# Lecture 03:
# Transformer and Large Model

# Notes

- <span style="color:red">Quiz 1 score can be found in Brightspace</span>
- Course website: https://www.saiqianzhang.com/COURSE/
- I use Brightspace to post announcements and grades
- I provide an online zoom meeting option for people interested in auditing the class. However, enrolled students are required to attend in person unless special condition.
- Discussion groups has been created in the Brightspace
- Course email: efficientaiaccelerator@gmail.com

NYU SAI LAB

# Recap

- Convolutional Neural Network
  - Basic building blocks
  - Popular CNN architectures
    - VGG
    - ResNet
    - MobileNet
    - ShuffleNet
    - SqueezeNet
    - DenseNet
    - EfficientNet
    - ConvNext
    - ShiftNet
  - CNN architectures for other vision tasks
    - Image Segmentation, Object Detection

NYU SAI LAB

# Topics

- <span style="color:red">Transformer basics</span>
- Bert
- Vision transformer
- Large Language Model
- Self-supervised learning

# Transformers

- Proposed in "Attention Is All You Need" in 2017
- The vanilla Transformer is a sequence-to-sequence model and consists of transformer blocks.



**Attention Is All You Need**

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
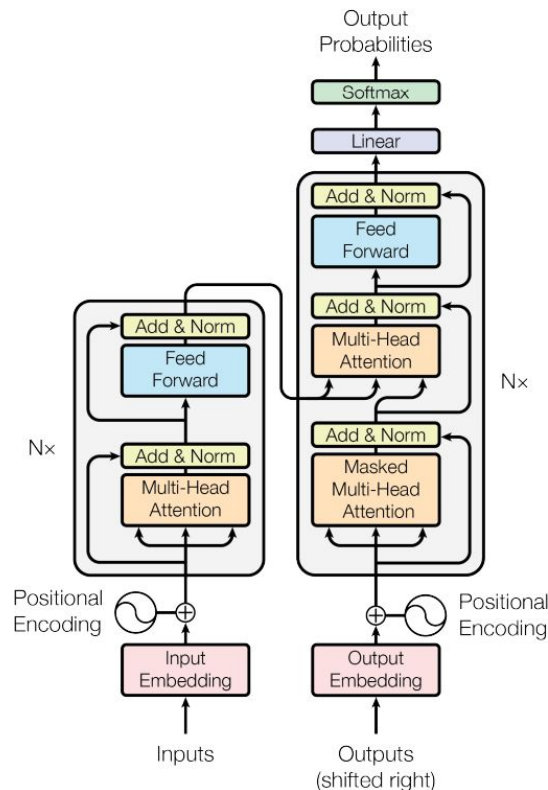noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
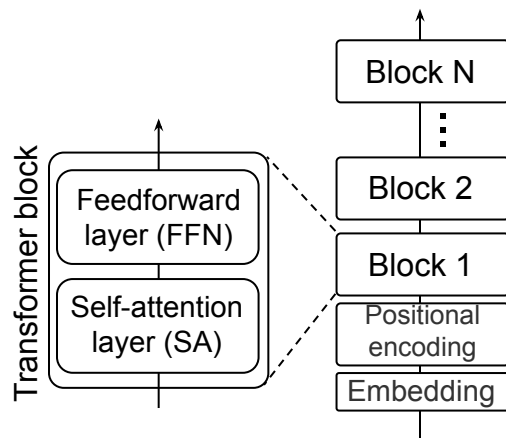usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** †
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
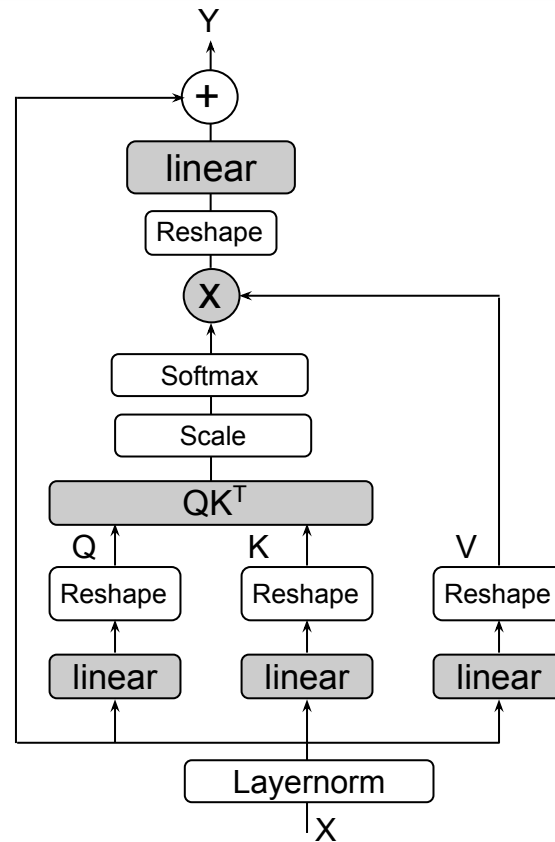Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** ‡
illia.polosukhin@gmail.com

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).
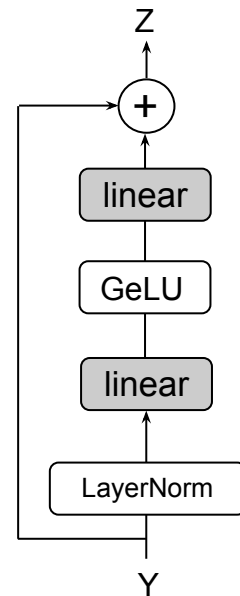
NYU SAI LAB

5

# Transformers



- Each transformer block includes a self-attention layer and a feedforward layer.
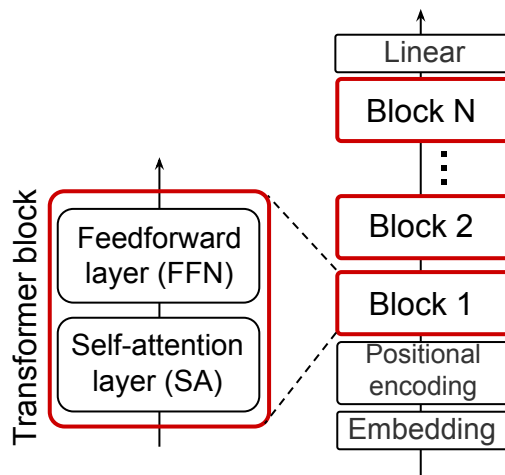
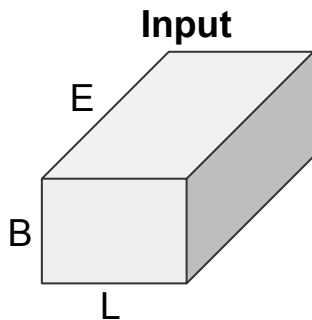**Self attention block (SA)**
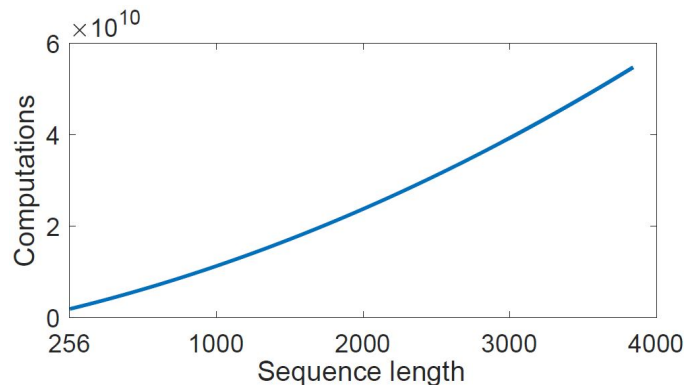
**Feed forward block (FFN)**

# Transformers: Transformer Block
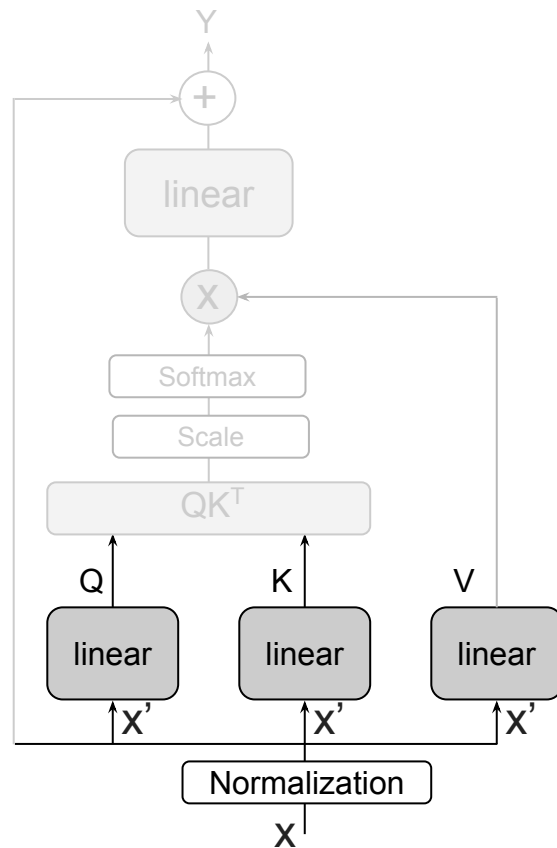
# Transformers



- The input contains three dimensions:
  - B: batch
  - L: token length
  - E: embeddings
- The amount of computation is closely related to the token length L.
- Longer sequences are disproportionately expensive because attention is quadratic to the sequence length.
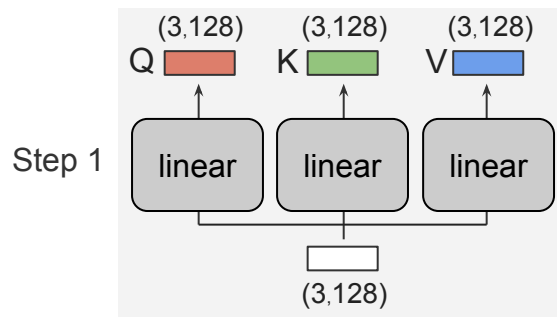
# Self-Attention Block

- The input x is first normalized, then the first step in calculating self-attention is to create three vectors from the input x', denoted as: Query (Q), Key (K), Value (V).
  - $(B,L,E) \times (E \times E) \rightarrow (B \times L \times E)$ $(BLE^2)$
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
  - $QK^T \rightarrow (B, L \times E) \times (B, E \times L) \rightarrow (B, L \times L)$
- Scale and normalize the score using softmax.
  - $Softmax(QK^T) \rightarrow (B, L \times L)$
- Multiply each value vector by the softmax score.
  - $Softmax(QK^T) \times V$
  - $(B, L \times L) \times (B, L \times E) \rightarrow (B, L \times E)$
- Pass the result to the linear layer, sum with the input.

