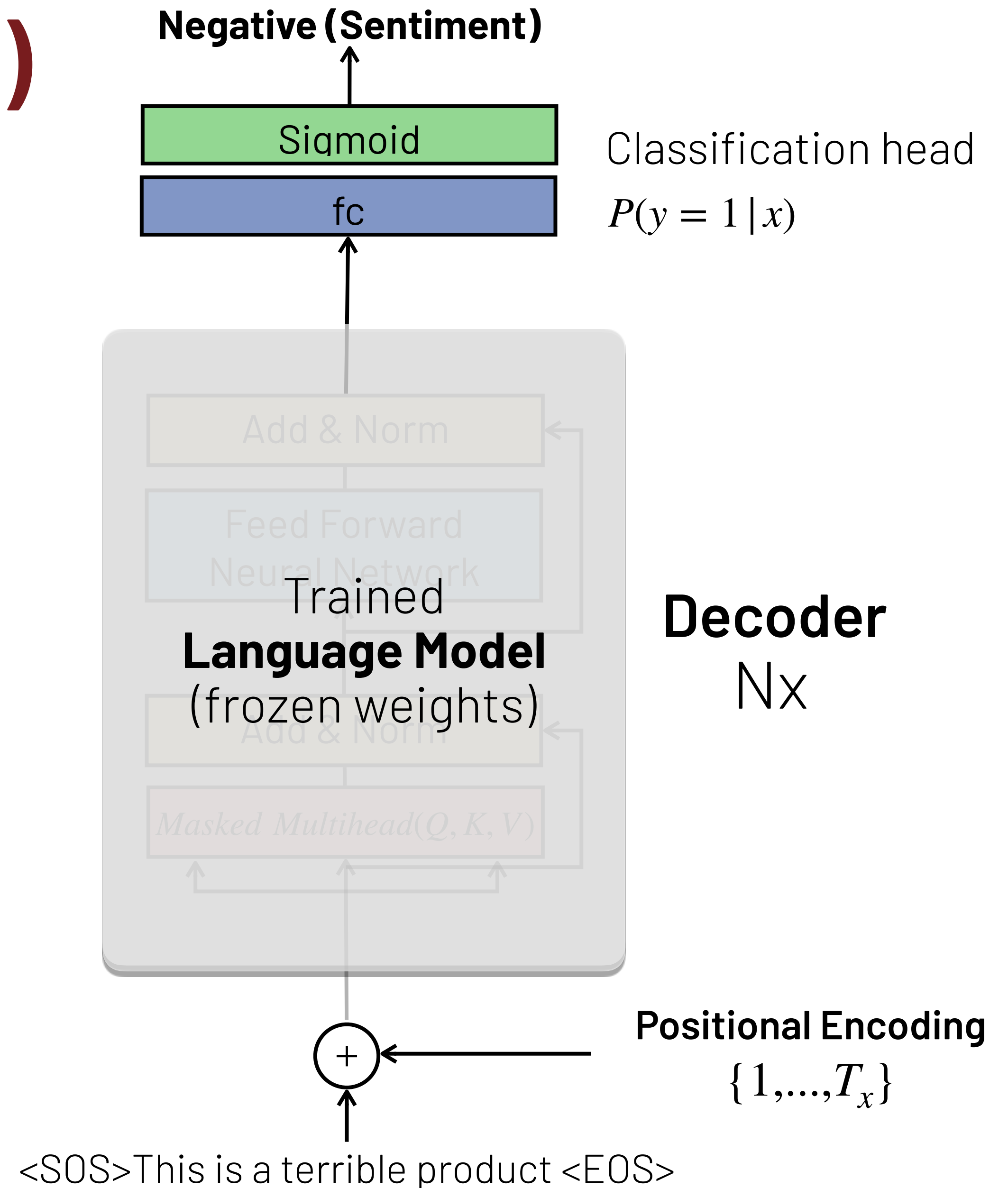
# General Pre-Training (GPT)

## Language Model based on the Transformer Decoder

1. Train a language model in unlabelled text
2. Freeze the model up to the last transformer block and use it to extract contextual embeddings
3. Substitute the Language Model head by a Classification head (e.g., sentiment analysis)

Step 1. is called **pre-training!**

Step 3. is called **fine-tuning!**



# Bytepair Encoding (BTE)

GPT (and many other transformer models) uses the Bytepair Encoding (BPE) algorithm to optimize the English Vocabulary, turning common patterns into single tokens.

