# Attention and sequential structure
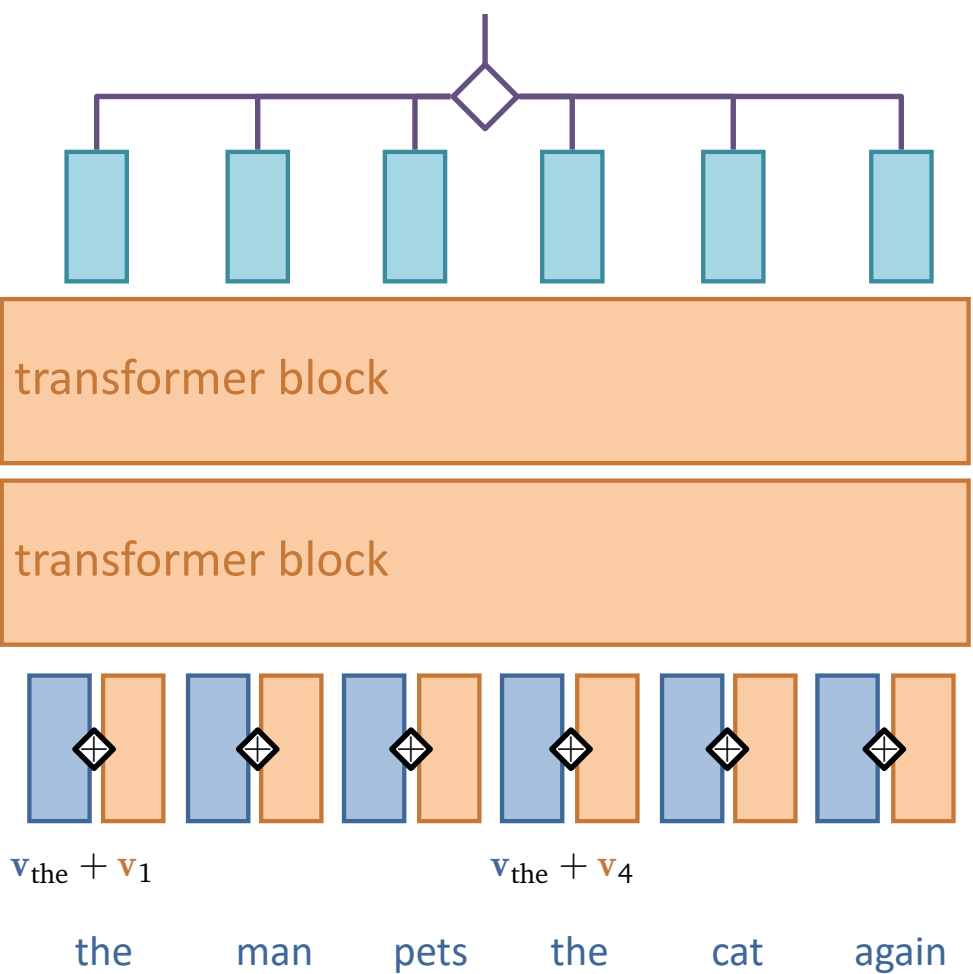
**Position Embeddings**

$\mathbf{v}_{the}, \mathbf{v}_{man}, \mathbf{v}_{pets}, \mathbf{v}_{cat}, \mathbf{v}_{again}$

Conceptually simple

Easy to implement

Adds set of learnable parameters $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4, \mathbf{v}_5, \ldots$

Maximum context length at test time limited to max sequence length in training set

