

Midterm Coverage

General Information

- In-class exam on Oct 29 7:30pm-8:45pm
- 75 mins
- Close book, you may bring calculator
- I will join online, TAs will invigilate



Recap (Lecture 1: P1-P32)

- Neural Network Basics
 - MLP
 - Forward and backward propagation of MLP
 - Weight decay, dropout
 - The training optimizer: SGD, RMSProp, Adam
 - Multistage learning rate scheduler
- Lecture 0 is not covered.

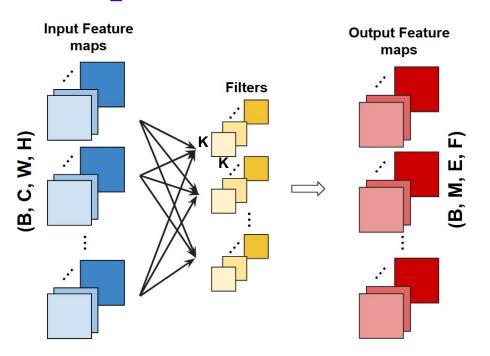


Recap (Lecture 2: P1-P71)

- 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



Computational Cost: Standard Convolution



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

B: batch size

C: number of input channels

H,W: size of the input feature maps

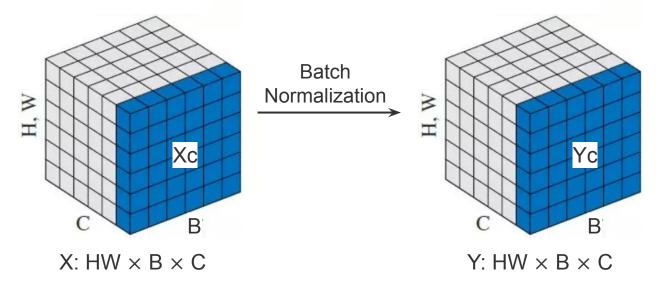
M: number of weight filters

K: weight kernel size

E,F: size of the output feature maps



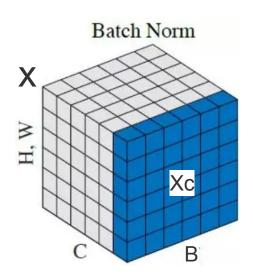
Batch Normalization



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



Batch Normalization



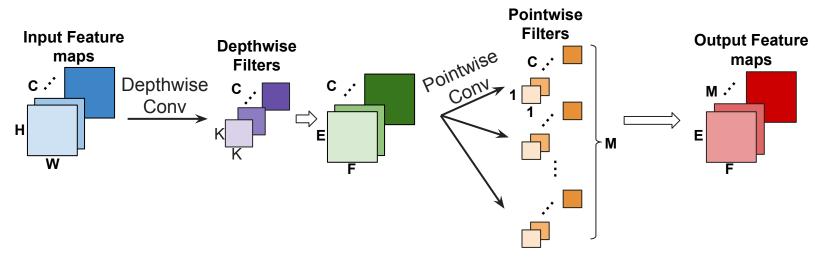
$$X: HW \times B \times C$$

$$egin{aligned} Y_c &= lpha_c rac{X_c - \mu_c}{\sigma_c} + eta_c & ext{For each c} \in \mathbb{C} \ lpha &= \{lpha_c\}, eta = \{eta_c\}, \mu = \{\mu_c\}, \sigma = \{\sigma_c\} \end{aligned}$$

- For each channel c, we have:
 - Xc: (HW x B)
 - \circ μ_c and δ_c are the mean and standard deviation of Xc.
 - o αc and βc are learnable parameters
 - αc, βc, μc, δc are scalers
- Overall, we have:
 - \circ μ, δ, α and β all have a length of C
 - \circ μ , δ , α and β are all fixed during the inference
 - \circ μ , δ are statistics based on the training dataset



