# Natural Language Processing

# 第八周 XFormer

庞彦

yanpang@gzhu.edu.cn

# Overview

**CONTENTS**

Spring 2023

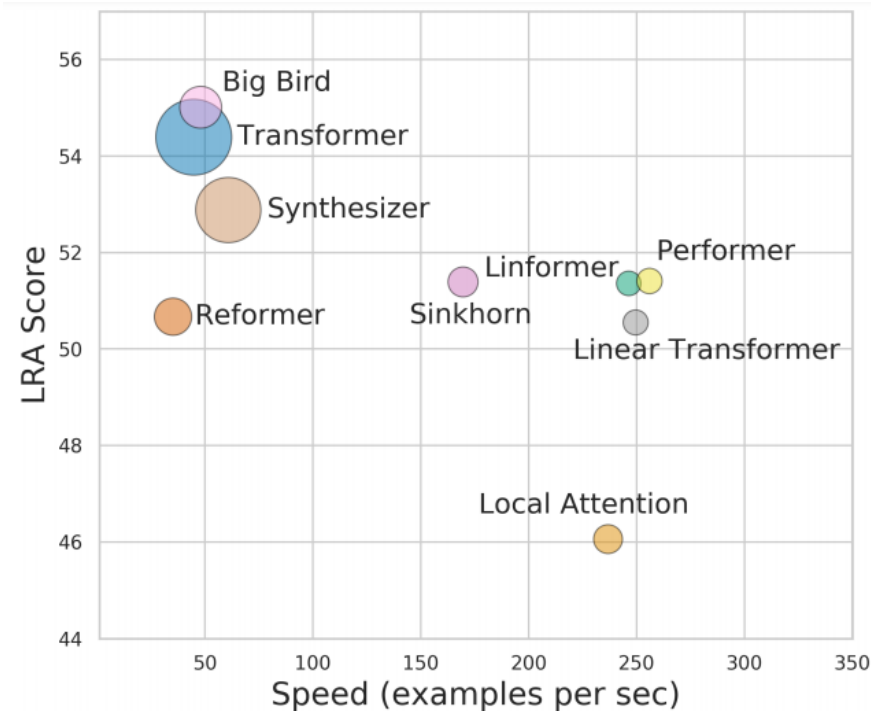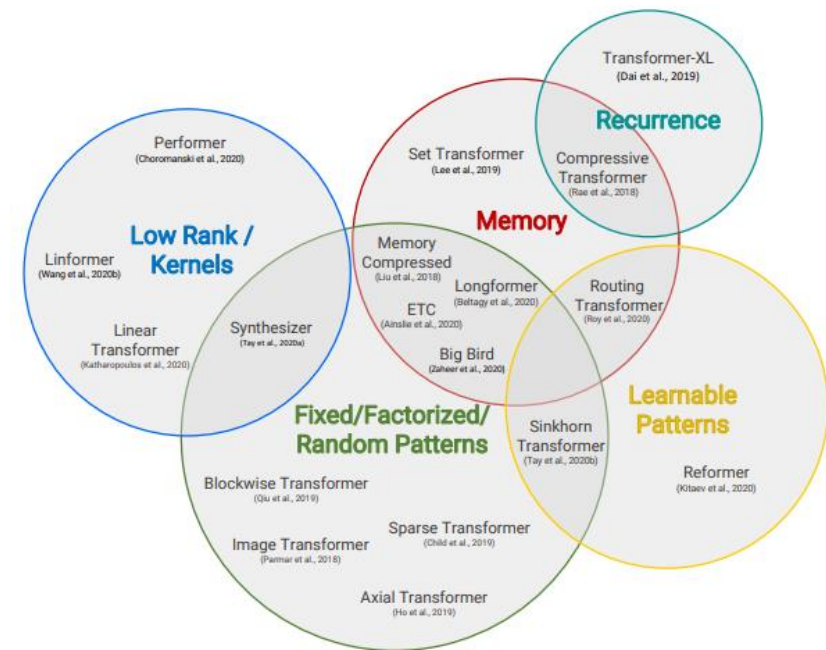**01** | More Attention Mechanisms