Y

+

linear

X

Softmax

Scale

$QK^T$

Q          K          V

linear     linear     linear

X'          X'          X'

Normalization

X

9

# Example

"I love AI" $\longrightarrow$ 3 $\overset{128}{\boxed{\phantom{xx}}}$

Step 1

Q (3,128) [🟥]   K (3,128) [🟩]   V (3,128) [🟦]
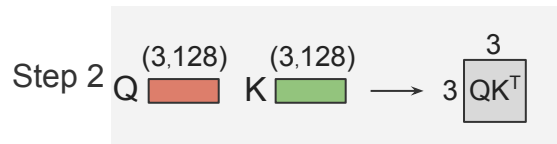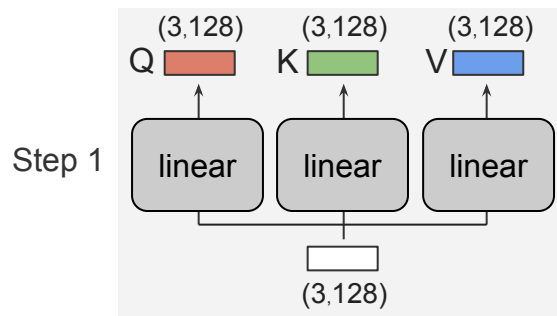
[linear]   [linear]   [linear]

[   ]
(3,128)

# 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) $\ast$ (E$\ast$E) $\rightarrow$ (B$\ast$L$\ast$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$\ast$E) $\ast$ (B, E$\ast$L) $\rightarrow$ (B, L$\ast$L) ($BL^2E$)
- Scale and normalize the score using softmax.
  - Softmax($QK^T$) $\rightarrow$ (B, L$\ast$L)
- Multiply each value vector by the softmax score.
  - Softmax($QK^T$) $\ast$ V
  - (B, L$\ast$L) $\ast$ (B, L$\ast$E) $\rightarrow$ (B, L$\ast$E)
- Pass the result to the linear layer, sum with the input.

# Example

"I love AI" $\longrightarrow$ 3 [ 128 ]

**Step 1**

(3,128)   (3,128)   (3,128)

Q [ ]   K [ ]   V [ ]

linear   linear   linear

[ ]

(3,128)

**Step 2**

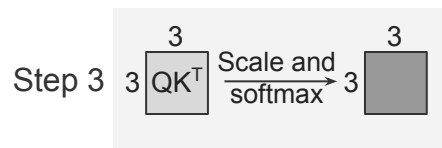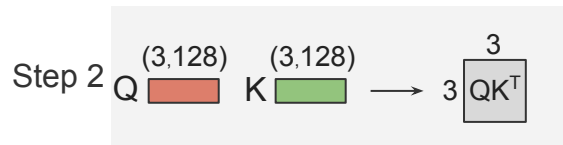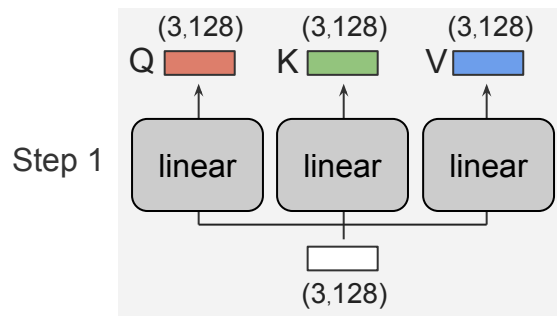Q (3,128) [ ]   K (3,128) [ ] $\longrightarrow$ 3 [ QK$^T$ ] 3

NYU SAI LAB

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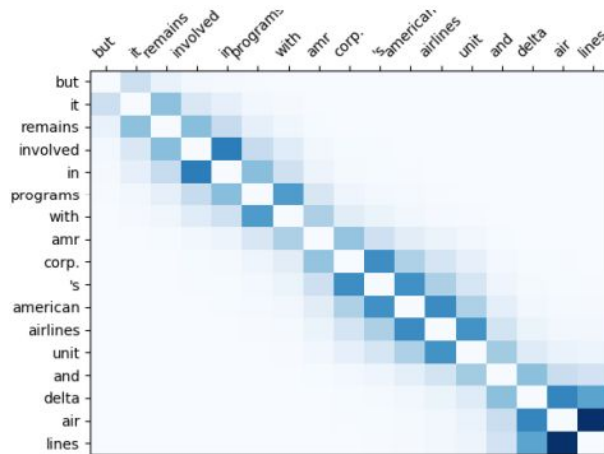- Pass the result to the linear layer, sum with the input.

# Example

"I love AI" $\longrightarrow$ 3 $\boxed{\phantom{xxx}}$ 128

**Step 1**

Q (3,128) $\blacksquare$    K (3,128) $\blacksquare$    V (3,128) $\blacksquare$

linear    linear    linear

(3,128)

**Step 2** Q (3,128) $\blacksquare$    K (3,128) $\blacksquare$ $\longrightarrow$ 3 $\boxed{QK^T}$ 3

**Step 3** 3 $\boxed{QK^T}$ 3 $\xrightarrow{\text{Scale and softmax}}$ 3 $\boxed{\phantom{xx}}$ 3
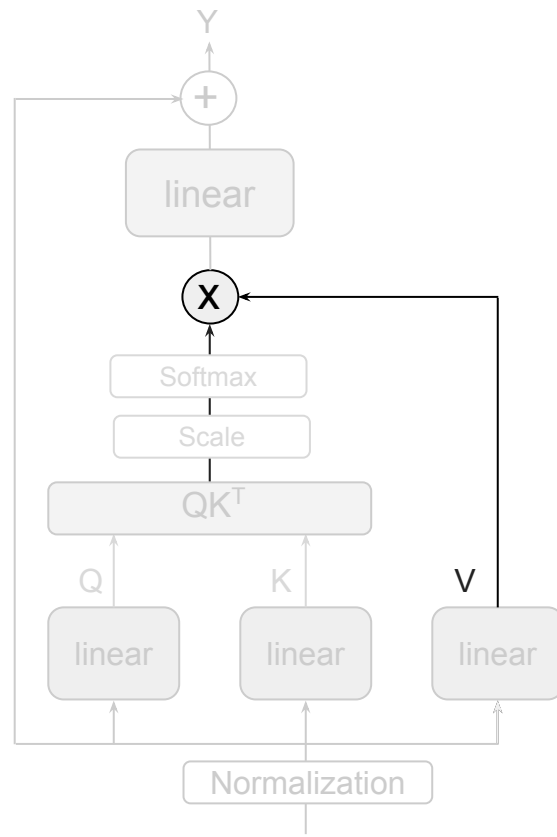
NYU SAI LAB

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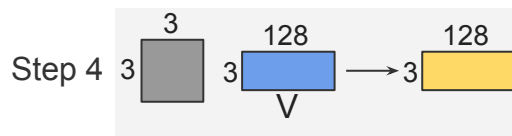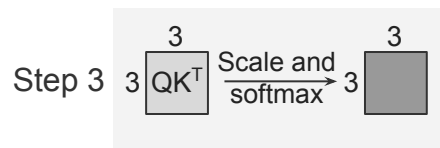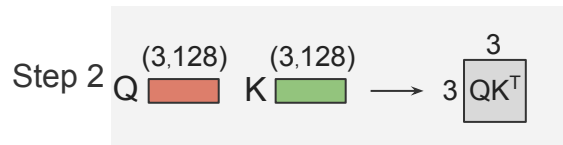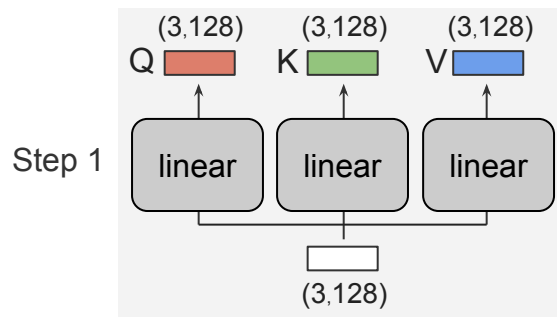- Pass the result to the linear layer, sum with the input.

NYU SAI LAB

# 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).
  - $(B,L,E) \ast (E \ast E) \rightarrow (B \ast L \ast E)$
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
  - $QK^{T} \rightarrow (B, L \ast E) \ast (B, E \ast L) \rightarrow (B, L \ast L)$
- Scale and normalize the score using softmax.
  - $Softmax(QK^{T}) \rightarrow (B, L \ast L)$
- Multiply each value vector by the softmax score.
  - $Softmax(QK^{T}) \ast V$
  - $(B, L \ast L) \ast (B, L \ast E) \rightarrow (B, L \ast E)$ $(BL^{2}E)$
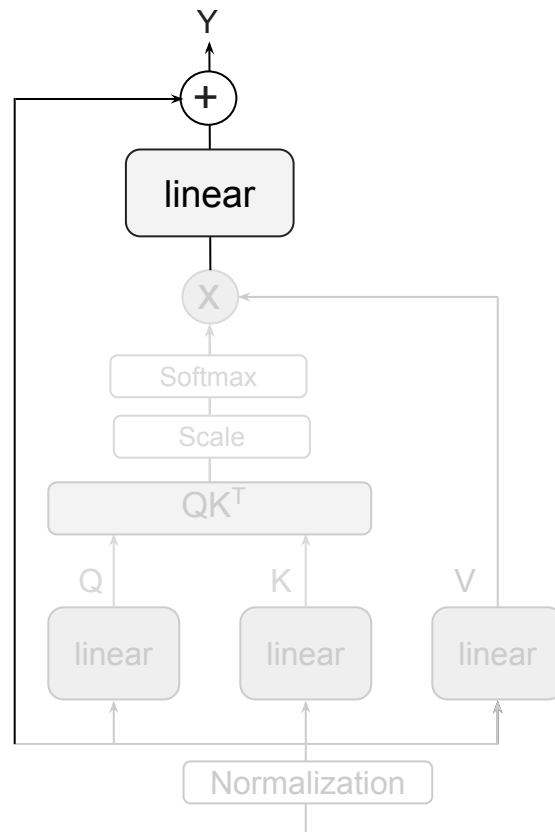- Pass the result to the linear layer, sum with the input.
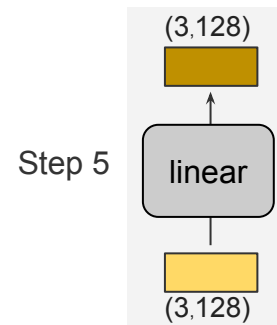
# Example

"I love AI" $\longrightarrow$ 3 $\boxed{\phantom{128}}$ 128

**Step 1**

(3,128) Q    (3,128) K    (3,128) V

linear   linear   linear

(3,128)

