



Midterm Coverage

General Information

- In-class exam on Mar 27 5:10pm-6:30pm
- 80 mins
- Close book, you can bring calculator
- No cheat sheet is allowed
- Five questions

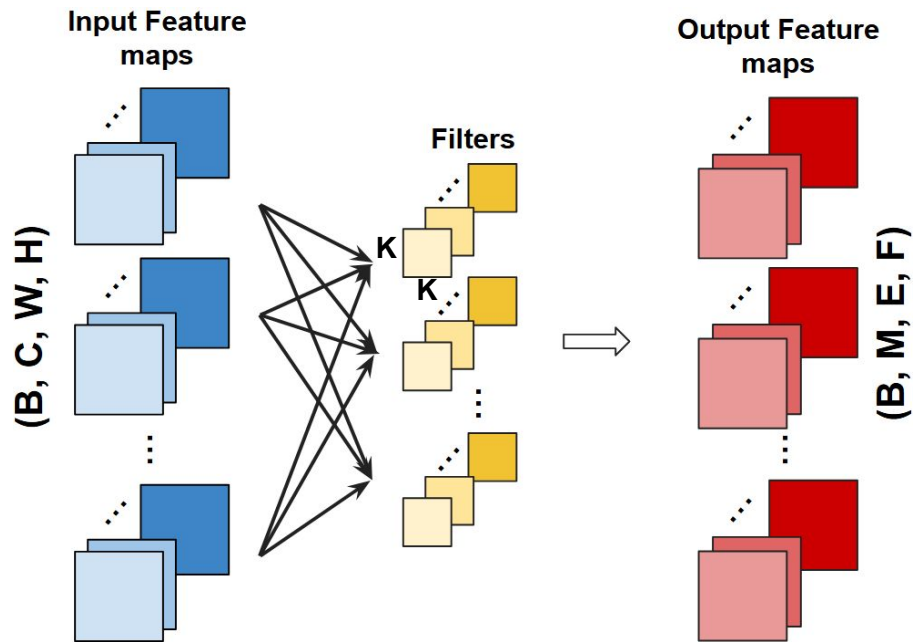
Recap (Lecture 1: P65-P98)

- Neural Network Basics
 - MLP
 - Forward and backward propagation of MLP
 - Weight decay, dropout
 - The training optimizer: SGD, RMSProp, Adam
 - Multistage learning rate scheduler
 - Initialization

Recap (Lecture 2: P1-P69, P75-P78)

- Conv2D operation
 - How the computation is performed
 - Input dimension, weight dimension, output dimension
 - Computational cost
- BatchNorm
 - Parameter folding-in during inference
- ResNet, MobileNet, ShuffletNet, SqueezeNet, DenseNet
 - Depthwise Separable Conv
 - Groupwise Convolution
- Other tasks: segmentation

Computational Cost: Standard Convolution



- Number of MACs: $B \times M \times K \times K \times C \times E \times F$
- Storage cost:
 $32 \times (M \times C \times K \times K + B \times C \times H \times W + B \times M \times E \times F)$

B: batch size

C: number of input channels

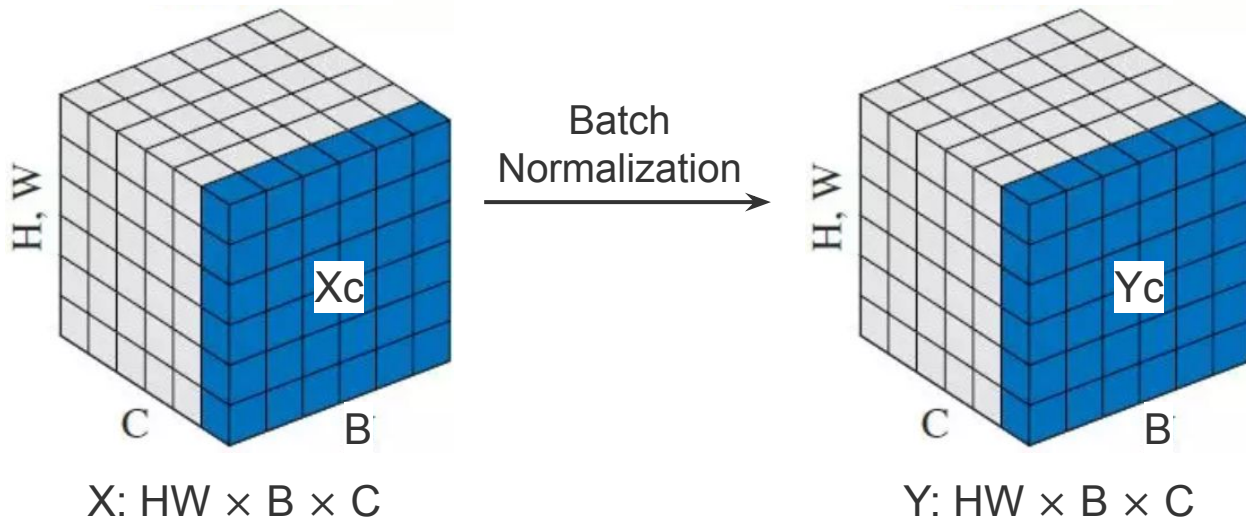
H,W: size of the input feature maps

M: number of weight filters

K: weight kernel size

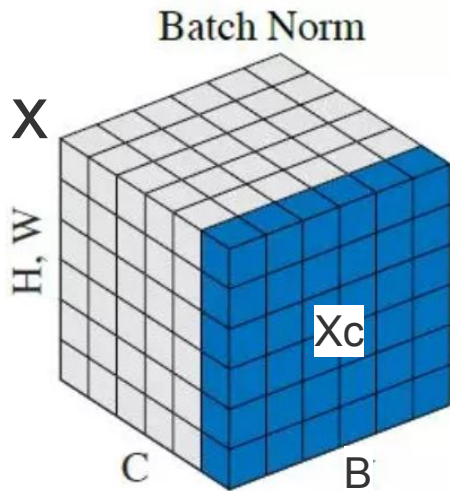
E,F: size of the output feature maps

Batch Normalization



- **Batch Normalization (BatchNorm)** is a technique used in deep learning to improve the training stability and performance of neural networks.