Embedding quality diminishes (in theory) with $t$

transformer block
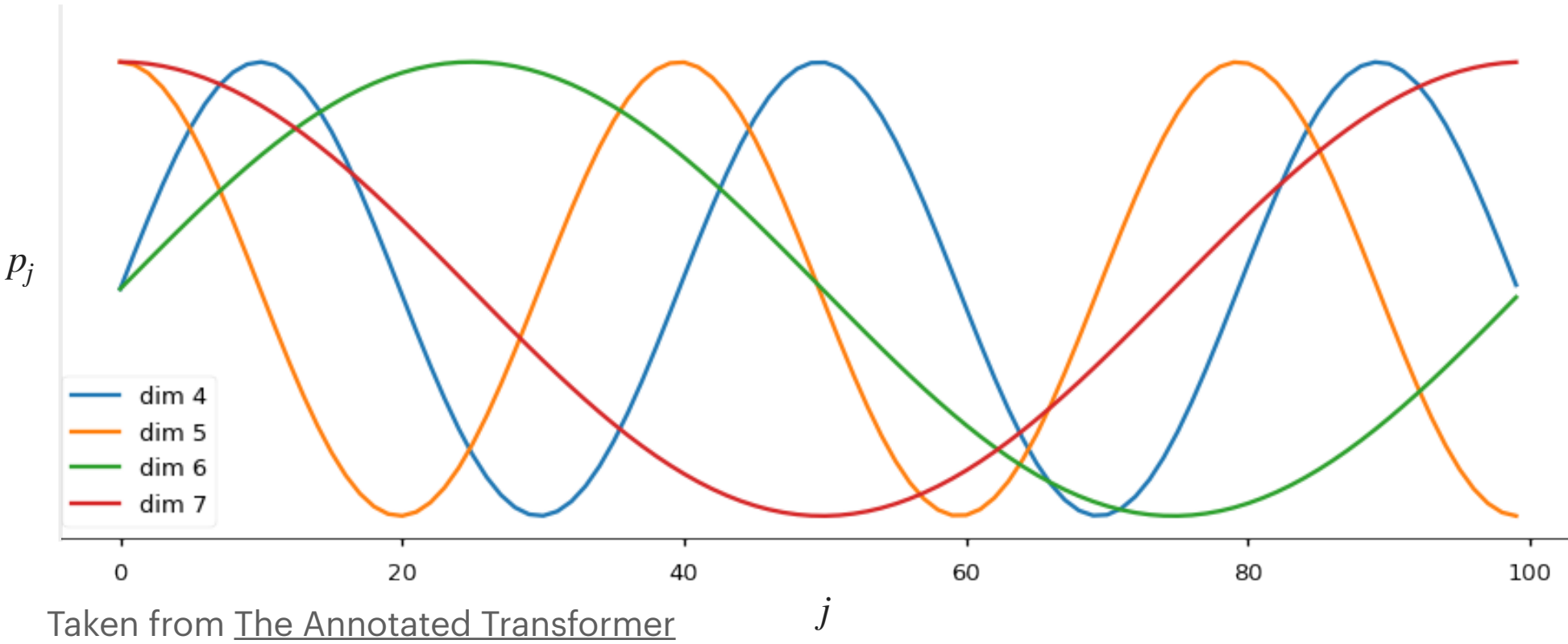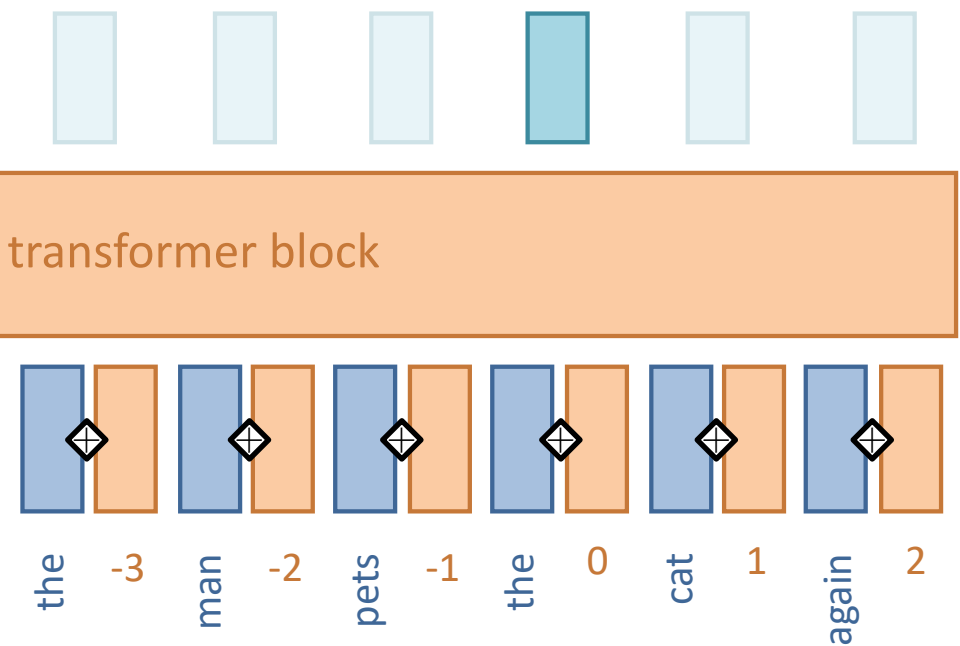
transformer block