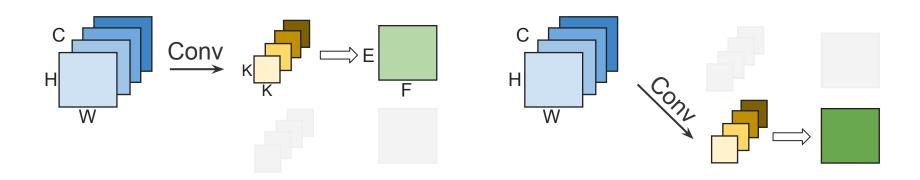
Depthwise Separable Convolution



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



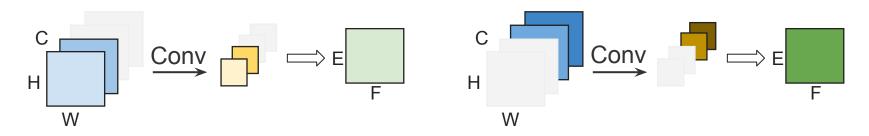
Group Convolution



The original MAC: E×F×K×K×C×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×F×K×K×C×M/G

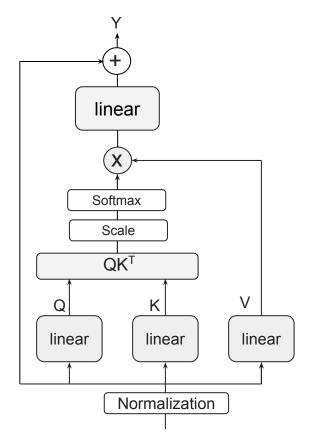
Recap (Lecture 3: P1-P69, P77-P79)

- 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).
 - $\circ \quad (B,L,E) * (E*E) \rightarrow (B*L*E)$
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
 - \circ QK^T \rightarrow (B, L*E) * (B,E*L) \rightarrow (B, L*L)
- Scale and normalize the score using softmax.
 - \circ Softmax(QK^T) \rightarrow (B, L*L)
- Multiply each value vector by the softmax score.
 - Softmax(QK^T) ***** V
 - \circ (B, L*L) * (B, L*E) \rightarrow (B, L*E)
- Pass the result to the linear layer, sum with the input.

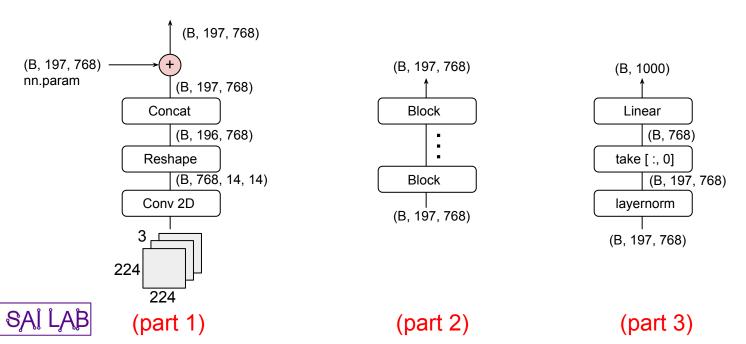




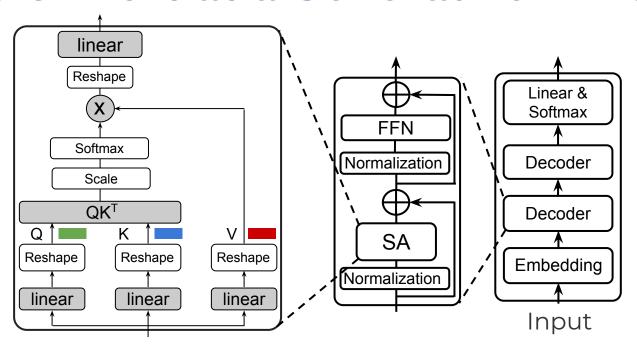
Vision Transformer

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

13



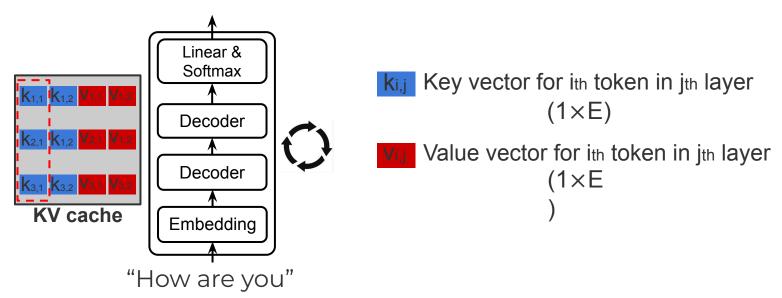
Transformers as a Generative AI Tool





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

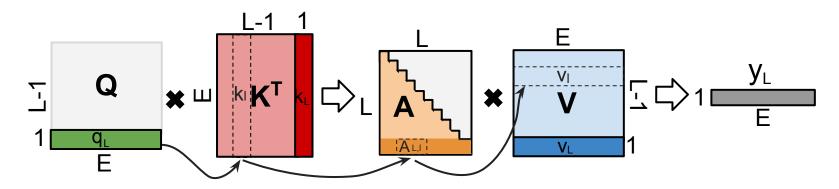
GPT-2: Prefilling



 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.

