



Institute of Computing
Technology, CAS



Research Center for Intelligent
Computing Systems

CTA: Hardware-Software Co-design for Compressed Token Attention Mechanism

Haoran Wang, Haobo Xu, Ying Wang, Yinhe Han

Research Center for Intelligent Computing Systems
Institute of Computing Technology, Chinese Academy of Sciences

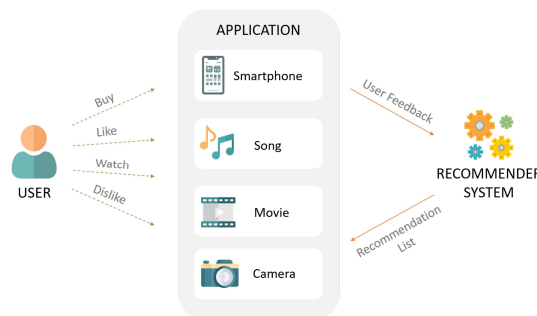
HPCA 2023

Attention mechanism significantly improves neural network capability

- Transformer-based models have been widely applied
 - Natural Language Processing(NLP): machine translation, question answering, text classification, language modeling...
 - Computer Vision(CV): image segmentation, video classification, zero-shot image classification...
 - Recommendation system



Machine translation:
Google translate



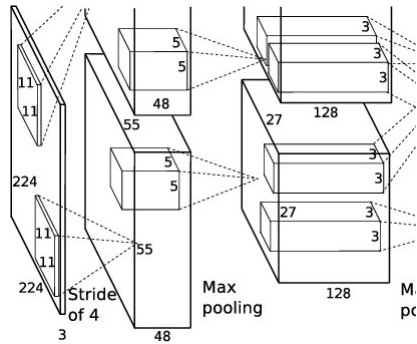
Online recommendations



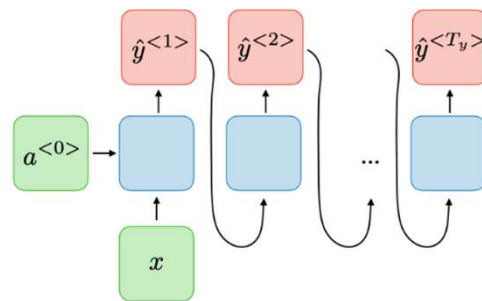
Language model for dialogue:
ChatGPT

Attention mechanism significantly improves neural network capability

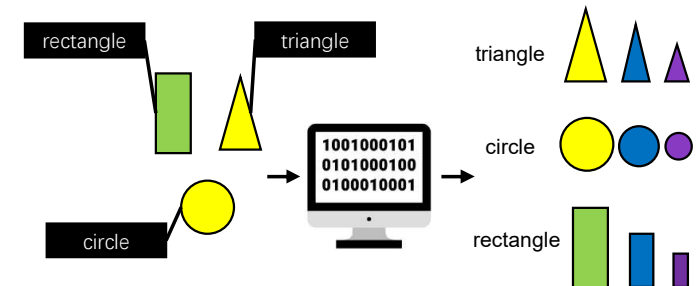
- The attention mechanism is the key enabler in Transformer-based models
 - Modeling relations among elements without regard to their distance in sequence
 - Parallelizing the sequence modeling
 - Enabling using larger amount of no-label data and achieves better performance



CNN models local relations by convolution but fail to model distant relations



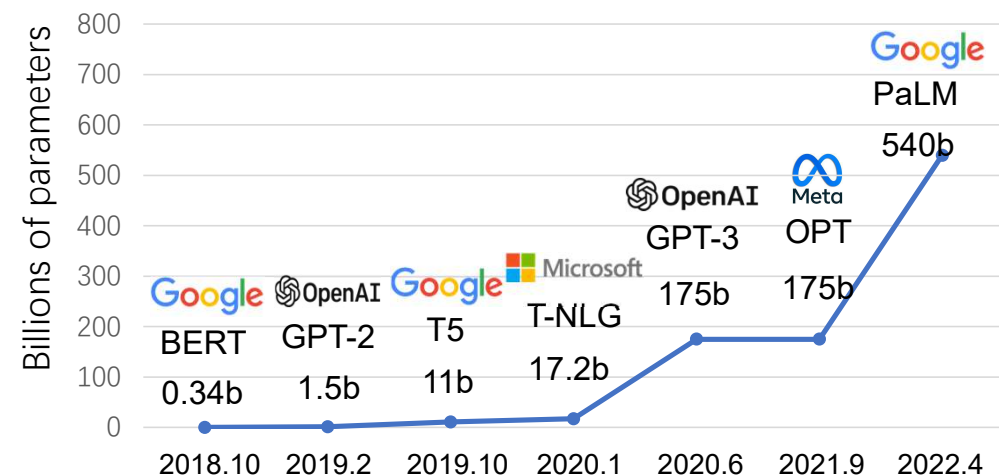
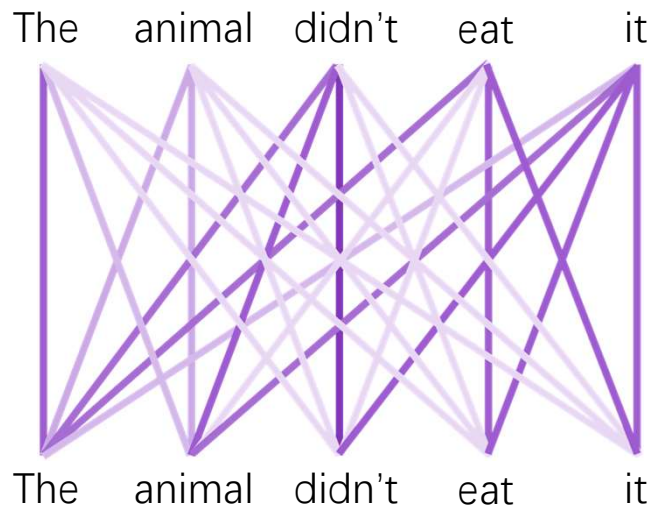
RNN utilizes sequential computation with less parallelism



Previous NN models need more labeling to achieve intelligence

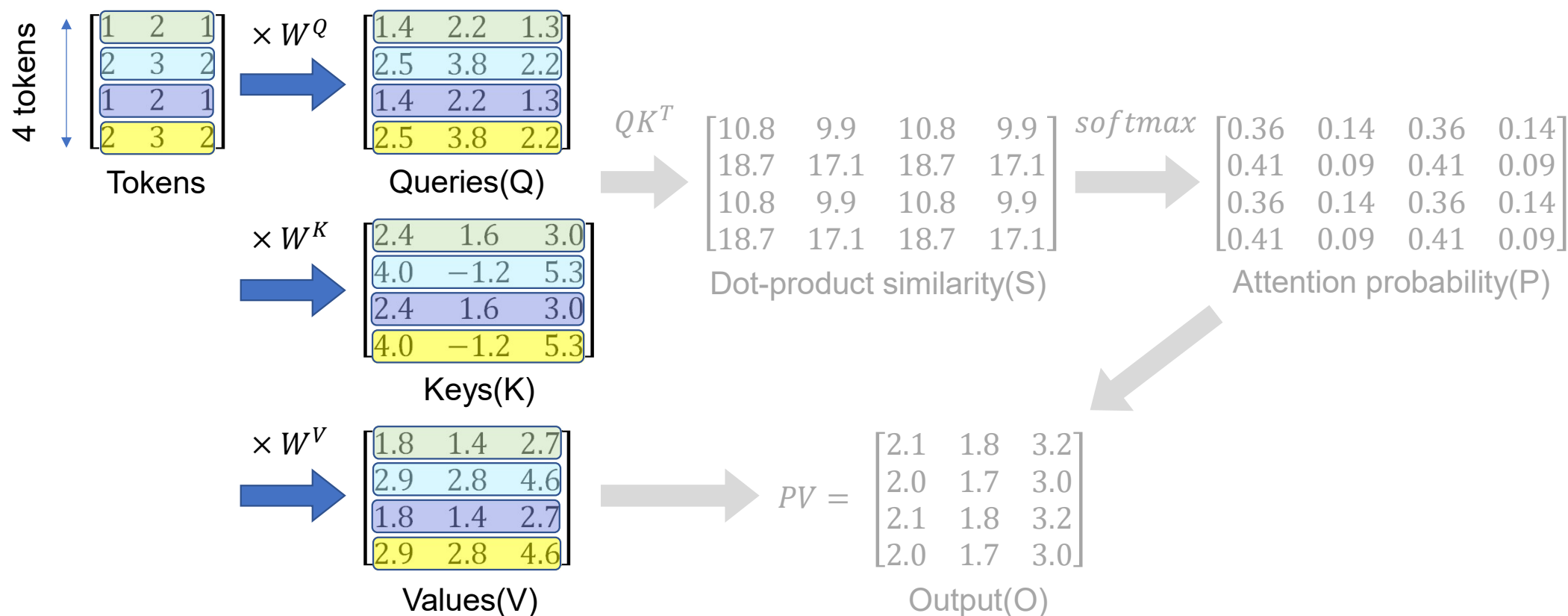
Efficient attention mechanism is important

- Attention mechanism incurs high overhead
 - Attention mechanism models dense relations among sequence elements: 30%-50% of the computation time is spent on attention mechanism in Transformer-based models
 - Transformer-based models are growing recklessly in size



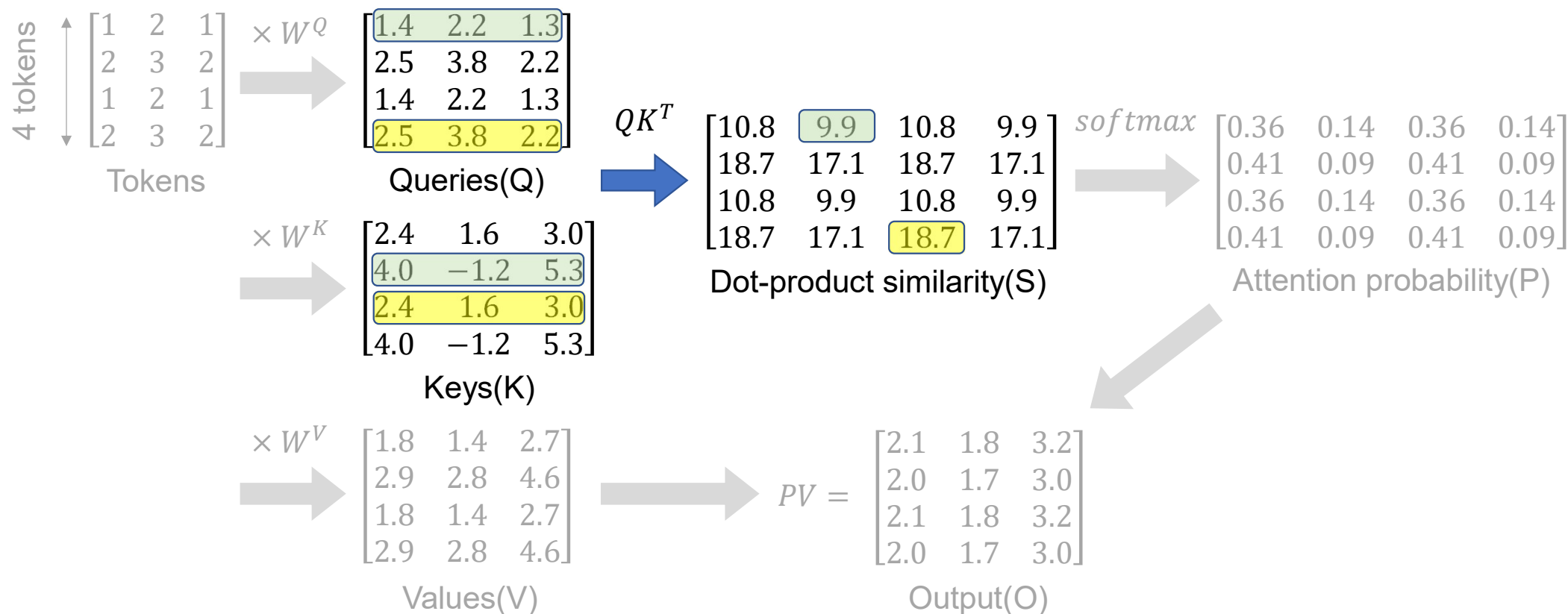
Attention mechanism overview

Step1: Three linear transformations are applied on tokens to generate queries, keys and values



