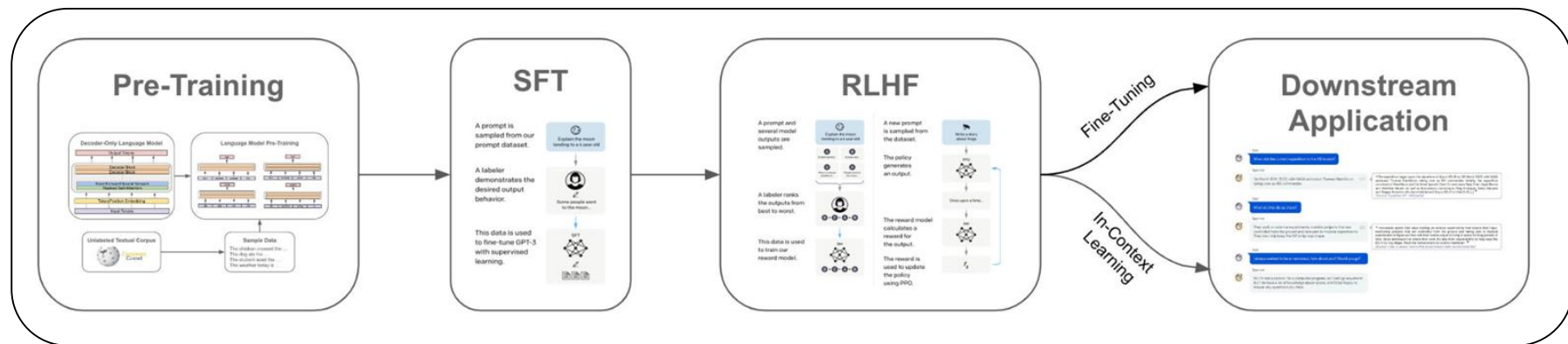


# Large Language Models

Reza Fayyazi

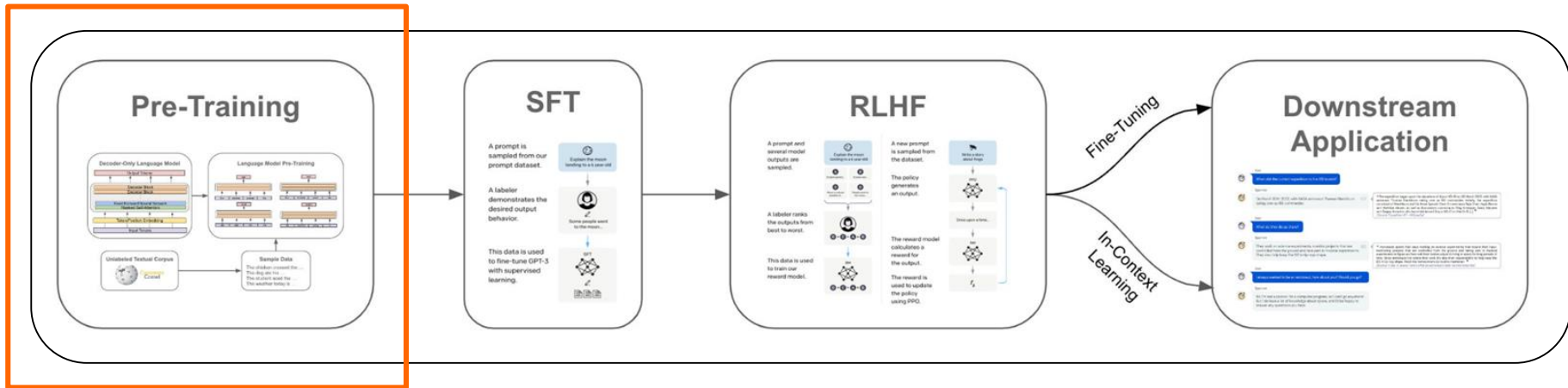
# Introduction

# The overall process of creating your own LLM

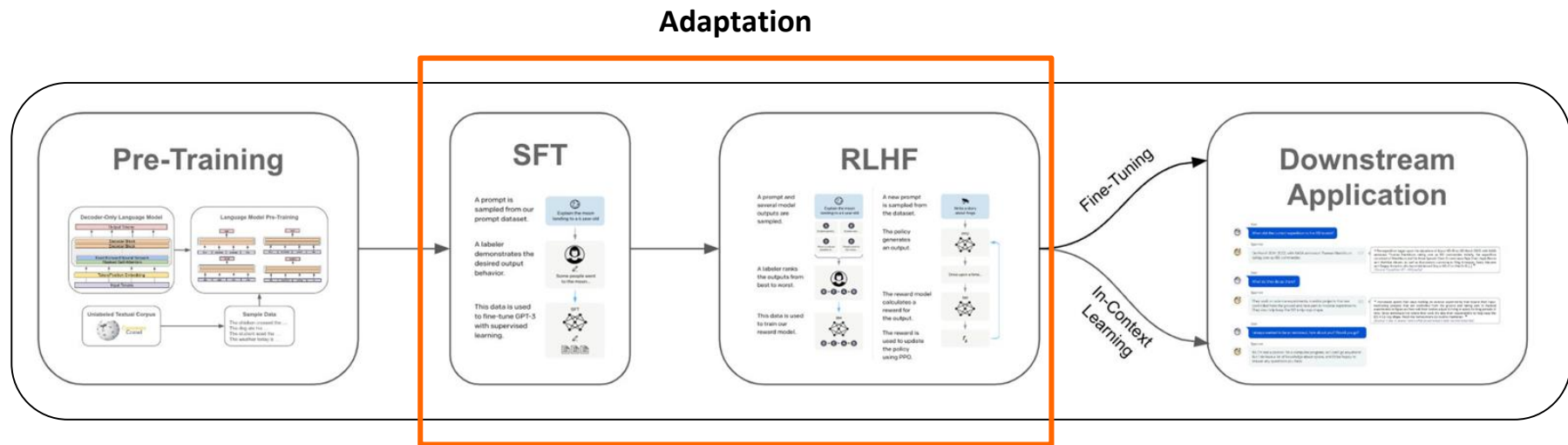


# The overall process of creating your own LLM

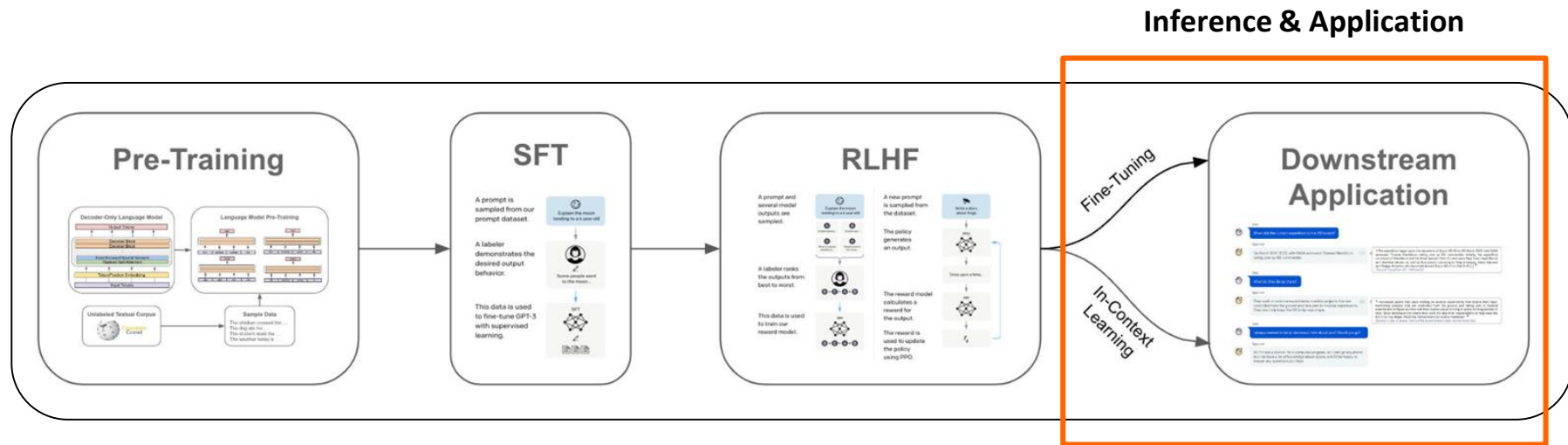
## Foundation Model



# The overall process of creating your own LLM



# The overall process of creating your own LLM



# Transformers

## Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

