

Transformer Architecture

Transformers are not only self attention, you also have FFN layers

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

usually simple, one hidden layer
but where most params are

input is added to output (represented by the arrows into "Add & Norm")

vectors: d_{model}

Queries/keys: d_k , always smaller than d_{model} (saves computation)

Values: separate dim d_v , output is multiplied by W^o which is $d_v \times d_{\text{model}}$ to get back to d_{model}

FFN can explode the dimension with W_1 and collapse it back with W_2

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

layers

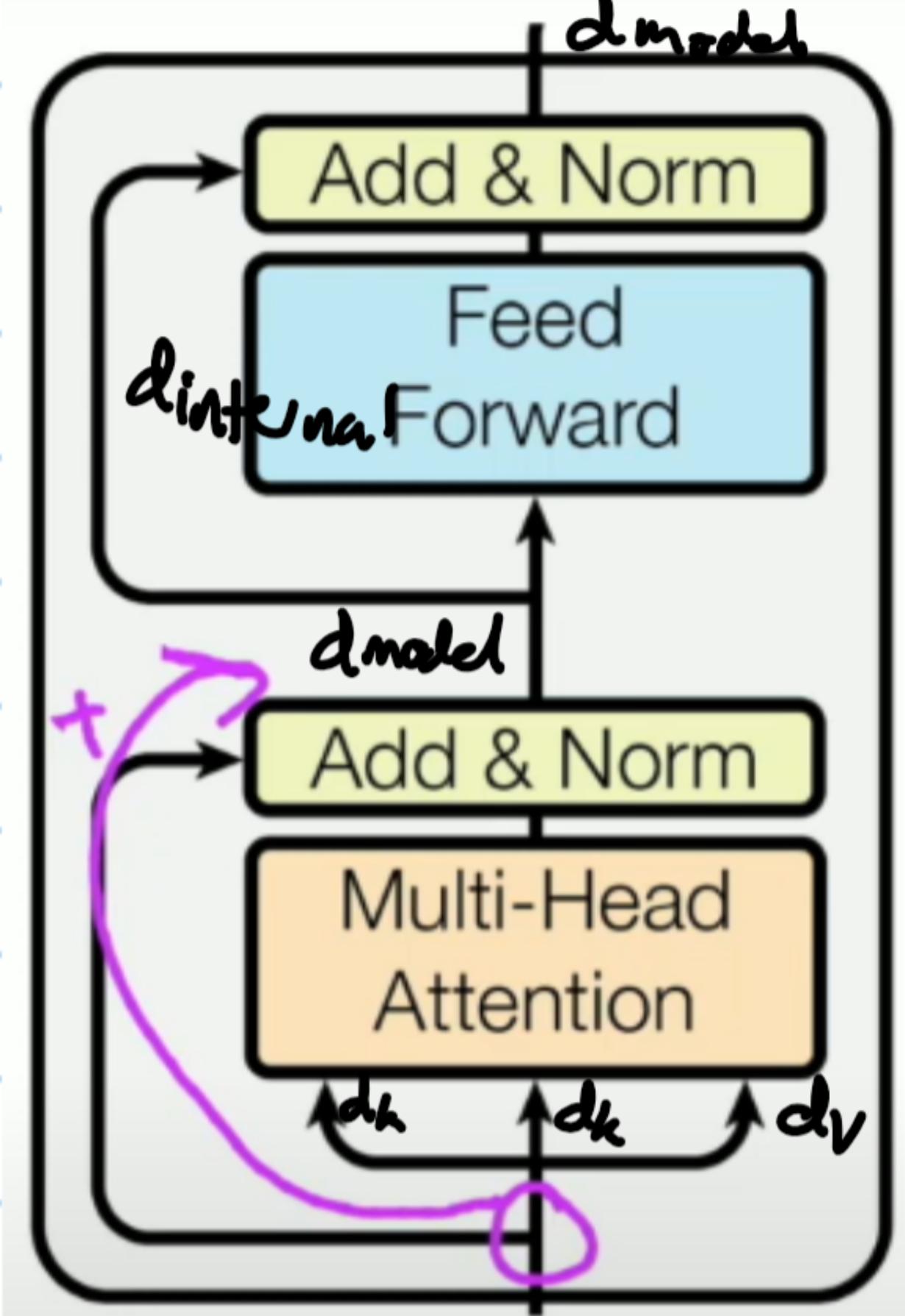
	N	d_{model}	d_{ff}	h	d_k	d_v
base	6	512	2048	8	64	64

gets expensive quickly

From Vaswani et al.

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}
GPT-3 Small	125M	12	768	12	64
GPT-3 Medium	350M	24	1024	16	64
GPT-3 Large	760M	24	1536	16	96
GPT-3 XL	1.3B	24	2048	24	128
GPT-3 2.7B	2.7B	32	2560	32	80
GPT-3 6.7B	6.7B	32	4096	32	128
GPT-3 13B	13.0B	40	5140	40	128
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128