Batch Normalization



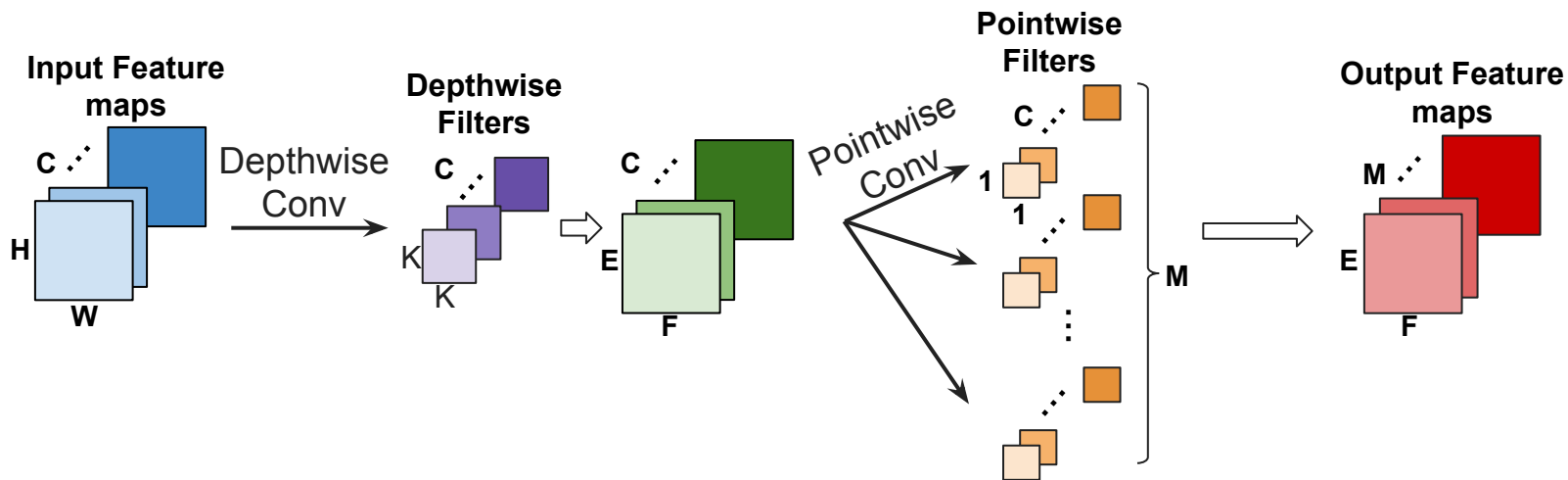
X: $HW \times B \times C$

$$Y_c = \alpha_c \frac{X_c - \mu_c}{\sigma_c} + \beta_c \quad \text{For each } c \in C$$

$$\alpha = \{\alpha_c\}, \beta = \{\beta_c\}, \mu = \{\mu_c\}, \sigma = \{\sigma_c\}$$

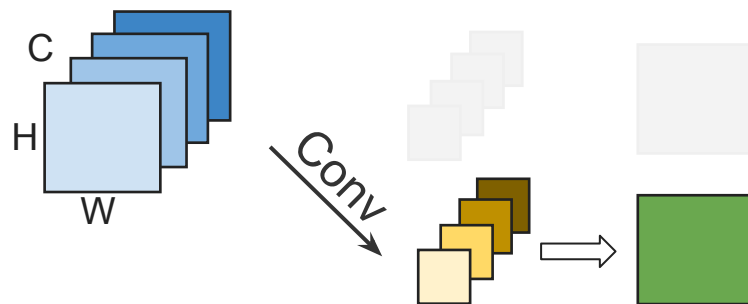
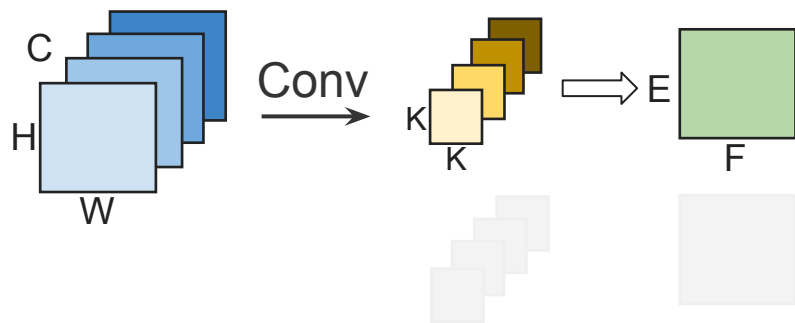
- For each channel c , we have:
 - X_c : ($HW \times B$)
 - μ_c and δ_c are the mean and standard deviation of X_c .
 - α_c and β_c are learnable parameters
 - $\alpha_c, \beta_c, \mu_c, \delta_c$ are scalars
- Overall, we have:
 - μ, δ, α and β all have a length of C
 - μ, δ, α and β are all fixed during the inference
 - μ, δ are statistics based on the training dataset

Depthwise Separable Convolution



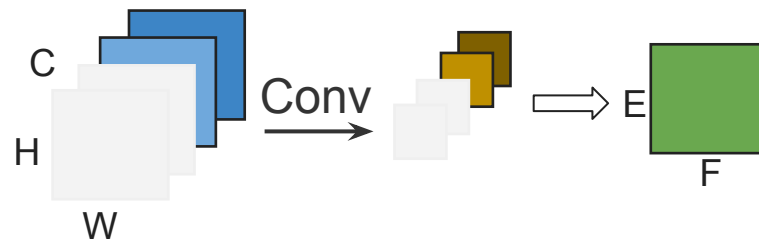
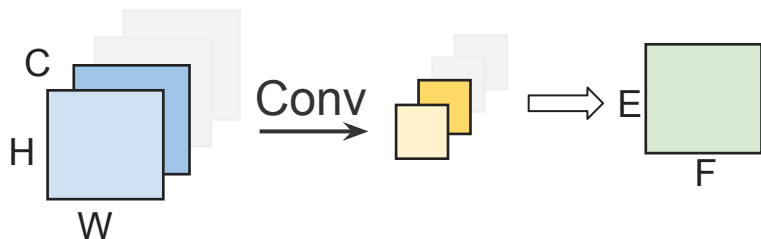
- Number of MACs: $K \times K \times C \times E \times F + M \times C \times E \times F$
- Storage cost: $32 \times (C \times H \times W + C \times K \times K + C \times E \times F + M \times C + M \times E \times F)$

