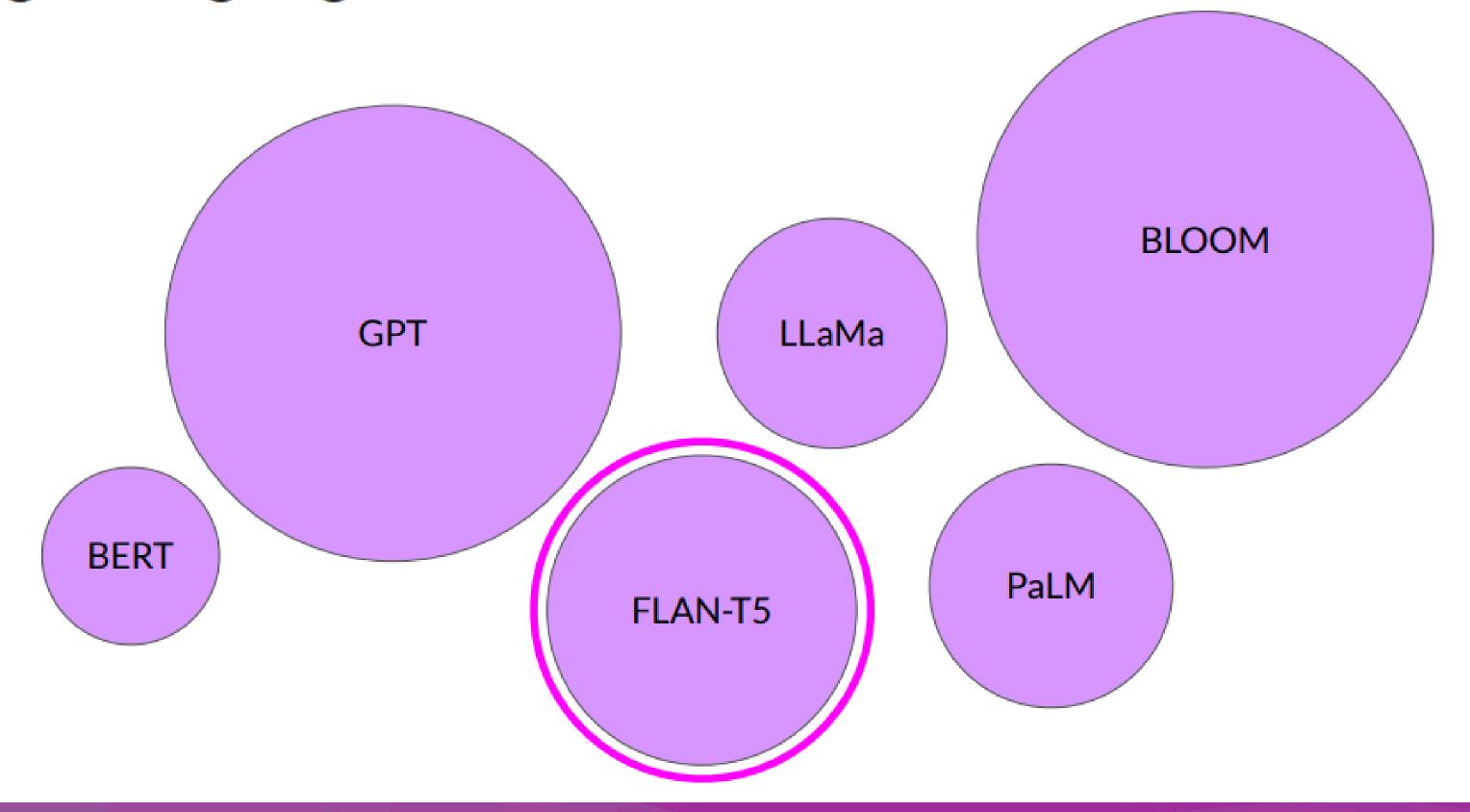
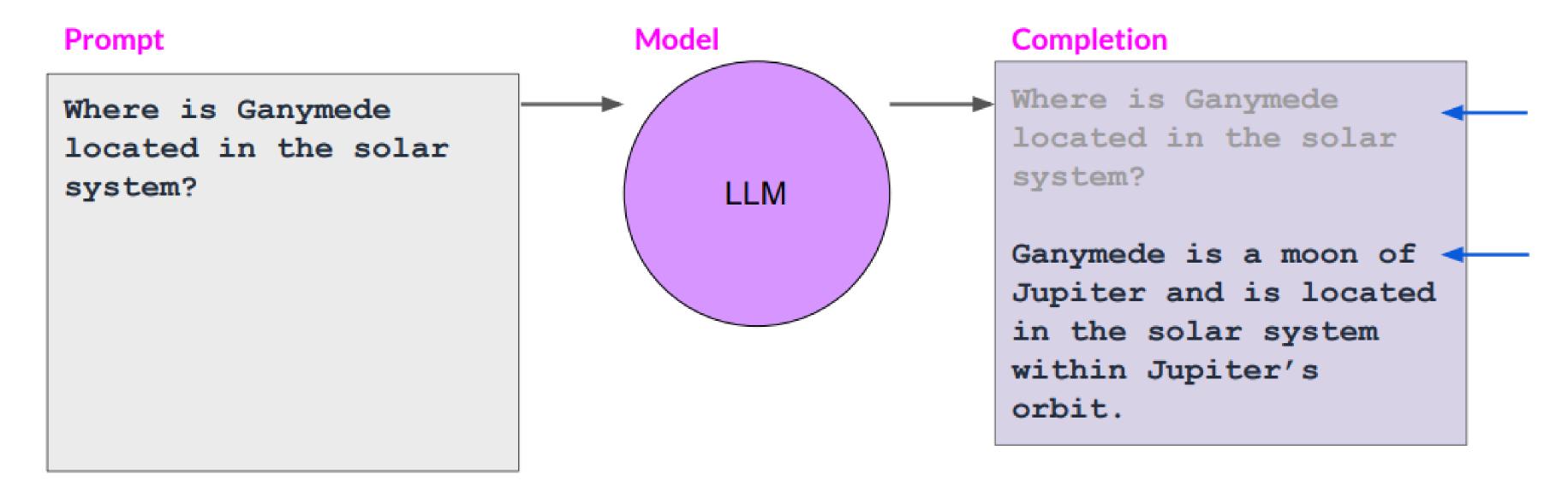
Large Language Models







Prompts and completions

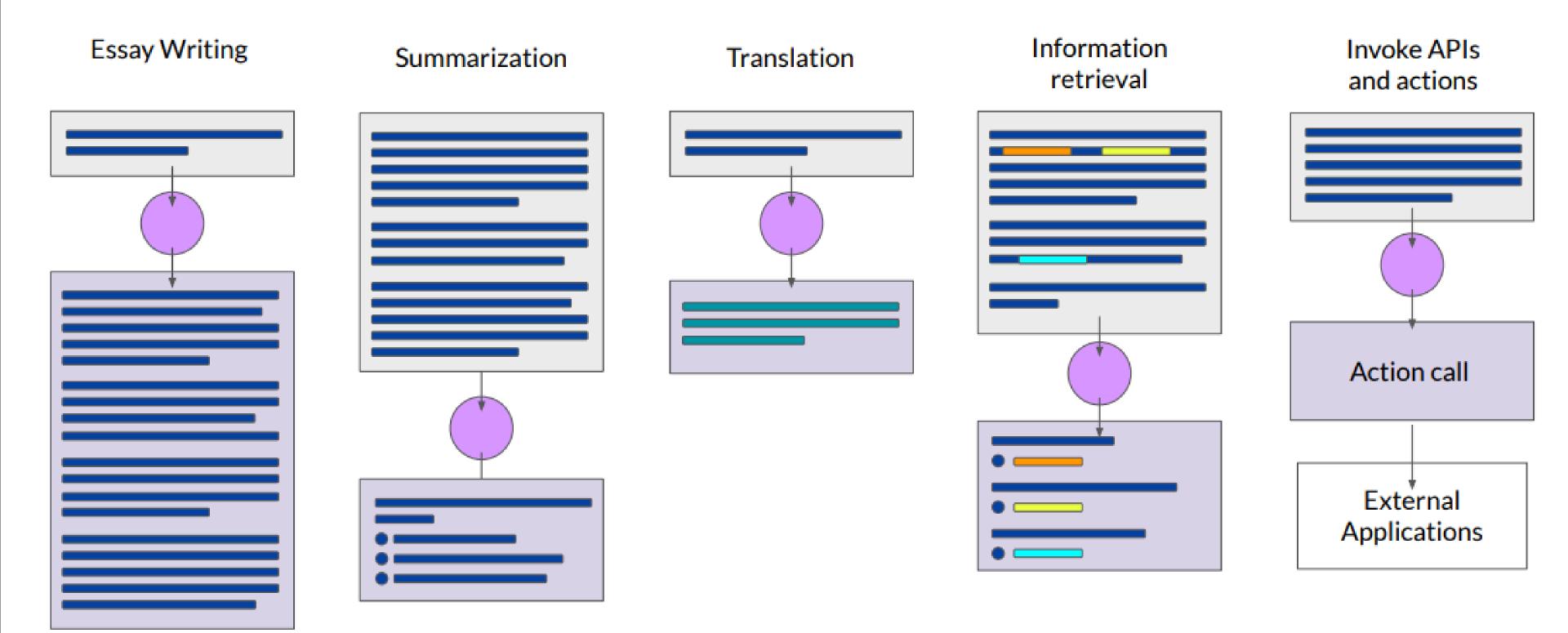


Context window

 typically a few 1000 words.