更多的自注意力机制

# To Learn More ...
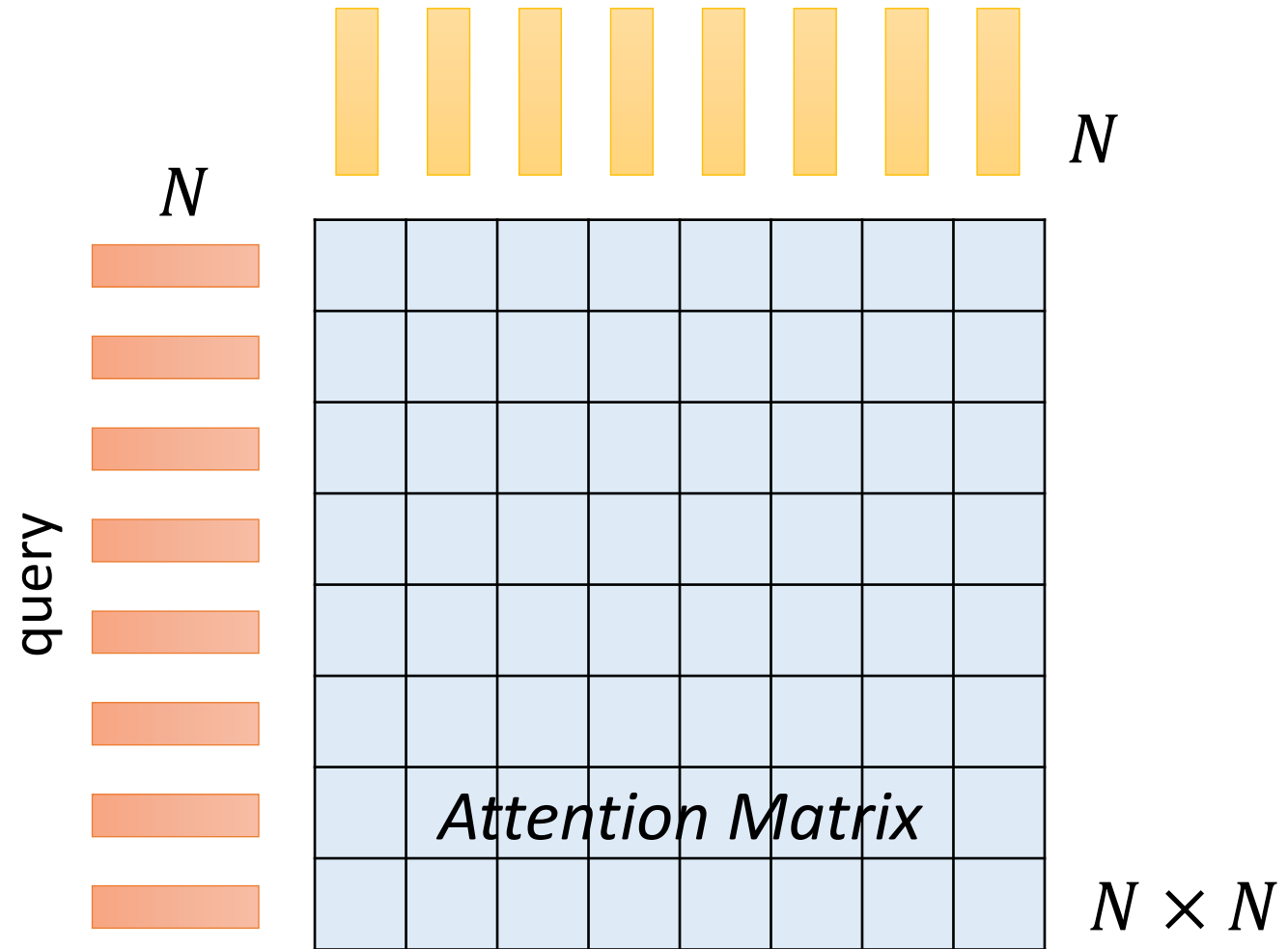


Long Range Arena: A Benchmark for Efficient Transformers

Efficient Transformers: A Survey
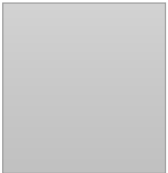
# How to make self-attention efficient?