- BPE is like building a dictionary of common word pieces by repeatedly combining the most frequent pairs of characters or subwords.

```
Text: "low lower lowest"
```

```
Initial: l o w _ l o w e r _ l o w e s t
```

```
Merges:
```

```
1. (l,o) → "lo": "lo w _ l o w e r _ l o w e s t"
```

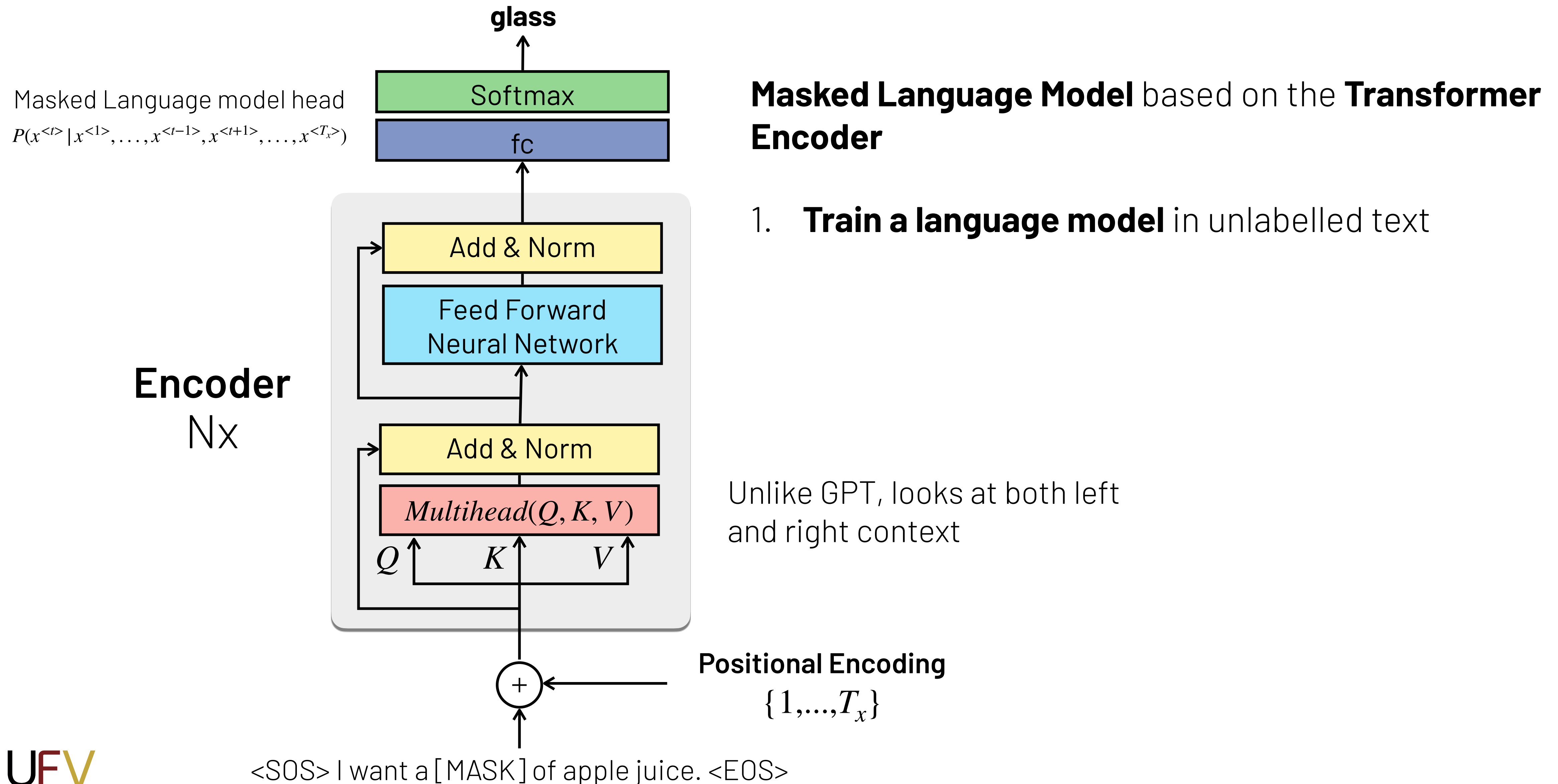
```
2. (lo,w) → "low": "low _ l o w e r _ l o w e s t"
```

```
3. (e,r) → "er": "low _ low er _ l o w e s t"
```

