

Lecture 03: Transformer and Large Model

Notes

- Quiz 1 score can be found in Brightspace
- Course website: https://www.saigianzhang.com/COURSE/
- I use Brightspace to post announcements and grades
- I provide an <u>online zoom meeting</u> option for people interested in auditing the class. However, enrolled students are required to attend in person unless special condition.
- A suggested reading list which contains interesting papers can be found <u>here</u>.
- Discussion groups has been created in the Brightspace
- Course email: efficientaiaccelerator@gmail.com

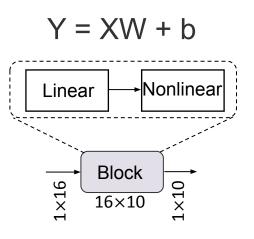


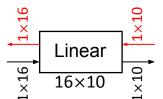
Quiz

During the training of a linear (fully connected) layer in a multi-layer perceptron (MLP), which step typically requires more computation, the forward pass or the backward pass? Why?



Fully-connected layers (Linear layers)





$$\frac{dL}{dX} = \frac{dL}{dY} \frac{dY^{\top}}{dX} = \frac{dL}{dY} W^{\top} \ \text{Derivative wrt data}$$

$$rac{dL}{db} = rac{dL}{dY}$$

$$\frac{dL}{dW} = X^{\top} \frac{dL}{dY}$$
16×10 16×1 1×10

Derivative wrt bias

Derivative wrt weight



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



Topics

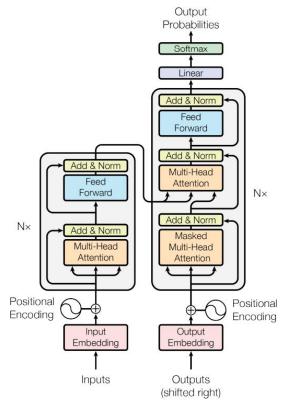
- Transformer basics
- Bert
- Vision transformer
- Large Language Model
- Self-supervised learning



Transformers

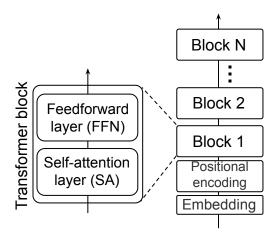
- Proposed in "Attention Is All You Need" in 2017
- The vanilla Transformer is a sequence-to-sequence model and consists of transformer blocks.

Attention Is All You Need Ashish Vaswani* Noam Shazeer* Niki Parmar* Jakob Uszkoreit* Google Brain Google Brain Google Research Google Research avaswani@google.com noam@google.com nikip@google.com usz@google.com Llion Jones* Aidan N. Gomez* † Łukasz Kaiser* Google Research Google Brain University of Toronto llion@google.com aidan@cs.toronto.edu lukaszkaiser@google.com Illia Polosukhin* ‡ illia.polosukhin@gmail.com

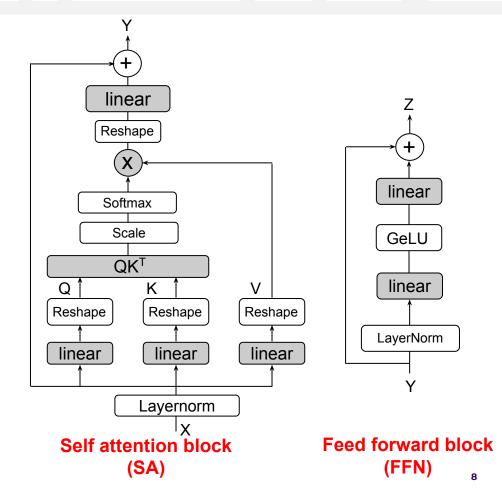




Transformers

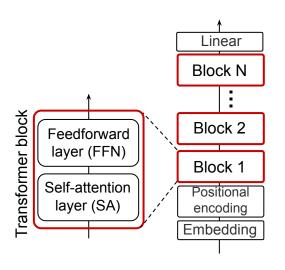


 Each transformer block includes a self-attention layer and a feedforward layer.



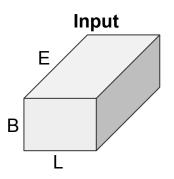


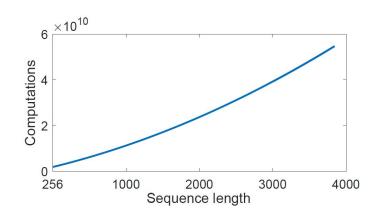
Transformers: Transformer Block





Transformers



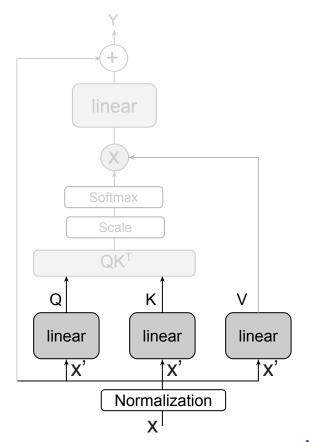


- The input contains three dimensions:
 - o B: batch
 - o L: token length
 - o E: embeddings
- The amount of computation is closely related to the token length L.
- Longer sequences are disproportionately expensive because attention is quadratic to the sequence length.