$\mathbf{v}_{the} + \mathbf{v}_1$   $\mathbf{v}_{the} + \mathbf{v}_4$

the   man   pets   the   cat   again

**Relative position embeddings**

Captures the property of absolute position invariance which is desirable

Like absolute position embeddings, limited in context length at inference

Conceptually and practically more complex; on the surface this results in (T-1)^2 different inputs because each input will have T-1 separate representations (there is a hack that turns this into 2T-1)

transformer block

the   -3   man   -2   pets   -1   the   0   cat   1   again   2

**Position encodings**

$\mathbf{v}_{the}, \mathbf{v}_{man}, \mathbf{v}_{pets}, \mathbf{v}_{cat}, \mathbf{v}_{again}$

Captures the property of absolute position invariance which is desirable

Conceptually more complex; on the surface this results in (T-1)^2 different inputs because each input will have T-1 separate representations (there is a hack that turns this into 2T-1)

transformer block

transformer block

$\mathbf{v}_{the}, \mathbf{v}_{man}, \mathbf{v}_{pets}, \mathbf{v}_{cat}, \mathbf{v}_{again}$

$\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4$

The idea behind position embeddings is simple. Just like we assign each word in our vocabulary an embedding vector, we also assign each *position* in our vocabulary an embedding vector. This way, the input vectors for the first "the" in the input sequence and the second "the" are different, because the first is added to the position embedding $\mathbf{v}_1$ and the second is added to the input embedding $\mathbf{v}_2$.

This break our equivariance: the position information becomes *part* of our embedding vectors, and is fed into the self attention. This is very effective, and very easy to implement. The only drawback is that we can't run the model very well on sequences that are longer than the largest

The idea behind relative position encodings is that it doesn't really matter so much where the word is

$p_j$

dim 4
dim 5
dim 6
dim 7

$j$

Taken from The Annotated Transformer

definition of self attention.

# BERT (2018)

**BERT**

The most popular/used/studied Transformer model to date

30K citations [1]

Trained on Wikipedia and Book Corpus (~10K books)

Wordpiece embedding [2]

Non-causal, uses masked LM procedure for pretraining

340 M parameters in total

LM prediction head: 4 ReLU (D=4096)