Sequence length = $N$
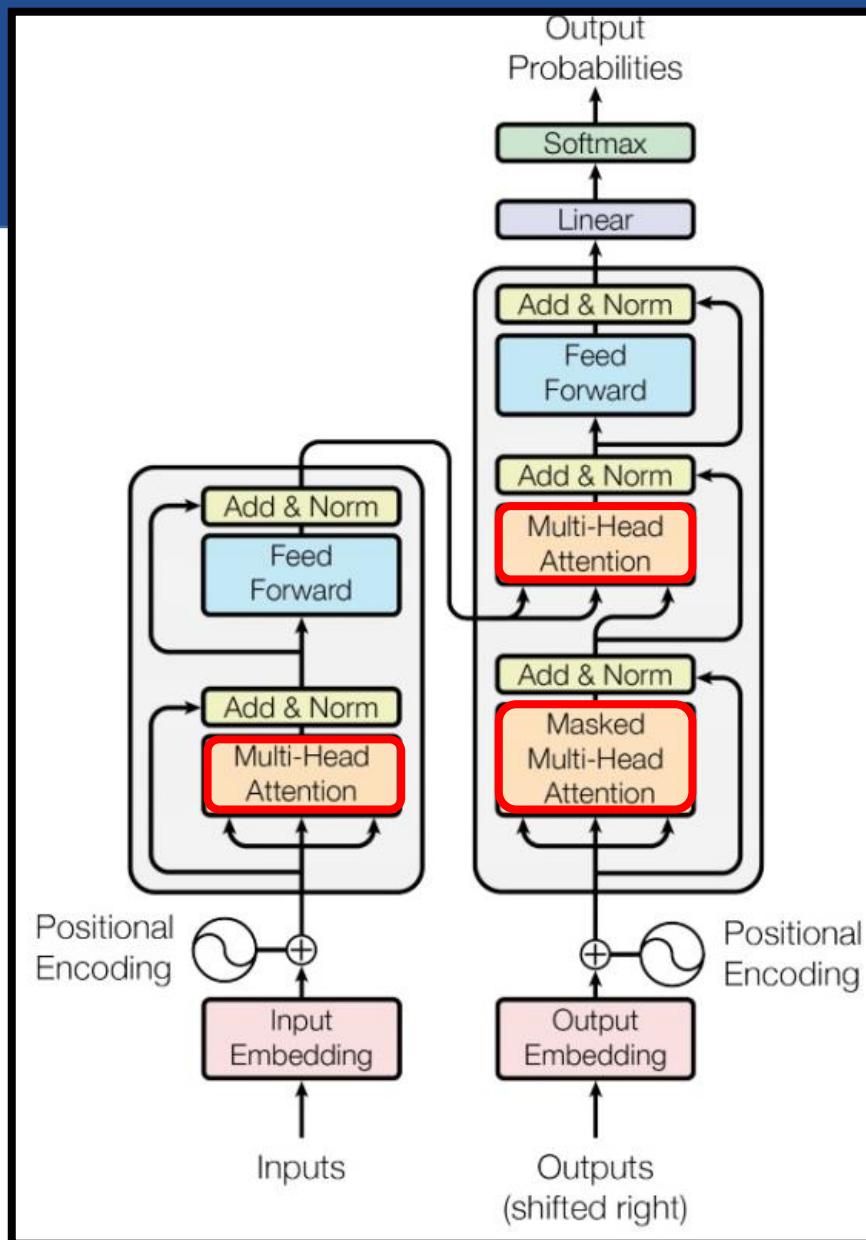
$N$

$N$

query

Attention Matrix

$N \times N$

# Notice

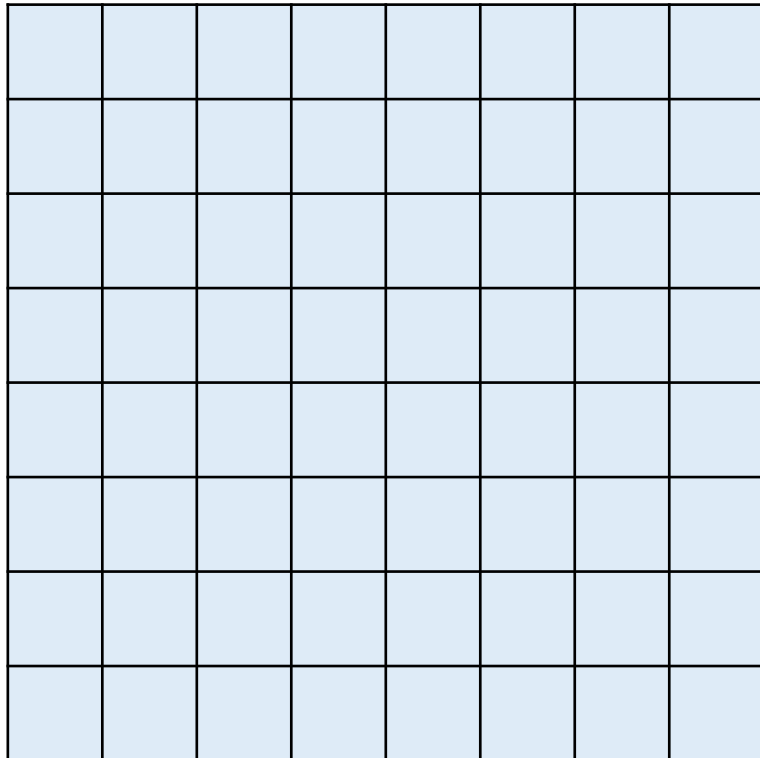- Self-attention is only a module in a larger network. 自注意力机制仅仅知识大网络的一个模块。