Group Convolution



- The original MAC: $E \times F \times K \times K \times C \times M$

Group Convolution



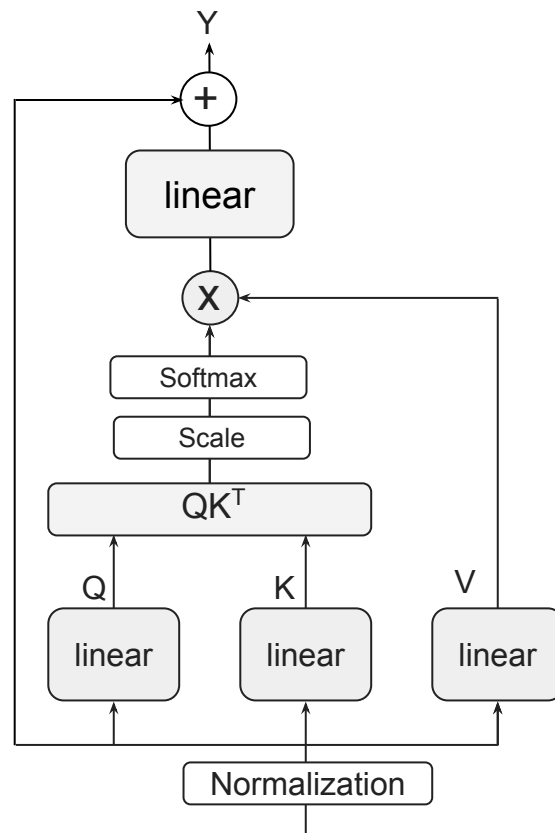
- Group size = 2
- Each group of feature maps within the input only convolved with partial weight kernels.
- This will lead to a large saving on memory consumption and computational cost.
- The number of MAC: $E \times F \times K \times K \times C \times M/G$

Recap (Lecture 3: P1-P69, P75-P78)

- Transformers
 - How the computation is performed and why
 - Multi-headed attention, FFN
 - LayerNorm, RMSNorm, GeLU
 - Positional embedding, Word embedding
- Vision Transformer
 - How to convert image into visual tokens
- LLM
 - Prefilling, decoding
 - KV cache
- SSL
 - Basics

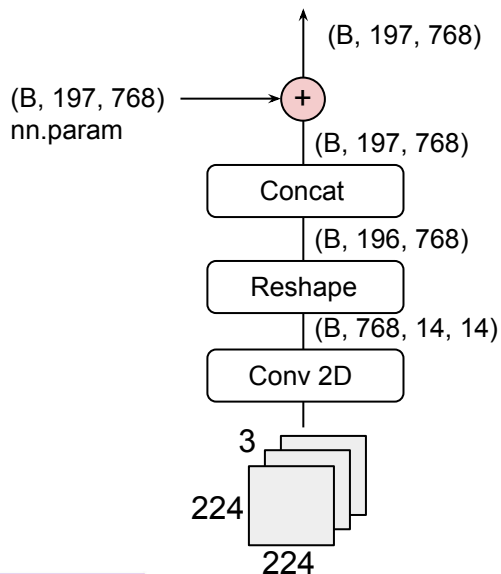
Self-Attention Block

- Given input x , the first step in calculating self-attention is to create three vectors from each of the input x , denoted as: Query (Q), Key (K), Value (V).
 - $(B, L, E) * (E * E) \rightarrow (B * L * E)$
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
 - $QK^T \rightarrow (B, L * E) * (B, E * L) \rightarrow (B, L * L)$
- Scale and normalize the score using softmax.
 - $\text{Softmax}(QK^T) \rightarrow (B, L * L)$
- Multiply each value vector by the softmax score.
 - $\text{Softmax}(QK^T) * V$
 - $(B, L * L) * (B, L * E) \rightarrow (B, L * E)$
- Pass the result to the linear layer, sum with the input.

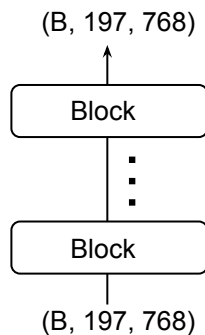


Vision Transformer

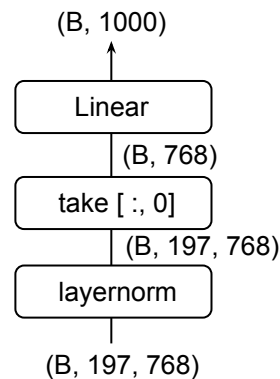
- Transformer architecture can also be applied over the computer vision tasks.



(part 1)

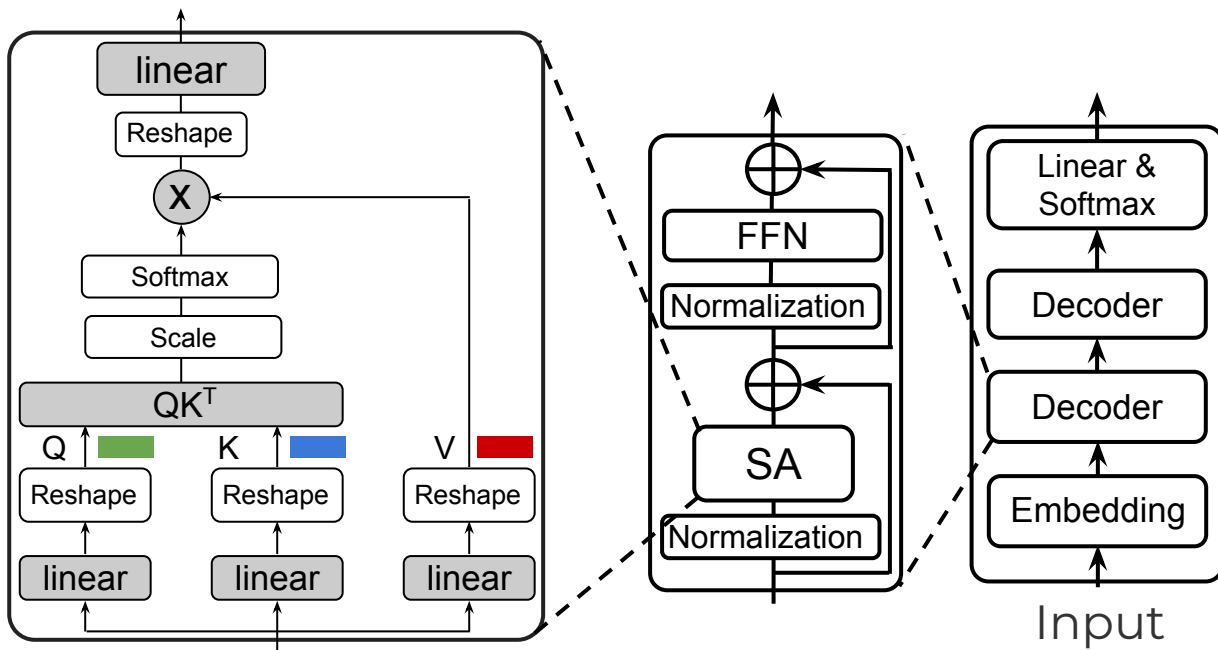


(part 2)



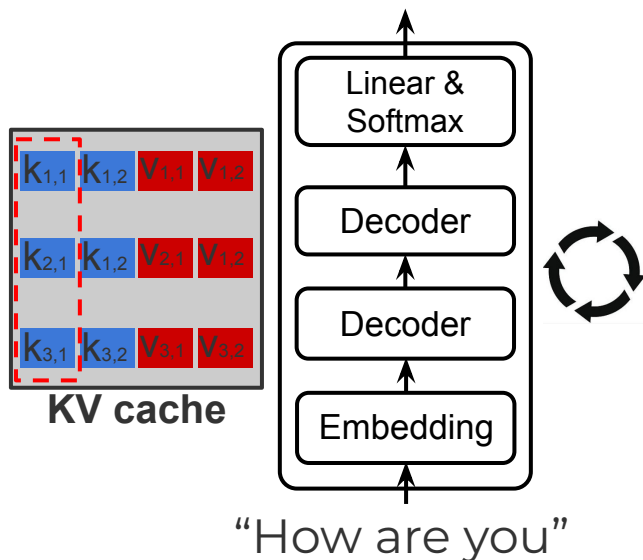
(part 3)

Transformers as a Generative AI Tool



- We need to buffer the v and k for later usage.