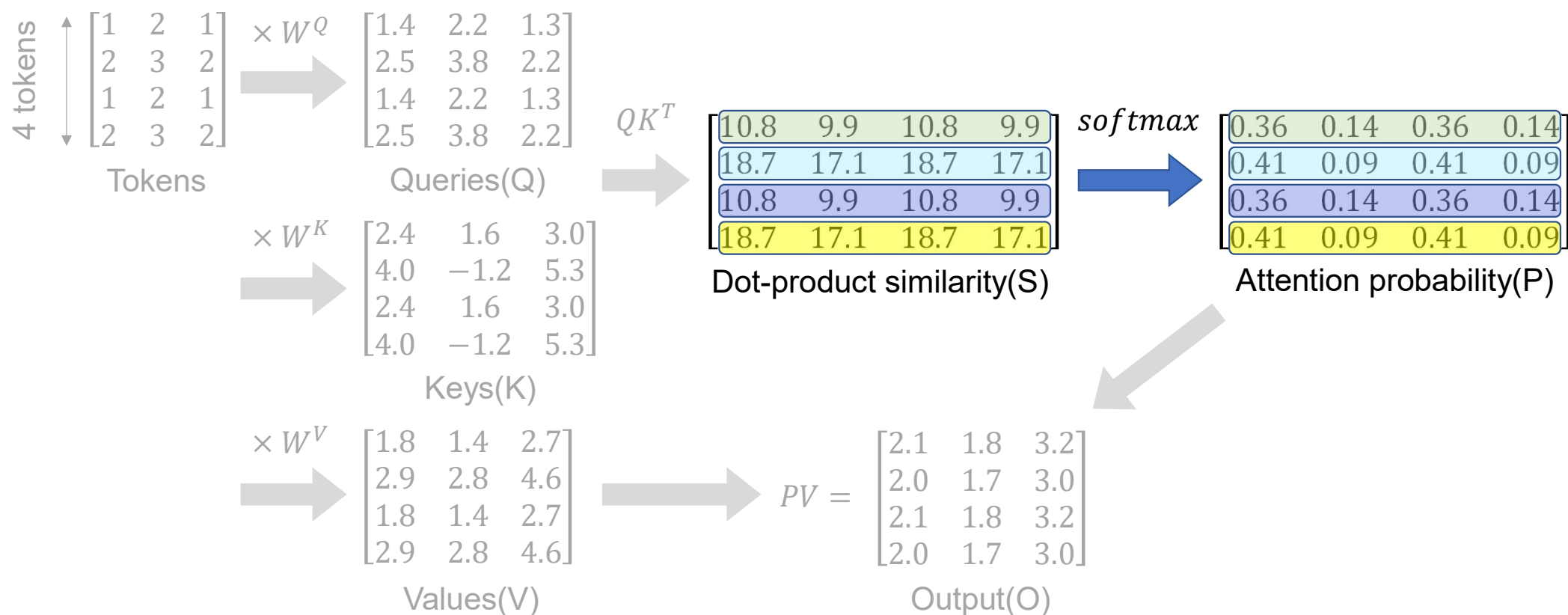
Attention mechanism overview

Step2: Calculate dot-product similarity score between each pair of query and key



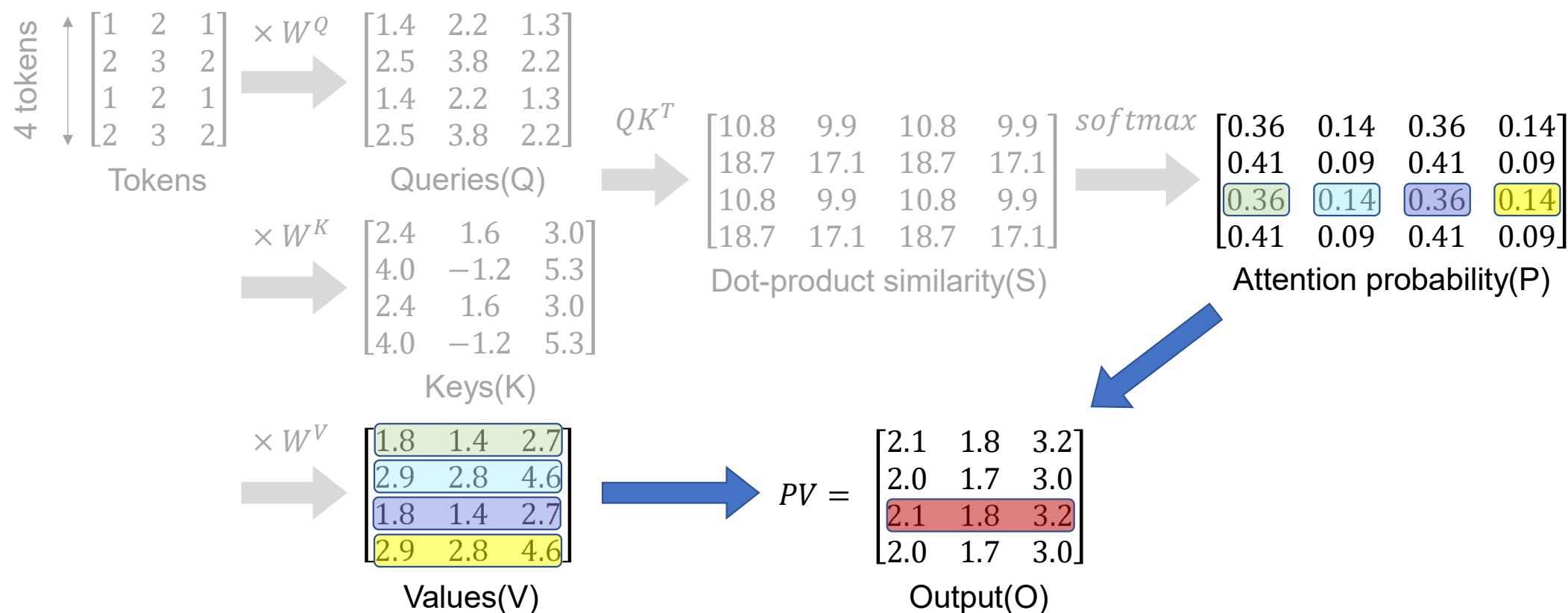
Attention mechanism overview

Step3: Apply softmax to normalize each row of similarity scores to generate attention probabilities



Attention mechanism overview

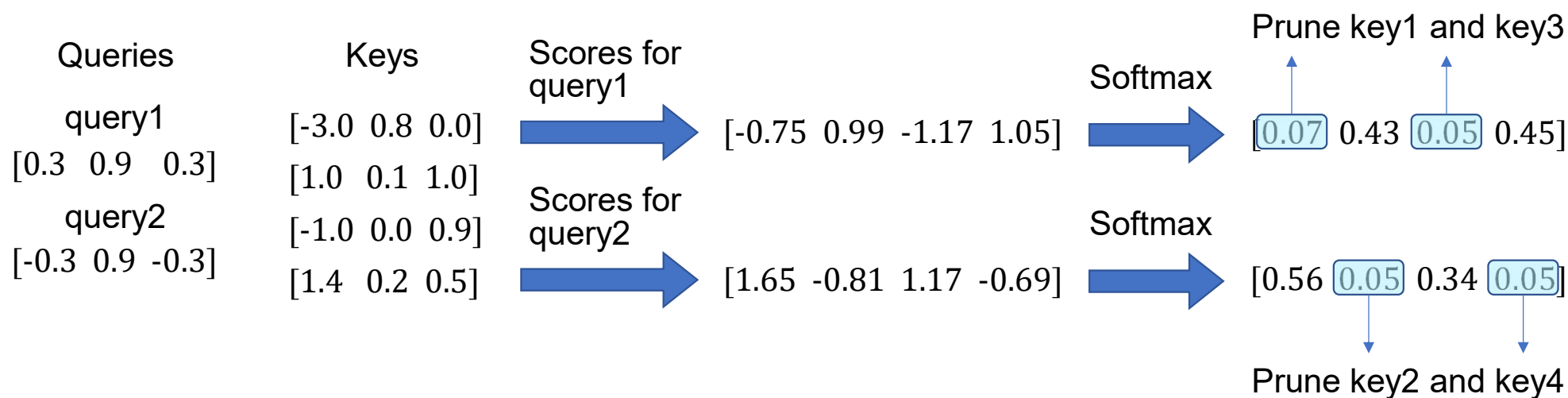
Step4: Using attention probabilities as weights to calculate weighted sum of values



Previous insights

KEY OBSERVATION: Given a query, keys do not contribute equally to the output.

After softmax normalization, for those keys that have small dot-product with this query, their attention probability is near 0 and can be pruned.

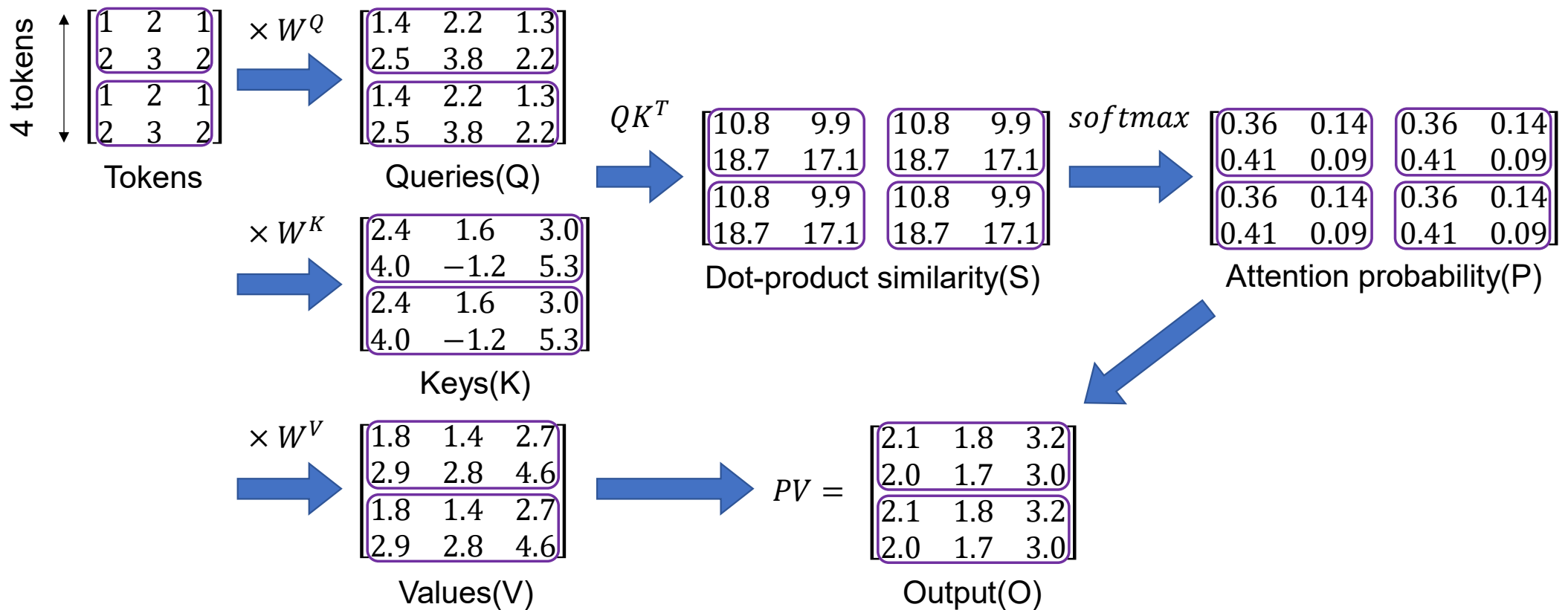


Query-specific pruning breaks the matrix multiplication formulation for multi-query processing.

Our work keeps the matrix multiplication formulation with shrinking matrix scale.

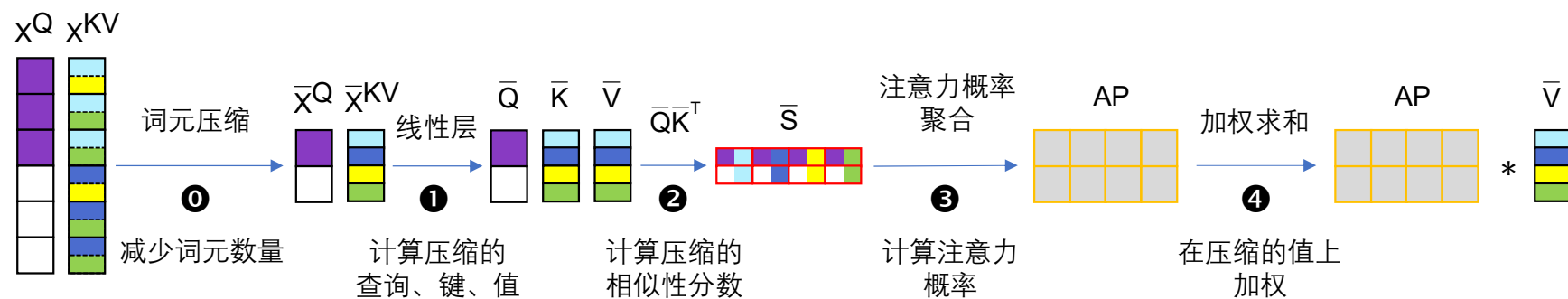
Motivation

KEY IDEA: Similar tokens lead to [similar patterns throughout attention computation](#). Large amount of computation can be avoided if we can remove semantic repetitions in tokens.

