# Acceleration

Nikita Kiselev

Deep Learning, Intelligent Systems
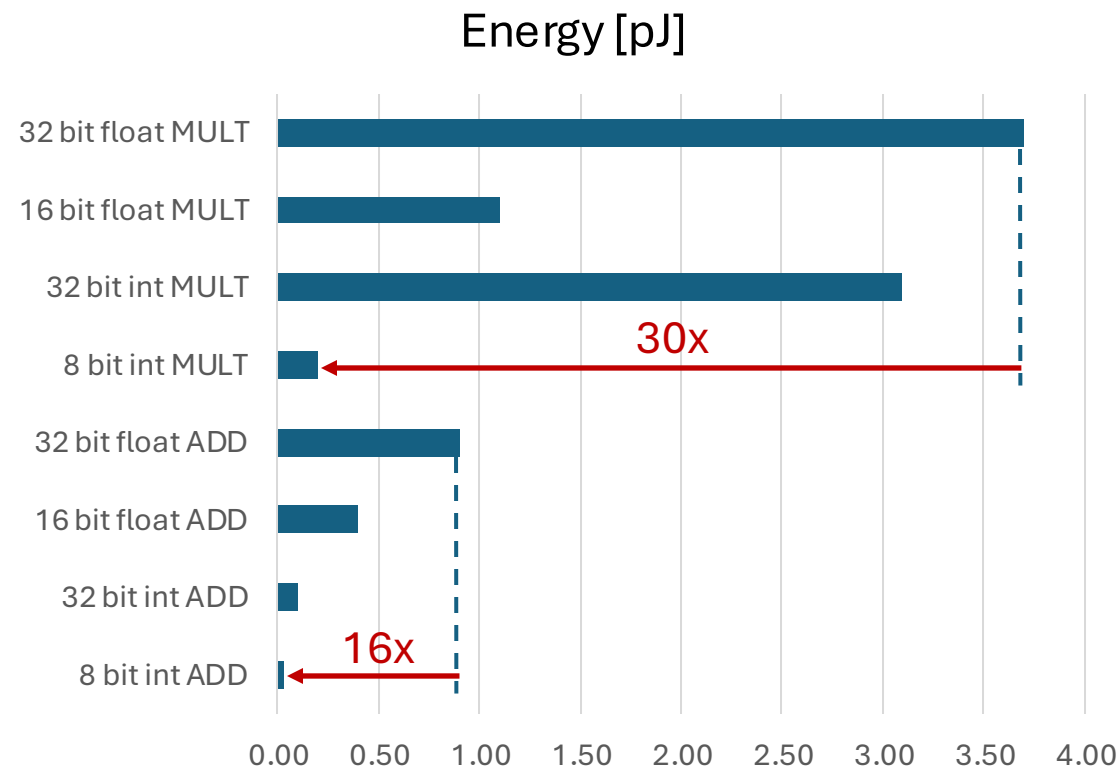
December 16, 2025

# Introduction

How to close this gap?



Energy [pJ]

| | |
|---|---|
| 32 bit float MULT | |
| 16 bit float MULT | |
| 32 bit int MULT | |
| 8 bit int MULT | 30x |
| 32 bit float ADD | |
| 16 bit float ADD | |
| 32 bit int ADD | |
| 8 bit int ADD | 16x |

**We should make deep learning more efficient…**

Xiao G. et al. Smoothquant: Accurate and efficient post-training quantization for large language models, 2023

# Contents

- Quantization

- Pruning

- Distillation

- KV-Cache

- Flash Attention

# Quantization

# Numerical Data Types: Integer

## Unsigned Integer
- $n$-bit range: $[0, 2^n - 1]$



$$2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = 49$$

## Signed Integer
- $n$-bit range: $[-2^{n-1} + 1, 2^{n-1} - 1]$
- Both 000…00 and 100…00 represent 0

**Sign Bit**



$$- \quad 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = -49$$

# Fixed-Point Number



Integer . Fraction

"Decimal" Point



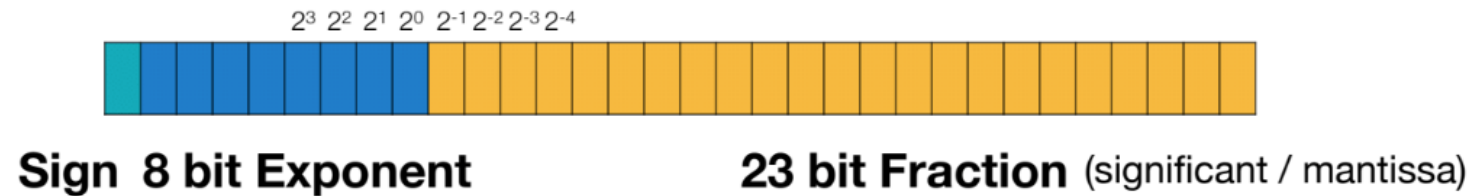$$-2^3 + 2^2 + 2^1 + 2^0 + 2^{-1} + 2^{-2} + 2^{-3} + 2^{-4} = 3.0625$$



$$( -2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 ) \times 2^{-4} = 49 \times 0.0625 = 3.0625$$

# Floating-Point Number

Example: 32-bit floating-point number in IEEE 754

$$2^3 \; 2^2 \; 2^1 \; 2^0 \; 2^{-1} 2^{-2} 2^{-3} 2^{-4}$$

**Sign  8 bit Exponent**                           **23 bit Fraction** (significant / mantissa)

$$(-1)^{sign} \times (1 + \textbf{Fraction}) \times 2^{\textbf{Exponent}-127} \quad \longleftarrow \quad \textbf{Exponent Bias} = 127 = 2^{8-1}-1$$

How to represent **0.265625**?

$$\textbf{0.265625} = 1.0625 \times 2^{-2} = (1 + \underline{0.0625}) \times 2^{\underline{125}-127}$$

| 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

125                           0.0625

# Floating-Point Number

Exponent Width → Range; Fraction Width → Precision

| | Exponent (bits) | Fraction (bits) | Total (bits) |
|---|---|---|---|
| **IEEE 754 Single Precision 32-bit Float (IEEE FP32)** | | | |
| | 8 | 23 | 32 |
| **IEEE 754 Half Precision 16-bit Float (IEEE FP16)** | | | |
| | 5 | 10 | 16 |
| **Google Brain Float (BF16)** | | | |
| | 8 | 7 | 16 |

# Dynamic Range and Precision



Maarten Grootendorst. A Visual Guide to Quantization, 2024
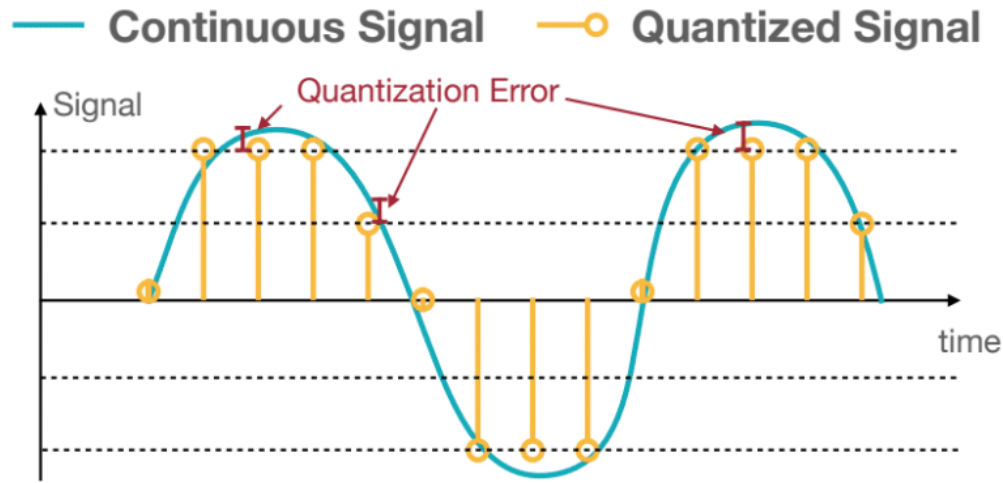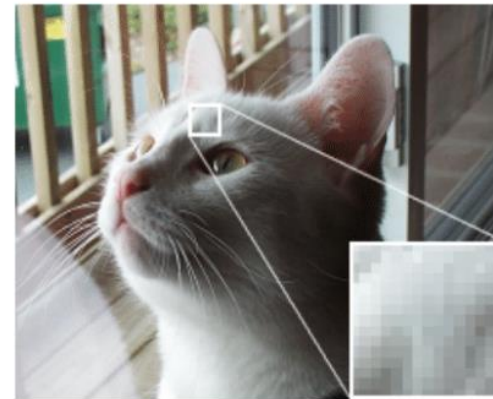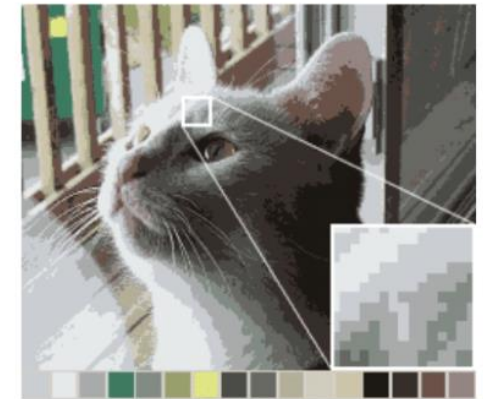
# What is Quantization?

*Quantization is the process of constraining an input from a continuous or otherwise large set of values to a discrete set*
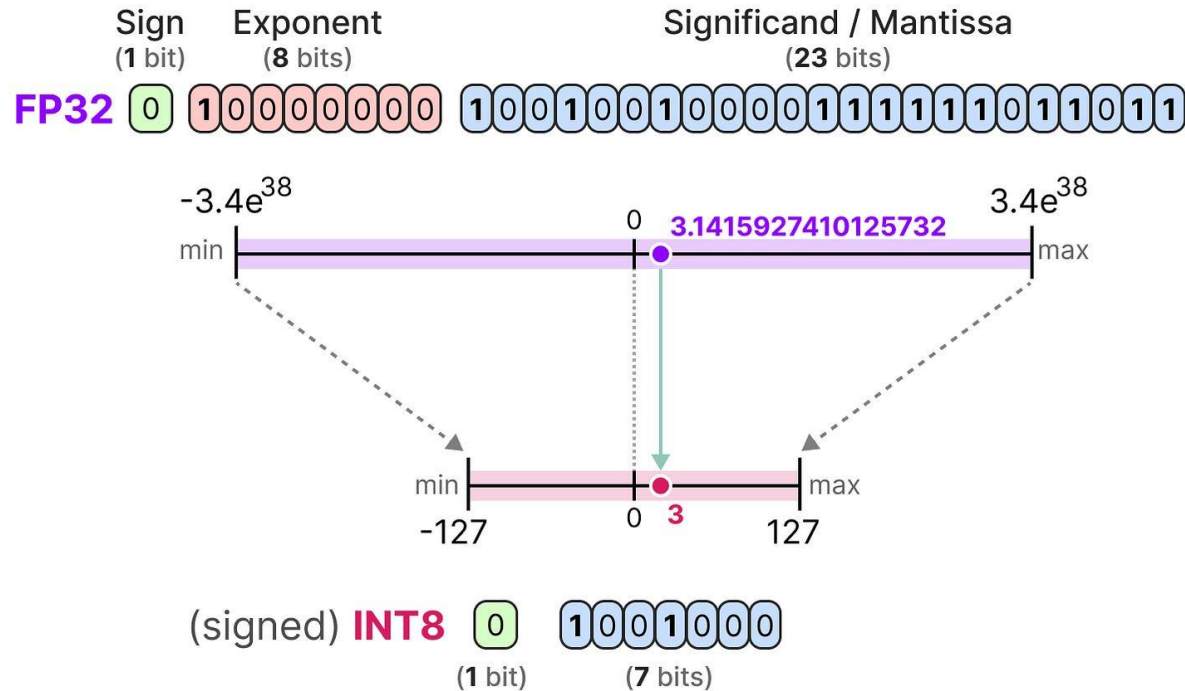
"Palettization"