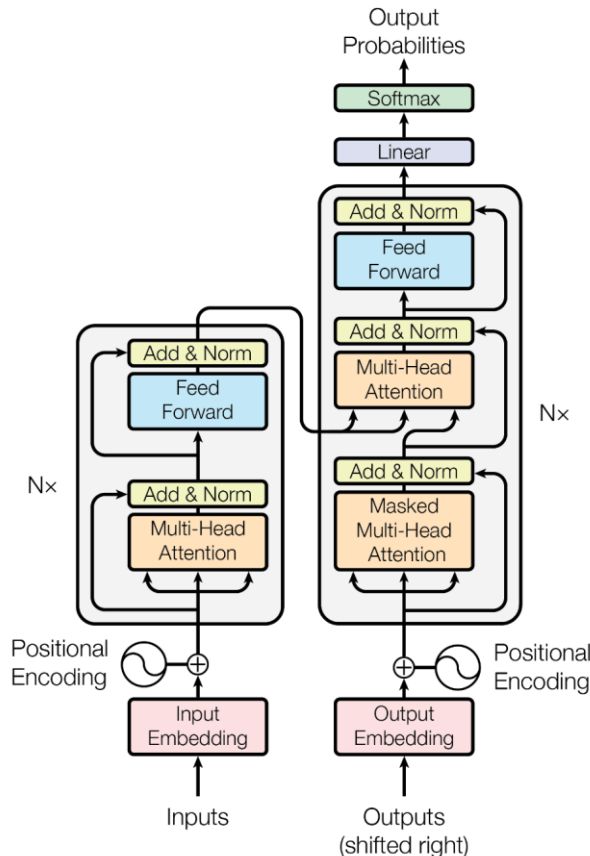
Aidan N. Gomez\*<sup>†</sup>  
University of Toronto  
aidan@cs.toronto.edu

Lukasz Kaiser\*  
Google Brain  
lukaszkaier@google.com

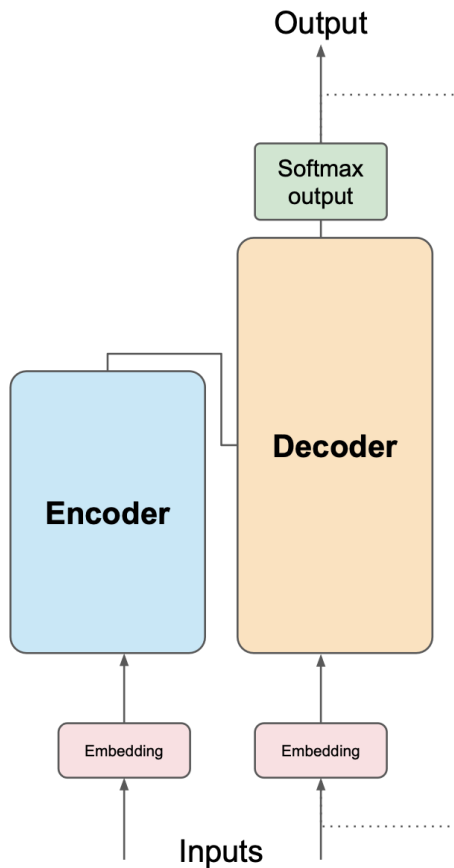
Illia Polosukhin\*<sup>‡</sup>  
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

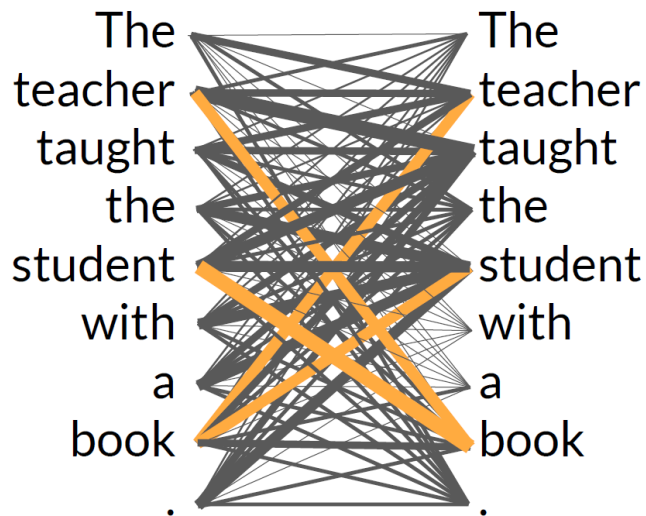


# Transformers

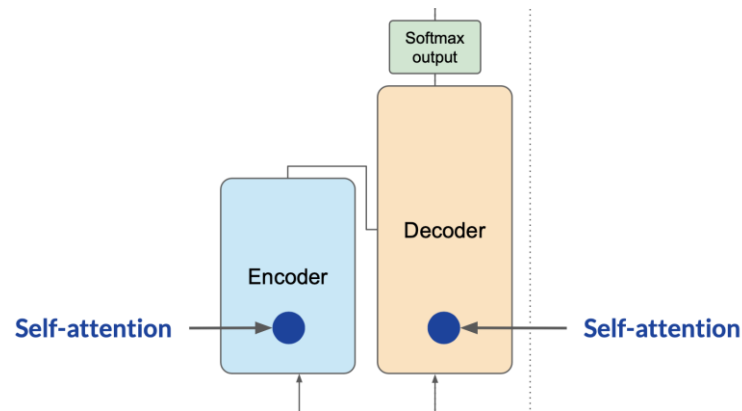




# Self-Attention

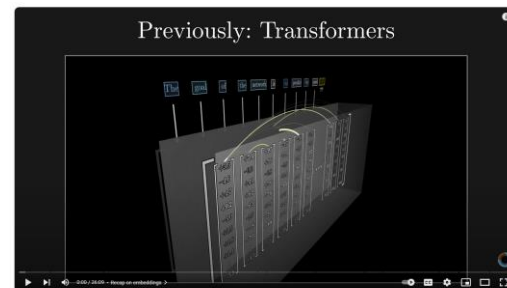


[3]