From GPT-3; d_{head} is our d_k

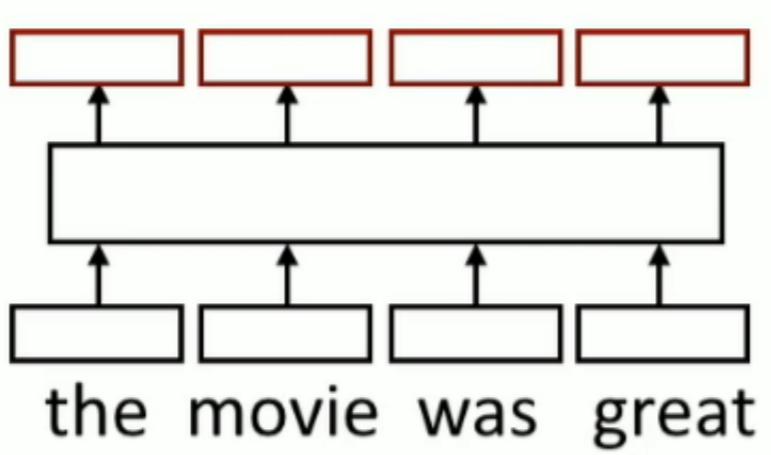


80% of the computation is in the FFN layers
other is attention

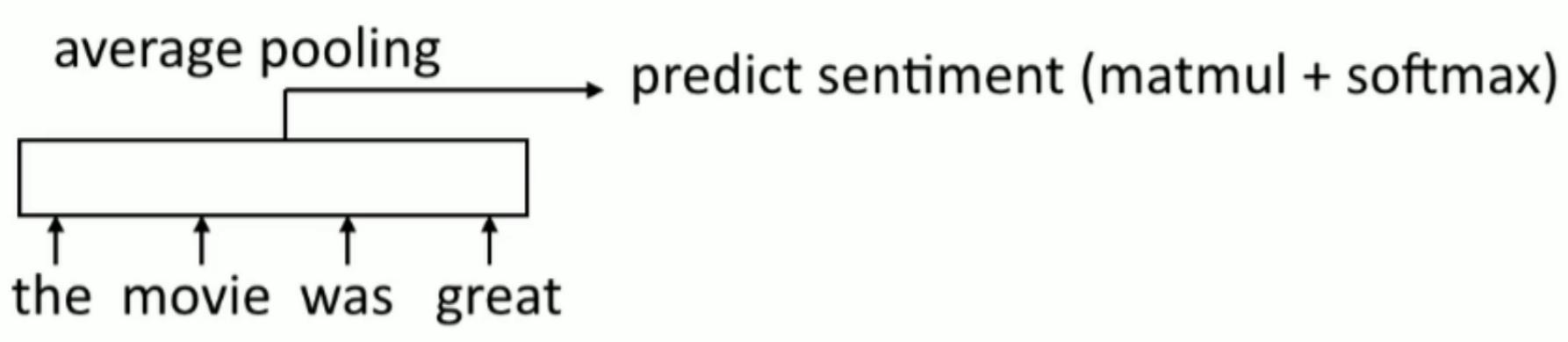
most of the information in models like GPT is from
the FFN params

Using Transformers

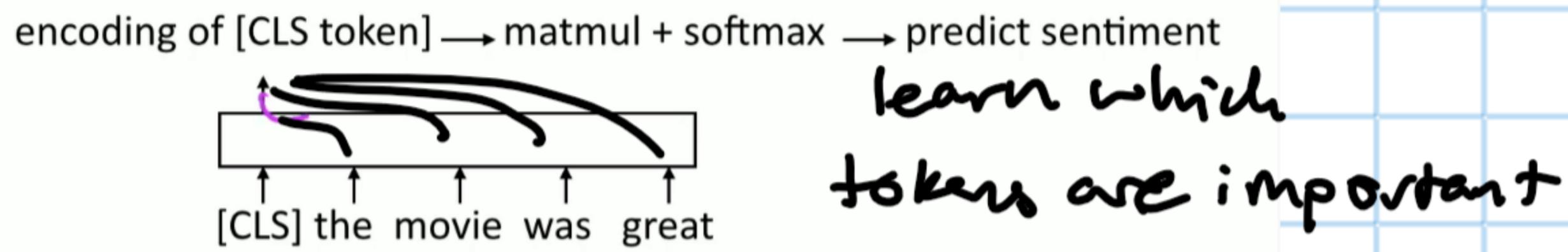
What do Transformers produce?



- ▶ **Encoding of each word** — can pass this to another layer to make a prediction (like predicting the next word for language modeling)
- ▶ Like RNNs, Transformers can be viewed as a transformation of a sequence of vectors into a sequence of context-dependent vectors



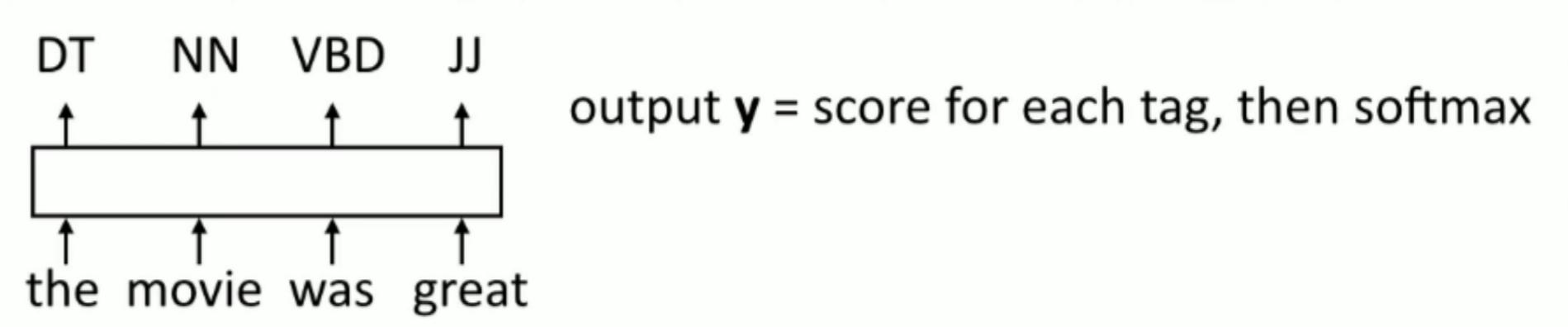
- ▶ Alternative: use a placeholder [CLS] token at the start of the sequence.



