

Transformer + Course Projects

Wei Xu

(many slides from Greg Durrett)

Administrivia

▶ Readings —

- ▶ “The Annotated Transformer” by Sasha Rush

<https://nlp.seas.harvard.edu/2018/04/03/attention.html>

- ▶ “The Illustrated Transformer” by Jay Lamar

<http://jalammar.github.io/illustrated-transformer/>

Transformers

Attention is All You Need

Attention Is All You Need

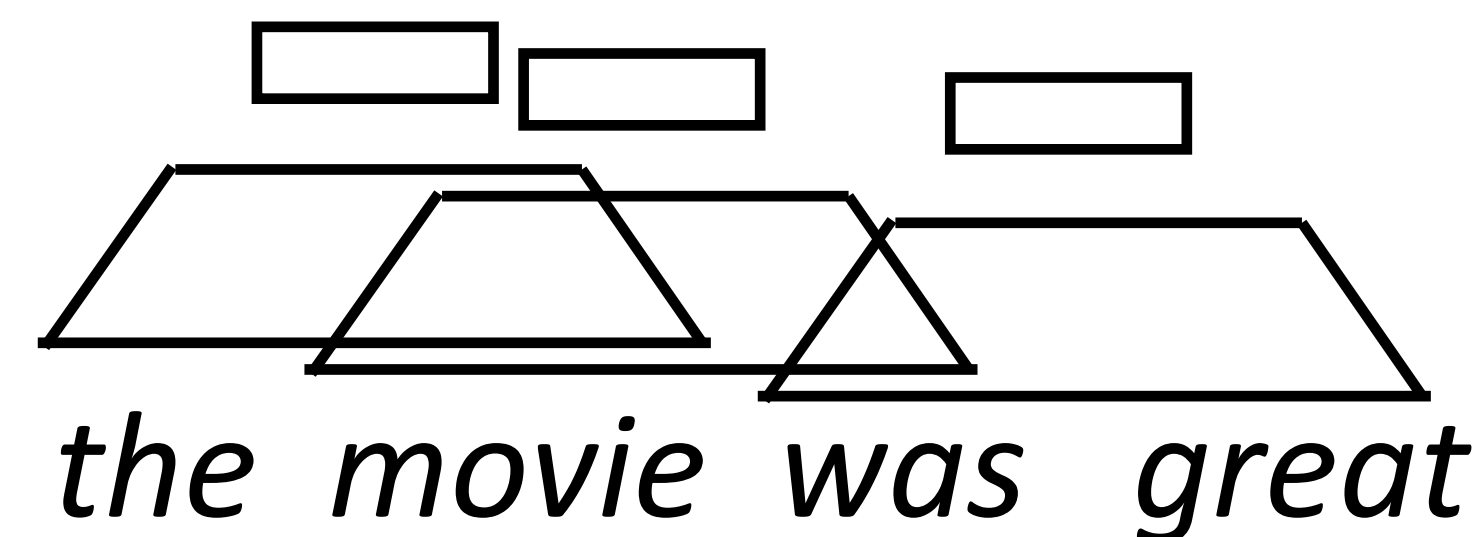
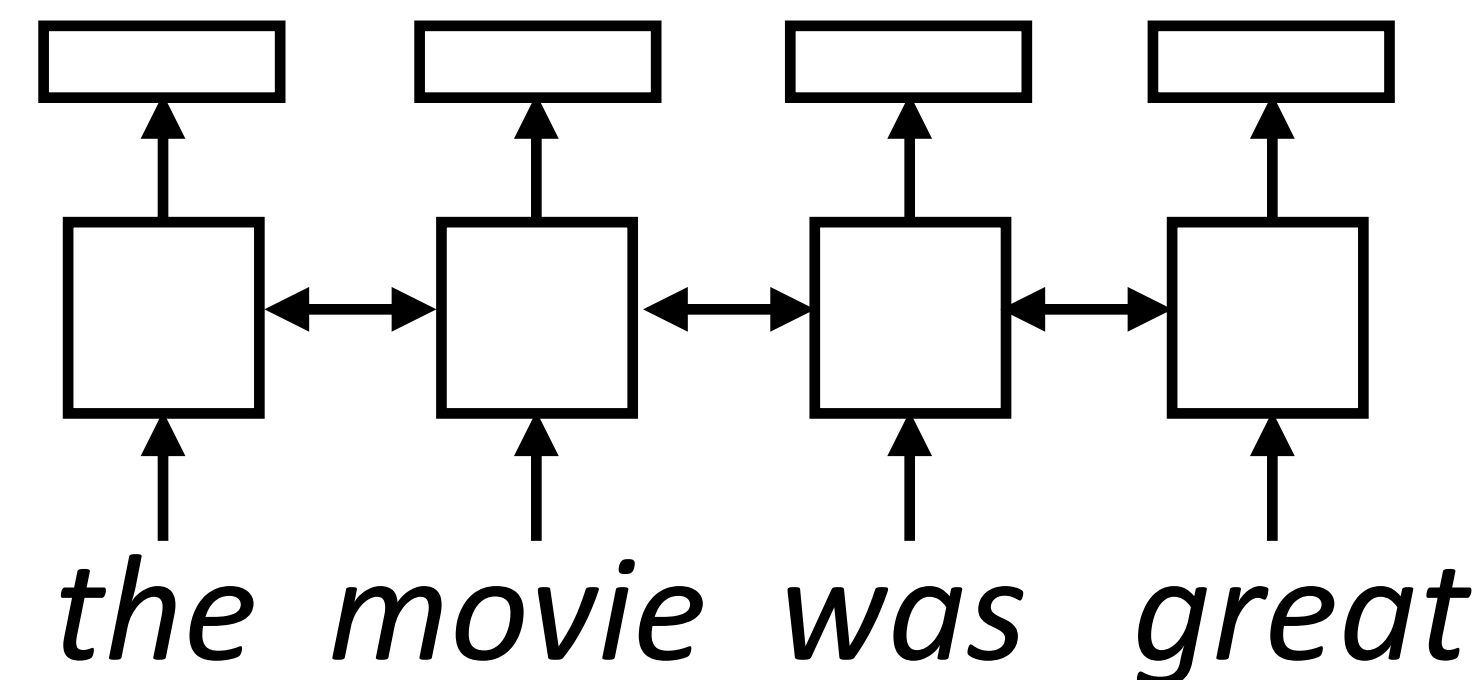
Ashish Vaswani* Google Brain avaswani@google.com	Noam Shazeer* Google Brain noam@google.com	Niki Parmar* Google Research nikip@google.com	Jakob Uszkoreit* Google Research usz@google.com
Llion Jones* Google Research llion@google.com	Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu	Lukasz Kaiser* Google Brain lukaszkaizer@google.com	
Illia Polosukhin* † 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 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Sentence Encoders

- ▶ LSTM abstraction: maps each vector in a sentence to a new, context-aware vector
- ▶ CNNs do something similar with filters
- ▶ Attention can give us a third way to do this



Self-Attention

- ▶ Assume we're using GloVe — what do we want our neural network to do?



*The ballerina is very excited that **she** will dance in the **show**.*

- ▶ What words need to be contextualized here?
 - ▶ Pronouns need to look at antecedents
 - ▶ Ambiguous words should look at context
 - ▶ Words should look at syntactic parents/children
- ▶ Problem: LSTMs and CNNs don't do this

Self-Attention

- ▶ Want:

*The ballerina is very excited that **she** will dance in the **show**.*



The diagram illustrates long-range dependencies in the sentence. A blue arc connects the word "that" to the word "she", and a red arc connects the word "will" to the word "show".

- ▶ LSTMs/CNNs: tend to look at local context

*The ballerina is very excited that **she** will dance in the **show**.*



The diagram illustrates local context dependencies. Multiple blue arcs connect adjacent words: "The" to "ballerina", "ballerina" to "is", "is" to "very", "very" to "excited", "excited" to "that", "that" to "she", "she" to "will", "will" to "dance", "dance" to "in", "in" to "the", and "the" to "show". Additionally, a red arc connects "will" to "show".

