

Lecture 04: Pruning Strategies for Efficient DNN Implementation

Notes

- Lab1 will be released this week!
- Lab0 will be posted to help you understand DNN pruning.



Recap

- Transformer basics
- Bert
- Vision transformer
- Large Language Model
- Self-supervised learning

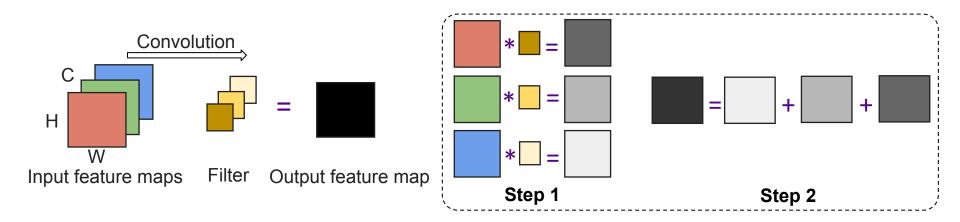


Topics

- Why pruning?
 - Reduce running cost
 - Reduce storage
- General pruning techniques
- Transformer pruning
- Large model pruning



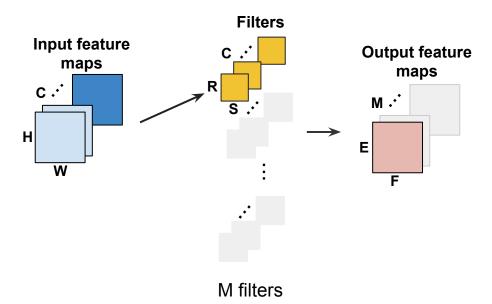
Convolutional Layers



Core building block of a CNN, it is also the most computational intensive layer.

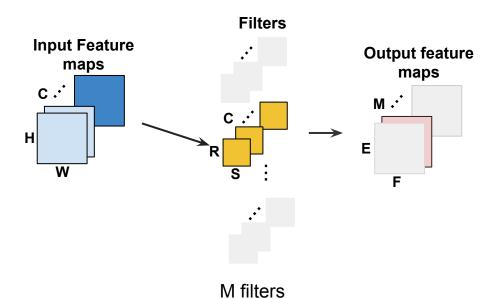


Convolution



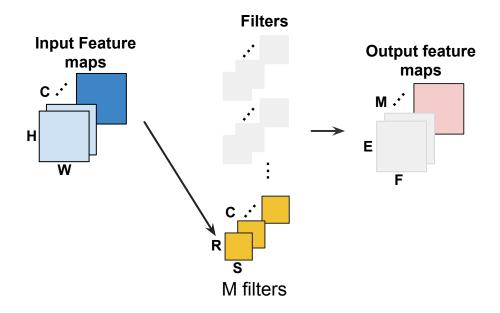


Convolution



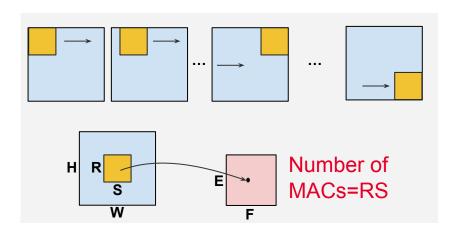


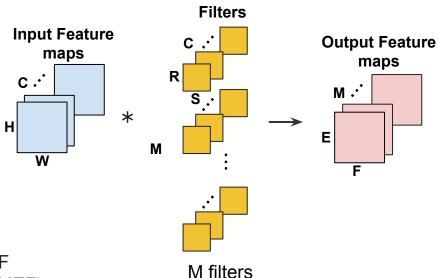
Convolution





Computational Cost of Convolution

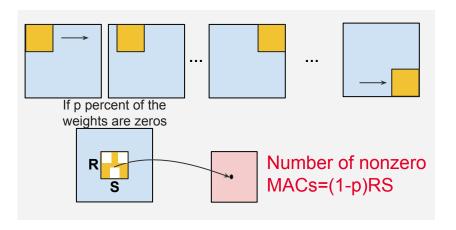




- Number of MACs: M×C×R×S×E×F
- Storage cost: 32×(MCRS+CHW+MEF)
- The input activation and output activations are transient storage, can be eliminated once this layer is finished processing.



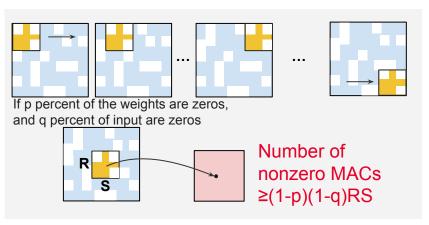
Convolution with Sparse Weight



- Number of MACs: (1-p)×M×C×R×S×E×F
- Weight pruning can reduce the computations.
- Sparse weight matrices can be stored more efficiently, which helps minimize memory usage.



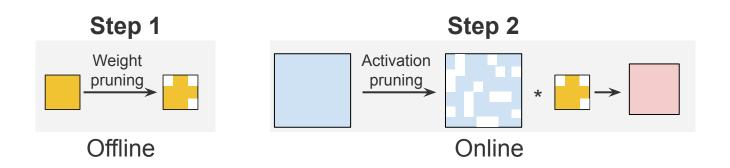
Convolution with Sparse Weight



- 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 usage.



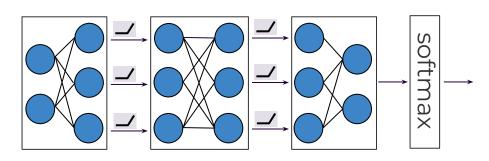
Convolution with Sparse Weight

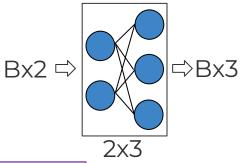


 Activation sparsity requires online pruning, which leads to additional overhead for sorting.



