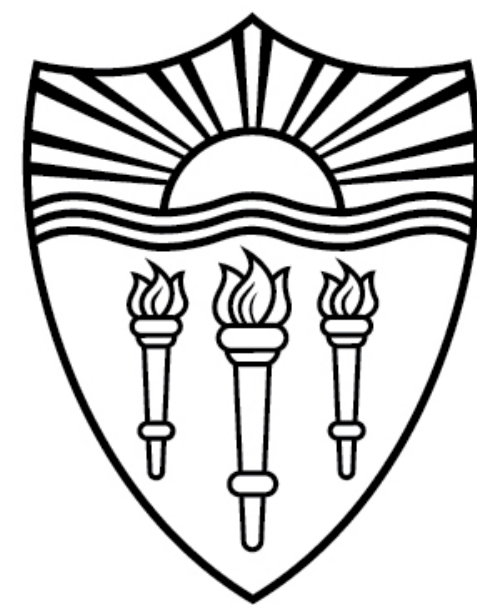


CSCI 544: Applied Natural Language Processing

Transformer-II

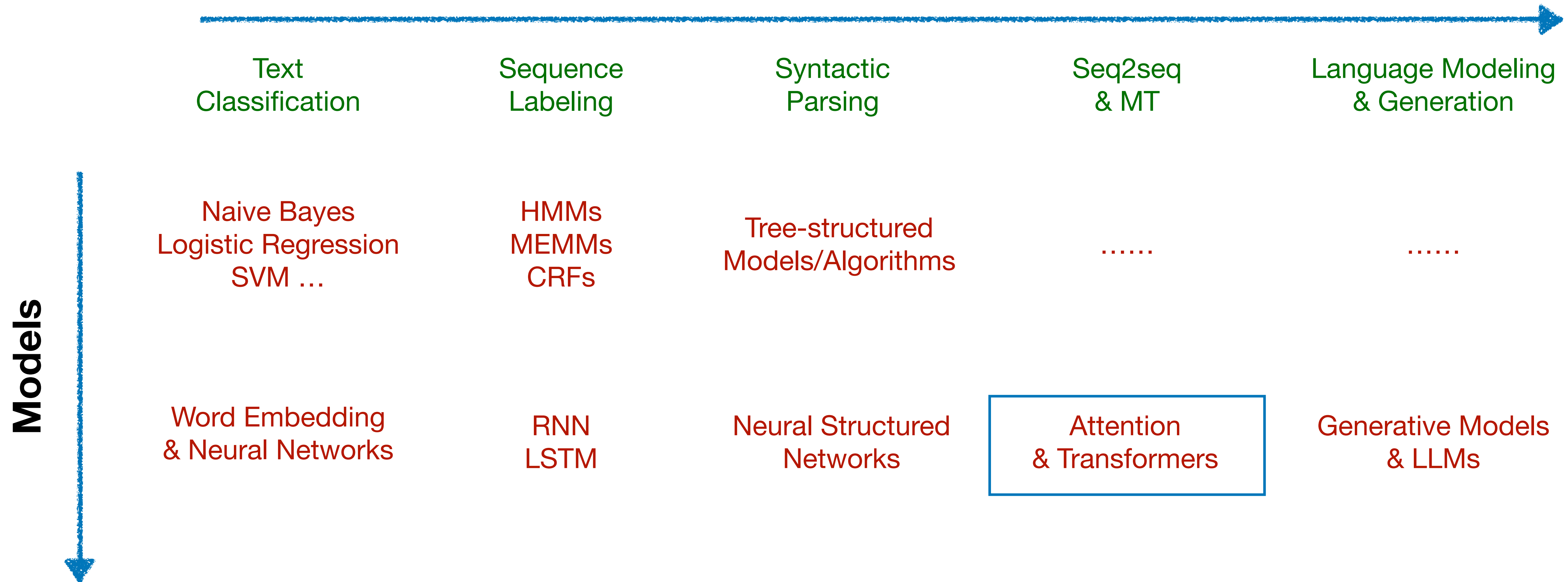
Xuezhe Ma (Max)