**Step 2**   (3,128) Q   (3,128) K $\longrightarrow$ 3 $QK^T$ 3

**Step 3**   3 $QK^T$ 3 $\xrightarrow[\text{softmax}]{\text{Scale and}}$ 3   3

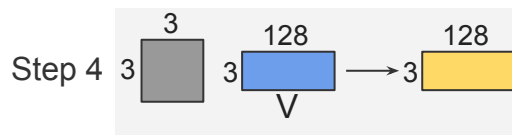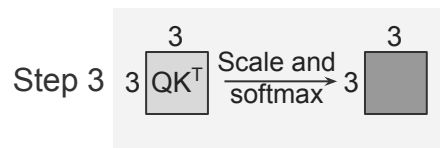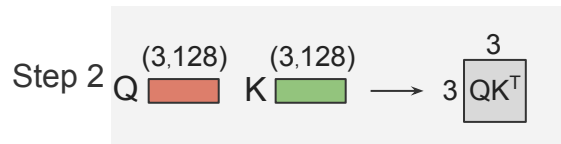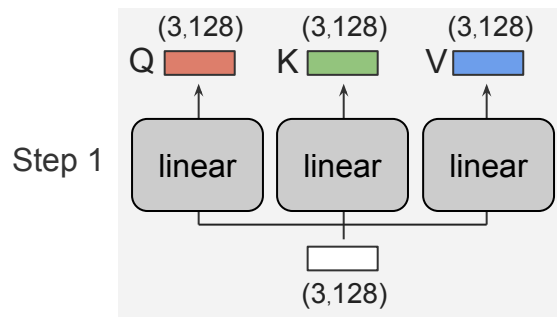**Step 4**   3   3   3    128 V $\longrightarrow$ 3   128

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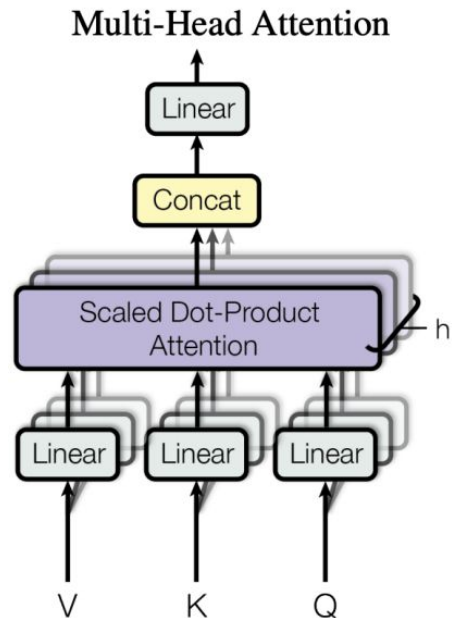- Pass the result to the linear layer, sum with the input.

# Example

"I love AI" $\longrightarrow$ 3 $\boxed{\phantom{128}}$ 128

**Step 1**

(3,128) Q
(3,128) K
(3,128) V

linear linear linear

(3,128)

**Step 2** Q (3,128) K (3,128) $\longrightarrow$ 3 $\boxed{QK^T}$ 3

**Step 3** 3 $\boxed{QK^T}$ 3 $\xrightarrow{\text{Scale and softmax}}$ 3 $\boxed{\phantom{x}}$ 3

**Step 4** 3 $\boxed{\phantom{x}}$ 3 · 3 $\boxed{\phantom{128}}$ 128 V $\longrightarrow$ 3 $\boxed{\phantom{128}}$ 128

**Step 5**

(3,128)

linear

(3,128)

# Multi-headed Attention

- Q, K, V tensors are broken into multiple components along the embedding dimension.
    - $(B,L,E) ✶ (E✶E) \rightarrow (B✶L✶E)$
    - $(B,L,E) \rightarrow (B, M, L, E/M) \rightarrow (B, M, L, D)$, where D=E/M
- All the following operations can be performed independently over each head M.
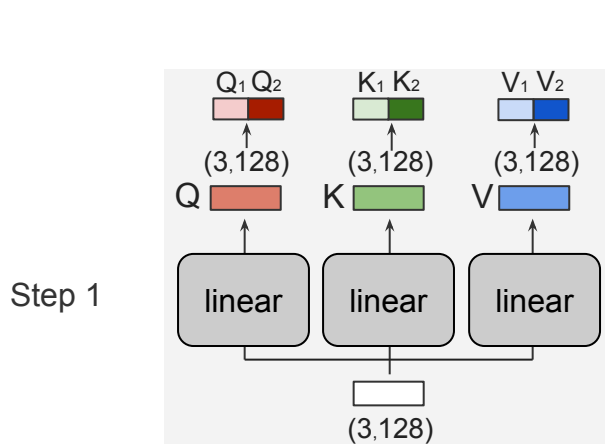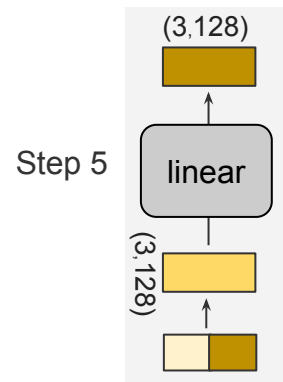    - $QK^T \rightarrow (B, M, L✶D) ✶ (B, M, D✶L) \rightarrow (B, M, L✶L)$
    - $Softmax(QK^T) \rightarrow (B, M, L✶L)$
    - $Softmax(QK^T) ✶ V \rightarrow (B, M, L✶L) ✶ (B, M, L✶D) \rightarrow (B, M, L✶D) \rightarrow (B✶L✶E)$



**Multi-Head Attention**

20

# Example

# Multi-headed Attention

- Why we need multiple heads?
  - Multiple attention heads in transformers are used to enhance the expressive power and modeling capabilities of the network.
  - By using multiple attention heads, transformers can capture different types of dependencies and relationships between words or elements in a sequence.
  - Having multiple heads allows the model to perform attention calculations in parallel, which can improve computational efficiency.



Multi-Head Attention

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

NYU SAI LAB

# Feed Forward Layer

Z

+

linear

GeLU

linear

LayerNorm

Y

- The two linear layers are big:
  - (Ex4E) and (4ExE), E can be large (e.g., 4096)
  - This is expensive to implement.
- GeLU:
  - $GeLU(x) = x\Phi(x)$
    $\Phi(x) = P(y \leq x), \text{ where } Y \sim N(0,1)$

### GELU activation function

Hendrycks, Dan, and Kevin Gimpel. "Gaussian error linear units (gelus)." *arXiv preprint arXiv:1606.08415* (2016).

NYU SAI LAB

# Example



Z

+

linear

GeLU

linear

LayerNorm

Y

"I love AI" ⟶ 3 □ 128

Step 6    Linear 1    (3,512)    (3,128)

Step 7    GeLU    (3,512)    (3,512)

Step 8    Linear 2    (3,128)    (3,512)

NYU SAI LAB

# Layer Normalization

**Step 1**

E



$\mu_0, \delta_0$
$\mu_1, \delta_1$

L     X

$\mu_{L-1}, \delta_{L-1}$

$$X' = \frac{X_r - \mu_r}{\sigma_r} \text{ For each } r \in L$$

**Step 2**

E



L     X'
...

$$Y_e = \alpha_e X'_e + \beta_e \text{ For each } e \in E$$

- LayerNorm is applied on each input sample.
- Both α and β have a length of E.

Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer normalization." *arXiv preprint arXiv:1607.06450* (2016).

# Layer Normalization

**Step 1**

E

$$X' = \frac{X_r - \mu_r}{\sigma_r}$$ For each r∈L

μ₀,δ₀ → $\mu_0, \delta_0$
μ₁,δ₁ → $\mu_1, \delta_1$

L    X

$\mu_{L-1}, \delta_{L-1}$

**Step 2**

E

L    X'
     …

$$Y_e = \alpha_e X'_e + \beta_e$$ For each e∈E

- Layer Norm does not store the running mean and running variance, so during the inference time, the mean and variance need to be computed.

Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer normalization." *arXiv preprint arXiv:1607.06450* (2016).

NYU SAI LAB

# RMS Normalization

**Step 1**

E

$$X'_r = \frac{X_r}{\sqrt{\frac{\sum_i X_{r,i}^2}{L}}}$$

For each r∈L

**Step 2**

E

$$Y_e = \alpha_e X'_e$$

For each e∈E

- The experiments demonstrate RMSNorm achieves similar and even better results than LayerNorm.

Zhang, Biao, and Rico Sennrich. "Root mean square layer normalization." *Advances in Neural Information Processing Systems* 32 (2019).

NYU SAI LAB

# Transpose & Reshape



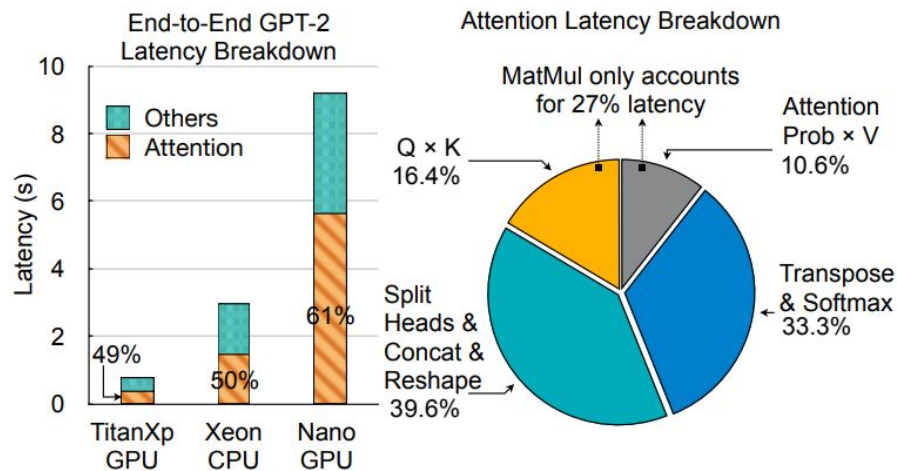Fig. 2. End-to-End GPT-2 latency breakdown on various platforms, and attention latency breakdown on TITAN Xp GPU. Attention accounts for over 50% of total latency. Data movements account for 73% of attention latency.