**Watch this YouTube video on how the self-attention mechanism works:**

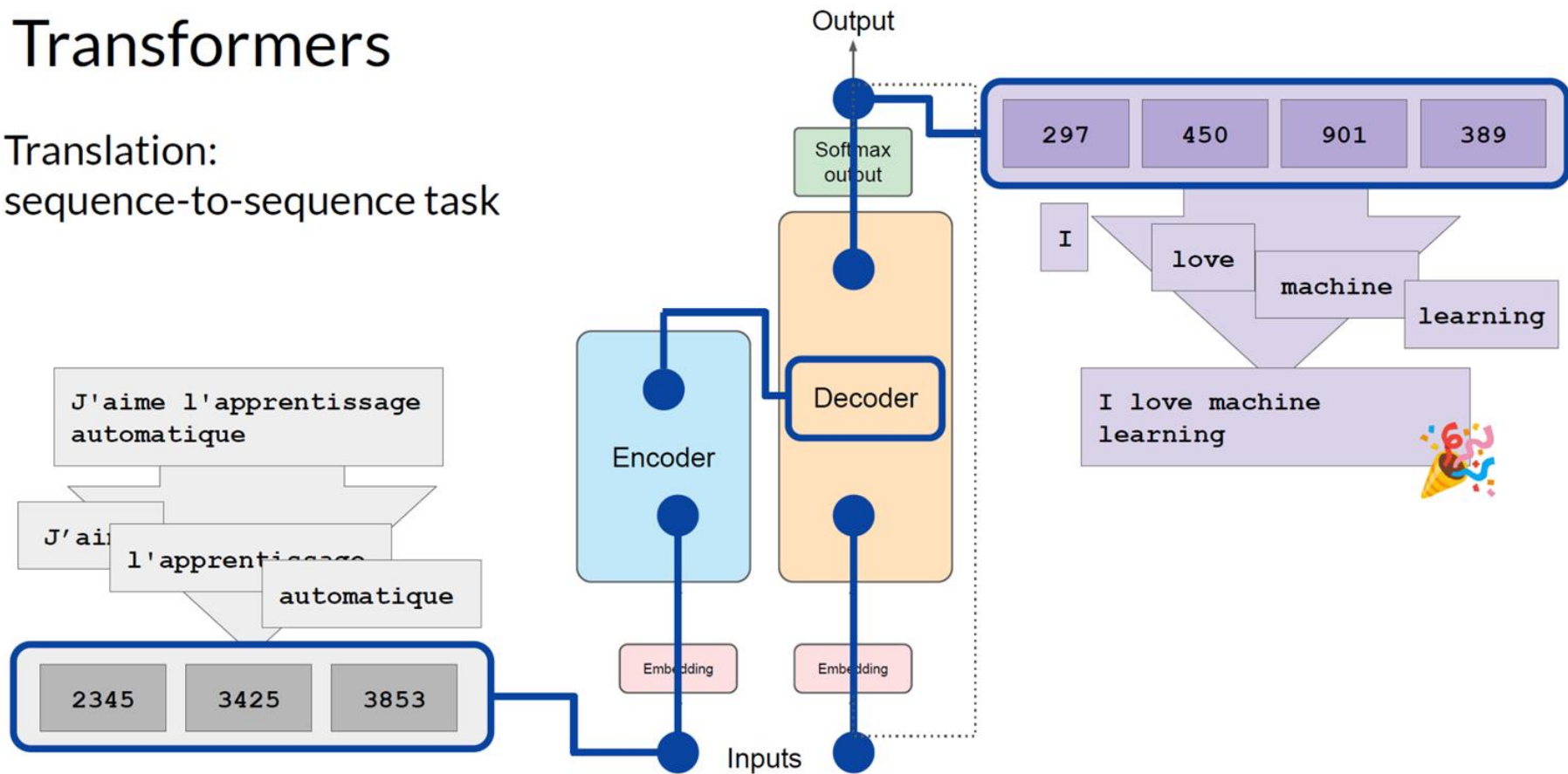
<https://www.youtube.com/watch?app=desktop&v=eMlx5fFNoYc>



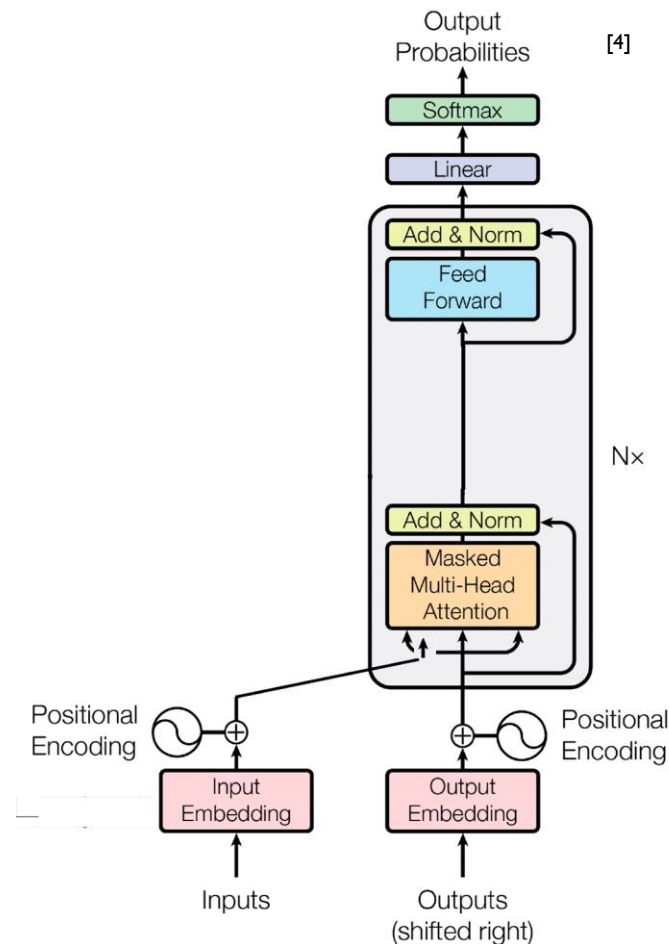
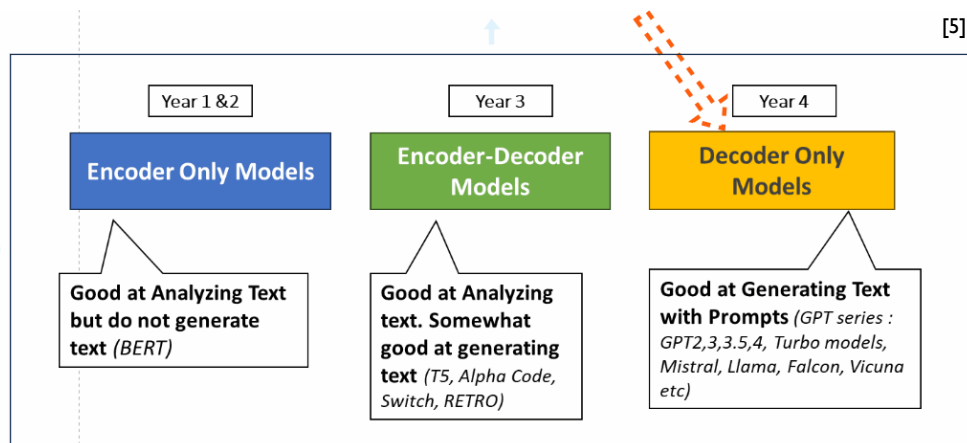
Attention in transformers, visually explained | Chapter 6, Deep Learning