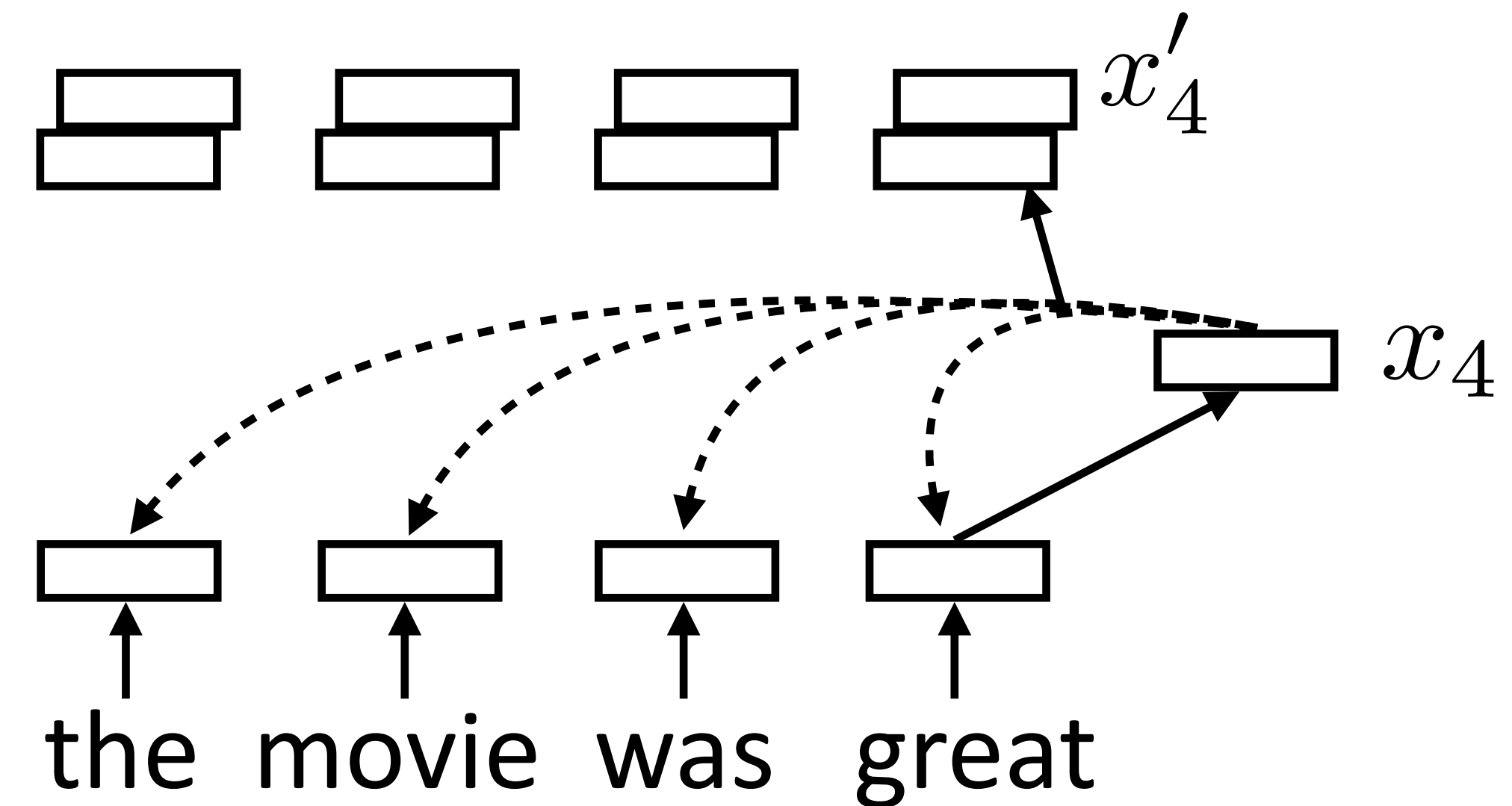
- ▶ To appropriately contextualize embeddings, we need to pass information over long distances dynamically for each word

Self-Attention

- Each word forms a “query” which then computes attention over each word

$$\alpha_{i,j} = \text{softmax}(x_i^\top x_j) \quad \text{scalar}$$

$$x'_i = \sum_{j=1}^n \alpha_{i,j} x_j \quad \text{vector} = \text{sum of scalar} * \text{vector}$$



- Multiple “heads” analogous to different convolutional filters. Use parameters W_k and V_k to get different attention values + transform vectors

$$\alpha_{k,i,j} = \text{softmax}(x_i^\top W_k x_j) \quad x'_{k,i} = \sum_{j=1}^n \alpha_{k,i,j} V_k x_j$$

What can self-attention do?