Sing Han et al. TinyML and Efficient Deep Learning Computing

https://efficientml.ai

# Model Weights Quantization



Sign
(**1** bit)

Exponent
(**8** bits)

Significand / Mantissa
(**23** bits)

**FP32** 0 10000000 10010010000011111110011011

-3.4e$^{38}$

min

0 **3.1415927410125732**

max

3.4e$^{38}$

min -127

0 **3**

max 127

(signed) **INT8** 0 1001000

(**1** bit) (**7** bits)

Depending on the hardware, integer-based calculations might be faster that floating-point calculations but this isn't always the case. However, computations are generally faster when using fewer bits.

$$memory = \frac{(\#\text{bits per number})}{8} \times (\#params)$$

**64-bits** $= \dfrac{64}{8} \times$ 70B $\approx$ **560** GB

$\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots$

**32-bits** $= \dfrac{32}{8} \times$ 70B $\approx$ **280** GB

$\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots$

**16-bits** $= \dfrac{16}{8} \times$ 70B $\approx$ **140** GB
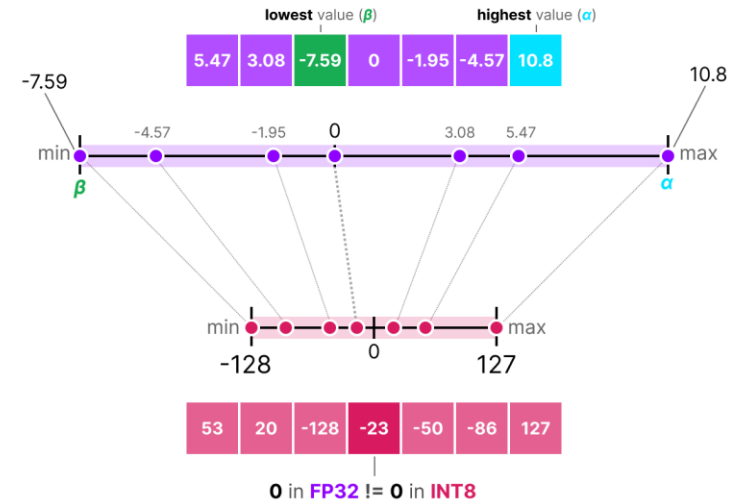
# Quantization Formula

## Symmetric



$$s = \frac{2^{b-1} - 1}{\alpha}$$

$$x_{\text{quantized}} = \text{round}(s \cdot x)$$

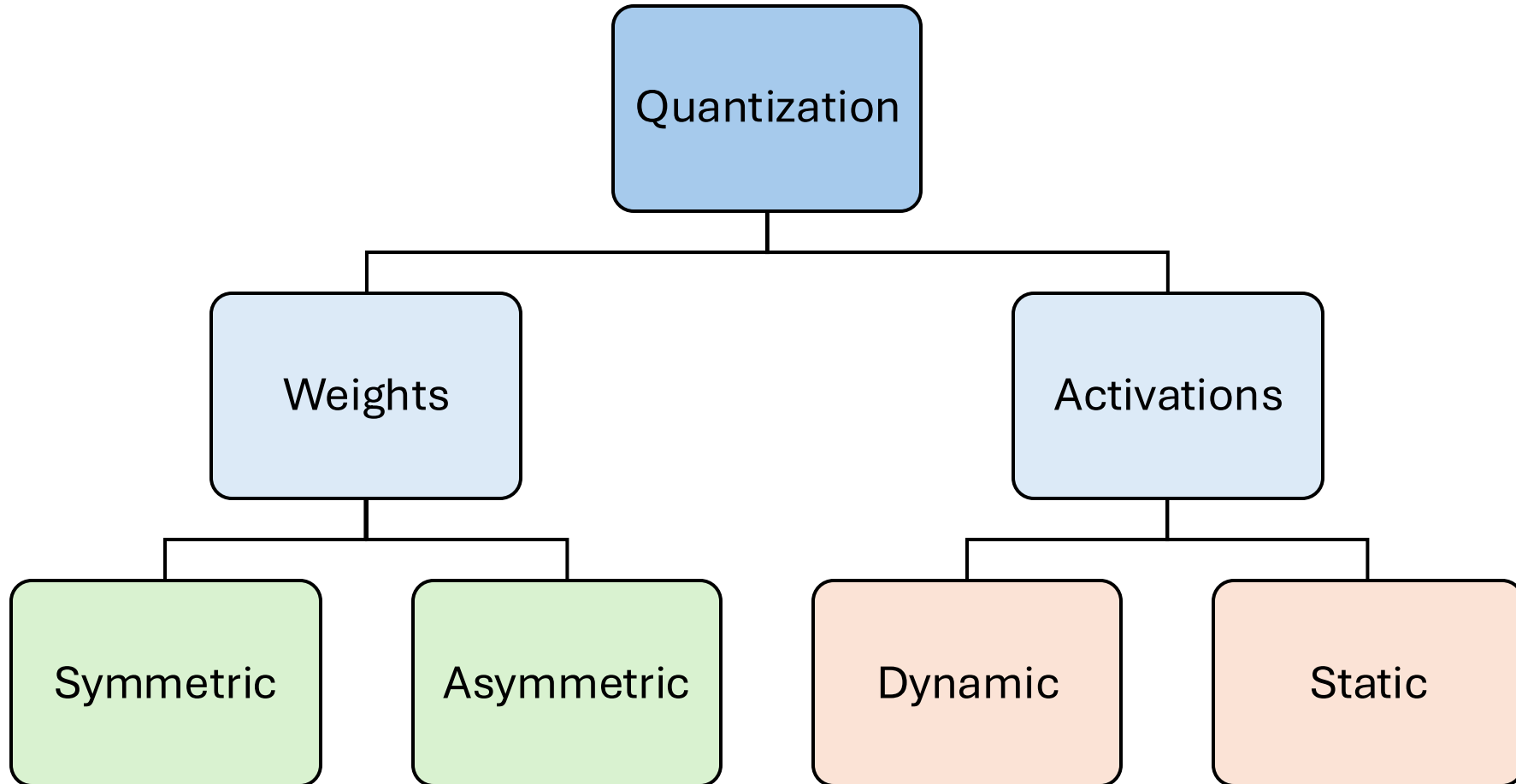$$x_{\text{dequantized}} = \frac{x_{\text{quantized}}}{s}$$

## Asymmetric



$$s = \frac{2^{b-1} - (-2^{b-1} + 1)}{\alpha - \beta}$$

$$z = \text{round}(-s \cdot \beta) - 2^{b-1}$$

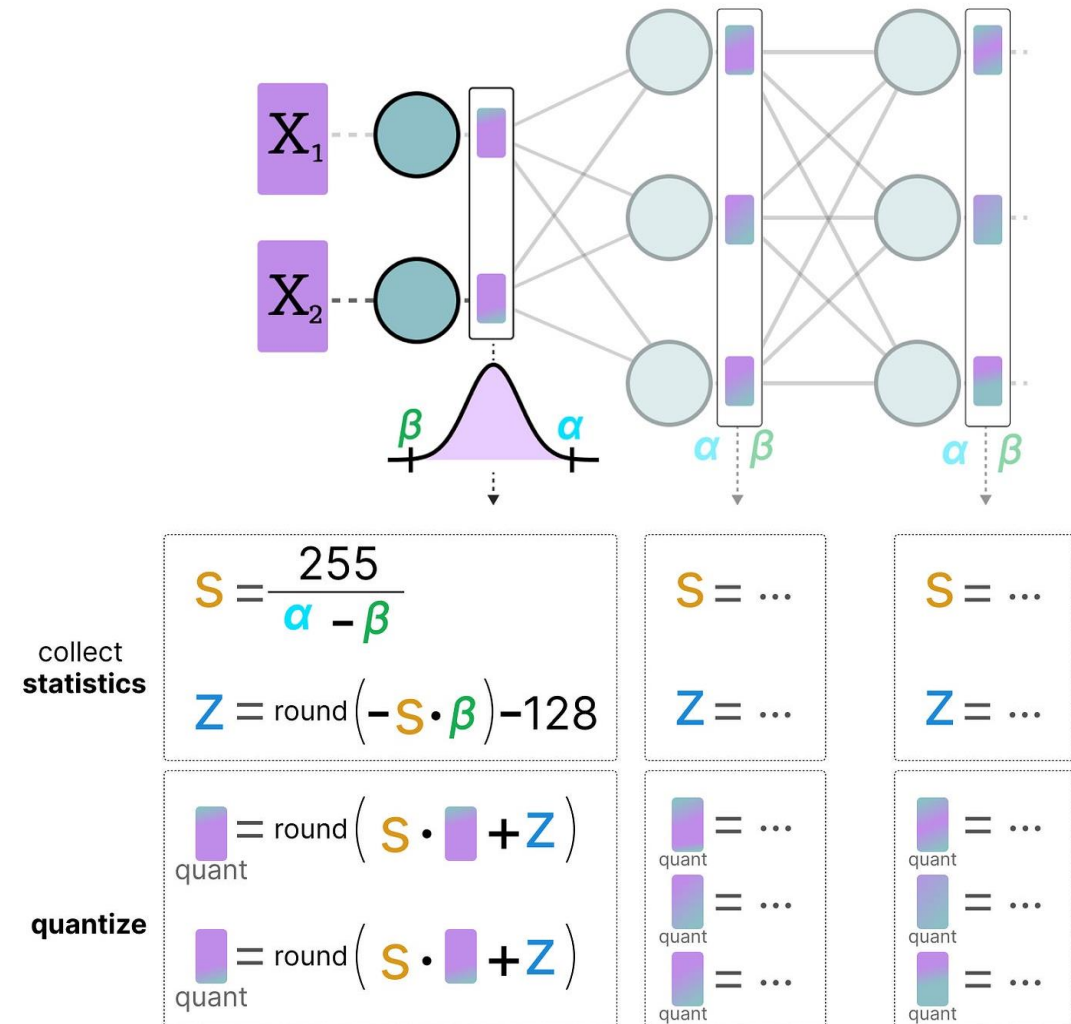$$x_{\text{quantized}} = \text{round}(s \cdot x + z)$$