USC University of
Southern California

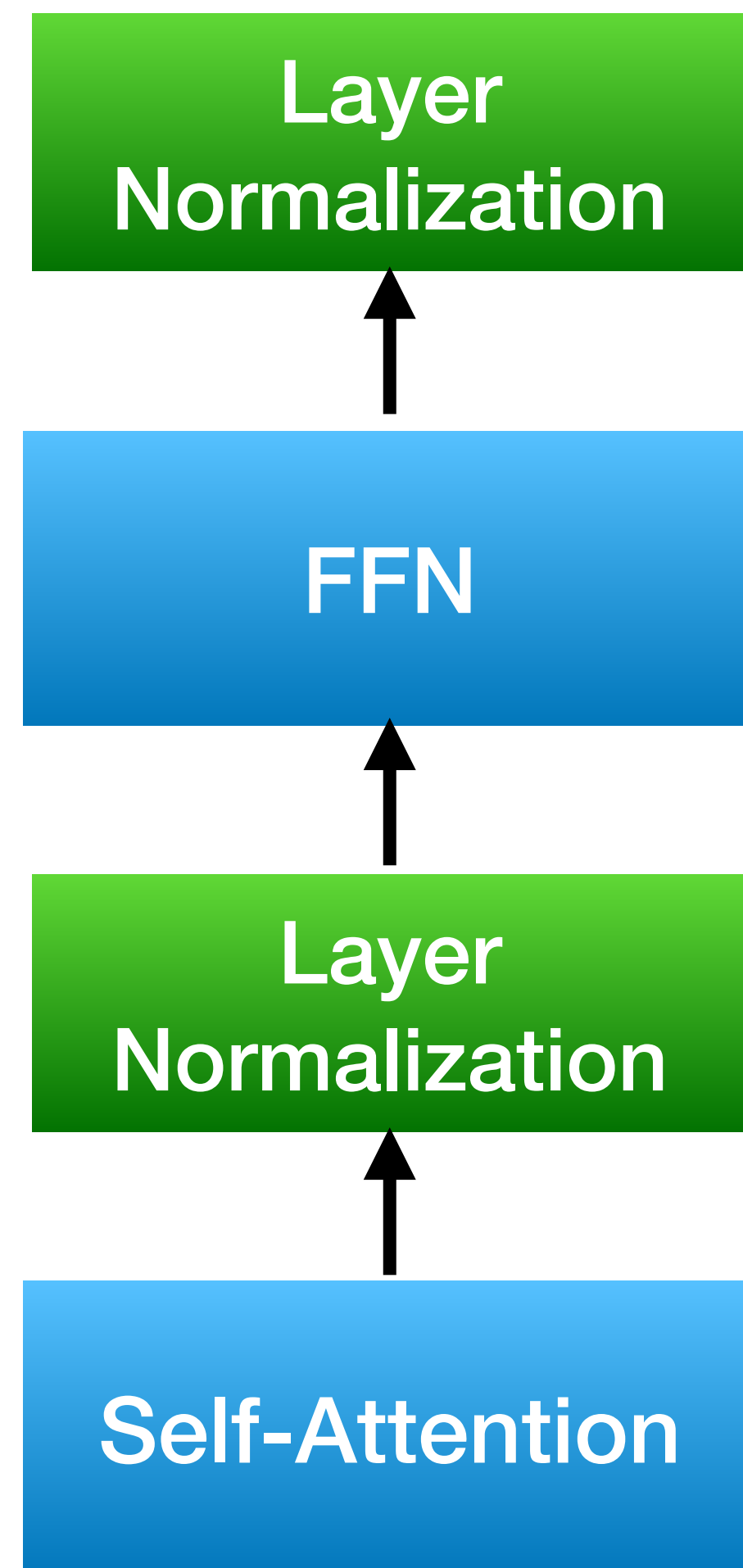
Course Organization

NLP Tasks



Transformers

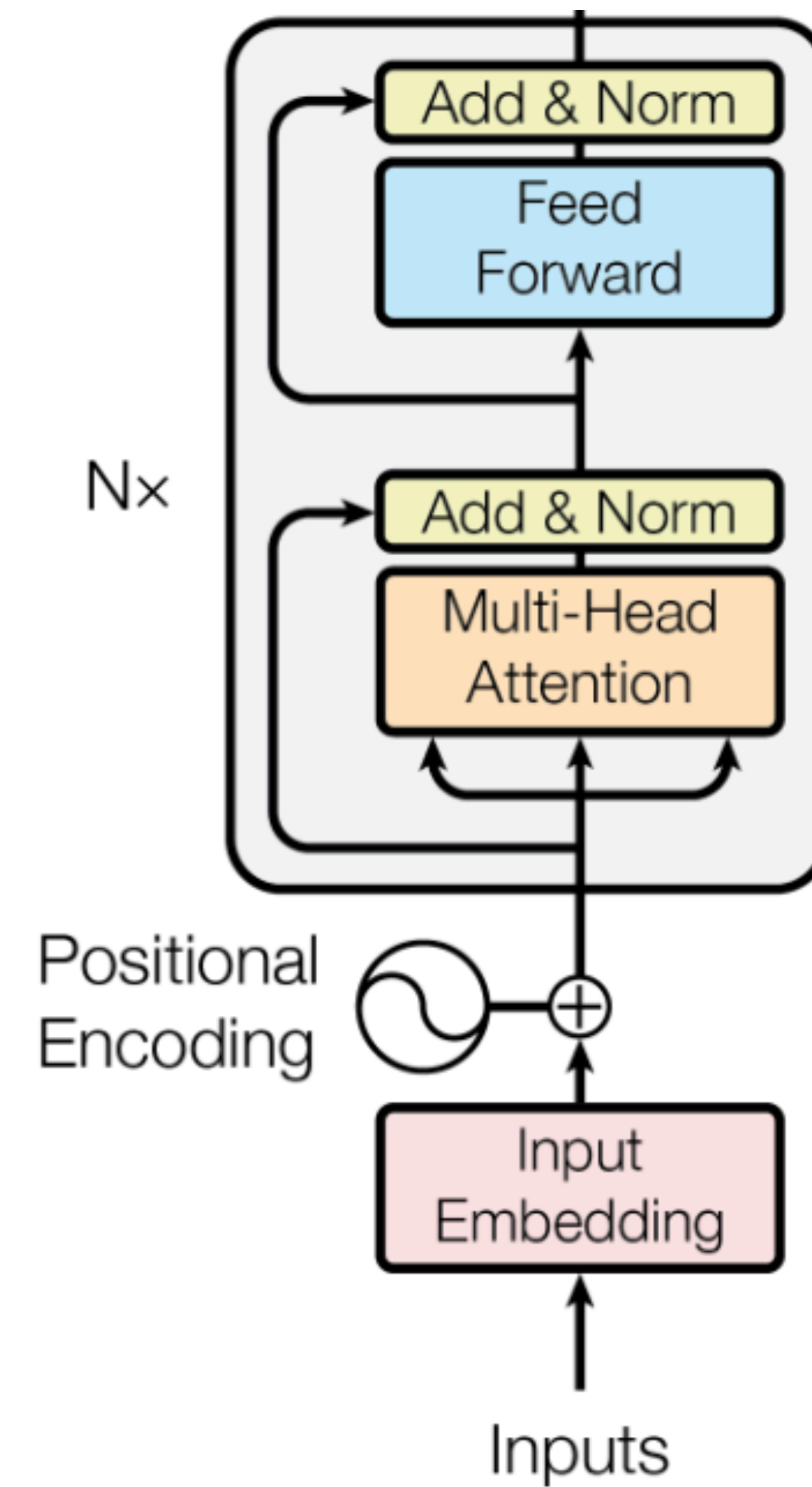
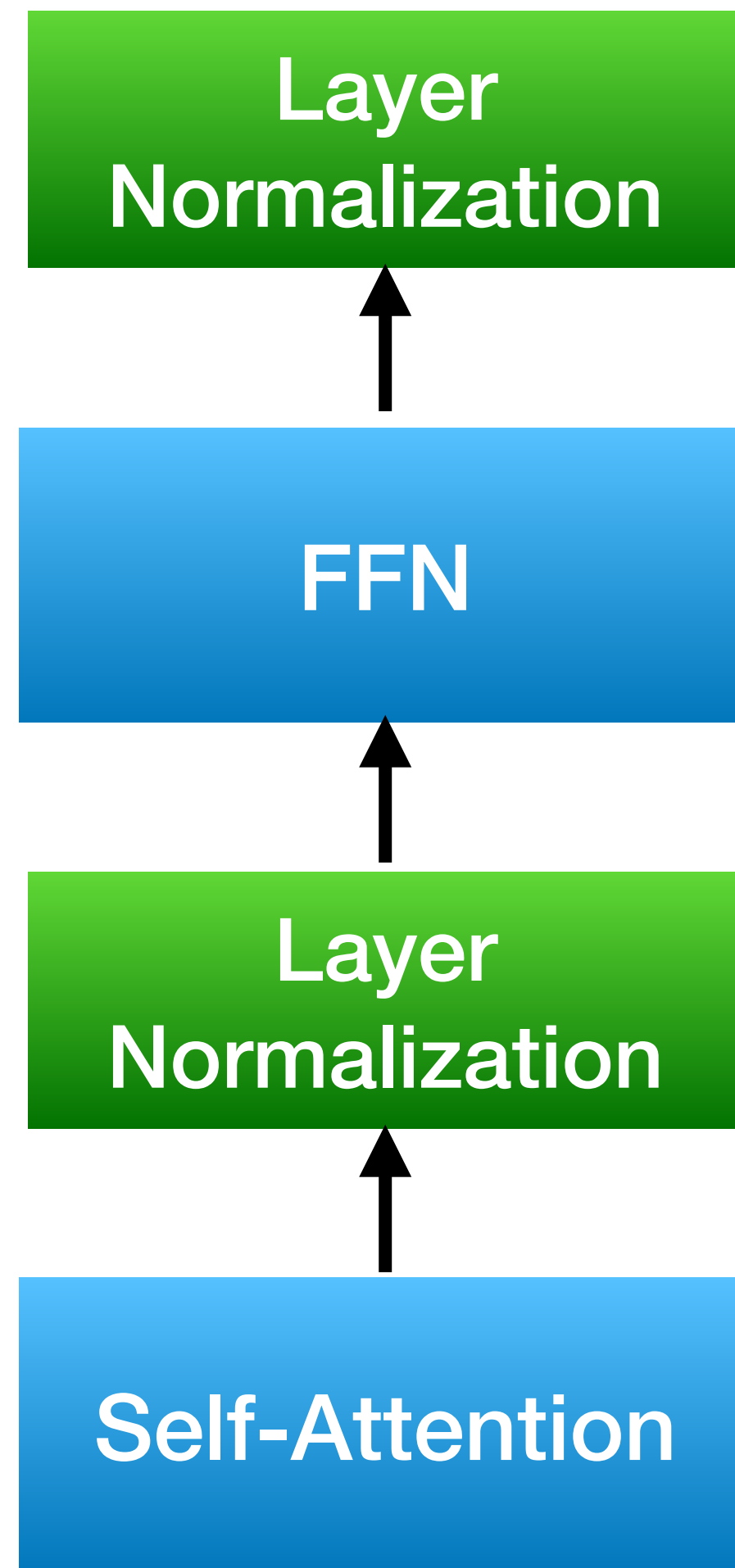
Transformer Encoder Block



- **Three Key Components**

- (Masked) Multi-head Self-Attention
- Layer Normalization
- Position-wise Feed-Forward Network

Transformers



Outline

- **Residual Connections**
- **Normalization Layers**
 - Variants of Normalization Layers
 - Pre-Norm vs. Post-Norm
- **Positional Embeddings**
 - Absolute Positional Embeddings
 - Rotary Positional Embeddings

Residual Connections

Residual Connections

Deep Residual Learning for Image Recognition

Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun

Microsoft Research

{kahe, v-xiangz, v-shren, jiansun}@microsoft.com

Deep residual learning for image recognition

[K He](#), [X Zhang](#), [S Ren](#), [J Sun](#) - ... and **pattern recognition**, 2016 - [openaccess.thecvf.com](#)

... **Deeper** neural **networks** are more difficult to train. We present a **residual learning** framework to ease the training of **networks** that are substantially **deeper** than those used previously. ...