*The ballerina is very excited that **she** will dance in the **show**.*

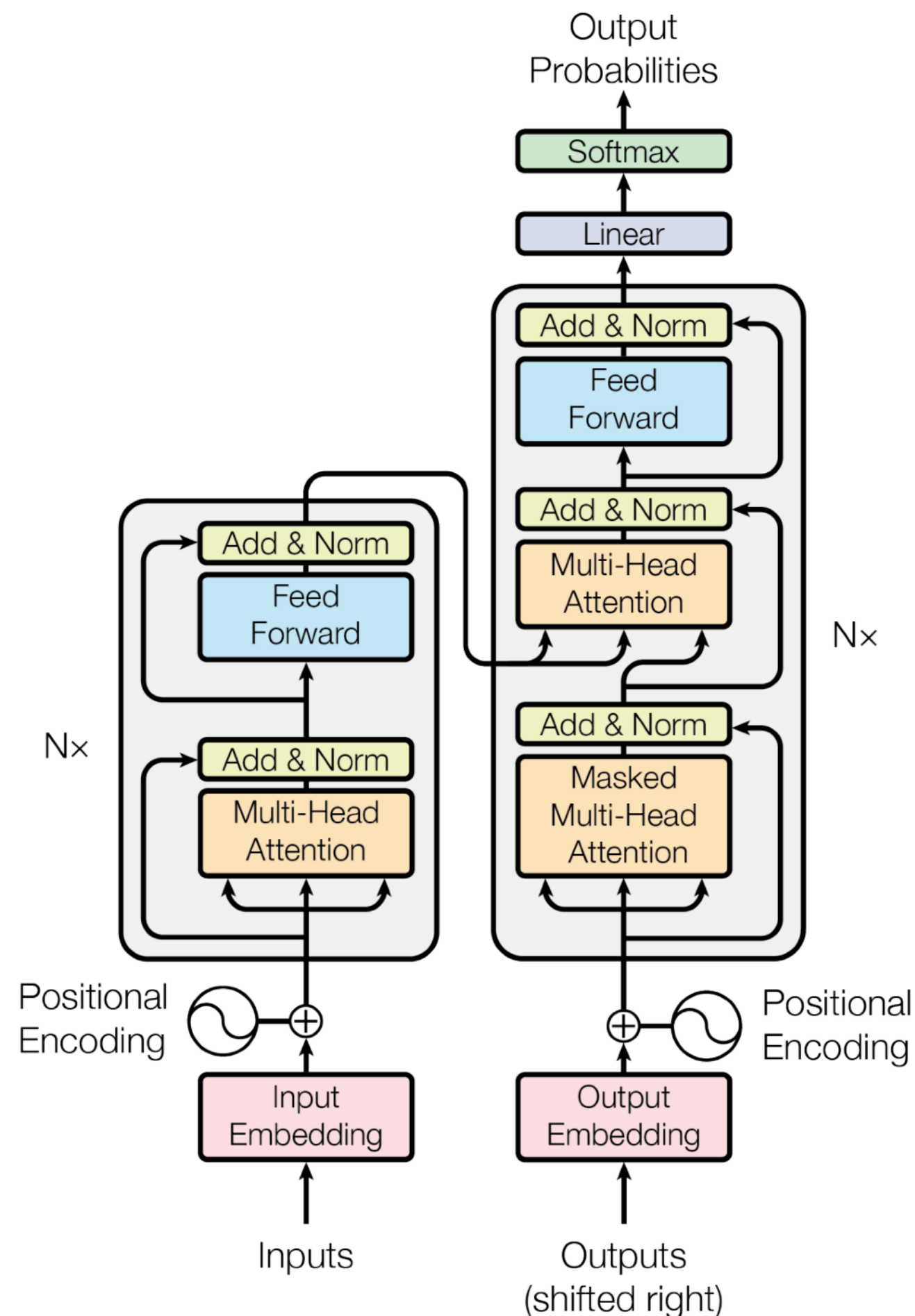


0	0.5	0	0	0.1	0.1	0	0.1	0.2	0	0	0
0	0.1	0	0	0	0	0	0	0.5	0	0.4	0

- ▶ Attend nearby + to semantically related terms
- ▶ Why multiple heads? Softmaxes end up being peaked, single distribution cannot easily put weight on multiple things

