



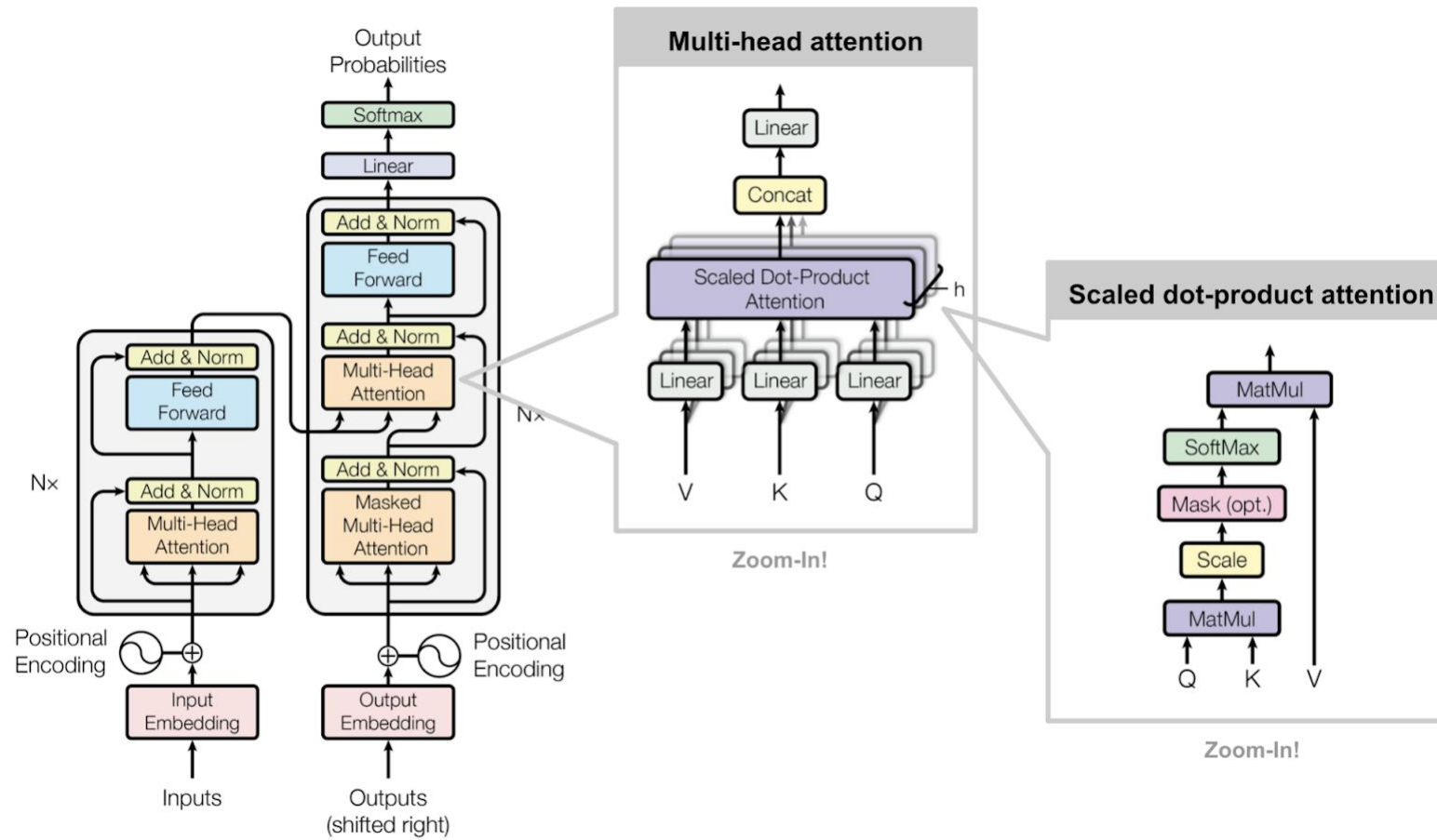
# Paper Sharing

**MEGA: MOVING AVERAGE EQUIPPED GATED ATTENTION**

**Lecturer: Yuxin Wu**

**2023.6.14**

# Transformer Architecture



# Limitations of Attention Mechanism

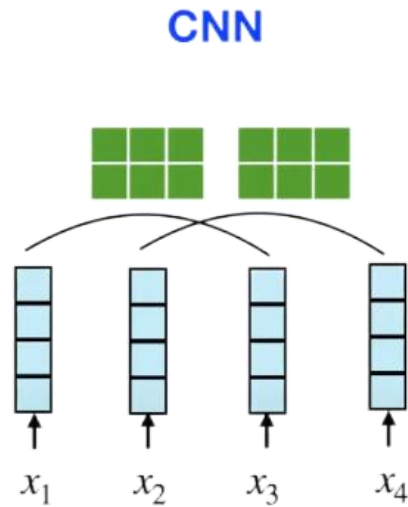
- **Weak Inductive Bias**

- Almost **no prior knowledge** of dependency patterns
  - Learning directly from data
- Position information only from **absolute/relative positional embeddings**