# Transformers

Translation:  
sequence-to-sequence task



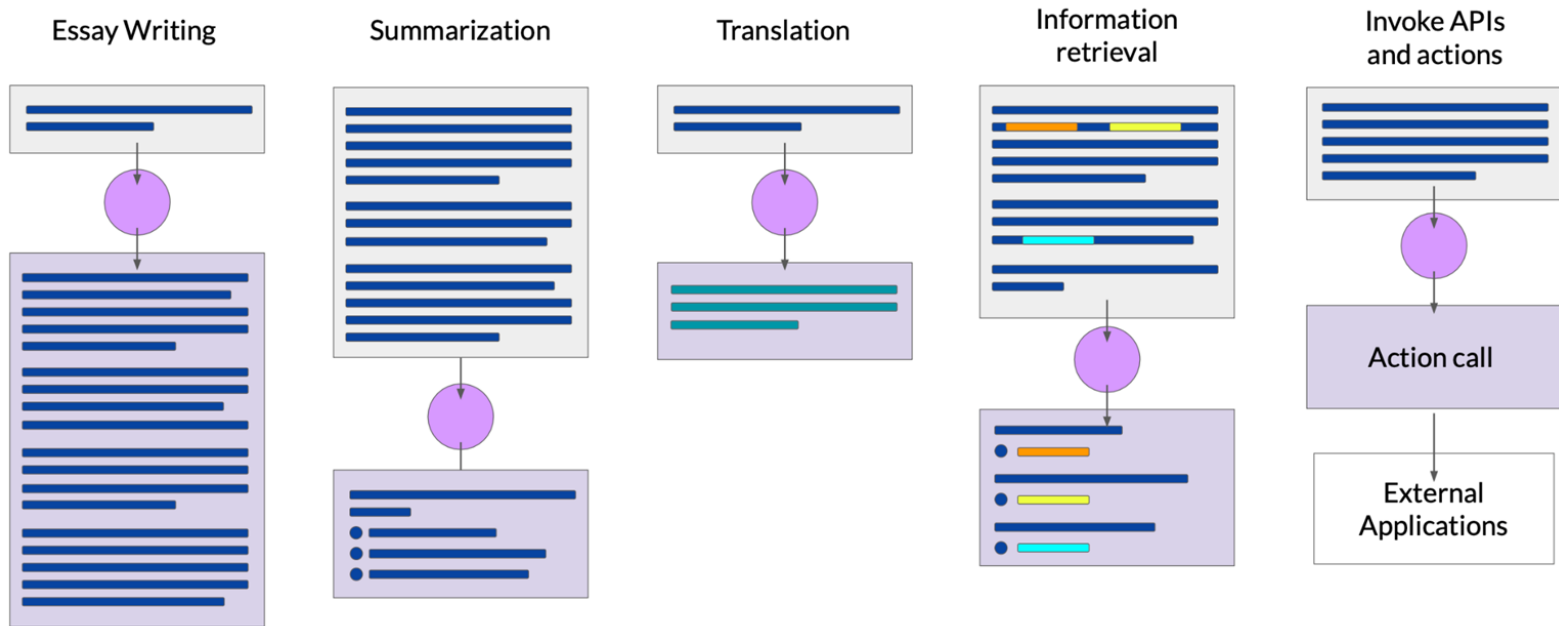
# State-of-the-art LLMs are Decoder-Only Models



[4] <https://stackoverflow.com/questions/75672816/how-does-gpt-like-transformers-utilize-only-the-decoder-to-do-sequence-generatio>

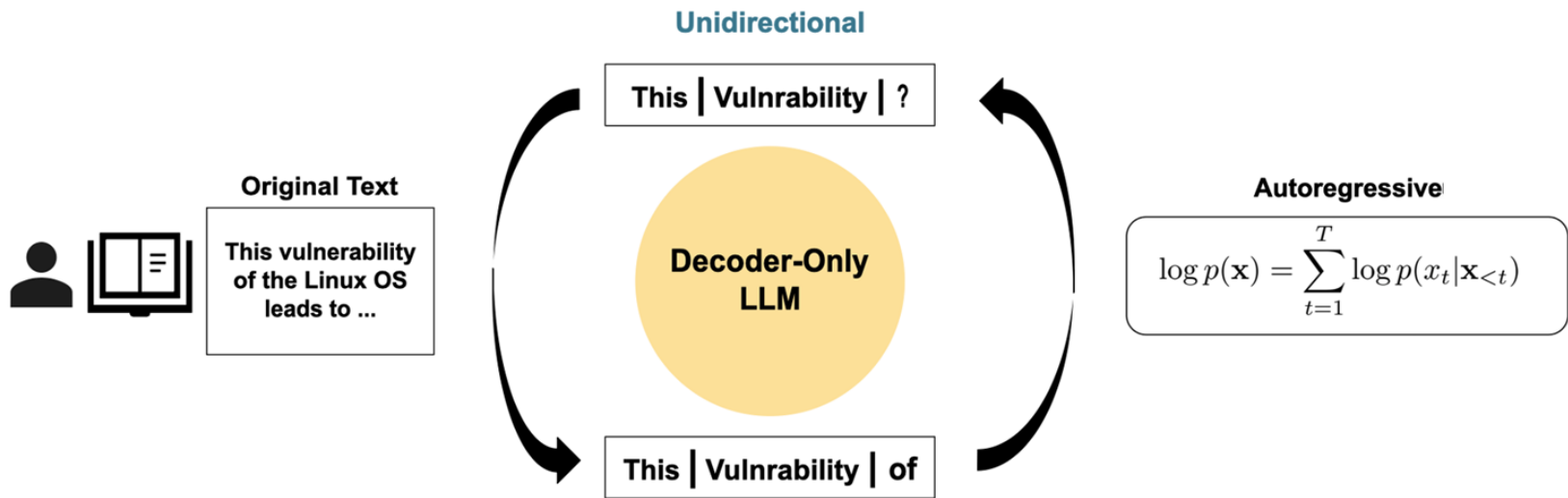
[5] <https://www.linkedin.com/pulse/transformer-architectures-dummies-part-2-decoder-only-qi6vc>