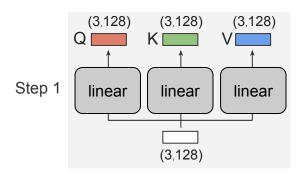
Self-Attention Block

- 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).
 - $\bigcirc \qquad (\mathsf{B},\mathsf{L},\mathsf{E}) * (\mathsf{E} * \mathsf{E}) \to (\mathsf{B} * \mathsf{L} * \mathsf{E}) \ (\mathsf{BLE}^2)$
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
 - $\bigcirc \qquad \mathsf{QK}^{\top} \rightarrow (\mathsf{B}, \mathsf{L} \times \mathsf{E}) \times (\mathsf{B}, \mathsf{E} \times \mathsf{L}) \rightarrow \ (\mathsf{B}, \mathsf{L} \times \mathsf{L})$
- 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^{\top}) * V
 - \circ (B, L*L) * (B, L*E) \rightarrow (B, L*E)
- Pass the result to the linear layer, sum with the input.





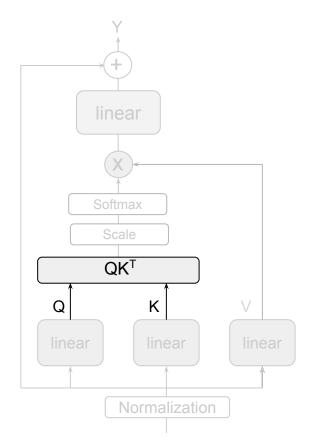
Example





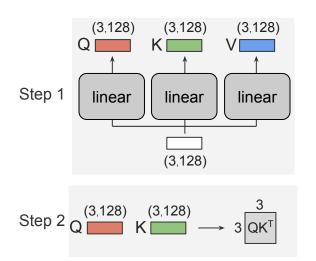
Self-Attention Block

- 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).
 - o $(B,L,E) \times (E \times E) \rightarrow (B \times L \times E)$
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
 - $\bigcirc \qquad \mathsf{QK}^{\mathsf{T}} \rightarrow (\mathsf{B}, \mathsf{L} \times \mathsf{E}) \times (\mathsf{B}, \mathsf{E} \times \mathsf{L}) \rightarrow (\mathsf{B}, \mathsf{L} \times \mathsf{L}) \ (\mathsf{BL}^2 \mathsf{E})$
- Scale and normalize the score using softmax.
 - \circ Softmax(QK^T) \rightarrow (B, L*L)
- Multiply each value vector by the softmax score.
 - \circ 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.





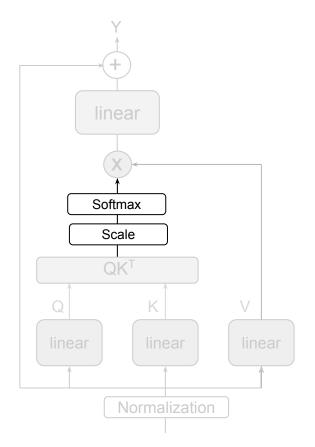
Example





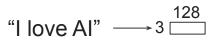
Self-Attention Block

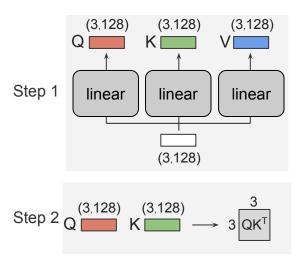
- Given input x, the first step in calculating self-attention is to create three vectors from each of the input x', denoted as: Query (Q), Key (K), Value (V).
 - \circ (B,L,E) * (E*E) \rightarrow (B*L*E)
- The second step in calculating self-attention. This will compute the attention score between each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