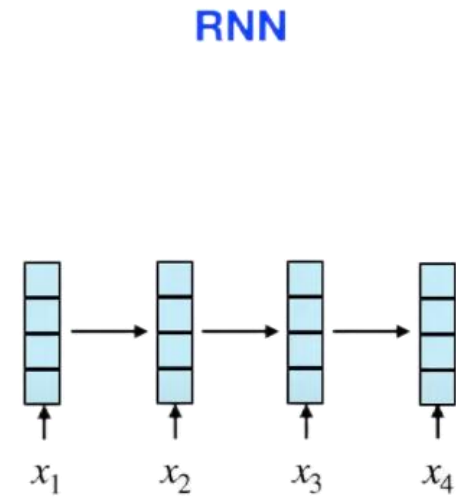
- **Quadratic Complexity**

- Both time and space
- $O(hn^2)$ :  $h$  heads and sequence length of  $n$

# Inductive Bias: CNN & RNN

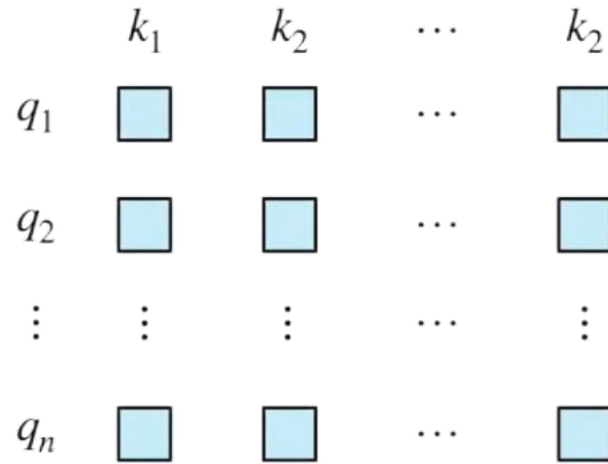


- Local dependencies
  - Window size is usually small (e.g. 3 or 5)
- Time-invariant kernel