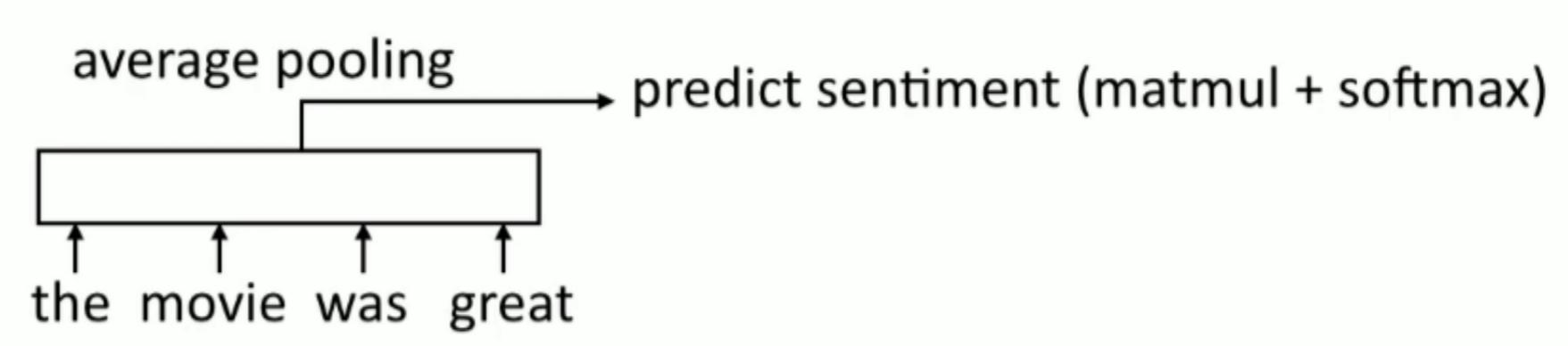
- ▶ Because [CLS] attends to everything with self-attention, it can do the pooling for you!

Transformer Uses

- ▶ Transducer: make some prediction for each element in a sequence



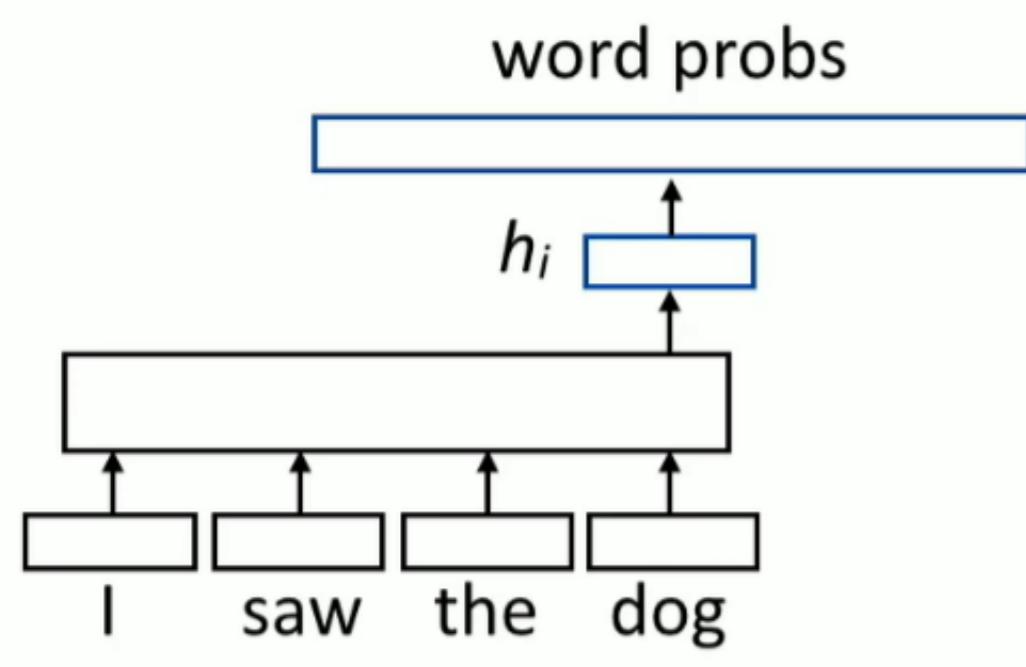
- ▶ Classifier: encode a sequence into a fixed-sized vector and classify that



Transformer Language Modeling

Goal: predict next word

Transformer Language Modeling



$$P(w|\text{context}) = \frac{\exp(\mathbf{w} \cdot \mathbf{h}_i)}{\sum_{w'} \exp(\mathbf{w}' \cdot \mathbf{h}_i)}$$