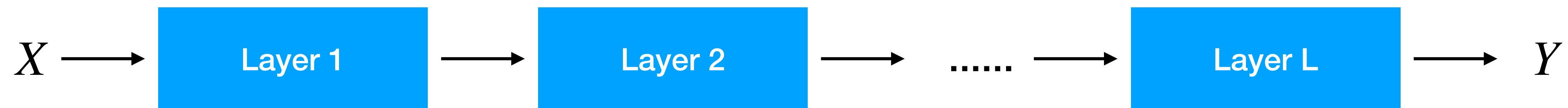
☆ Save 📄 Cite Cited by 200865 Related articles All 76 versions 🔗

Residual Connections

- Deep Neural Networks

$$Y_1 = f_1(X)$$

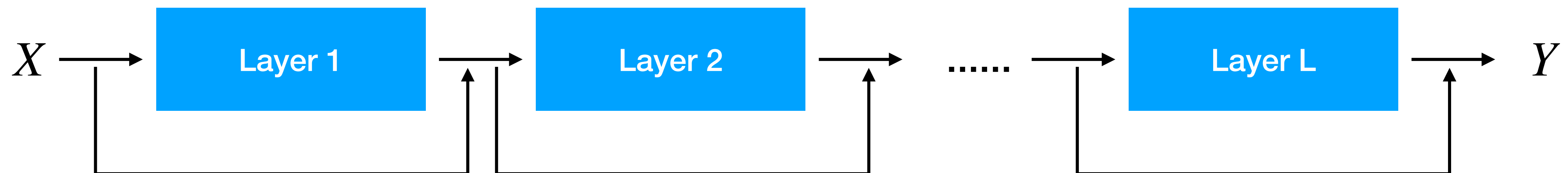
$$Y_2 = f_2(Y_1)$$



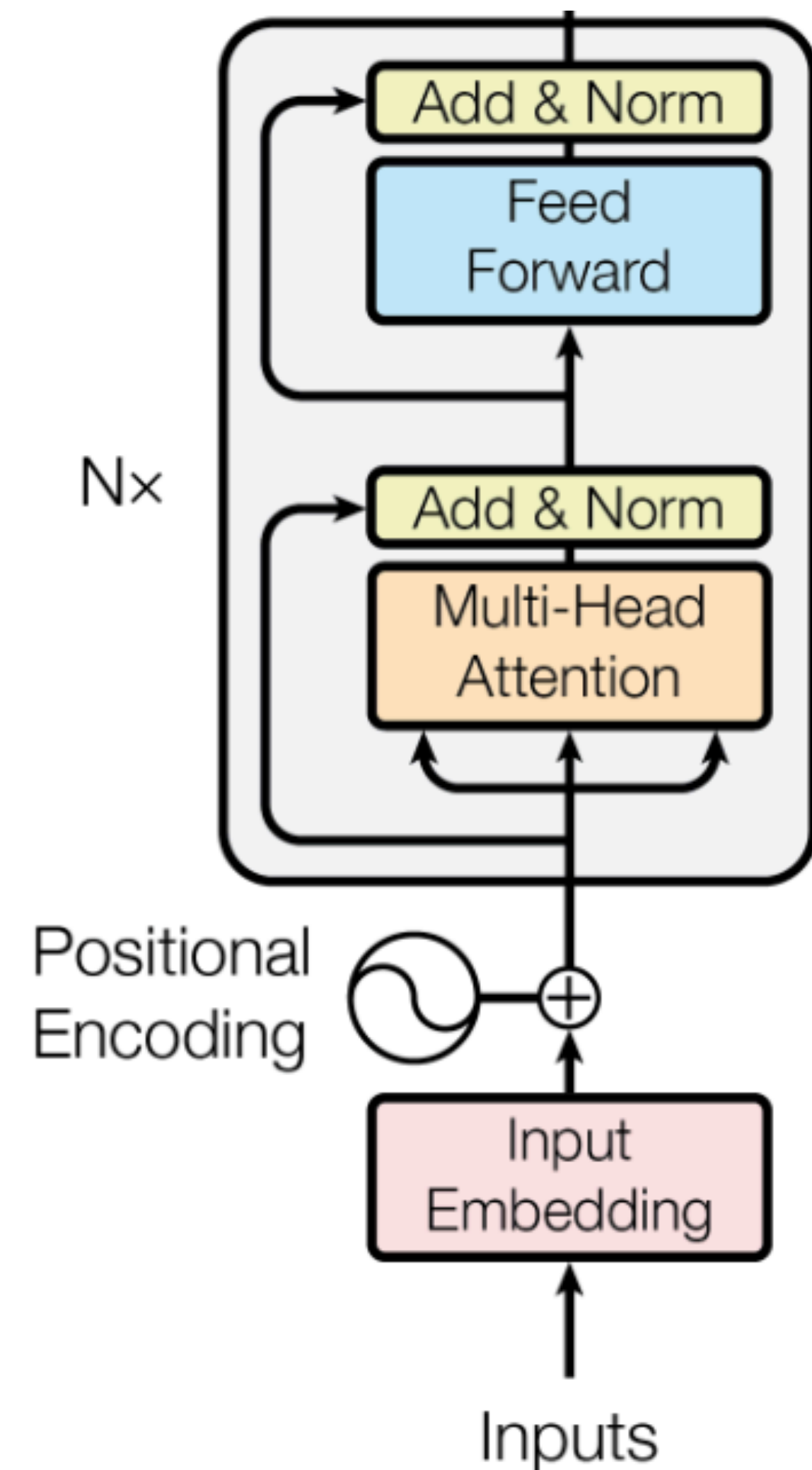
- Residual Connections

$$Y = f_1(X) + X$$

$$Y_2 = f_2(Y_1) + Y_1$$



Residual Connection in Transformers



$$Y = \text{LayerNorm} (\text{FFN}(X) + X)$$

$$Y = \text{LayerNorm} (\text{SelfAttn}(X) + X)$$

Outline

- **Residual Connections**
- **Normalization Layers**
 - Variants of Normalization Layers
 - Pre-Norm vs. Post-Norm
- **Positional Embeddings**
 - Absolute Positional Embeddings
 - Rotary Positional Embeddings

Normalization Layers

Normalization Layers

- **General Form:**

$$Y = \frac{X - E[X]}{\sqrt{\text{Var}[X] + \epsilon}} * \gamma + \beta$$

- **Variants of Normalization Layers**

- Along which dimensions to compute expectation and variance
- Batch Normalization
- Layer Normalization
-

Batch Normalization

- **First Proposed for CV models**
 - Lofe and Szegedy, 2014

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Sergey Ioffe
Christian Szegedy
Google, 1600 Amphitheatre Pkwy, Mountain View, CA 94043