Computational Cost for MLP

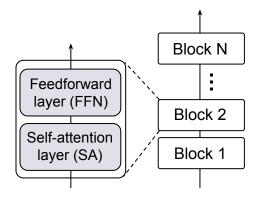




- B is the batch size
- Number of MACs:
 - \circ Bx2x3 = 6B
- Storage cost:
 - 6 x 32 = 192 bits (Weights)
 - $(2B + 3B) \times 32 \text{ bits (Activation)}$

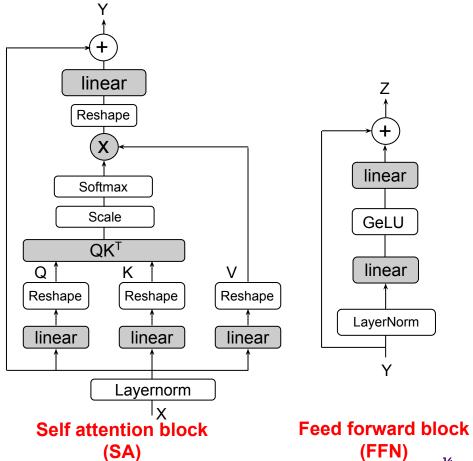


Transformers

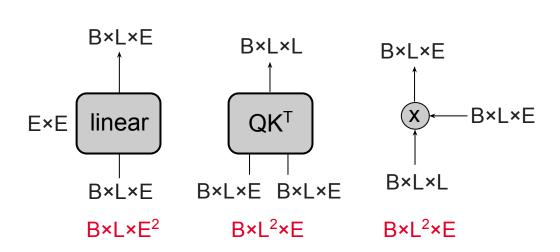


- The input sentence has three dimensions:
 - B: batch
 - L: sequence length (number of words)
 - E: embeddings 0



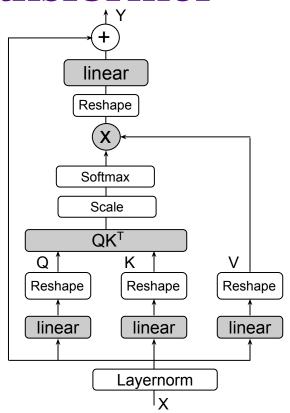


Computational Cost of Transformer

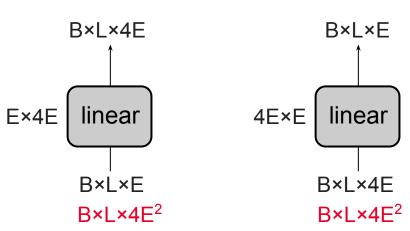


Total =
$$4B \times L \times E^2 + 2B \times L^2 \times E$$



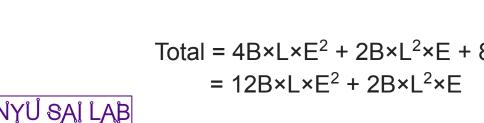


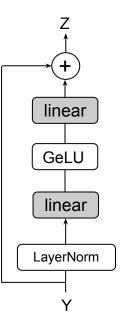
Computational Cost of Transformer



Total =
$$4B \times L \times E^2 + 2B \times L^2 \times E + 8B \times L \times E^2$$

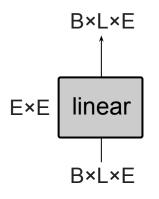
= $12B \times L \times E^2 + 2B \times L^2 \times E$

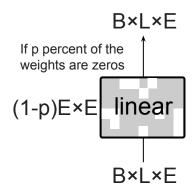


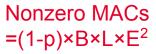


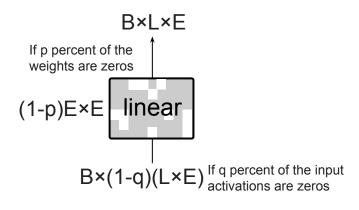


Computational Cost of Transformer









Nonzero MACs
$$\geq (1-p)(1-q)\times B\times L\times E^2$$

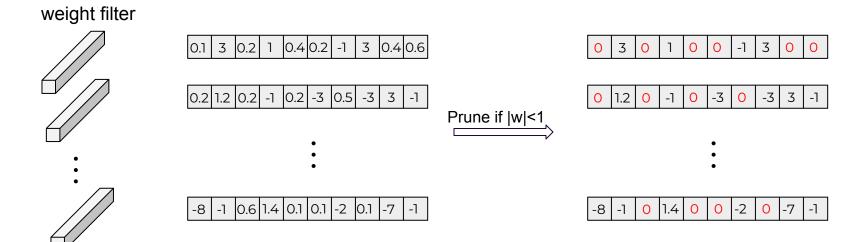


Topics

- Why pruning?
 - Reduce running cost
 - Reduce storage
- General pruning techniques
- Transformer pruning
- Large model pruning



Pruning

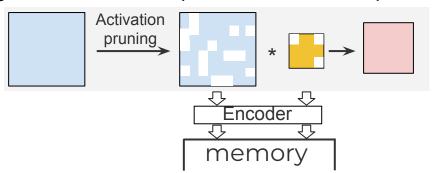


Pruning reduces both computational demands and storage costs.



Benefit of Pruning

- Reduce computational complexity
 - To support sparse matrix with random sparsity pattern, specialized hardware is required.
- Reduce the storage complexity
 - To achieve it, we need to encode the sparse weights.
 - Encoding the activations requires additional computation.



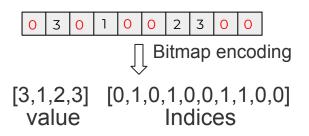


Sparse Matrices Encodings

- Efficient encoding scheme for sparse matrix storage.
 - o Bitmap
 - Run Length Encoding (RLE)
 - Coordinate format (COO)
 - Compressed sparse row (CSR), Compressed sparse column (CSC)



Bitmap Encoding

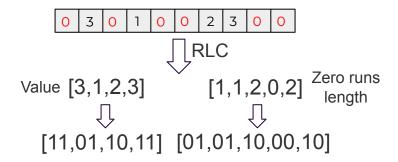


- In summary, the storage cost of bitmap encoding (in bits) is: (1-p)×L×n + L
- n: number bits per value
- L: number of elements
- p: sparsity (%)
- Bitmap is effective for compressing the tensors of low or moderate sparsity.
- Encoding cost is low.



Run Length Encoding (RLC)

- Record the values and length of zero runs between the values.
- Assume 2 bits are used to encode the length of zero runs (0-3).
- Each value requires 2 bits.

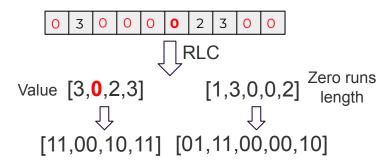


RLC can reduce storage requirement when sparsity is moderate.



Run Length Encoding (RLC)

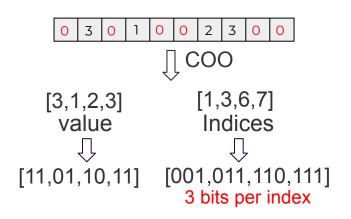
- Record the values and length of zero runs between the values.
- Assume 2 bits are used to encode the length of zero runs (0-3).
- Each value requires 2 bits.



- RLC can reduce storage requirement when sparsity is moderate.
- Hard to formulate it analytically.