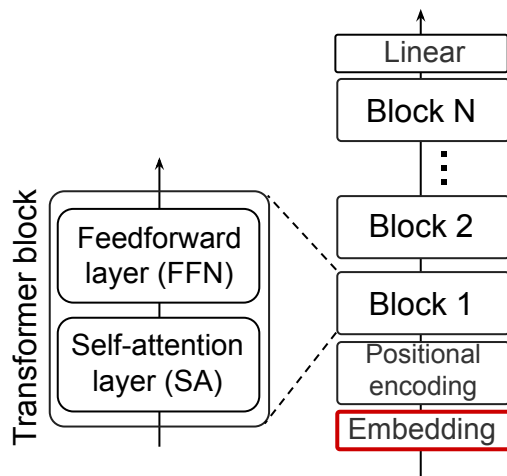**Reshaping operation**

$(B, L, E) \longrightarrow (B, M, L, E/M)$

**Transpose operation**

$(B, L, E) \longrightarrow (B, E, L)$

Wang, Hanrui, Zhekai Zhang, and Song Han. "Spatten: Efficient sparse attention architecture with cascade token and head pruning." *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2021.
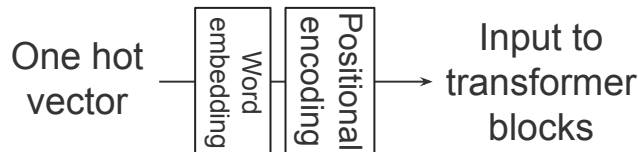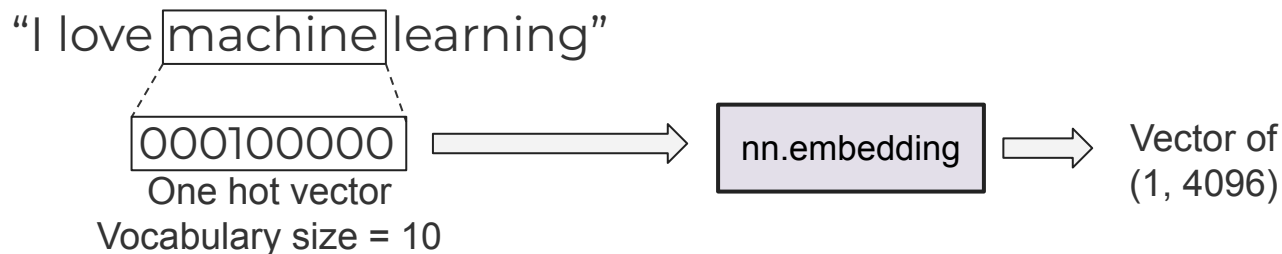
# Transformers: Word Embedding

# Transformers: Word Embedding

- For each word, we can convert them into a one-hot vector.
- We use the embedding layer to encode these one-hot vectors, which acts like a trainable lookup table.
- Dictionary: {are, is, machine, good, awesome, learning, love, ….}

"I love machine learning"

000100000
One hot vector
Vocabulary size = 10

nn.embedding

Vector of
(1, 4096)

One hot vector → Word embedding → Positional encoding → Input to transformer blocks

Mikolov, Tomas. "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781* (2013).

NYU SAI LAB

# Transformers: Positional Encoding

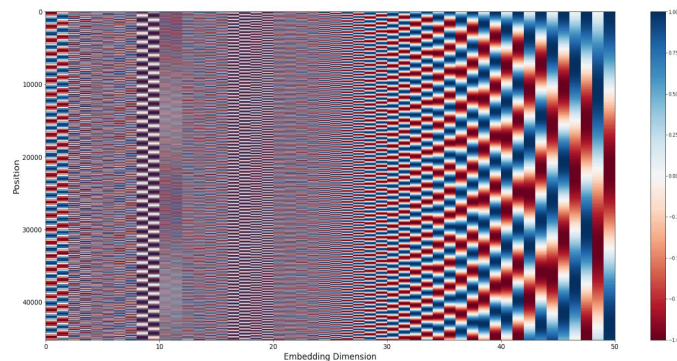# Transformers: Positional Encoding

- One thing that's missing from the model as we have described it so far is a way to account for the order of the words in the input sequence.
- Transformer adds a vector to each input embedding. These vectors follow a specific pattern that the model learns, which helps it determine the position of each word.
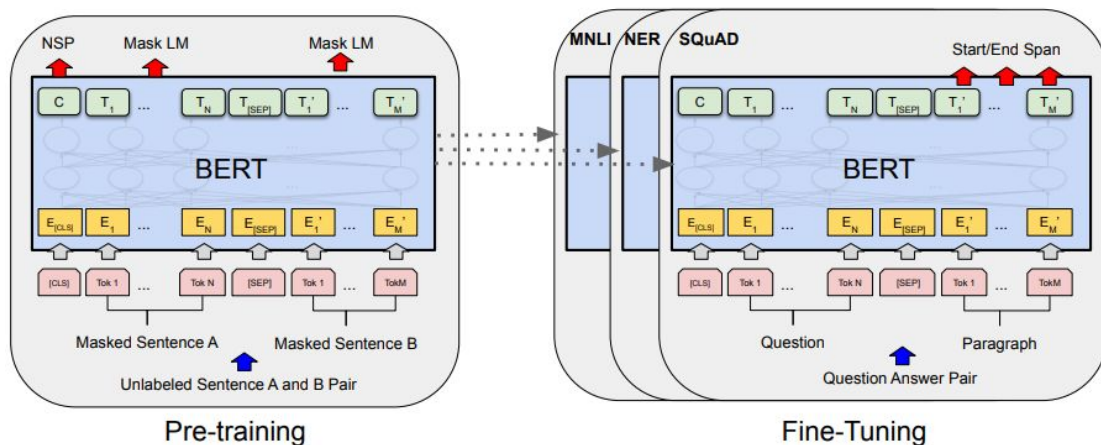
$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{\text{model}}})$$

- pos is the positional of the token, i is the index of the embedding.

# Case Study: BERT

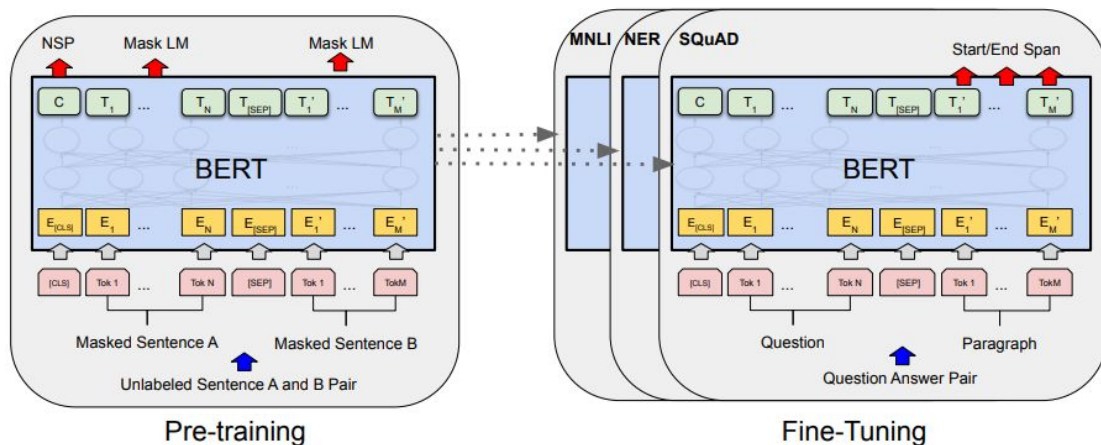- Bidirectional Encoder Representations from Transformers.
- BERT is designed to understand the text by considering both the words before and after it.
- BERT consists of transformer encoders and takes the entire sentence as the input.
- BERT can not generate new text.



Devlin, Jacob. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).
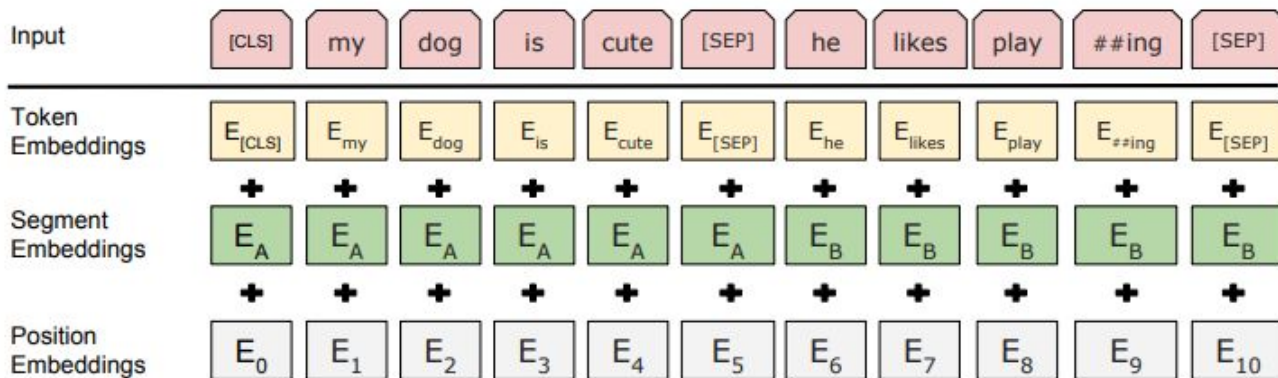
# Case Study: BERT

- Use the encoder's output embeddings as input features for downstream tasks.
- Represent a sentence as a vector for semantic similarity, clustering, or search.
- The pre-trained BERT model can be finetuned with just one additional output layer to create state-of-the-art models for a wide range of tasks.



Devlin, Jacob. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).

# BERT

- A [CLS] token is inserted at the start of every sequence, and the two sentences in the sequence are separated by a [SEP] token.
- The final hidden state corresponding to this token is used as the aggregate sequence representation for classification tasks.
- In addition to the positional information, BERT contains a segment embeddings to differentiate the sentences.

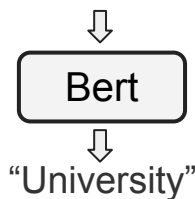| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

NYU SAI LAB

# BERT Pretraining

- BERT is pretrained using two unsupervised tasks:
  - Masked Language Modeling (MLM)
    - We simply mask some percentage of the input tokens at random, and then predict those masked tokens.
  - Next Sentence Prediction (NSP)
    - Given two sentences, A and B, predict whether B is A's following sentence.
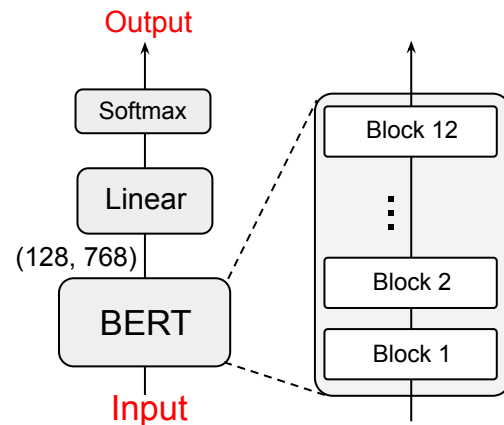
"New York University is a great school."          "New York University is a great school."