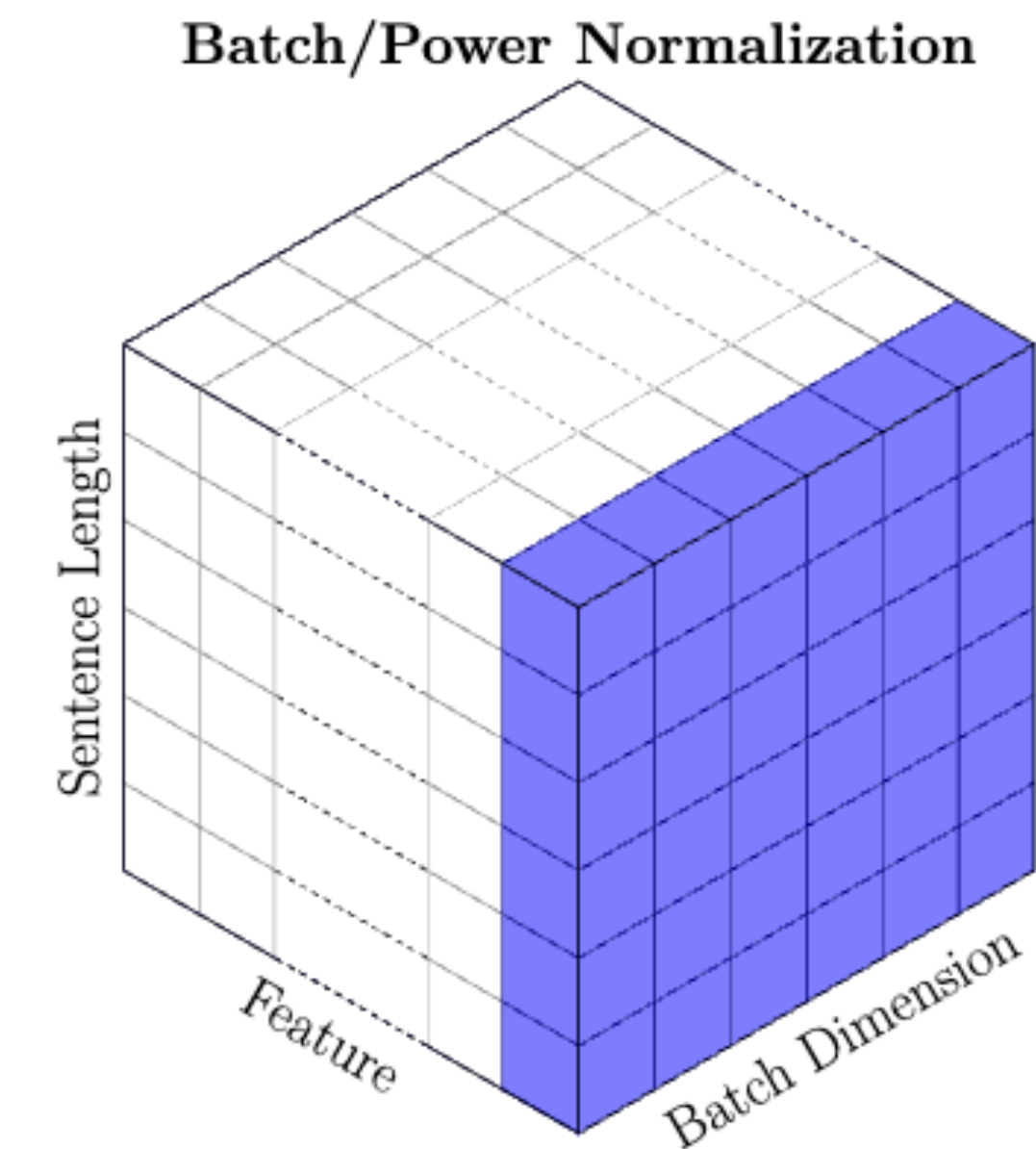
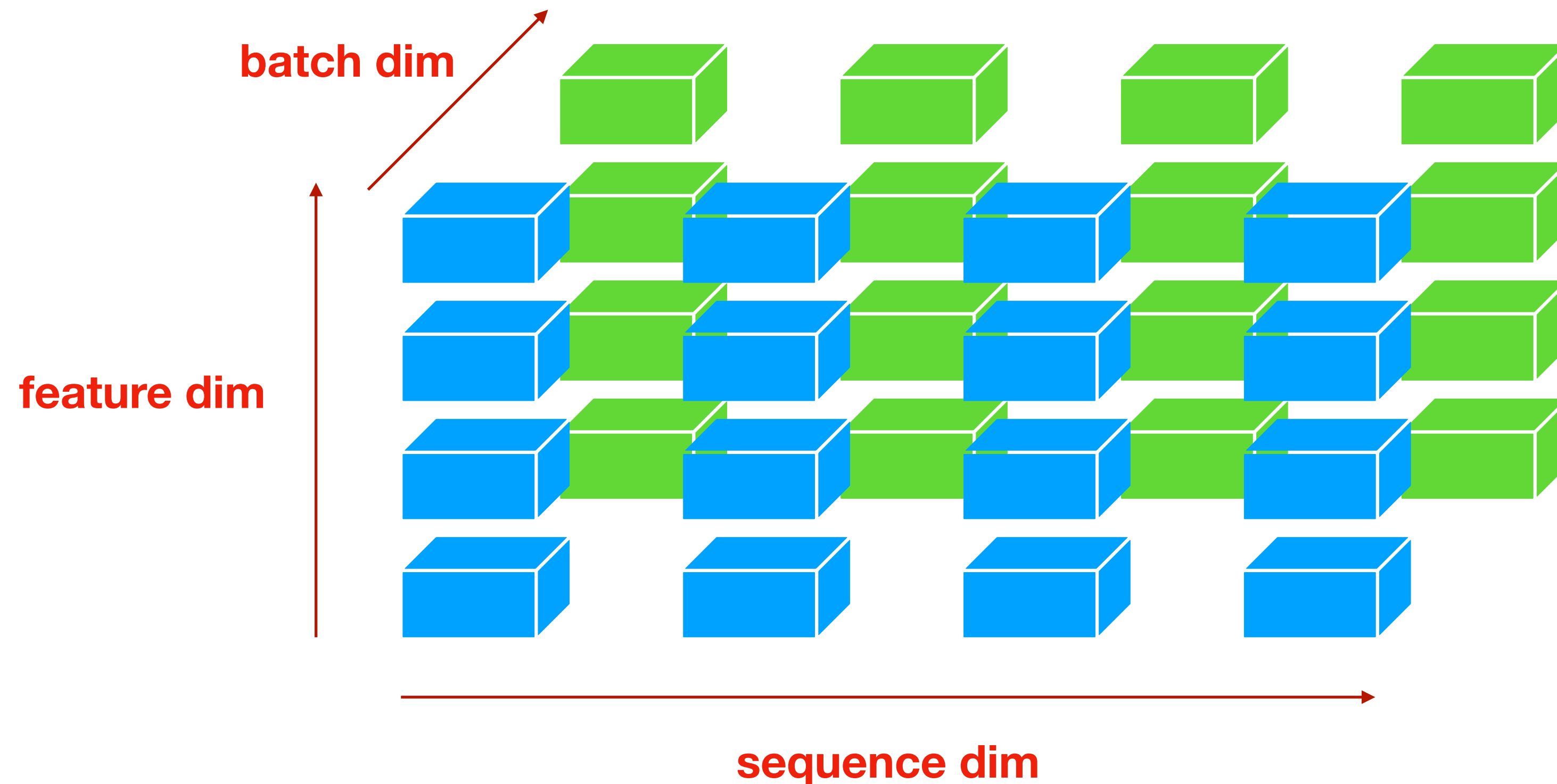
SIOFFE@GOOGLE.COM
SZEGEDY@GOOGLE.COM

- **Two Basic Motivations**
 - Reducing internal covariate shift (bias) from data
 - Controlling output range/scale of each layer

Batch Normalization

- Normalize along both **batch & spatial/sequential** dimensions

$$Y = \frac{X - E[X]}{\sqrt{\text{Var}[X] + \epsilon}} * \gamma + \beta$$



Issues of Batch Normalization

- **Normalization involves batch dimension**
 - Batch size cannot be too small
 - When batch size < 8 , batch normalization works poorly in practice
- **What are the expectation/variance in test time?**
 - Running stats for expectation/variance
 - At each mini-batch in training:
 - Compute mean and variance in this mini-batch