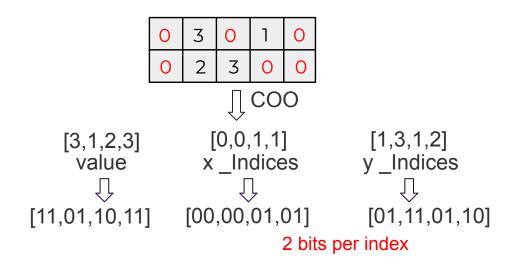
Coordinate Format (COO)



- COO is efficient with the sparsity level is extremely high.
- The storage cost (in bits) is: (1-p)×L×n + (1-p)
 ×L×Ceil(log₂L)
- n: number bits per value
- L: number of elements
- p: sparsity (%)



Coordinate Format (COO)



COO is efficient with the sparsity level is extremely high.



Compressed Sparse Row/Column (CSR/CSC)

```
\begin{pmatrix}
10 & 20 & 0 & 0 & 0 & 0 \\
0 & 30 & 0 & 40 & 0 & 0 \\
0 & 0 & 50 & 60 & 70 & 0 \\
0 & 0 & 0 & 0 & 0 & 80
\end{pmatrix}

v = [10 & 20 & 30 & 40 & 50 & 60 & 70 & 80 ]

COL_INDEX = [0 & 1 & 1 & 3 & 2 & 3 & 4 & 5 ]

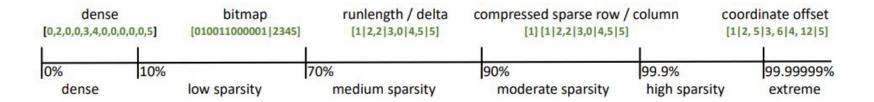
ROW_INDEX = [0 & 2 & 4 & 7 & 8 ]
```

Encoded form

- CSR/CSC is also suitable for matrices with high sparsity.
- The row index specifies the amount of nonzero values within each row.



Encoding Approaches Tailored to Different Sparsity Levels



• Different encoding scheme can be applied for different sparsity levels.



Topics

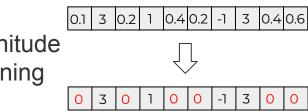
- Why pruning?
 - Reduce running cost
 - Reduce storage
- General pruning techniques
- Transformer pruning
- Large model pruning



Pruning Criteria: Magnitude Pruning

- We can prune the weights using the importance score:
 - Magnitude
 - Gradient
 - Hessian

Magnitude pruning



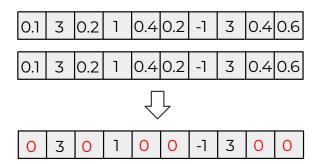
$$m_i = \begin{cases} 1 : if \ ||w_i||_1 \ge a \\ 0 : if \ ||w_i||_1 < a \end{cases}$$



Drawbacks of Magnitude Pruning

 The major drawback of magnitude based pruning is that it does not consider the impact of the input when making the pruning decision.

Magnitude pruning

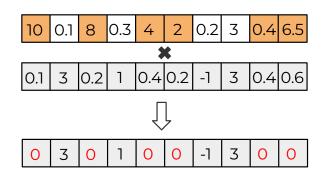




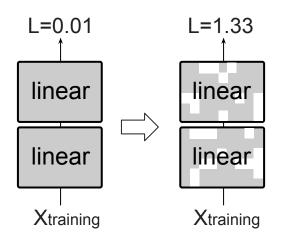
Drawbacks of Magnitude Pruning

 The major drawback of magnitude based pruning is that it does not consider the impact of the input when making the pruning decision.

Magnitude pruning







$$\mathcal{L}(w) = \sum_{(x,y) \in D_{training}} l(y, F_w(x))$$

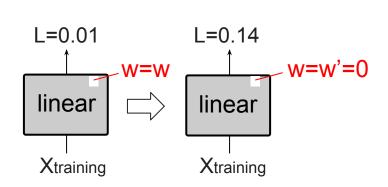
Dtraining: is the training dataset

Fw(.): neural network function with parameter w.

L(.): Training loss function

- Another pruning principle is to minimize the impact on training loss as much as possible.
- For a trained DNN, L(.) remains low.
- Training loss is typically computed over either a subset or the entire training dataset.





$$\mathcal{L}(w) = \sum_{(x,y) \in D_{training}} l(y, F_w(x))$$

Dtraining: is the training dataset

Fw(.): neural network function with parameter w.

I(.): loss function

- Another pruning principle is to minimize the impact on training loss as much as possible.
- For a trained DNN, L(.) remains low.



- Pruning criteria:
 - Keep the change on training loss as small as possible
 - Let L(.) denote the training loss
 - For trained DNN, L(.) will be low.

$$\mathcal{L}(w') = \mathcal{L}(w) +
abla \mathcal{L}(w)^T (w-w')$$

• If w is pruned, then we have w'=0:

$$\mathcal{L}(0) - \mathcal{L}(w) = rac{d\mathcal{L}}{dw}w$$

- We can use it as the pruning criteria
 - Sort the weight based on the product of gradient and its value.



$$egin{aligned} \mathcal{L}(w') &= \mathcal{L}(w) +
abla \mathcal{L}(w)^T (w-w') \ &+ rac{1}{2} (w-w')^T
abla^2 \mathcal{L}(w) (w-w') \end{aligned}$$

 When reflecting on each individual value, the pruning criteria becomes:

$$\mathcal{L}(0) - \mathcal{L}(w) = rac{d\mathcal{L}(w)}{dw}w + rac{1}{2}rac{d^2\mathcal{L}(w)}{dw^2}w^2$$

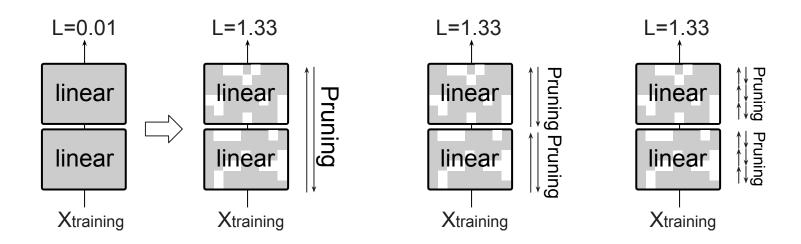
The gradient is usually estimated to zero.