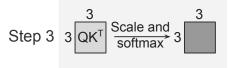




Example



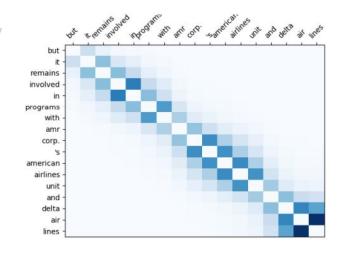






Self-Attention Block

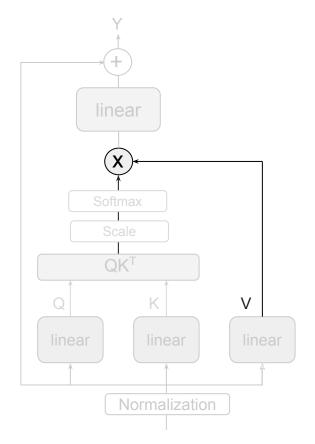
- Given input x, the first step in calculating self-attention is to create three vectors from each of the input x', denoted as: Query (Q), Key (K), Value (V).
 - \circ (B,L,E) * (E*E) \rightarrow (B*L*E)
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
 - o $QK^{\top} \rightarrow (B, L \times E) \times (B, E \times L) \rightarrow (B, L \times L)$
- 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^{\top}) * V
 - \circ (B, L*L) * (B, L*E) \rightarrow (B, L*E)
- Pass the result to the linear layer, sum with the input.





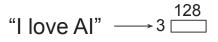
Self-Attention Block

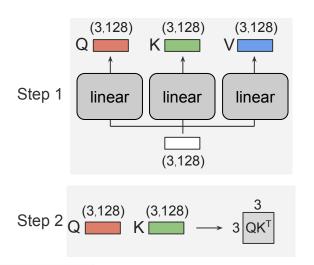
- Given input x, the first step in calculating self-attention is to create three vectors from each of the input x', denoted as: Query (Q), Key (K), Value (V).
 - \circ (B,L,E) * (E*E) \rightarrow (B*L*E)
- The second step in calculating self-attention. This will compute the attention score between each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

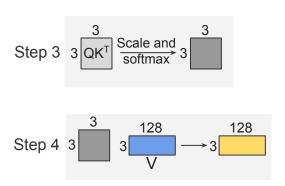




Example



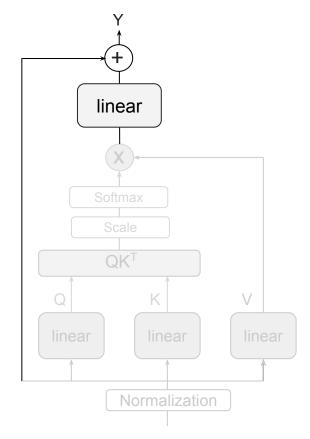






Self-Attention Block

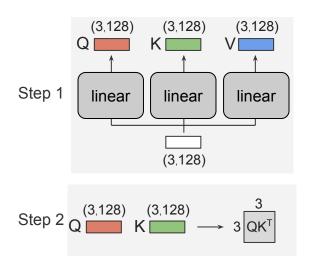
- Given input x, the first step in calculating self-attention is to create three vectors from each of the input x', denoted as: Query (Q), Key (K), Value (V).
 - \circ (B,L,E) * (E*E) \rightarrow (B*L*E)
- The second step in calculating self-attention. This will compute the attention score between each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

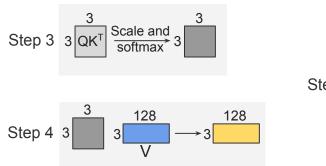


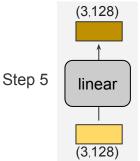


Example





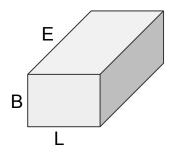


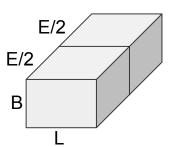


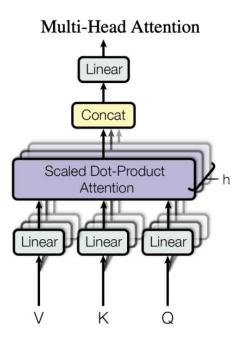


Multi-headed Attention

- 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)
 - $\bigcirc \qquad \text{Softmax}(\mathsf{QK}^\top) \star \mathsf{V} \to (\mathsf{B}, \mathsf{M}, \mathsf{L} \star \mathsf{L}) \star (\mathsf{B}, \mathsf{M}, \mathsf{L} \star \mathsf{D}) \to (\mathsf{B}, \mathsf{M}, \mathsf{L} \star \mathsf{D}) \to (\mathsf{B} \star \mathsf{L} \star \mathsf{E})$