LLM use cases & tasks





The significance of scale: language understanding



BLOOM ___

*Bert-base

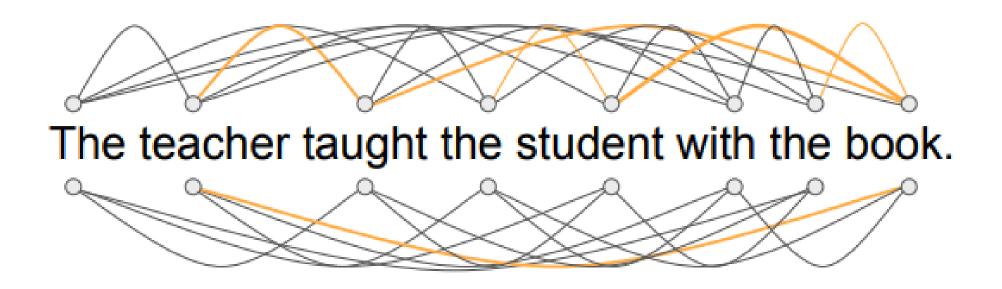


How LLMs work -Transformers architecture



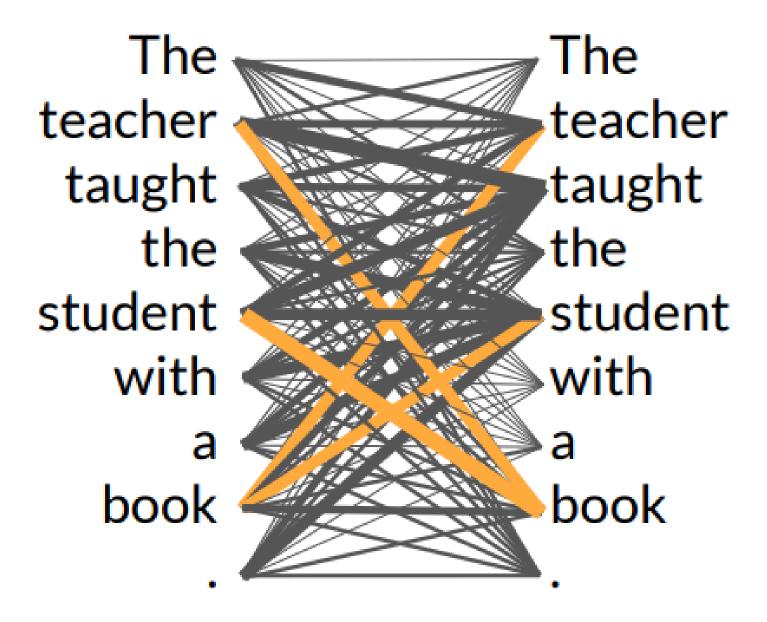


Transformers



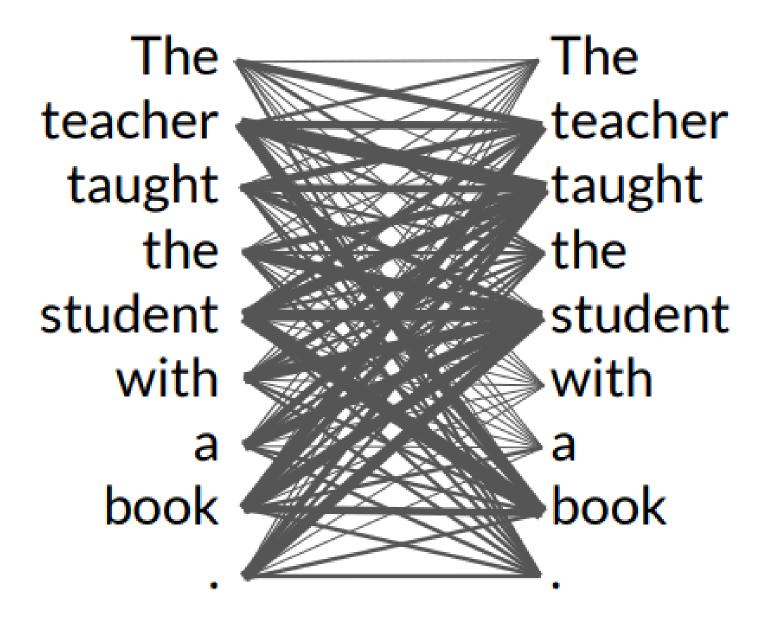


Self-attention





Self-attention