Hard to scale up

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

- Update running mean and variance

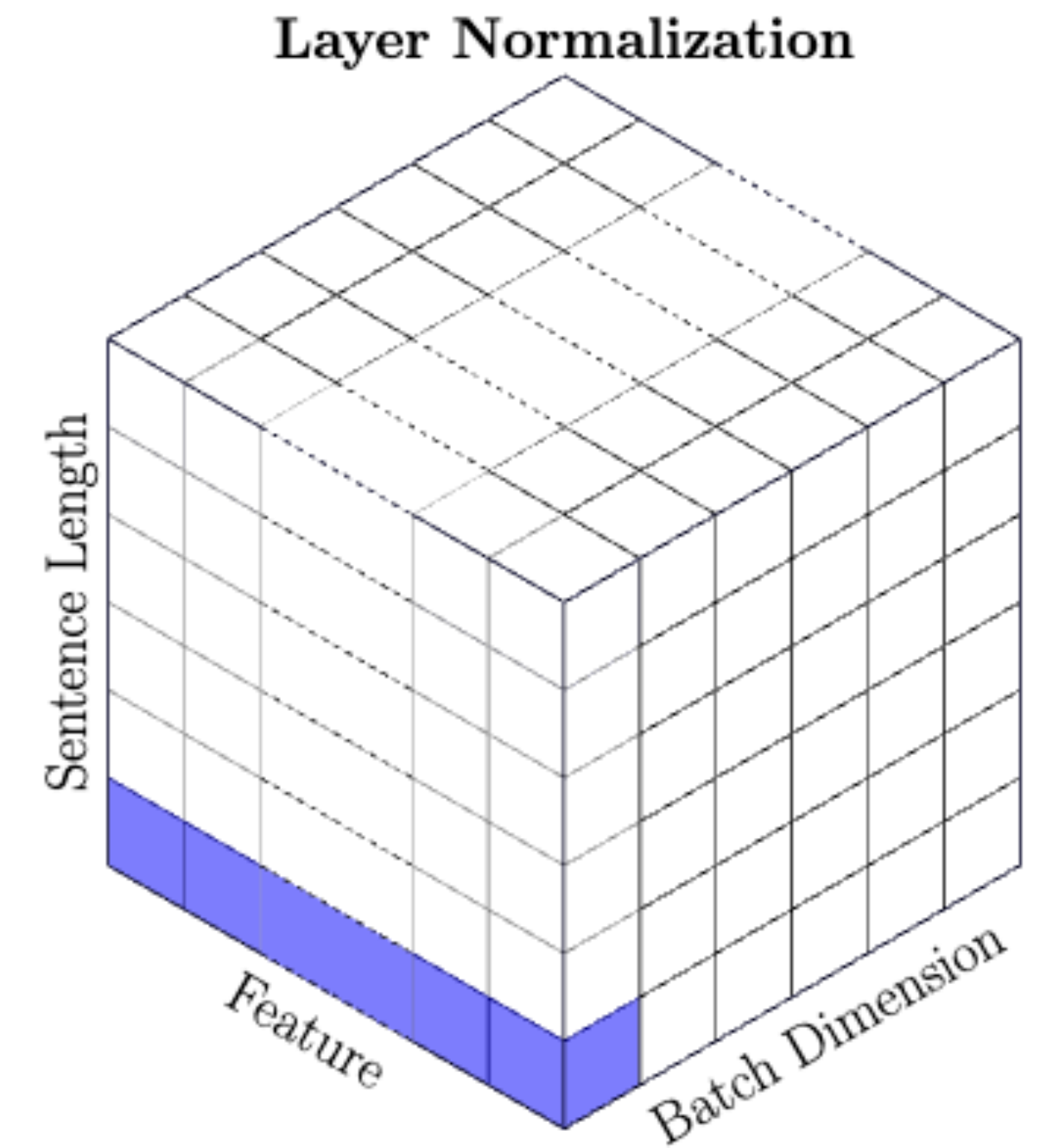
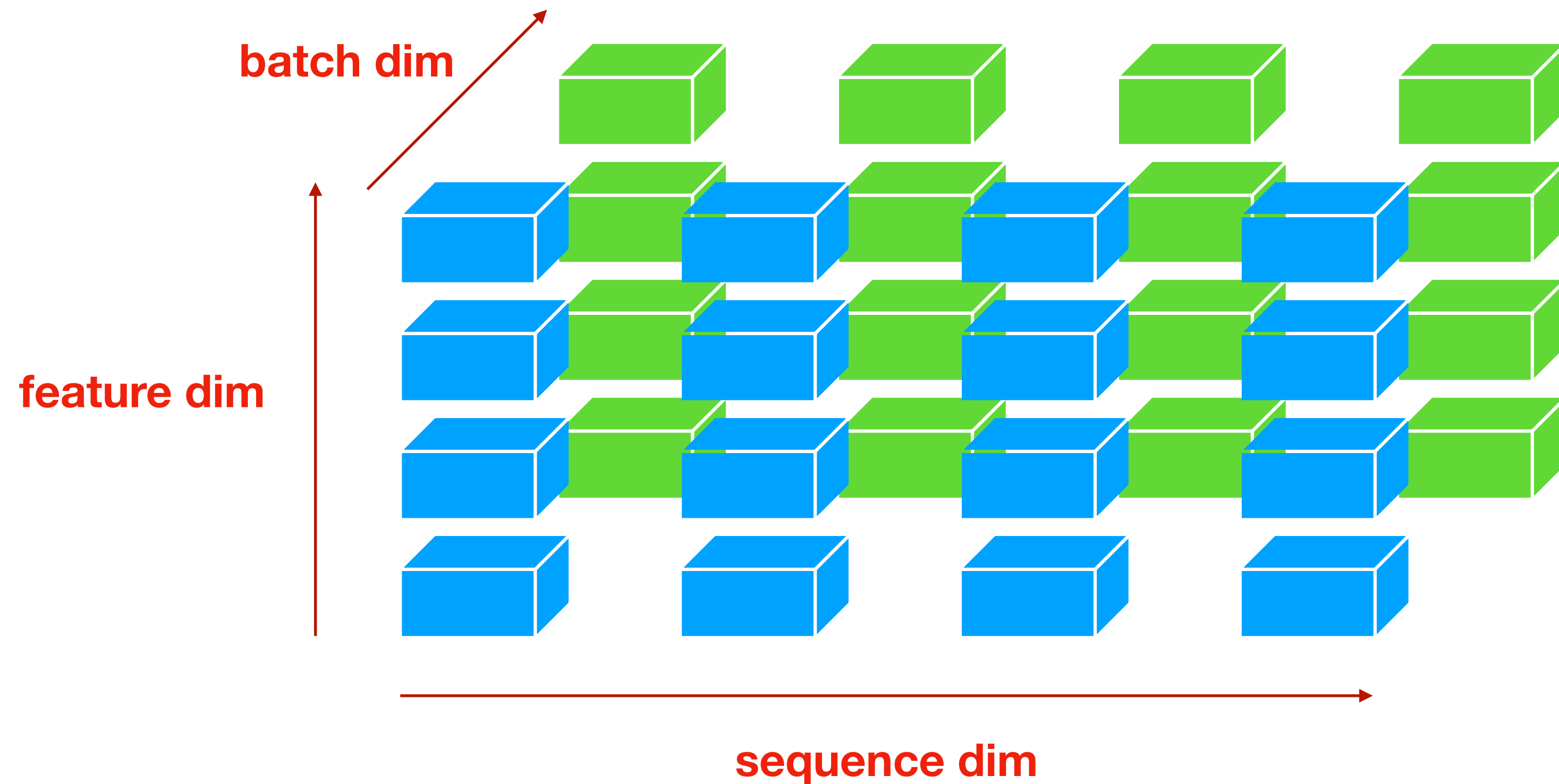
$$\mu \leftarrow \alpha \mu + (1 - \alpha) \mu_{\mathcal{B}}$$
$$\sigma^2 \leftarrow \alpha \sigma^2 + (1 - \alpha) \sigma_{\mathcal{B}}^2$$

Brittle for data from different domains

Layer Normalization

- Normalize along the **feature** dimension

$$Y = \frac{X - E[X]}{\sqrt{\text{Var}[X] + \epsilon}} * \gamma + \beta$$



Layer Normalization

- **Pros:**

- Does not involve the batch dimension (no need for running stats)
- Does not involve sequential dimension
 - Easy to parallel
 - Working well for data w. various lengths (language)

- **Cons:**

- Cannot reduce internal covariate shift/bias from data
 - Bias of data usually raises across different instances or steps
- Root Mean Square (RMS) Normalization
 - Used in LLama

$$Y = \frac{X}{\sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2 + \epsilon}} * \gamma$$

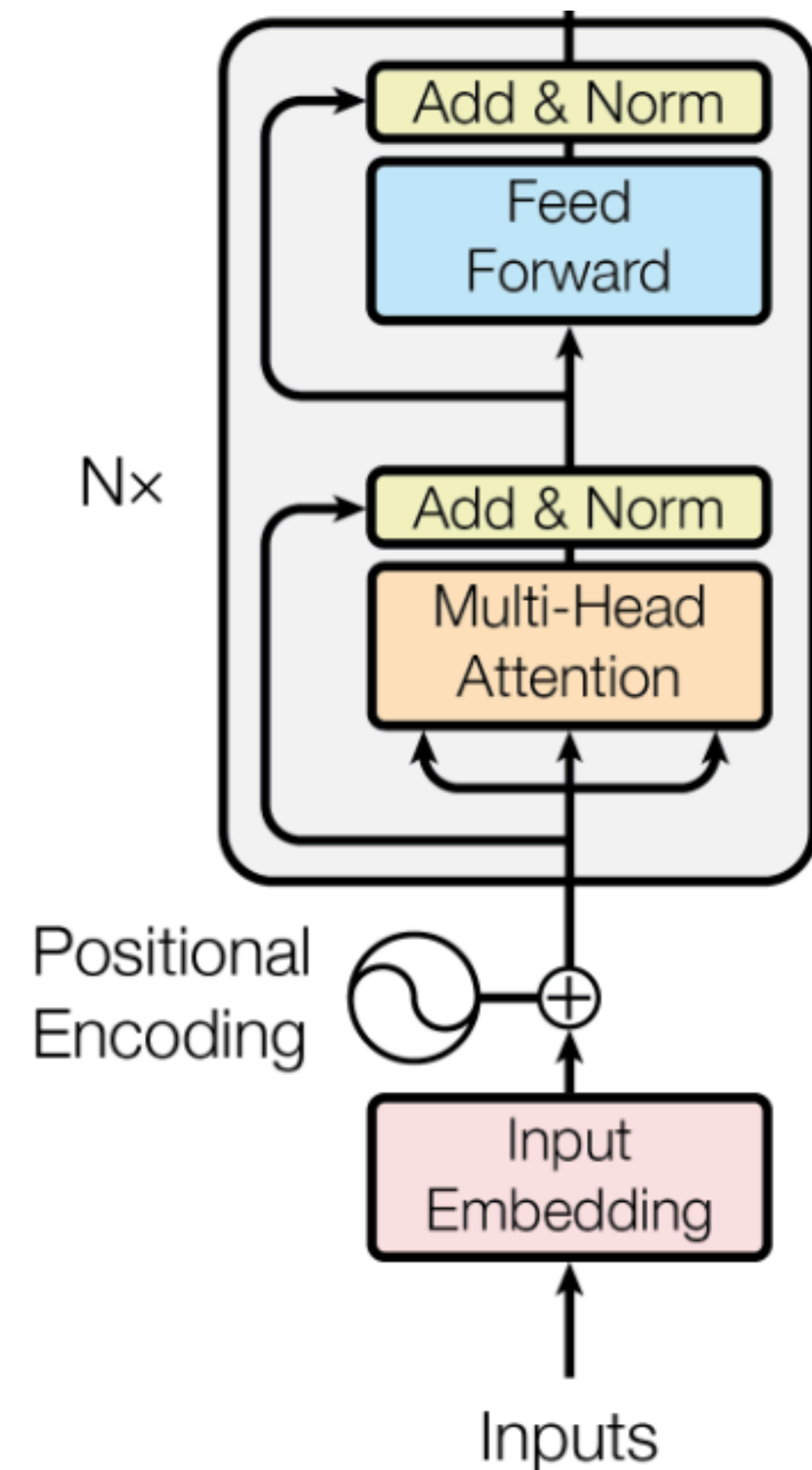
Outline

- **Residual Connections**
- **Normalization Layers**
 - Variants of Normalization Layers
 - Pre-Norm vs. Post-Norm
- **Positional Embeddings**
 - Absolute Positional Embeddings
 - Rotary Positional Embeddings