- Sequential dependencies
- Time-invariant recurrence

# Inductive Bias: Attention

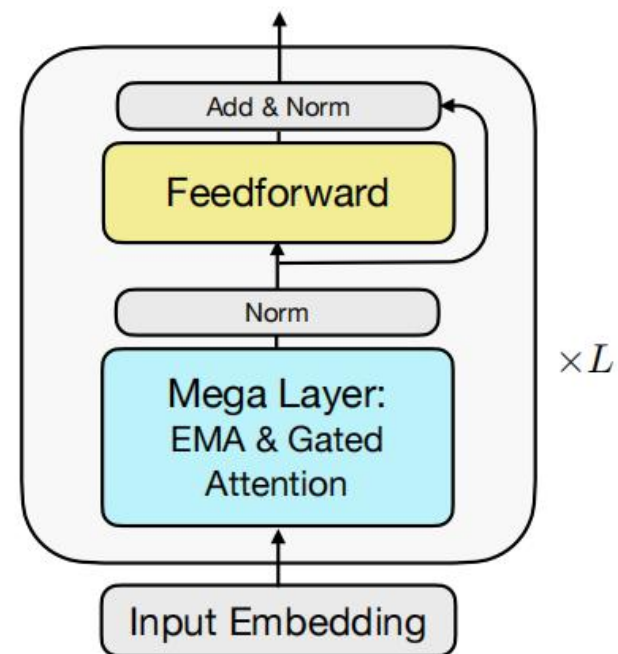


- Pair-wise interaction
- Order-invariant if no positional embeddings

Neither accurate nor efficient for long sequence modeling

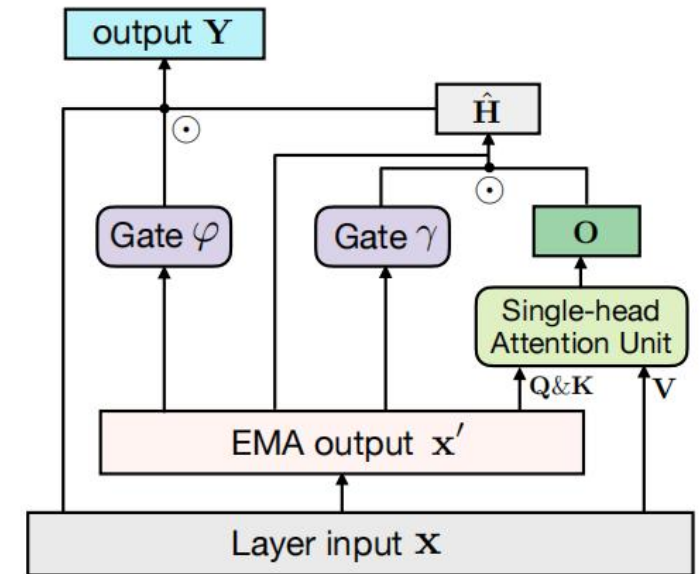
# Mega: Overview

- **Effective and efficient drop-in replacement of attention for long sequence modeling**
  - Outstanding results on various types
    - text, images and audios
  - Exponential Moving Average (EMA)
  - Mega-chunk: **linear complexity** of time and space



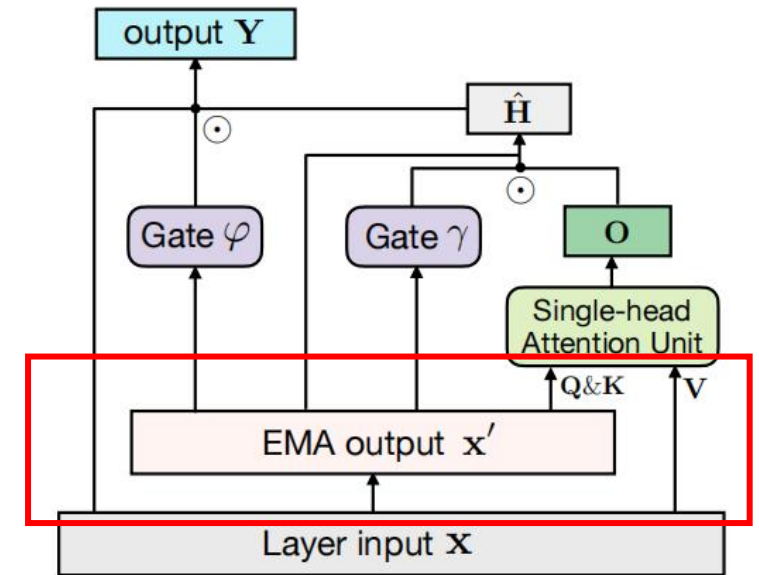
# Mega Architecture: Outline

- **Exponential Moving Average (EMA)**
  - Local dependencies that decaying exponentially over time
- **Single-head Gated Attention**
  - Adding a reset gate to the attention output
  - Theoretically proving that **single-head** gated attention is as expressive as multi-head one
- **Mega-Chunk**
  - Applying attention to local chunks of fixed length
  - Reducing quadratic complexity to **linear**



# Mega Architecture: Outline

- **Exponential Moving Average (EMA)**
  - Local dependencies that decaying exponentially over time
- **Single-head Gated Attention**
  - Adding a reset gate to the attention output
  - Theoretically proving that **single-head** gated attention is as expressive as multi-head one
- **Mega-Chunk**
  - Applying attention to local chunks of fixed length
  - Reducing quadratic complexity to **linear**