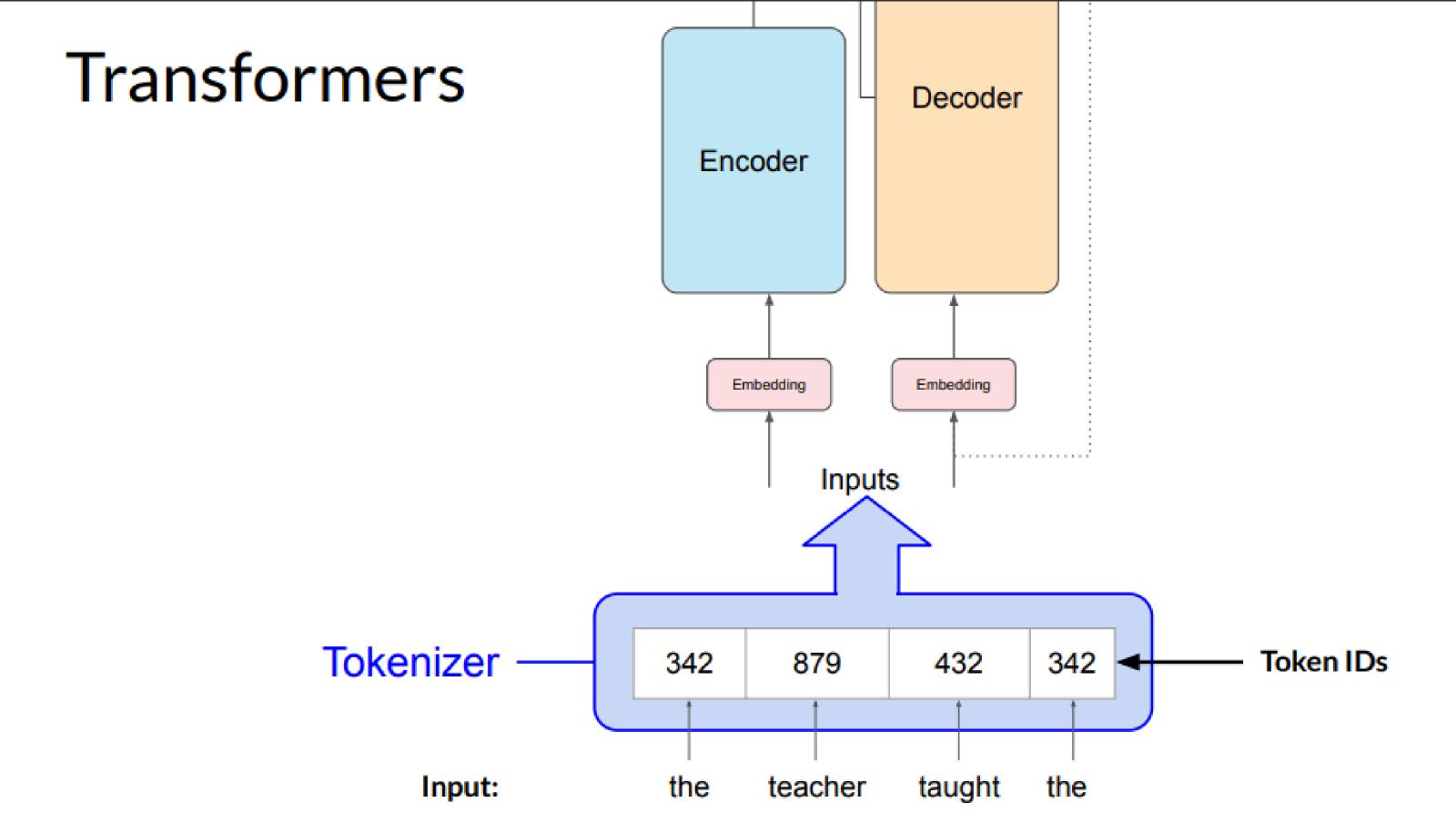


Output **Transformers** Softmax output Decoder Encoder Embedding Embedding

Inputs



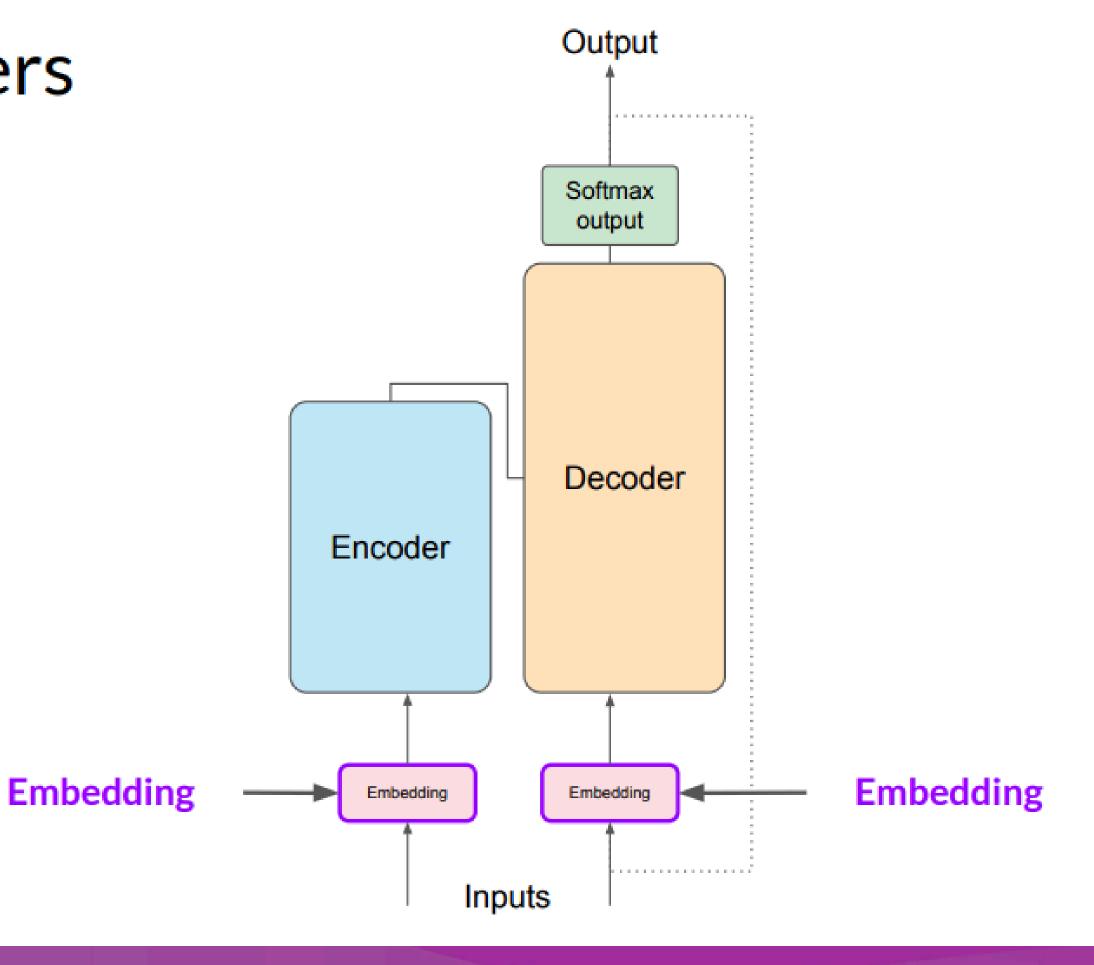








Transformers





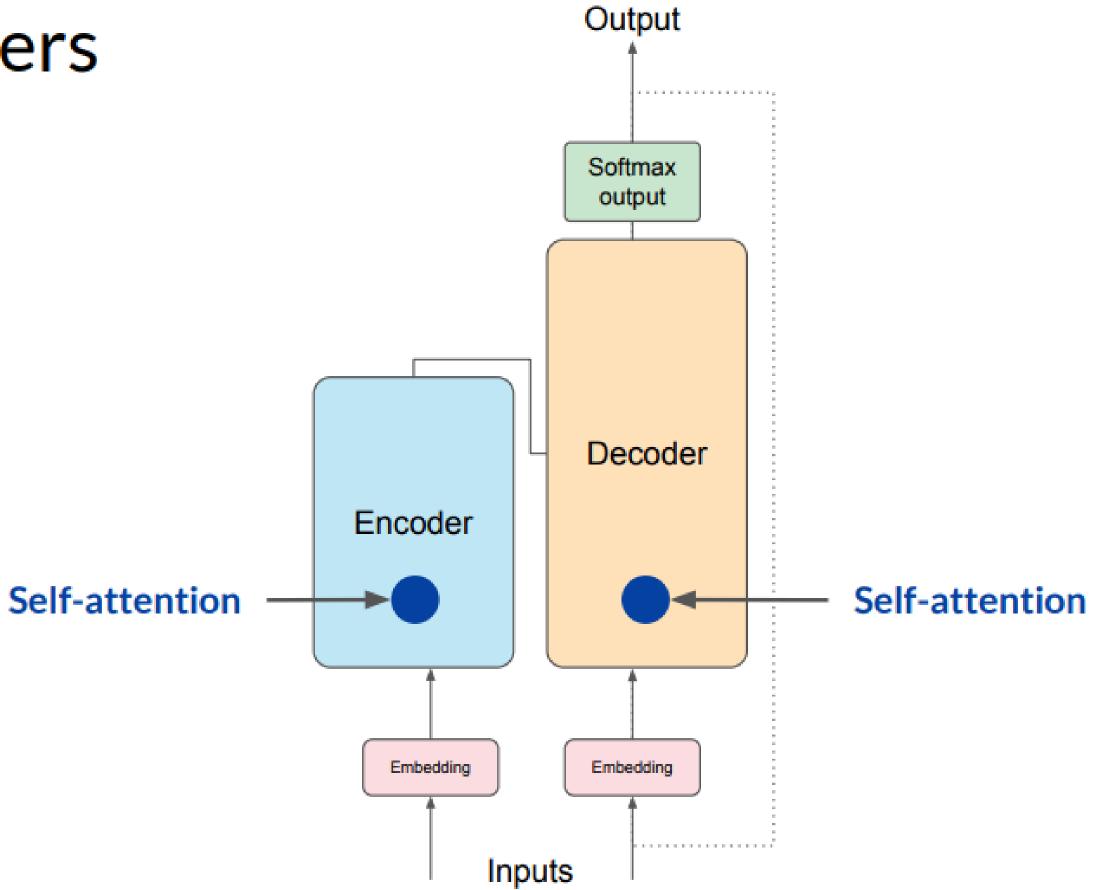


Output **Transformers** X₁ X_3 X_4 X_2 e.g. 512 342 879 432 342 **Embedding Embedding** Embedding Inputs