GPT-2: Prefilling

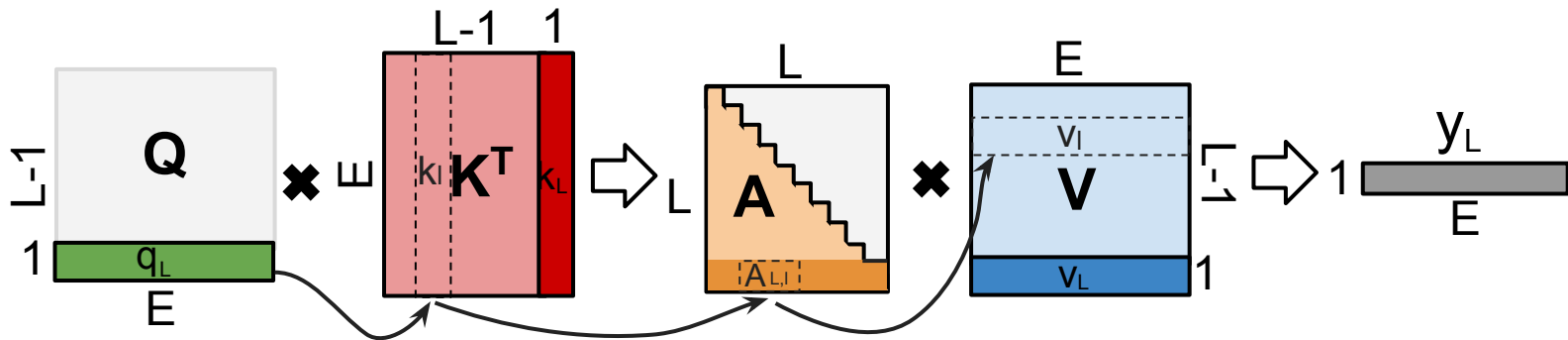


$k_{i,j}$ Key vector for i th token in j th layer
($1 \times E$)

$v_{i,j}$ Value vector for i th token in j th layer
($1 \times E$)

- During the prefilling stage, LLM processes the entire prompt, or context tokens jointly, saving the KV vectors into the memory.

Why KV Cache Saves Computation?

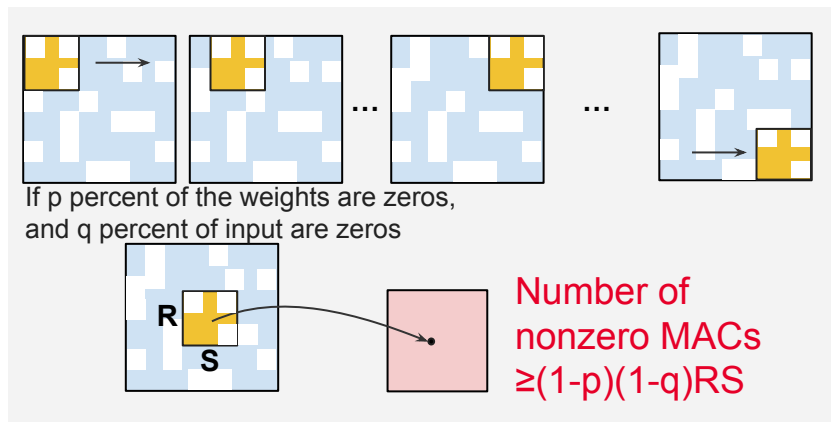


- During the decoding phase, new tokens are continuously generated and must be processed using the buffered K and V vectors to generate subsequent tokens.
- Without a KV cache, all previous K and V vectors must be recomputed, resulting in significant computational overhead.

Recap (Lecture 4: P1-P65, P68-P75)

- Computational cost saving with pruning
 - CNN & Transformer
- Sparse matrix encoding
 - Bitmap, Run-length encoding, COO
- General pruning techniques
 - Magnitude pruning, gradient-based, Hessian-based pruning
 - Lasso
 - Taxonomy of Pruning
 - Network Slimming, N:M sparsity
 - Cascade effect of pruning
- Transformer pruning
 - Token pruning
 - Head pruning

Convolution with Sparse Weight



- Number of MACs $\geq (1-p) \times (1-q) \times M \times C \times R \times S \times E \times F$
- Input pruning can also reduce the computations.
- Sparse input and weight matrices can be stored more efficiently, which helps minimize memory storage.

Sparse Matrices Encodings

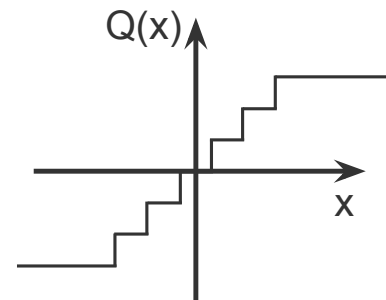
- Efficient encoding scheme for sparse matrix storage.
 - Bitmap
 - Run Length Encoding (RLE)
 - Coordinate format (COO)

Recap (Lecture 5: P1-P70)

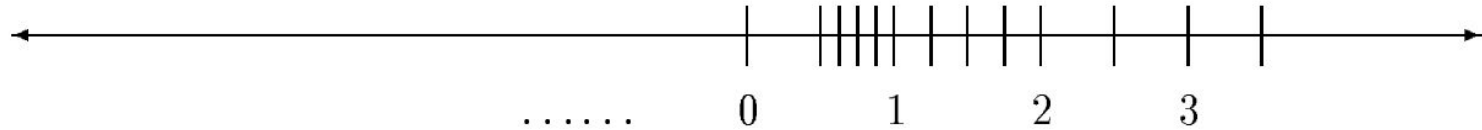
- Basic Data Formats
 - Fixed point (INT), Floating point (FP), Block floating point (BFP)
- Quantization
 - Unsymmetrical & Symmetrical
 - Why fixed, FP, BFP & logarithm quantization can save computation?
- STE
- Taxonomy of Quantization
- Quantization during training
 - Stochastic quantization
- Learnable adaptive quantization scheme