# GPT Evolution Over Time

	GPT 1 (2018)	GPT 2 (2019)	GPT 3 (2020)	GPT 3.5 (2022)	GPT 4 (2023)
Parameters	117M	1.5B	175B	~175B	Estimated ~1.8T
Nº of Decoder Blocks	12	48	96	~96	Unknown
Pre-training Dataset	5GB (~0.5B Tokens)	40GB (~8B tokens)	570 GB (~300B tokens)	570 GB (~300B tokens)	Unknown
Main Contribution	Introduced generative pretraining for transformers, showing decent results on text classification and sentiment analysis.	Generated coherent long-form text and exhibited surprising zero-shot learning abilities, like translation and summarization without specific task training.	Showed impressive few-shot, zero-shot, and multi-task learning capabilities.  Could perform tasks like question-answering, text generation, and even code generation.	Improved performance in language understanding and generation.  Reduced bias, hallucination, and increased coherence in long outputs.	Demonstrated strong multi-modal capabilities (image and text input), exhibited higher reasoning abilities, and improved complex task handling.

# Bidirectional Encoder Representations from Transformers (BERT)

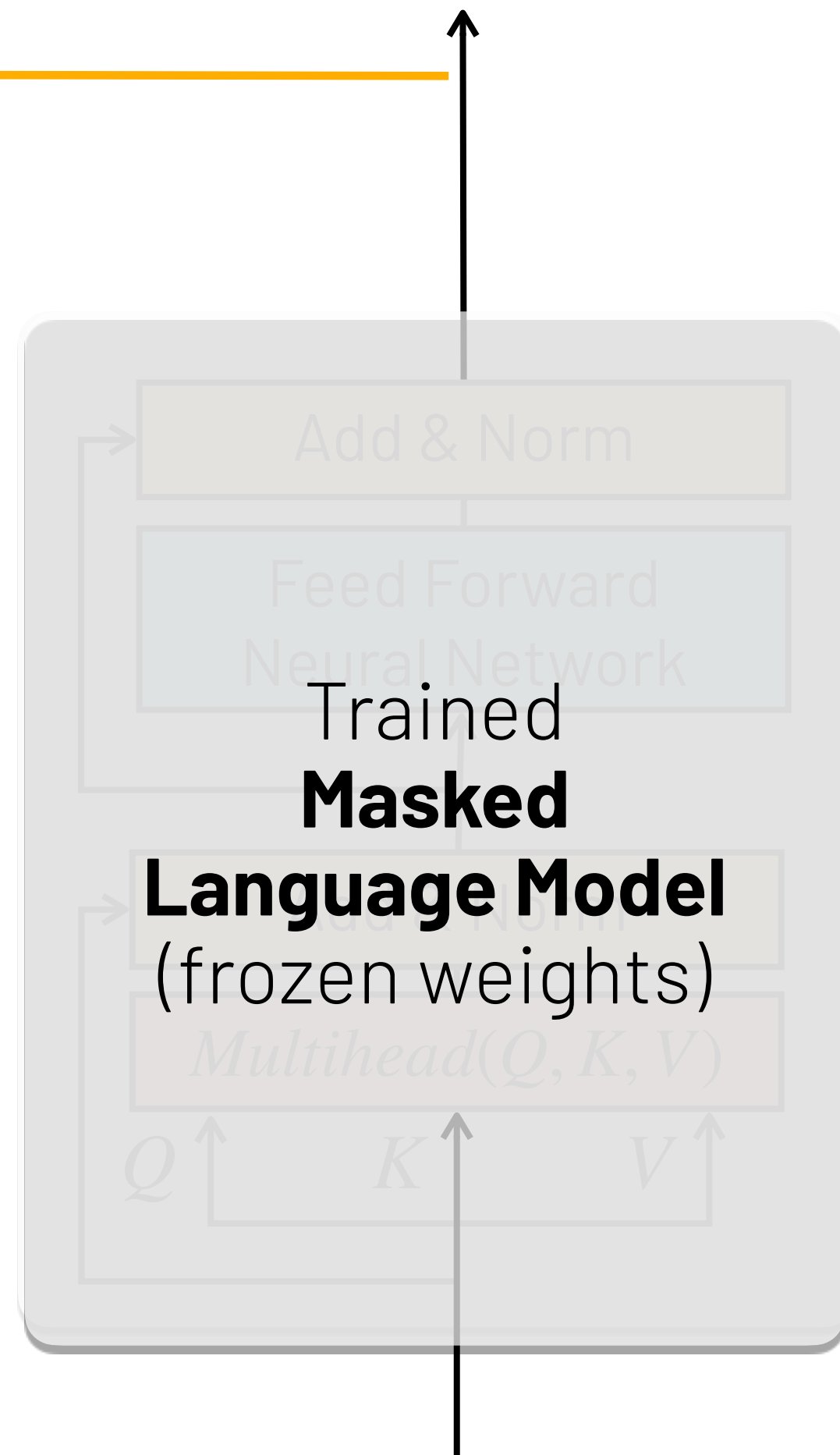