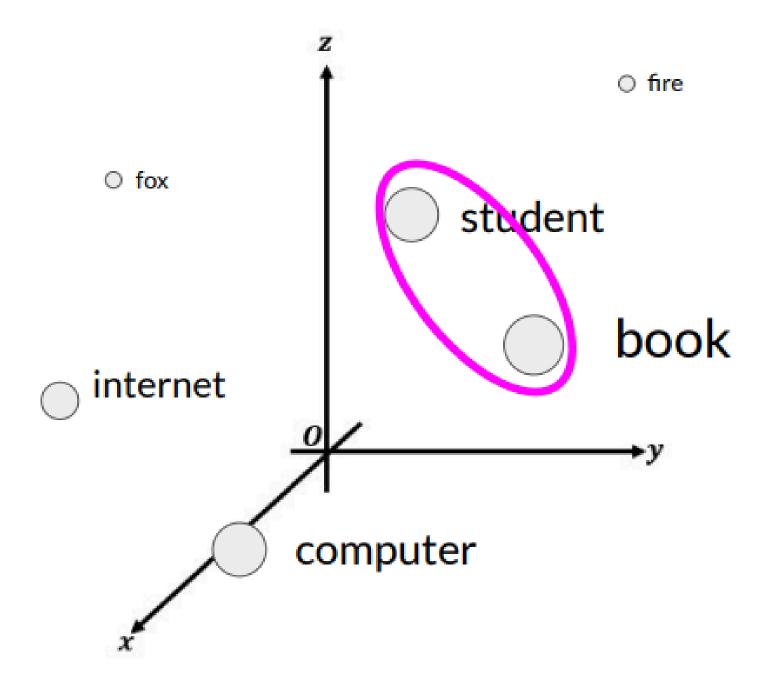


Transformers





Transformers







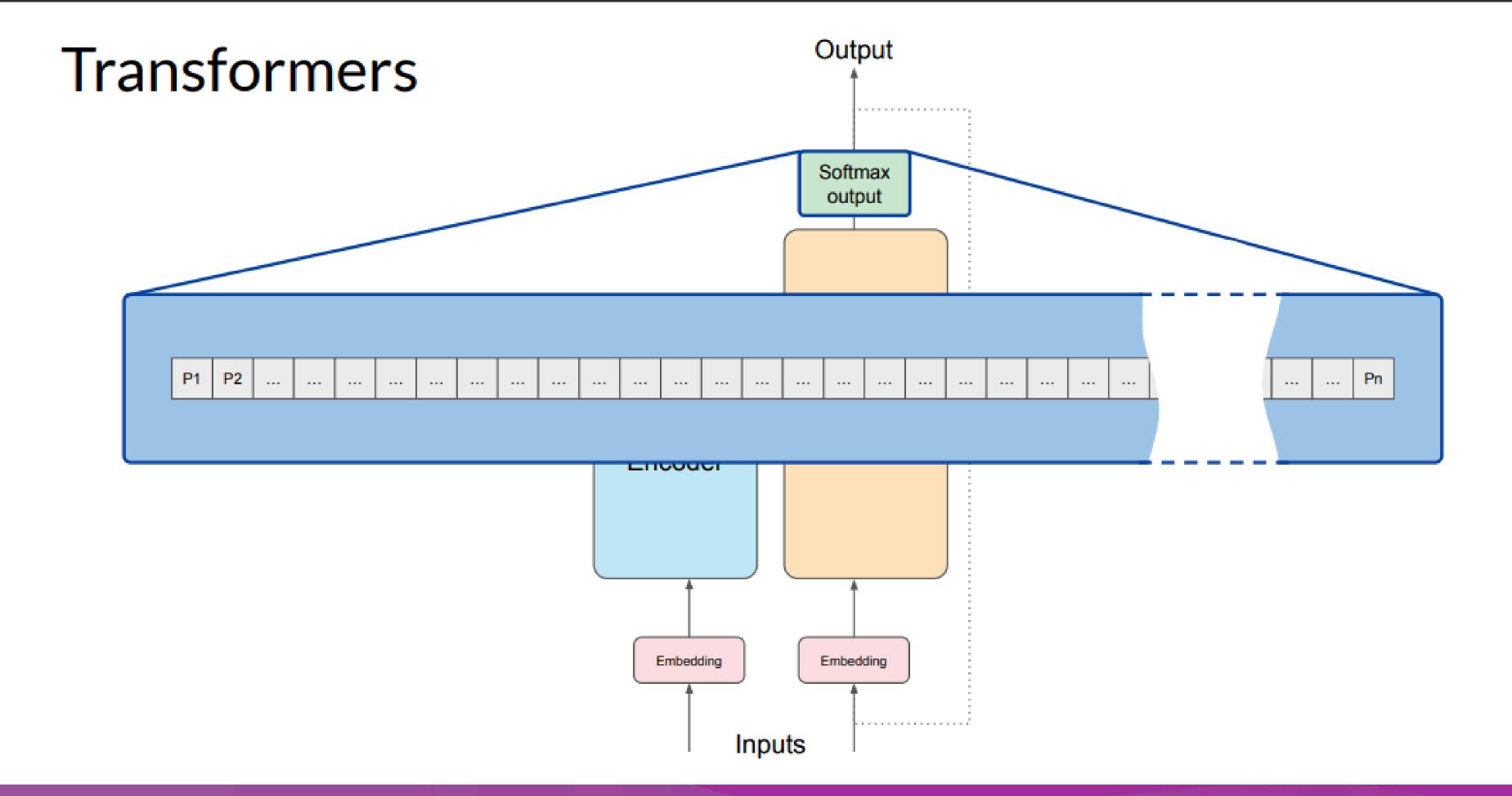
Output **Transformers** Softmax output Multi-headed Multi-headed **Self-attention Self-attention** Embedding Embedding Inputs



Output **Transformers** Softmax output **Feed forward** network **Feed forward** Decoder network Encoder Embedding Embedding