- Multiple approaches have been propose to estimate the Hessian:
 - Empirical Fisher (Outer Product of Gradients)

$$H(heta) pprox rac{1}{N} \sum_{i=1}^N g_i g_i^ op, \quad g_i =
abla_ heta \log p(x_i \mid heta)$$

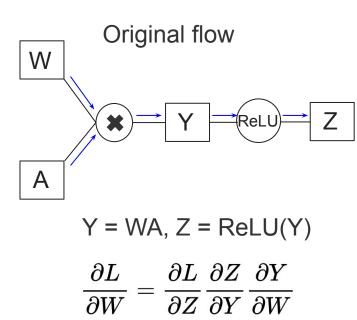


Granularity of Pruning

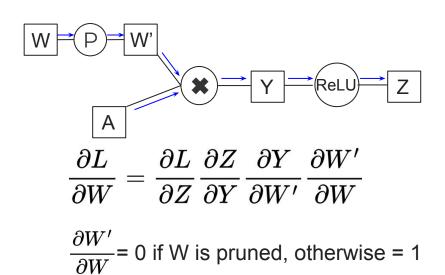




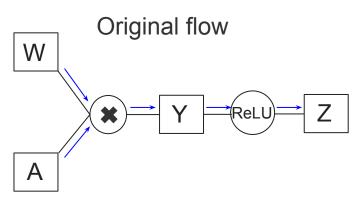
Computational Flow of Pruning



Flow with pruning



Computational Flow of Pruning

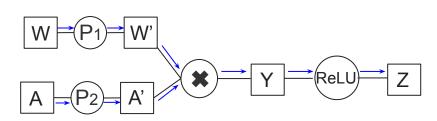


$$Y = WA, Z = ReLU(Y)$$

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

NYU SAI LAB

Flow with pruning

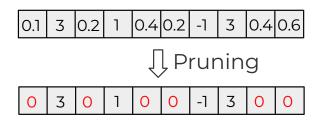


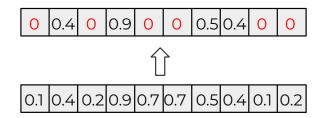
$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial Y} \frac{\partial Y}{\partial W'} \frac{\partial W'}{\partial W}$$

$$\frac{\partial W'}{\partial W}$$
 = 0 if W is pruned, otherwise = 1

$$rac{dL}{dA} = rac{dL}{dZ}rac{dZ}{dY}rac{dY}{dA}rac{dA'}{dA}$$

Backward Pass of Pruning Operation





• During backward pass, the gradients of the pruned elements are masked.



Pytorch Implementation of Iterative Pruning

```
def forward(self, x):
    y = F.conv2d(self.w, x)
    return y
```

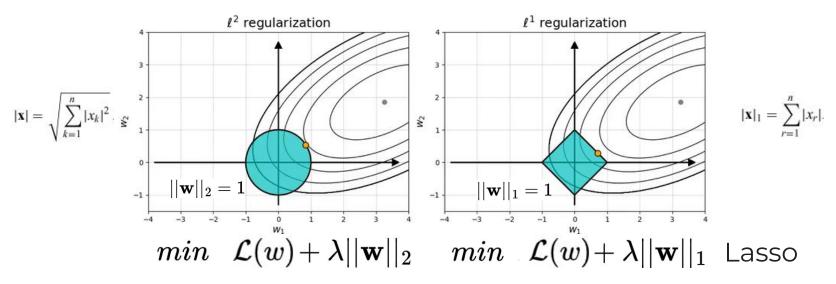
Fake pruning to stimulate the impact of sparse weight on the model accuracy!

```
mask = nn.Parameter(... requires_grad=False)
For w in each layer:
    mask = mask_Prune(w, mask, percent)
    w = w * mask
...
def forward(self, x):
    y = F.conv2d(self.w, x)
    return y
```



Regularization-based Pruning

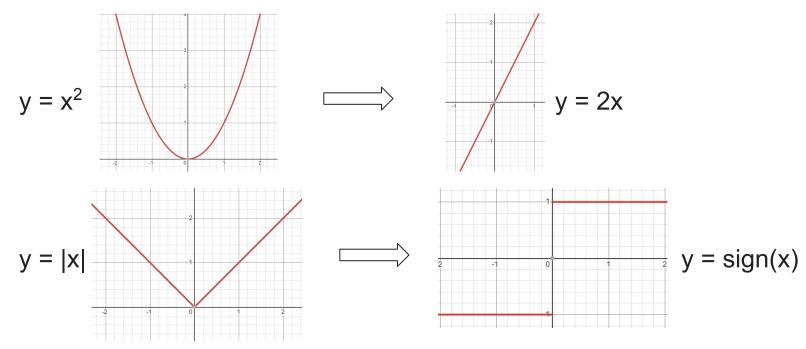
 ℓ^1 induces sparse solutions for least squares



Add this term can make DNN naturally select the unimportant weight during the training process.



Regularization-based Pruning



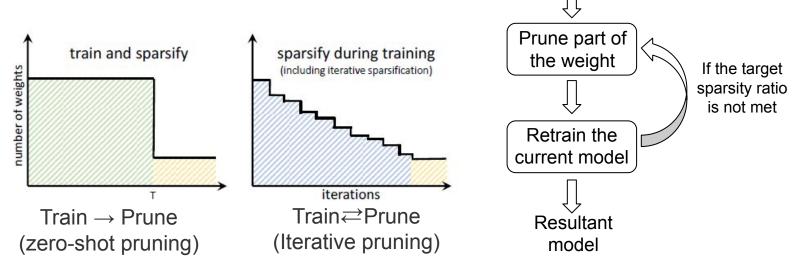


Taxonomy of Pruning

- Pruning techniques can be classified from different perspectives
 - Iterative pruning, zero-shot pruning
 - Structured pruning, unstructured pruning, N:M pruning
 - Weight pruning, activation pruning
 - Static pruning and dynamic pruning
 - Pruning for inference, pruning for training



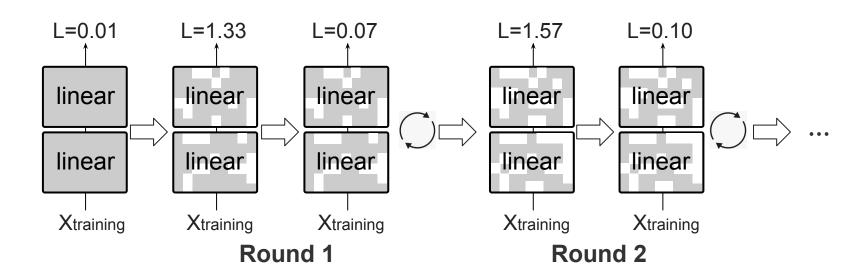
When to Prune?



- Usually interactive pruning has the best accuracy performance, however, it also requires multiple rounds of training and computational cost.
- Zero-shot pruning also termed post-training pruning.