⇩

Bert

⇩

"University"

# Downstream Tasks

- For text classification task, Bert will return a binary output.
  - Single sentence task (SST-2): The task is to predict the sentiment of a given sentence.
- The input may contain a single sentence, or a pair of sentences.
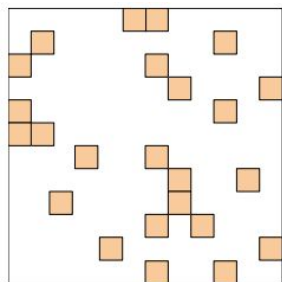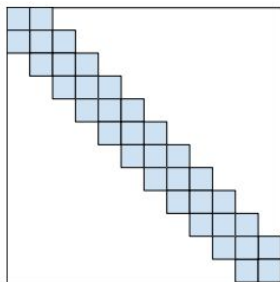  - Similarity and Paraphrase tasks (MRPC): Is the second sentence a paraphrase of the first sentence.

Output

Softmax

Linear

(128, 768)

BERT

Input

Block 12

Block 2

Block 1

# BERT Performance

| System | MNLI-(m/mm) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTE | Average |
|---|---|---|---|---|---|---|---|---|---|
| | 392k | 363k | 108k | 67k | 8.5k | 5.7k | 3.5k | 2.5k | - |
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.8 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 87.4 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.1 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| **BERT$_{LARGE}$** | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

- BERT can achieve better performance over GLUE datasets than GPT-1.
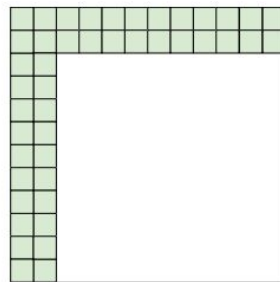
NYU SAI LAB

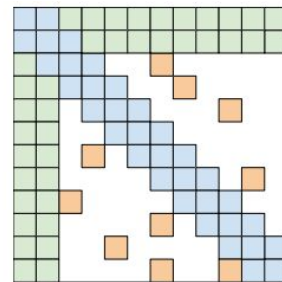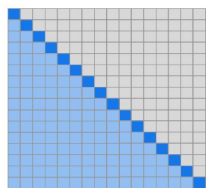# Efficient Self-Attention Design



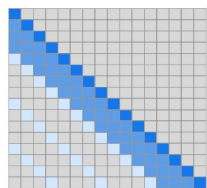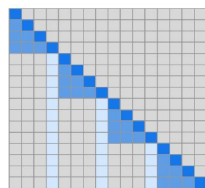(a) Random attention  (b) Window attention  (c) Global Attention  (d) BIGBIRD
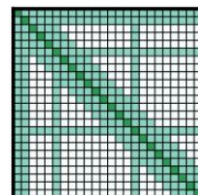
Sparse Transformer (2019)          Longformer (2020)
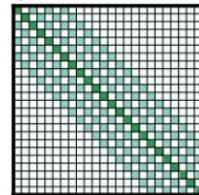
Original          Strided          Fixed

Beltagy, Iz, Matthew E. Peters, and Arman Cohan. "Longformer: The long-document transformer." *arXiv preprint arXiv:2004.05150* (2020).

Child, Rewon, et al. "Generating long sequences with sparse transformers." *arXiv preprint arXiv:1904.10509* (2019).

# Efficient Self-Attention Design

"I love AI" $\longrightarrow$ 3 $\overset{128}{\boxed{\phantom{xxx}}}$



Step 1

Q (3,128) K (3,128) V (3,128)

linear   linear   linear

(3,128)

Step 2  Q (3,128)  K (3,128)  $\longrightarrow$  3 $\overset{3}{\boxed{QK^T}}$

Step 3  3 $\overset{3}{\boxed{QK^T}}$  $\xrightarrow[\text{softmax}]{\text{Scale and}}$  3 $\overset{3}{\boxed{\phantom{x}}}$  $\xrightarrow{\text{Pruning}}$  3 $\overset{3}{\boxed{\phantom{x}}}$

Step 4  3 $\overset{3}{\boxed{\phantom{x}}}$  3 $\overset{128}{\boxed{\phantom{xx}}}$  $\longrightarrow$  3 $\overset{128}{\boxed{\phantom{xx}}}$
V

# Token Merging



Pooling layer will reduce the number of tokens.

- We can reduce the number of tokens by merging them together.

Bolya, Daniel, et al. "Token merging: Your vit but faster." *arXiv preprint arXiv:2210.09461* (2022).

# Token Merging



- x,y are two token vectors with length of E

Max pooling
$$z_i = \begin{cases} x_i & \text{if } (|x_i| >= |y_i|) \\ y_i & \text{otherwise} \end{cases}$$

Average pooling
$$z_i = (x_i + y_i)/2$$

Interleaved merging
$$z_i = \begin{cases} x_i & i \text{ is odd} \\ y_i & \text{otherwise} \end{cases}$$

Bolya, Daniel, et al. "Token merging: Your vit but faster." *arXiv preprint arXiv:2210.09461* (2022).
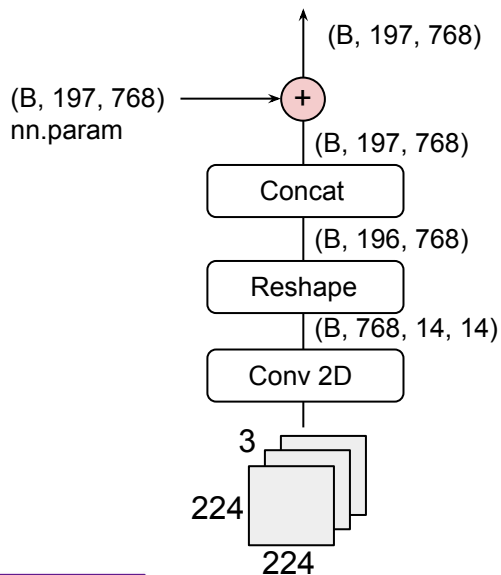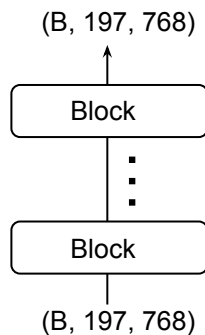
# Topics

- Transformer basics
- Bert
- <span style="color:red">Vision transformer</span>
- Large Language Model
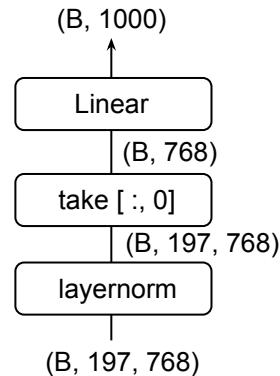- Self-supervised learning

# Vision Transformer

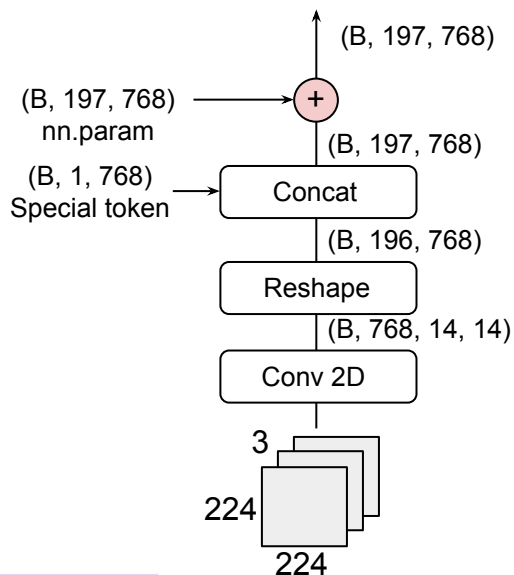- Transformer architecture can also be applied over the computer vision tasks.

(B, 197, 768)

(B, 197, 768)
nn.param

(B, 197, 768)

Concat

(B, 196, 768)

Reshape

(B, 768, 14, 14)

Conv 2D

3

224

224

(part 1)

(B, 197, 768)

Block

.
.
.

Block

(B, 197, 768)

(part 2)

(B, 1000)

Linear

(B, 768)

take [ :, 0]

(B, 197, 768)
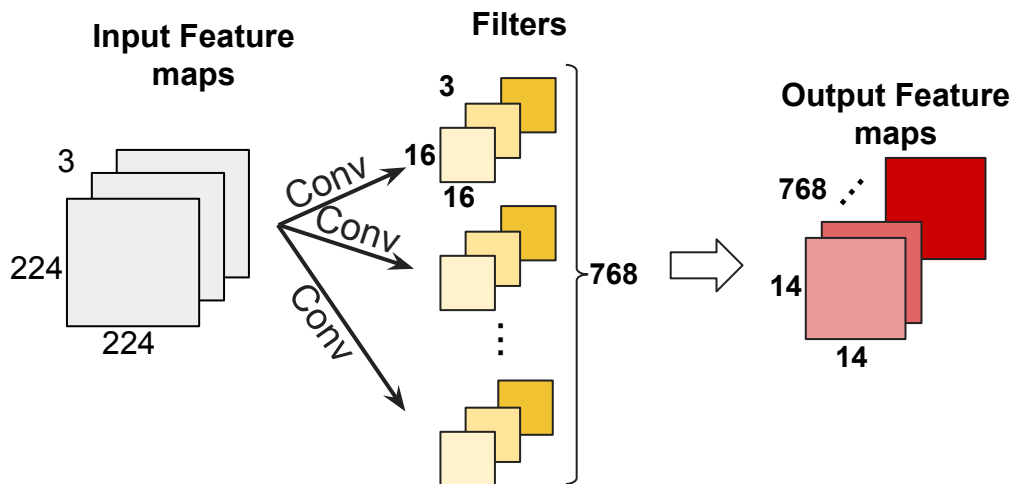
layernorm

(B, 197, 768)

(part 3)