# LLM Use Cases

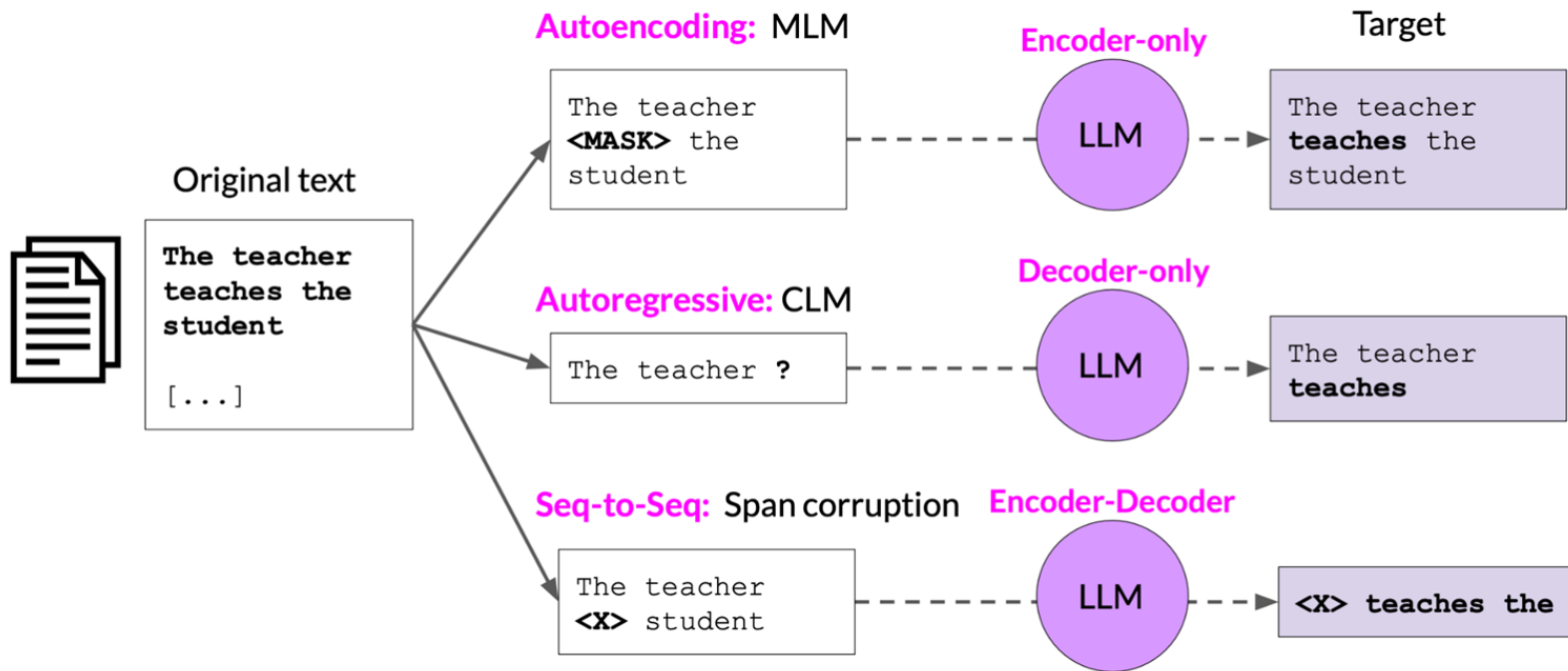


# LLM Pre-Training

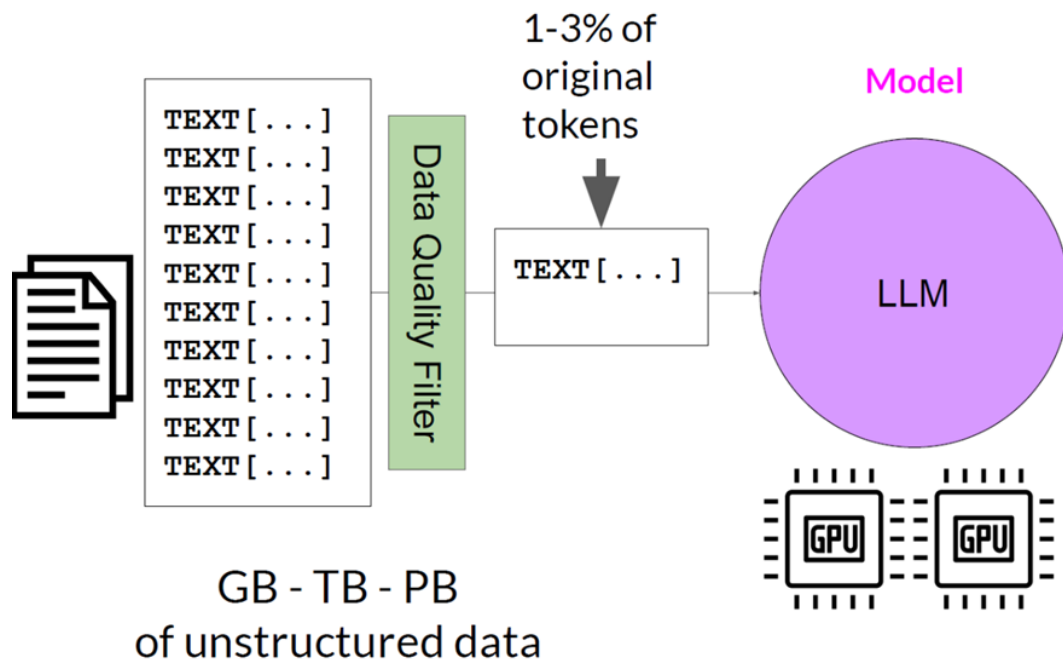
# LLMs work on predicting the next probable token!



# Pre-training Objectives



# LLM Pre-training

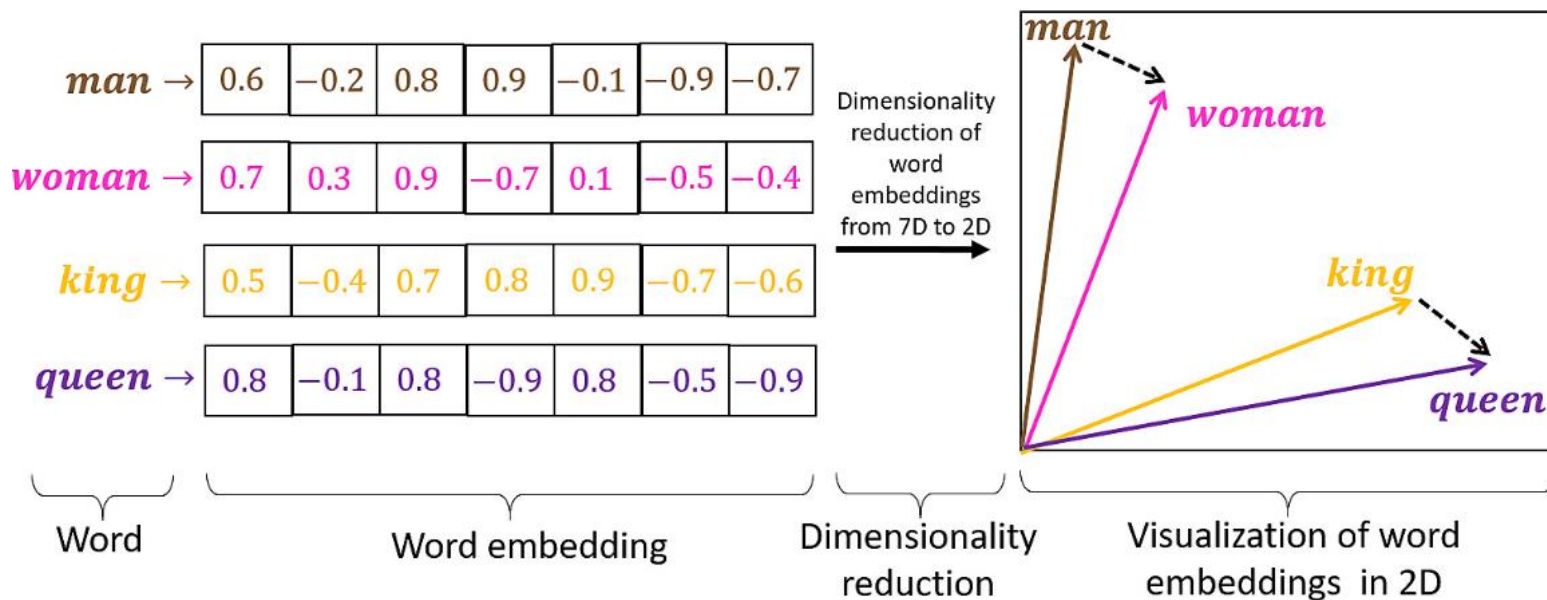


Token String	Token ID	Embedding / Vector Representation
'_The'	37	[-0.0513, -0.0584, 0.0230, ...]
'_teacher'	3145	[-0.0335, 0.0167, 0.0484, ...]
'_teaches'	11749	[-0.0151, -0.0516, 0.0309, ...]
'_the'	8	[-0.0498, -0.0428, 0.0275, ...]
'_student'	1236	[-0.0460, 0.0031, 0.0545, ...]
...	...	...