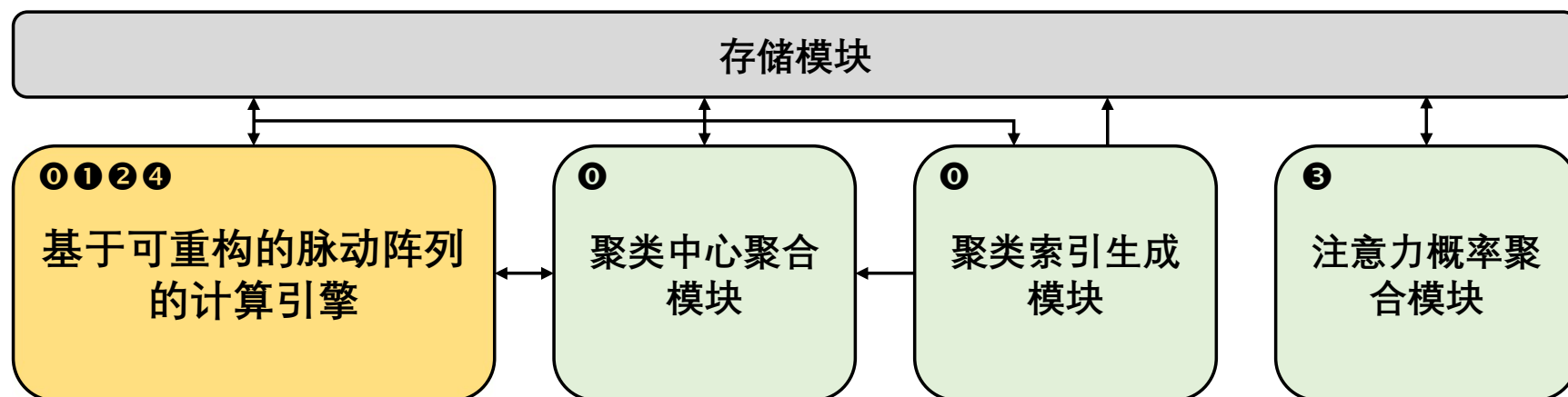


Overview: Compressed Token Attention Mechanism (CTA)

算法

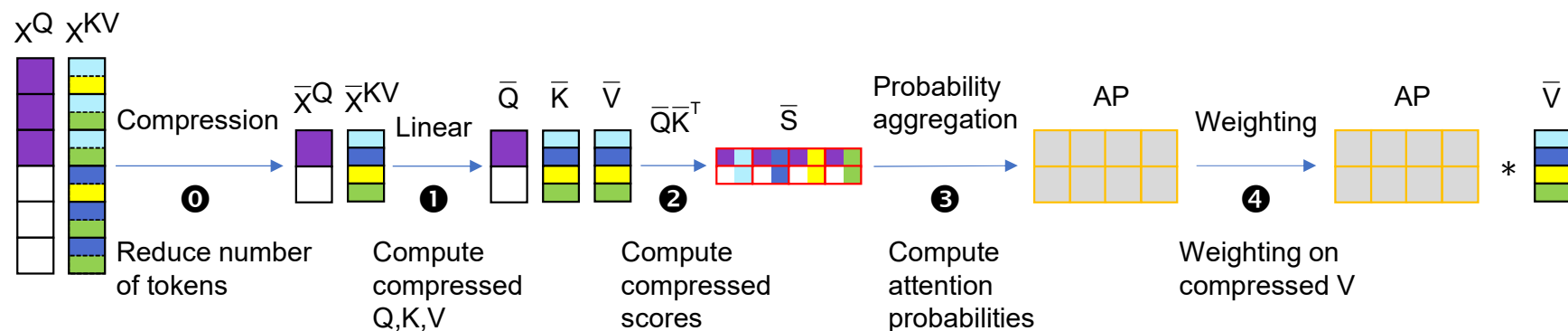


架构

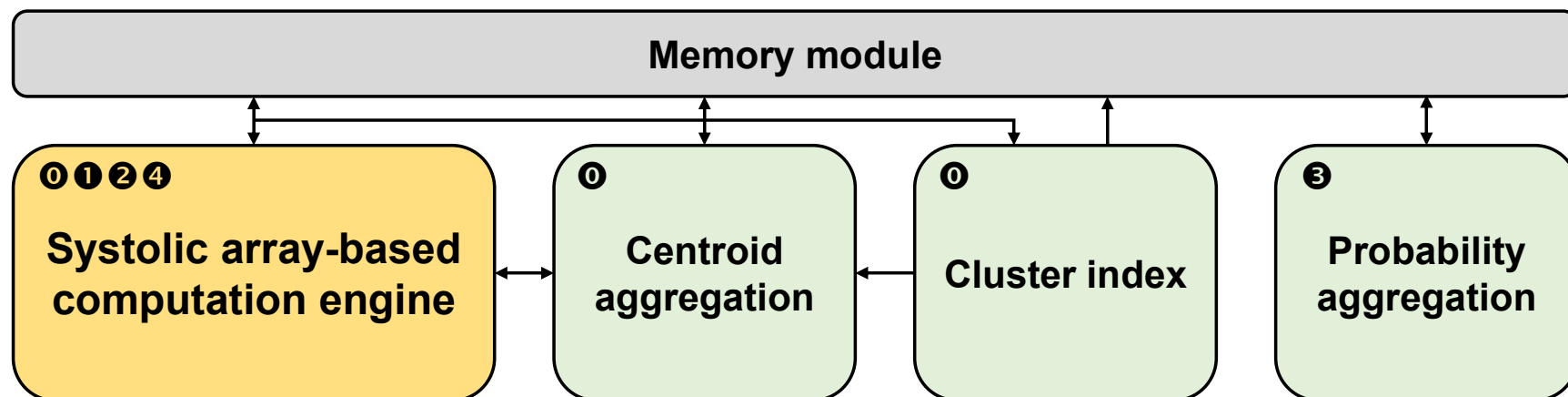


Overview: Compressed Token Attention Mechanism (CTA)

Algorithm



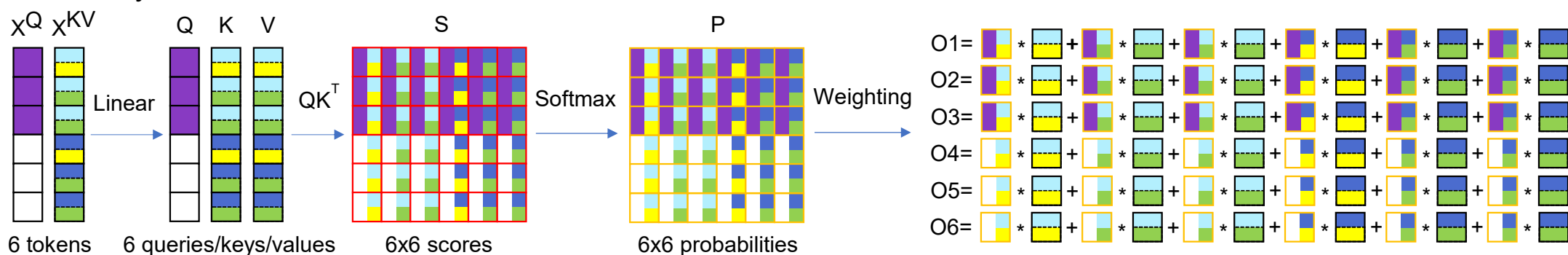
Architecture



Algorithm overview

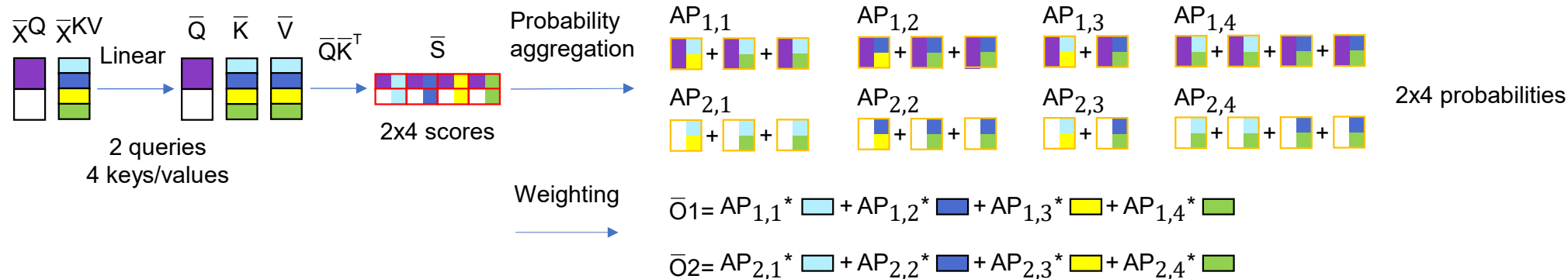
Compress tokens to remove repeating semantics, saving attention computation.

Ordinary attention



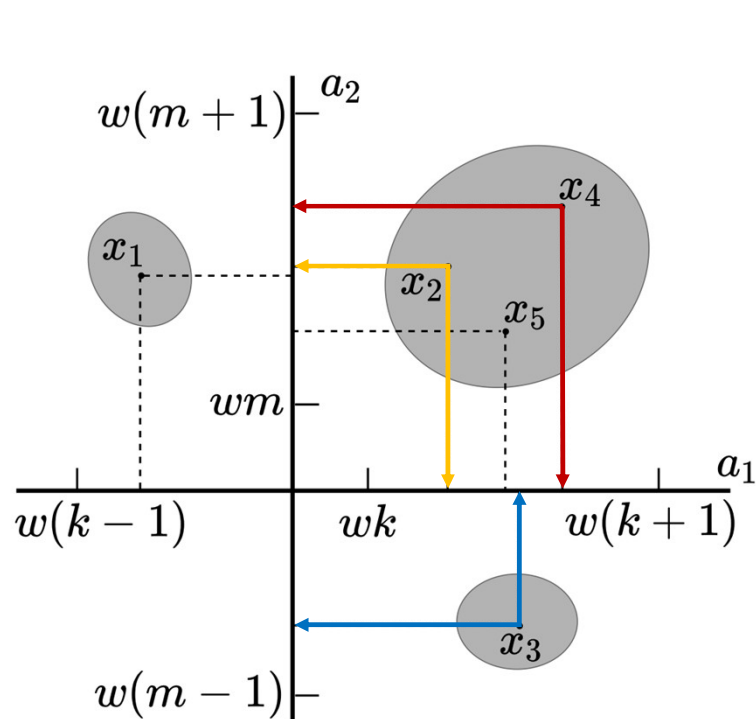
Compression

CTA attention



Efficiently remove semantic repetitions in tokens

Goal: group similar tokens and find appropriate representation for grouped tokens



$$\begin{array}{cc} a_1 & a_2 \\ \downarrow & \downarrow \\ h(x_2) = (k, m) \\ h(x_3) = (k, m-1) \\ h(x_4) = (k, m) \end{array}$$

Tokens (X) projects to which interval on direction axes (A)

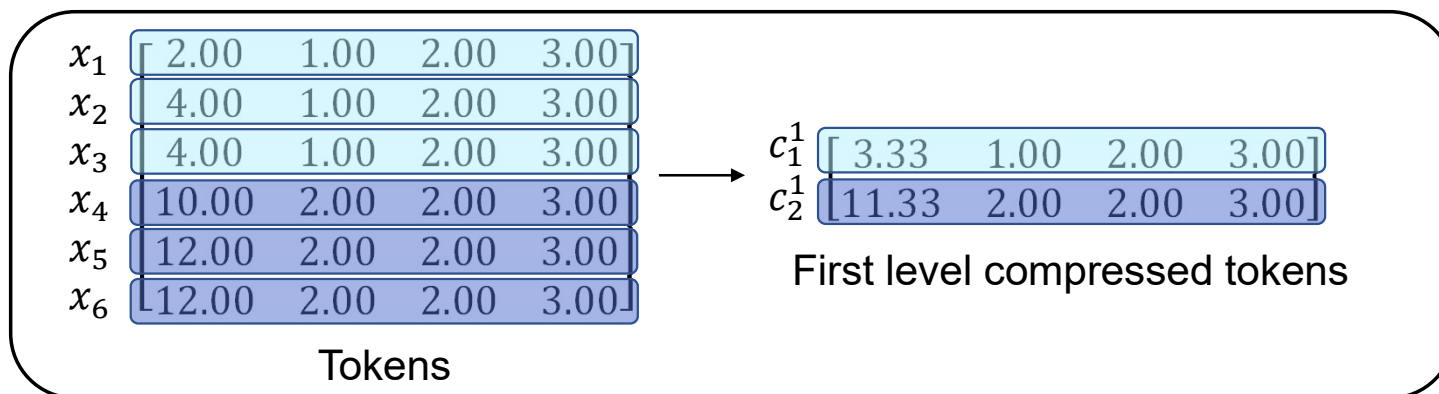
$$H = \lfloor (A \cdot X^T + B) / w \rfloor$$

- Group tokens with the same hash value into a cluster
- Compress tokens to their cluster centroids - average of tokens of their cluster