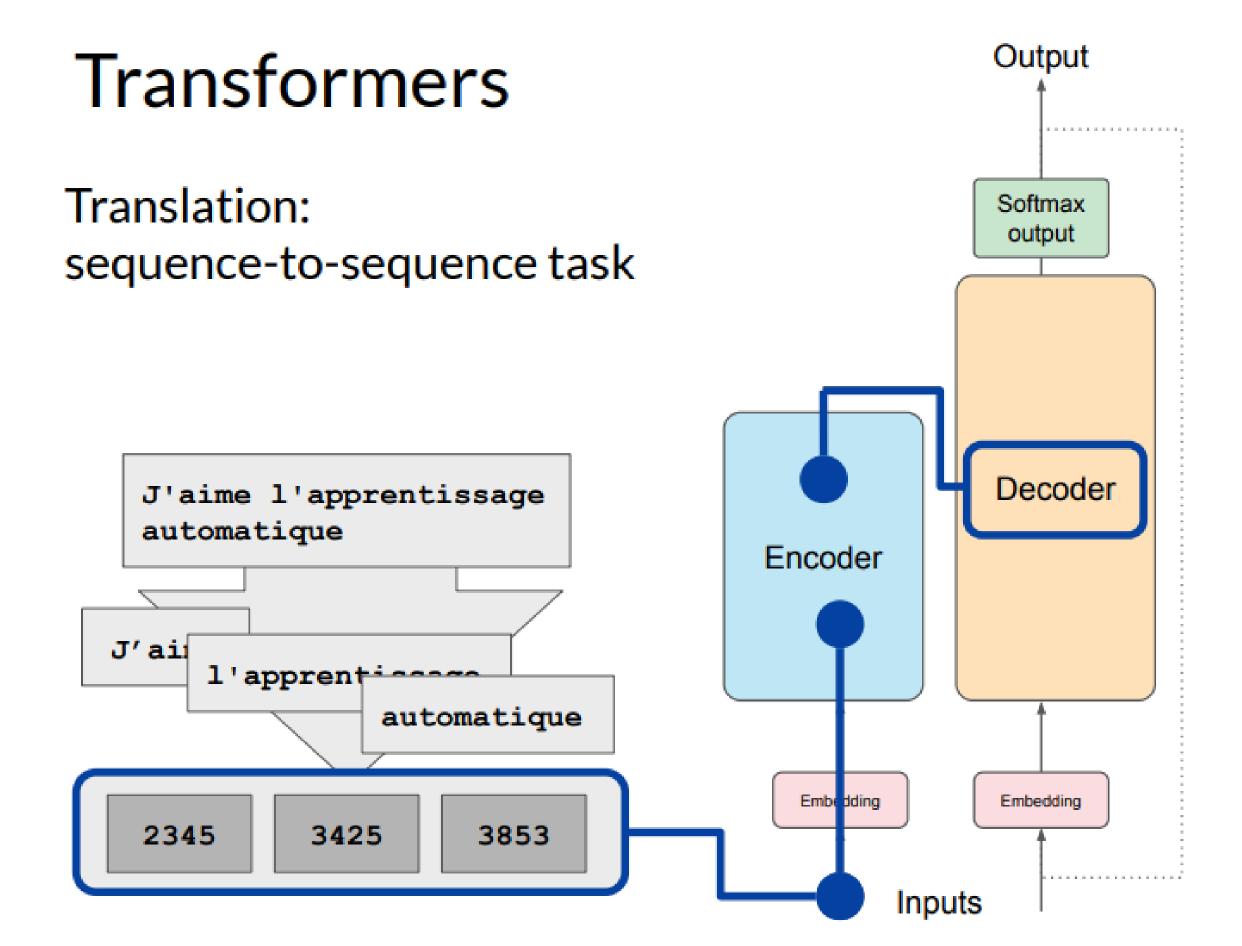
Inputs





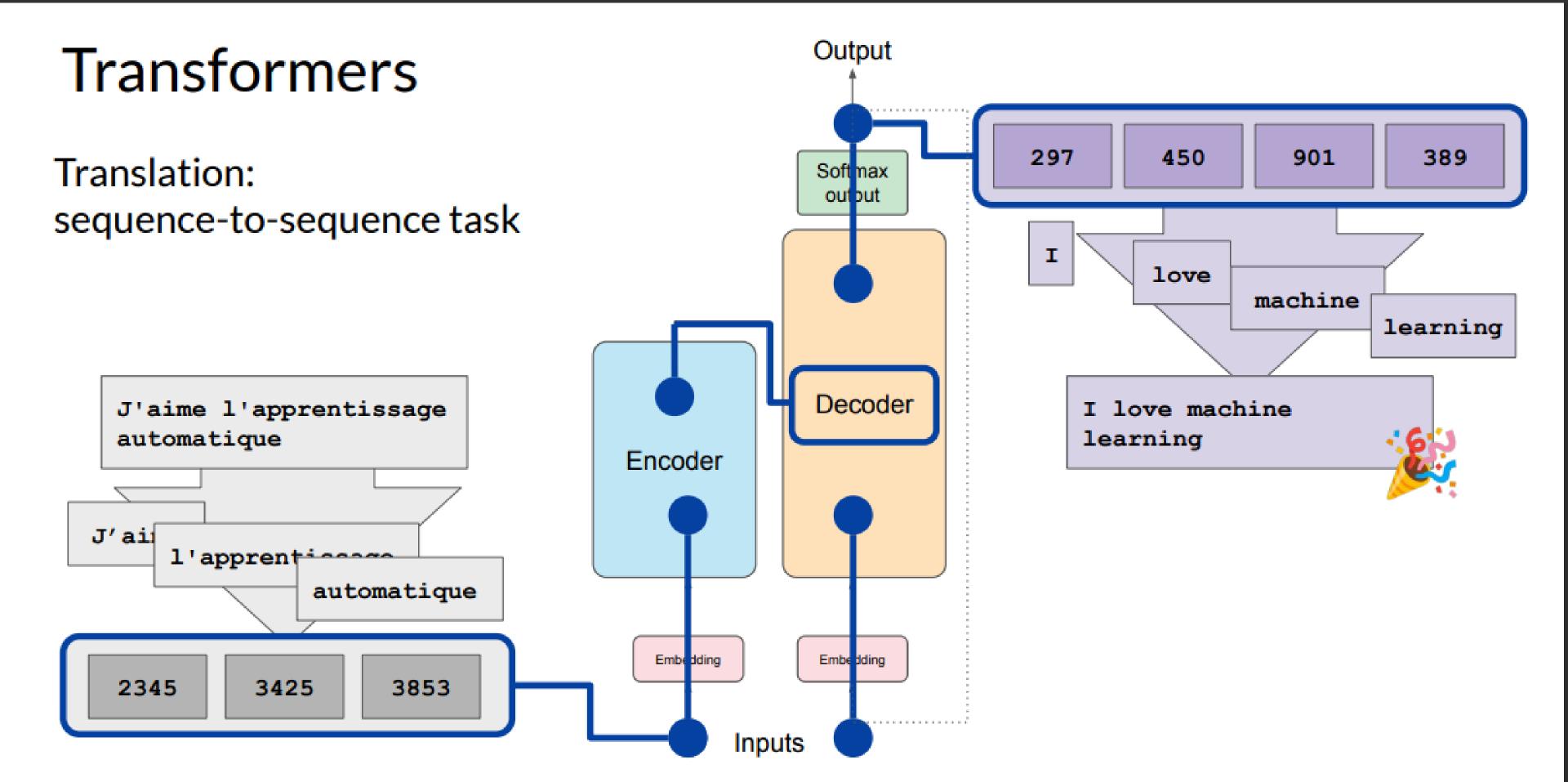












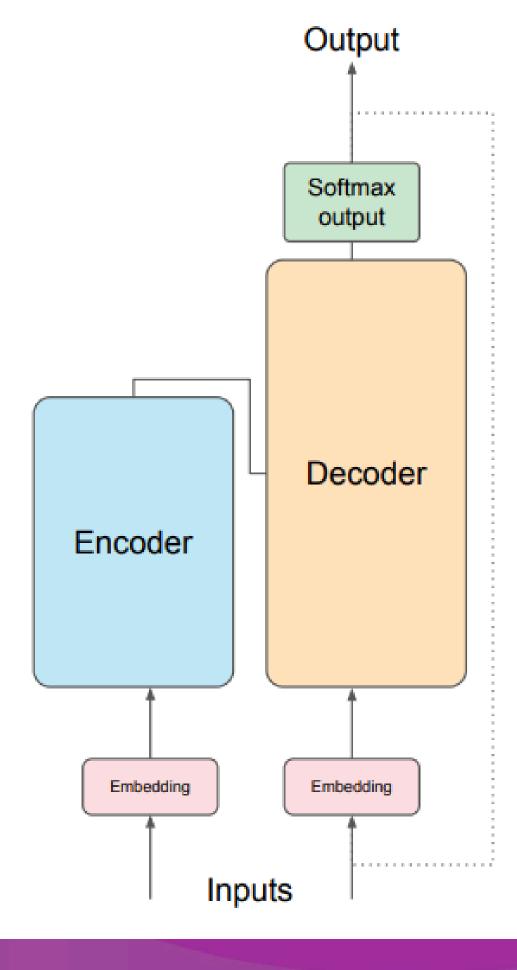




Transformers

Encoder

Encodes inputs ("prompts") with contextual understanding and produces one vector per input token.



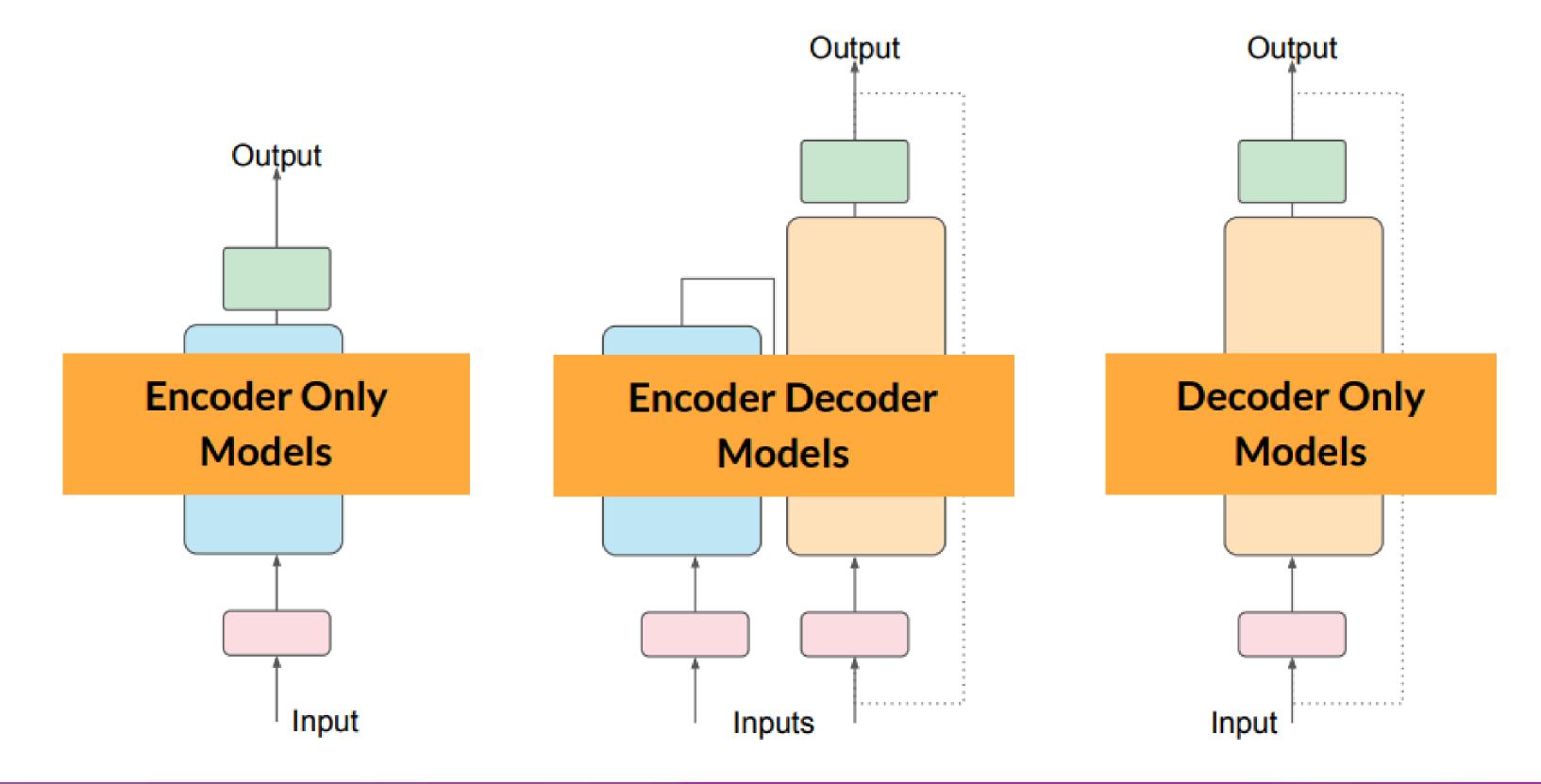
Decoder

Accepts input tokens and generates new tokens.