- Self-attention dominates computation when $N$ is large. 当$N$非常大的时候，自注意力机制的计算量显著上升。

- Usually developed for image processing

256 ▢ 256

$N = 256 * 256$

Can we fill in some values
with human knowledge?

# Local Attention / Truncated Attention

Set to 0

Calculate attention weight
计算权重

Similar with CNN
近似CNN

key

query

# Stride Attention

# Global Attention

Add special token into original sequence在原句中新增特殊的token

- Attend to every token → collect global information搜集全局信息
- Attended by every token → it knows global information



*No attention between non-special token*

# Many Different Choices …



**Different heads use different patterns.**
不同的头利用不用的特征。

# Many Different Choices ...

- LongFormer
  https://arxiv.org/abs/2004.05150



(b) Sliding window attention    (c) Dilated sliding window    (d) Global+sliding window
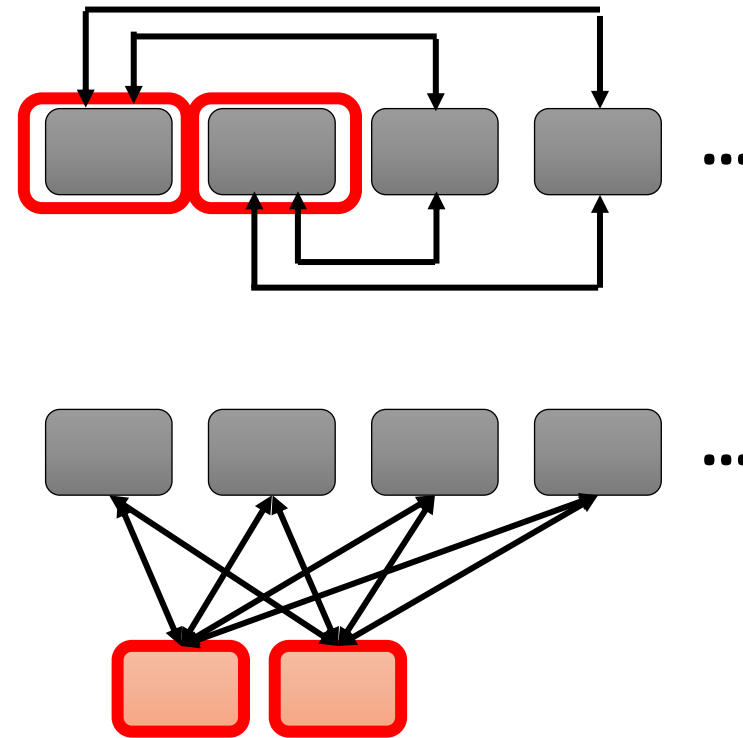
- Big Bird
  https://arxiv.org/abs/2007.14062



(a) Random attention    (b) Window attention    (c) Global Attention    (d) BIGBIRD

# Can we only focus on Critical Parts?



small value小值

- Directly set to 0直接归零
- Smaller influence on results 对结果影响较小

large value大值

**How to quickly estimate the portion with small attention weights?**

# Clustering

***Step 1***

query                                        key

Clustering
based on similarity        (approximate & fast)