Pre-Norm vs. Post-Norm

- **Position of Normalizations**
 - Normalization layers block gradients in practice
 - The original Transformer architecture is hard to train w. more layers, even w. residual connections
 - Post-Norm -> Pre-Norm

Post-Normalization



$$Y = \text{LayerNorm} (\text{FFN}(X) + X)$$

$$Y = \text{LayerNorm} (\text{SelfAttn}(X) + X)$$

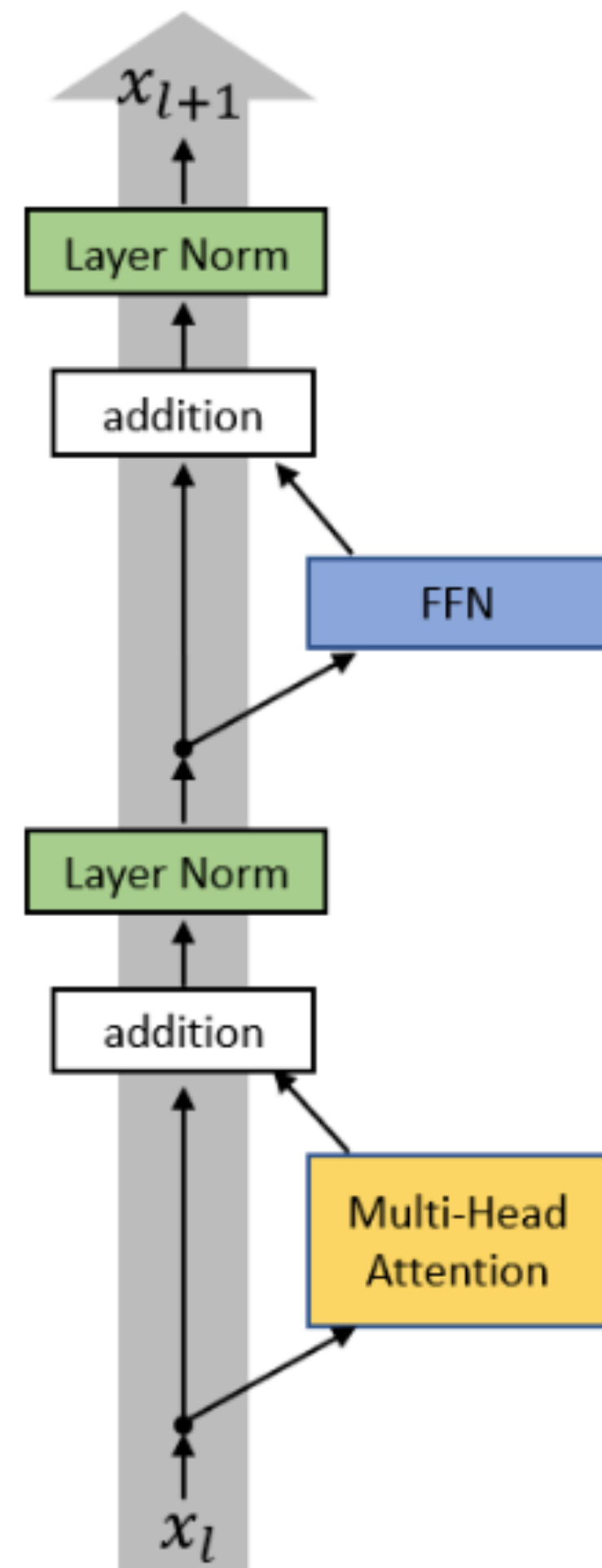
Detailed analysis

$$Y_1 = \text{LayerNorm} (f_1(X) + X)$$

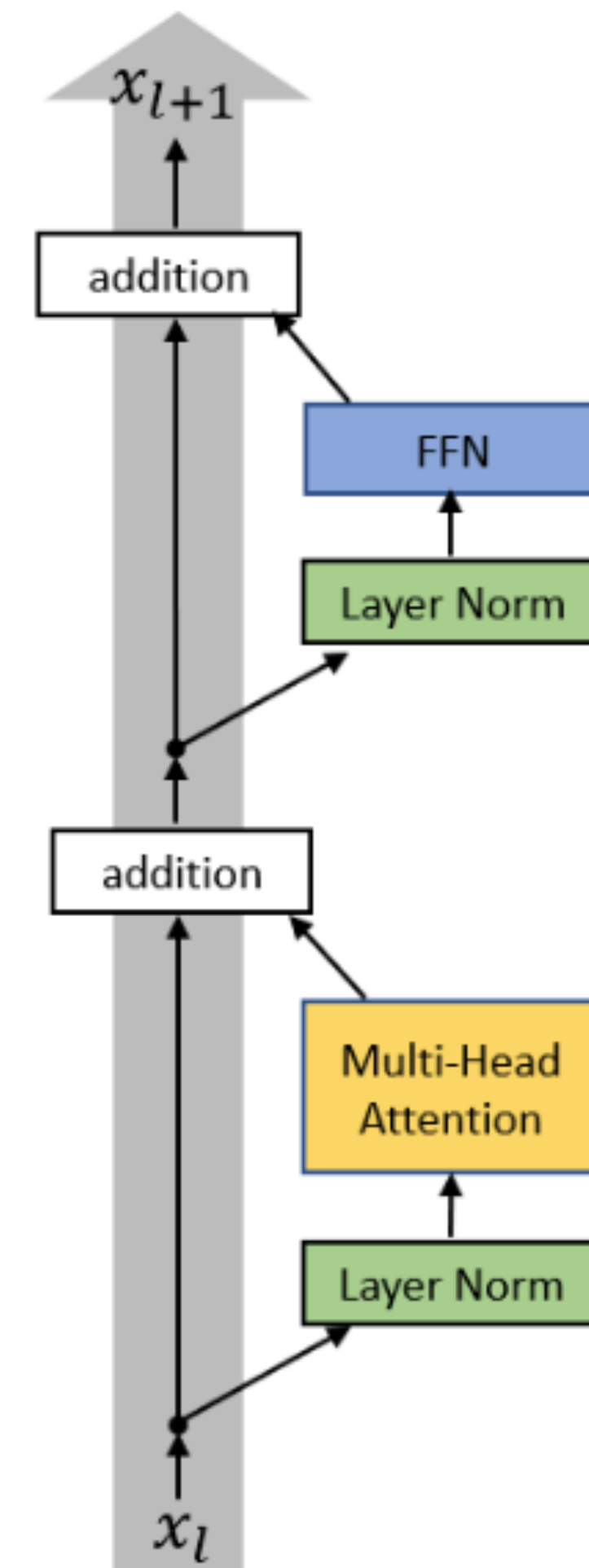
$$Y_2 = \text{LayerNorm} (f_2(Y_1) + Y_1) = \text{LayerNorm} (f_2(Y_1) + \text{LayerNorm} (f_1(X) + X))$$

$$Y_l = \text{LayerNorm} (f_l(Y_{l-1}) + Y_{l-1}) \quad l \text{ layer-normalization blocks}$$