Iterative Pruning





Lottery Ticket Hypothesis

"A randomly-initialized, dense neural network contains a subnetwork that is initialized such that—when trained in isolation—it can match the test accuracy of the original network after training for at most the same number of iterations."



How to Find the Winning Tickets?

Iterative Magnitude Pruning (IMP):

- Initialized DNN with random weights wo.
- While the sparsity level has not reached:
 - Train the DNN with k epochs until convergence
 - prune p% of the nonzero weights.
 - Reinitialize the remaining weights using the values in w₀, finetune the remaining weights for k epochs (Rewind).
- Return the weights.
- Later work has shown that rewind to wi (i is small) works better for larger networks.

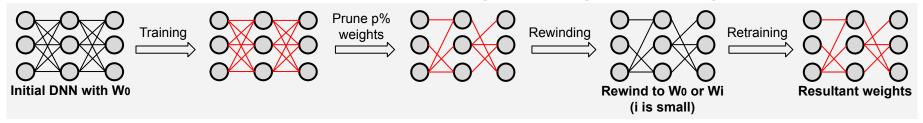


Weight Rewinding

Conventional iterative pruning



Conventional iterative pruning with weight rewinding



• The pruned architecture itself, rather than a set of inherited "important" weights, is more crucial to the accuracy in the final model, which suggests that in some cases pruning can be useful as an architecture search paradigm.

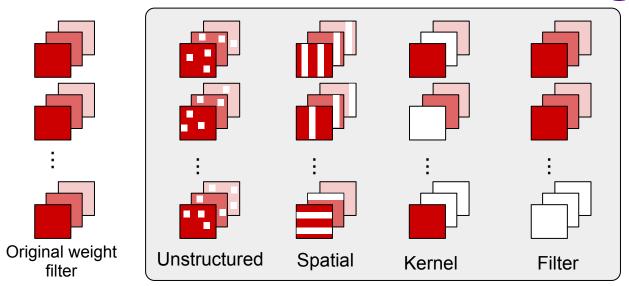


Taxonomy of Pruning

- Pruning techniques can be classified from different perspectives
 - Iterative pruning, zero-shot pruning
 - Structured pruning, unstructured pruning, N:M pruning
 - Weight pruning, activation pruning
 - Static pruning and dynamic pruning
 - Pruning for inference, pruning for training



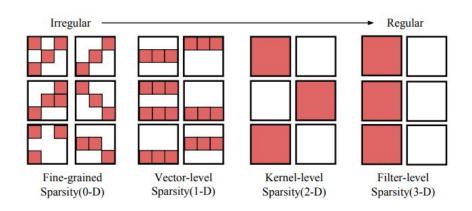
Unstructured/Structured Pruning

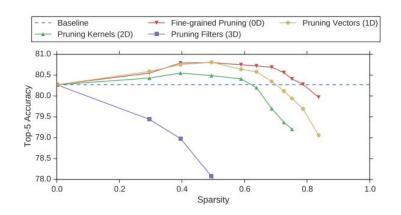


• Structured pruning is amenable to hardware performance, due to the regular sparsity distribution.



Unstructured/Structured Pruning





- Unstructured sparsity has a better accuracy than structured sparsity.
- We can apply the same method as the unstructured pruning to prune a group of parameters.

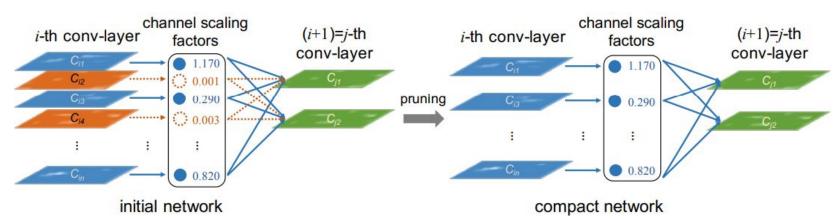


Structured Pruning

- In magnitude-based weight filter pruning, we first compute the sum of absolute weight values within each filter, then rank the filters and prune those with the smallest sums.
- Similarly, structured pruning can be performed by assessing the importance of a weight group through the sum of their importance scores.



Network Slimming

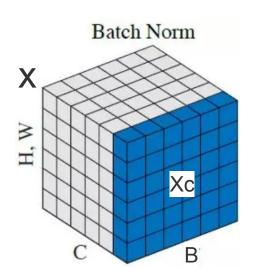


• We associate a scaling factor (from a batch normalization layer) with each filter in convolutional layers. Sparsity regularization is imposed on these scaling factors during training to automatically identify unimportant filters.

$$L = \sum_{(x,y)} l(f(x,W),y) + \lambda \sum_{\mathbf{i}} |\mathbf{pi}|$$



Batch Normalization



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

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

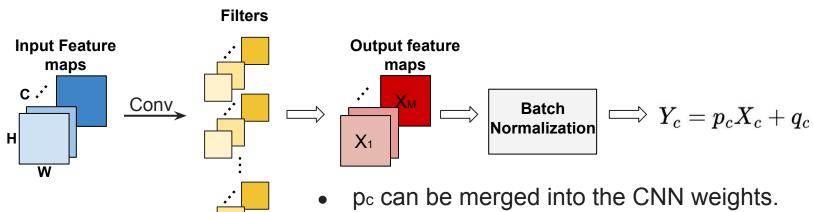
- For each channel c, we have:
 - o Xc: (HW x B)
 - \circ μ_c and δ_c are the mean and standard deviation of Xc.
 - o αc and βc are learnable parameters
 - αc, βc, μc, δc are scalers
- Overall, we have:
 - \circ μ, δ, α and β all have a length of C
 - \circ μ , δ , α and β are all fixed during the inference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



Batch Normalization: During Inference

• Given all the parameters are fixed, for each channel c, we have:

$$Y_c = lpha_c rac{X_c - \mu_c}{\sigma_c} + eta_c = rac{lpha_c}{\sigma_c} X_c + (eta_c - rac{lpha_c \mu_c}{\sigma_c}) \implies Y_c = p_c X_c + q_c$$



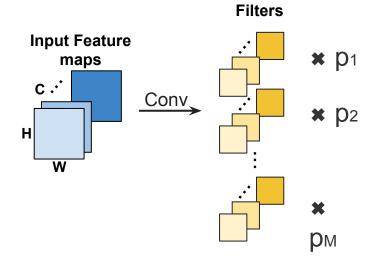
qc can be merged into the CNN bias.