# Bidirectional Encoder Representations from Transformers (BERT)

$$c = \{c^{<1>}, \dots, c^{<7>}\}$$

The output of the last encoder block is the **contextual embeddings**

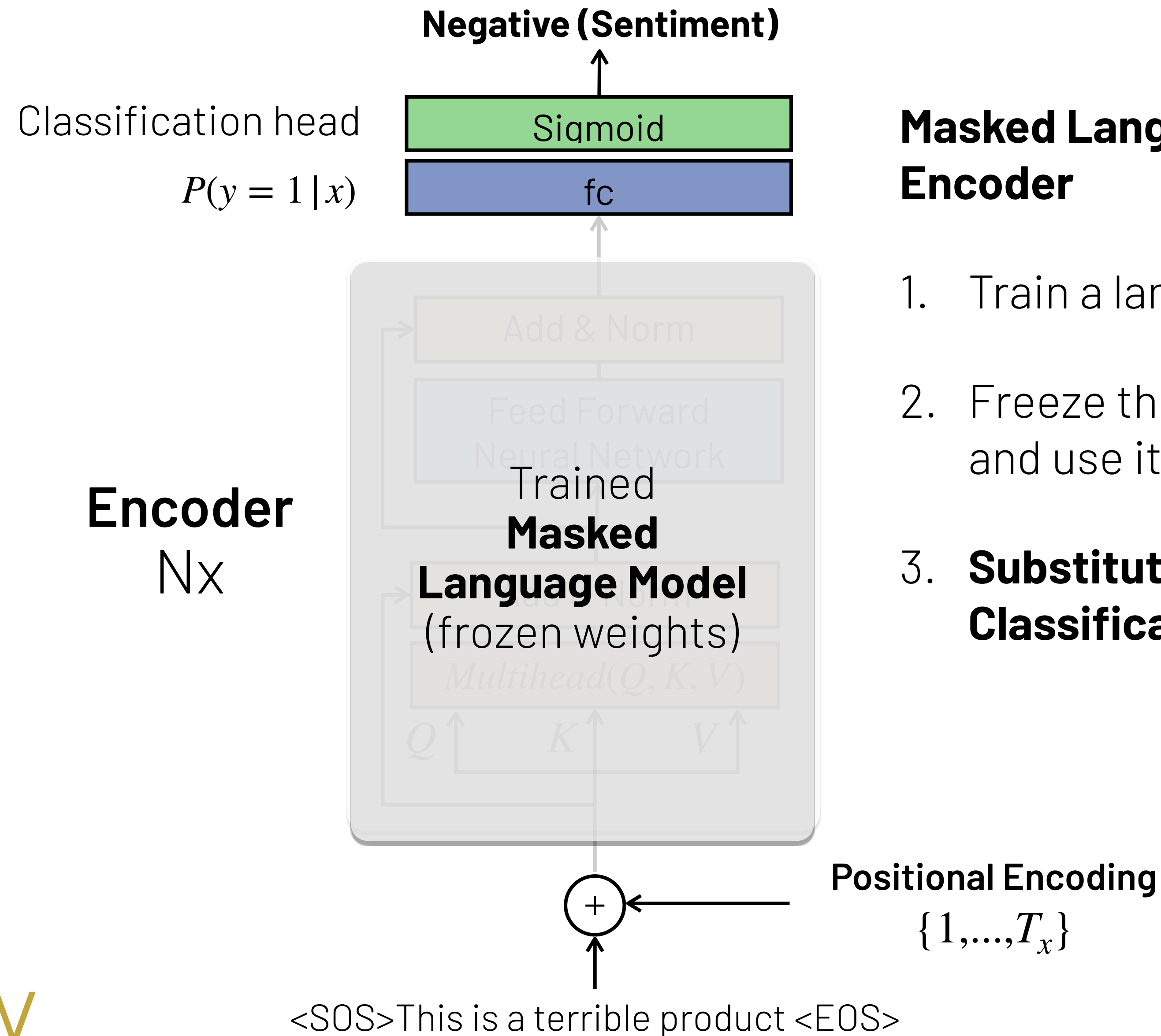
**Encoder**  
Nx



## **Masked Language Model** based on the **Transformer Encoder**

1. Train a language model in unlabelled text
2. **Freeze the model** up to the last transformer block and use it to **extract contextual embeddings**

# Bidirectional Encoder Representations from Transformers (BERT)



## Masked Language Model based on the Transformer Encoder

1. Train a language model in unlabelled text
2. Freeze the model up to the last transformer block and use it to extract contextual embeddings
3. **Substitute** the Language Model head by a **Classification** head (e.g., sentiment analysis)



# BERT Evolution Over Time

	BERT (Base)	BERT (Large)	RoBERTa	DistilBERT	ALBERT (Large)
Parameters	110M	340M	355M	66M	18M
Nº of Decoder Blocks	12 encoder layers, 12 attention heads	24 encoder layers, 16 attention heads	24 encoder layers, 16 attention heads	6 encoder layers, 12 attention heads	12 encoder layers, 12 attention heads (shared weights)
Pre-training Dataset	16 GB (~3.3B tokens)	16 GB (~3.3B tokens)	160 GB (~33 B tokens)	16 GB (~3.3B tokens)	16 GB (~3.3B tokens)
Pre-training Dataset	Introduced bidirectional pretraining, greatly improved performance on NLP benchmarks like GLUE, SQuAD, and others. It became a foundation for many downstream tasks.	Same as BERT Base but with higher capacity, resulting in improved performance across NLP tasks, though with greater computational cost.	Tweaked BERT's training process (e.g., longer training), resulting in better performance on NLP benchmarks. Achieved state-of-the-art results on many tasks.	A distilled version of BERT, with 40% fewer parameters and 60% faster inference, while retaining 97% of BERT's performance on downstream tasks. Efficient for real-time applications.	Optimized for parameter efficiency by sharing layers and factorizing embedding parameters. Achieved performance close to BERT Large with significantly fewer parameters.