Transformer Uses

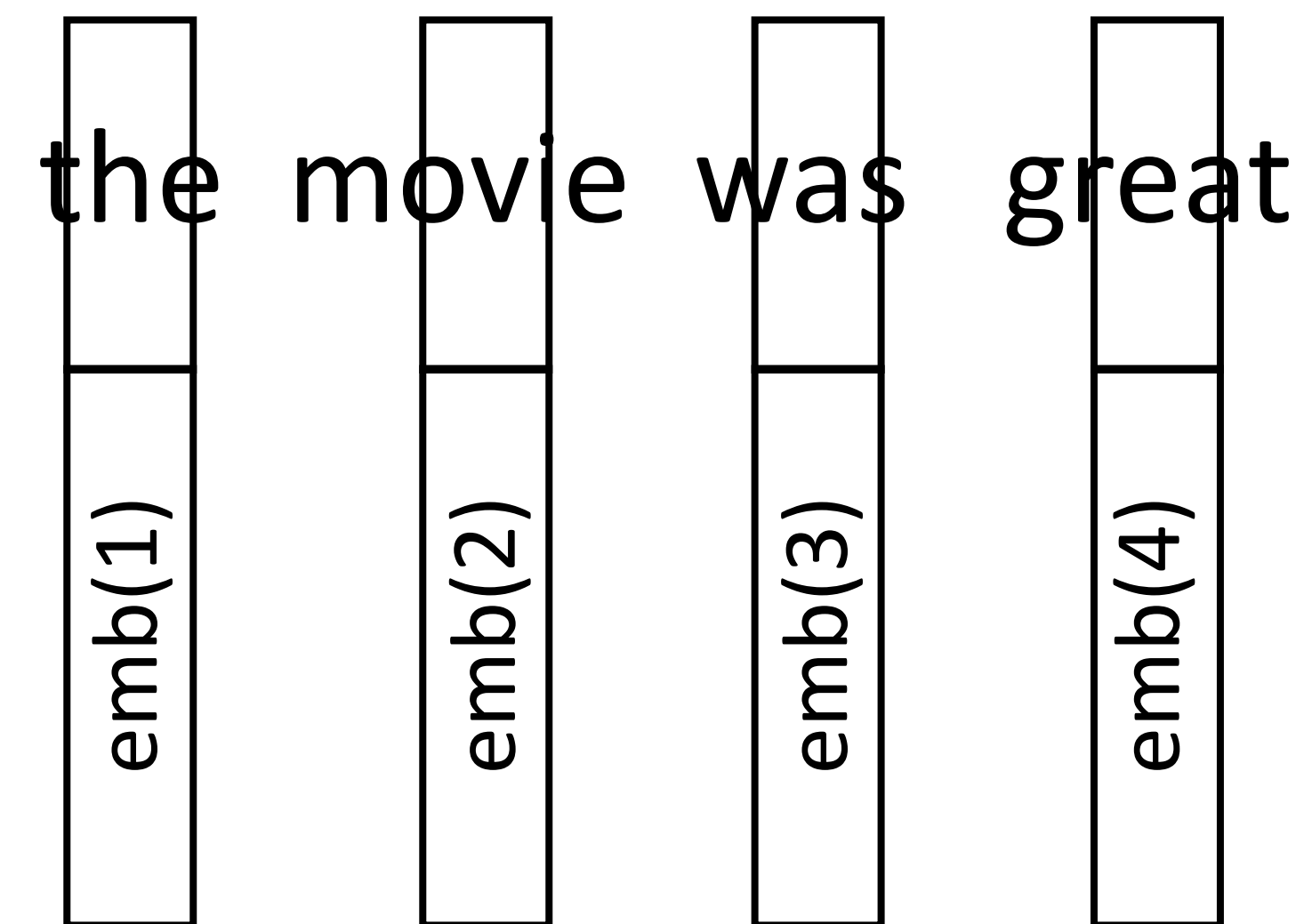
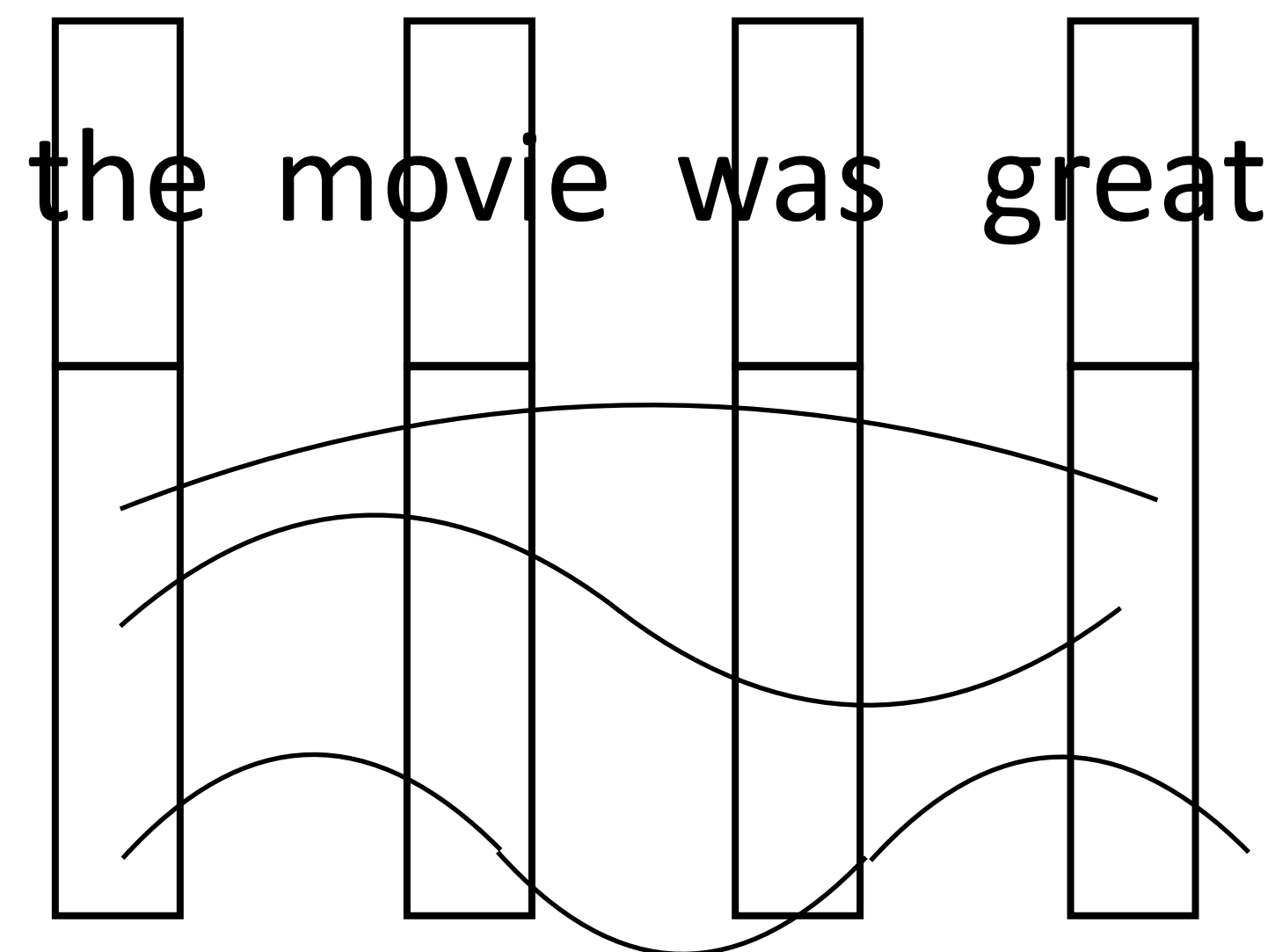
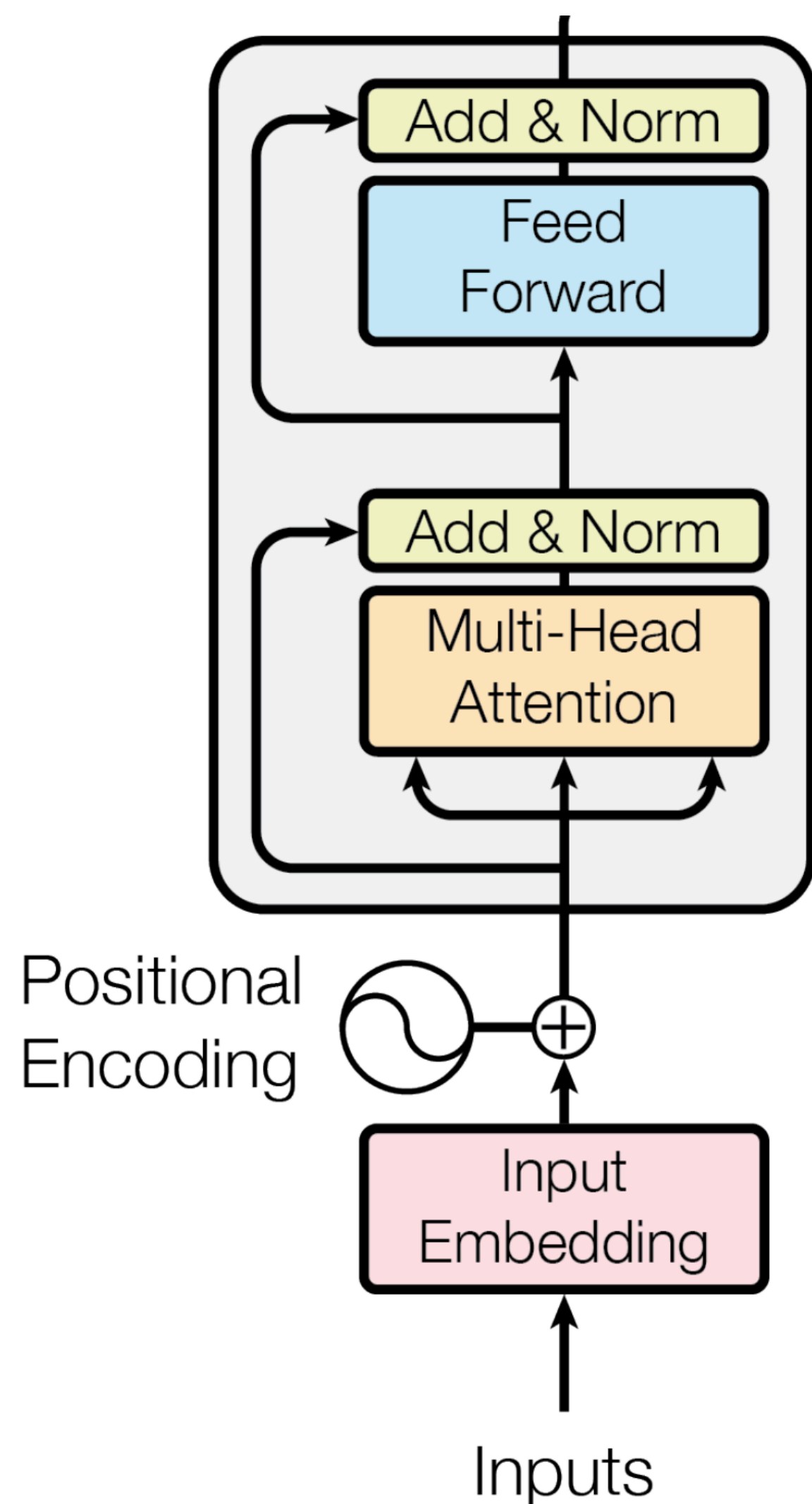
- ▶ Supervised: transformer can replace LSTM as encoder, decoder, or both; such as in machine translation and natural language generation tasks.