Batch Normalization

For each channel c, we have:

$$Y_c = lpha_c rac{X_c - \mu_c}{\sigma_c} + eta_c = rac{lpha_c}{\sigma_c} X_c + (eta_c - rac{lpha_c \mu_c}{\sigma_c}) \;\; \Longrightarrow Y_c = p_c X_c + q_c$$



We can fold in the p and q to the weights and bias of convolutional layer during inference and reduce the online computational cost.

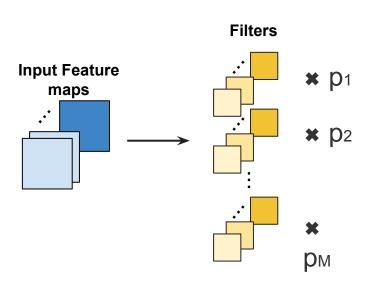


Network Slimming

- Lasso regularization is imposed on the scaling factors of batch normalization during training to automatically identify unimportant channels.
- g(.) is the lasso I1-norm g(.) = Σ|pi|

$$L = \sum_{(x,y)} l(f(x,W),y) + \lambda \sum_{\mathbf{i}} |\mathbf{pi}|$$

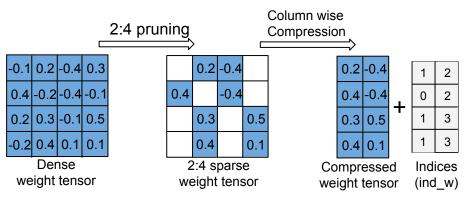
 The unimportant channel are naturally eliminated during the training process.



Add a lasso loss on pi



N:M Sparsity

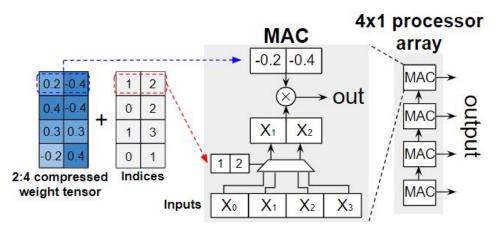


- DNN with structured sparsity can be easily adopted for acceleration, but incur low accuracy.
- On the other hand, DNN with unstructured sparsity is hard to accelerate.

- N:M sparsity is proposed as a middleground between structured and unstructured sparsity.
- 2:4 sparsity is supported in Nvidia V100 GPUs.



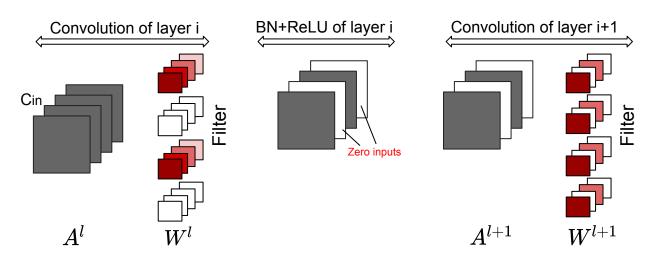
N:M Sparsity



- N:M sparsity is proposed as a middleground between structured and unstructured sparsity.
- 2:4 sparsity is supported in Nvidia V100 GPUs.



Cascade Effect of Filterwise Pruning in CNN



- Assume the bias of the batch normalization is zero.
- Filter pruning at layer I can also result in weight and input sparsity in layer I+1.
- When the bias is not zero, the feature maps of layer i+1 will contain a uniform constant value.



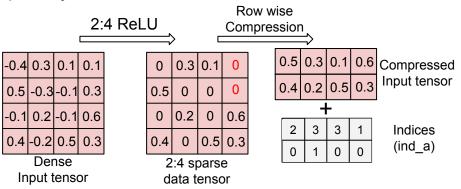
Taxonomy of Pruning

- Pruning techniques can be classified from different perspectives
 - Iterative pruning, zero-shot pruning
 - Structured pruning, unstructured pruning, N:M pruning
 - Weight pruning, activation pruning
 - Static pruning and dynamic pruning
 - Pruning for inference, pruning for training



Pruning on Input Activation

- Why pruning can not be applied to input activation?
 - Large computing cost to determine the importance scores.
 - Due to the usage of ReLU, activation in CNN are 50% sparse, but with irregular sparsity distributions.





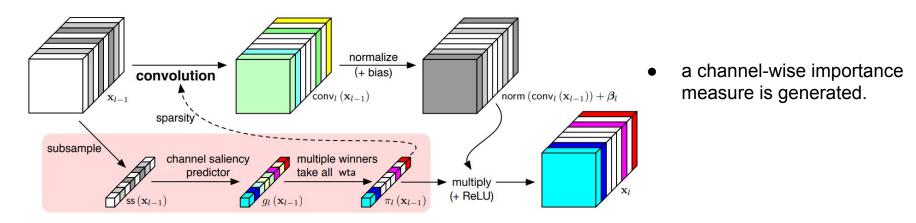
Taxonomy of Pruning

- Pruning techniques can be classified from different perspectives
 - Iterative pruning, zero-shot pruning
 - Structured pruning, unstructured pruning, N:M pruning
 - Weight pruning, activation pruning
 - Static pruning and dynamic pruning
 - Pruning for inference, pruning for training



Static vs Dynamic Pruning

- Conventional pruning adopts static pruning criteria and permanently removes components.
- Dynamic pruning exploits input-specific characteristic pruning criteria and preserves the entire network structures and accelerates the networks by dynamically skipping unimportant components.



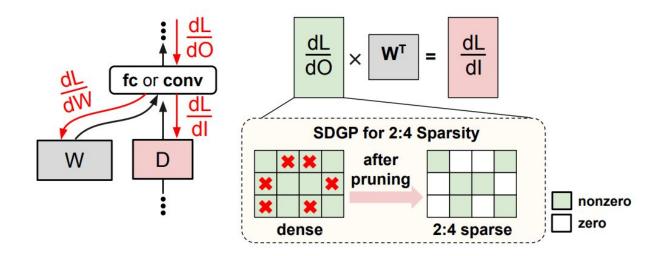


Taxonomy of Pruning

- Pruning techniques can be classified from different perspectives
 - Iterative pruning, zero-shot pruning
 - Structured pruning, unstructured pruning, N:M pruning
 - Weight pruning, activation pruning
 - Static pruning and dynamic pruning
 - Pruning for inference, pruning for training



Pruning during DNN Training



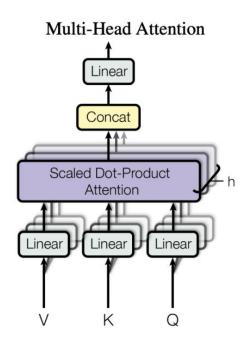
Topics

- Why pruning?
 - Reduce running cost
 - Reduce storage
- General pruning techniques
- Transformer pruning
- Large model pruning