Fixed-Point Format (Symmetrical)

- How to convert a number x to INT representation?
 - Set the clipping range: $(-L, L)$, bitwidth: b
 - Compute the scale: $s = 2L / (2^b - 2)$
 - Clip the input x : $x_c = \text{Clip}(x, L, -L)$
 - Calculate the INT representation: $x_{int} = \text{round}(x_c / s)$
 - Rescale: $x_q = Sx_{int}$
- Have a uniform representation power within the clipping range.
- All the computations can be performed using x_{int}



Floating Point Arithmetic



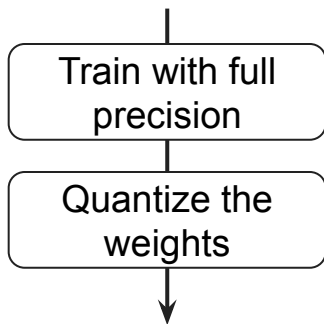
- Have better representation power for values with small magnitudes.
- How to convert a real number x to FP representation?

$$x = |x| \quad s = \text{sign}(x)$$

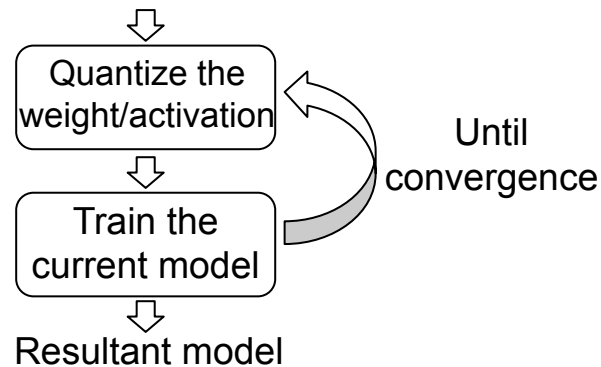
$$a = \lfloor \log_2 x \rfloor \quad e = a + \text{bias} \quad m = \frac{x}{2^a} - 1$$

When to Quantize?

Post-training quantization (PTQ)

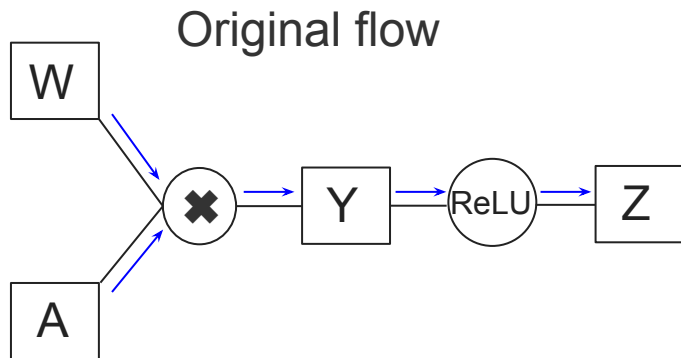


Quantization-aware Training (QAT)



- PTQ has lower computational cost, but accuracy is also lower.
- For the model which is expensive to train (LLM), PTQ is applied to facilitate their implementations.

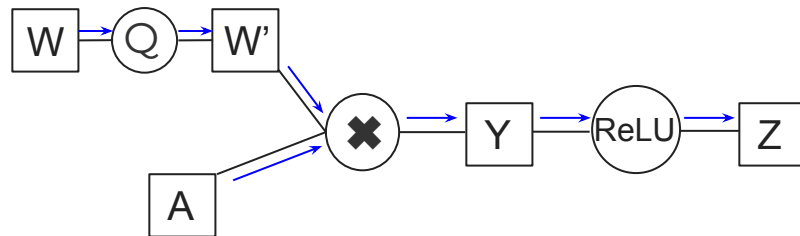
Another Way to Look at Quantization



$$Y = WA, Z = \text{ReLU}(Y)$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W}$$

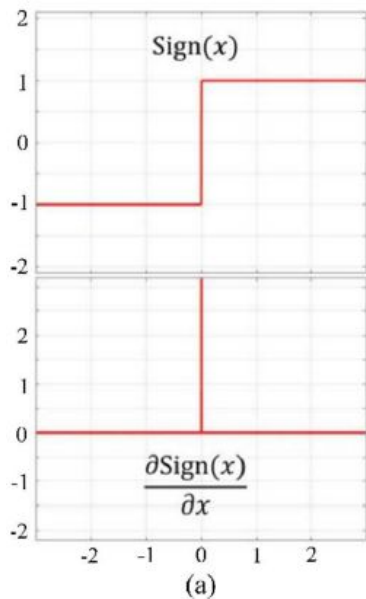
Flow with quantization



$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W'} \frac{\partial W'}{\partial W}$$

How to compute $\frac{\partial W'}{\partial W}$?

Straight Through Estimator (STE)



- Staircase function has a derivative of 0 at most of the values. This will make the DNN not trainable.
- We instead use STE to estimate the gradient of a non-differentiable quantized function in the backward pass.

$$\frac{\partial W'}{\partial W} = 1$$

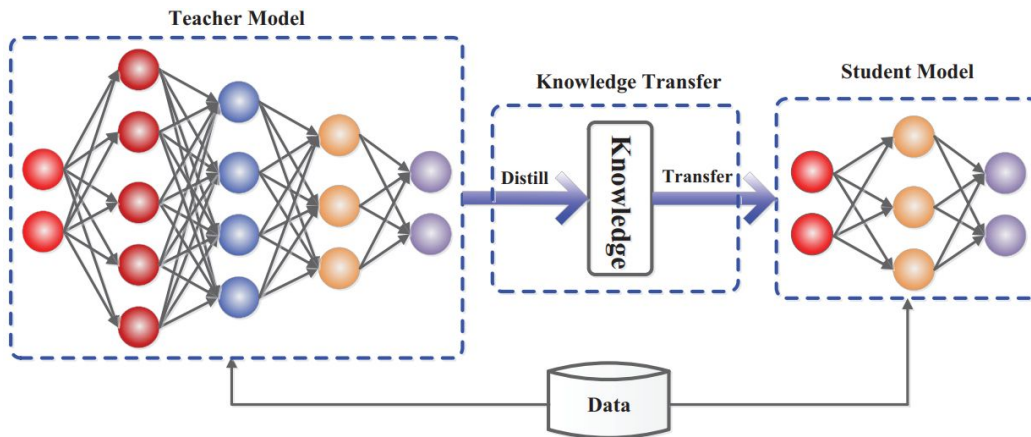
- During the forward pass, apply quantization, for backprop, ignore it.