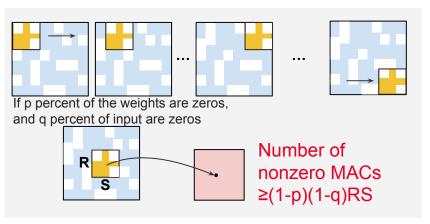


Recap (Lecture 4: P1-P46, P50-P65, P68-P79)

- 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≥(1-p)×(1-q)×M×C×R×S×E×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)



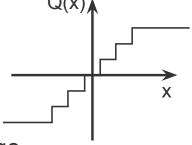
Recap (Lecture 5: P1-P65)

- 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
- Qunatization during training
 - Stochastic quantization
- Learnable adaptive quantization scheme



Fixed-Point Format (Symmetrical)

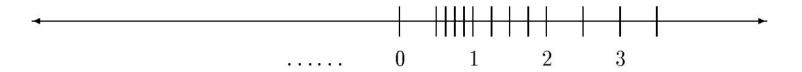
- How to convert a number x to INT representation?
 - Set the clipping range: (-L, L), bitwidth: b
 - Compute the scale: $s = 2L/(2^b 2)$
 - Clip the input x: $x_c = Clip(x, L, -L)$
 - \circ Calculate the INT representation: $x_{int} = round(x_c/s)$
 - \circ Rescale: $x_q = sx_{int}$



- Have a uniform representation power within the clipping range.
- All the computations can be performed using x_{int}



Floating Point Arithmetic



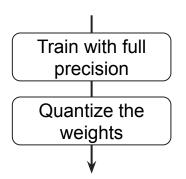
- Have better representation power for values with small magnitudes.
- How to convert a real number x to FP representation?

$$egin{aligned} \mathsf{x} = |\mathsf{x}| & \mathsf{s} = \mathsf{sign}(\mathsf{x}) \ a = \lfloor log_2 x
floor & e = a + bias & m = rac{x}{2^a} - 1 \end{aligned}$$

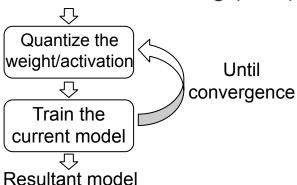


When to Quantize?

Post-training quantization (PTQ)



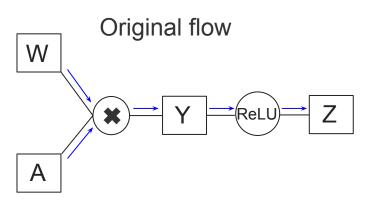
Quantization-aware Training (QAT)



- 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.



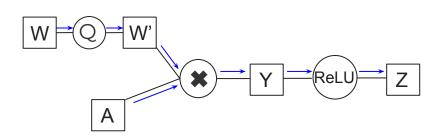
Another Way to Look at Quantization



$$Y = WA, Z = ReLU(Y)$$

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

Flow with quantization

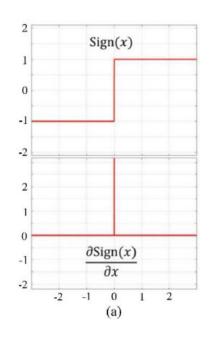


$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W'} \frac{\partial W'}{\partial W}$$

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 makes 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.