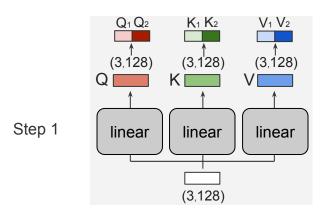




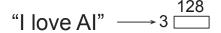


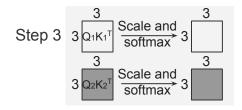


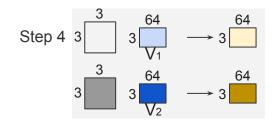
Example

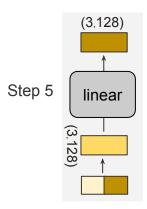








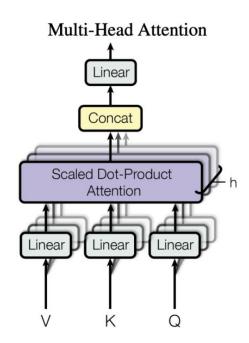






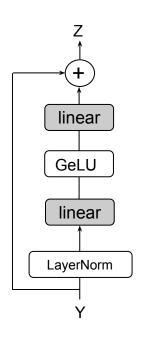
Multi-headed Attention

- 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.



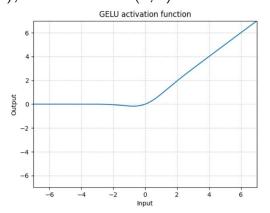


Feed Forward Layer



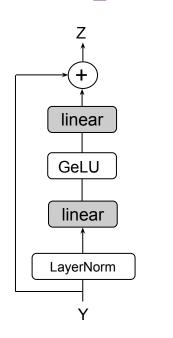
- The two linear layers are big:
 - (Ex4E) and (4ExE), E can be large (e.g., 4096)
 - This is expensive to implement.
- GeLU:

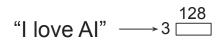
$$egin{aligned} & GeLU(x) = x\Phi(x) \ & \Phi(x) = P(y \leq x), \ ext{where Y} \sim N(0,1) \end{aligned}$$

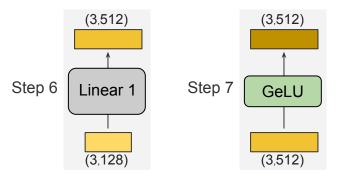


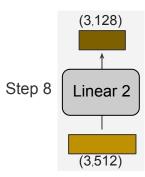


Example



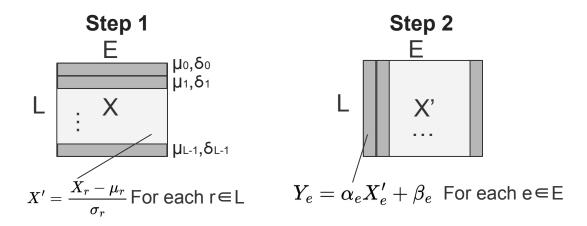








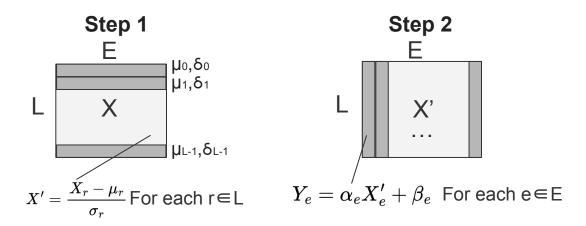
Layer Normalization



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

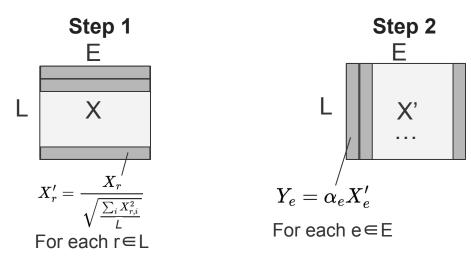


Layer Normalization



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

RMS Normalization



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

Transpose & Reshape

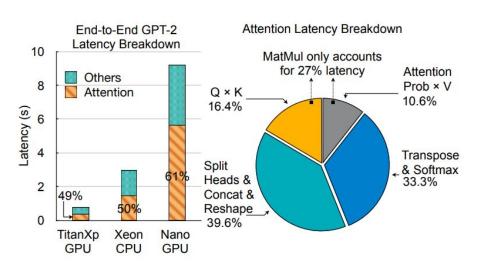


Fig. 2. End-to-End GPT-2 latency breakdown on various platforms, and attention latency breakdown on TITAN Xp GPU. Attention accounts for over 50% of total latency. Data movements account for 73% of attention latency.