# Vision Transformer

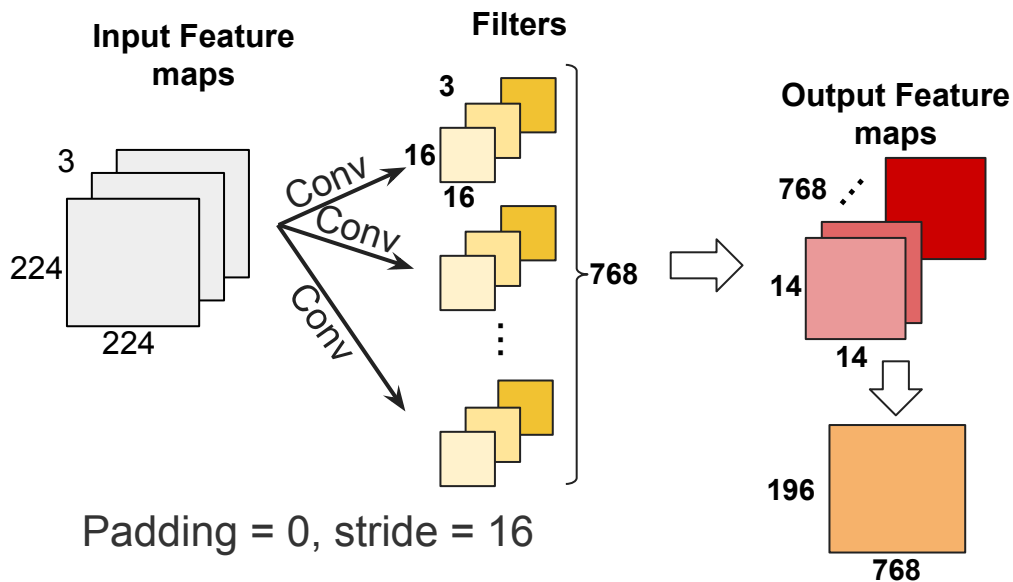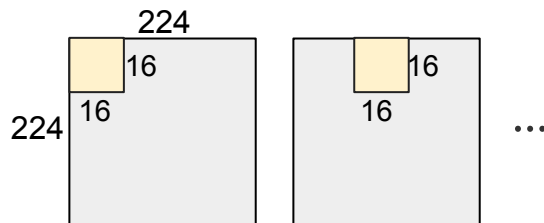- A special placeholder is introduced to aggregate global information about the whole image.
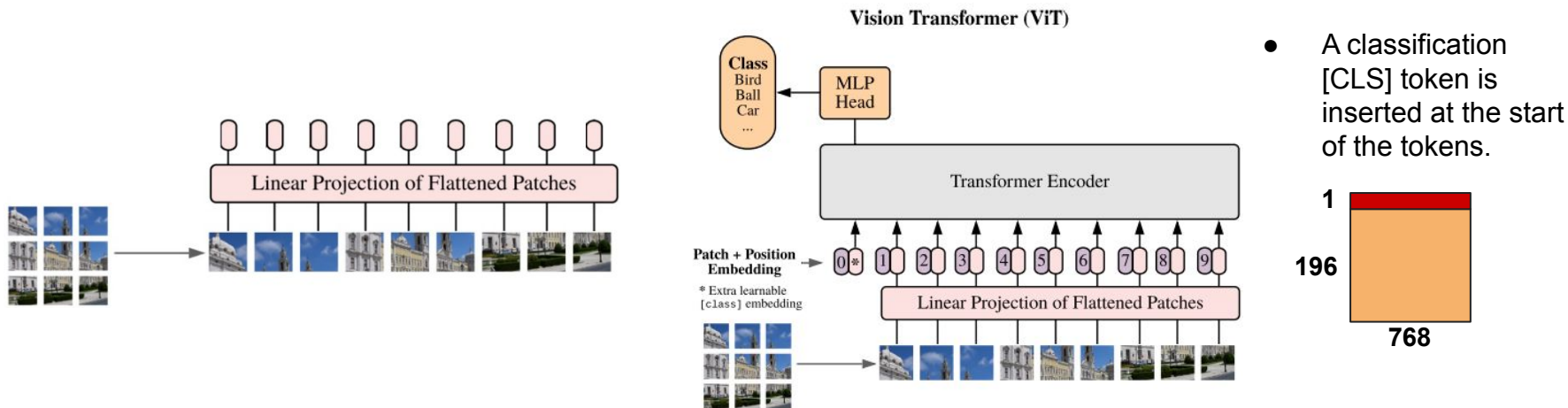


(part 1)

Padding = 0, stride = 16

# Vision Transformer

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



Padding = 0, stride = 16

# Vision Transformer

- An image C×H×W is divided into patches of C×P×P. P is 16✖16 in the previous example.
- They are then flattened and linearly projected to E (e.g., 768) dimensions for a sequence of (H/P) × (W/P) tokens.



- A classification [CLS] token is inserted at the start of the tokens.

Dosovitskiy, Alexey. "An image is worth 16x16 words: Transformers for image recognition at scale." *arXiv preprint arXiv:2010.11929* (2020).

# Vision Transformer

ViT-H/14

ViT-L/16

ViT-B/16

ResNet-50

ResNet-50
23M

ViT-Base
86M

ViT-Large
307M

ViT-Huge
632M

NYU SAI LAB

# Vision Transformer



- Accuracy increases as the training dataset grow.
- ResNet 50 can achieves an accuracy between 75-80% on ImageNet.
- ViT does not strike an efficient balance between parameter count and accuracy when dataset is small.

# Topics

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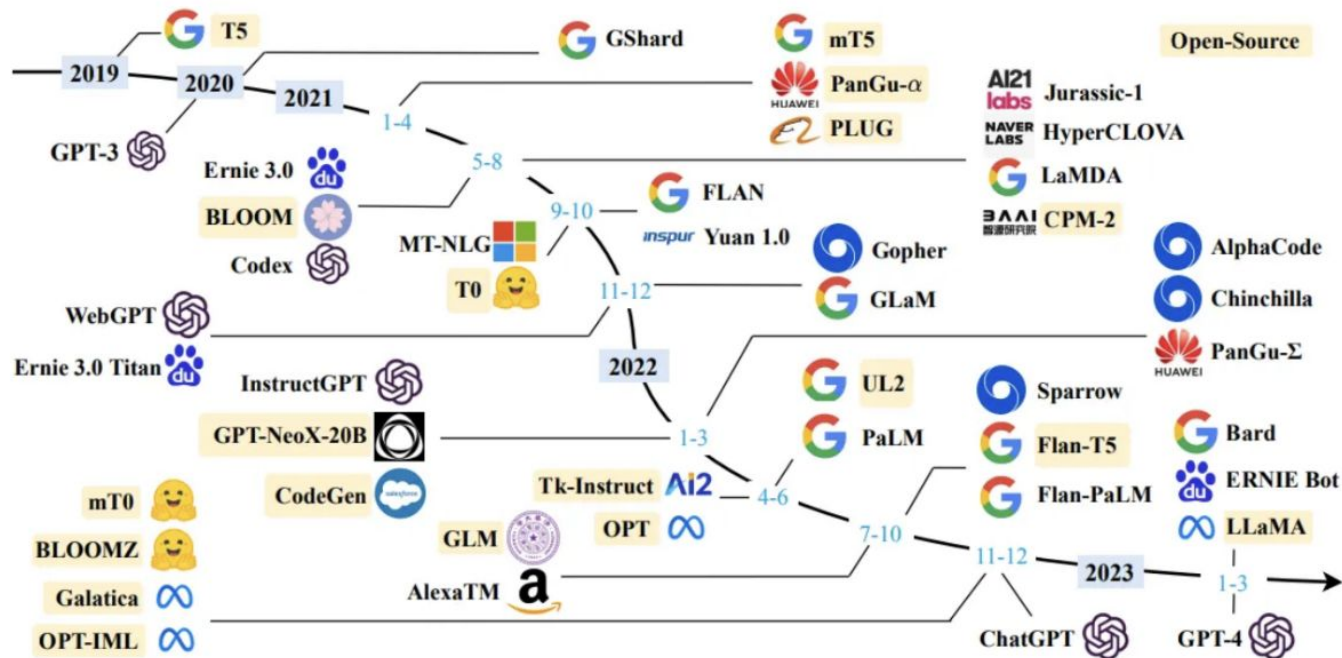- Bert
- Vision transformer
- Large Language Model
- Self-supervised learning

NYU SAI LAB

# Large Language Models (LLMs)

# **Transformers as a Generative AI Tool**



Transformer ⟳ Processing

"Where is new york university?"

**Step 1: Prefilling**

"New" — "York"

Transformer        Transformer ⟳ Generating

"Where is new york university?"    "Where is new york university?"

**Step 2: Decoding**

# Transformers as a Generative AI Tool



- Each token is generated in an autoregressive manner.

Radford, Alec, et al. "Language models are unsupervised multitask learners." *OpenAI blog* 1.8 (2019): 9.

# Transformers as a Generative AI Tool



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

# GPT-2: Prefilling
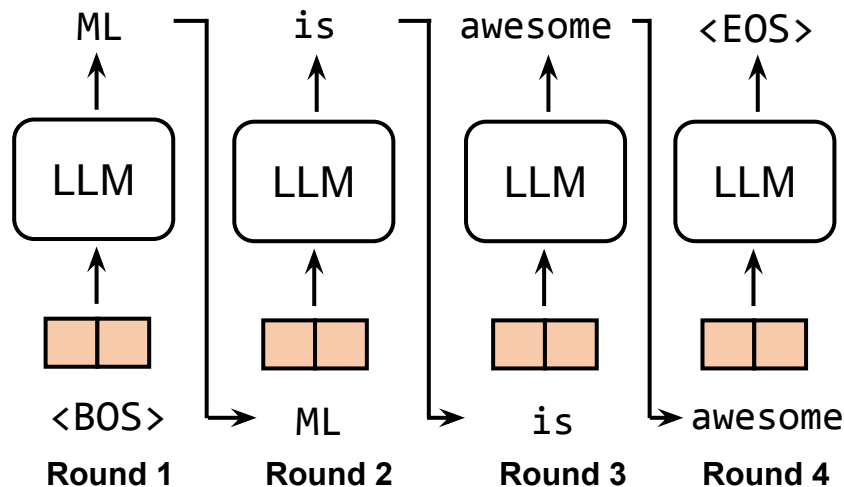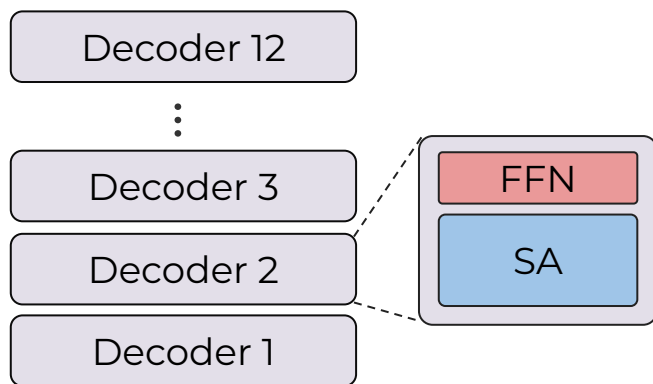


**KV cache**

"How are you"

$k_{i,j}$ 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.

# GPT-2: Decoding



**KV cache**

"How are you"

**KV cache**

I

"How are you"
**Round 1**

- During the decoding stage, LLM generates the responses in an autoregressive way.

# GPT-2: Decoding



am

KV cache

Linear & Softmax

Decoder

Decoder

Embedding

"How are you I"

**Round 1**

doing

KV cache

Linear & Softmax

Decoder

Decoder

Embedding

"How are you I am"

**Round 2**

well

KV cache

Linear & Softmax

Decoder

Decoder

Embedding

"How are you I am doing"

**Round 3**

NYU SAI LAB

57

# GPT-2: Decoding

well

KV cache

well
doing
am
I

Linear &
Softmax

Decoder

Decoder

Embedding

"How are you I am doing"

good

KV cache

good
doing
am
I

Linear &
Softmax

Decoder

Decoder

Embedding

"How are you I am doing"

- We can simply select the token with the highest score. But better results are achieved if the model considers other words as well. So a better strategy is to sample a word from the entire list using the score as the probability of selecting that word.