1   4   1   2   3   3   3   2   2   1   3   3   1   4   1   4

Routing Transformer
https://arxiv.org/abs/2003.05997

# Clustering

key

**Step 2**

query

Belong to the same cluster, then calculate attention weight

Not the same cluster, set to 0

# Learnable Patterns

Input sequence

Sinkhorn Sorting Network

key

NN ← **Jointly learned**

query



A grid should be skipped or not is decided by another learned module

# Do we need full attention matrix?

Many redundant columns



query

key          ***Low Rank***

Linformer
https://arxiv.org/abs/2006.04768

# Clustering



$N$ value

$N$ key

$K$

$K$ Representative keys

output

Can we reduce the number of queries? 可以减少Q的数量吗？

query

change output sequence length 改变序列长度

# Reduce Number of Keys

## Compressed Attention
https://arxiv.org/abs/1801.10198

## Linformer
https://arxiv.org/abs/2006.04768



Linear combination of N vectors N个矢量的线性组合

# Attention Mechanism is three-matrix Multiplication

*Review*

$$d \times N \quad \boxed{Q} \quad = \quad \boxed{W^q} \quad \boxed{I}$$

$$d \times N \quad \boxed{K} \quad = \quad \boxed{W^k} \quad \boxed{I}$$

$$d' \times N \quad \boxed{V} \quad = \quad \boxed{W^v} \quad \boxed{I}$$

$$A' \quad \longleftarrow \quad A \quad = \quad \boxed{K^T} \quad \boxed{Q}$$

softmax
ignore

$$A'$$

$$\boxed{O} \quad = \quad \boxed{V}$$

# Attention Mechanism is three-matrix Multiplication

*Review*

$$d \times N \quad \boxed{Q} = \boxed{W^q} \; \boxed{I}$$

$$d \times N \quad \boxed{K} = \boxed{W^k} \; \boxed{I}$$

$$d' \times N \quad \boxed{V} = \boxed{W^v} \; \boxed{I}$$

$$\boxed{O} \approx \boxed{V} \; \boxed{K^T} \; \boxed{Q}$$

$$d' \times N \qquad\quad d' \times N \qquad\quad d \times N$$

$$N \times d$$

# Attention Mechanism is three-matrix Multiplication

$$O \approx V \quad K^T \quad Q$$

$d' \times N \qquad d' \times N \qquad \qquad d \times N$

$N \times d$

$$O \approx V \quad K^T \quad Q$$

$d' \times N \qquad d' \times N \qquad \qquad d \times N$

$N \times d$

What is the difference?
区别?

# Review Linear Algebra

## Practical Issue

$k=1$     $m=1000$

$n=1$     $p=1000$

- Let A and B be k x m matrices, C be an m x n matrix, and P and Q be n x p matrices
  - A(CP) = (AC)P

| A | C | P |
|---|---|---|
| k X m | m X n | n X p |

| A | CP |
|---|---|
| k X m | m X p |

m X n X p

$10^6$

k X m X p

$10^6$

| A | C | P |
|---|---|---|
| k X m | m X n | n X p |

| AC | P |
|---|---|
| k X n | n X p |

k X m X n

$10^3$

k X n X p

$10^3$

# Review Linear Algebra

$$O \approx V \; K^T \; Q$$

$$(d + d')N^2$$

$d' \times N$  $d' \times N$  $d \times N$

$N \times d$

$N \times d \times N$

$$V \quad A$$

Attention Matrix

$d' \times N$  $N \times N$

$d' \times N \times N$

# Review Linear Algebra

$$O \approx V \quad K^T \quad Q$$

$d' \times N$     $d' \times N$     $d \times N$     $\boxed{2d'dN}$

$N \times d$

$d' \times N \times d$

$d \times N$

$d' \times d$   ?   Q

$d' \times d \times N$

# Review Linear Algebra

$$O \approx V \quad K^T \quad Q$$

$$d \times N \qquad d \times N \qquad d \times N$$

$$N \times d$$

$$(d + d')N^2$$

$$\vee$$

$$O \approx V \quad K^T \quad Q$$

$$d \times N \qquad d \times N \qquad d \times N$$

$$N \times d$$

$$2d'dN$$

# Review Linear Algebra

Let's put softmax back …
把Softmax激活函数放回来

# Softmax

$$b^1 = \sum_{i=1}^{N} \boxed{\alpha'_{1,i}} v^i = \sum_{i=1}^{N} \boxed{\frac{exp(q^1 \cdot k^i)}{\sum_{j=1}^{N} exp(q^1 \cdot k^j)}} v^i$$