Reshaping operation

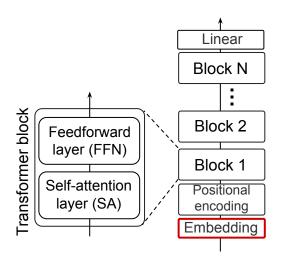
 $(B, L, E) \longrightarrow (B, M, L, E/M)$

Transpose operation

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



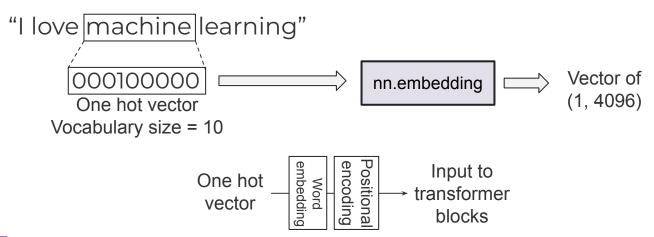
Transformers: Word Embedding



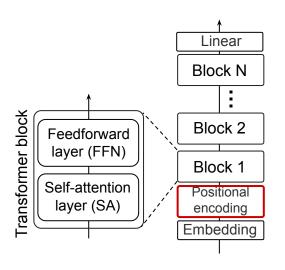


Transformers: Word Embedding

- 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,}



Transformers: Positional Encoding





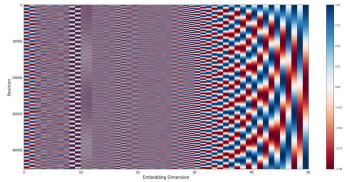
Transformers: Positional Encoding

- 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.

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{\text{model}}})$$

 $PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{\text{model}}})$

For long sequences, the indices can become quite large.
 Normalizing these index values to a range between 0 and 1 can cause issues with variable-length sequences, as each sequence would be normalized differently.

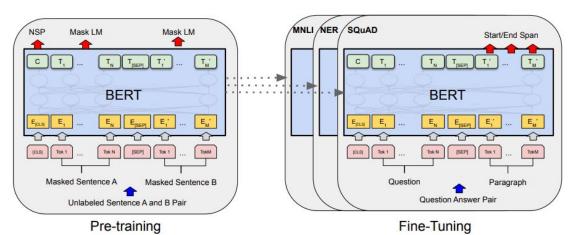


pos is the positional of the token, i is the index of the embedding.



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.

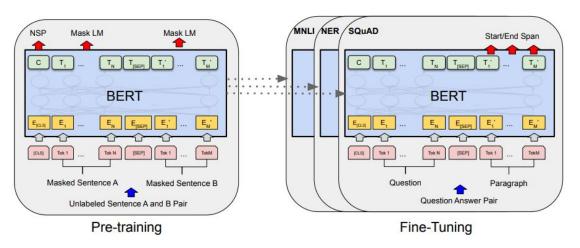




Devlin, Jacob. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).

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.

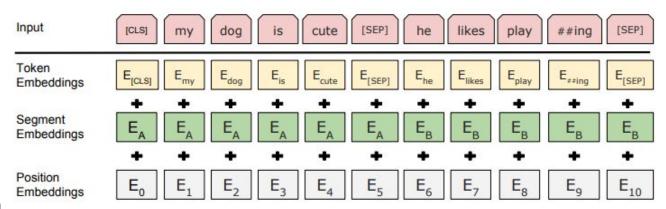




Devlin, Jacob. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).

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.



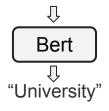


BERT Pretraining

- 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."

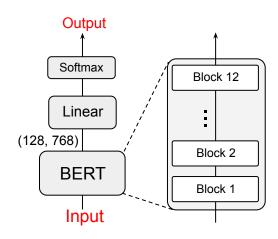
"New York University is a great school."





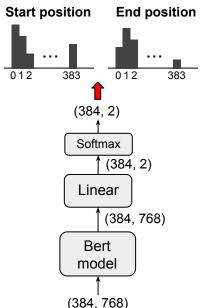
Downstream Tasks

- 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.





Downstream Tasks



Context paragraph

"Similarly, movies and television often revert to standard, clichéd snatches of classical music to convey refinement or opulence: some of the most-often heard pieces in this category include Bach's Cello Suite No. 1, Mozart's Eine kleine Nachtmusik, Vivaldi's Four Seasons, Mussorgsky's Night on Bald Mountain (as orchestrated by Rimsky-Korsakov), and Rossini's William Tell Overture."

Question

"Who wrote William Tell Overture?"

Answer

Classical_music (12, 13)

- In Question Answering tasks, an input sample consists of a context paragraph and a question with a total length of 384.
- The goal is to find, for each question, a span of text in the context paragraph that answers that question.
- BERT produces the answers by generating the start and end positions of the text span.