# BERTimbau

BERTimbau is a bert model especially pretrained for the Brazillian Portuguese language:

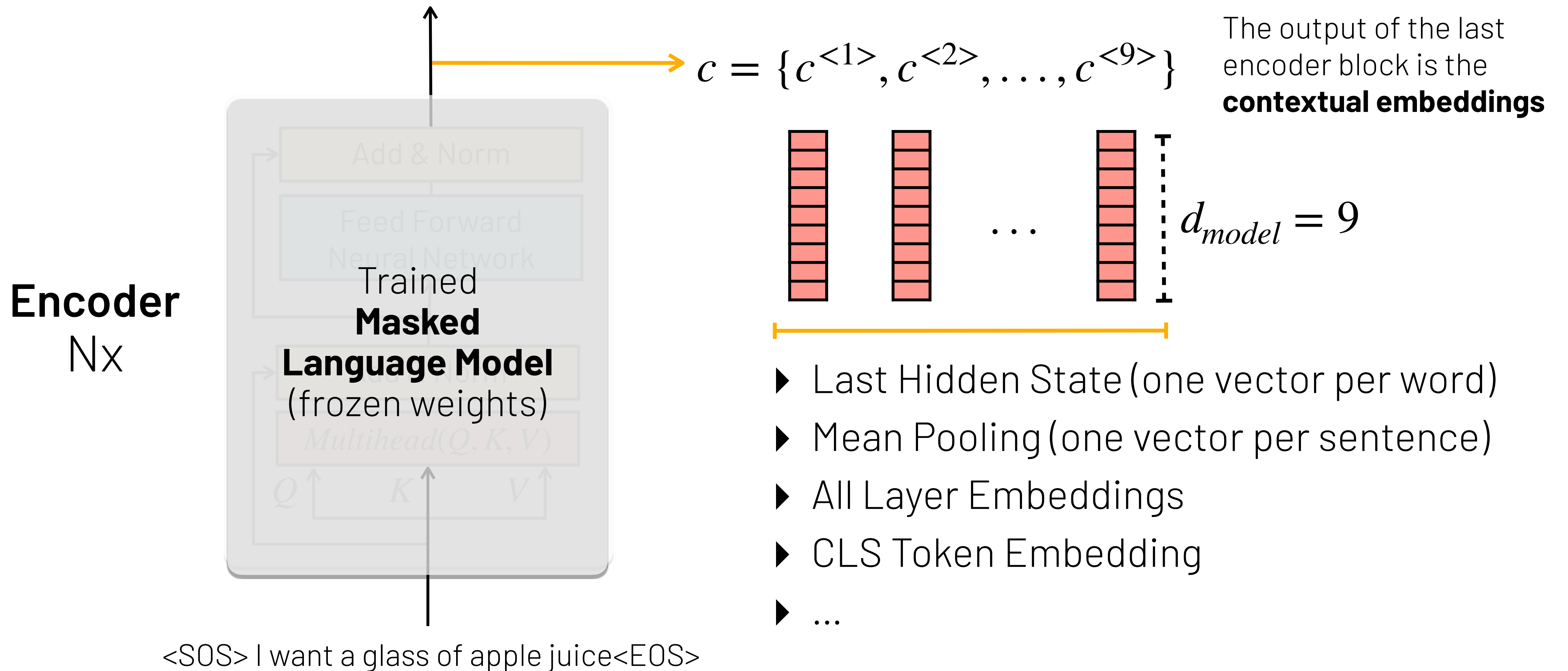
	BERTimbau Base	BERTimbau Large
Parameters	117M	1.5B
Nº of Decoder Blocks	12	48
Pre-training Dataset	5GB (~0.5B Tokens)	40GB (~8B tokens)
Pre-training Dataset	First large-scale pre-trained model for Brazilian Portuguese. Comparable to BERT base for Portuguese language tasks.	Larger version with increased capacity, achieving better performance on tasks requiring more linguistic nuance and complex text understanding.

Vocabulary specifically optimized for Portuguese morphology (and includes Portuguese-specific tokens and accents)



# Extracting BERT Contextual Embeddings

There are different methods to extract contextual word embeddings from BERT



# Next Lecture

## **L19:** Transfer Learning

Exploiting large unlabelled dataset to improve performance of supervised learning models