# Exponential Moving Average (EMA)

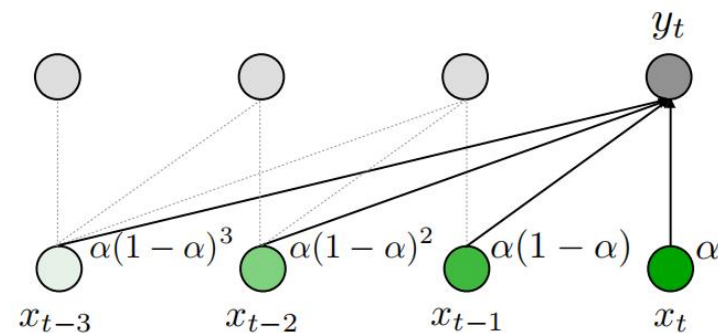
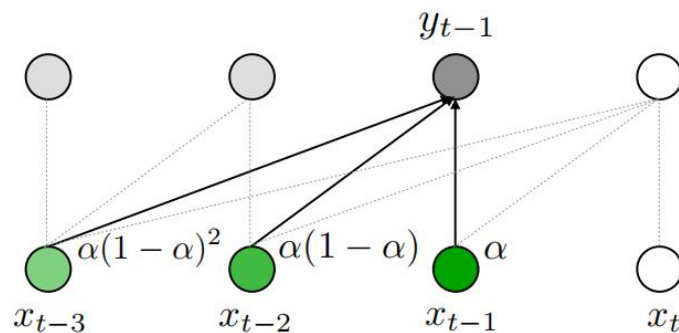
- **Notations:** Assuming 1-dim input sequence  $X = [X_1, X_2, \dots, X_n]$ ,  $X \in \mathbb{R}$

- **EMA:** 
$$y_t = \alpha \odot x_t + (1 - \alpha) \odot y_{t-1} \quad \alpha \in (0, 1)$$

Learnable Parameters

- **Damped EMA:** 
$$y_t = \alpha \odot x_t + (1 - \alpha \odot \delta) \odot y_{t-1} \quad \delta \in (0, 1)$$

Relaxing the coupled weights



# Multi-dimensional Damped EMA

- Expanding  $x_t$  to  $h$  dimensions

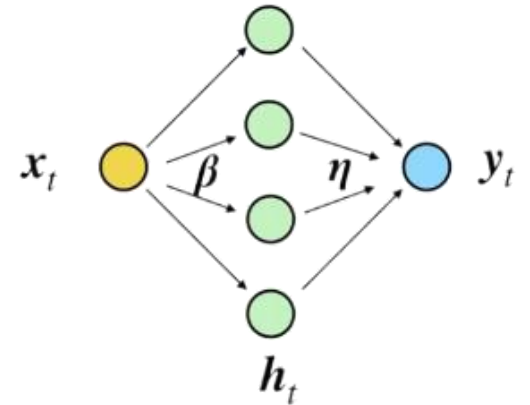
$$\mathbf{u}_t = \beta \mathbf{x}_t \in \mathbb{R}^h$$

- Applying damped EMA **individually** to each dimension

$$\mathbf{h}_t = \alpha \odot \mathbf{u}_t + (1 - \alpha \odot \delta) \odot \mathbf{h}_{t-1} \in \mathbb{R}^h$$

- Mapping the  $h$ -dimensional vector back to 1 dimension

$$y_t = \eta^T \mathbf{h}_t \in \mathbb{R}$$

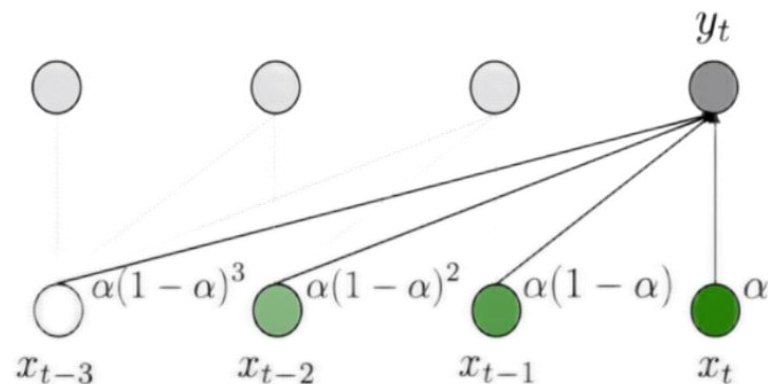
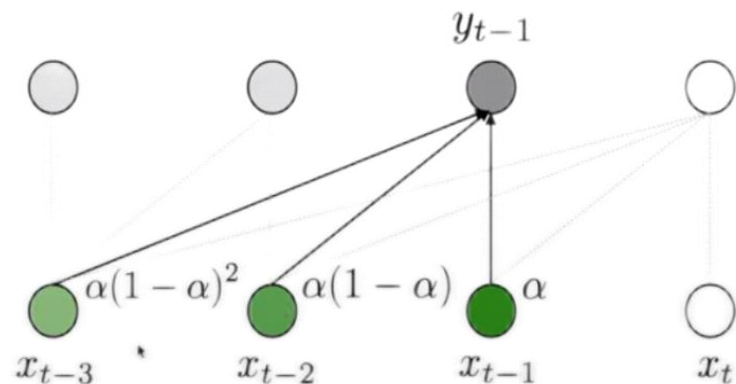