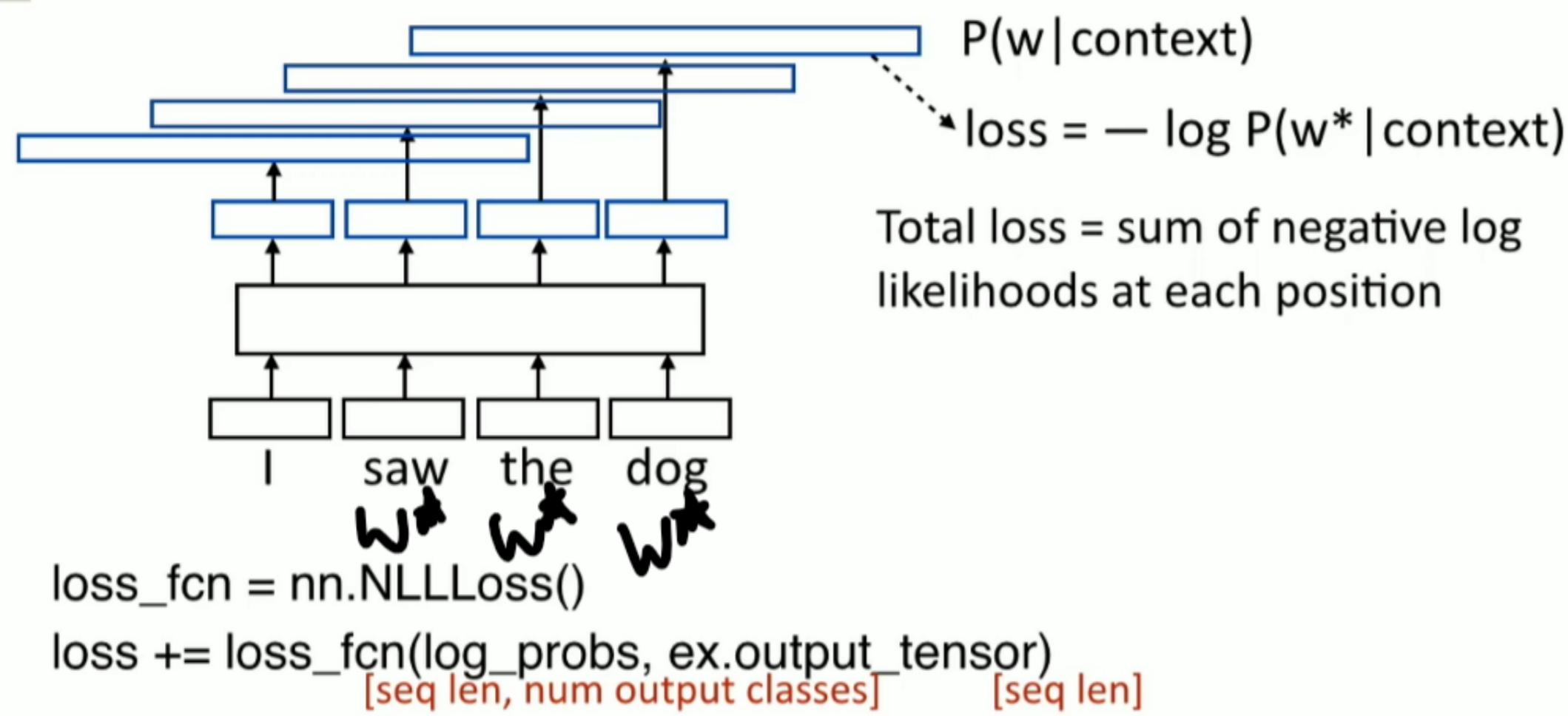
equivalent to

$$P(w|\text{context}) = \text{softmax}(W\mathbf{h}_i)$$

- h_i is the embedding of dog produced by the Transformer

- W is a (vocab size) x (hidden size) matrix; linear layer in PyTorch (rows are word embeddings)

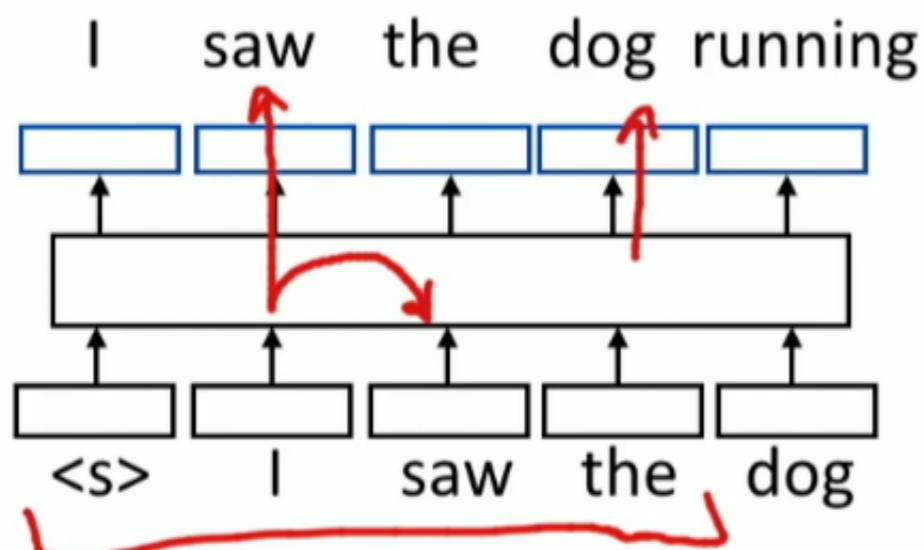
Training Transformer LMs



- Batching is a little tricky with NLLLoss: need to collapse [batch, seq len, num classes] to [batch * seq len, num classes]. You do not need to batch

A Small Problem with Transformer LMs

- This Transformer LM as we've described it will *easily* achieve perfect accuracy. Why?



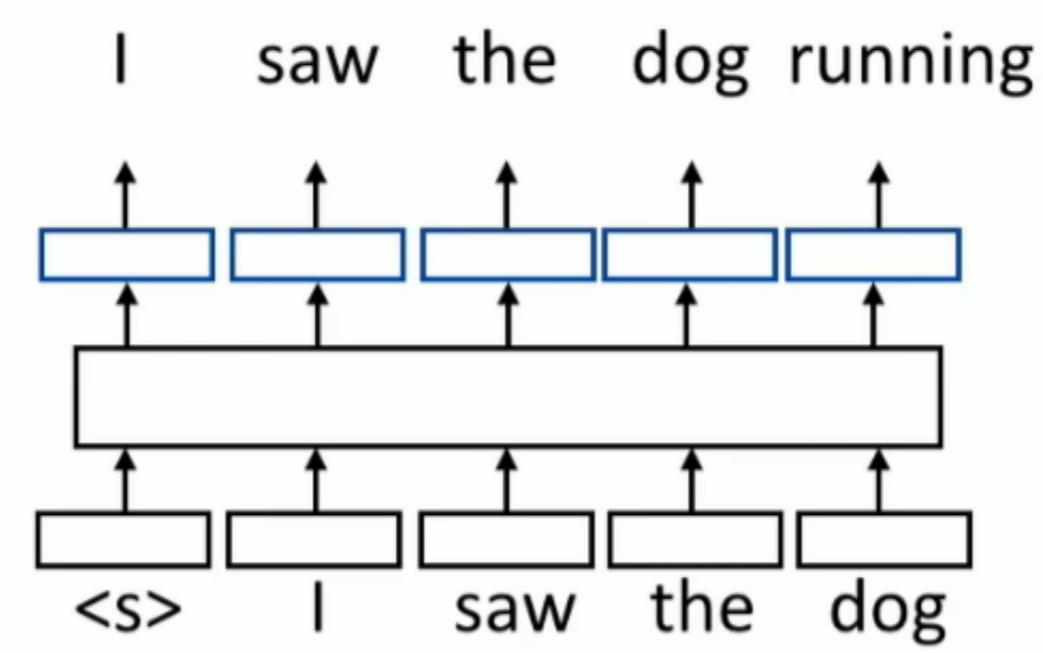
- With standard self-attention: "I" attends to "saw" and the model is "cheating". How do we ensure that this doesn't happen?