- ▶ Encoder and decoder are both transformers
- ▶ Decoder consumes the previous generated token (and attends to input), but has *no recurrent state*
- ▶ Many other details to get it to work: residual connections, layer normalization, positional encoding, optimizer with learning rate schedule, label smoothing

Vaswani et al. (2017)

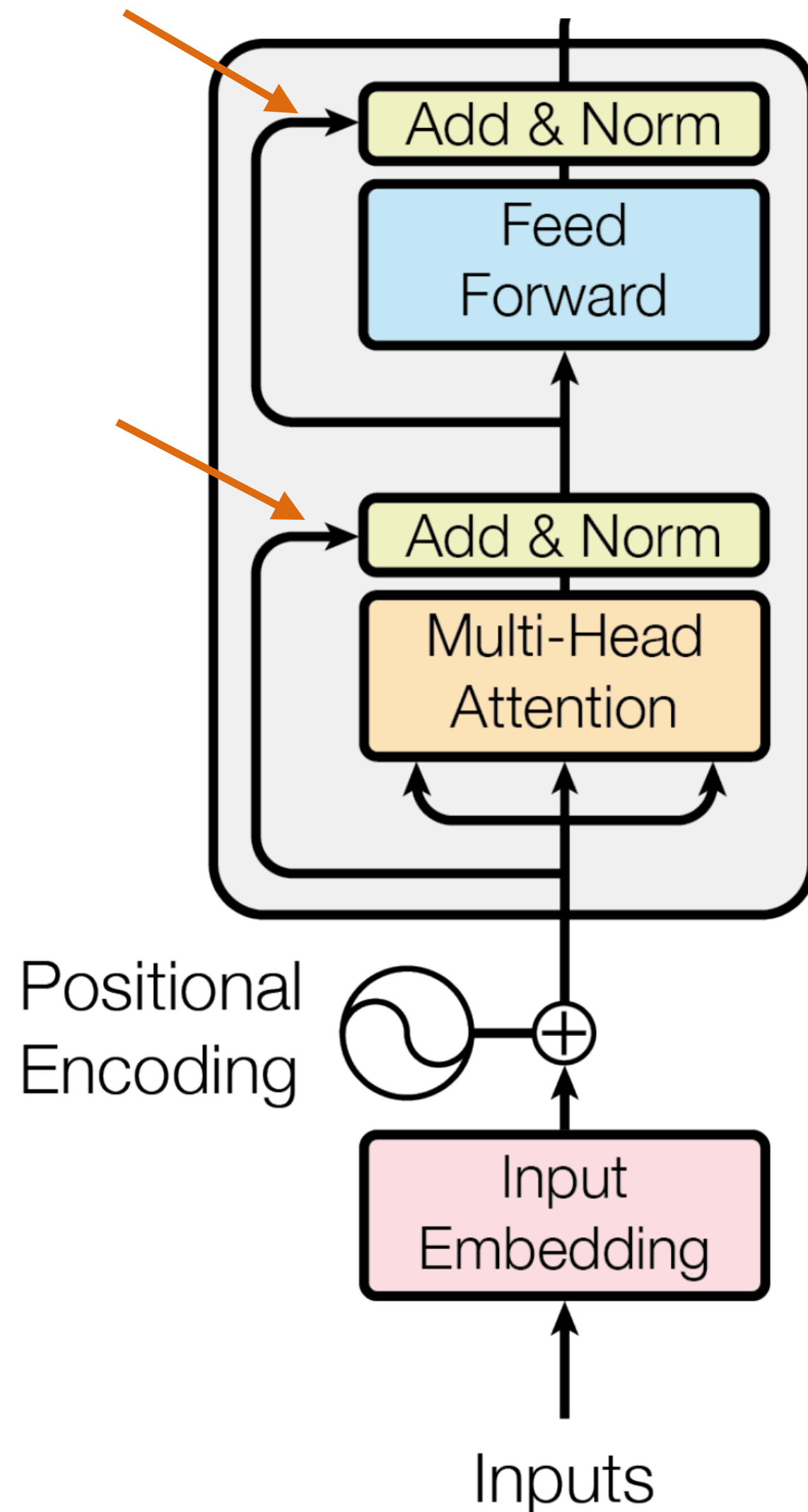
Transformers



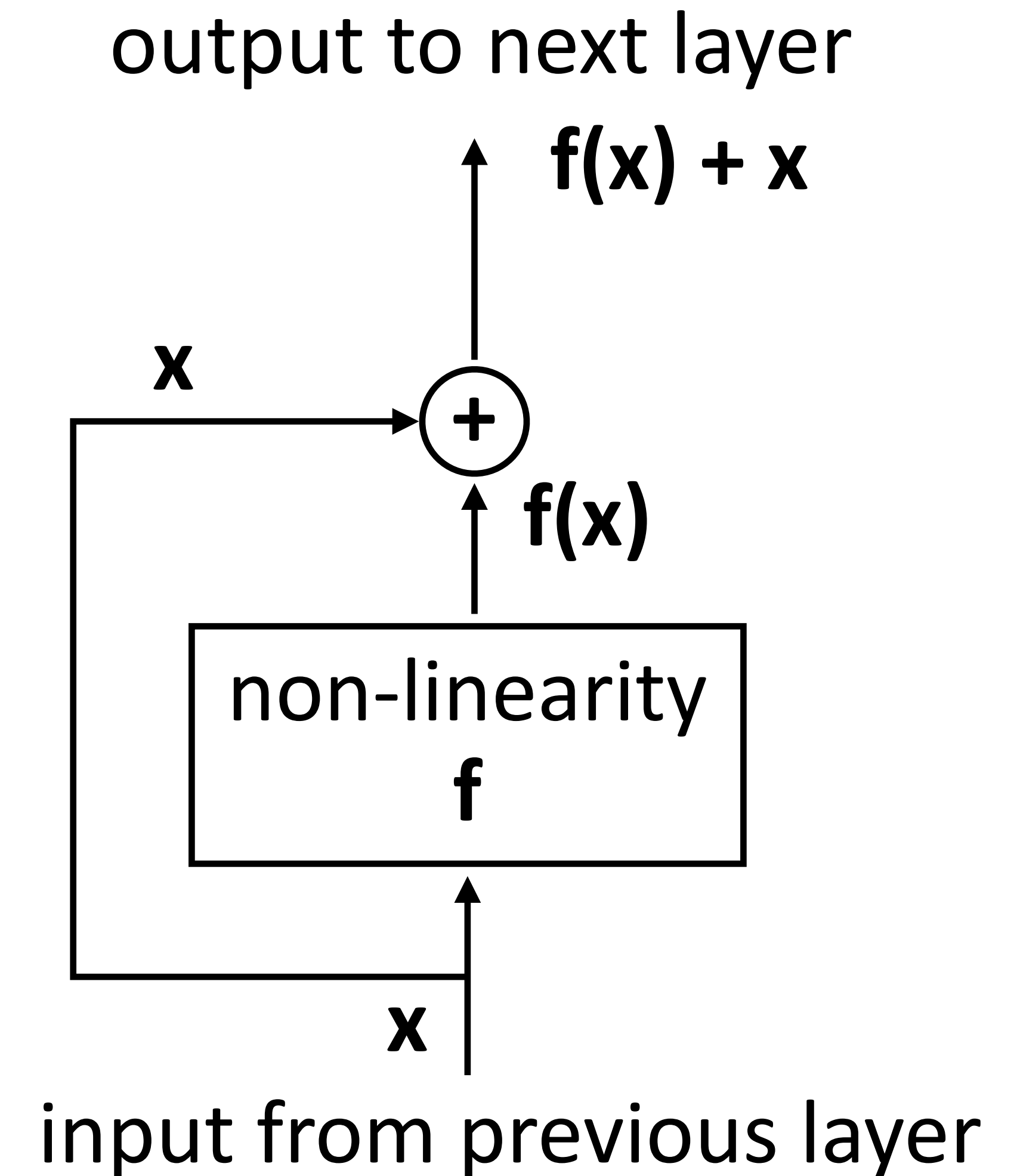
- ▶ Augment word embedding with position embeddings, each dim is a sine/cosine wave of a different frequency. Closer points = higher dot products
- ▶ Works essentially as well as just encoding position as a one-hot vector

Vaswani et al. (2017)

Residual Connections



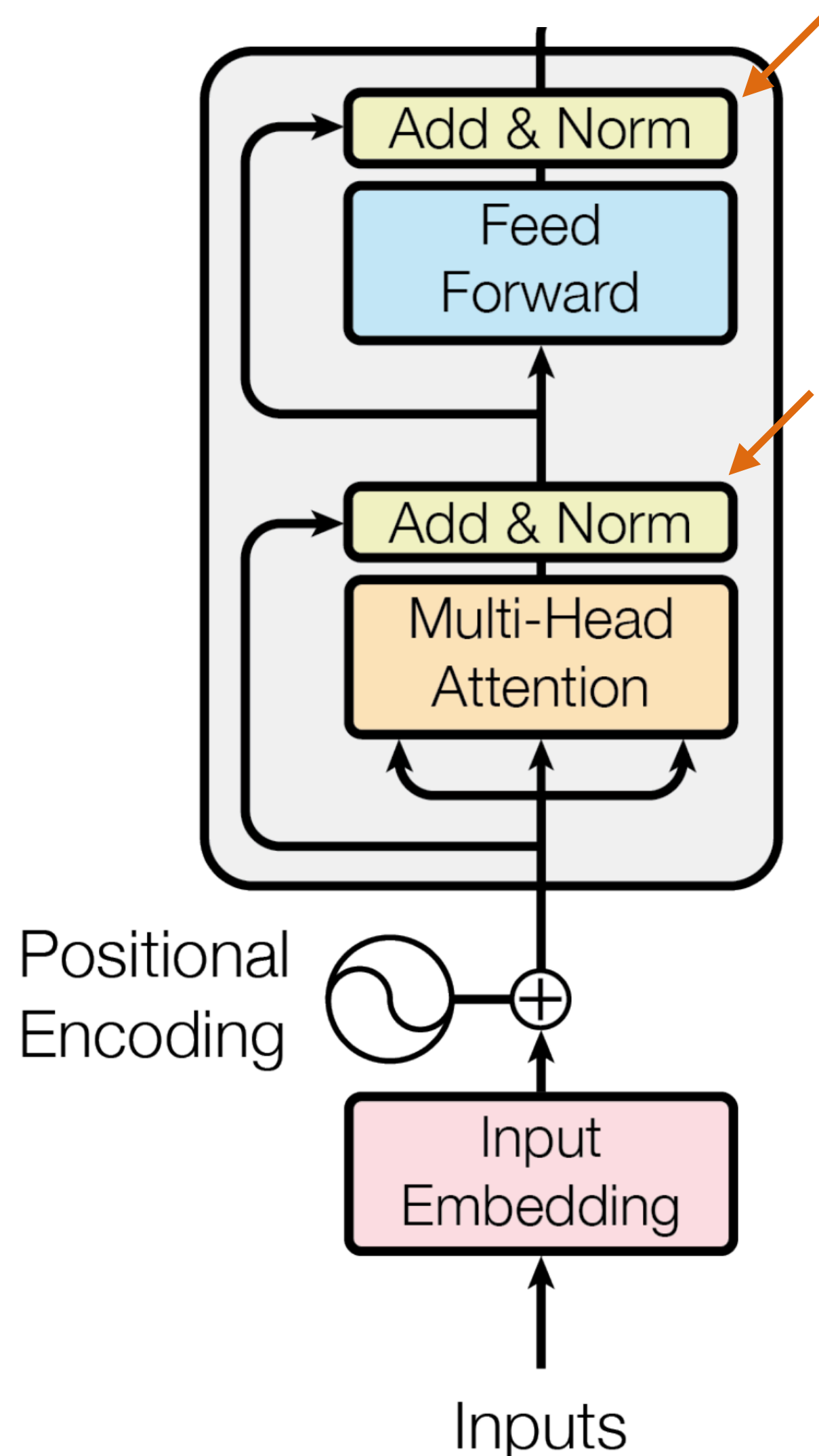
- ▶ allow gradients to flow through a network directly, without passing through non-linear activation functions