Vocabulary



# Word Embeddings



# LLM Comparison

Model	Provider	Context Window	Open-Source	Price / Million	Price	Quality	Speed
GPT-4o	OpenAI	128k	No	7.5	★☆☆	★★★	★★☆
GPT-4 Turbo	OpenAI	128k	No	15	★☆☆	★★★	★★☆
GPT-4	OpenAI	8k	No	37.5	★☆☆	★★★	★★☆
GPT-3.5 Turbo	OpenAI	16k	No	0.75	★★☆	★★☆	★★☆
Gemini 1.5 Pro	Google	1m	No	5.25	★★☆	★★★	★★☆
Gemini 1.5 Flash	Google	1m	No	0.53	★★★	★★☆	★★★
Gemma 7B	Google	8k	Yes	0.2	★★★	★★☆	★★★
Claude 3 Opus	Anthropic	200k	No	30	★☆☆	★★★	★★☆
Claude 3 Sonnet	Anthropic	200k	No	6	★☆☆	★★☆	★★☆
Claude 3 Haiku	Anthropic	200k	No	0.5	★★★	★★☆	★★★
Command R +	Cohere	128k	Yes	6	★☆☆	★★☆	★★☆
Command R	Cohere	128k	Yes	0.75	★★☆	★★☆	★★☆
Llama 3 70B	Meta AI	8k	Yes	0.93	★★☆	★★★	★★☆
Llama 3 8B	Meta AI	8k	Yes	0.2	★★★	★★☆	★★★
Code Llama	Meta AI	16k	Yes	0.9	★★☆	★★★	★★☆
Mistral Large	Mistral AI	32k	No	12	★☆☆	★★★	★★☆
Mistral Medium	Mistral AI	32k	No	4.05	★★☆	★★☆	★★☆
Mistral Small	Mistral AI	32k	No	2.25	★★☆	★★☆	★★☆
Mixtral 8x22B	Mistral AI	65k	Yes	1.2	★★☆	★★☆	★★☆
Mixtral 8x7B	Mistral AI	32k	Yes	0.5	★★★	★★☆	★★★
Mistral 7B	Mistral AI	32k	Yes	0.2	★★★	★★☆	★★☆
DBRX	Databricks	32k	Yes	1.4	★★☆	★★☆	★★★

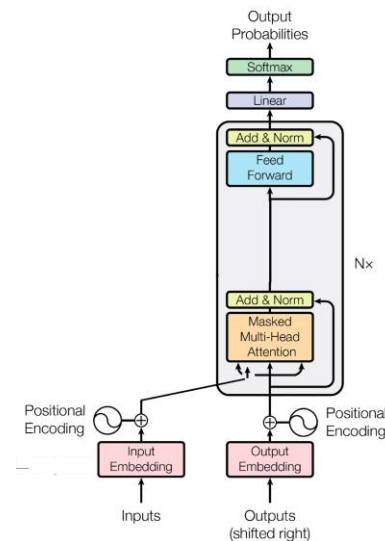
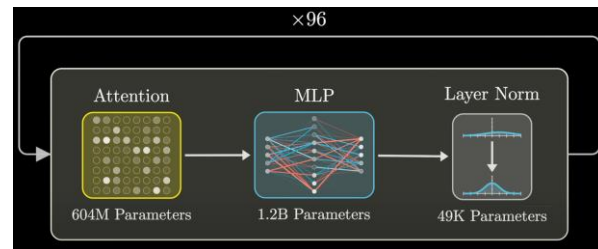
# High Number of Parameters



Total weights:

175,181,291,520

Embedding	$12,288 \times 50,257$ $d_{\text{embed}} * n_{\text{vocab}}$	$= 617,558,016$
Key	$128 \times 12,288 \times 96 \times 96$ $d_{\text{query}} * d_{\text{embed}} * n_{\text{heads}} * n_{\text{layers}}$	$= 14,495,514,624$
Query	$128 \times 12,288 \times 96 \times 96$ $d_{\text{query}} * d_{\text{embed}} * n_{\text{heads}} * n_{\text{layers}}$	$= 14,495,514,624$
Value	$128 \times 12,288 \times 96 \times 96$ $d_{\text{value}} * d_{\text{embed}} * n_{\text{heads}} * n_{\text{layers}}$	$= 14,495,514,624$
Output	$12,288 \times 128 \times 96 \times 96$ $d_{\text{embed}} * d_{\text{value}} * n_{\text{heads}} * n_{\text{layers}}$	$= 14,495,514,624$
Up-projection	$49,152 \times 12,288 \times 96$ $n_{\text{neurons}} * d_{\text{embed}} * n_{\text{layers}}$	$= 57,982,058,496$
Down-projection	$12,288 \times 49,152 \times 96$ $d_{\text{embed}} * n_{\text{neurons}} * n_{\text{layers}}$	$= 57,982,058,496$
Unembedding	$50,257 \times 12,288$ $n_{\text{vocab}} * d_{\text{embed}}$	$= 617,558,016$



# Computational Challenges

[6]

## CUDA Out of Memory

```
RuntimeError: CUDA out of memory. Tried to allocate 200.00 MiB (GPU 0; 15.78 GiB total capacity; 14.56 GiB already allocated; 38.44 MiB free; 14.80 GiB reserved in total by PyTorch) If reserved memory is >> allocated memory try setting max_split_size_mb to avoid fragmentation. See documentation for Memory Management and PYTORCH_CUDA_ALLOC_CONF
```

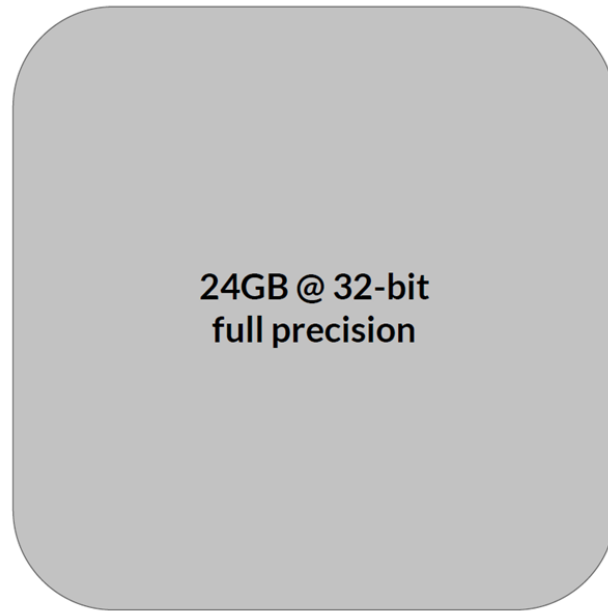
# Computational Challenges for 1B Param Model

Memory needed to store model



**4GB @ 32-bit  
full precision**

Memory needed to train model



**24GB @ 32-bit  
full precision**

# Computational Challenges

**1B param  
model**

**175B param  
model**

**500B param  
model**

4,200 GB @ 32-bit  
full precision

12,000 GB @ 32-bit  
full precision



# Computational Challenges

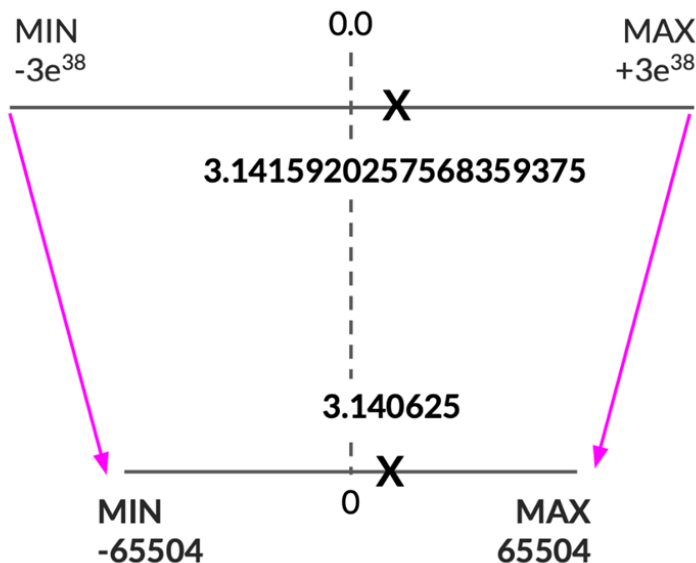
	Bytes per parameter
<b>Model Parameters (Weights)</b>	4 bytes per parameter
<b>Adam optimizer (2 states)</b>	+8 bytes per parameter
<b>Gradients</b>	+4 bytes per parameter
<b>Activations and temp memory (variable size)</b>	+8 bytes per parameter (high-end estimate)
<b>TOTAL</b>	<b>=4 bytes per parameter +20 extra bytes per parameter</b>