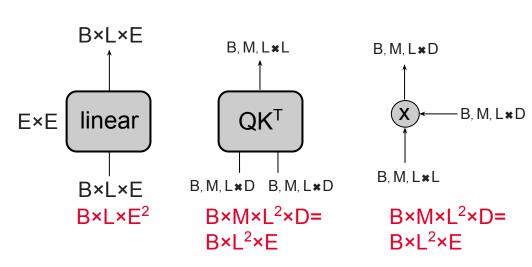
Multi-headed Attention

- Q, K, V tensors are broken into multiple components along the embedding dimension.
 - \circ (B,L,E) * (E*E) \rightarrow (B*L*E)
 - \circ (B,L,E) \rightarrow (B, M, L, E/M) \rightarrow (B, M, L, D), where D=E/M
- All the following operations can be performed independently over each head M.
 - $\bigcirc \qquad \mathsf{QK}^{\mathsf{T}} \rightarrow (\mathsf{B}, \mathsf{M}, \mathsf{L} \times \mathsf{D}) \times (\mathsf{B}, \mathsf{M}, \mathsf{D} \times \mathsf{L}) \rightarrow (\mathsf{B}, \mathsf{M}, \mathsf{L} \times \mathsf{L})$
 - \circ Softmax(QK^T) \rightarrow (B, M, L*L)
 - $\bigcirc \qquad \text{Softmax}(\mathsf{QK}^\top) \star \mathsf{V} \to (\mathsf{B}, \mathsf{M}, \mathsf{L} \star \mathsf{L}) \star (\mathsf{B}, \mathsf{M}, \mathsf{L} \star \mathsf{D}) \to (\mathsf{B}, \mathsf{M}, \mathsf{L} \star \mathsf{D}) \to (\mathsf{B} \star \mathsf{L} \star \mathsf{E})$

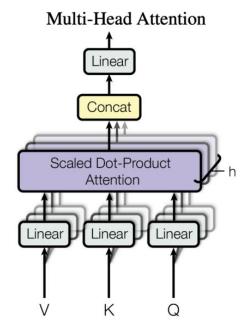




Multi-Head Attention



 The introduction of multiple heads do not change the computational cost of the transformer.





Pruning on Transformers: Token Pruning

 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).

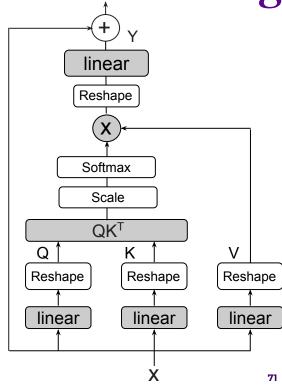
$$\circ \quad (B,L,E) * (E*E) \rightarrow (B*L*E)$$

• The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.

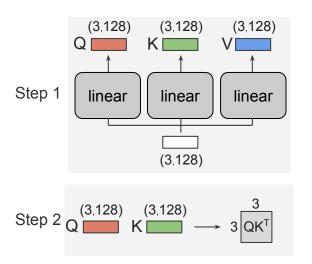
$$\circ$$
 QK^T \rightarrow (B, L*E) * (B,E*L) \rightarrow (B, L*L)

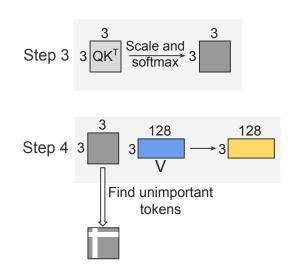
- Scale and normalize the score using softmax.
 - \circ Softmax(QK^T) \rightarrow (B, L*L)
- Multiply each value vector by the softmax score.
 - Softmax(QK^T) * V
 - \circ (B, L*L) * (B, L*E) \rightarrow (B, L*E)
- Pass the result to the linear layer, sum with the input.

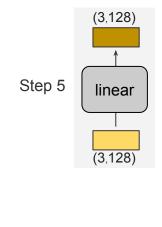




Pruning on Transformers: Token Pruning

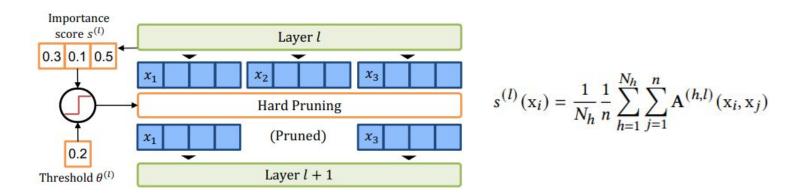






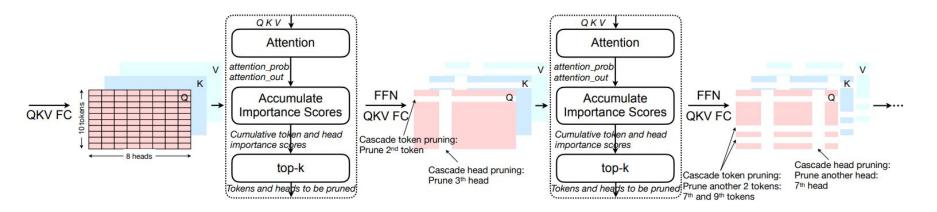


Pruning on Transformers: Token Pruning



• One simple approach involves computing the importance score of each token, and remove the tokens whose importance score is lower than a predefined threshold.

Pruning on Transformers: Token Pruning



 Tokens and heads can be pruned jointly, the removed tokens and heads will result in a much reduced computation for all the following layers.



Pruning on Transformers: Token Pruning



- Not all the tokens are necessary to generate the final results.
- Unimportant tokens can be removed progressively as an input sequence passes through transformer layers.

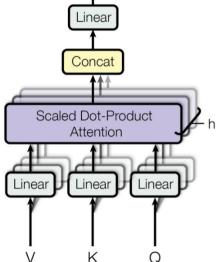


New Pruning Dimension: Head

Pribladdition to the valuewise and channelwise pruning, transformer allows for a new type of pruning: multi-head pruning.

$$(1, 197, 768) \rightarrow (1, 12, 197, 64) \rightarrow (1, 4, 197, 64)$$

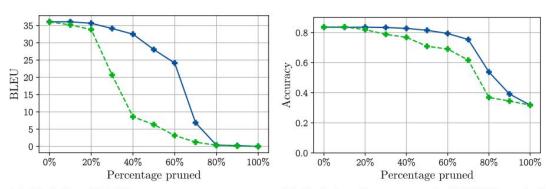
Input Input with 12 heads Input with 4 heads





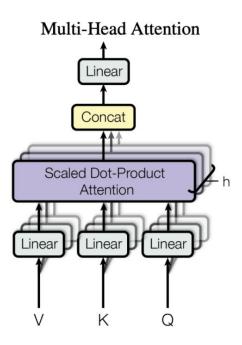
Multi-headed Attention

• We observe that the majority of attention heads can be removed without deviating too much from the original score. Surprisingly, in some cases removing an attention head results in an increase in BLEU/accuracy.