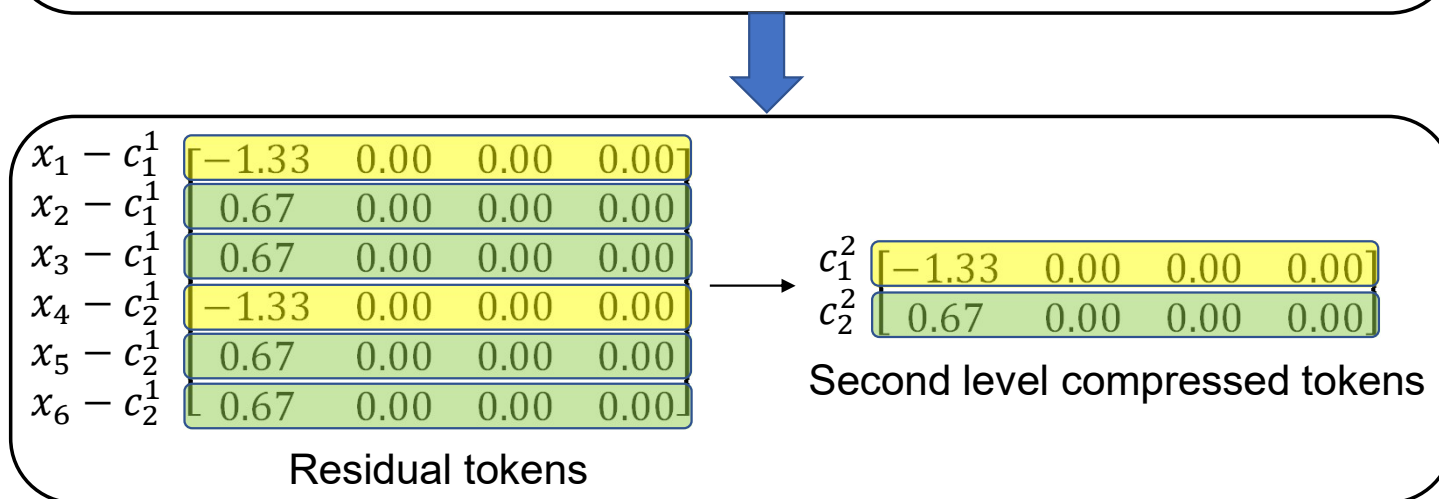
$$x_1, x_2, x_3, x_4, x_5 \xrightarrow{\text{compression}} x_1, \frac{x_2 + x_4 + x_5}{3}, x_3$$

Step ①: Token compression

Efficiently extracting semantic features in tokens, yielding compressed tokens with **little information loss** compared to original tokens



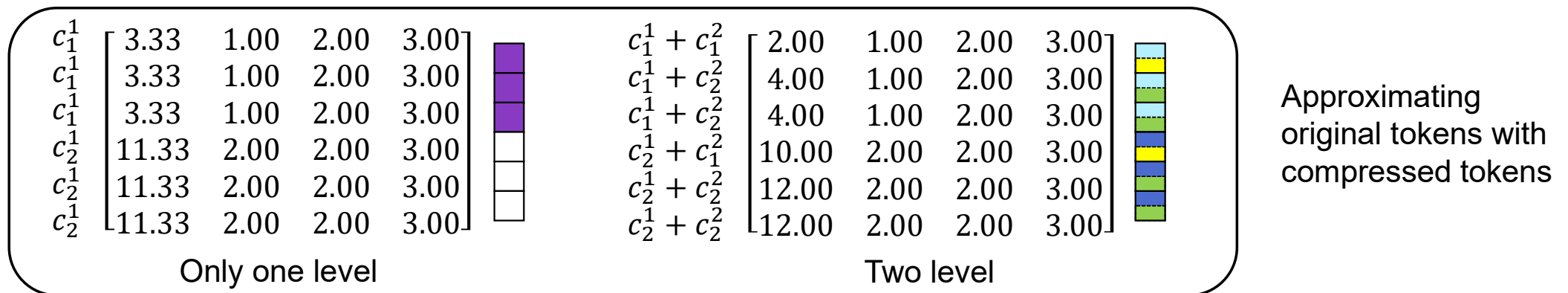
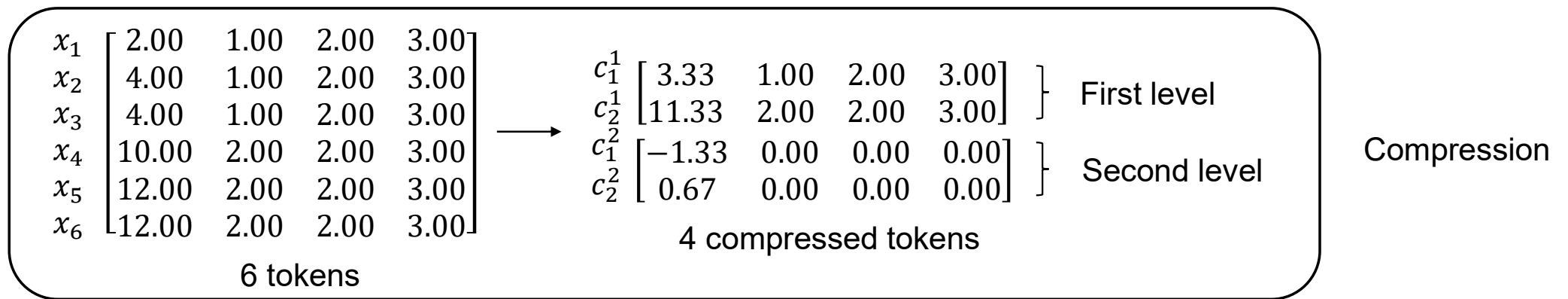
First level compression:
tokens are compressed



Second level compression:
residual tokens are further
compressed

Step ①: Token compression

Approximate original tokens with compressed tokens utilizing only first level compressed tokens or the sum of two levels' compressed tokens

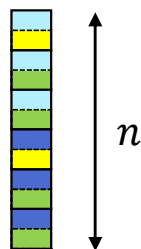


Step ①: Token compression

In this example,

$$X^Q = X^{KV} = \begin{bmatrix} 2.00 & 1.00 & 2.00 & 3.00 \\ 4.00 & 1.00 & 2.00 & 3.00 \\ 4.00 & 1.00 & 2.00 & 3.00 \\ 10.00 & 2.00 & 2.00 & 3.00 \\ 12.00 & 2.00 & 2.00 & 3.00 \\ 12.00 & 2.00 & 2.00 & 3.00 \end{bmatrix} \approx$$

or



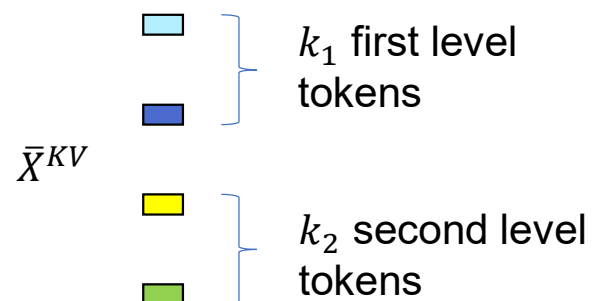
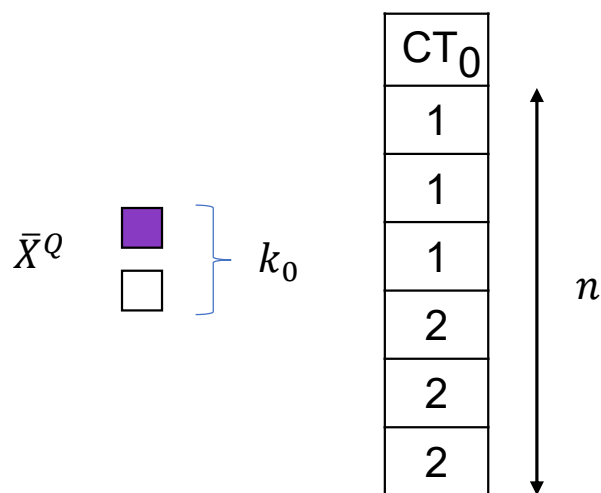
METHOD

1-level compression for query tokens

2-level compression for key-value tokens

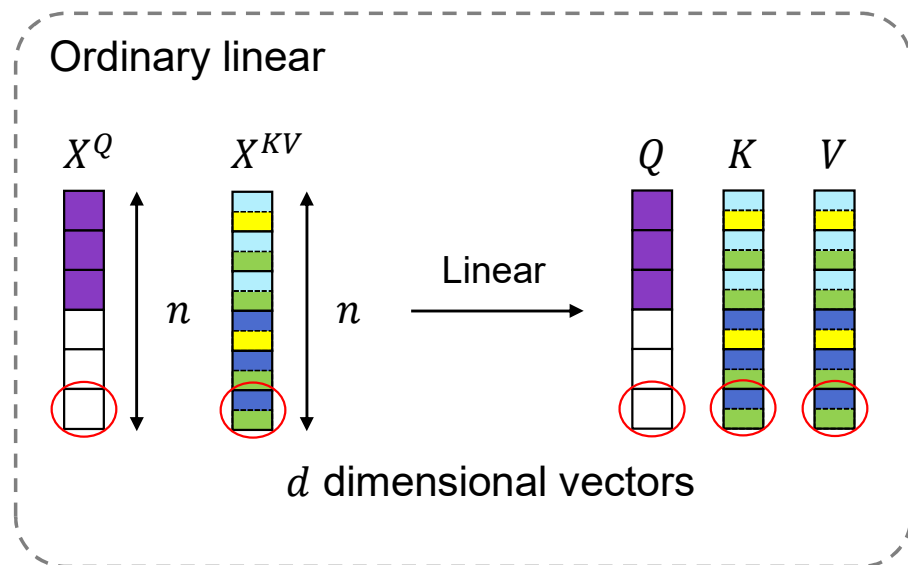
1-level compression
for X^Q

2-level compression
for X^{KV}

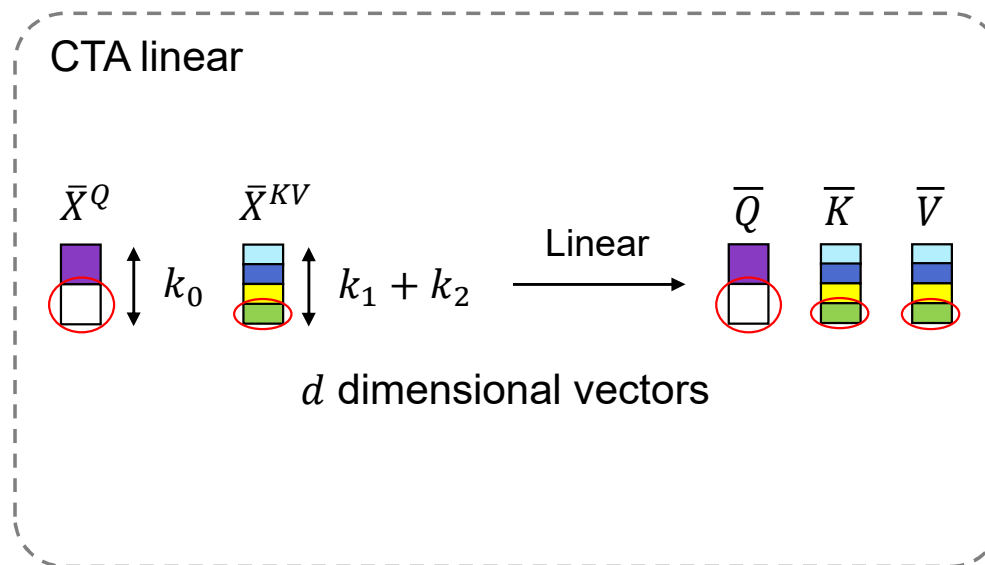


Step ❶: Linear transformations on compressed tokens

Linear transformations maps compressed tokens to compressed queries, keys and values



$$Q = X^Q \cdot W^Q, K = X^{KV} \cdot W^K, V = X^{KV} \cdot W^V$$



$$\bar{Q} = \bar{X}^Q \cdot W^Q, \bar{K} = \bar{X}^{KV} \cdot W^K, \bar{V} = \bar{X}^{KV} \cdot W^V$$

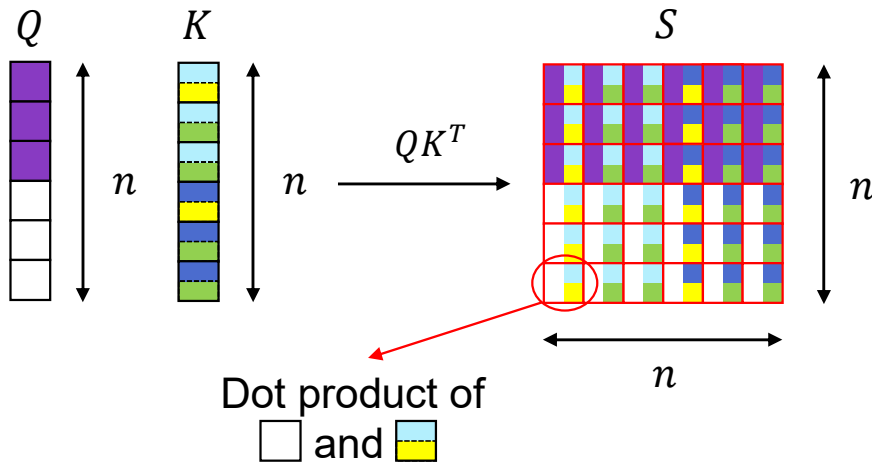
Computation for queries is reduced from nd^2 to k_0d^2