Implementation in PyTorch

- nn.TransformerEncoder can be built out of nn.TransformerEncoderLayers, can accept an input and a mask for language modeling:

```
# Inside the module; need to fill in size parameters
layers = nn.TransformerEncoderLayer([...])
transformer_encoder = nn.TransformerEncoder(encoder_layers, num_layers=[...])
[...]
# Inside forward(): puts negative infinities in the red part
mask = torch.triu(torch.ones(len, len) * float('-inf'), diagonal=1)
output = transformer_encoder(input, mask=mask)
```

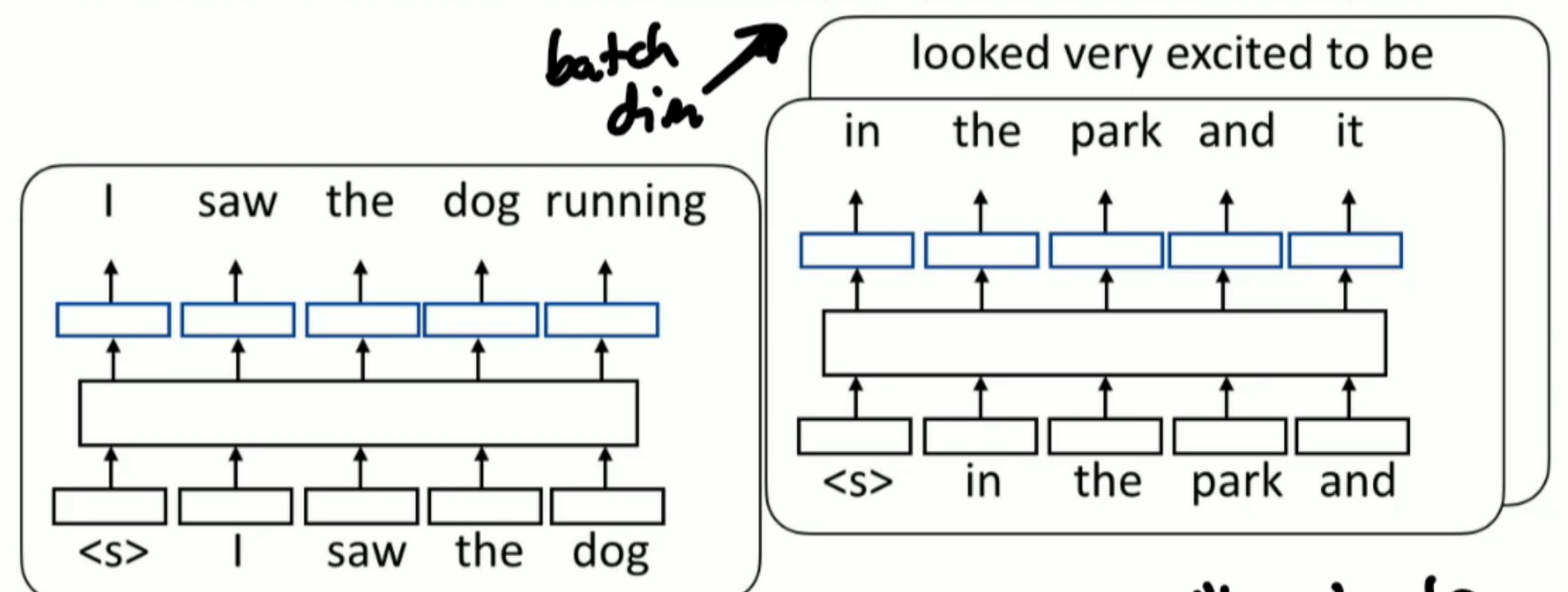
Training Transformer LMs



- Input is a sequence of words, output is those words shifted by one
- Allows us to train on predictions across several timesteps simultaneously (similar to batching but this is NOT what we refer to as batching)

Batched LM Training

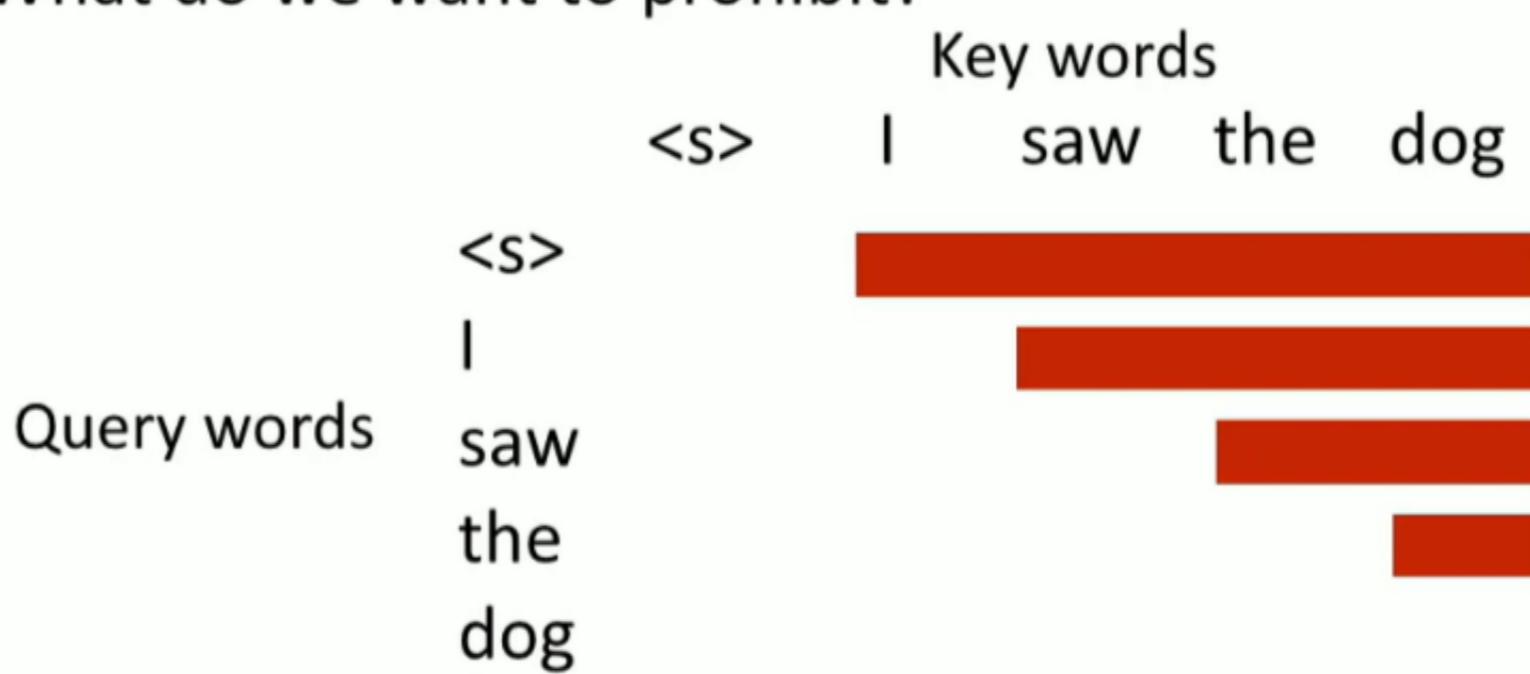
I saw the dog running in the park and it looked very excited to be there