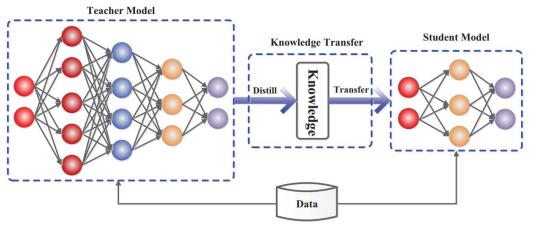
Recap (Lecture 6: P1-P35, P56-P63)

- Distillation
 - Feature-Based Knowledge Distillation
 - Online distillation
 - Self distillation
 - Multi-teacher, multi-student, cross-modal
- Low-rank Decomposition



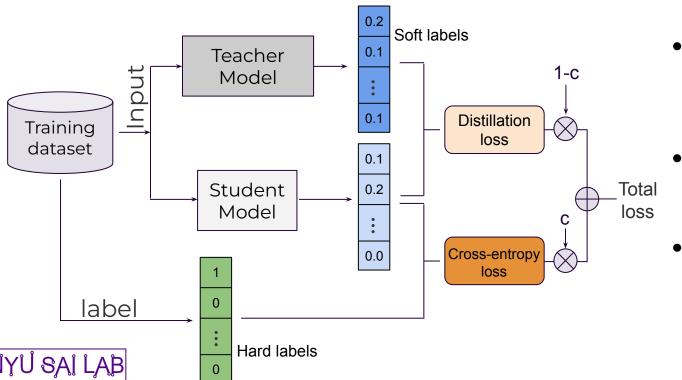
Knowledge Distillation Basics

 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 DKL(P // 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.

$$D_{\mathrm{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \; \log iggl(rac{P(x)}{Q(x)}iggr)$$

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



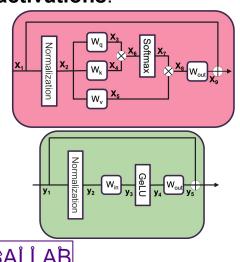
Recap (Lecture 7: P1-P41)

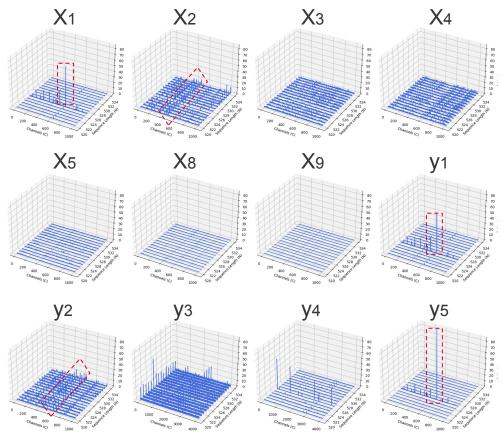
- Outlier Distribution in LLM
 - Massive activation
 - Channelwise outlier
- Quantization and smoothing techniques for large models



Outlier Study: CLIP Activations

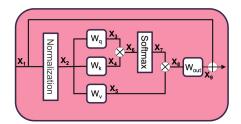
 Outliers with large magnitudes appear at positions x1, y1, and y5, referred to as massive activations.



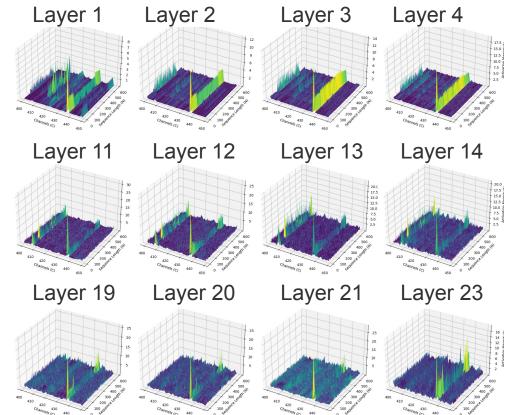


Outlier Study: CLIP Activations

 3D plots of x2 across layers.



 x2 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
Original	5.567	10.790	6.941	9.046	4.355	6.405	14.795	19.460	6.520	11.773
TMAs to mean at y7	1124111.75	21046.82	21281.49	1301562.25	1301562.25	6469.42	14.841	19.560	71216.17	66588.86
TMAs to zeroes at y7	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.



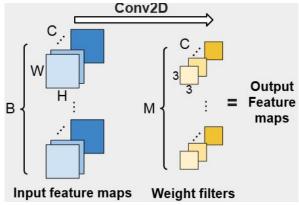
Recap (Lecture 8: P1-P45)

- Efficient training of DNNs
 - Efficient computing
 - Efficient storage
- Parameter efficient finetuning



Example Question

Figure below describes a standard Conv2D operation where the input has dimensions B x C x H x W, and the weight filters have sized M x C x 3 x 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. Will removing weight filters reduce the computational cost of the subsequent Conv2D layers, assuming that only ReLU layers exist between consecutive Conv2D layers? Why?