# Efficient computation of EMA

- Efficiently compute EMA outputs of all tokens in parallel

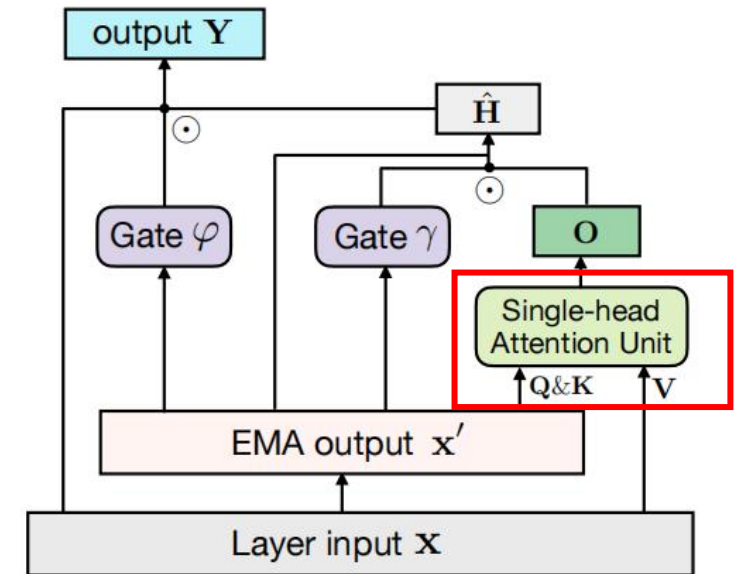
$$y_t = \alpha \odot \mathbf{x}_t + (1 - \alpha) \odot y_{t-1}$$

EMA weights are input independent



# Mega Architecture: Outline

- **Exponential Moving Average (EMA)**
  - Local dependencies that decaying exponentially over time
- **Single-head Gated Attention**
  - Adding a reset gate to the attention output
  - Theoretically proving that **single-head** gated attention is as expressive as multi-head one
- **Mega-Chunk**
  - Applying attention to local chunks of fixed length
  - Reducing quadratic complexity to **linear**



# Single-head Gated Attention

- Gated Attention Unit (GAU; Hua et al. (2022)) as the backbone
- first use the output from the EMA to compute the shared representation in GAU

$$\begin{aligned}\mathbf{X}' &= \text{EMA}(\mathbf{X}) && \in \mathbb{R}^{n \times d} \\ \mathbf{Z} &= \phi_{\text{silu}}(\mathbf{X}'\mathbf{W}_z + b_z) && \in \mathbb{R}^{n \times z}\end{aligned}$$

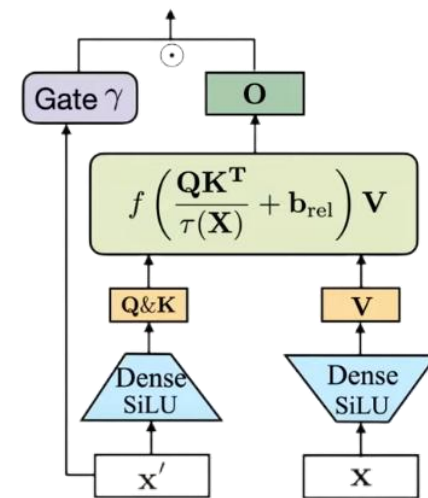
- the query and key sequences are computed by applying per-dimension scalars and offsets to  $\mathbf{Z}$ , and the value sequence is from the original  $\mathbf{X}$ :

$$\begin{aligned}\mathbf{Q} &= \boldsymbol{\kappa}_q \odot \mathbf{Z} + \boldsymbol{\mu}_q \\ \mathbf{K} &= \boldsymbol{\kappa}_k \odot \mathbf{Z} + \boldsymbol{\mu}_k \\ \mathbf{V} &= \phi_{\text{silu}}(\mathbf{X}\mathbf{W}_v + b_v)\end{aligned}$$

# Single-head Gated Attention

- output of attention:

$$\mathbf{O} = f \left( \frac{\mathbf{Q}\mathbf{K}^T}{\tau(\mathbf{X})} + \mathbf{b}_{\text{rel}} \right) \mathbf{V}$$