if you have pos enc., the model will not learn to attend to stuff far away.

Attention Masking

- What do we want to prohibit?



- This is called a **causal** mask (also, causal self-attention / causal Transformers). Only things in the "past" can influence the "present"

Transformers Extension

Scaling Laws

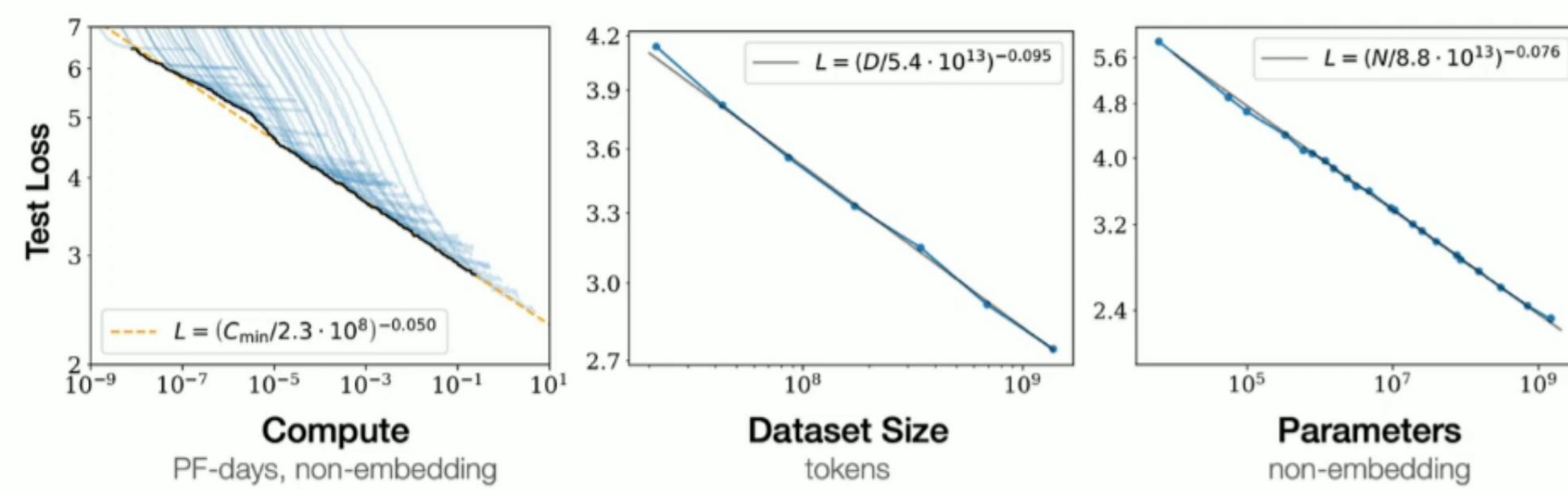


Figure 1 Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute² used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