Computation for keys/values is reduced from nd^2 to $(k_1 + k_2)d^2$

Step ②: Similarities between compressed queries and keys

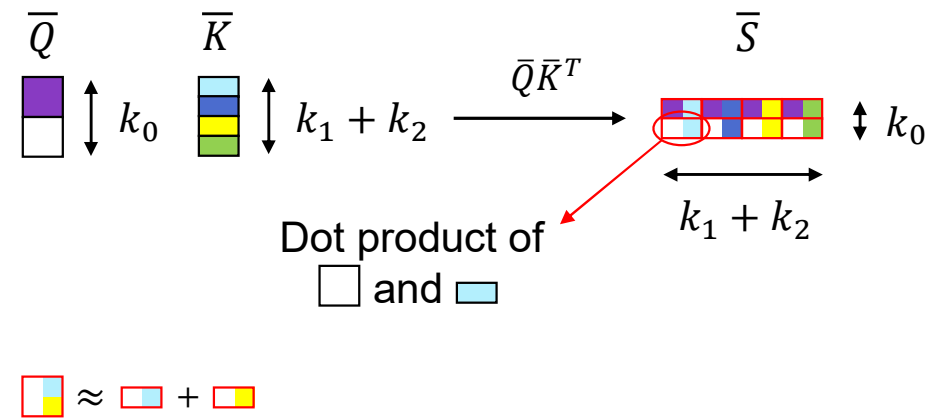
Compressed queries and keys generates compressed similarities

Ordinary similarities



$$S = (Q \cdot K^T) / \sqrt{d}$$

CTA similarities

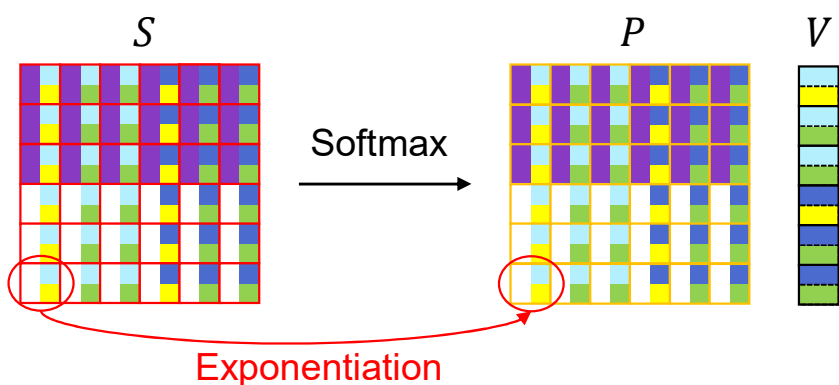


$$\bar{S} = (\bar{Q} \cdot \bar{K}^T) / \sqrt{d}$$

Computation is reduced from n^2d to $k_0(k_1 + k_2)d$

Step ③: Attention probability aggregation

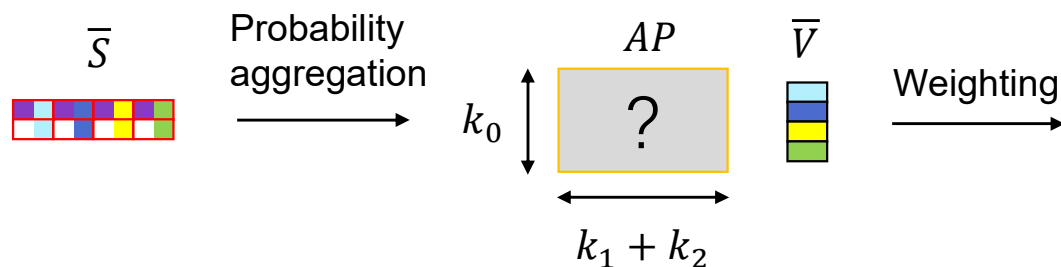
Ordinary Softmax and weighting



We omit softmax division in step ③, because it can be equivalently put off to the end of step ④

$$O_1 = \begin{matrix} \text{[colored square]} & * & \text{[colored square]} & + & \text{[colored square]} & * & \text{[colored square]} & + & \text{[colored square]} & * & \text{[colored square]} & + & \text{[colored square]} & * & \text{[colored square]} & + & \text{[colored square]} & * & \text{[colored square]} & + & \text{[colored square]} & * & \text{[colored square]} \end{matrix}$$

CTA attention probability aggregation and weighting



Key idea:
Find a way of weighting
compressed values that can
approximate ordinary output

Step ③: Attention probability aggregation

Reform output into the form of weighted sum of compressed values

Observation

$$O_1 = \begin{matrix} \text{[purple, yellow]} * \text{[light blue, yellow]} & + & \text{[purple, green]} * \text{[light blue, green]} & + & \text{[purple, green]} * \text{[light blue, green]} & + & \text{[purple, yellow]} * \text{[dark blue, yellow]} & + & \text{[purple, green]} * \text{[dark blue, green]} & + & \text{[purple, green]} * \text{[dark blue, green]} \end{matrix}$$

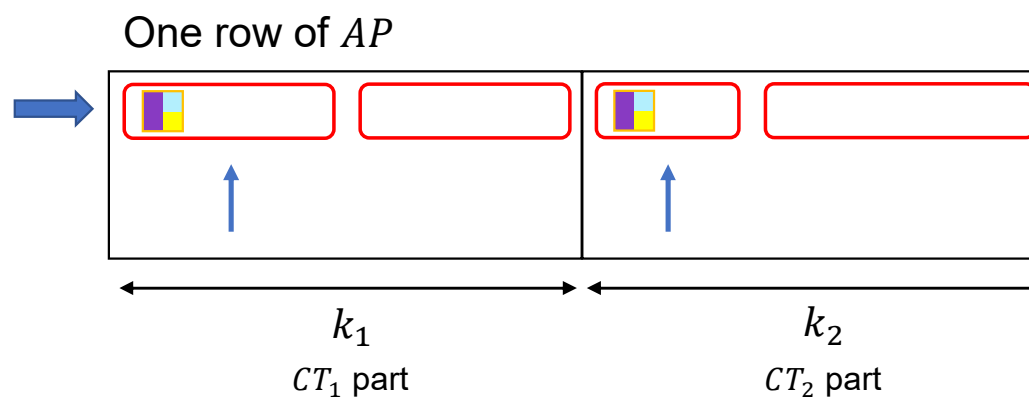
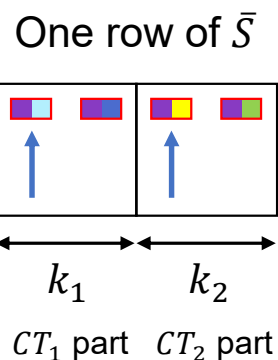
$$= (\text{[purple, yellow]} + \text{[purple, green]} + \text{[purple, green]}) * \text{[light blue]} + (\text{[purple, green]} + \text{[purple, green]} + \text{[purple, green]}) * \text{[dark blue]} + (\text{[purple, yellow]} + \text{[purple, yellow]}) * \text{[yellow]} + (\text{[purple, green]} + \text{[purple, green]} + \text{[purple, green]} + \text{[purple, green]}) * \text{[green]}$$

- Split value vectors into sum of compressed value vectors
- Merge coefficients

Aggregation algorithm

| CT ₁ | CT ₂ |
|-----------------|-----------------|
| 1 | 1 |
| 1 | 2 |
| 1 | 2 |
| 2 | 1 |
| 2 | 2 |
| 2 | 2 |

Each time select two scores from \bar{S} according to CT_1 and CT_2 , take the exponent of their sum and add to the same positions in AP



Step ③: Attention probability aggregation

Reform output into the form of weighted sum of compressed values

Observation

$$O_1 = \begin{matrix} \text{[yellow]} * \text{[light blue]} & + & \text{[yellow]} * \text{[green]} & + & \text{[yellow]} * \text{[green]} & + & \text{[yellow]} * \text{[blue]} & + & \text{[yellow]} * \text{[blue]} & + & \text{[yellow]} * \text{[blue]} \end{matrix}$$

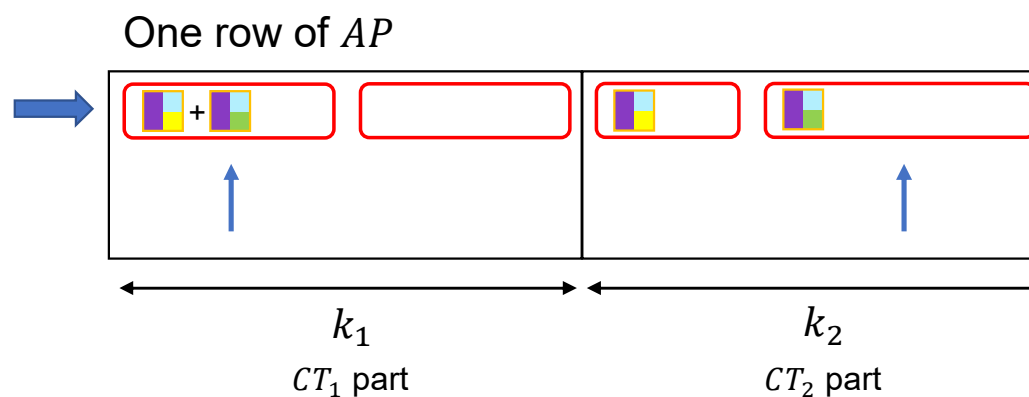
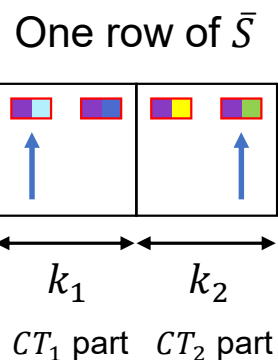
$$= (\text{[yellow]} + \text{[yellow]} + \text{[yellow]}) * \text{[light blue]} + (\text{[yellow]} + \text{[yellow]} + \text{[yellow]}) * \text{[blue]} + (\text{[yellow]} + \text{[yellow]}) * \text{[yellow]} + (\text{[yellow]} + \text{[yellow]} + \text{[yellow]} + \text{[yellow]}) * \text{[green]}$$

- Split value vectors into sum of compressed value vectors
- Merge coefficients

Aggregation algorithm

| CT ₁ | CT ₂ |
|-----------------|-----------------|
| 1 | 1 |
| 1 | 2 |
| 1 | 2 |
| 2 | 1 |
| 2 | 2 |
| 2 | 2 |

Each time select two scores from \bar{S} according to CT_1 and CT_2 , take the exponent of their sum and add to the same positions in AP



Step ③: Attention probability aggregation

Reform output into the form of weighted sum of compressed values

Observation

$$O_1 = \begin{matrix} \text{[yellow]} * \text{[light blue]} + \text{[yellow]} * \text{[green]} + \text{[yellow]} * \text{[green]} + \text{[yellow]} * \text{[blue]} + \text{[yellow]} * \text{[blue]} + \text{[yellow]} * \text{[blue]} \end{matrix}$$

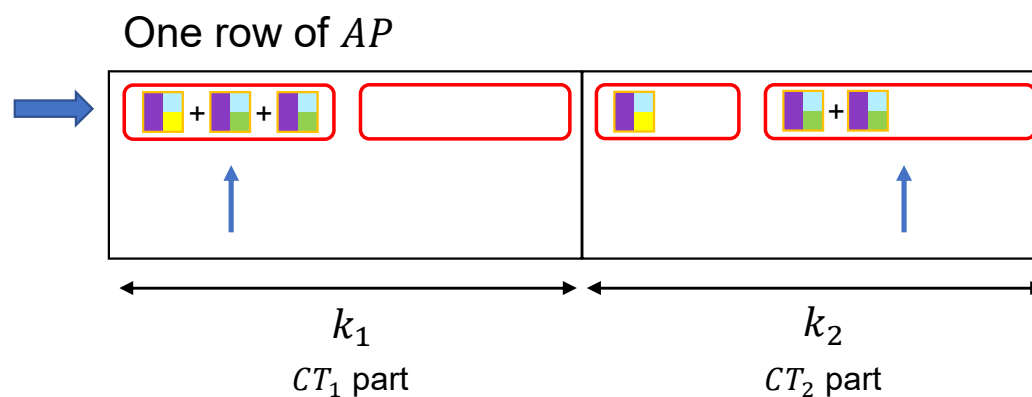
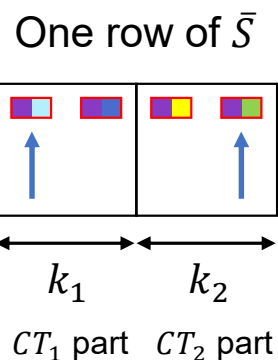
$$= (\text{[yellow]} + \text{[yellow]} + \text{[yellow]}) * \text{[light blue]} + (\text{[yellow]} + \text{[yellow]} + \text{[yellow]}) * \text{[blue]} + (\text{[yellow]} + \text{[yellow]}) * \text{[yellow]} + (\text{[yellow]} + \text{[yellow]} + \text{[yellow]} + \text{[yellow]}) * \text{[green]}$$

- Split value vectors into sum of compressed value vectors
- Merge coefficients

Aggregation algorithm

| CT ₁ | CT ₂ |
|-----------------|-----------------|
| 1 | 1 |
| 1 | 2 |
| 1 | 2 |
| 2 | 1 |
| 2 | 2 |
| 2 | 2 |

Each time select two scores from \bar{S} according to CT_1 and CT_2 , take the exponent of their sum and add to the same positions in AP



Step ③: Attention probability aggregation

Reform output into the form of weighted sum of compressed values

Observation

$$O_1 = \begin{matrix} \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[green]} & + & \text{[red]} * \text{[green]} & + & \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[green]} \end{matrix}$$