NYU SAI LAB

58

# 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.

# Why KV Cache Saves Computation?



- Given the input, the $q_L$, $k_L$, $v_L$ are first computed by passing through the linear layers.
- After that, the $k_l$, $v_l$ vectors (l=1,...,L-1) are also loaded from the memory.

# Why KV Cache Saves Computation?



- K and V are loaded from the memory, the q vector of the current token ($q_L$) is multiplied with the each of the key vector $k_i$, ($i=1\ldots L$) to produce the result $A_{L,i}$

# Why KV Cache Saves Computation?



- K and V are loaded from the memory, the q vector of the current token ($q_L$) is multiplied with the each of the key vector $k_i$, (i=1…L) to produce the result $A_{L,i}$

# Why KV Cache Saves Computation?



- K and V are loaded from the memory, the q vector of the current token ($q_L$) is multiplied with the each of the key vector $k_i$, ($i=1…L$) to produce the result $A_{L,i}$

# Why KV Cache Saves Computation?



- Afterwards, the computed $A_{L,i}$ (i=1…L) will then passed to softmax function.
- Then each element of $A_L$ will then multiplied with vi (i=1…L) and elementwise sum together.

# Why KV Cache Saves Computation?



- Afterwards, the computed $A_{L,i}$ (i=1…L) will then passed to softmax function
- Then each element of $A_L$ will then multiplied with vi (i=1…L) and elementwise sum together.

# Why KV Cache Saves Computation?



- Afterwards, the computed $A_{L,i}$ (i=1…L) will then passed to softmax function
- Then each element of $A_L$ will then multiplied with vi (i=1…L) and elementwise sum together.

# Deepseek V3



$$c_t^Q = W^{DQ}\mathbf{h}_t, \tag{37}$$

$$[\mathbf{q}_{t,1}^C; \mathbf{q}_{t,2}^C; ...; \mathbf{q}_{t,n_h}^C] = \mathbf{q}_t^C = W^{UQ}c_t^Q, \tag{38}$$

$$[\mathbf{q}_{t,1}^R; \mathbf{q}_{t,2}^R; ...; \mathbf{q}_{t,n_h}^R] = \mathbf{q}_t^R = \text{RoPE}(W^{QR}c_t^Q), \tag{39}$$

$$\mathbf{q}_{t,i} = [\mathbf{q}_{t,i}^C; \mathbf{q}_{t,i}^R], \tag{40}$$

$$c_t^{KV} = W^{DKV}\mathbf{h}_t, \tag{41}$$

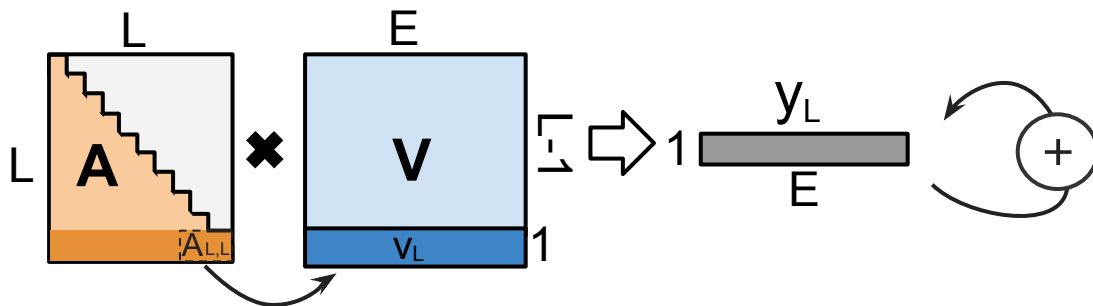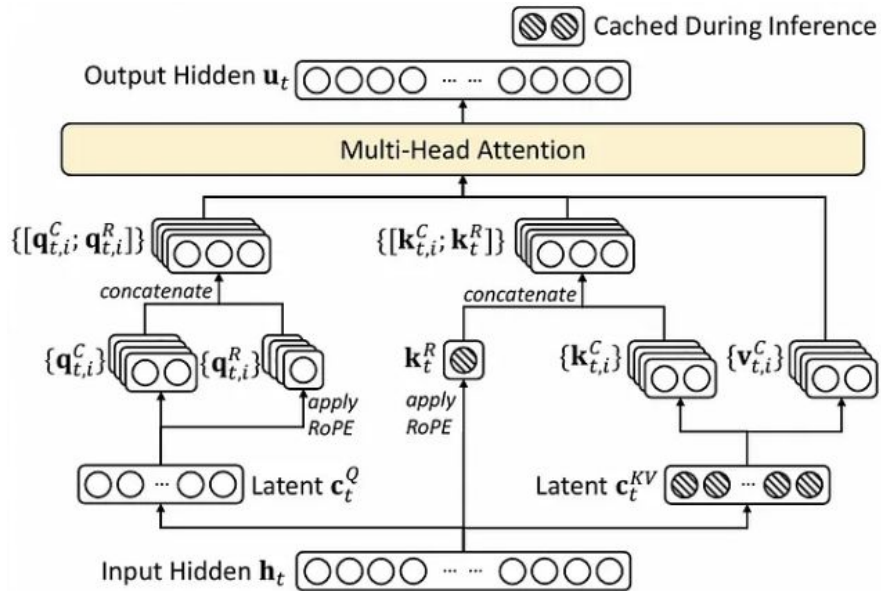$$[\mathbf{k}_{t,1}^C; \mathbf{k}_{t,2}^C; ...; \mathbf{k}_{t,n_h}^C] = \mathbf{k}_t^C = W^{UK}c_t^{KV}, \tag{42}$$

$$\mathbf{k}_t^R = \text{RoPE}(W^{KR}\mathbf{h}_t), \tag{43}$$

$$\mathbf{k}_{t,i} = [\mathbf{k}_{t,i}^C; \mathbf{k}_t^R], \tag{44}$$

$$[\mathbf{v}_{t,1}^C; \mathbf{v}_{t,2}^C; ...; \mathbf{v}_{t,n_h}^C] = \mathbf{v}_t^C = W^{UV}c_t^{KV}, \tag{45}$$

$$\mathbf{o}_{t,i} = \sum_{j=1}^{t} \text{Softmax}_j\left(\frac{\mathbf{q}_{t,i}^T \mathbf{k}_{j,i}}{\sqrt{d_h + d_h^R}}\right)\mathbf{v}_{j,i}^C, \tag{46}$$

$$\mathbf{u}_t = W^O[\mathbf{o}_{t,1}; \mathbf{o}_{t,2}; ...; \mathbf{o}_{t,n_h}], \tag{47}$$

Liu, Aixin, et al. "Deepseek-v3 technical report." *arXiv preprint arXiv:2412.19437* (2024).

# LLaMA

- LLaMA has a similar architecture as GPT-2, with some minor differences:
  - RMSNorm is used to replace LayerNorm
  - SwiGLU
  - MLPs with gating

Touvron, Hugo, et al. "Llama: Open and efficient foundation language models." *arXiv preprint arXiv:2302.13971* (2023).

NYU SAI LAB

68

# MLPs with Gating



- A lot of LLM models applies gated feed-forward network to replace the conventional FFN in the transformer.

Liu, Hanxiao, et al. "Pay attention to mlps." *Advances in neural information processing systems* 34 (2021): 9204-9215.

# How LLM is Trained?

- The loss function consists of two parts:

$$L_1(\mathcal{U}) = \sum \log P(u_i | u_{i-k}, \ldots, u_{i-1}; \Theta)$$

"A newspaper article should contain these five main components: a headline, a byline, a lead/lede paragraph, an explanation, and any other additional information."

A newspaper article should contain these five main **xxx** &longrightarrow; "components" (GPT)

A newspaper article should **xxx** these five main components: a headline, a byline, a lead/lede paragraph, an explanation, and any other additional information. &longrightarrow; "contain" (BERT)

Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018).
Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).

NYU SAI LAB

# How LLM is Trained?



$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018).

71

# Vision Language Model



"An image of two golden retrievers"

Language model

⊕

Embedding

Projection

Visual Encoder

**Text** "Describe the image"

**Image**

- A **Vision-Language Model (VLM)** is an large model that jointly processes from visual data (e.g., images, video) and textual data to understand, align, and generate multimodal content.

NYU SAI LAB

72

# LLaVA



- In Llava, the visual encoder takes the input images and produced the visual embeddings.

- The visual embeddings and textual embeddings are concatenated, which is then forwarded to the fusion model.

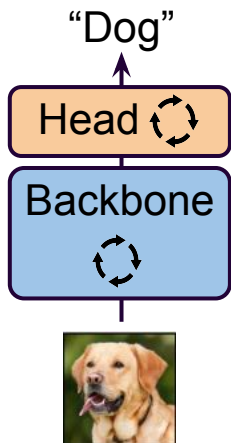Liu, Haotian, et al. "Visual instruction tuning." *Advances in neural information processing systems* 36 (2023): 34892-34916.

# Topics

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

NYU SAI LAB

# Self-Supervised Learning

- Self-Supervised Learning (SSL) is a paradigm that leverages intrinsic structures within unlabeled data to create pretext tasks, enabling models to learn meaningful representations that can be fine-tuned for downstream applications.
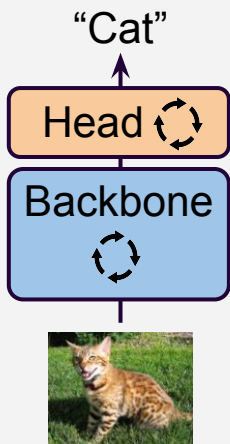
# Self-Supervised Learning