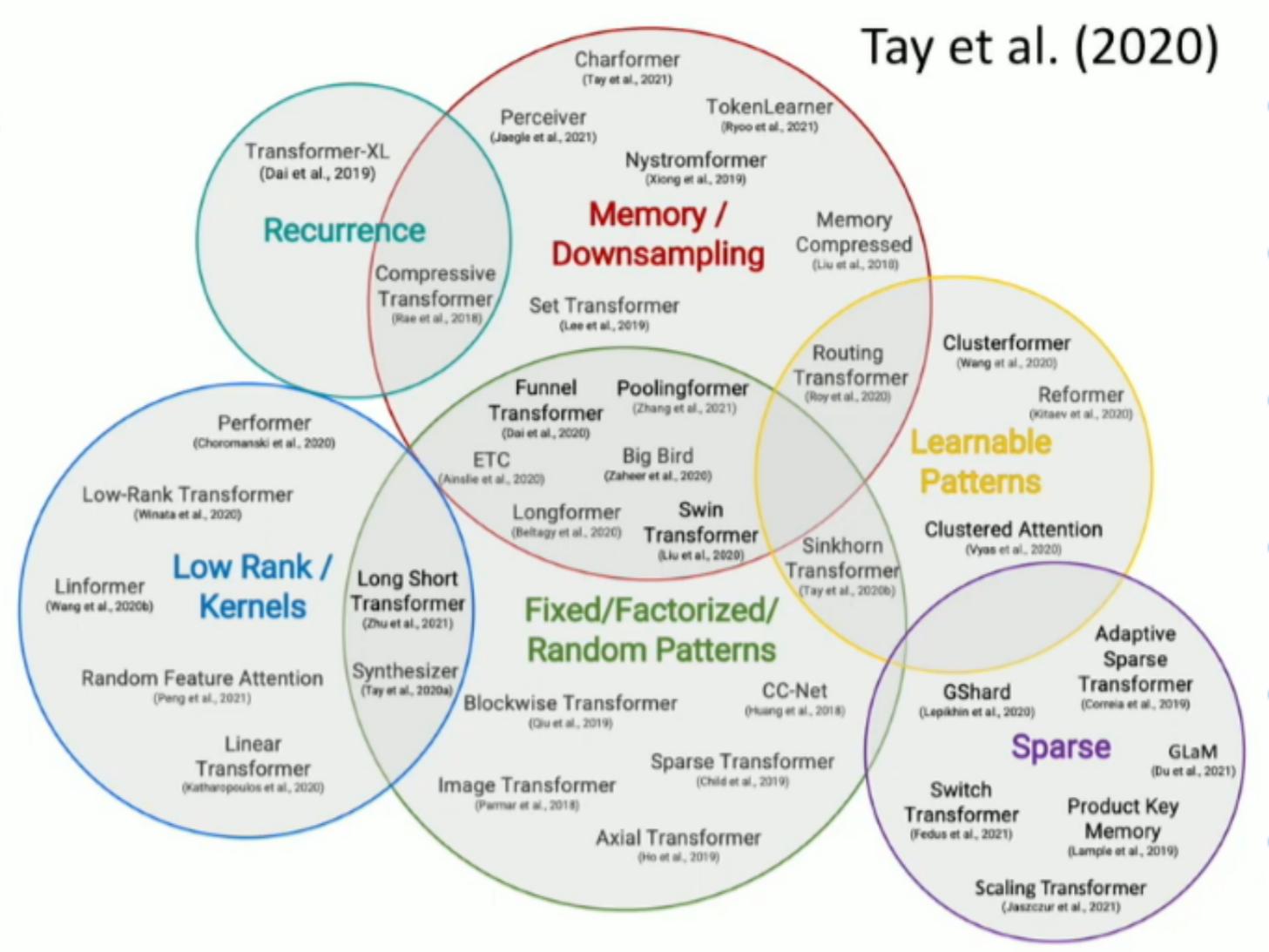
- ▶ Transformers scale really well!

just increasing training probably will decrease loss, which is rare for such models

Kaplan et al. (2020)

Transformer Runtime

- ▶ Even though most parameters and FLOPs are in feedforward layers, Transformers are still limited by quadratic complexity of self-attention
- ▶ Many ways proposed to handle this



Performers

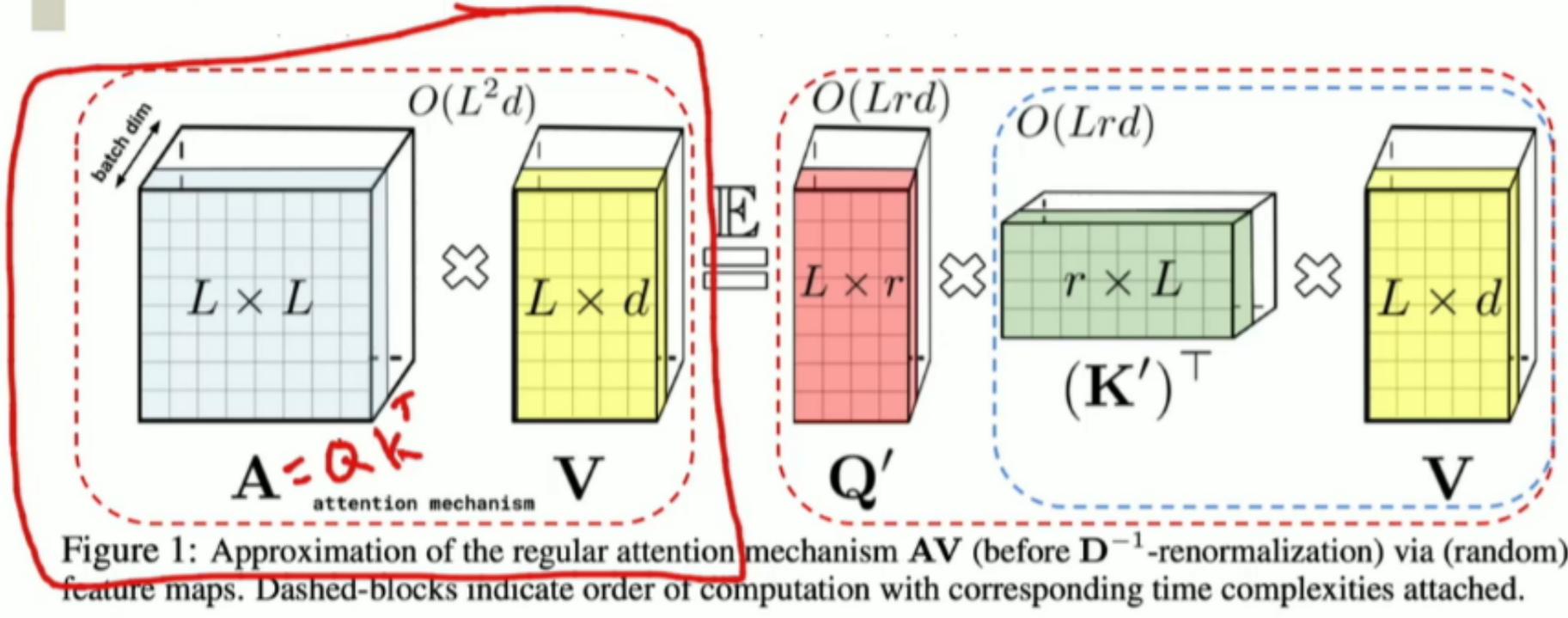


Figure 1: Approximation of the regular attention mechanism AV (before D^{-1} -renormalization) via (random) feature maps. Dashed-blocks indicate order of computation with corresponding time complexities attached.