Recap (Lecture 6: P1-P32, P50-P60, P63-P68)

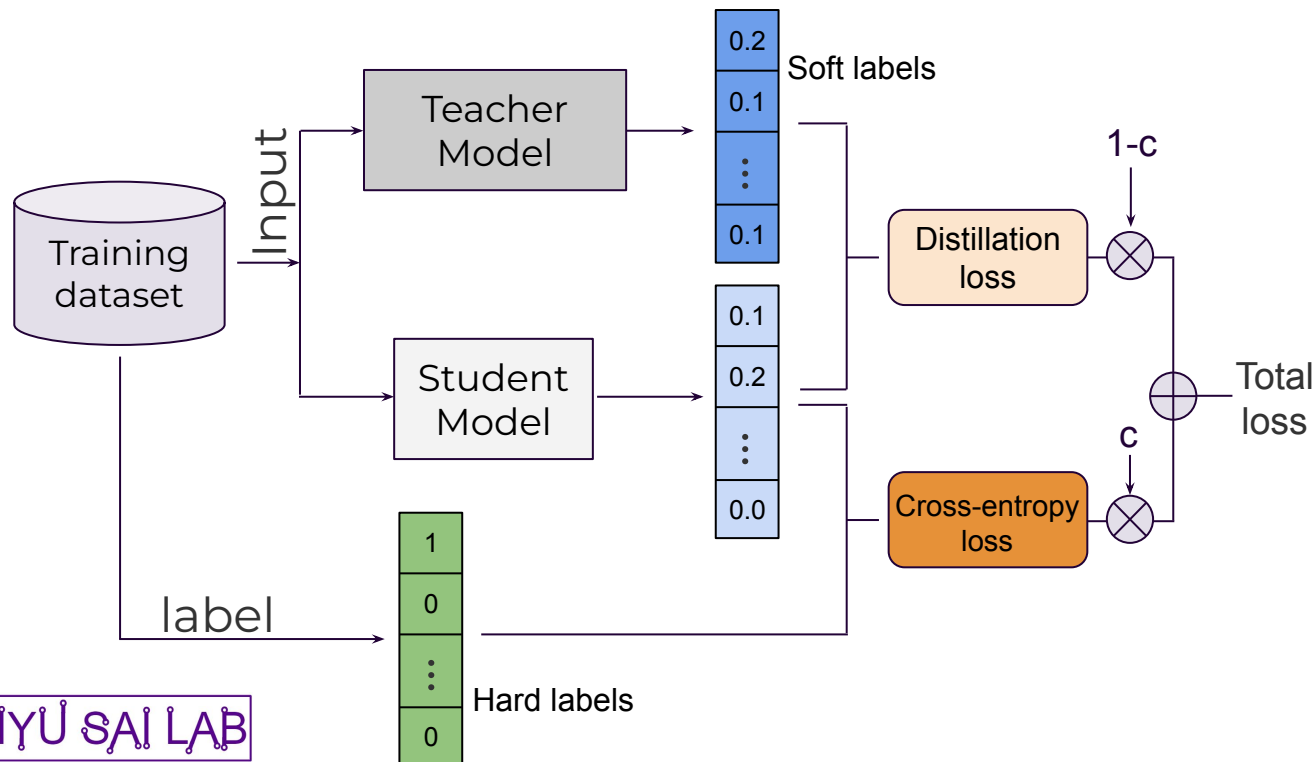
- Distillation
 - Feature-Based Knowledge Distillation
 - Online distillation
 - Self distillation
 - Multi-teacher, multi-student, cross-modal
- Once-for-all NAS
- Low-rank Decomposition
- Dynamic / Conditional Computing

Knowledge Distillation Basics

- Transferring knowledge from a large and complex model or set of models to a single, smaller model that can be effectively deployed in real-world scenarios with practical limitations.



Response-Based Knowledge Distillation



- During the training process, only the weights within the student model got updated.
- The teacher model can be either pretrained or trained together with student.
- c tends to be close to 1: ~ 0.9 .

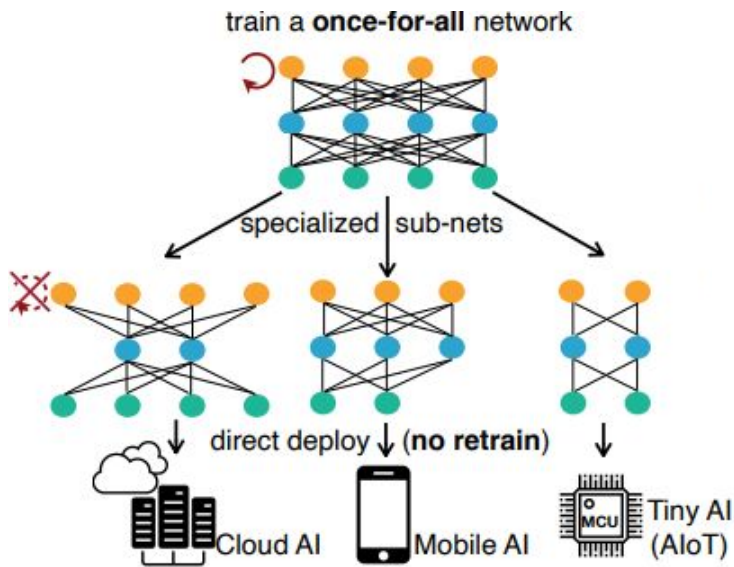
KL Divergence

- In mathematical statistics, the Kullback–Leibler (KL) divergence, denoted $D_{\text{KL}}(P \parallel Q)$, is a measure of how much a model probability distribution Q is different from a true probability distribution P , where both P and Q are discrete.

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

- If $P(x) = Q(x)$ for all x , then the KL divergence equals 0.
- Otherwise, KL divergence is greater than 0.

Once-For-All: Train One Network And Specialize It for Efficient Deployment



- A supernet is train which contained multiple smaller subnetworks.
- All the subnetworks are trained jointly.

Singular Value Decomposition

$$\begin{matrix} & n \\ m & \boxed{W} \end{matrix} = \begin{matrix} & r \\ m & \boxed{U} \end{matrix} \times \begin{matrix} & r \\ r & \boxed{R} \end{matrix} \times \begin{matrix} & n \\ r & \boxed{V^T} \end{matrix} = \begin{matrix} & r \\ m & \boxed{W_1} \end{matrix} \times \begin{matrix} & n \\ r & \boxed{W_2} \end{matrix}$$

- W: Input matrix
 - $m \times n$ matrix
- V ($n \times r$ matrix) contains the orthonormal eigenvectors of $W^T W$
- U ($m \times r$ matrix) contains the orthonormal eigenvectors of $W W^T$
- R: Singular value matrix
 - $r \times r$ diagonal matrix, r is the rank of W