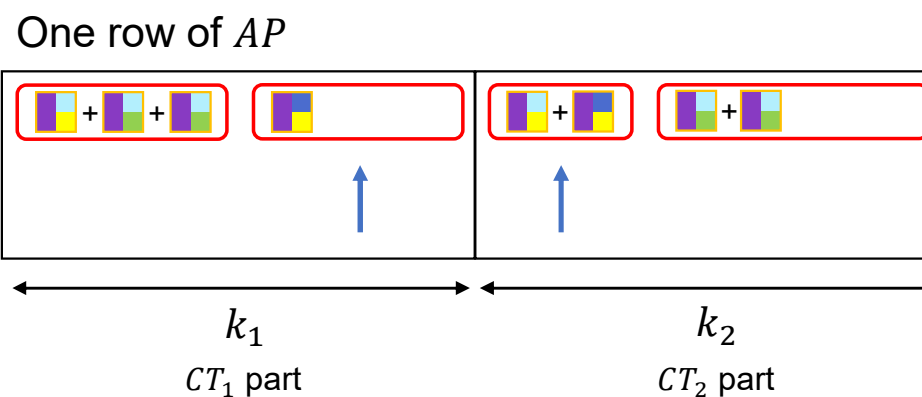
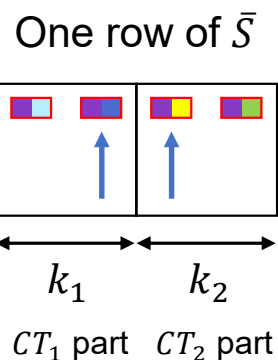
$$= (\text{[red]} + \text{[red]} + \text{[red]}) * \text{[blue]} + (\text{[red]} + \text{[red]} + \text{[red]}) * \text{[green]} + (\text{[red]} + \text{[red]}) * \text{[blue]} + (\text{[red]} + \text{[red]} + \text{[red]} + \text{[red]}) * \text{[green]}$$

- Split value vectors into sum of compressed value vectors
- Merge coefficients

Aggregation algorithm

| CT ₁ | CT ₂ |
|-----------------|-----------------|
| 1 | 1 |
| 1 | 2 |
| 1 | 2 |
| 2 | 1 |
| 2 | 2 |
| 2 | 2 |

Each time select two scores from \bar{S} according to CT_1 and CT_2 , take the exponent of their sum and add to the same positions in AP



Step ③: Attention probability aggregation

Reform output into the form of weighted sum of compressed values

Observation

$$O_1 = \begin{matrix} \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[green]} & + & \text{[red]} * \text{[green]} & + & \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[green]} \end{matrix}$$

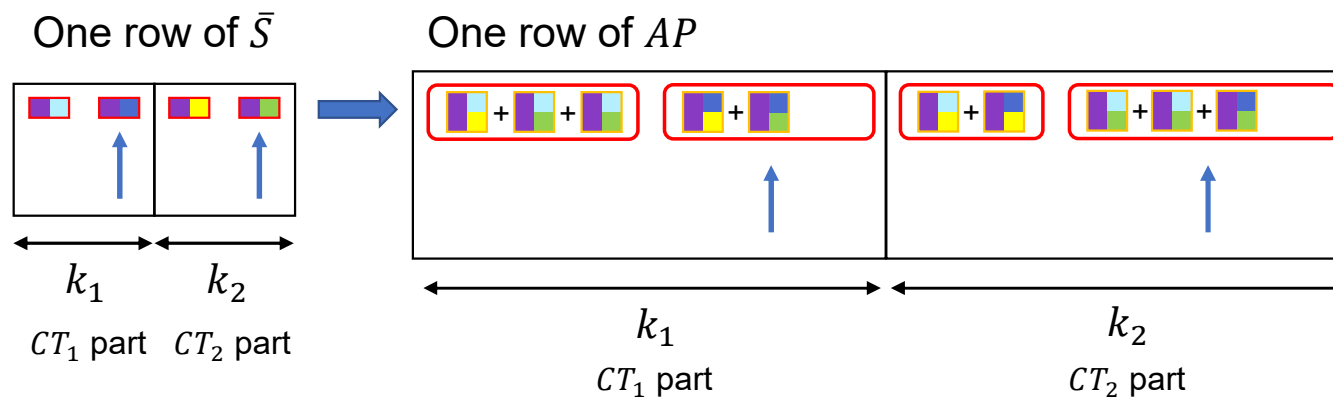
$$= (\text{[red]} + \text{[red]} + \text{[red]}) * \text{[blue]} + (\text{[red]} + \text{[red]} + \text{[red]}) * \text{[green]} + (\text{[red]} + \text{[red]}) * \text{[blue]} + (\text{[red]} + \text{[red]} + \text{[red]} + \text{[red]}) * \text{[green]}$$

- Split value vectors into sum of compressed value vectors
- Merge coefficients

Aggregation algorithm

| CT ₁ | CT ₂ |
|-----------------|-----------------|
| 1 | 1 |
| 1 | 2 |
| 1 | 2 |
| 2 | 1 |
| 2 | 2 |
| 2 | 2 |

Each time select two scores from \bar{S} according to CT_1 and CT_2 , take the exponent of their sum and add to the same positions in AP



Step ③: Attention probability aggregation

Reform output into the form of weighted sum of compressed values

Observation

$$O_1 = \begin{matrix} \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[green]} & + & \text{[red]} * \text{[green]} & + & \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[blue]} & + & \text{[red]} * \text{[green]} \end{matrix}$$

$$= (\text{[red]} + \text{[red]} + \text{[red]}) * \text{[blue]} + (\text{[red]} + \text{[red]} + \text{[red]}) * \text{[green]} + (\text{[red]} + \text{[red]}) * \text{[blue]} + (\text{[red]} + \text{[red]} + \text{[red]} + \text{[red]}) * \text{[green]}$$

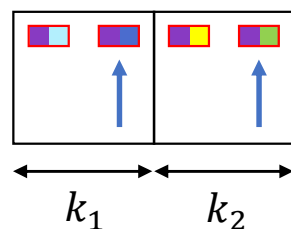
- Split value vectors into sum of compressed value vectors
- Merge coefficients

Aggregation algorithm

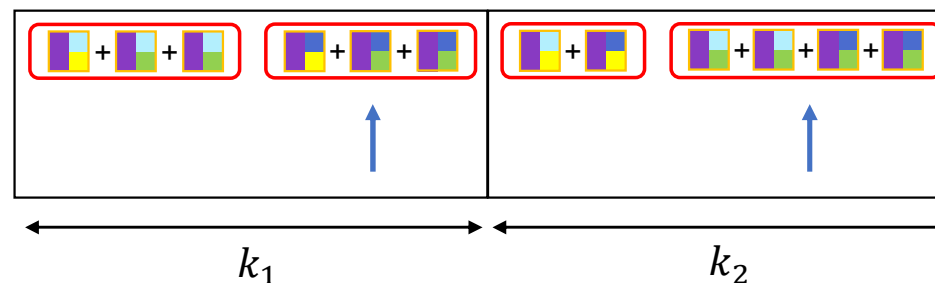
| CT ₁ | CT ₂ |
|-----------------|-----------------|
| 1 | 1 |
| 1 | 2 |
| 1 | 2 |
| 2 | 1 |
| 2 | 2 |
| 2 | 2 |
| 2 | 2 |

Each time select two scores from \bar{S} according to CT_1 and CT_2 , take the exponent of their sum and add to the same positions in AP

One row of \bar{S}



One row of AP

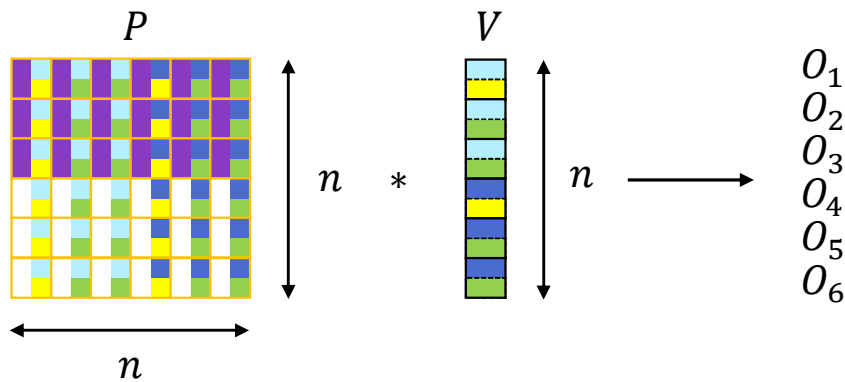


n exponent operations per row, the total number is reduced from n^2 to $k_0 n$

Step ④: Output weighting

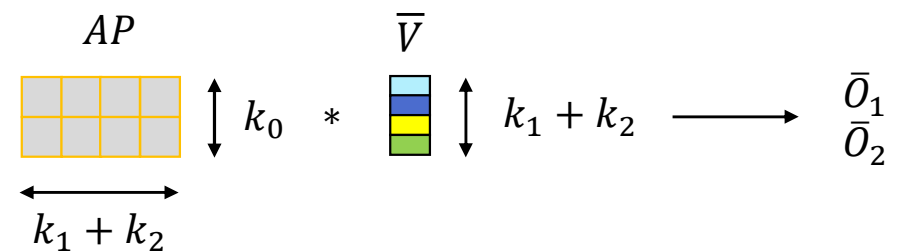
Given the aggregated attention probability matrix AP , output can be calculated as $AP \cdot \bar{V}$

Ordinary weighting



$$O = P \cdot V$$

CTA weighting



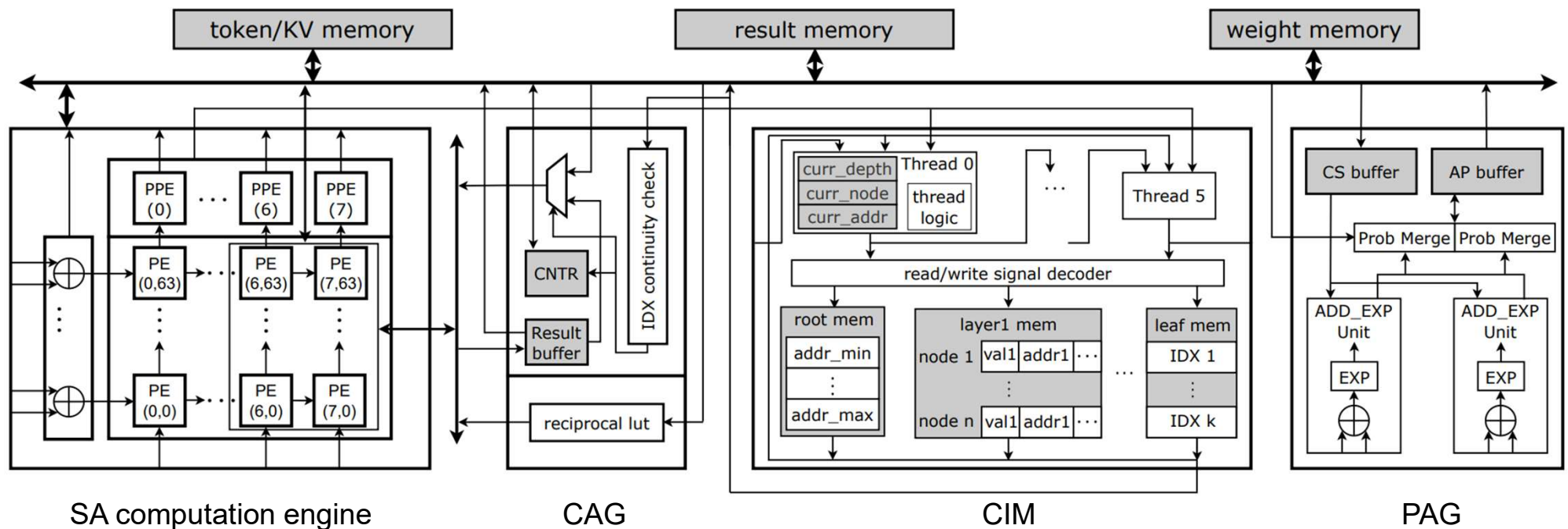
$$\bar{O} = AP \cdot \bar{V}$$

Computation is reduced from n^2d to $k_0(k_1 + k_2)d$

Architecture overview: Specialized accelerator for CTA

Four modules:

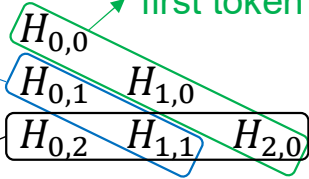
- Systolic array computation engine: accommodate major **matrix multiplication** operations
- CAG: averaging tokens assigned to the same cluster to **compute compressed tokens**
- CIM: **separate tokens** into clusters by assigning cluster index according to their hash values
- PAG: calculate **aggregated attention probabilities** from compressed attention scores