$$x_{\text{dequantized}} = \frac{x_{\text{quantized}} - z}{s}$$
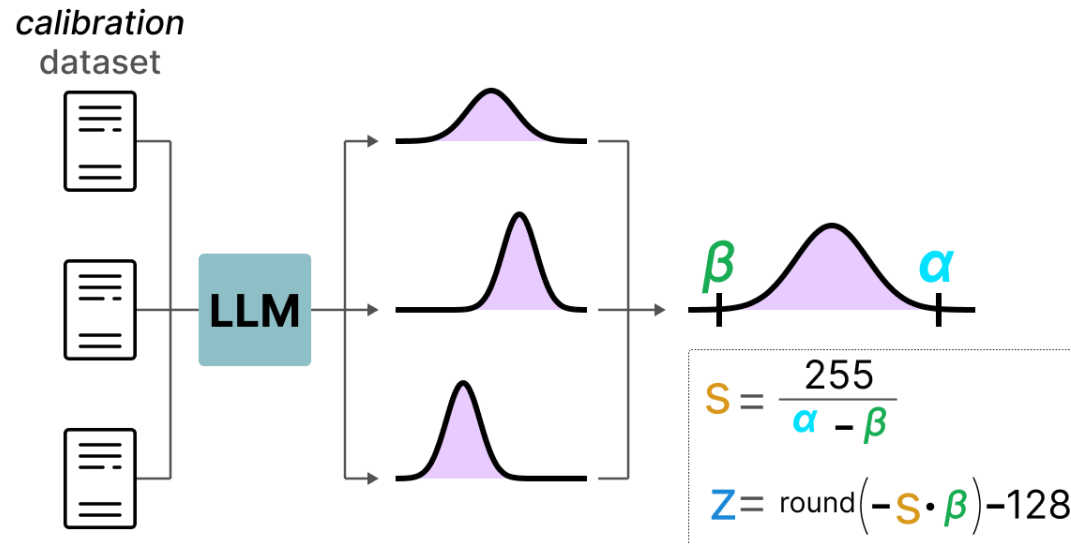
# Post-Training Quantization

# Dynamic Quantization

1. After data passes a hidden layer, its activations are collected

2. This distribution of activations is then used to calculate the zeropoint and scale factor

3. The process is repeated each time data passes through new later



$$S = \frac{255}{\alpha - \beta}$$

$$Z = \text{round}\left(-S \cdot \beta\right) - 128$$

collect statistics

quantize

$$\blacksquare_{\text{quant}} = \text{round}\left(S \cdot \blacksquare + Z\right)$$

$$\blacksquare_{\text{quant}} = \text{round}\left(S \cdot \blacksquare + Z\right)$$

$$S = \dots$$

$$Z = \dots$$

$$\blacksquare_{\text{quant}} = \dots$$

$$\blacksquare_{\text{quant}} = \dots$$

$$\blacksquare_{\text{quant}} = \dots$$

$$S = \dots$$

$$Z = \dots$$

$$\blacksquare_{\text{quant}} = \dots$$

$$\blacksquare_{\text{quant}} = \dots$$

$$\blacksquare_{\text{quant}} = \dots$$

# Static Quantization



calibration dataset

$S = \dfrac{255}{\alpha - \beta}$

$Z = \text{round}(-S \cdot \beta) - 128$

1. Use **calibration dataset** to collect these potential distributions

2. Dynamic quantization is more accurate, but increase compute time

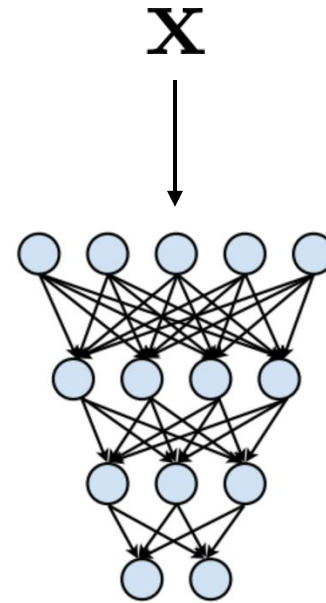3. In contrast, static quantization is less accurate, but is faster

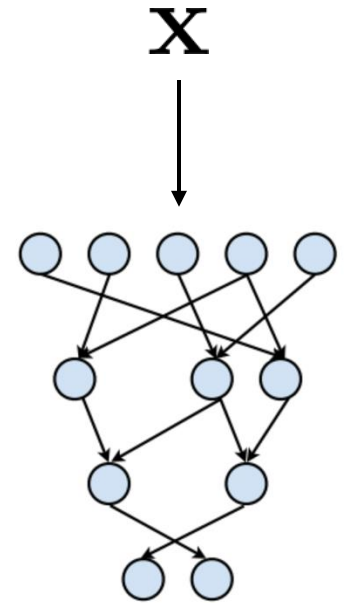# Pruning

# Neural Network Pruning

$$\arg\min_{\boldsymbol{\theta}_p} \mathcal{L}_{\boldsymbol{\theta}_p}(\mathbf{x})$$

$$s.t. \|\boldsymbol{\theta}_p\|_0 \leq N$$

- $\mathcal{L}$ represents objective function for neural network training
- $\mathbf{x}$ is input, $\boldsymbol{\theta}$ is original weights, $\boldsymbol{\theta_p}$ is pruned weights
- $\|\boldsymbol{\theta_p}\|_0$ calculates the #nonzeros in $\boldsymbol{\theta_p}$, and N is the target of #nonzeros



$$\arg\min_{\boldsymbol{\theta}} \mathcal{L}_{\boldsymbol{\theta}}(\mathbf{x})$$

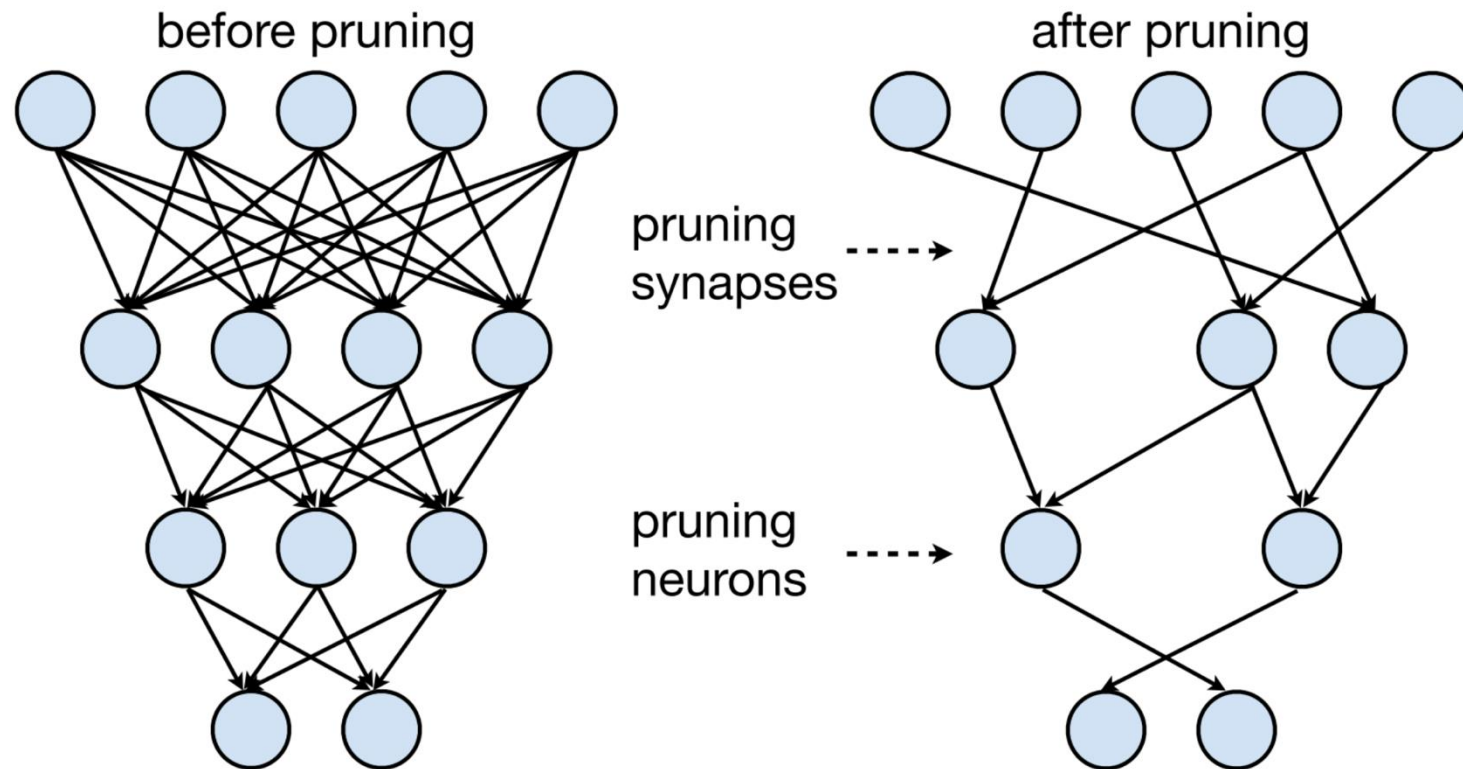

$$\arg\min_{\boldsymbol{\theta}_p} \mathcal{L}_{\boldsymbol{\theta}_p}(\mathbf{x})$$

$$s.t. \|\boldsymbol{\theta}_p\|_0 \leq N$$

Han S. et al. Learning both weights and connections for efficient neural network, 2015

# Neural Network Pruning

Make neural network smaller by removing synapses and neurons



Han S. et al. Learning both weights and connections for efficient neural network, 2015

# Neural Network Pruning



Train Connectivity

Han S. et al. Learning both weights and connections for efficient neural network, 2015

# Neural Network Pruning



Train Connectivity

⇩

Prune Connections

Pruning

Han S. et al. Learning both weights and connections for efficient neural network, 2015

# Neural Network Pruning



Train Connectivity

Prune Connections

Train Weights

Han S. et al. Learning both weights and connections for efficient neural network, 2015

# Neural Network Pruning



Train Connectivity → Prune Connections → Train Weights

Pruning · Pruning+Finetuing · Iterative Pruning and Finetuing

Accuracy Loss vs Pruning Ratio (Parameters Pruned Away)

Han S. et al. Learning both weights and connections for efficient neural network, 2015

# Pruning at Different Granularities

A simple example of 2D weight matrix



**Fine-grained / Unstructured**

- More flexible pruning index choice
- Hard to accelerate (irregular)

**Coarse-grained / Structured**

- Less flexible pruning index choice (a subset of the fine-grained case)
- Easy to accelerate (just a smaller matrix!

# Pruning at Different Granularities

The case of convolutional layers



Fine-grained Pruning  Pattern-based Pruning  Vector-level Pruning  Kernel-level Pruning  Channel-level Pruning

https://efficientml.ai

# Pruning Criterions

# Distillation

# Challenge: limited hardware resources



| | Cloud AI | Tiny AI |
|---|---|---|
| Computation (fp32) | 19.5 TFLOPS | MFLOPs |
| Memory | 80GB | 256kB |
| Neural Network | ResNet<br>ViT-Large<br>… | MCUNet<br>MobileNetV2-Tiny<br>… |

**Neural network must be tiny to run efficiently on tiny edge devices.**
**How to train tiny model with the help of large model?**

# Tiny models are hard to train

Tiny models underfit large datasets, how to help them...?