Pre-Norm vs. Post-Norm

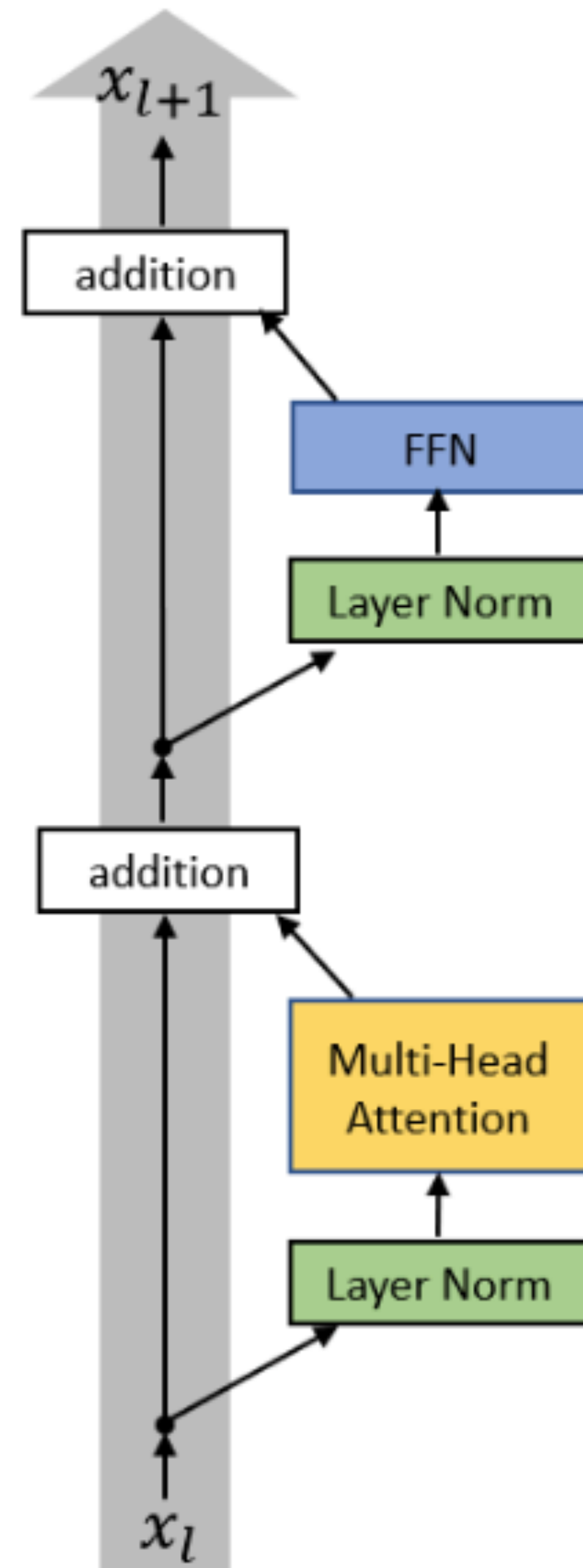


Post-Norm



Pre-Norm

Pre-Normalization



Pre-Norm

$$Y = \text{FFN} (\text{LayerNorm}(X)) + X$$

$$Y = \text{SelfAttn} (\text{LayerNorm}(X)) + X$$

Detailed analysis

$$Y_1 = f_1 (\text{LayerNorm}(X)) + X$$

$$Y_2 = f_2 (\text{LayerNorm}(Y_1)) + Y_1 = f_2 (\text{LayerNorm}(Y_1)) + f_1 (\text{LayerNorm}(X)) + X$$

Pre-Normalization

- **Pros:**

- Keeping (almost) all good properties of post-norm
- Similar performance w. Post-norm
- Able to train very deep Transformer models

- **Cons:**

- Numerically unstable (training spikes)

What is the range of Y_2 in post-norm and pre-norm?

Post-norm $Y_2 = \mathbf{LayerNorm} (f_2(Y_1) + Y_1) = \mathbf{LayerNorm} \left(f_2(Y_1) + \mathbf{LayerNorm} (f_1(X) + X) \right)$

Pre-norm $Y_2 = f_2 (\mathbf{LayerNorm}(Y_1)) + Y_1 = f_2 (\mathbf{LayerNorm}(Y_1)) + f_1 (\mathbf{LayerNorm}(X)) + X$

Outline

- **Residual Connections**
- **Normalization Layers**
 - Variants of Normalization Layers
 - Pre-Norm vs. Post-Norm
- **Positional Embeddings**
 - Absolute Positional Embeddings
 - Rotary Positional Embeddings

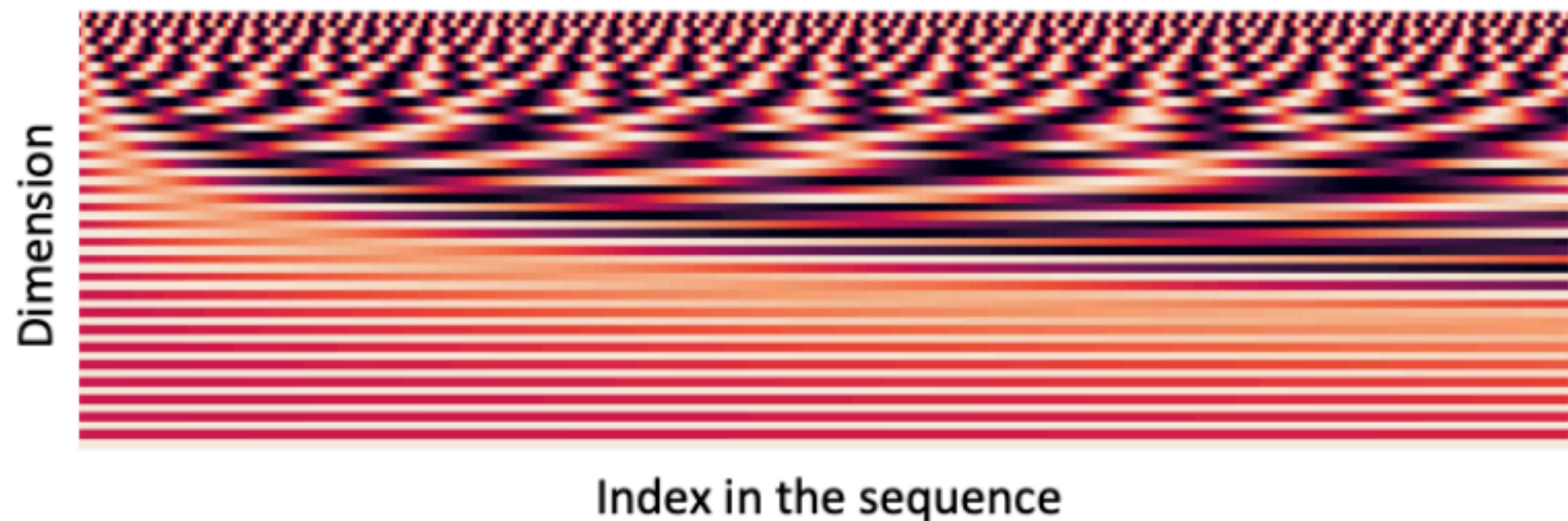
Positional Information

- Unlike RNNs, self-attention does **not** build in order information
 - Encode the order of the sentence into the input x_1, \dots, x_n
- Solution: add **positional encoding** to the input embeddings

$$x_i \leftarrow x_i + p_i$$

- Use sine and cosine functions of different frequencies (not learnable)

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



Sinusoidal Positional Embedding

$$p_{t,j} = \sin \left(t \cdot \theta^{\frac{j}{2d}} \right) \quad j \% 2 = 0$$

$$p_{t,j} = \cos \left(t \cdot \theta^{\frac{j-1}{2d}} \right) \quad j \% 2 = 1$$

Position $t = 0, 1, \dots, n$

$$\theta = \frac{1}{10000}$$

dim $j = 0, \dots, d$

	t=0	t=1	t=2	...	t=n
j=0	$\sin(0 \cdot \theta^0)$	$\sin(1 \cdot \theta^0)$	$\sin(2 \cdot \theta^0)$		$\sin(n \cdot \theta^0)$
j=1	$\cos(0 \cdot \theta^0)$	$\cos(1 \cdot \theta^0)$	$\cos(2 \cdot \theta^0)$		$\cos(n \cdot \theta^0)$
j=2	$\sin(0 \cdot \theta^{\frac{1}{2d}})$	$\sin(1 \cdot \theta^{\frac{1}{2d}})$	$\sin(2 \cdot \theta^{\frac{1}{2d}})$		$\sin(n \cdot \theta^{\frac{1}{2d}})$
j=3	$\cos(0 \cdot \theta^{\frac{1}{2d}})$	$\cos(1 \cdot \theta^{\frac{1}{2d}})$	$\cos(2 \cdot \theta^{\frac{1}{2d}})$		$\cos(n \cdot \theta^{\frac{1}{2d}})$

Sinusoidal Positional Embedding

$$p_{t,j} = \sin \left(t \cdot \theta^{\frac{j}{2d}} \right) \quad j \% 2 = 0$$

$$p_{t,j} = \cos \left(t \cdot \theta^{\frac{j-1}{2d}} \right) \quad j \% 2 = 1$$

Position $t = 0, 1, \dots, n$

$$\theta = \frac{1}{10000}$$

t=0

t=1

t=2

...

t=n

j=0	$\sin(0 \cdot \theta^0)$	$\sin(1 \cdot \theta^0)$	$\sin(2 \cdot \theta^0)$	$\sin(n \cdot \theta^0)$
j=1	$\cos(0 \cdot \theta^0)$	$\cos(1 \cdot \theta^0)$	$\cos(2 \cdot \theta^0)$	$\cos(n \cdot \theta^0)$

dim $j = 0, \dots, d$

j=2	$\sin(0 \cdot \theta^{\frac{1}{2d}})$	$\sin(1 \cdot \theta^{\frac{1}{2d}})$	$\sin(2 \cdot \theta^{\frac{1}{2d}})$	$\sin(n \cdot \theta^{\frac{1}{2d}})$
j=3	$\cos(0 \cdot \theta^{\frac{1}{2d}})$	$\cos(1 \cdot \theta^{\frac{1}{2d}})$	$\cos(2 \cdot \theta^{\frac{1}{2d}})$	$\cos(n \cdot \theta^{\frac{1}{2d}})$