Subsequently, MEGA introduces the reset gate  $\gamma$ , the update gate  $\varphi$ , and computes the candidate activation output  $\hat{\mathbf{H}}$ :

$$\gamma = \phi_{\text{silu}}(\mathbf{X}'W_{\gamma} + b_{\gamma}) \in \mathbb{R}^{n \times v} \quad (12)$$

$$\varphi = \phi_{\text{sigmoid}}(\mathbf{X}'W_{\varphi} + b_{\varphi}) \in \mathbb{R}^{n \times d} \quad (13)$$

$$\hat{\mathbf{H}} = \phi_{\text{silu}}(\mathbf{X}'W_h + (\gamma \odot \mathbf{O})U_h + b_h) \in \mathbb{R}^{n \times d} \quad (14)$$

The final output  $\mathbf{Y}$  is computed with the update gate  $\varphi$ :

$$\mathbf{Y} = \varphi \odot \hat{\mathbf{H}} + (1 - \varphi) \odot \mathbf{X} \in \mathbb{R}^{n \times d} \quad (15)$$



# Single-head Gated Attention

- Single-head gated attention is as expressive as multi-head one

$$\mathbf{O}_{\text{SHA}} = \mathbf{a}^T \mathbf{V} = \begin{bmatrix} \mathbf{a}^T \mathbf{V}^{(1)} \\ \vdots \\ \mathbf{a}^T \mathbf{V}^{(h)} \end{bmatrix}, \quad \mathbf{O}_{\text{MHA}} = \begin{bmatrix} \mathbf{a}^{(1)T} \mathbf{V}^{(1)} \\ \vdots \\ \mathbf{a}^{(h)T} \mathbf{V}^{(h)} \end{bmatrix} \quad (19)$$

It is straightforward to see that  $\mathbf{O}_{\text{MHA}}$  is more expressive than  $\mathbf{O}_{\text{SHA}}$ , because  $\mathbf{O}_{\text{MHA}}$  leverages  $h$  sets of attention weights.

In the single-head gated attention, we introduce a gate vector  $\gamma = \mathcal{G}(\mathbf{X})$  for each  $\mathbf{q}$ , and the output of single-head gated attention is  $\mathbf{O}_{\text{SHGA}} = \mathbf{O}_{\text{SHA}} \odot \gamma$ . The following theorem reveals the equivalence of  $\mathbf{O}_{\text{SHGA}}$  and  $\mathbf{O}_{\text{MHA}}$  w.r.t expressiveness (proof in Appendix B):

**Theorem 1** *Suppose the transformation  $\mathcal{G}$  is a universal approximator. Then,  $\forall \mathbf{X}, \exists \gamma = \mathcal{G}(\mathbf{X})$  s.t.*

$$\mathbf{O}_{\text{SHGA}} = \mathbf{O}_{\text{MHA}} \quad (20)$$

Theorem 1 indicates that by simply introducing the gate vector,  $\mathbf{O}_{\text{SHGA}}$  is as expressive as  $\mathbf{O}_{\text{MHA}}$ . In practice,  $\mathcal{G}$  is commonly modeled by a (shallow) neural network, whose universality of approximation has been extensively studied (Hornik et al., 1989; Yarotsky, 2017; Park et al., 2020).

# Single-head Gated Attention

- Single-head gated attention is as expressive as multi-head one

**Proof** We split  $\gamma$  into  $h$  heads in the same way as  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$ :

$$\gamma = \begin{bmatrix} \gamma^{(1)} \\ \vdots \\ \gamma^{(h)} \end{bmatrix}$$

Then we have

$$\mathbf{O}_{\text{SHGA}} = \mathbf{a}^T \mathbf{V} \odot \gamma = \begin{bmatrix} \mathbf{a}^T \mathbf{V}^{(1)} \odot \gamma^{(1)} \\ \vdots \\ \mathbf{a}^T \mathbf{V}^{(h)} \odot \gamma^{(h)} \end{bmatrix}$$