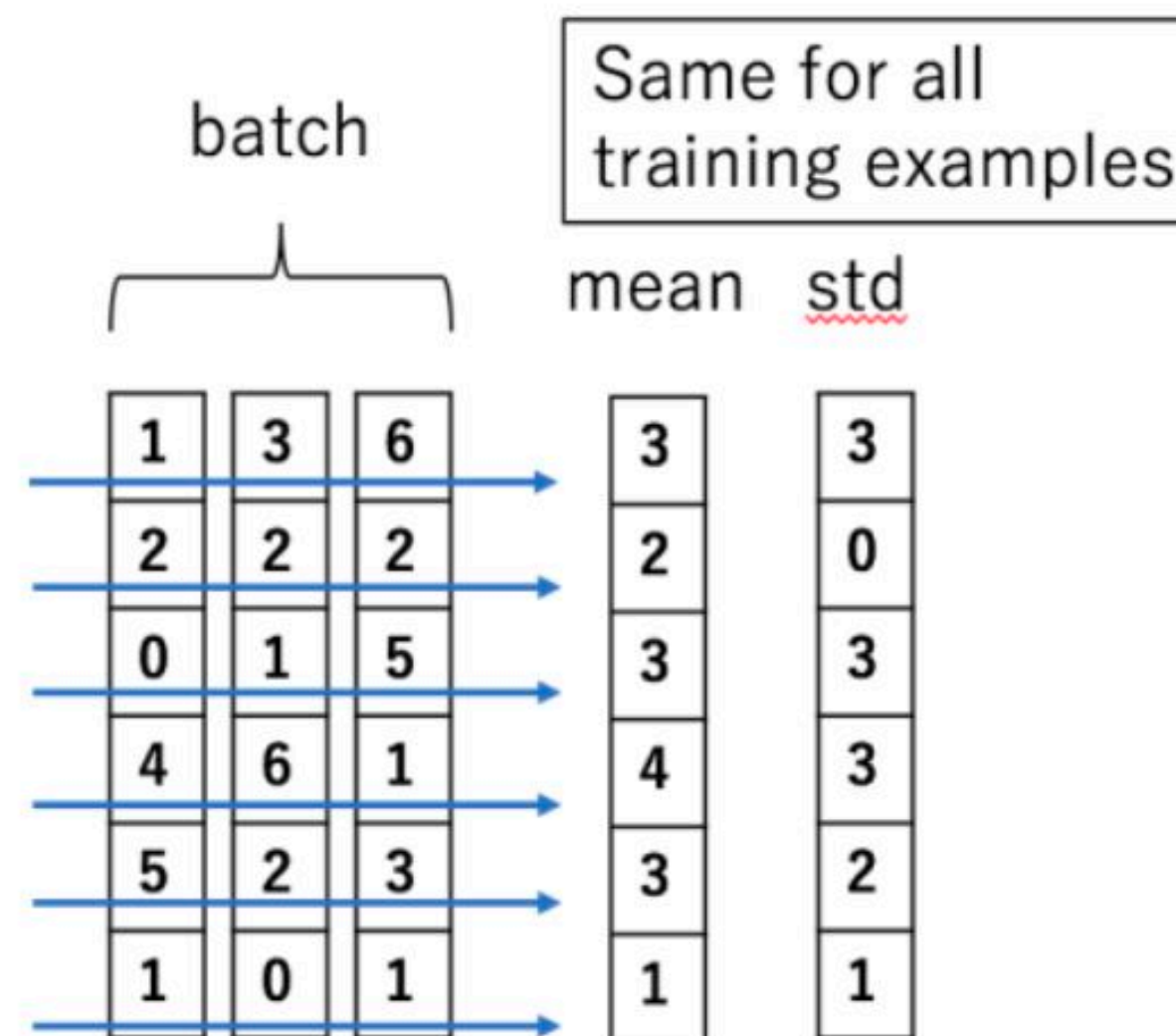
He et al. (2015)