Sinusoidal Positional Embedding

$$p_{t,j} = \sin \left(t \cdot \theta^{\frac{j}{2d}} \right) \quad j \% 2 = 0$$

$$p_{t,j} = \cos \left(t \cdot \theta^{\frac{j-1}{2d}} \right) \quad j \% 2 = 1$$

Position $t = 0, 1, \dots, n$

$$\theta = \frac{1}{10000}$$

	t=0	t=1	t=2	...	t=n
j=0	$\sin(0 \cdot \theta^0)$	$\sin(1 \cdot \theta^0)$	$\sin(2 \cdot \theta^0)$		$\sin(n \cdot \theta^0)$
j=1	$\cos(0 \cdot \theta^0)$	$\cos(1 \cdot \theta^0)$	$\cos(2 \cdot \theta^0)$		$\cos(n \cdot \theta^0)$
j=2	$\sin(0 \cdot \theta^{\frac{1}{2d}})$	$\sin(1 \cdot \theta^{\frac{1}{2d}})$	$\sin(2 \cdot \theta^{\frac{1}{2d}})$		$\sin(n \cdot \theta^{\frac{1}{2d}})$
j=3	$\cos(0 \cdot \theta^{\frac{1}{2d}})$	$\cos(1 \cdot \theta^{\frac{1}{2d}})$	$\cos(2 \cdot \theta^{\frac{1}{2d}})$		$\cos(n \cdot \theta^{\frac{1}{2d}})$

dim $j = 0, \dots, d$

Sinusoidal Positional Embedding

$$p_{t,j} = \sin \left(t \cdot \theta^{\frac{j}{2d}} \right) \quad j \% 2 = 0$$

$$p_{t,j} = \cos \left(t \cdot \theta^{\frac{j-1}{2d}} \right) \quad j \% 2 = 1$$

Position $t = 0, 1, \dots, n$

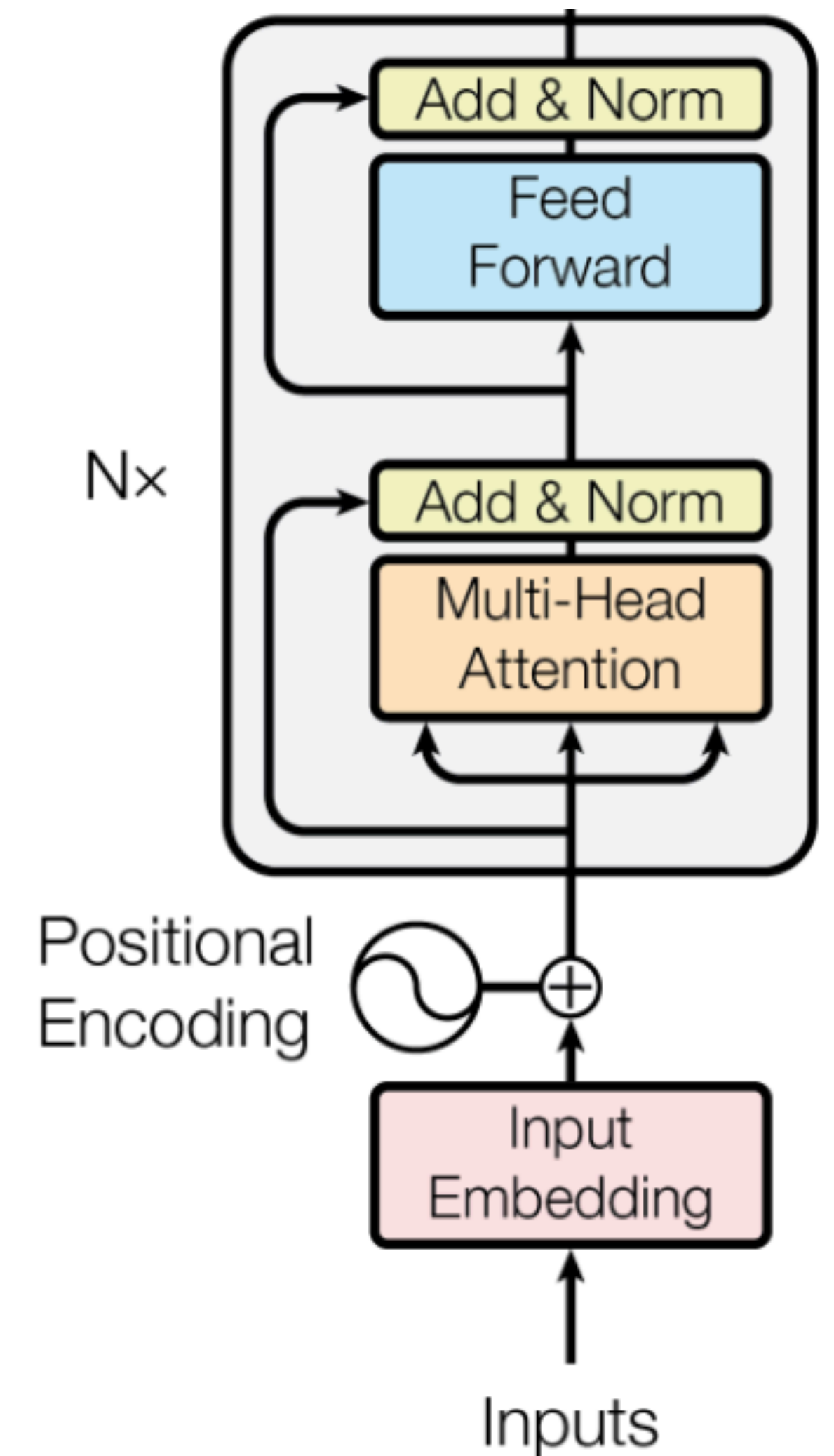
$$\theta = \frac{1}{10000}$$

dim $j = 0, \dots, d$

	t=0	t=1	t=2	...	t=n
j=0	$\sin(0 \cdot \theta^0)$	$\sin(1 \cdot \theta^0)$	$\sin(2 \cdot \theta^0)$		$\sin(n \cdot \theta^0)$
j=1	$\cos(0 \cdot \theta^0)$	$\cos(1 \cdot \theta^0)$	$\cos(2 \cdot \theta^0)$		$\cos(n \cdot \theta^0)$
j=2	$\sin(0 \cdot \theta^{\frac{1}{2d}})$	$\sin(1 \cdot \theta^{\frac{1}{2d}})$	$\sin(2 \cdot \theta^{\frac{1}{2d}})$		$\sin(n \cdot \theta^{\frac{1}{2d}})$
j=3	$\cos(0 \cdot \theta^{\frac{1}{2d}})$	$\cos(1 \cdot \theta^{\frac{1}{2d}})$	$\cos(2 \cdot \theta^{\frac{1}{2d}})$		$\cos(n \cdot \theta^{\frac{1}{2d}})$

Sinusoidal Positional Embedding

- Sinusoidal positional embedding depends on absolute positional information
- Only added to the input word embeddings
- Ability for length extension
 - Training on up to k words
 - Test on $> k$ words



Rotary Positional Embedding

- **Absolute positional index is not meaningful**
 - I don't know
 - In fact, I don't know
- **How about relative positional embedding?**
 - Positional information should be invariant w.r.t relative distance of two positions
 - Rotary Position Embedding
 - Used in [Llama models](#)

Rotary Positional Embedding

Rotary Operation

$$X \in \mathbb{R}^d \quad \theta_j = \left(\frac{1}{10000} \right)^{\frac{j}{2d}}$$

$$\mathbf{Rotary}(X, m) = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_d \\ x_{d-1} \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{pmatrix}$$

Rotary Positional Embedding

Attention

$$\text{attn}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$

Applying rotary operation to each query and key vector

$$q'_i = \mathbf{Rotary}(q_i, i)$$

$$k'_j = \mathbf{Rotary}(k_j, j)$$

Replacing Q, K with Q', K' to compute attention

Rotary Positional Embedding

- A simple case: $d = 2$

$$\mathbf{Rotary}(X, m) = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_d \\ x_{d-1} \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{pmatrix}$$

$$q'_i = \mathbf{Rotary}(q_i, i) = \begin{pmatrix} q_{i,1} \cos(i\theta) - q_{i,2} \sin(i\theta) \\ q_{i,1} \sin(i\theta) + q_{i,2} \cos(i\theta) \end{pmatrix}$$

$$k'_j = \mathbf{Rotary}(k_j, j) = \begin{pmatrix} k_{j,1} \cos(j\theta) - k_{j,2} \sin(j\theta) \\ k_{j,1} \sin(j\theta) + k_{j,2} \cos(j\theta) \end{pmatrix}$$

$\langle q'_i, k'_j \rangle$ only depends on $(i - j)\theta$

Q&A