Before: mn

After: $mr + rn = r(m+n) = \min(m,n)(m+n)$
assume W is full rank

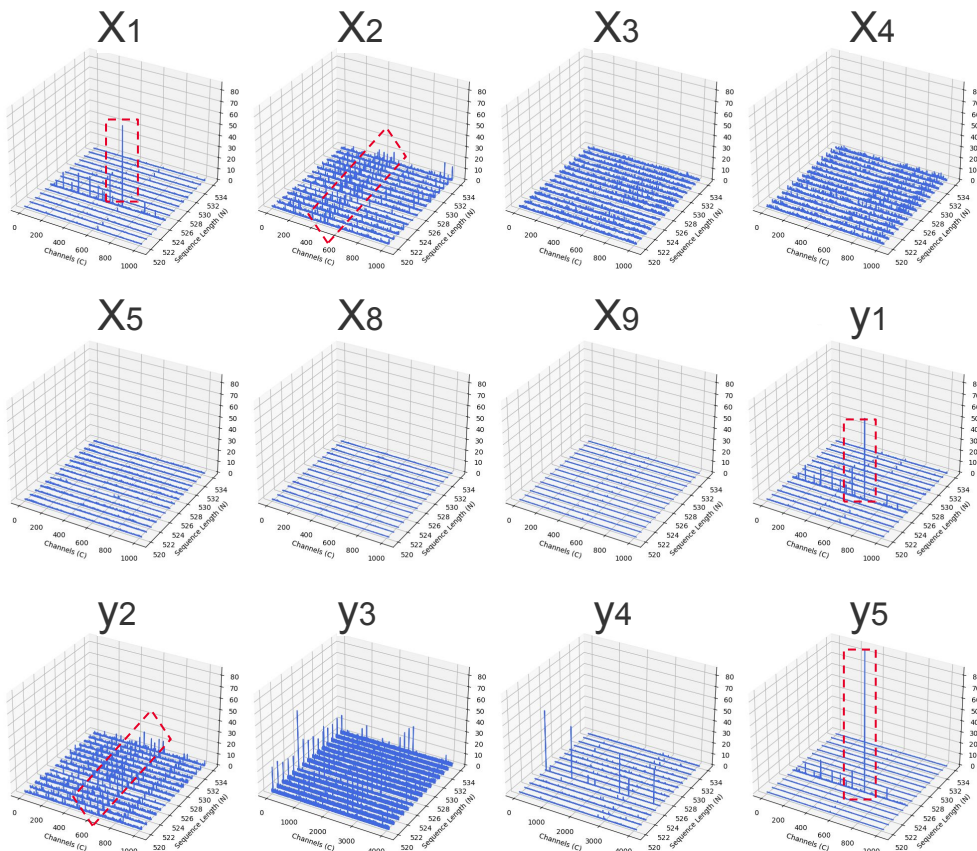
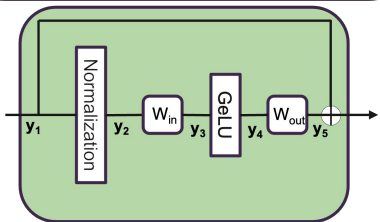
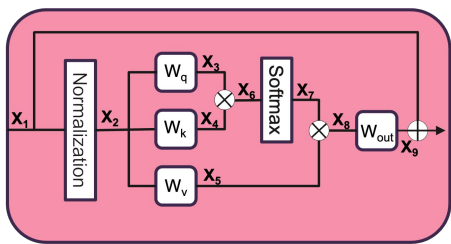
- Without rank truncation, the number of parameters increases.

Recap (Lecture 7: P1-P41, P53-P68)

- Outlier Distribution in LLM
 - Massive activation
 - Channelwise outlier
- Quantization and smoothing techniques for large models
- LLM Pruning
- Low-rank decomposition for LLM

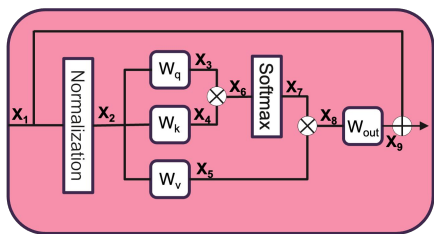
Outlier Study: CLIP Activations

- Outliers with large magnitudes appear at positions x_1 , y_1 , and y_5 , referred to as **massive activations**.

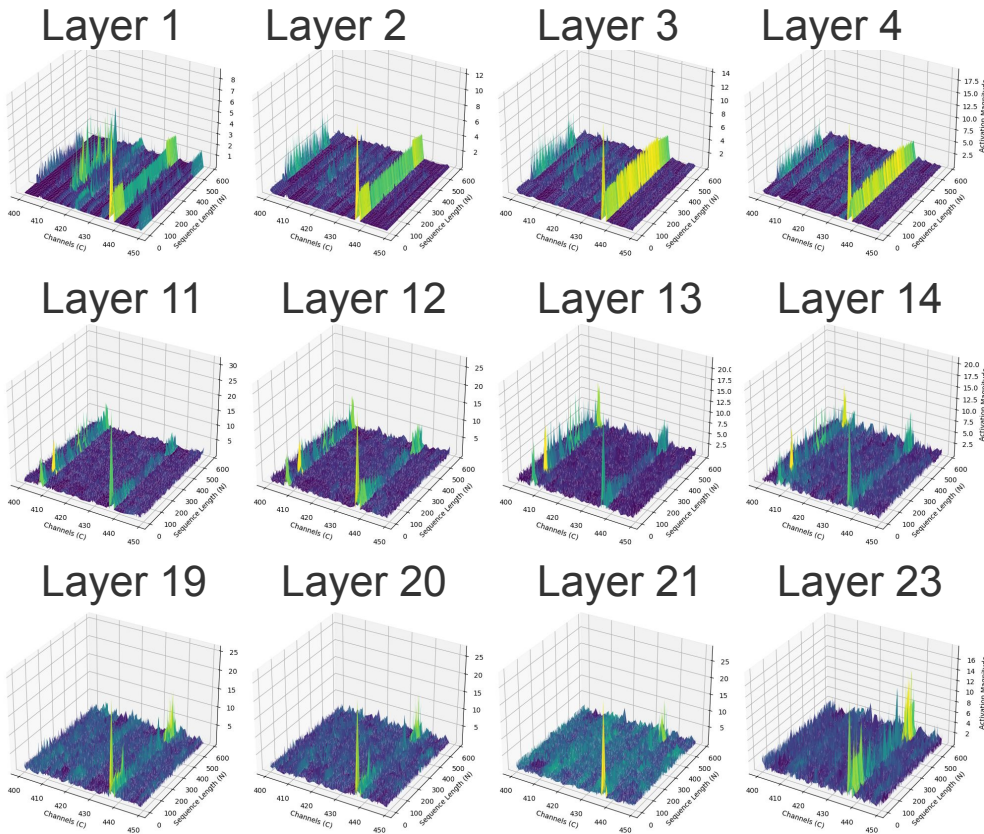


Outlier Study: CLIP Activations

- 3D plots of x_2 across layers.