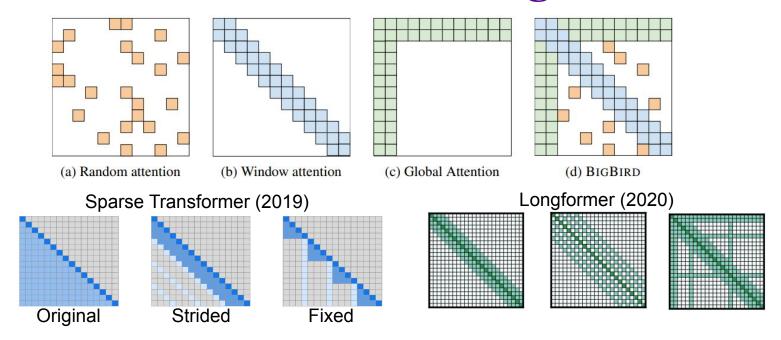
BERT Performance

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
BERTBASE	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERTLARGE	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.



Efficient Self-Attention Design



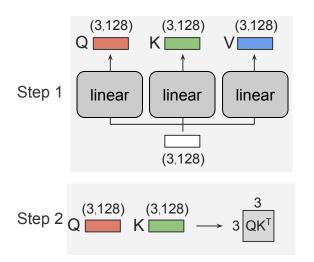


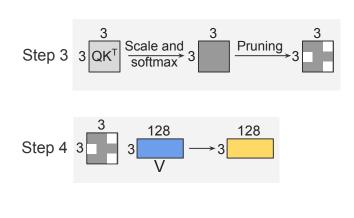
Beltagy, Iz, Matthew E. Peters, and Arman Cohan. "Longformer: The long-document transformer." *arXiv preprint arXiv:2004.05150* (2020).

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

Efficient Self-Attention Design

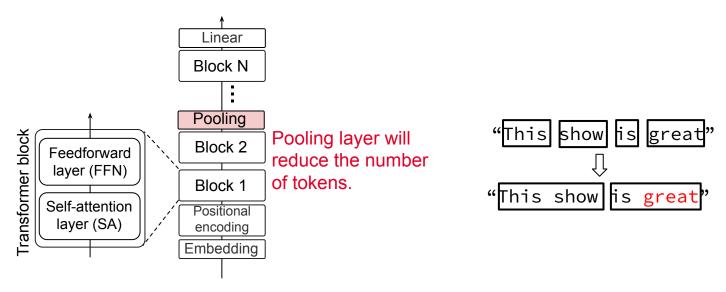
"I love AI"
$$\longrightarrow$$
 3 $\stackrel{128}{\square}$







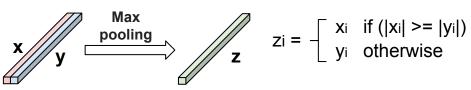
Token Merging



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

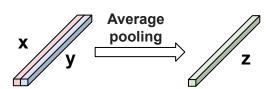


Token Merging



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

x,y are two token vectors with length of E



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

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

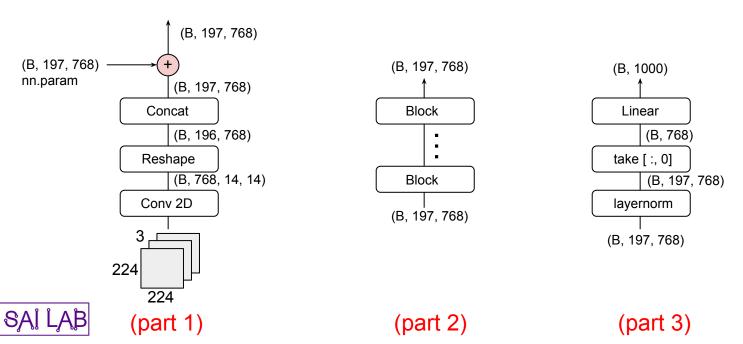


Topics

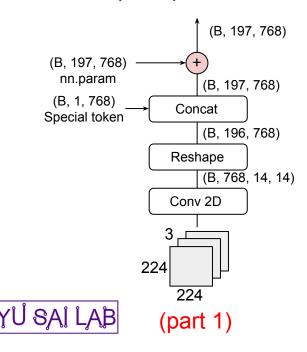
- Transformer basics
- Bert
- Vision transformer
- Large Language Model
- Self-supervised learning

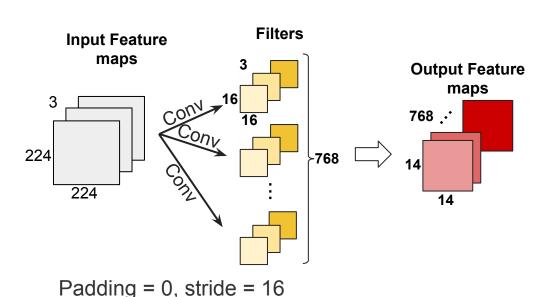


Transformer architecture can also be applied over the computer vision tasks.