[3] <https://www.coursera.org/learn/generative-ai-with-llms>

[11] [https://huggingface.co/docs/transformers/v4.20.1/en/perf\\_train\\_gpu\\_one#anatomy-of-models-memory](https://huggingface.co/docs/transformers/v4.20.1/en/perf_train_gpu_one#anatomy-of-models-memory)

[12] <https://www.youtube.com/live/g68qlo9lf0>

# One Solution: Quantization from FP32 to FP16



**FP32** 4 bytes memory

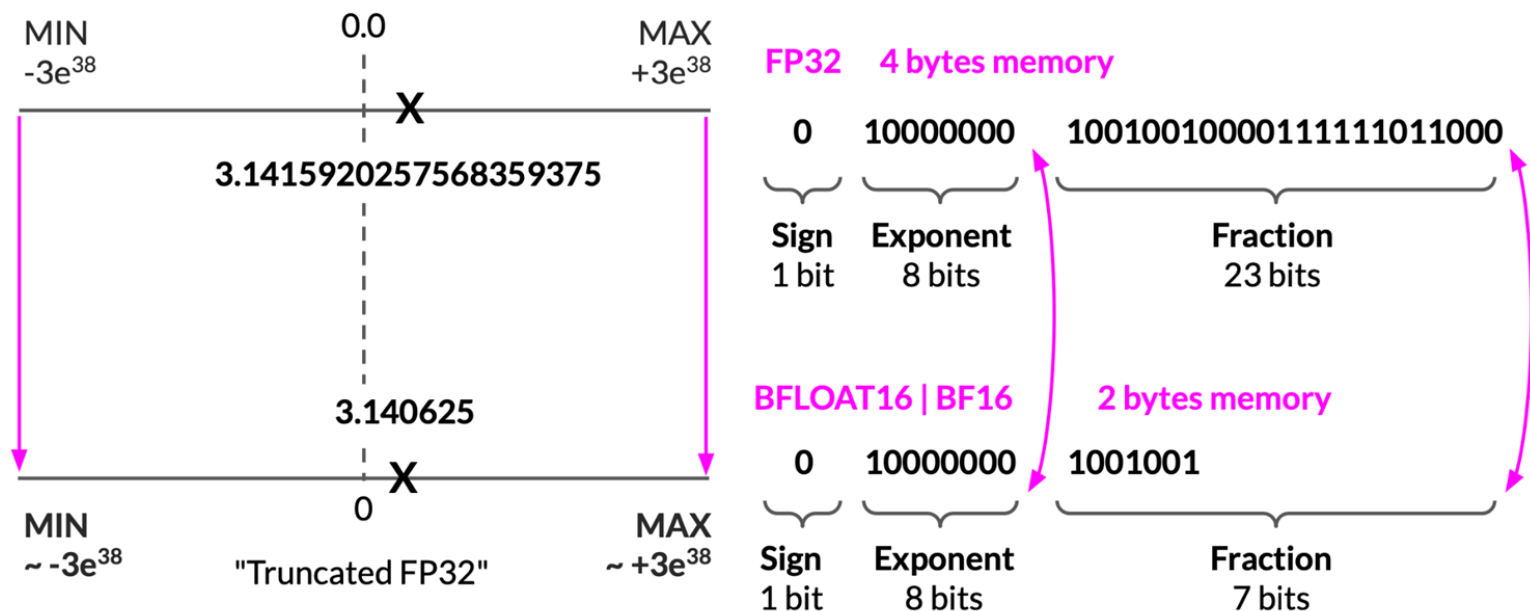
0	10000000	10010010000111111011000
<hr/>		
Sign	Exponent	Fraction
1 bit	8 bits	23 bits

**FP16** 2 bytes memory

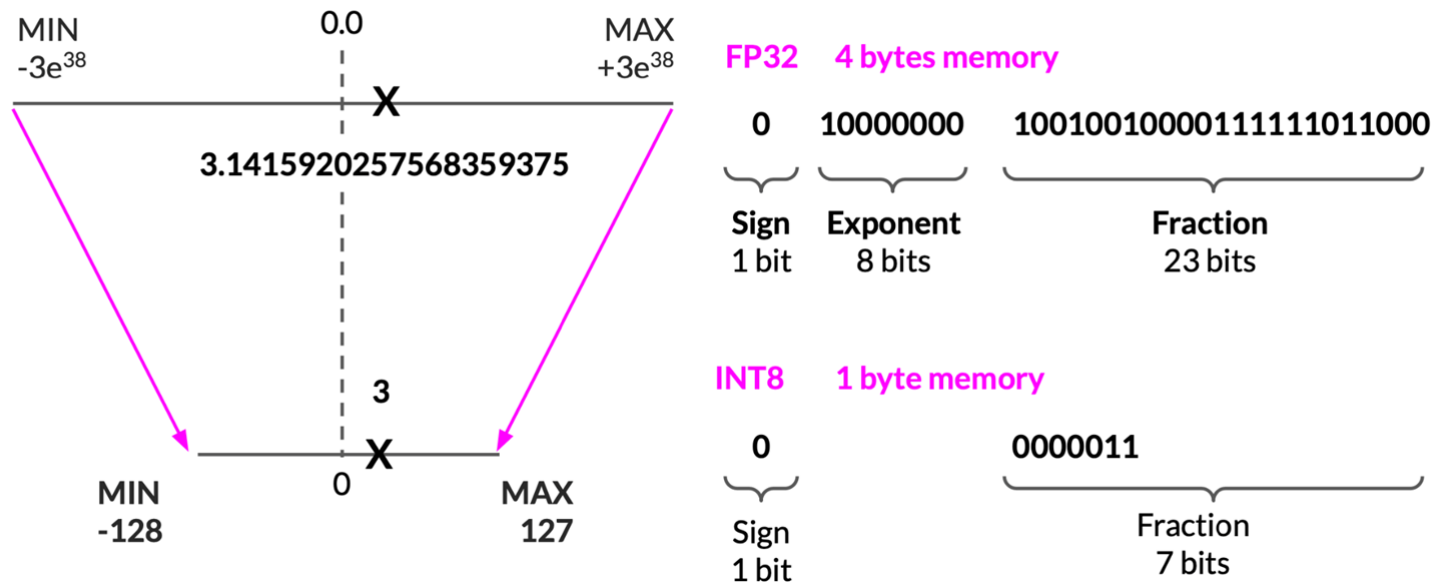
0	10000	1001001000
<hr/>		
Sign	Exponent	Fraction
1 bit	5 bits	10 bits