LM training took 4 days on 64 TPU cores
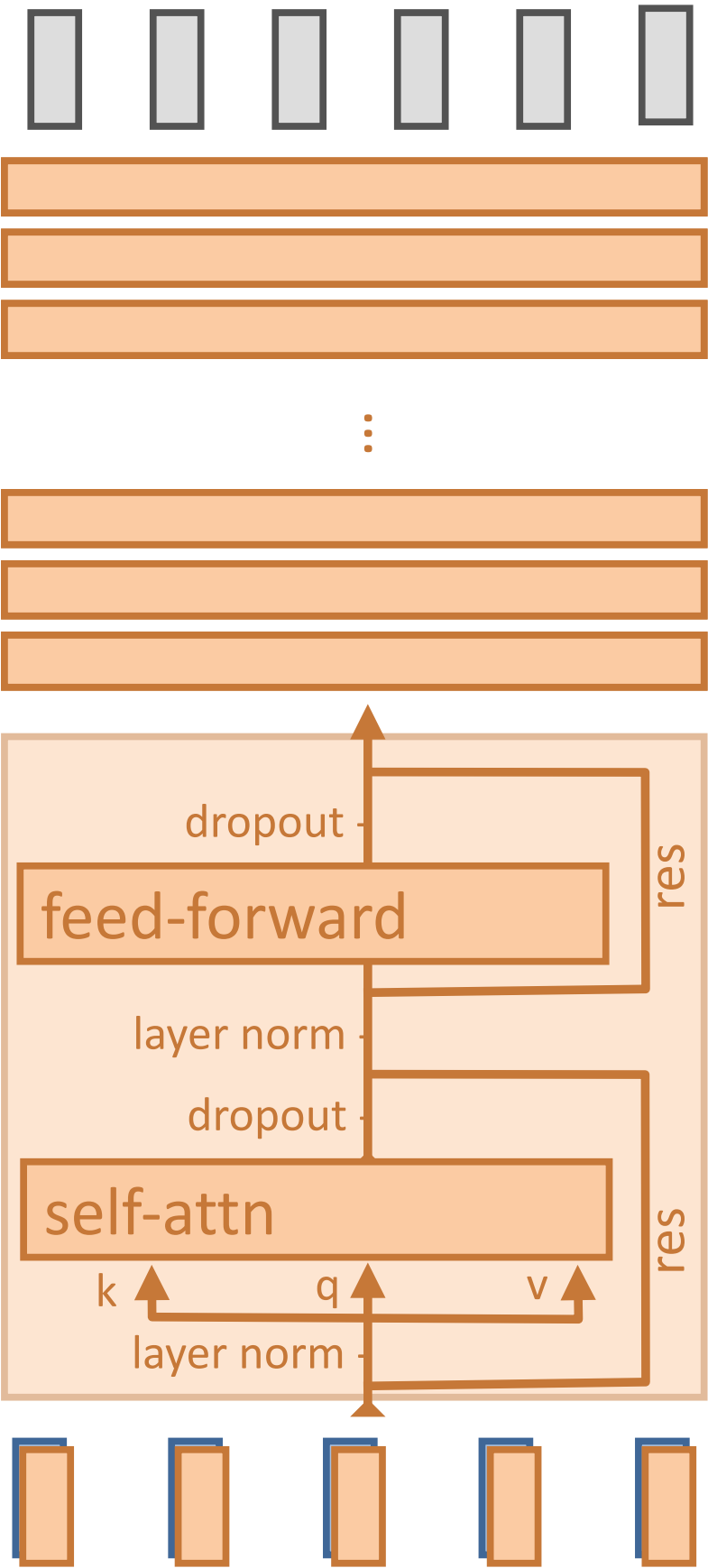
$D = 1024$

$N_H = 16$

$N_B = 24$

$T_{max} = 1024$



Figure taken from Bloem, 2020



BERT Input Features [1]

| Hyperparams | | | | Dev Set Accuracy | | |
|---|---|---|---|---|---|---|
| #L | #H | #A | LM (ppl) | MNLI-m | MRPC | SST-2 |
| 3 | 768 | 12 | 5.84 | 77.9 | 79.8 | 88.4 |
| 6 | 768 | 3 | 5.24 | 80.6 | 82.2 | 90.7 |
| 6 | 768 | 12 | 4.68 | 81.9 | 84.8 | 91.3 |
| 12 | 768 | 12 | 3.99 | 84.4 | 86.7 | 92.9 |
| 12 | 1024 | 16 | 3.54 | 85.7 | 86.9 | 93.3 |
| 24 | 1024 | 16 | 3.23 | 86.6 | 87.8 | 93.7 |

Ablation Study [1]

[1] Devlin et al., *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (2018)*, [2] Wu et al., *Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation (2016)*