- x_2 exhibits channel wise outlier

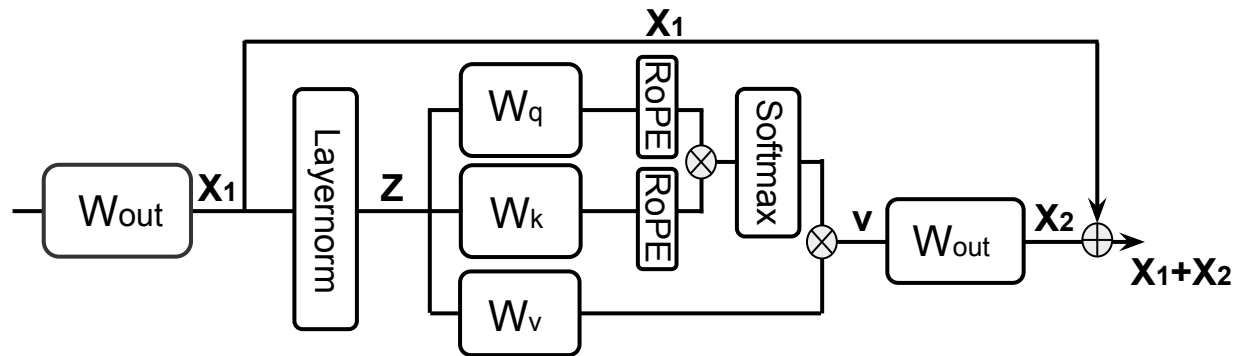


Impact of Massive Activation

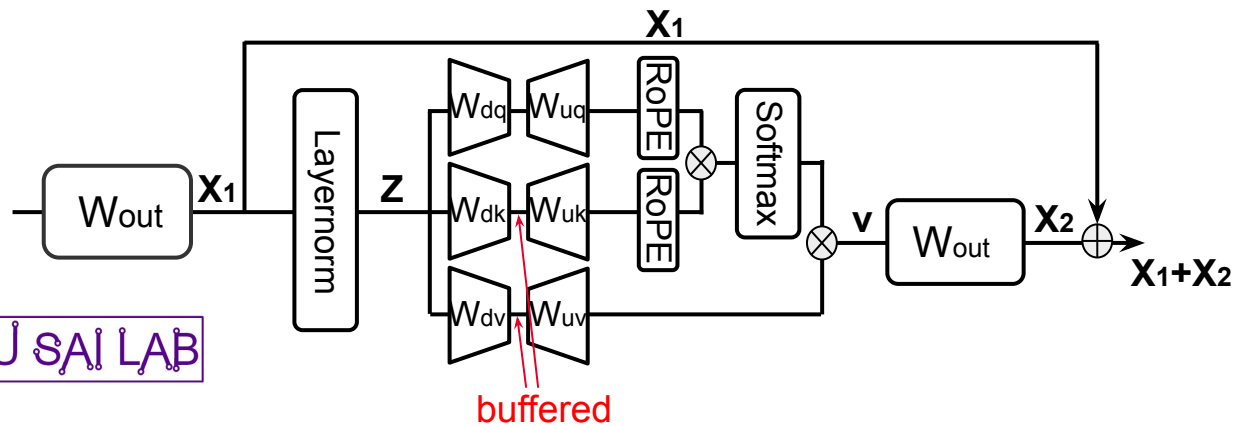
Intervention	LLaMA3.2-3B		LLaMA3.1-8B		LLaMA2-13B		GPT-2		Qwen2.5-7B	
	WikiText	C4	WikiText	C4	WikiText	C4	WikiText	C4	WikiText	C4
<i>Original</i>	5.567	10.790	6.941	9.046	4.355	6.405	14.795	19.460	6.520	11.773
<i>TMA to mean at y_7</i>	1124111.75	21046.82	21281.49	1301562.25	1301562.25	6469.42	14.841	19.560	71216.17	66588.86
<i>TMA to zeroes at y_7</i>	1138151.23	21951.41	21601.10	1302018.53	1309211.61	7128.32	14.911	19.928	71835.61	67518.35

- The truncation of massive activation will cause the significant accuracy degradation of the LLM.
- Massive activations also occur in other types of foundation models that utilize attention-based architectures.

Singular Value Decomposition



SVD decomposition can potentially save the MAC operations, memory storage and KV cache size.

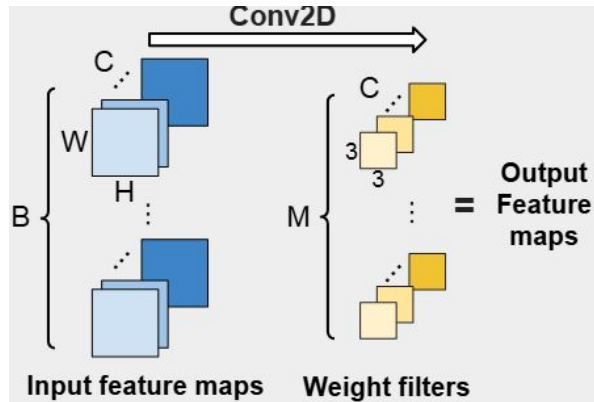


Recap (Lecture 8: P1-P50)

- Efficient training of DNNs
 - Efficient computing
 - Efficient storage
- Parameter efficient finetuning
- Basics of Speculative Decoding

Example Question

Figure below describes a standard Conv2D operation where the input has dimensions $B \times C \times H \times W$, and the weight filters have sized $M \times C \times 3 \times 3$. Here, B represents the batch size, C is the number of input channels, and H and W denote the spatial dimensions of the input feature map. Assume a padding size and stride of 1. No bias terms is involved.



1. How many multiply-accumulate (MAC) operations are required for conventional Conv2D?
2. If $p\%$ of the filters within the weight are pruned (meaning every elements within the weight filters are set to zero), how does this affect the computational cost of the convolution? Additionally, what is the new total of MAC operations for the convolution?
3. What is the size of the output generated by the batch normalization operation? During inference, is the BN operation explicitly computed? What if changing to layer normalization?

Example Question

as illustrated in figure below. The input to the SA layer is shaped as $B \times L \times E$, where B represents the batch size, L denotes the token length, and E is the embedding size. Each of the weight matrices W_Q , W_K , and W_V has dimensions $E \times E$, and the number of attention heads is assumed to be 1.

1. Could you explain the purpose and underlying intuition of the QK^T operation in the self-attention mechanism?
2. What is the total count of MAC operations in this SA block? Please express it using the terms B , L , and E . How does the computational complexity evolve as the token length L increases? You don't need to account for the cost of nonlinear operations within softmax, layernorm, and scaling operations.