Transformers







Summary of in-context learning (ICL)

Prompt // Zero Shot

Classify this review:
I loved this movie!
Sentiment:

Context Window (few thousand words)

Prompt // One Shot

Classify this review:
I loved this movie!
Sentiment: Positive

Classify this review:
I don't like this chair.
Sentiment:

Prompt // Few Shot >5 or 6 examples

```
Classify this review:
I loved this movie!
Sentiment: Positive
Classify this review:
I don't like this
chair.
Sentiment: Negative
Classify this review:
Who would use this
product?
Sentiment:
```



The significance of scale: task ability

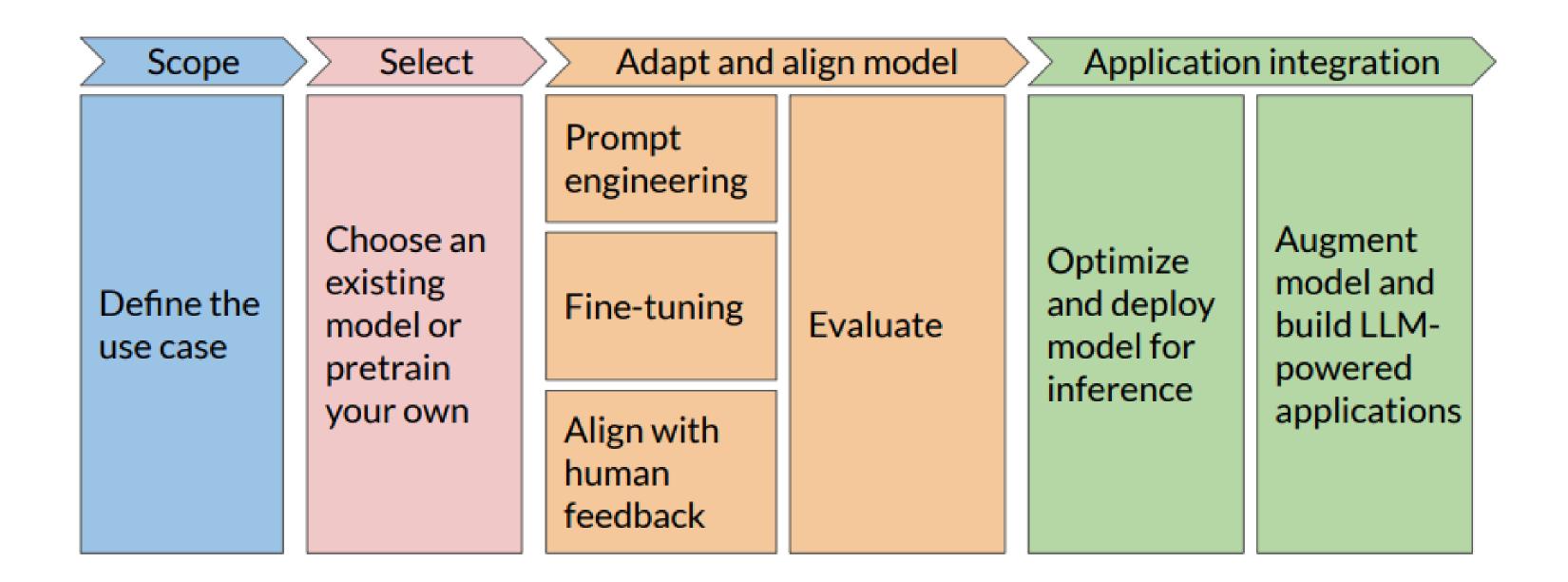


BLOOM ₋

*Bert-base



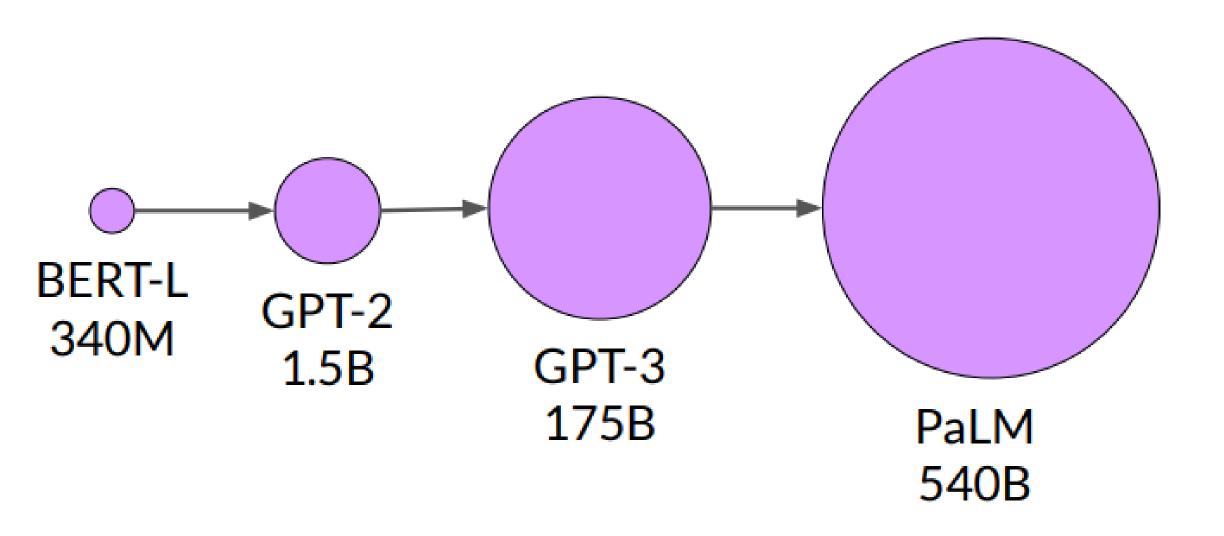
Generative Al project lifecycle







Model size vs. time



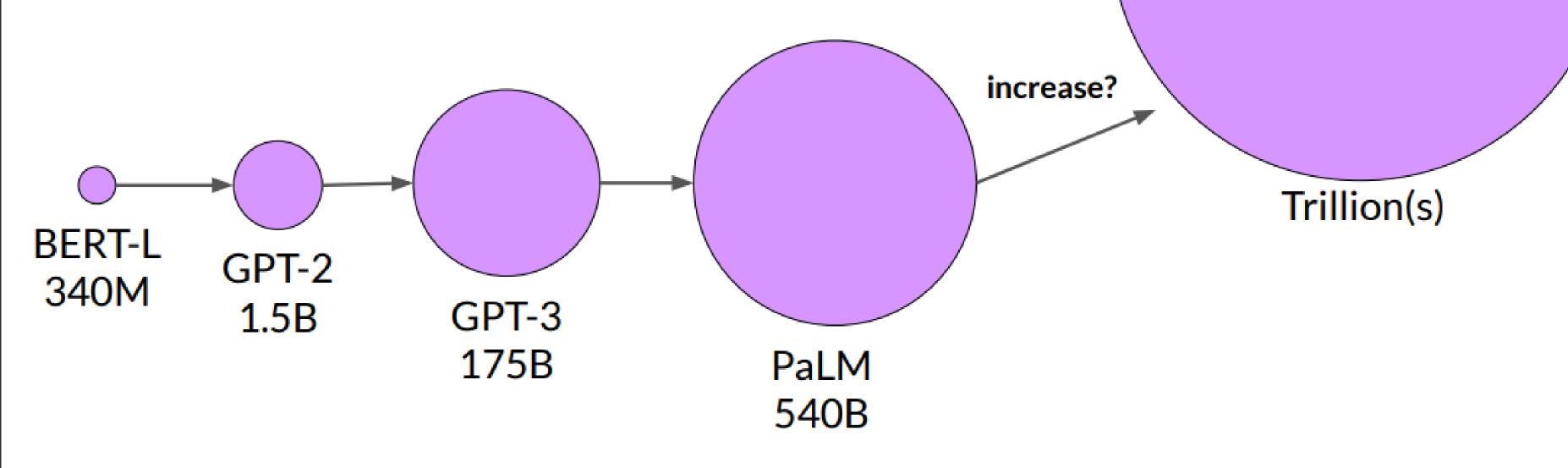
Growth powered by:

- Introduction of transformer
- Access to massive datasets
- More powerful compute resources

2018 2022 2023



Model size vs. time



2018 2022





Computational challenges

OutOfMemoryError: CUDA out of memory.