Training curve for ResNet50

Training curve for MobileNetV2-Tiny

# Illustration of Knowledge Distillation

# Intuition of Knowledge Distillation

Matching prediction probabilities between teacher and student



| | | Logits | Probabilities |
|---|---|---|---|
| Cat | | 5 | 0.982 |
| Dog | | 1 | 0.017 |

| | | Logits | Probabilities |
|---|---|---|---|
| Cat | | 3 | 0.731 |
| Dog | | 2 | 0.269 |

The student model is less confident

# Intuition of Knowledge Distillation

Concept of temperature

$$\frac{\exp(z_i)}{\sum_{j=1}^{K} \exp(z_j)}$$



| | Logits | Probabilities (T=1) | Probabilities (T=10) |
|---|---|---|---|
| Cat | 5 | 0.982 | 0.599 |
| Dog | 1 | 0.017 | 0.401 |

$$\frac{\exp(z_i/T)}{\sum_{j=1}^{K} \exp(z_j/T)}$$

# Formal Definition of KD

Neural networks typically use a softmax function to generate the **logits** $z_i$ to class **probabilities** $p(z_i, T) = \dfrac{\exp(z_i/T)}{\sum_{j=1}^{k} \exp(z_j/T)}$

Temperature is normally set to 1

The goal of knowledge distillation is to **align the class probability distribution from teacher and student networks**

# What to match?

1. Output logits
2. Intermediate weights
3. Intermediate features

# Matching output logits



Teacher Model

Input → Layer 1 → Layer 2 → Layer N → Logits

Student Model

Input → Layer 1 → Layer 2 → Layer N → Logits → Classification Loss

Logits → Distillation Loss

Cross Entropy Loss
$$\mathbb{E}(-p_t \log p_s)$$

L2 Loss
$$\mathbb{E}\|p_t - p_s\|_2^2$$

# What to match?

1. Output logits

2. Intermediate weights

3. Intermediate features

# Matching intermediate weights

# What to match?

1. Output logits
2. Intermediate weights
3. Intermediate features

# Matching intermediate features



Intuition: teacher and student networks should have similar feature distributions, not just output probability distributions

# TIME FOR A BREAK

# KV-Cache

# Preliminaries



- GPT uses a Transformer Decoder that generates text autoregressively – one token at a time
- Each step takes all previous tokens as input and predicts the next token's probability

**Problem:** without optimization, GPT must recompute attention for all previous tokens at every step

# Transformer Decoder



**Self-Attention** (masked): tokens look at the previous tokens
queries, keys, values are computed from hidden states

# Transformer Decoder



**Feed-Forward**: after taking information from other tokens, take a moment to think and process this information

**Self-Attention** (masked): tokens look at the previous tokens queries, keys, values are computed from hidden states

# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



$Q$     $K^T$     $QK^T$     $V$     $O$

Softmax

[**seq_len**, emb_size]    [emb_size, **seq_len**]    [**seq_len, seq_len**]    [**seq_len**, emb_size]    [**seq_len**, emb_size]

# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$Q$ $K^T$ $QK^T$ $V$ $O$

Softmax

$\times$ $=$ $\times$ $=$

$[1, \text{emb\_size}]$ $[\text{emb\_size}, 1]$ $[1, 1]$ $[1, \text{emb\_size}]$ $[1, \text{emb\_size}]$

<bos> Transformers

# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$Q$        $K^T$        $QK^T$        $V$        $O$

Softmax $\left[ \quad \times \quad = \quad \right] \times \quad =$

$[2, \text{emb\_size}]$    $[\text{emb\_size}, 2]$    $[2, 2]$    $[2, \text{emb\_size}]$    $[2, \text{emb\_size}]$

<bos>   Transformers   predict

# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



$Q$      $K^T$      $QK^T$      $V$      $O$

Softmax

$\times$      $=$      $\times$      $=$

$[3, \text{emb\_size}]$      $[\text{emb\_size}, 3]$      $[3, 3]$      $[3, \text{emb\_size}]$      $[3, \text{emb\_size}]$

\<bos\>    Transformers    predict    tomorrow's

# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



Softmax

$Q$      $K^T$      $QK^T$      $V$      $O$

$[4, \text{emb\_size}]$    $[\text{emb\_size}, 4]$    $[4, 4]$    $[4, \text{emb\_size}]$    $[4, \text{emb\_size}]$

\<bos\>   Transformers   predict   tomorrow's   words

# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



Softmax

$Q$      $K^T$      $QK^T$      $V$      $O$

$[5, \text{emb\_size}]$      $[\text{emb\_size}, 5]$      $[5, 5]$      $[5, \text{emb\_size}]$      $[5, \text{emb\_size}]$

\<bos\>    Transformers    predict    tomorrow's    words    from
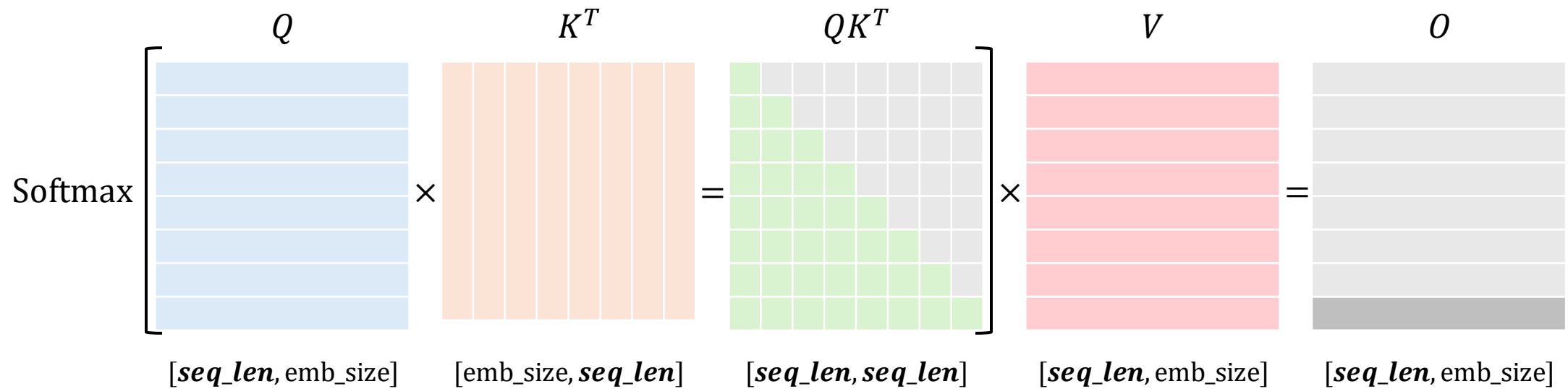
# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



$Q$       $K^T$       $QK^T$       $V$       $O$

Softmax

$[6, \text{emb\_size}]$    $[\text{emb\_size}, 6]$    $[6, 6]$    $[6, \text{emb\_size}]$    $[6, \text{emb\_size}]$

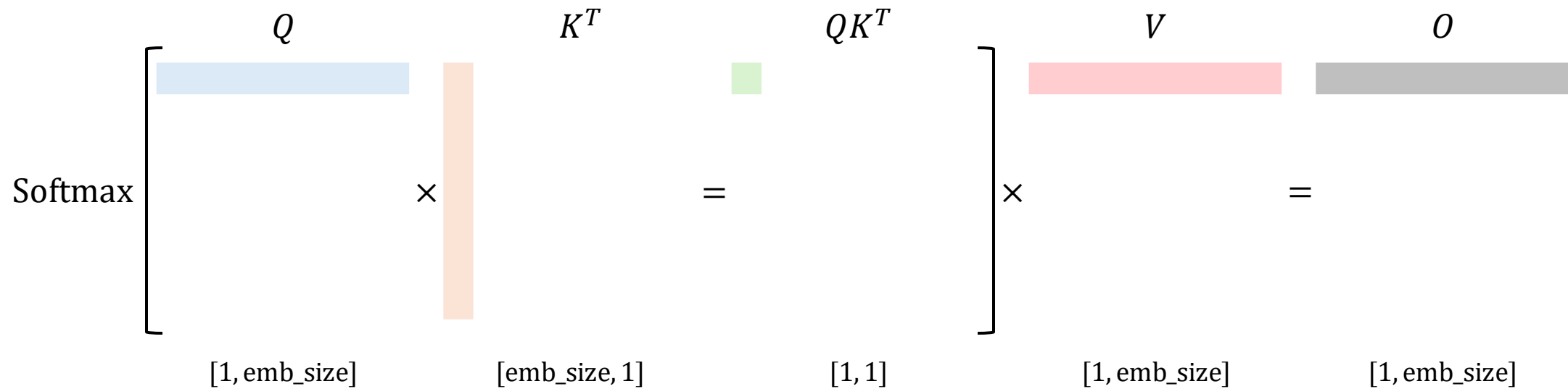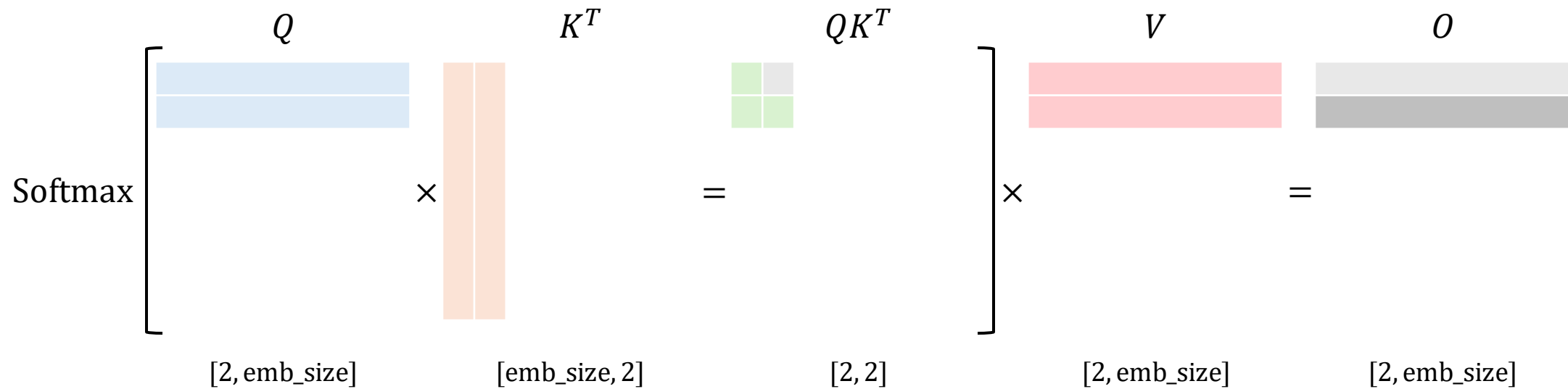<bos>   Transformers   predict   tomorrow's   words   from   today's

# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



$Q$      $K^T$      $QK^T$      $V$      $O$

Softmax

[7, emb_size]    [emb_size, 7]    [7, 7]    [7, emb_size]    [7, emb_size]