# Softmax

$$b^1 = \sum_{i=1}^{N} \alpha'_{1,i} v^i = \sum_{i=1}^{N} \frac{exp(q^1 \cdot k^i)}{\sum_{j=1}^{N} exp(q^1 \cdot k^j)} v^i$$

$$exp(q \cdot k) \approx \phi(q) \cdot \phi(k)$$

$$q \rightarrow \boxed{\phi} \rightarrow \phi(q)$$

$$= \sum_{i=1}^{N} \frac{\phi(q^1) \cdot \phi(k^i)}{\sum_{j=1}^{N} \phi(q^1) \cdot \phi(k^j)} v^i$$

$$= \frac{\sum_{i=1}^{N} [\phi(q^1) \cdot \phi(k^i)] v^i}{\sum_{j=1}^{N} \phi(q^1) \cdot \phi(k^j)}$$

$$\phi(q^1) \cdot \sum_{j=1}^{N} \phi(k^j) \quad\cdots\cdots\rightarrow\quad \phi(q^1)$$

$$b^1 = \sum_{i=1}^{N} \alpha'_{1,i} v^i = \frac{\sum_{i=1}^{N} [\phi(q^1) \cdot \phi(k^i)] v^i}{\phi(q^1) \cdot \sum_{j=1}^{N} \phi(k^j)}$$

$$\boxed{\sum_{i=1}^{N} [\phi(q^1) \cdot \phi(k^i)] v^i} \qquad \phi(q^1) = \begin{bmatrix} q_1^1 \\ q_2^1 \\ \vdots \end{bmatrix} \qquad \phi(k^1) = \begin{bmatrix} k_1^1 \\ k_2^1 \\ \vdots \end{bmatrix}$$

$$= [\phi(q^1) \cdot \phi(k^1)] v^1 + [\phi(q^1) \cdot \phi(k^2)] v^2 + \cdots$$

$$= (q_1^1 k_1^1 + q_2^1 k_2^1 + \cdots) v^1 + (q_1^1 k_1^2 + q_2^1 k_2^2 + \cdots) v^2 + \cdots$$

$$= \underline{q_1^1 k_1^1 v^1} + \underline{q_2^1 k_2^1 v^1} + \cdots + \underline{q_1^1 k_1^2 v^2} + \underline{q_2^1 k_2^2 v^2} + \cdots + \cdots$$

$$= q_1^1 (k_1^1 v^1 + k_1^2 v^2 + \cdots) + q_2^1 (k_2^1 v^1 + k_2^2 v^2 + \cdots)$$

# Softmax

$$\boldsymbol{b^1} = \sum_{i=1}^{N} \alpha'_{1,i} \boldsymbol{v^i} = \frac{\sum_{i=1}^{N}[\phi(\boldsymbol{q^1}) \cdot \phi(\boldsymbol{k^i})]\boldsymbol{v^i}}{\phi(\boldsymbol{q^1}) \cdot \sum_{j=1}^{N}\phi(\boldsymbol{k^j})}$$

$$\sum_{i=1}^{N}[\phi(\boldsymbol{q^1}) \cdot \phi(\boldsymbol{k^i})]\boldsymbol{v^i}$$

$$\phi(\boldsymbol{q^1}) = \begin{bmatrix} q_1^1 \\ q_2^1 \\ \vdots \end{bmatrix} \qquad \phi(\boldsymbol{k^1}) = \begin{bmatrix} k_1^1 \\ k_2^1 \\ \vdots \end{bmatrix}$$

M dim

$$= \boxed{q_1^1}(k_1^1 \boldsymbol{v^1} + k_1^2 \boldsymbol{v^2} + \cdots) + \boxed{q_2^1}(k_2^1 \boldsymbol{v^1} + k_2^2 \boldsymbol{v^2} + \cdots)$$

$$\sum_{j=1}^{N} k_1^j \boldsymbol{v^j}$$

$$\sum_{j=1}^{N} k_2^j \boldsymbol{v^j}$$

$$\phi(\boldsymbol{q^1})$$

M vectors

# Softmax

$$v^1 \quad v^2 \quad v^3 \quad \cdots \quad v^N \qquad \phi(\boldsymbol{k^1}) \quad \phi(\boldsymbol{k^2}) \quad \phi(\boldsymbol{k^3}) \qquad \cdots\cdots \qquad \phi(\boldsymbol{k^N})$$