Layer Normalization

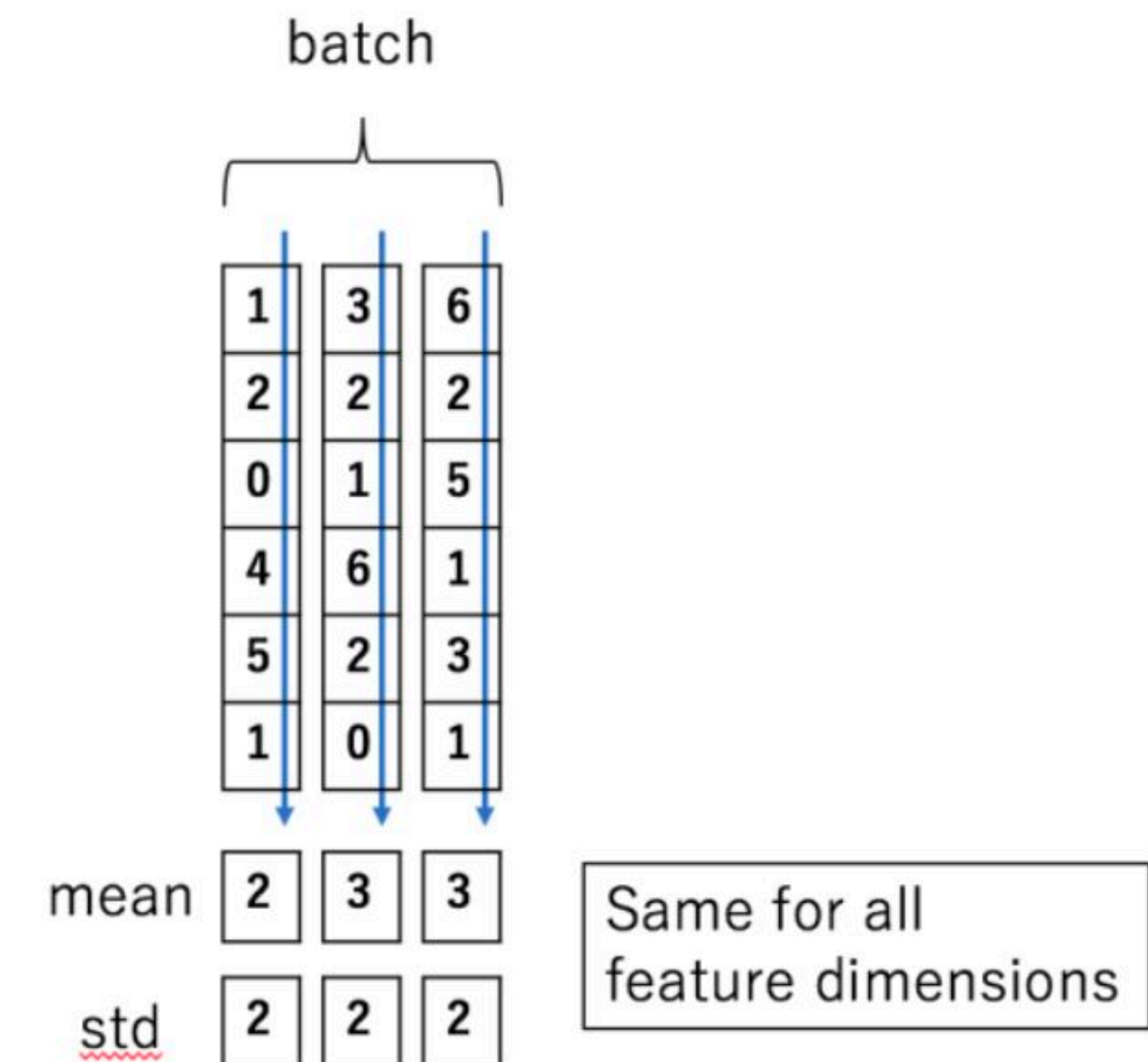
- ▶ subtract mean, divide by variance



Batch Normalization



Layer Normalization



Label Smoothing

- ▶ Instead of using a one-hot target distribution, create a distribution that has “confidence” of the correct word and the rest of the “smoothing” mass distributed throughout the vocabulary.
- ▶ Implemented by minimizing KL-divergence between smoothed ground truth probabilities and the probabilities computed by model.

I went to class and took _____

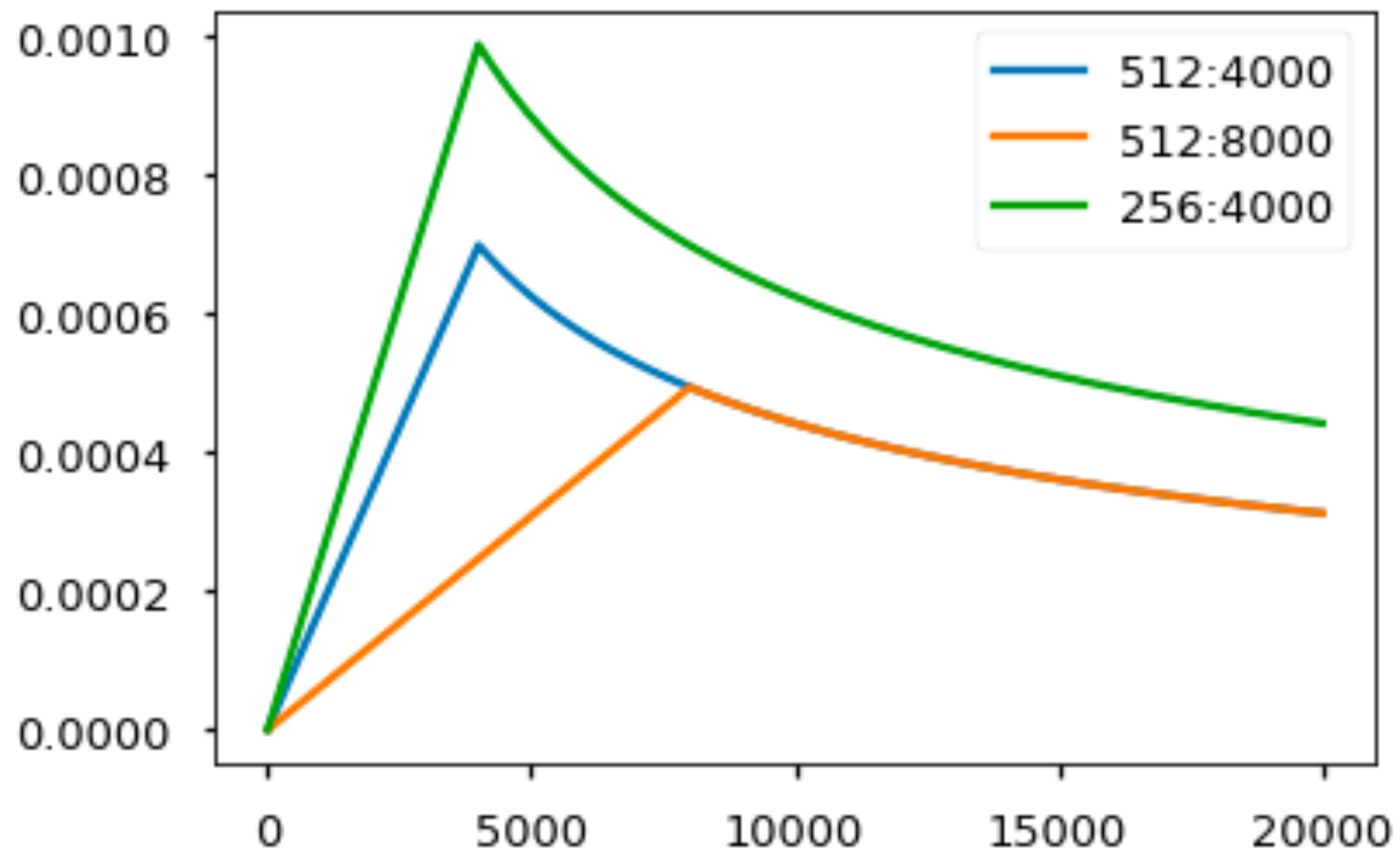
<i>cats</i>	<i>TV</i>	<i>notes</i>	<i>took</i>	<i>sofa</i>
0	0	1	0	0
0.025	0.025	0.9	0.025	0.025

← one-hot

← with label smoothing

Szegedy et al. (2015)