&lt;bos&gt;    Transformers    predict    tomorrow's    words    from    today's    hidden
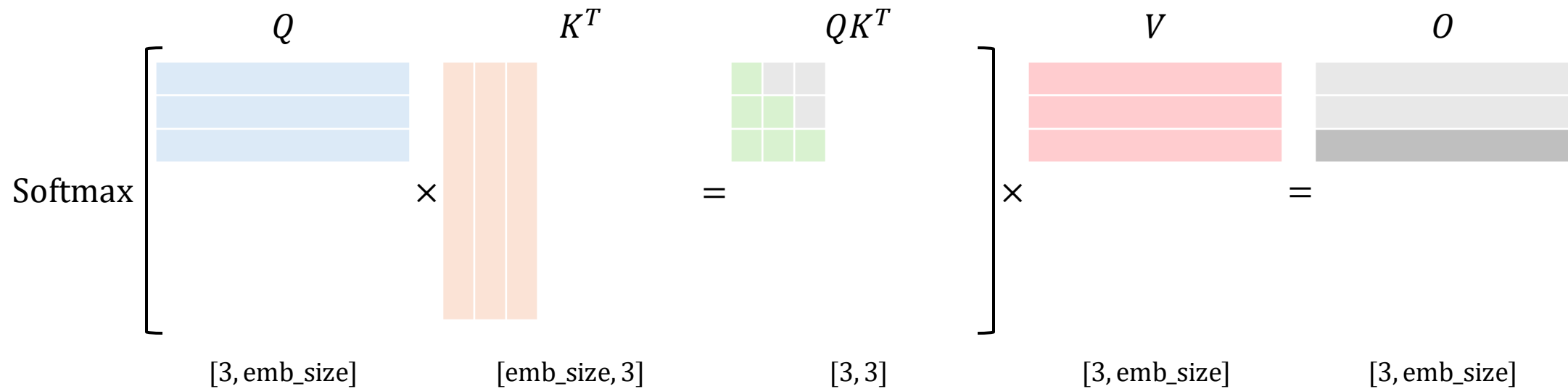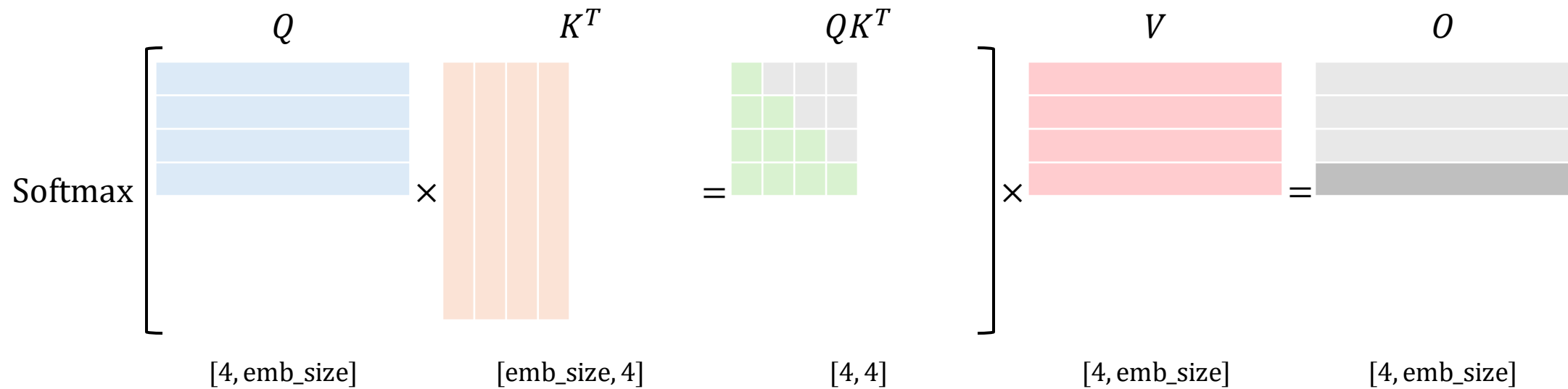
# Scaled dot-product Attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



| | $Q$ | $K^T$ | $QK^T$ | $V$ | $O$ |
|---|---|---|---|---|---|
| Softmax | $[8, \text{emb\_size}]$ | $[\text{emb\_size}, 8]$ | $[8, 8]$ | $[8, \text{emb\_size}]$ | $[8, \text{emb\_size}]$ |

<bos>  Transformers  predict  tomorrow's  words  from  today's  hidden  context
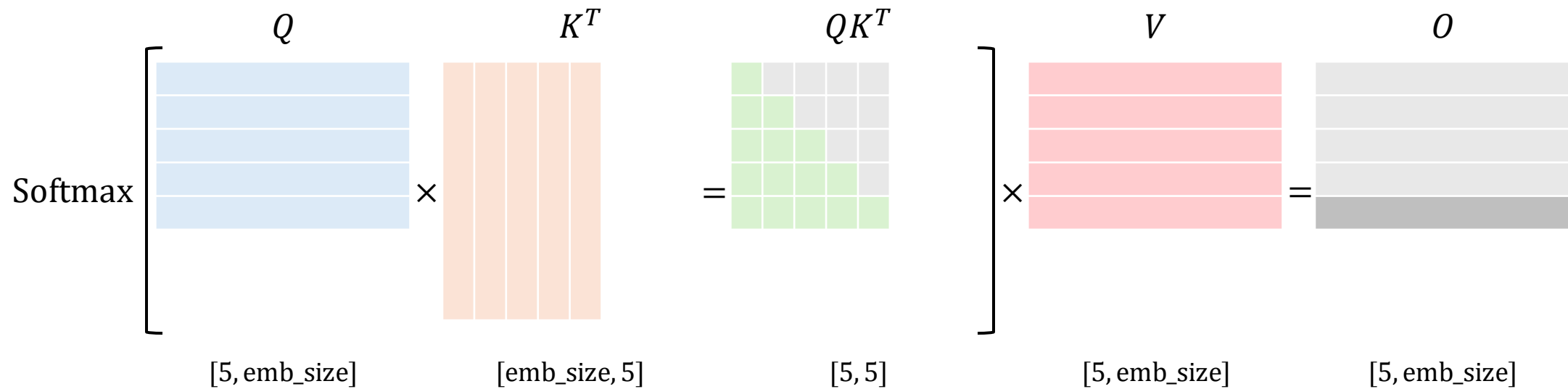
# Scaled dot-product Attention

$$\mathrm{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathrm{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$$Q \qquad\qquad K^T \qquad\qquad QK^T \qquad\qquad V \qquad\qquad O$$

Softmax

$$[8, \mathrm{emb\_size}] \qquad [\mathrm{emb\_size}, 8] \qquad [8, 8] \qquad [8, \mathrm{emb\_size}] \qquad [8, \mathrm{emb\_size}]$$

We don't use these vectors to generate the next token!

# KV-Cache

Cache the past keys and values once → reuse them for all future queries:

$$\mathbf{K}_{\text{cache}} = [k_1, k_2, \ldots, k_{t-1}]^T, \quad \mathbf{V}_{\text{cache}} = [v_1, v_2, \ldots, v_{t-1}]^T$$

Then for the next token:

$$o_t = \text{Softmax} \left( \frac{q_t^T \cdot [\mathbf{K}_{\text{cache}}^T, k_t]}{\sqrt{d_k}} \right) [\mathbf{V}_{\text{cache}}, v_t]^T$$

Complexity per step:

$$\mathcal{O}(t^2) \rightarrow \mathcal{O}(t)$$

# Scaled dot-product Attention w/ KV-Cache
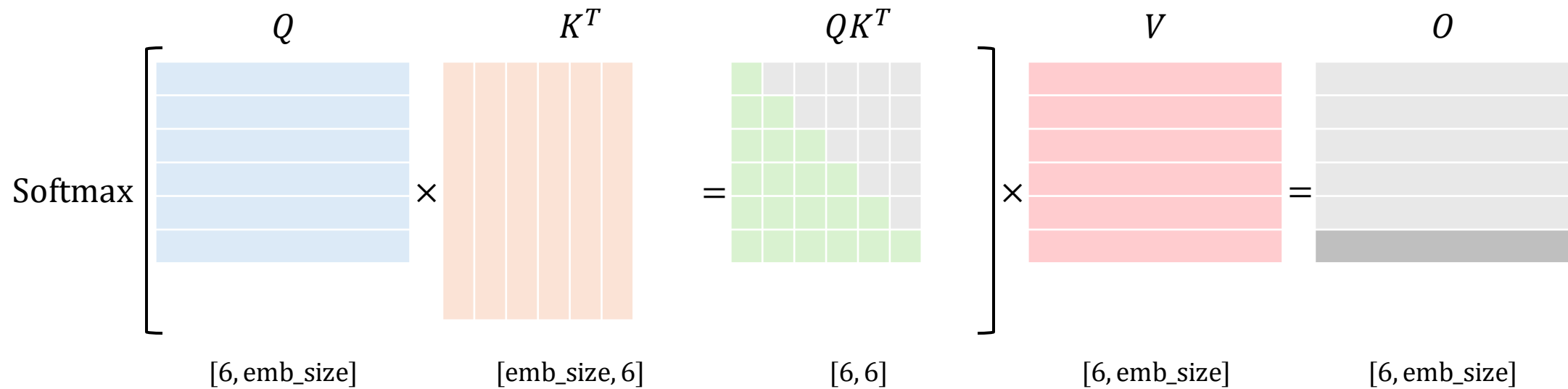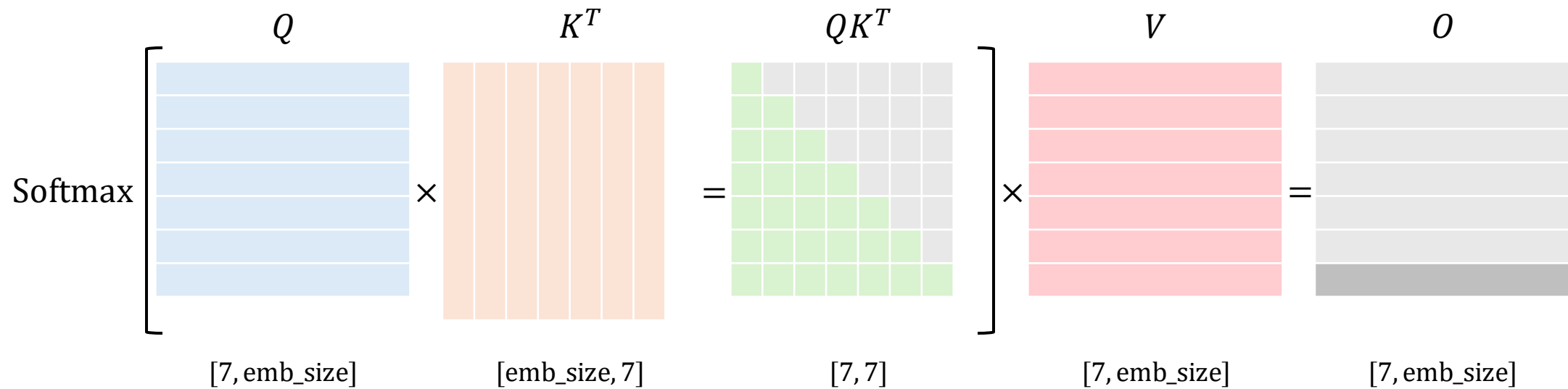
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$Q$       $K^T$       $QK^T$       $V$       $O$