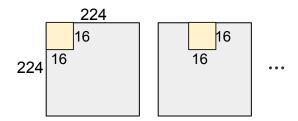


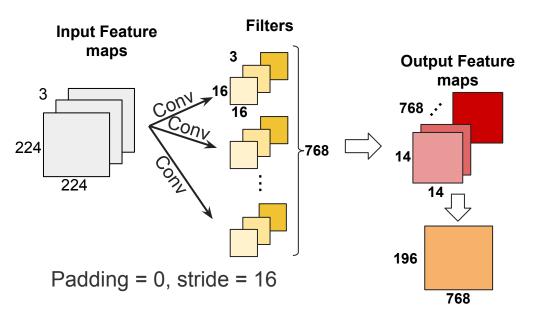
A special placeholder is introduced to aggregate global information about the whole image.





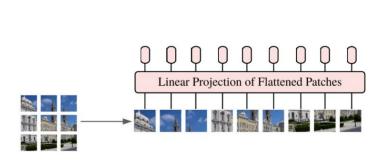
Transformer architecture can also be applied over the computer vision tasks.

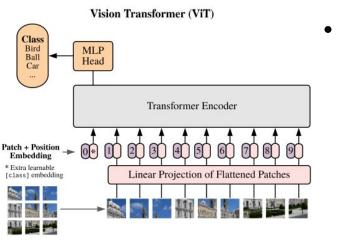




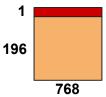


- An image C×H×W is divided into patches of C×P×P. P is 16x16 in the previous example.
- They are then flattened and linearly projected to E (e.g., 768) dimensions for a sequence of (H/P) × (W/P) tokens.



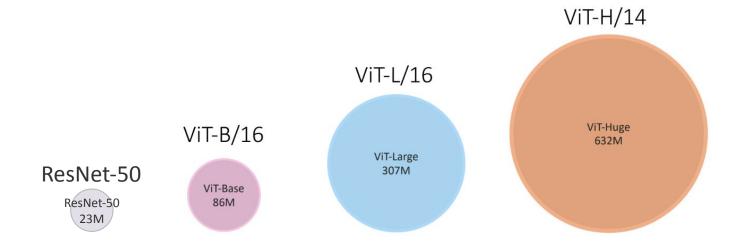


A classification
[CLS] token is
inserted at the start
of the tokens.

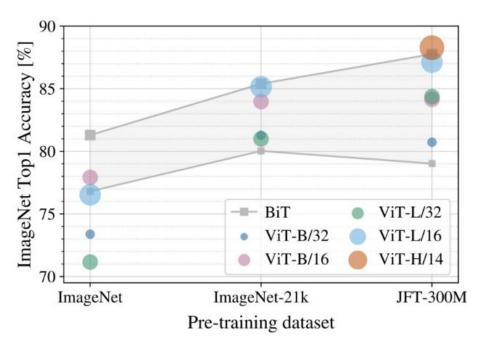




Dosovitskiy, Alexey. "An image is worth 16x16 words: Transformers for image recognition at scale." *arXiv preprint* arXiv:2010.11929 (2020).







- 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.

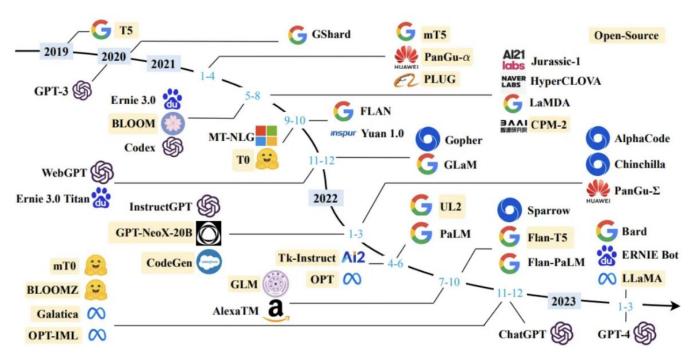


Topics

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

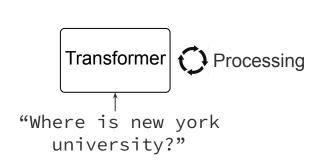


Large Language Models (LLMs)