Approximate GPU RAM needed to store 1B parameters

1 parameter = 4 bytes (32-bit float)

1B parameters = 4×10^9 bytes = 4GB

4GB @ 32-bit full precision

Sources: https://huggingface.co/docs/transformers/v4.20.1/en/perf_train_gpu_one#anatomy-of-models-memory, https://github.com/facebookresearch/bitsandbytes

Additional GPU RAM needed to train 1B parameters

	Bytes per parameter		
Model Parameters (Weights)	4 bytes per parameter		

~20 extra bytes per parameter

Sources: https://huggingface.co/docs/transformers/v4.20.1/en/perf train gpu one#anatomy-of-models-memory, https://github.com/facebookresearch/bitsandbytes





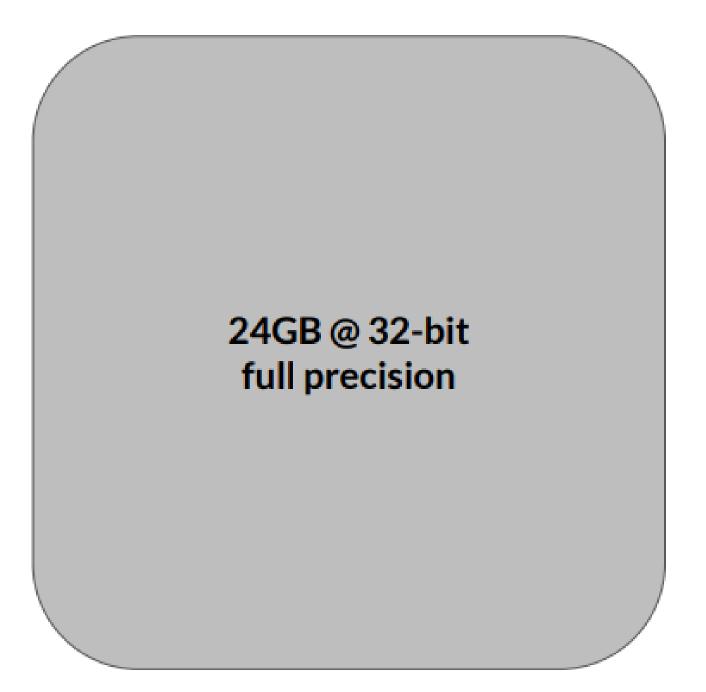
Approximate GPU RAM needed to train 1B-params

Memory needed to store model



4GB @ 32-bit full precision

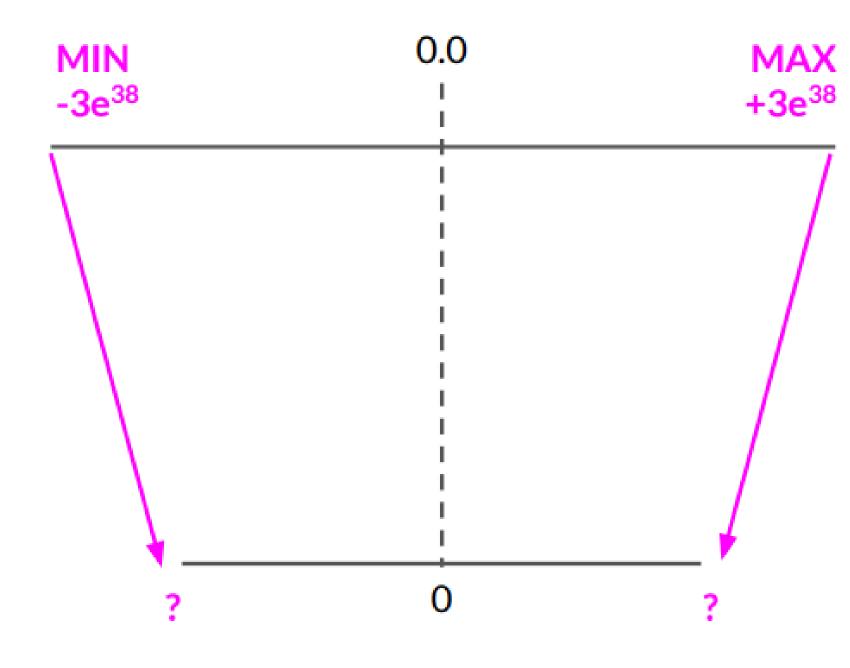
Memory needed to train model







Quantization



FP32

32-bit floating point

Range:

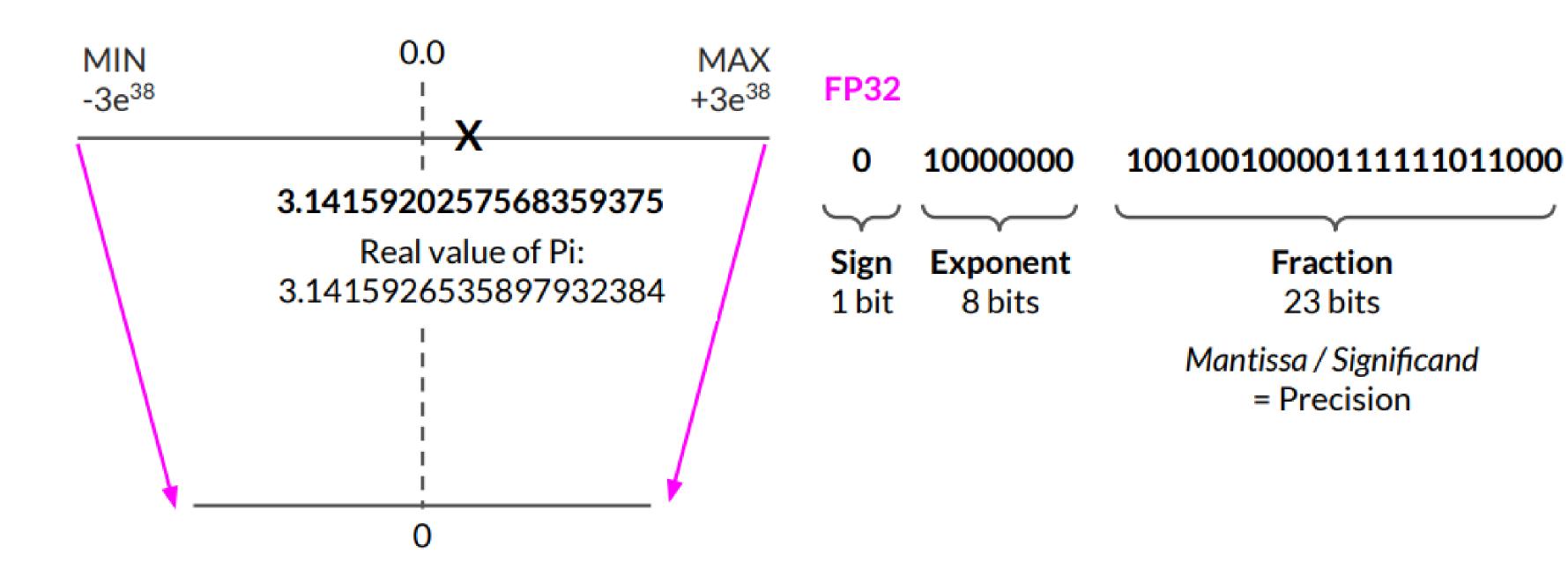
From $-3e^{38}$ to $+3e^{38}$

FP16 | BFLOAT16 | INT8

16-bit floating point | 8-bit integer

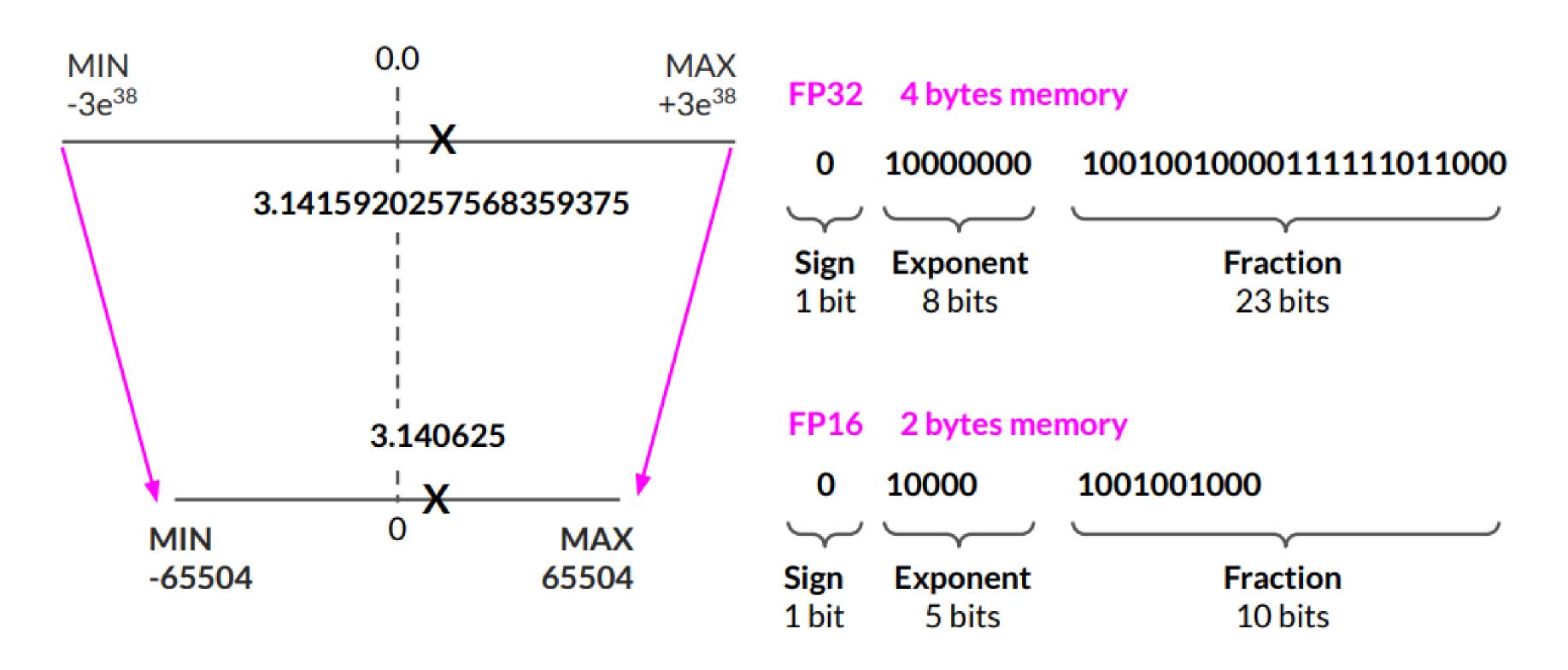


Quantization: FP32



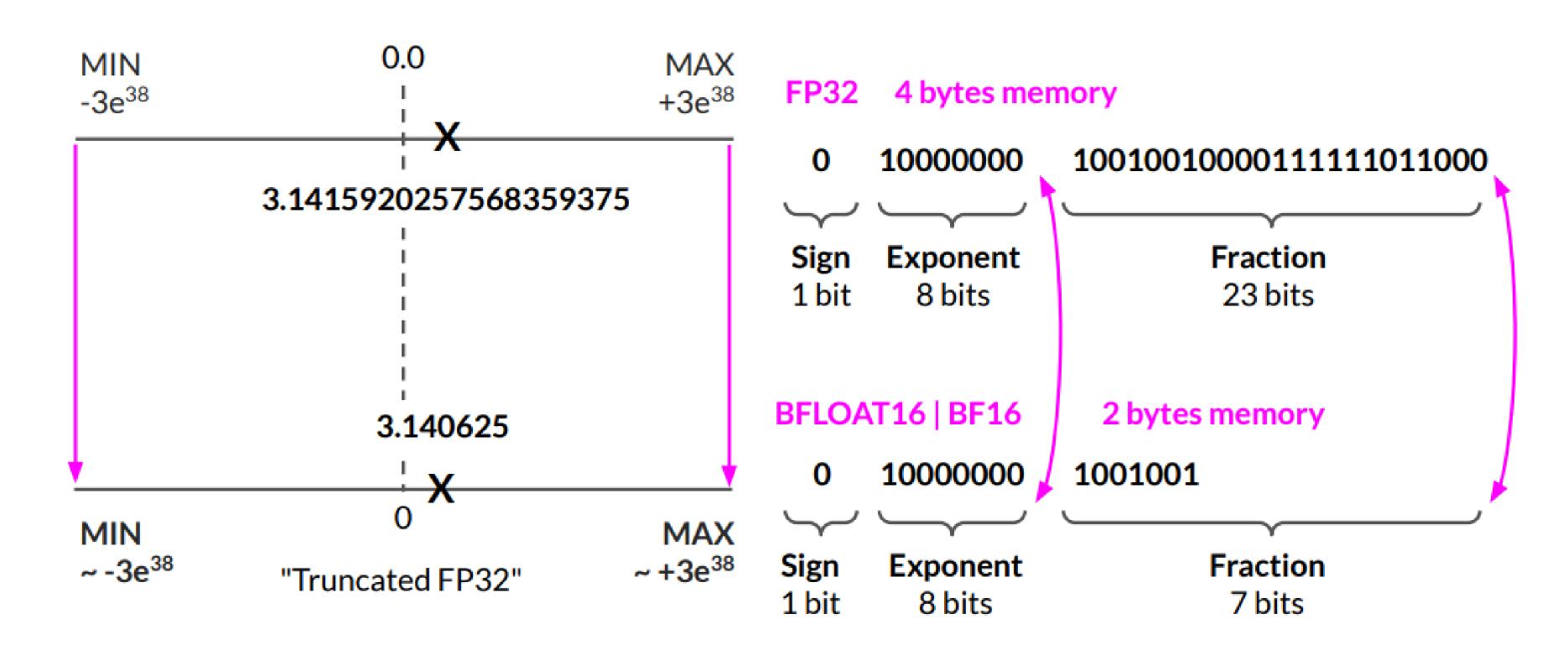


Quantization: FP16



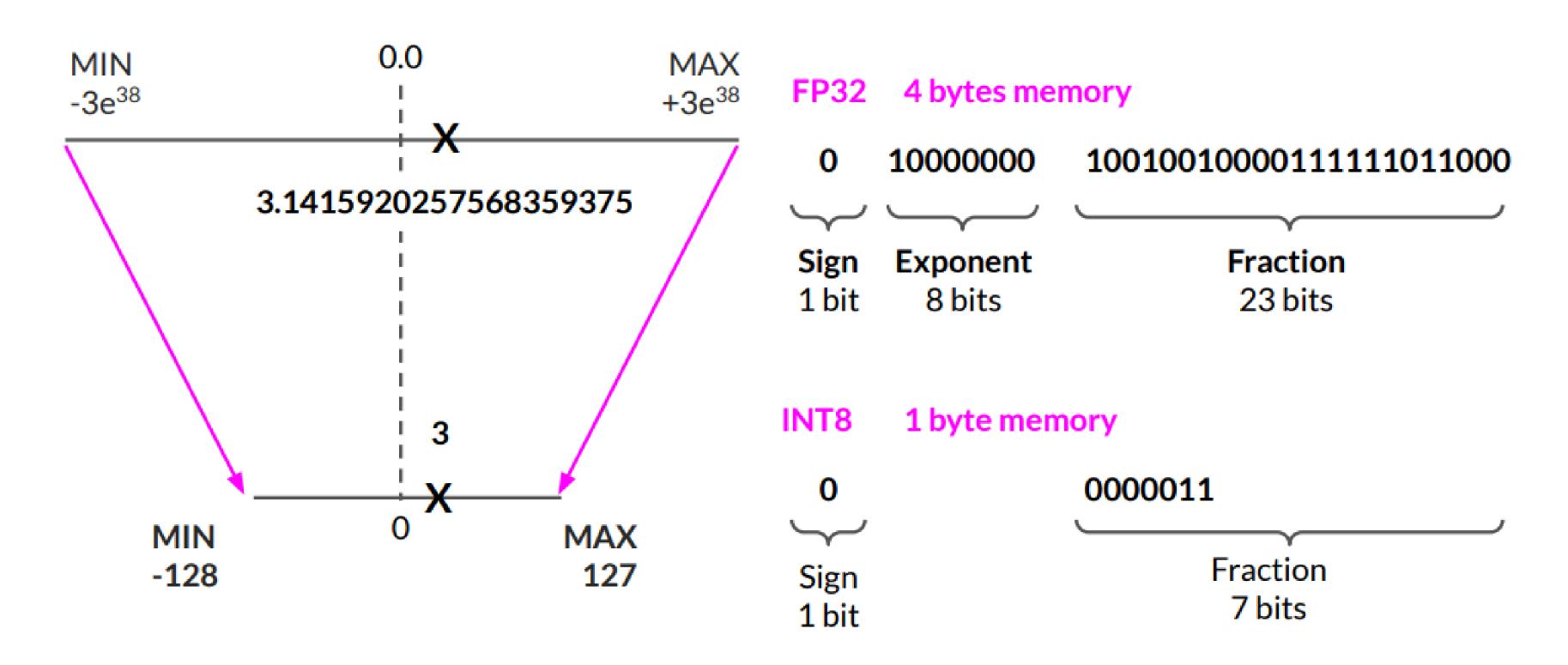


Quantization: BFLOAT16





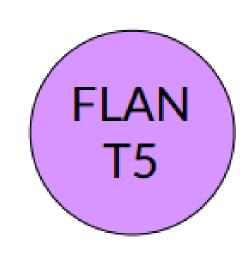
Quantization: INT8





Quantization: Summary

	Bits	Exponent	Fraction	Memory needed to store one value
FP32	32	8	23	4 bytes
FP16	16	5	10	2 bytes
BFLOAT16	16	8	7	2 bytes
INT8	8	-/-	7	1 byte



- Reduce required memory to store and train models
- Projects original 32-bit floating point numbers into lower precision spaces
- Quantization-aware training (QAT) learns the quantization scaling factors during training
- BFLOAT16 is a popular choice





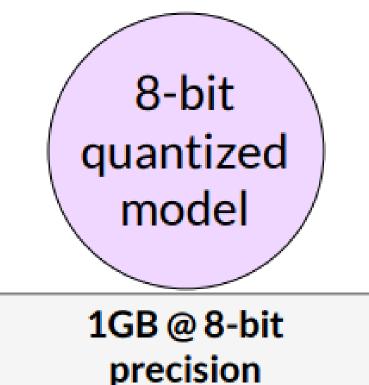
Approximate GPU RAM needed to store 1B parameters

Fullprecision model

4GB @ 32-bit full precision

16-bit quantized model

2GB @ 16-bit half precision



Sources: https://huggingface.co/docs/transformers/v4.20.1/en/perf train gpu one#anatomy-of-models-memory, https://github.com/facebookresearch/bitsandbytes





GPU RAM needed to train larger models

1B param model

175B param model

4,200 GB @ 32-bit full precision

500B param model

12,000 GB @ 32-bit full precision