Softmax    $\times$    $=$    $\times$    $=$

$[1, \text{emb\_size}]$    $[\text{emb\_size}, 1]$    $[1, 1]$    $[1, \text{emb\_size}]$    $[1, \text{emb\_size}]$

`<bos>` Transformers

# Scaled dot-product Attention w/ KV-Cache
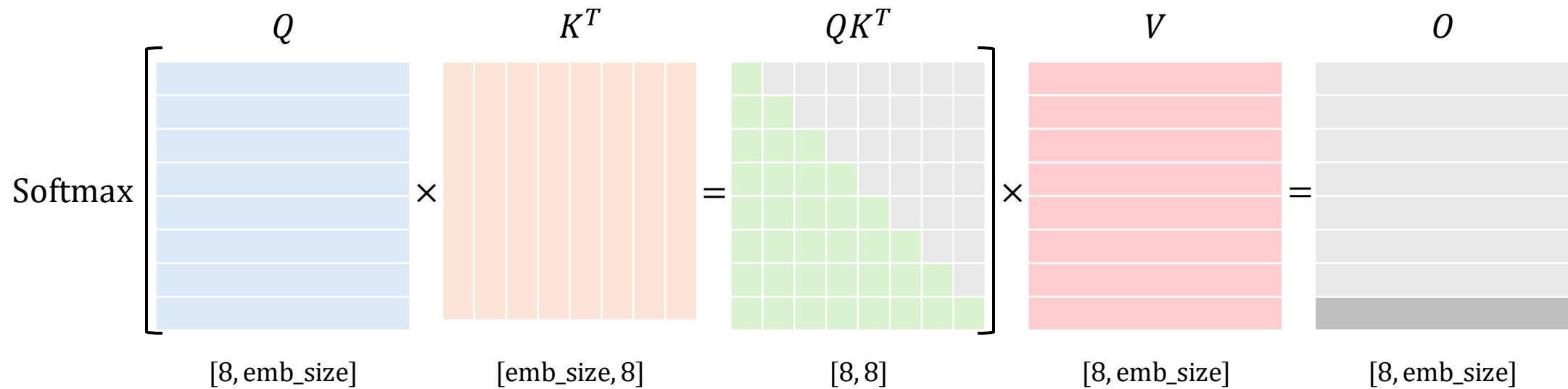
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$Q$      $K^T$      $QK^T$      $V$      $O$

Softmax $\left[ \qquad \times \qquad = \qquad \right] \times \qquad =$

$[2, \text{emb\_size}]$    $[\text{emb\_size}, 2]$    $[2, 2]$    $[2, \text{emb\_size}]$    $[2, \text{emb\_size}]$

\<bos\>    Transformers    predict

# Scaled dot-product Attention w/ KV-Cache
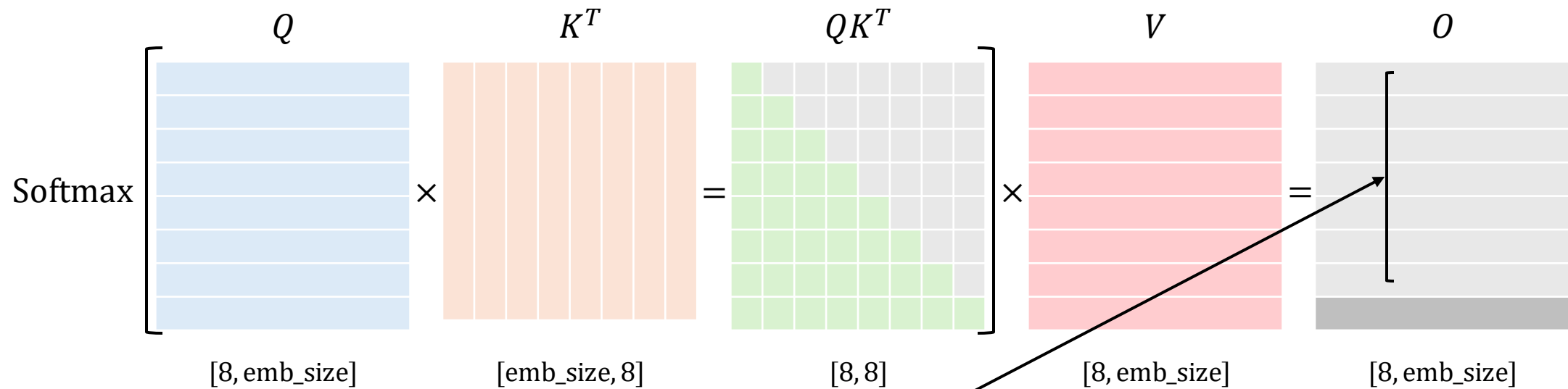
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$Q$      $K^T$      $QK^T$      $V$      $O$

Softmax

$\times$     $=$     $\times$     $=$

$[3, \text{emb\_size}]$     $[\text{emb\_size}, 3]$     $[3, 3]$     $[3, \text{emb\_size}]$     $[3, \text{emb\_size}]$

\<bos\>     Transformers     predict     tomorrow's

# Scaled dot-product Attention w/ KV-Cache

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



| $Q$ | $K^T$ | $QK^T$ | $V$ | $O$ |
|---|---|---|---|---|
| $[4, \text{emb\_size}]$ | $[\text{emb\_size}, 4]$ | $[4, 4]$ | $[4, \text{emb\_size}]$ | $[4, \text{emb\_size}]$ |

Softmax

<bos>   Transformers   predict   tomorrow's   words

# Scaled dot-product Attention w/ KV-Cache

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$Q$ $K^T$ $QK^T$ $V$ $O$

Softmax

$[5, \text{emb\_size}]$  $[\text{emb\_size}, 5]$  $[5, 5]$  $[5, \text{emb\_size}]$  $[5, \text{emb\_size}]$

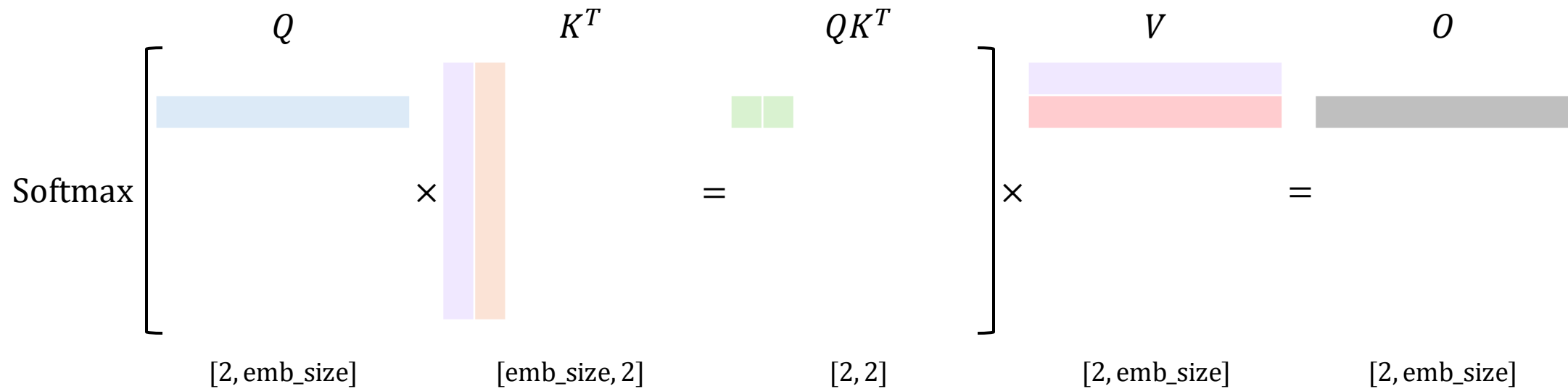\<bos\>   Transformers   predict   tomorrow's   words   from

# Scaled dot-product Attention w/ KV-Cache
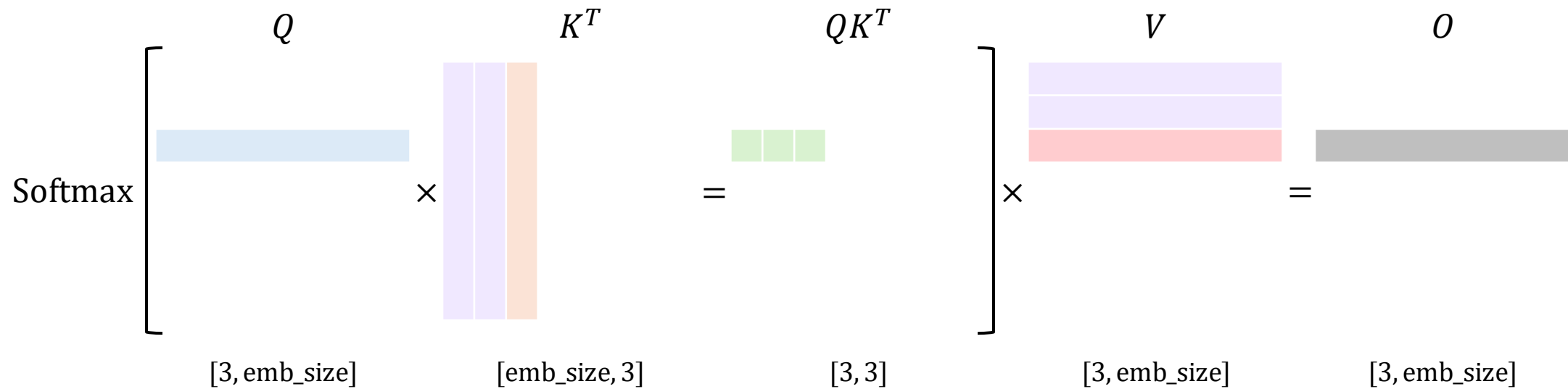
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



$Q$     $K^T$     $QK^T$     $V$     $O$

Softmax

$[6, \text{emb\_size}]$     $[\text{emb\_size}, 6]$     $[6, 6]$     $[6, \text{emb\_size}]$     $[6, \text{emb\_size}]$