To prove Theorem 1, we need to find  $\gamma$  such that

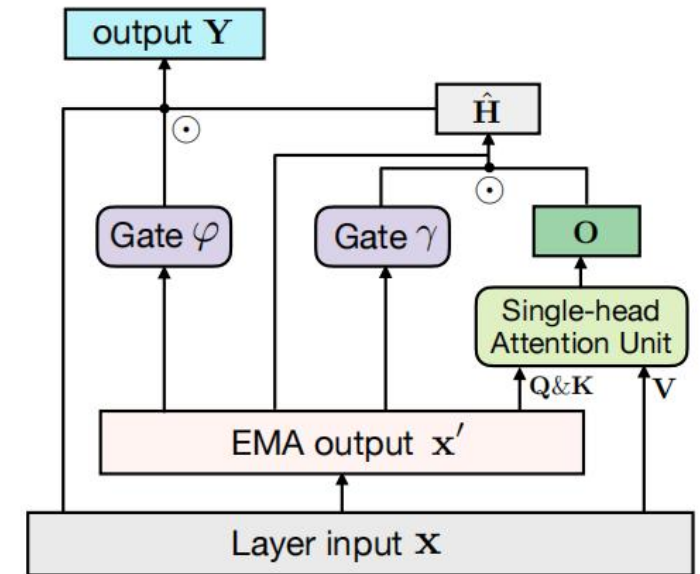
$$\mathbf{a}^T \mathbf{V}^{(i)} \odot \gamma^{(i)} = \mathbf{a}^{(i)T} \mathbf{V}^{(i)} \iff \gamma^{(i)} = \mathbf{a}^{(i)T} \mathbf{V}^{(i)} \oslash \mathbf{a}^T \mathbf{V}^{(i)}, \forall i \in \{1, \dots, h\},$$

where  $\oslash$  is the element-wise divide operation. Since  $\mathcal{G}(\mathbf{X})$  is a universal approximator and  $\mathbf{Q}$ ,  $\mathbf{K}$ ,  $\mathbf{V}$  and  $\mathbf{a}$  are all transformed from  $\mathbf{X}$ ,  $\gamma$  can theoretically recover  $\mathbf{a}^{(i)T} \mathbf{V}^{(i)} \oslash \mathbf{a}^T \mathbf{V}^{(i)}, \forall \mathbf{X}$ . ■



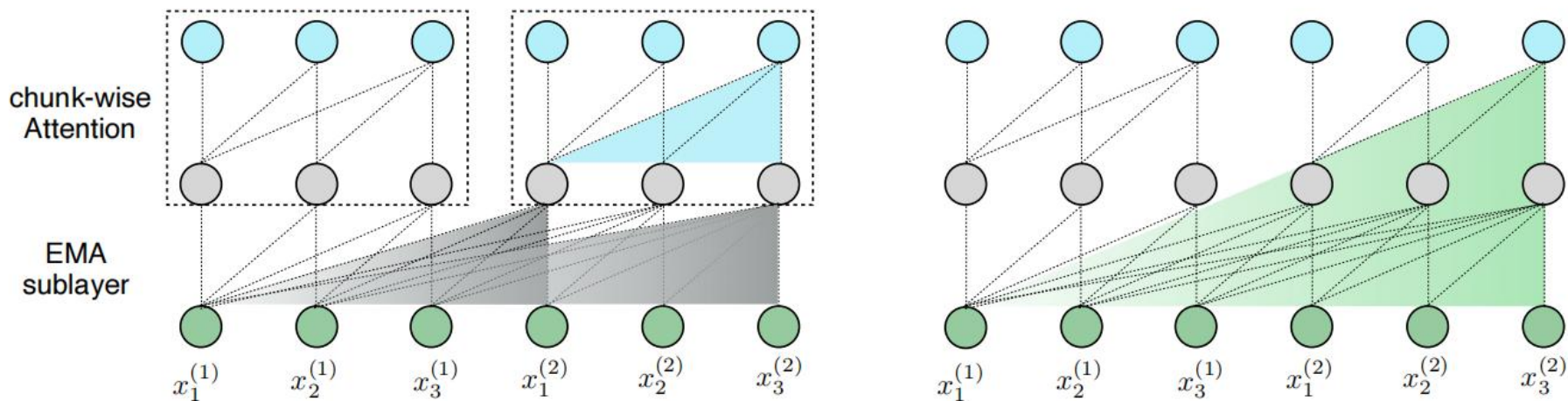
# Mega Architecture: Outline

- **Exponential Moving Average (EMA)**
  - Local dependencies that decaying exponentially over time
- **Single-head Gated Attention**
  - Adding a reset gate to the attention output
  - Theoretically proving that **single-head** gated attention is as expressive as multi-head one
- **Mega-Chunk**
  - Applying attention to local chunks of fixed length
  - Reducing quadratic complexity to **linear**



# Mega - Chunk: Efficient Mega

- Split input sequences into multiple chunks with fixed length
- Applying attention individually to each chunk
  - Linear complexity & easy implementation
  - Losing contextual information between chunks
  - **EMA preserves the information from previous chunks**



# Experiments

- **Long Range Arena (LRA)**
  - 3 tasks on byte-level text classification
  - 3 tasks on pixel-level image classification
- **Language Modeling**
  - Enwiki8 (character-level)
  - WikiText-103(Word-level)
- **Machine Translation**
  - WMT'14 English - German
- **Image Classification**
  - ImageNet-1K
- **Raw Speech Classification**
  - Speech Commands

# Experimental Results

Table 2: (**Long Range Arena**) Accuracy on the full suite of long range arena (LRA) tasks, together with training speed and peak memory consumption comparison on the Text task with input length of 4K. ‡ indicates results replicated by us.