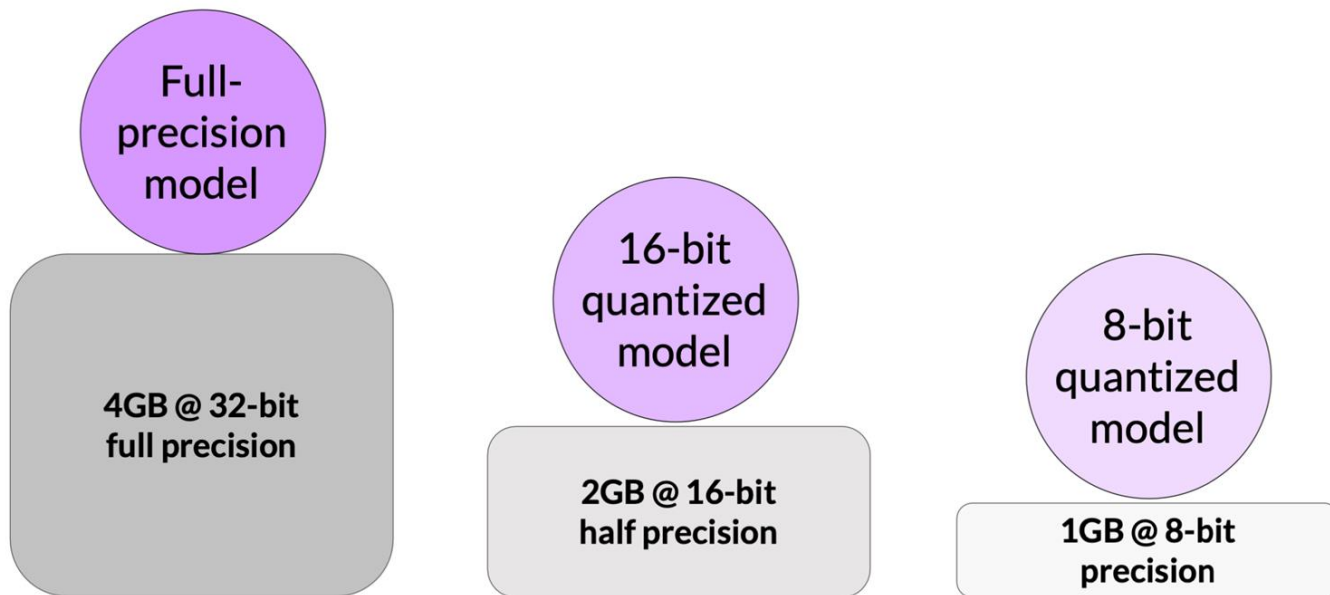
# One Solution: Quantization from FP32 to BFLOAT16



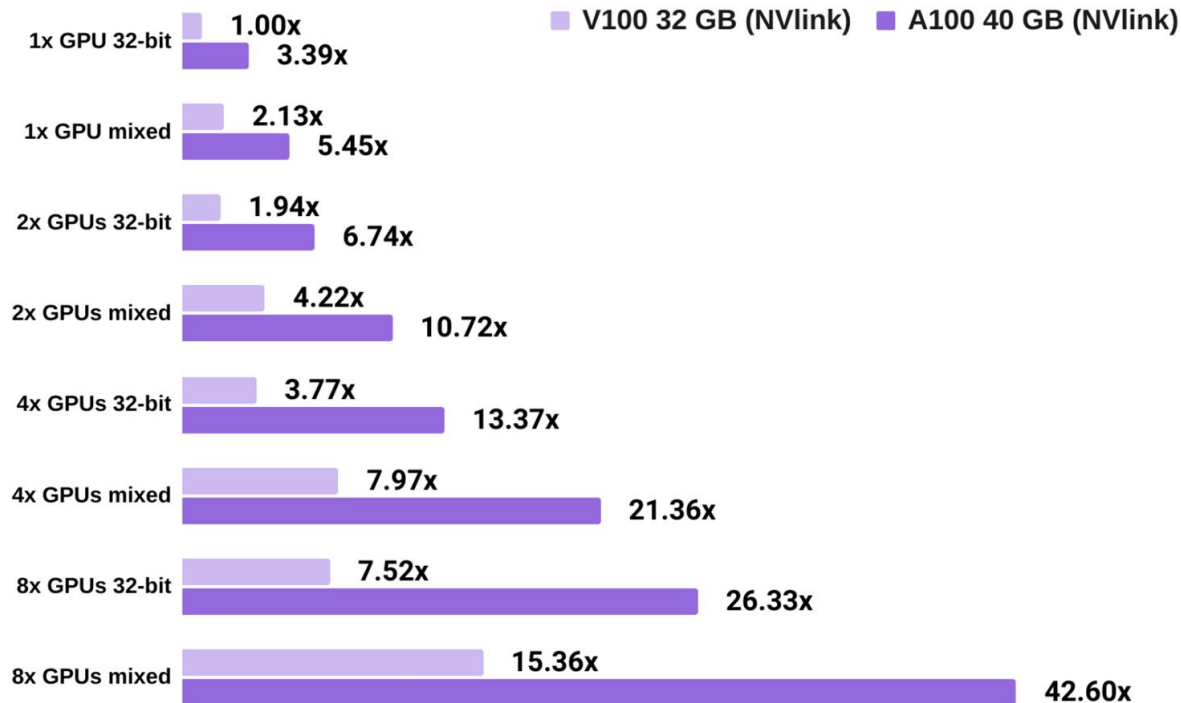
# One Solution: Quantization from FP32 to INT8



# Approximate GPU Needed for a 1B Parameter



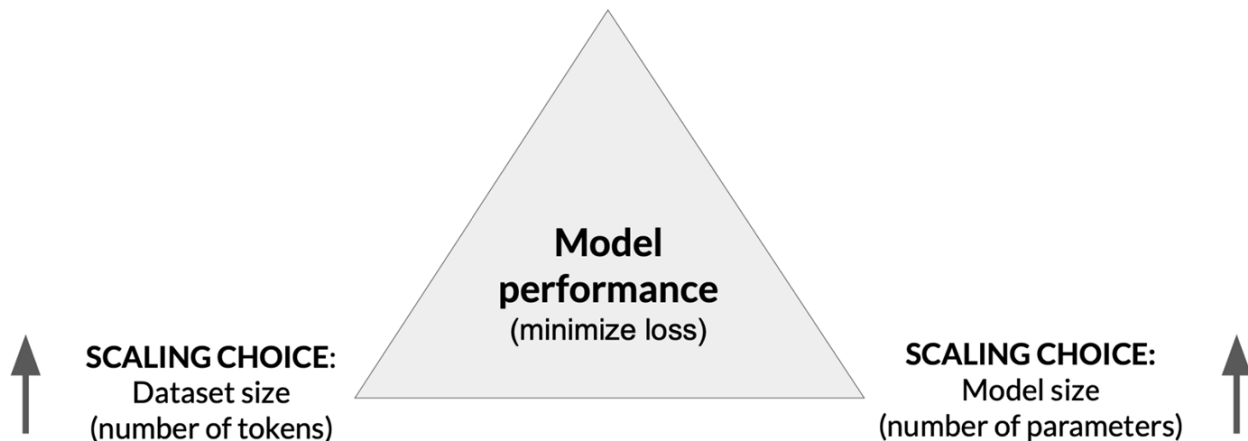
# LLM Training Speed of V100 and A100



# Scaling Choices for Pre-training

**Goal: maximize model performance**

**CONSTRAINT:**  
Compute budget  
(GPUs, training time, cost)



# Thank you!

Reza Fayyazi  
rf1679@rit.edu