Transformers



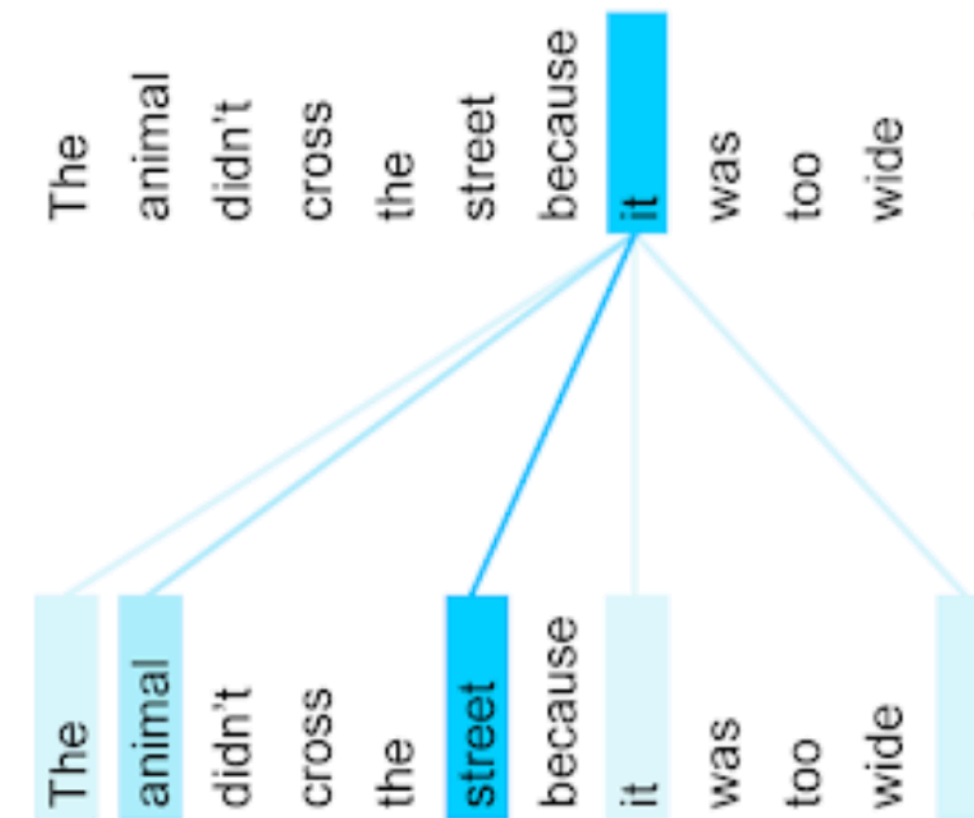
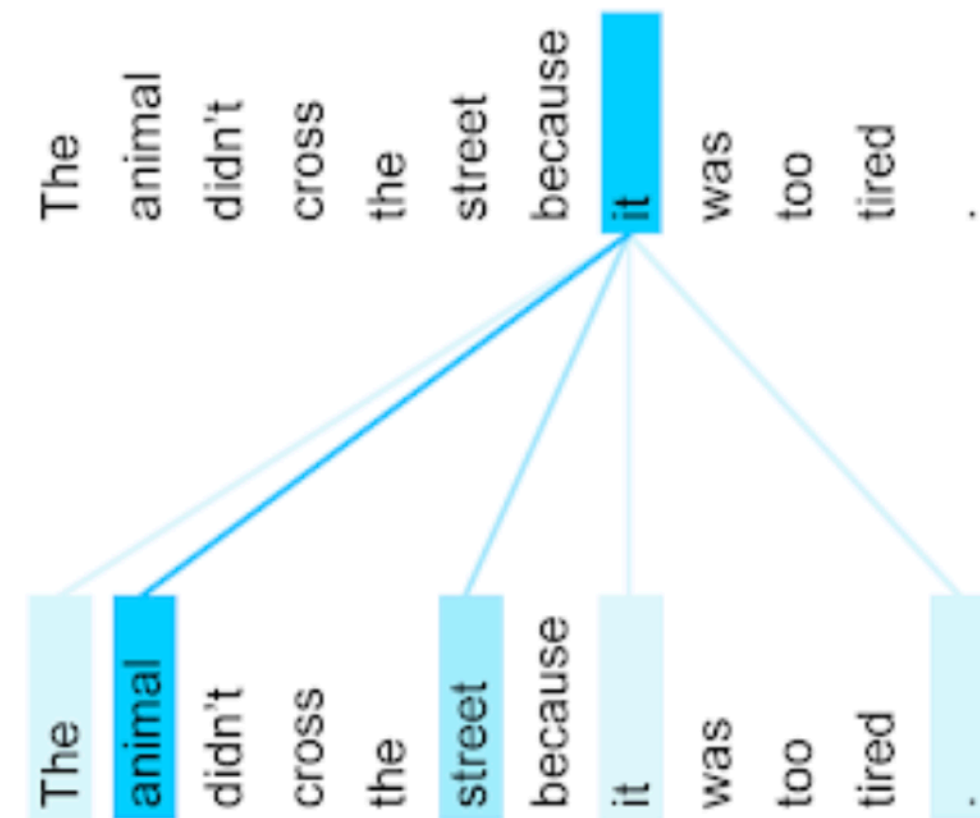
- ▶ Adam optimizer with varied learning rate over the course of training
- ▶ Linearly increase for warmup, then decay proportionally to the inverse square root of the step number
- ▶ This part is very important!

Transformers

Model	BLEU	
	EN-DE	EN-FR
ByteNet [18]	23.75	
Deep-Att + PosUnk [39]		39.2
GNMT + RL [38]	24.6	39.92
ConvS2S [9]	25.16	40.46
MoE [32]	26.03	40.56
Deep-Att + PosUnk Ensemble [39]		40.4
GNMT + RL Ensemble [38]	26.30	41.16
ConvS2S Ensemble [9]	26.36	41.29
Transformer (base model)	27.3	38.1
Transformer (big)	28.4	41.8

- Big = 6 layers, 1000 dim for each token, 16 heads, base = 6 layers + other params halved

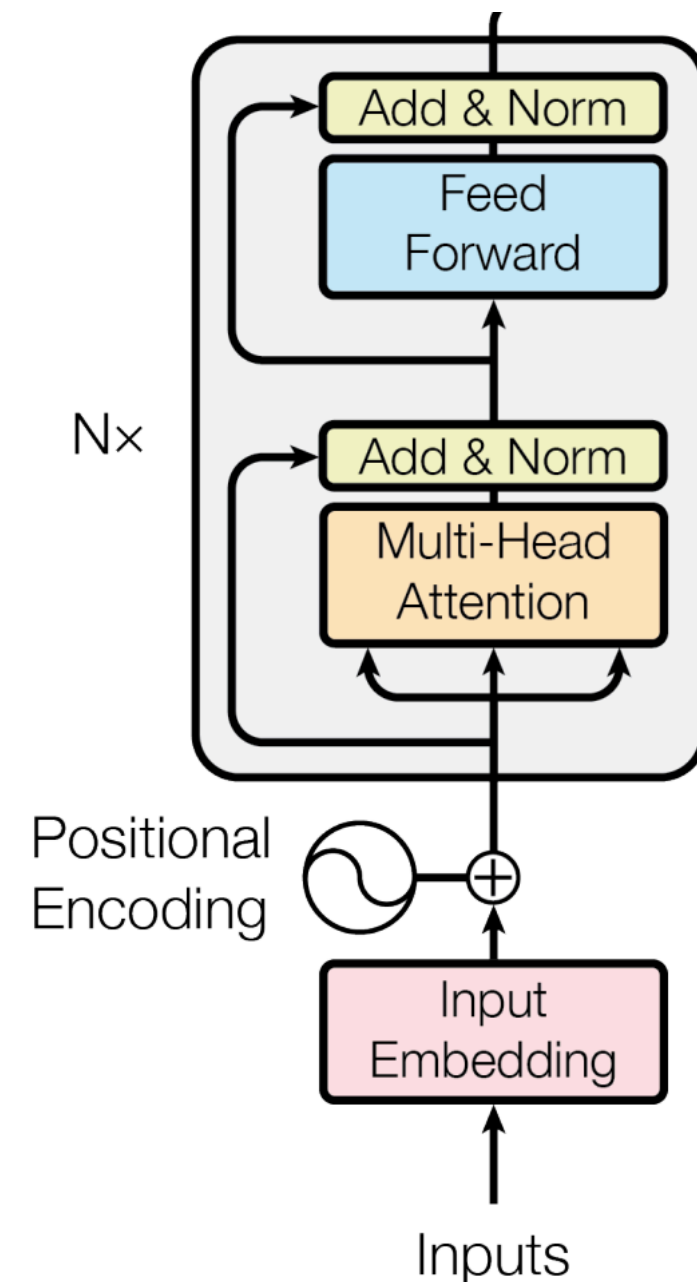
Visualization



<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

Other Transformer Variations

- ▶ Multilayer transformer networks consist of interleaved self-attention and feedforward sublayers.
- ▶ Could ordering the sublayers in a different pattern lead to better performance?



s f s f s f s f s f s f s f s f s f s f s f s f

(a) Interleaved Transformer

s s s s s s s f s f s f s f s f s f s f s f f f f f f f

(b) Sandwich Transformer

Figure 1: A transformer model (a) is composed of interleaved self-attention (green) and feedforward (purple) sublayers. Our sandwich transformer (b), a reordering of the transformer sublayers, performs better on language modeling. Input flows from left to right.