M dim

$$\sum_{j=1}^{N} k_1^j \boldsymbol{v^j} \qquad \sum_{j=1}^{N} k_2^j \boldsymbol{v^j}$$

M vectors

$$\boldsymbol{b^1} = \frac{\phi(\boldsymbol{q^1}) \quad \text{M dim}}{\displaystyle\sum_{j=1}^{N} \phi(\boldsymbol{k^j}) \quad \phi(\boldsymbol{q^1})}$$

# Softmax

$$b^1 = \frac{\overbrace{\begin{array}{cccc} \| & \| & \| & \cdots & \| \end{array}}^{\phi(q^1)}}{\rule{6cm}{0.4pt}}$$

$\phi(q^1)$

Don't compute again

$$b^2 = \frac{\overbrace{\begin{array}{cccc} \| & \| & \| & \cdots & \| \end{array}}^{\phi(q^2)}}{\rule{6cm}{0.4pt}}$$

$\phi(q^2)$

$\phi(q^2)$

$$b^N = \ \ldots$$

# Softmax

weighted sum

$$= \bullet\, v^1 + \bullet\, v^2 + \bullet\, v^3 + \bullet\, v^4$$

$$\phi(\boldsymbol{q}^1) = \begin{bmatrix} q_1^1 \\ q_2^1 \\ \vdots \end{bmatrix}$$

$$= \bullet\, v^1 + \bullet\, v^2 + \bullet\, v^3 + \bullet\, v^4$$

M vectors

M dimensions

$\phi(\boldsymbol{k}^1)$  $\phi(\boldsymbol{k}^2)$  $\phi(\boldsymbol{k}^3)$  $\phi(\boldsymbol{k}^4)$

$\boldsymbol{q}^1$  $\boldsymbol{k}^1$  $\boldsymbol{v}^1$   $\boldsymbol{k}^2$  $\boldsymbol{v}^2$   $\boldsymbol{k}^3$  $\boldsymbol{v}^3$   $\boldsymbol{k}^4$  $\boldsymbol{v}^4$

$\boldsymbol{a}^1$   $\boldsymbol{a}^2$   $\boldsymbol{a}^3$   $\boldsymbol{a}^4$

# Softmax

weighted sum

$$\phi(\boldsymbol{q^1}) = \begin{bmatrix} q_1^1 \\ q_2^1 \\ \vdots \end{bmatrix}$$

$$\bullet = \bullet\, v^1 + \bullet\, v^2 + \bullet\, v^3 + \bullet\, v^4$$

$$\bullet = \bullet\, v^1 + \bullet\, v^2 + \bullet\, v^3 + \bullet\, v^4$$

M vectors

$$\boldsymbol{b^1} = \frac{\begin{bmatrix} \quad \cdots \quad \end{bmatrix} \; \phi(\boldsymbol{q^1})}{\displaystyle\sum_{j=1}^{N} \phi(\boldsymbol{k^j}) \; \phi(\boldsymbol{q^1})}$$

# Realization

$$exp(\boldsymbol{q} \cdot \boldsymbol{k})$$
$$\approx \phi(\boldsymbol{q}) \cdot \phi(\boldsymbol{k})$$

$$\boldsymbol{q} \rightarrow \boxed{\phi} \rightarrow \phi(\boldsymbol{q})$$

- Efficient attention
  https://arxiv.org/pdf/1812.01243.pdf

- Linear Transformer
  https://linear-transformers.com/

- Random Feature Attention
  https://arxiv.org/pdf/2103.02143.pdf

- Performer
  https://arxiv.org/pdf/2009.14794.pdf

$$b^1 = \sum_{i=1}^{N} \alpha'_{1,i} v^i$$

softmax

| $\alpha_{1,1}$ | $\alpha_{1,2}$ | $\alpha_{1,3}$ | $\alpha_{1,4}$ |
|---|---|---|---|
| $\alpha_{1,2}$ | $\alpha_{2,2}$ | $\alpha_{2,3}$ | $\alpha_{2,4}$ |
| $\alpha_{1,3}$ | $\alpha_{2,3}$ | $\alpha_{3,3}$ | $\alpha_{3,4}$ |
| $\alpha_{1,4}$ | $\alpha_{2,4}$ | $\alpha_{3,4}$ | $\alpha_{4,4}$ |