- ▶ No more len^2 term, but we are fundamentally approximating the self-attention mechanism (cannot form A and take the softmax)

Choromanski et al. (2020)

Longformer

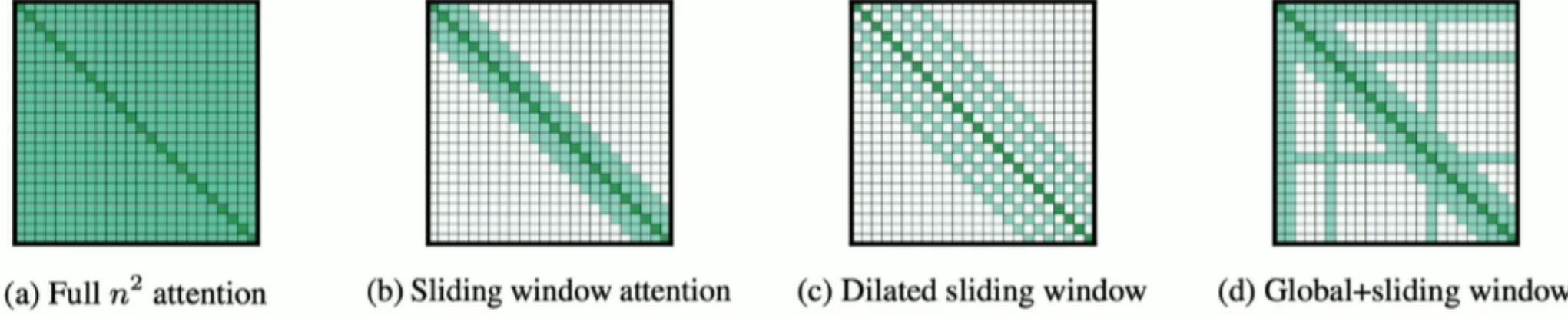
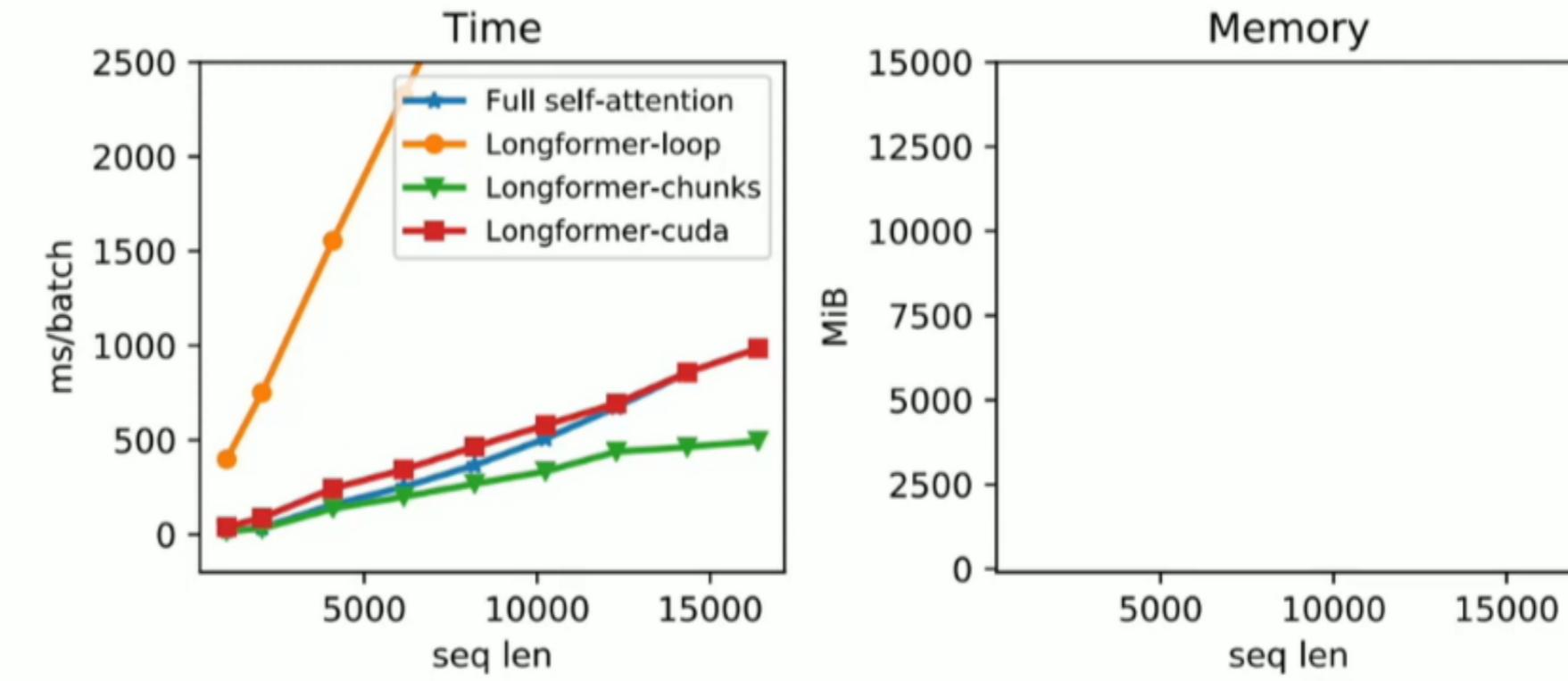


Figure 2: Comparing the full self-attention pattern and the configuration of attention patterns in our Longformer.

- ▶ Use several pre-specified self-attention patterns that limit the number of operations while still allowing for attention over a reasonable set of things
- ▶ Scales to 4096-length sequences

Beltagy et al. (2021)

Attention Maps



- ▶ Loop = non-vectorized version
- ▶ What will the memory profile look like?

Beltagy et al. (2021)