<bos>  Transformers  predict  tomorrow's  words  from  today's

# Scaled dot-product Attention w/ KV-Cache
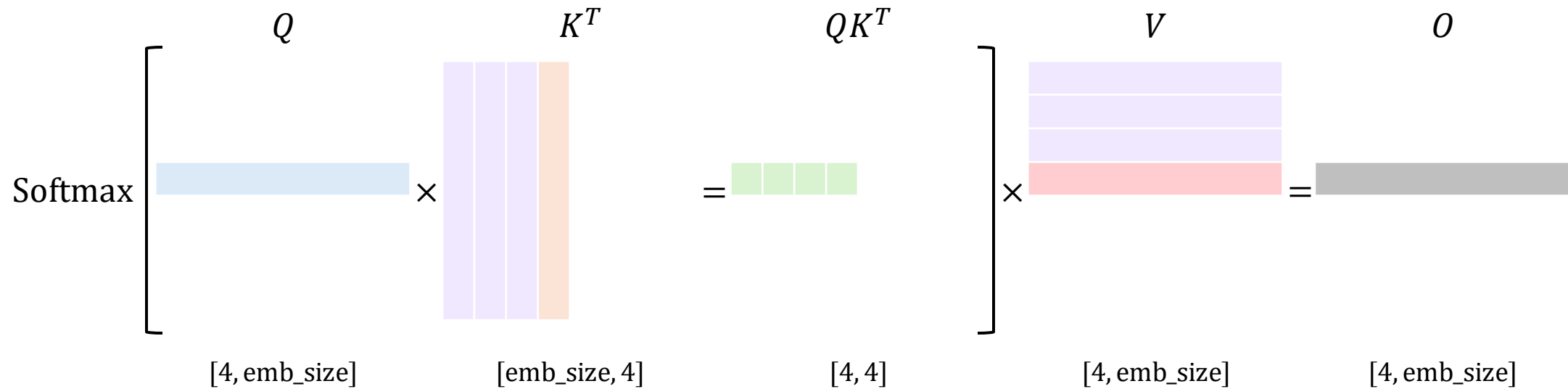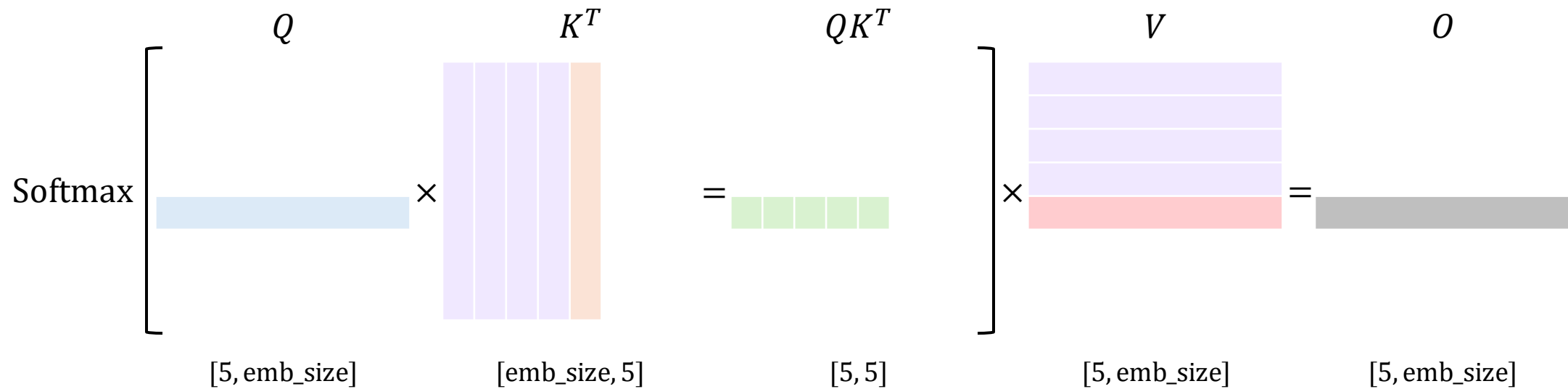
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$Q$      $K^T$      $QK^T$      $V$      $O$

Softmax

$\times$    $=$    $\times$    $=$

$[7, \text{emb\_size}]$    $[\text{emb\_size}, 7]$    $[7, 7]$    $[7, \text{emb\_size}]$    $[7, \text{emb\_size}]$

<bos>   Transformers   predict   tomorrow's   words   from   today's   hidden

# Scaled dot-product Attention w/ KV-Cache

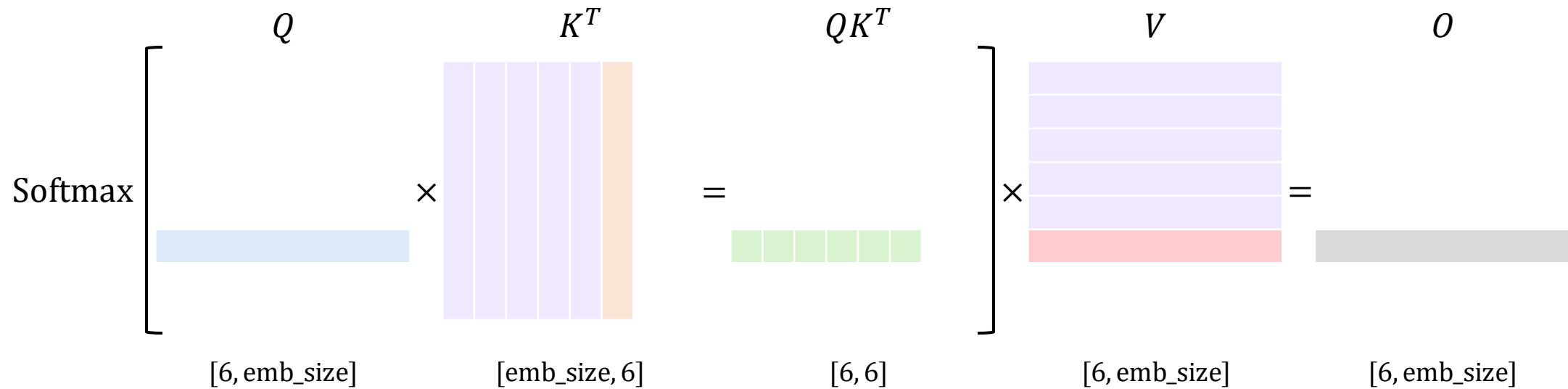$$\mathrm{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathrm{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



$Q$      $K^T$      $QK^T$      $V$      $O$

Softmax    $\times$    $=$    $\times$    $=$

$[8, \mathrm{emb\_size}]$    $[\mathrm{emb\_size}, 8]$    $[8, 8]$    $[8, \mathrm{emb\_size}]$    $[8, \mathrm{emb\_size}]$

<bos>   Transformers   predict   tomorrow's   words   from   today's   hidden   context

# Implementation Example

```python
from transformers import AutoModelForCausalLM, AutoTokenizer

tokenizer = AutoTokenizer.from_pretrained("HuggingFaceTB/SmolLM2-1.7B")
model = AutoModelForCausalLM.from_pretrained("HuggingFaceTB/SmolLM2-1.7B").cuda()

tokens = tokenizer.encode("The red cat was", return_tensors="pt").cuda()
output = model.generate(
    tokens, max_new_tokens=300, use_cache=True # by default is set to True
)
output_text = tokenizer.batch_decode(output, skip_special_tokens=True)[0]
```

| With KV-Cache | Standard Inference | Speedup |
|---------------|--------------------|---------|
| 11.7 s | 1 min 1 s | ~5.21x times faster |

# Flash Attention

# Memory Types

# Attention Step-by-Step
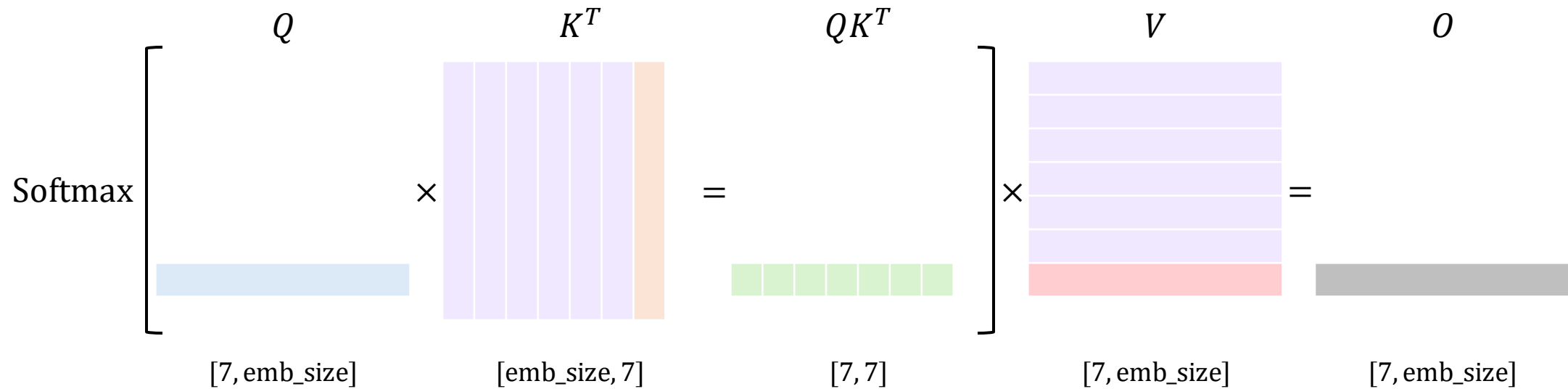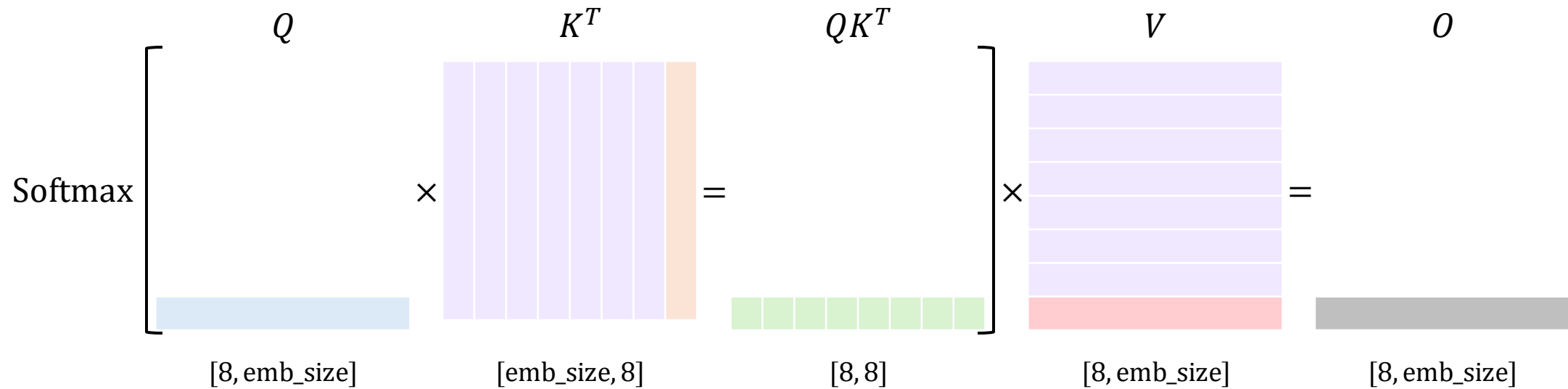
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$$S = QK^T$$

$$A = \text{Softmax}(S)$$

$$O = AV$$

# Why is Self-Attention slow?

High
Bandwidth
Memory
(HBM)

$$S = QK^T$$

$$A = \text{Softmax}(S)$$

$$O = AV$$

# Why is Self-Attention slow?



High
Bandwidth
Memory
(HBM)

Load $Q$, $K$ by blocks from HBM

Write $S$ to HBM

Compute

$S = QK^T$

$A = \text{Softmax}(S)$

$O = AV$

# Why is Self-Attention slow?



**H**igh
**B**andwidth
**M**emory
(HBM)

Load $Q, K$ by blocks from HBM

$$S = QK^T$$

Compute

Write $S$ to HBM

Load $S$ from HBM

$$A = \text{Softmax}(S)$$

Compute

Write $A$ to HBM

$$O = AV$$

# Why is Self-Attention slow?



High Bandwidth Memory (HBM)

$$S = QK^T$$

Load $Q$, $K$ by blocks from HBM

Compute

Write $S$ to HBM

$$A = \text{Softmax}(S)$$

Load $S$ from HBM

Compute

Write $A$ to HBM

$$O = AV$$

Load $A$, $V$ by blocks from HBM

Compute

Write $O$ to HBM

# IO-aware Algorithm – Tiling

# IO-aware Algorithm – Tiling

$B$

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

$A$

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

| $c_{1,1}$ | $c_{1,2}$ | | |
|---|---|---|---|
| $c_{2,1}$ | $c_{2,2}$ | | |
| | | | |
| | | | |