~~From q and k?~~

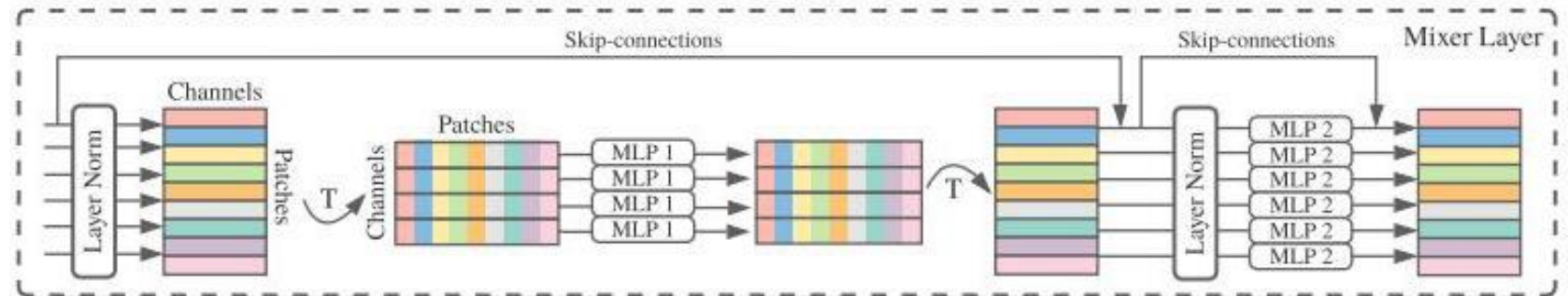They are network parameters!

# Attention-Free

- Fnet: Mixing tokens with fourier transforms

  https://arxiv.org/abs/2105.03824

- Pay Attention to MLPs  https://arxiv.org/abs/2105.08050

- MLP-Mixer: An all-MLP Architecture for Vision
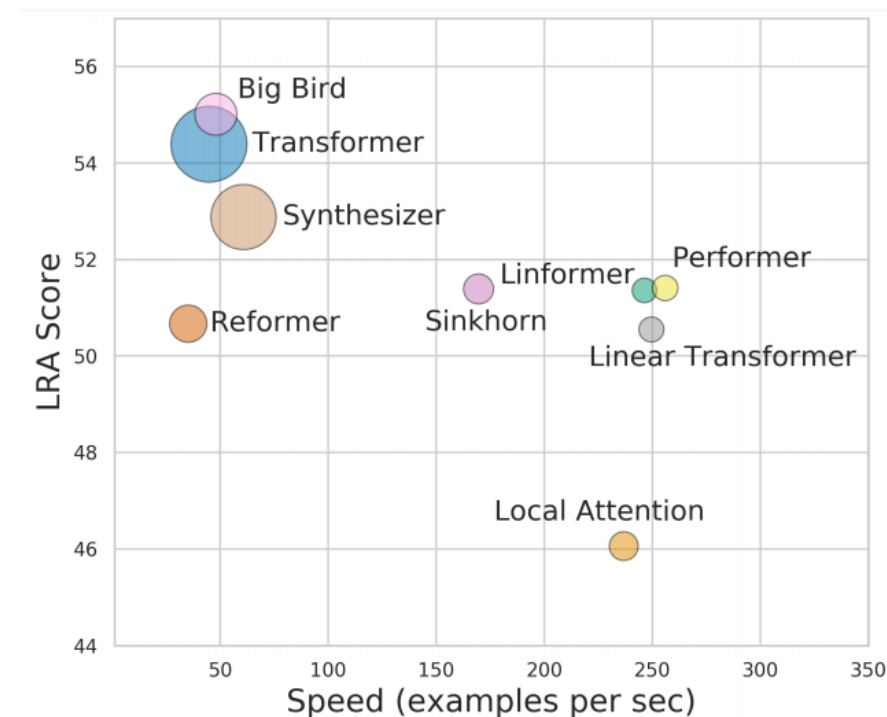
  https://arxiv.org/abs/2105.01601

# Summary

- Human knowledge
  - Local Attention, Big Bird
- Clustering
  - Reformer
- Learnable Pattern
  - SinFform
- Representative key
  - LinFormer
- k,q first → v,k first
  - Linear Transformer, Performer
- New framework
  - Synthesizer

# Q&A