Models	ListOps	Text	Retrieval	Image	Pathfinder	Path-X	Avg.	Speed	Mem.
XFM	36.37	64.27	57.46	42.44	71.40	<b>X</b>	54.39	–	–
XFM‡	37.11	65.21	79.14	42.94	71.83	<b>X</b>	59.24	1×	1×
Reformer	37.27	56.10	53.40	38.07	68.50	<b>X</b>	50.67	0.8×	0.24×
Linformer	35.70	53.94	52.27	38.56	76.34	<b>X</b>	51.36	5.5×	0.10×
BigBird	36.05	64.02	59.29	40.83	74.87	<b>X</b>	55.01	1.1×	0.30×
Performer	18.01	65.40	53.82	42.77	77.05	<b>X</b>	51.41	<b>5.7×</b>	<b>0.11×</b>
Luna-256	37.98	65.78	79.56	47.86	78.55	<b>X</b>	61.95	4.9×	0.16×
S4-v1	58.35	76.02	87.09	87.26	86.05	88.10	80.48	–	–
S4-v2	59.60	86.82	90.90	88.65	94.20	96.35	86.09	–	–
S4-v2‡	59.10	86.53	90.94	88.48	94.01	96.07	85.86	4.8×	0.14×
MEGA	<b>63.14</b>	<b>90.43</b>	<b>91.25</b>	<b>90.44</b>	<b>96.01</b>	<b>97.98</b>	<b>88.21</b>	2.9×	0.31×
MEGA-chunk	58.76	90.19	90.97	85.80	94.41	93.81	85.66	5.5×	0.13×

# Experimental Results

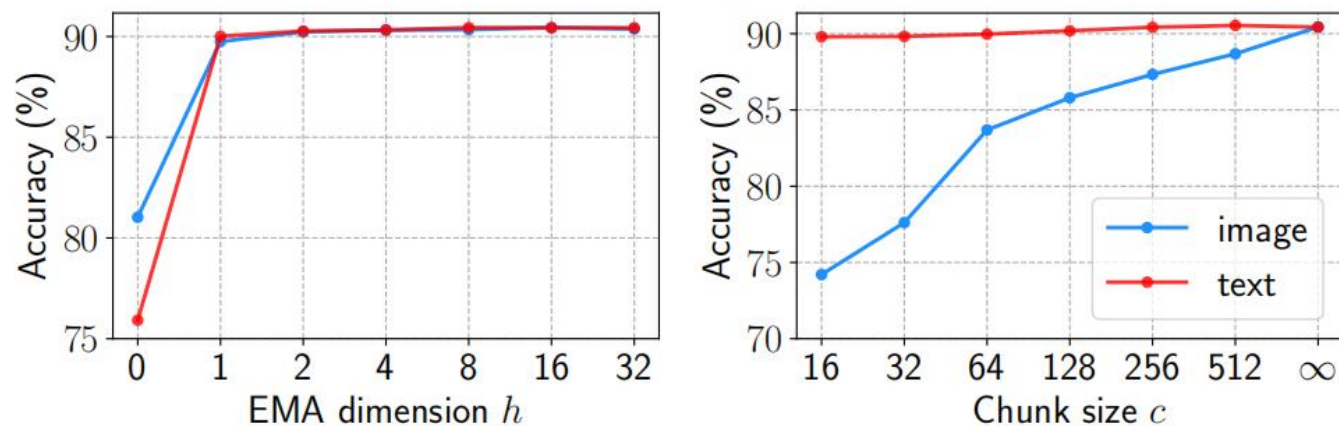


Figure 4: Ablations on EMA dimension and chunk size.

	Text	Image
softmax	<b>90.43</b>	89.87
relu <sup>2</sup>	90.08	90.22
laplace	90.22	<b>90.43</b>

Table 3: Attention functions.





# Thank you for listening

主講人：吳雨欣

2023.6.14