(a) Evolution of BLEU score on newstest2013 when heads are pruned from WMT.

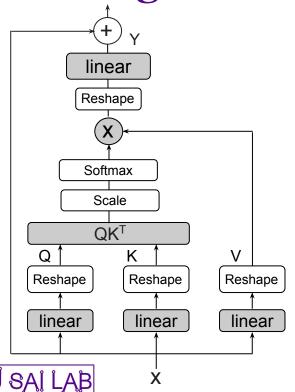
(b) Evolution of accuracy on the MultiNLI-matched validation set when heads are pruned from BERT.





Michel, Paul, Omer Levy, and Graham Neubig. "Are sixteen heads really better than one?" *Advances in neural information processing systems* 32 (2019).

Pruning on Transformers: Head Pruning



- Xi is an embedded vector of ith token.
- There are in total n tokens.
- The output vector of qth token can be expressed as:

$$\begin{aligned} \text{Att}_{W_k,W_q,W_v,W_o}(\mathbf{x},q) &= W_o \sum_{i=1}^n \alpha_i W_v x_i \\ \text{where } \alpha_i &= \operatorname{softmax} \left(\frac{q^\intercal W_q^\intercal W_k x_i}{\sqrt{d}} \right) \end{aligned}$$

 If we further expressed with multi-head attention, the output vector can be expressed as:

$$\mathrm{MHAtt}(\mathbf{x},q) = \sum_{h=1}^{N_h} \underline{\xi_h} \mathrm{Att}_{W_k^h,W_q^h,W_v^h,W_o^h}(\mathbf{x},q)$$

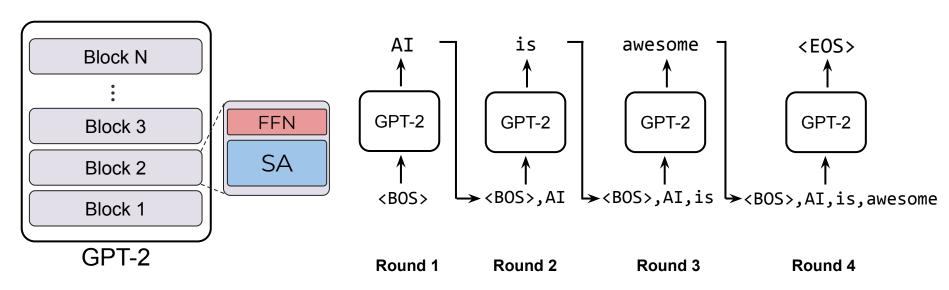
Where the ε_h are mask variables with values in {0, 1}.

$$I_h = \mathbb{E}_{x \sim X} \left| \frac{\partial \mathcal{L}(x)}{\partial \mathcal{E}_h} \right|$$

Topics

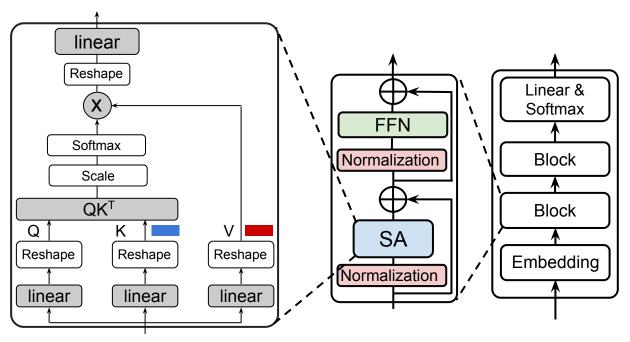
- Pruning in CNN and transformers
 - Reduce running cost
 - Reduce storage
- Sparsity encoding
- General pruning techniques
- Transformer pruning
- Large model pruning





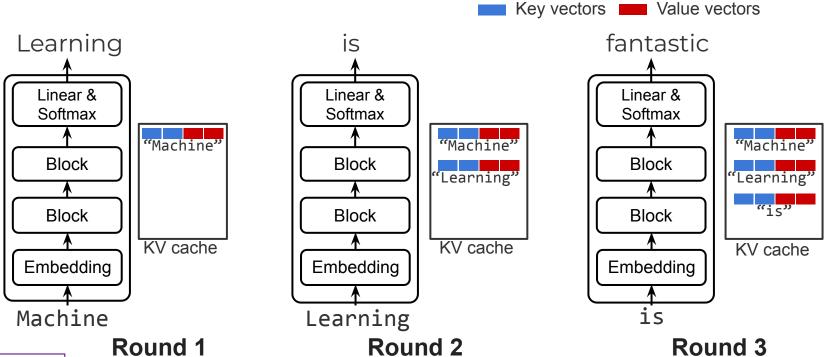
Each token is generated in an autoregressive manner.



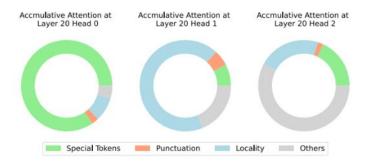




• We need to buffer the v and k for later usage.





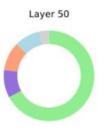


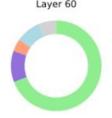






- We show the attention score of each token.
- Different attention heads usually have different importance scores on KV vectors.
- The importance of KV vectors also varies across layers.

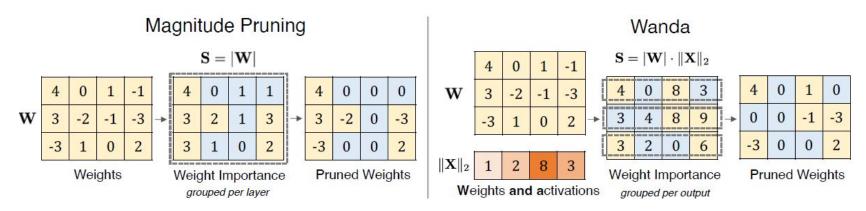








LLM Pruning: Wanda



- A zero-shot pruning method.
- Prune the weights by considering the input statistics.
- For each weight, if the corresponding input's magnitude is large, the output will also be large.
- Need some training samples for calibration.