- Self-Supervised Learning (SSL) is a paradigm that leverages intrinsic structures within unlabeled data to create pretext tasks, enabling models to learn meaningful representations that can be fine-tuned for downstream applications.
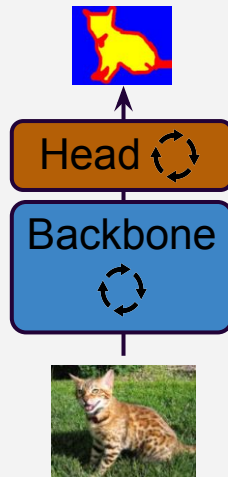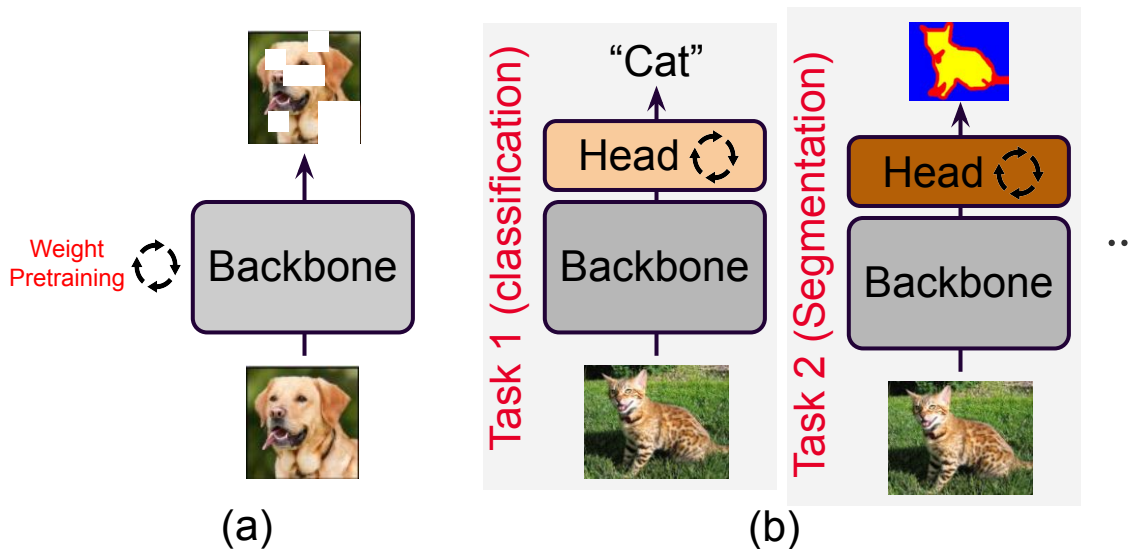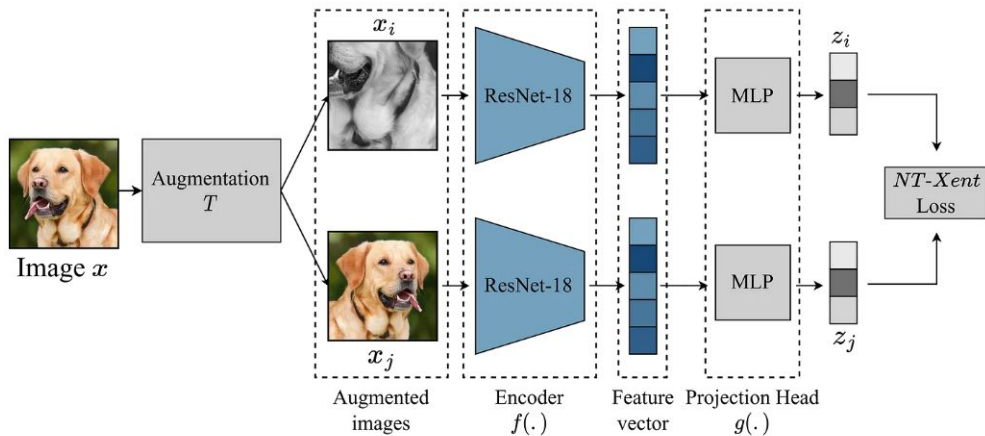


(a)　　　(b)

# Popular SSL Methods: Contrastive Learning

- Contrastive Learning is a framework in which models learn meaningful representations by contrasting positive pairs (similar data points) with negative pairs (dissimilar data points), encouraging the embedding space to capture semantic similarities and differences.



**SimCLR framwork.** Image by author.

Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." *International conference on machine learning*. PMLR, 2020.

# Popular SSL Methods: Contrastive Learning



(a) Original   (b) Crop and resize   (c) Crop, resize (and flip)   (d) Color distort. (drop)   (e) Color distort. (jitter)

(f) Rotate $\{90°, 180°, 270°\}$   (g) Cutout   (h) Gaussian noise   (i) Gaussian blur   (j) Sobel filtering

- The current augmentation approaches adopted by the AI community including random crop (with flip and resize), color distortion, and Gaussian blur.

NYU SAI LAB

# Momentum Contrast SSL (MoCo)



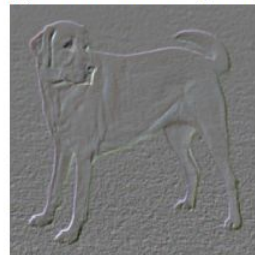- In the training process, only query encoder is updated, momentum encoder doesn't change.

$$\theta_{\mathrm{k}} \leftarrow m\theta_{\mathrm{k}} + (1 - m)\theta_{\mathrm{q}}.$$

- Only the encoder is updated.

He, Kaiming, et al. "Momentum contrast for unsupervised visual representation learning." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.

# Knowledge Distillation with No Labels (DINO)

**Algorithm 1** DINO PyTorch pseudocode w/o multi-crop.

```
# gs, gt: student and teacher networks
# C: center (K)
# tps, tpt: student and teacher temperatures
# l, m: network and center momentum rates
gt.params = gs.params
for x in loader: # load a minibatch x with n samples
    x1, x2 = augment(x), augment(x) # random views

    s1, s2 = gs(x1), gs(x2) # student output n-by-K
    t1, t2 = gt(x1), gt(x2) # teacher output n-by-K

    loss = H(t1, s2)/2 + H(t2, s1)/2
    loss.backward() # back-propagate

    # student, teacher and center updates
    update(gs) # SGD
    gt.params = l*gt.params + (1-l)*gs.params
    C = m*C + (1-m)*cat([t1, t2]).mean(dim=0)

def H(t, s):
    t = t.detach() # stop gradient
    s = softmax(s / tps, dim=1)
    t = softmax((t - C) / tpt, dim=1) # center + sharpen
    return - (t * log(s)).sum(dim=1).mean()
```
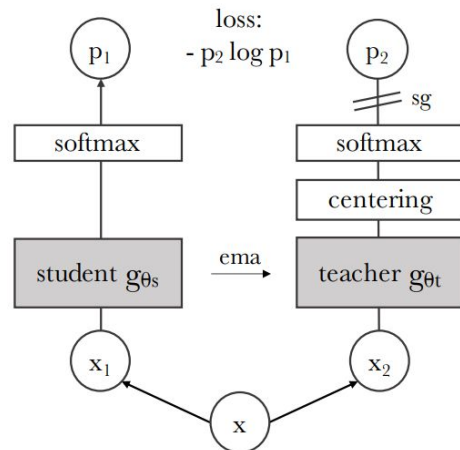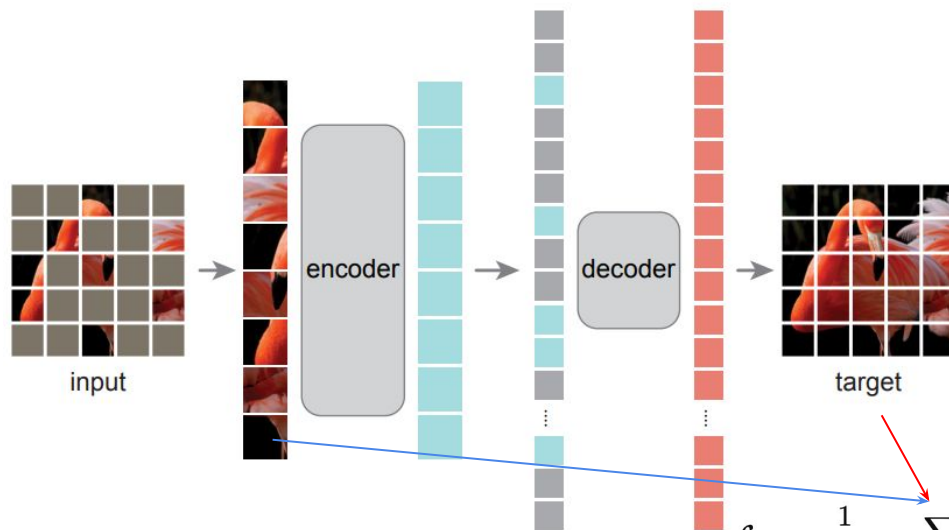


**DINO (2021)**

- During the backpropagation, only the student DNN is updated.
- The teacher updates its weight periodically using the following formula:

$$gt.params = l*gt.params + (1-l)*gs.params$$

# Masked AutoEncoder (MAE)



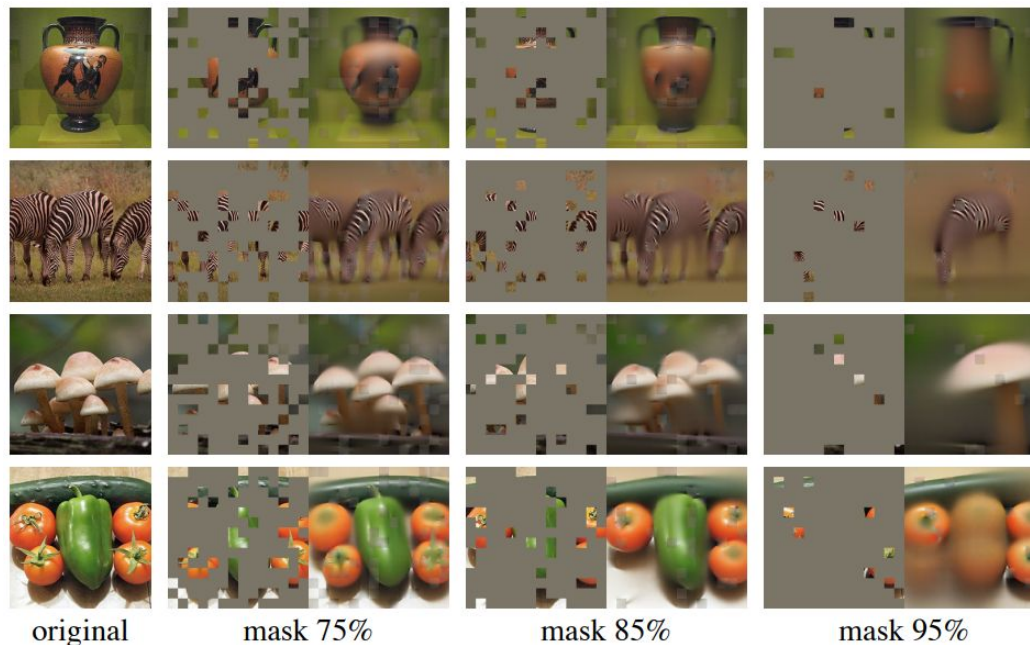- The input image is masked, the unmasked image patches will be sent to the encoder.
- The decoder will infer the masked portion of the image.

$$\mathcal{L} = \frac{1}{N_{\text{masked}}} \sum_{i \in \text{masked}} \| \hat{x}_i - x_i \|^2$$

He, Kaiming, et al. "Masked autoencoders are scalable vision learners." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.

# Self-Supervised Learning: Masked Autoencoder



original      mask 75%      mask 85%      mask 95%

He, Kaiming, et al. "Masked autoencoders are scalable vision learners." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.

# JEPA

Assran, Mahmoud, et al. "Self-supervised learning from images with a joint-embedding predictive architecture." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.

NYU SAI LAB