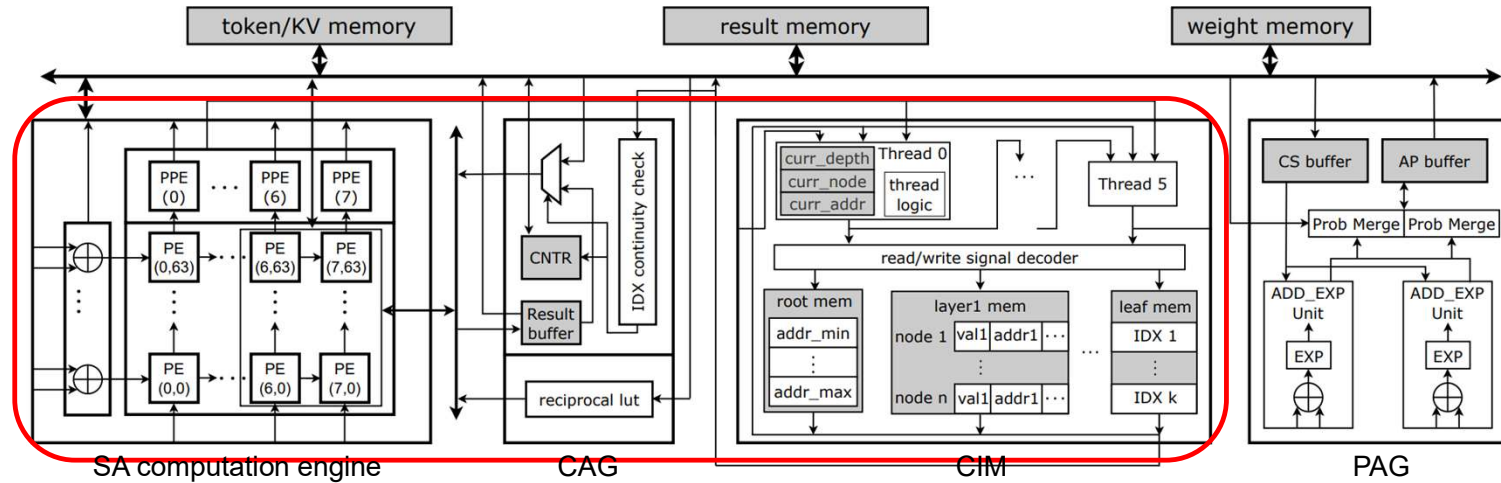
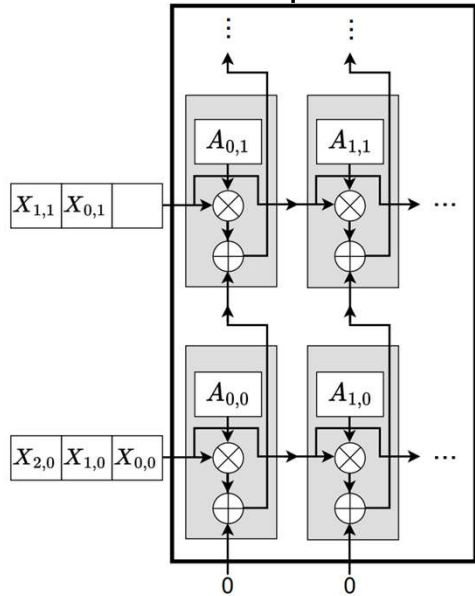
Dedicated accelerator for step ①

Hash values
of second
token

Hash values of
first token



Hash values
in one cycle



$$H = \lfloor (A \cdot X^T + B) / w \rfloor$$

- Each column of PEs accumulate the inner products between a direction vector and tokens.
- Each PPE adds bias to the result and divide it by w , then take floor.
- CIM separates tokens according to their hash values and CAG computes the average.

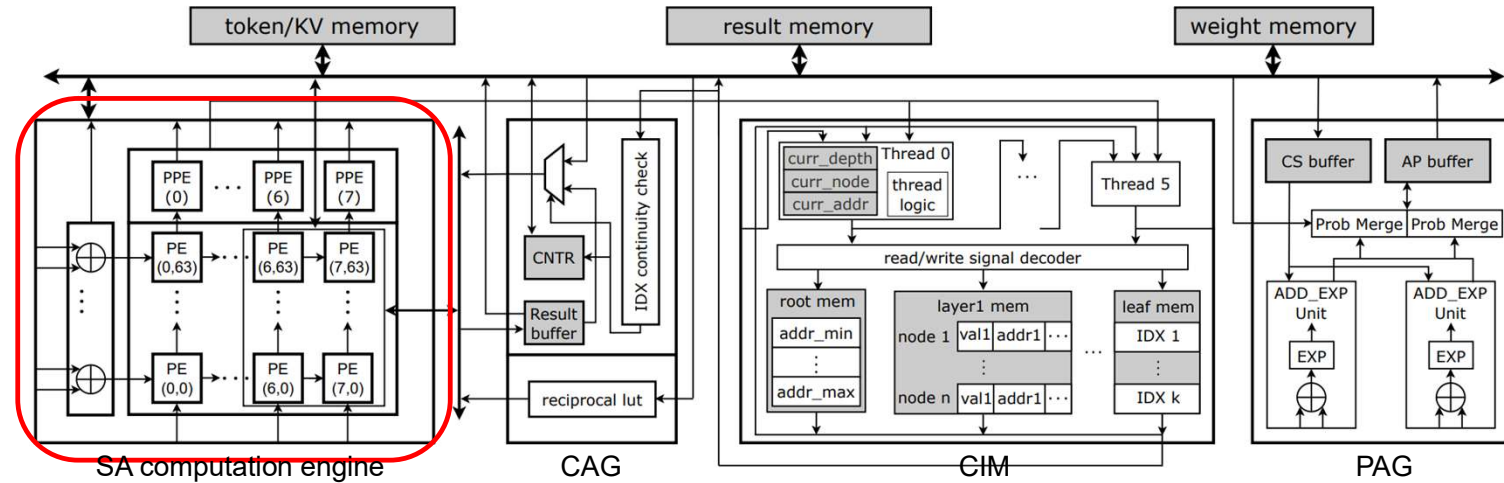
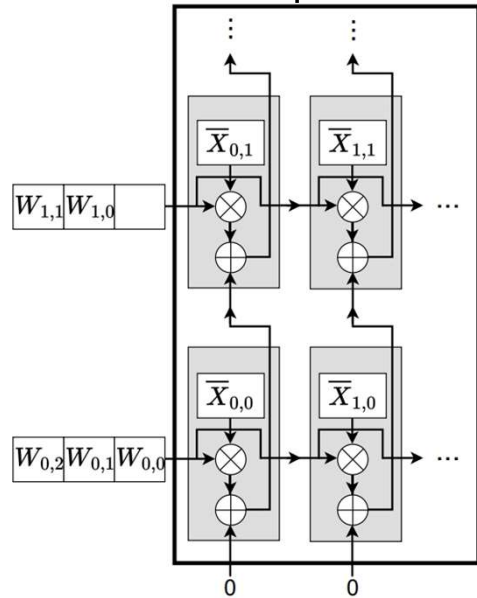
Dedicated accelerator for step ①

Query of first token

Or \bar{K}, \bar{V}

$\bar{Q}_{0,0}$
 $\bar{Q}_{0,1}$
 $\bar{Q}_{0,2}$

$\bar{Q}_{1,0}$
 $\bar{Q}_{1,1}$
 $\bar{Q}_{2,0}$



$$\bar{Q} = \bar{X}^Q \cdot W^Q, \bar{K} = \bar{X}^{KV} \cdot W^K, \bar{V} = \bar{X}^{KV} \cdot W^V$$

- Each column of PEs accumulate the inner products between a compressed token and weight matrix columns to compute the transformed query/key/value for this compressed token.

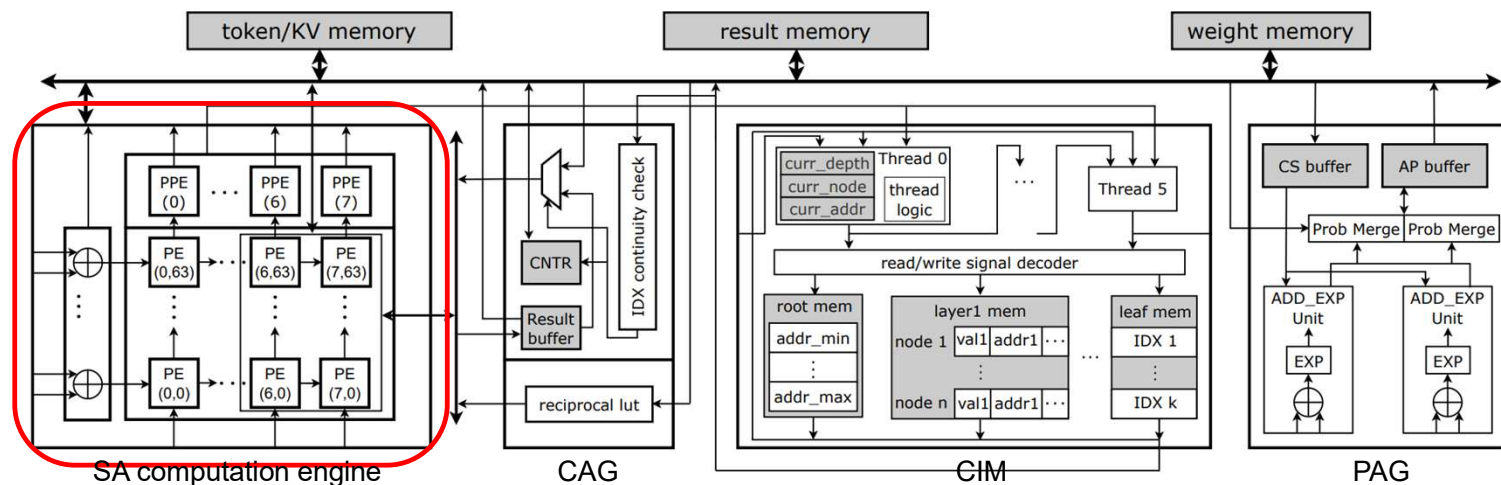
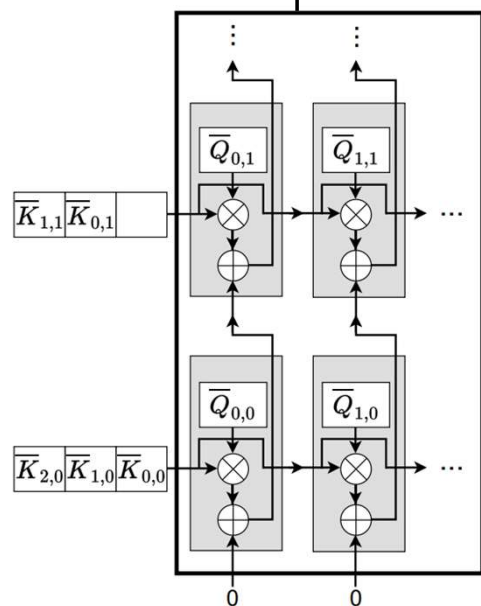
Dedicated accelerator for step ②

Similarity scores of first query

$\bar{S}_{0,0}$
 $\bar{S}_{0,1}$
 $\bar{S}_{0,2}$

$\bar{S}_{1,0}$
 $\bar{S}_{1,1}$

$\bar{S}_{2,0}$



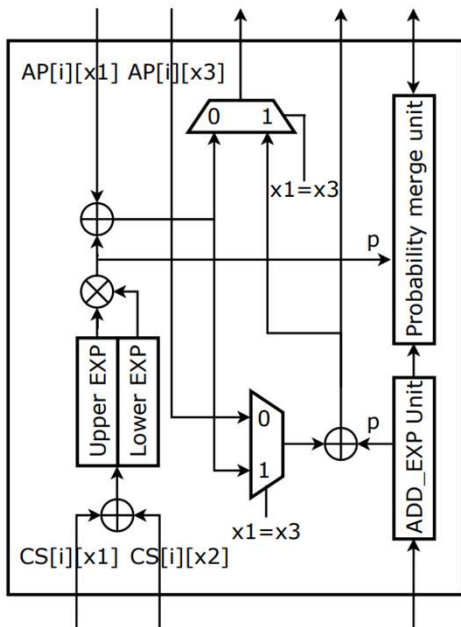
$$\bar{S} = (\bar{Q} \cdot \bar{K}^T) / \sqrt{d}$$

- Each column of PEs accumulate the inner products between a compressed query and compressed keys to compute the similarity scores of this query.
- Each PPE counts the largest first level scores (first k_1 scores), then subtract the maximum from the following scores.