$C_{1,1} = 1 \times 1 + 2 \times 5 + 3 \times 9 + 4 \times 13$

$C_{1,2} = 1 \times 2 + 2 \times 6 + 3 \times 10 + 4 \times 14$

$C_{2,1} = 5 \times 1 + 6 \times 5 + 7 \times 9 + 8 \times 13$

$C_{2,2} = 5 \times 2 + 6 \times 6 + 7 \times 10 + 8 \times 14$

$C = A \times B$

# IO-aware Algorithm – Tiling



$B$

| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

$A$

| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

| $c_{1,1}$ | $c_{1,2}$ | | |
| $c_{2,1}$ | $c_{2,2}$ | | |
| | | | |
| | | | |

**Without tiling:** 32 memory accesses

$C_{1,1} = 1 \times 1 + 2 \times 5 + 3 \times 9 + 4 \times 13$

$C_{1,2} = 1 \times 2 + 2 \times 6 + 3 \times 10 + 4 \times 14$

$C_{2,1} = 5 \times 1 + 6 \times 5 + 7 \times 9 + 8 \times 13$

$C_{2,2} = 5 \times 2 + 6 \times 6 + 7 \times 10 + 8 \times 14$

$C = A \times B$

# IO-aware Algorithm – Tiling

$B$

| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

**Without tiling:** 32 memory accesses
**With tiling:** 16 memory accesses

$N \times N$ block $\rightarrow 1/N$ memory access

$$\begin{bmatrix} c_{1,1} & c_{1,2} \\ c_{2,1} & c_{2,2} \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix}\begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix} + \begin{bmatrix} 3 & 4 \\ 7 & 8 \end{bmatrix}\begin{bmatrix} 9 & 10 \\ 13 & 14 \end{bmatrix}$$

$A$

| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

| $c_{1,1}$ | $c_{1,2}$ | | |
| $c_{2,1}$ | $c_{2,2}$ | | |

$$C = A \times B$$

# IO-aware Algorithm – Tiling



$B$

$$B_{1,1} \quad B_{1,2}$$
$$B_{2,1} \quad B_{2,2}$$

How to make efficient?

$$S = QK^T$$
$$A = \text{Softmax}(S)$$
$$O = AV$$

$A$

$$A_{1,1} \quad A_{1,2}$$
$$A_{2,1} \quad A_{2,2}$$

$$C_{1,1} \quad C_{1,2}$$
$$C_{2,1} \quad C_{2,2}$$

$$C = A \times B$$

# Softmax



Runtime Warning: overflow encountered in exp y = torch.exp(x)

# Safe Softmax

$$m = max(x_1, x_2, ..., x_N)$$

$$d_N = \sum_{i=1}^{N} e^{x_i - m}$$

# Safe Softmax

$$m_N = max(x_1, x_2, \ldots, x_N)$$

$$d_N = \sum_{i=1}^{N} e^{x_i - m_N}$$

$m_0 = -\infty$

**for** $i = 1, \ldots, N$ **do**

$\quad m_i = \max(m_{i-1}, x_i)$

$d_0 = 0$

**for** $i = 1, \ldots, N$ **do**

$\quad d_i = d_{i-1} + e^{x_i - m_N}$

**for** $i = 1, \ldots, N$ **do**

$\quad a_i = e^{x_i - m_N}/d_N$

**3 loops**

# Online Softmax

$m_0 = -\infty$

**for** $i = 1, \ldots, N$ **do**

$\quad m_i = \max(m_{i-1}, x_i)$

$d_0 = 0$

**for** $i = 1, \ldots, N$ **do**

$\quad d_i = d_{i-1} + e^{x_i - m_N}$

**for** $i = 1, \ldots, N$ **do**

$\quad a_i = e^{x_i - m_N}/d_N$

$$d_i = \sum_{j=1}^{i} e^{x_j - m_N} \qquad \textcolor{red}{d'_i = \sum_{j=1}^{i} e^{x_j - m_i}} \qquad \textcolor{red}{d'_N = d_N}$$

$$\textcolor{red}{d'_i} = \underbrace{\left(\sum_{j=1}^{i} e^{x_j - \textcolor{red}{m_{i-1}}}\right)}_{\textcolor{red}{d'_{i-1}}} e^{\textcolor{red}{m_{i-1} - m_i}} + e^{\textcolor{red}{x_i - m_i}}$$

# Online Softmax

$m_0 = -\infty$

**for** $i = 1, \dots, N$ **do**

$\quad m_i = \max(m_{i-1}, x_i)$

$d_0 = 0$

**for** $i = 1, \dots, N$ **do**

$\quad d_i = d_{i-1} + e^{x_i - m_N}$

**for** $i = 1, \dots, N$ **do**

$\quad a_i = e^{x_i - m_N} / d_N$

$\Longrightarrow$

$m_0 = -\infty$

$d_0 = 0$

**for** $i = 1, \dots, N$ **do**

$\quad m_i = \max(m_{i-1}, x_i)$

$\quad d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$

**for** $i = 1, \dots, N$ **do**

$\quad a_i = e^{x_i - m_N} / d'_N$

# Online Softmax

$m_0 = -\infty$

$d_0 = 0$

**for** $i = 1, \ldots, N$ **do**

$\quad x_i = q k_i^T$

$\quad m_i = \max(m_{i-1}, x_i)$

$\quad d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$

$o_0 = 0$

**for** $i = 1, \ldots, N$ **do**

$\quad a_i = e^{x_i - m_N} / d'_N$

$\quad o_i = o_{i-1} + a_i v_i$

$S = Q K^T$

$A = \text{Softmax}(S)$

$O = A V$

# Online Softmax

$m_0 = -\infty$

$d_0 = 0$

**for** $i = 1, \dots, N$ **do**

   $x_i = q k_i^T$

   $m_i = \max(m_{i-1}, x_i)$

   $d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$

$o_0 = 0$

**for** $i = 1, \dots, N$ **do**

   $a_i = e^{x_i - m_N} / d'_N$

   $o_i = o_{i-1} + a_i v_i$

$$S = Q K^T$$

$$A = \text{Softmax}(S)$$

$$O = AV$$

We can do better!

Let's do the same trick:

$$o_i = \sum_{j=1}^{i} \frac{e^{x_j - m_N}}{d'_N} v_j \longrightarrow o'_i = \sum_{j=1}^{i} \frac{e^{x_j - m_i}}{d'_i} v_j$$

$$o_N = o'_N$$

# Flash Attention ⚡

$m_0 = -\infty$

$d_0 = 0$

$o_0 = 0$

**for** $i = 1, \ldots, N$ **do**

$\quad x_i = qk_i^T$

$\quad m_i = \max(m_{i-1}, x_i)$

$\quad d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$

$\quad o'_i = o'_{i-1} \dfrac{d'_{i-1}}{d'_i} e^{m_{i-1} - m_i} + \dfrac{e^{x_i - m_i}}{d'_i} v_i$

**return** $o'_N$

$S = QK^T$

$A = \text{Softmax}(S)$

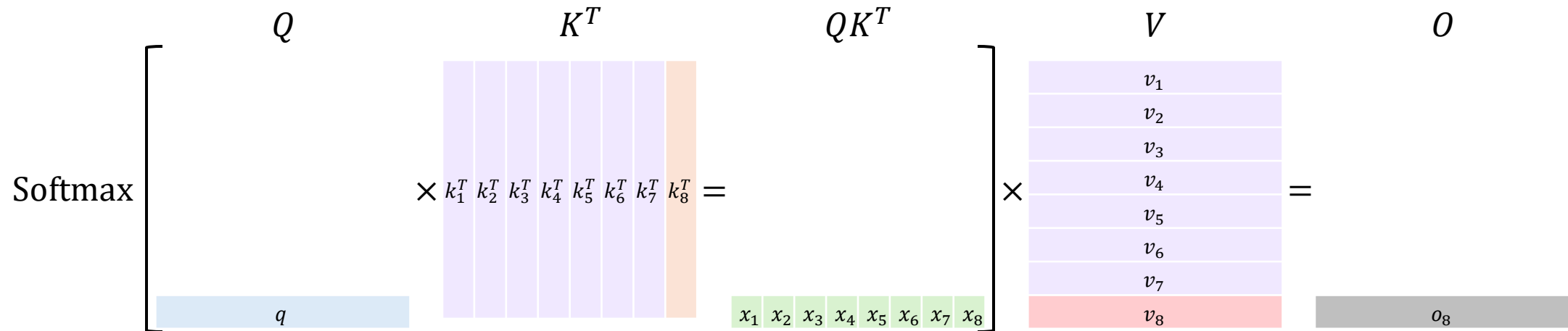$O = AV$

# Flash Attention ⚡

**for** $i = 1, \ldots, N$ **do**

$\quad x_i = q k_i^T$

$\quad m_i = \max(m_{i-1}, x_i)$

$\quad d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$

$\quad o'_i = o'_{i-1} \dfrac{d'_{i-1}}{d'_i} e^{m_{i-1} - m_i} + \dfrac{e^{x_i - m_i}}{d'_i} v_i$

**return** $o'_N$

$$Q \qquad K^T \qquad QK^T \qquad V \qquad O$$

Softmax

$\times \quad k_1^T \; k_2^T \; k_3^T \; k_4^T \; k_5^T \; k_6^T \; k_7^T \; k_8^T \quad =$

$q$

$x_1 \; x_2 \; x_3 \; x_4 \; x_5 \; x_6 \; x_7 \; x_8$

$\times$

$v_1$
$v_2$
$v_3$
$v_4$
$v_5$
$v_6$
$v_7$
$v_8$

$=$

$o_8$

# Flash Attention ⚡



Dao T. et al. FlashAttention: Fast and memory-efficient exact attention with io-awareness, 2022

# Recap

- Quantization

- Pruning

- Distillation

- KV-Cache

- Flash Attention

# Course Overview

| Introduction to Neural Networks | Natural Language Processing | Computer Vision | Reinforcement Learning | Generative Models | Multimodality | Acceleration |
|---|---|---|---|---|---|---|
| MLP | Word Embeddings | Image Classification | Multi-armed Bandits | Autoregression | CLIP | Quantization |
| Backpropagation | RNN, LTSM | Object Detection | Monte Carlo Methods | VAE | BLIP | Pruning |
| Initialization | Attention | Segmentation | Policy Improvement | GAN | VLM | Distillation |
| Regularization | Transformer | | | Diffusion | | KV-Cache |
| CNN | | | | Flow Matching | | Flash Attention |

# Game Rules

- 5 Homeworks = **70 points**

- Oral Exam = **30 points**

- Maximum Points: 70 + 30 = **100 points**

**Final Grade:** min(round(#points, 10), 10)

Thanks for your $\text{Softmax}\left(\dfrac{\mathbf{QK}^T}{\sqrt{d_k}}\right)\mathbf{V}$