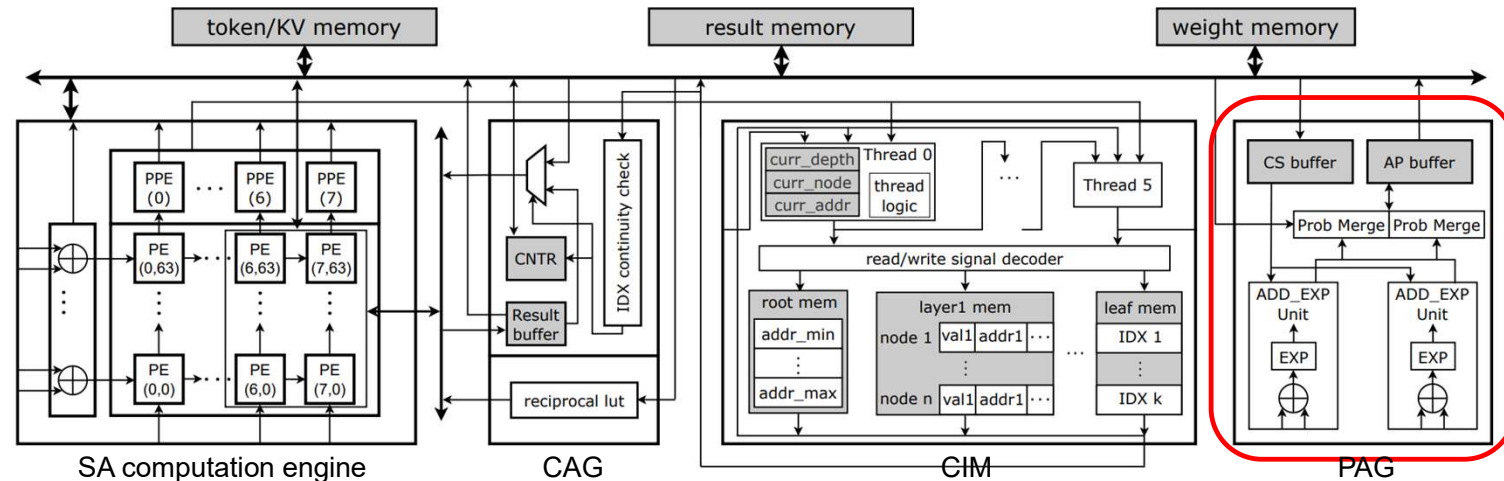
Dedicated accelerator for step ③

One tile of PAG

Aggregated probabilities



Compressed similarities

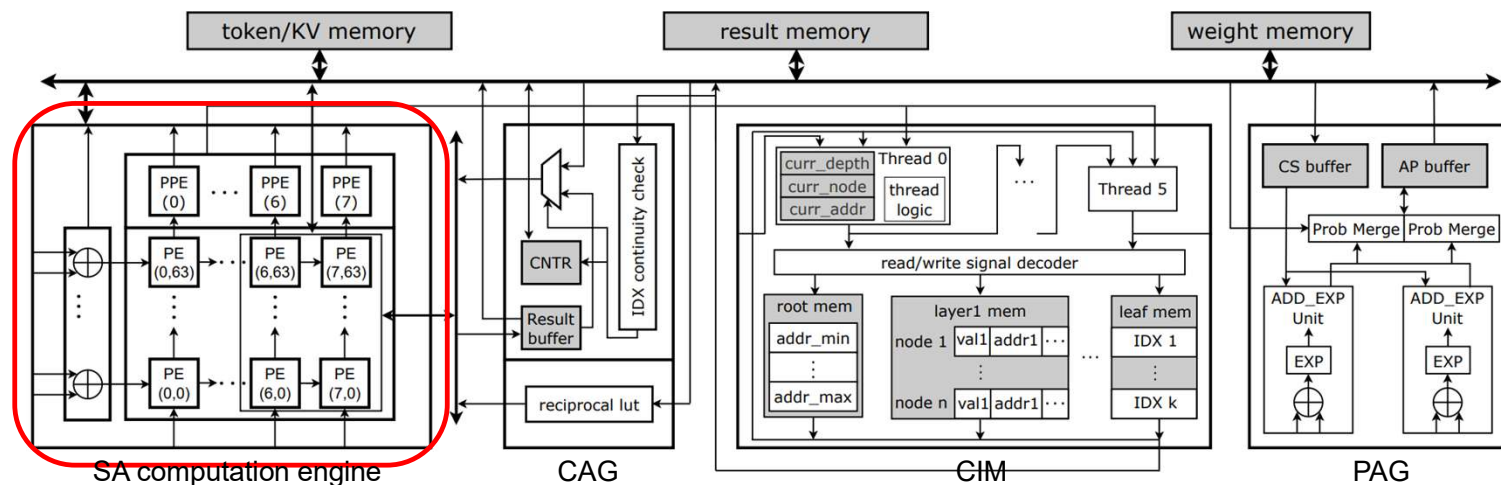
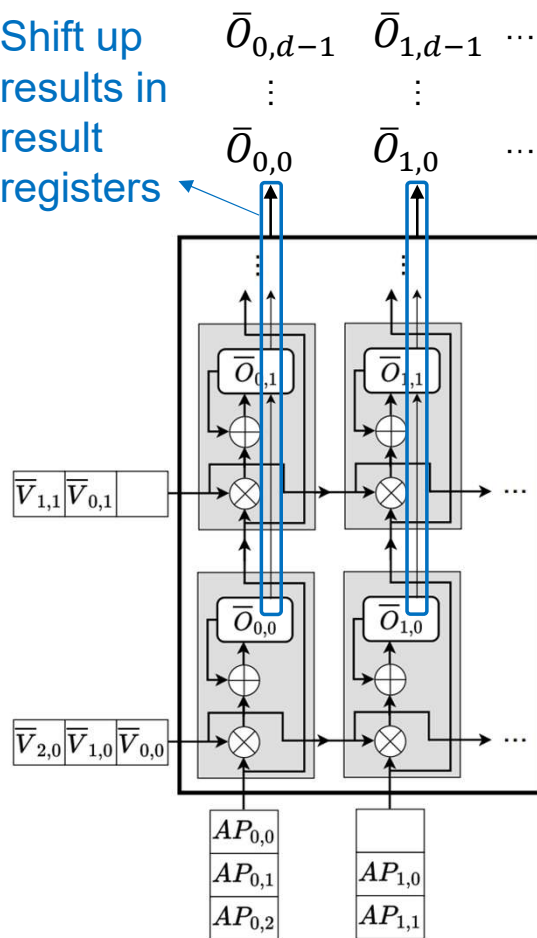


$$AP = PAG(\bar{S}, CT_1, CT_2)$$

- PAG reads pairs of compressed similarity scores indexed by CT_1 and CT_2 , sum the scores and take the exponent.
- PAG accumulates the result to the same positions of aggregated probabilities indexed by CT_1 and CT_2 .

Dedicated accelerator for step ④

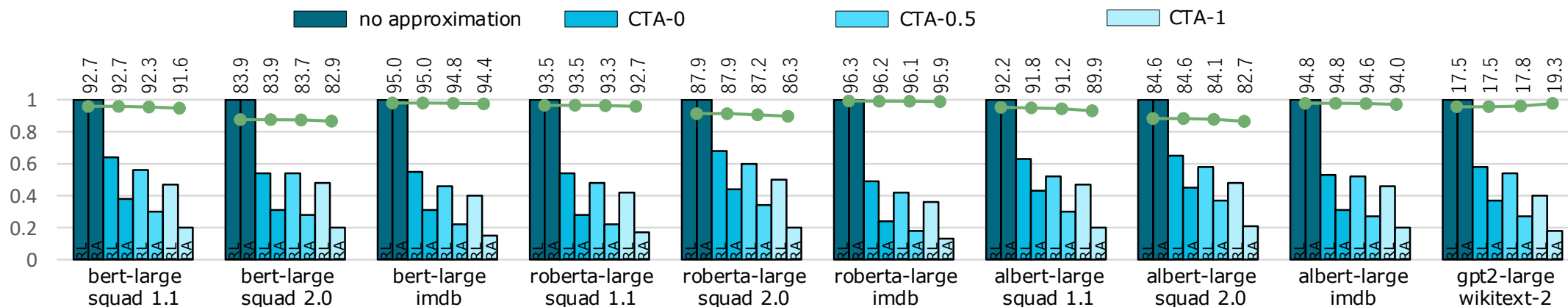
Shift up
results in
result
registers



$$\bar{O} = AP \cdot \bar{V}$$

- Each column of PEs accumulate an output vector in their result registers.
- Each PPE sums the aggregated probabilities, divides the results by half of the sum when shifting up the results.

Evaluation: accuracy and computation reduction



Green line: accuracy

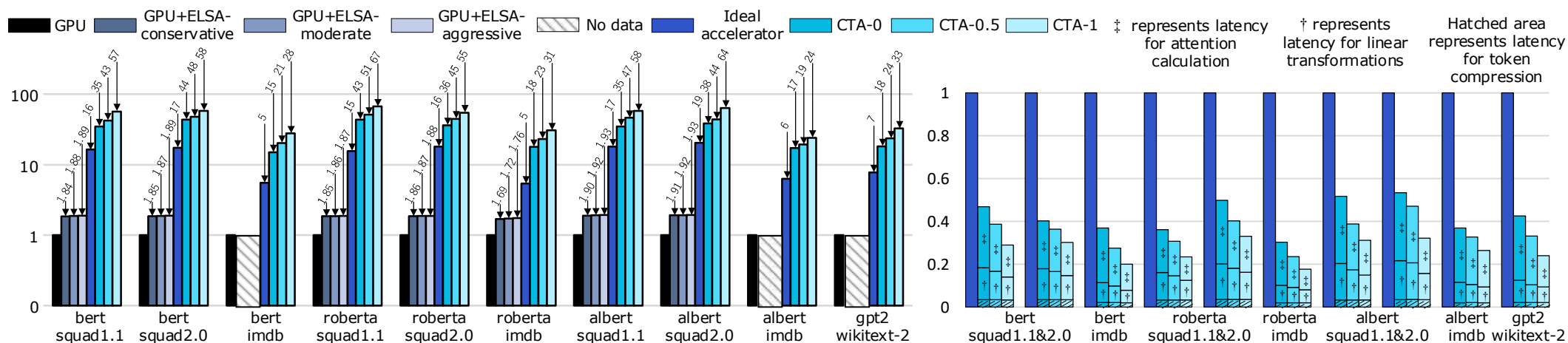
RL, RA: (computation after compression) / (computation before compression)

RL is for linear transformations (step1), RA is for the rest steps in attention (step2-4)

Average results

| | | | |
|----------|----------------------------|---|--------------------|
| CTA-0: | virtually no accuracy loss | & | RL=58.3%, RA=35.2% |
| CTA-0.5: | <0.5% accuracy loss | & | RL=52.2%, RA=27.5% |
| CTA-1: | <1% accuracy loss | & | RL=44.4%, RA=18.4% |

Evaluation: throughput and latency



Average results

GPU: V100-SXM2 32GB

ELSA: previous state-of-the-art attention accelerator

speedup over GPU

CTA-0: 27.7×
CTA-0.5: 33.8×
CTA-1: 44.2×

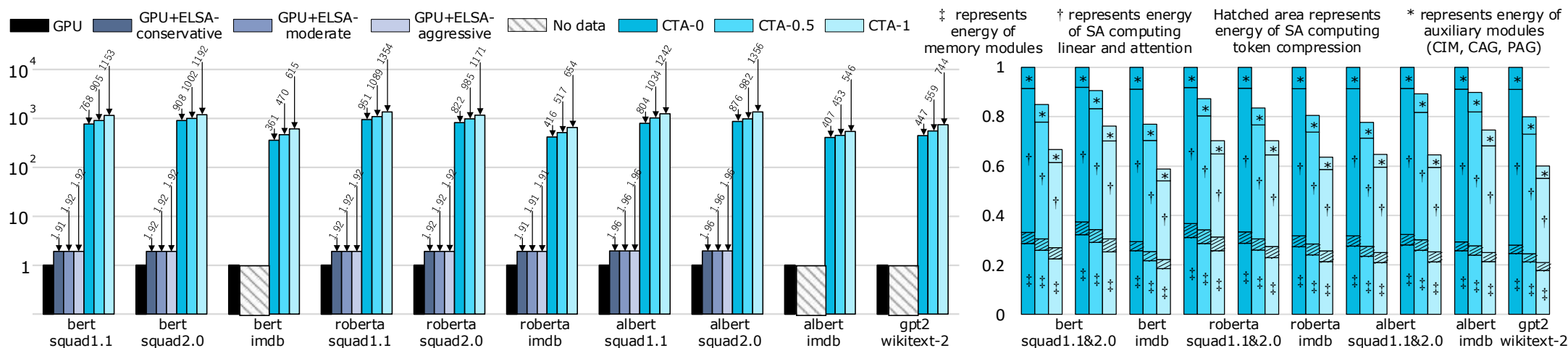
speedup over ELSA-aggressive+GPU

CTA-0: 18.3×
CTA-0.5: 22.1×
CTA-1: 28.7×

latency breakdown

attention computation: 59%
linear computation: 34%
token compression: 7%

Evaluation: energy efficiency and breakdown



Average results

GPU: V100-SXM2 32GB

ELSA: previous state-of-the-art attention accelerator

efficiency over GPU

CTA-0: 634×
CTA-0.5: 756×
CTA-1: 950×

efficiency over ELSA-aggressive+GPU

CTA-0: 399×
CTA-0.5: 471×
CTA-1: 587×

energy breakdown

Memory modules: 29%
SA module: 62%
Others: 9%

Energy spent on hashing for token compression: 5%

Conclusion

- CTA efficiently removes token semantic feature repetition
 - Token compression and attention operations on compressed tokens
- CTA keeps matrix multiplication formulation of major computation stage
 - Preserve inherent parallelism of attention mechanism
 - Support acceleration with more general architecture
- CTA architecture
 - Reconfigurable systolic array-based computation engine
 - Light-weight auxiliary modules

Thank you!