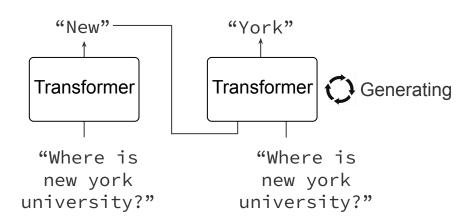




Transformers as a Generative AI Tool



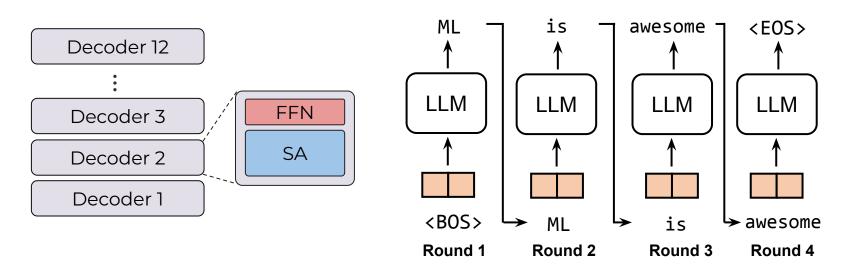
Step 1: Prefilling



Step 2: Decoding



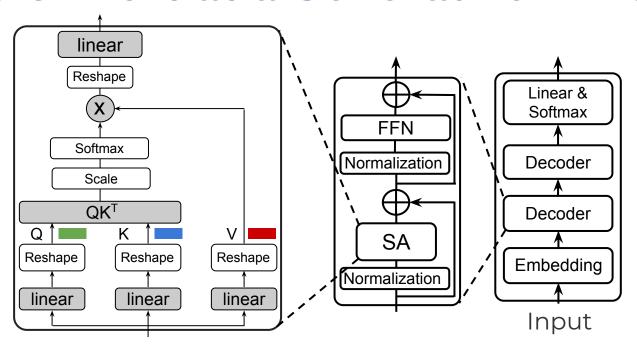
Transformers as a Generative AI Tool



Each token is generated in an autoregressive manner.



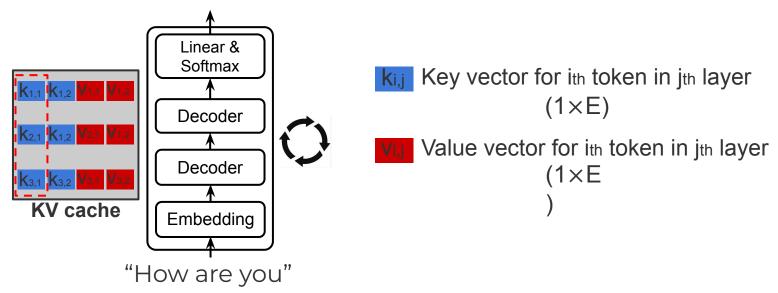
Transformers as a Generative AI Tool





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

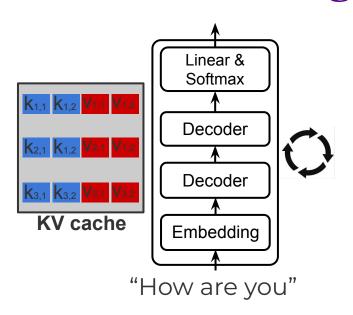
GPT-2: Prefilling

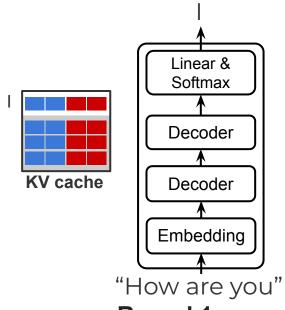


 During the prefilling stage, LLM processes the entire prompt, or context tokens jointly, saving the KV vectors into the memory.



GPT-2: Decoding



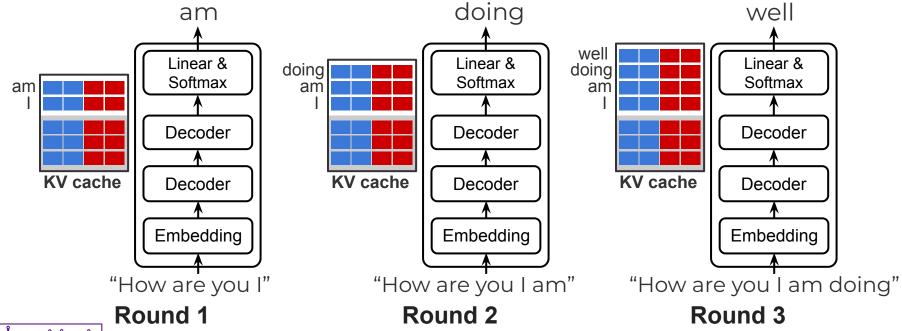


Round 1

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

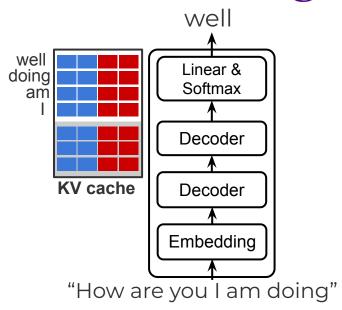


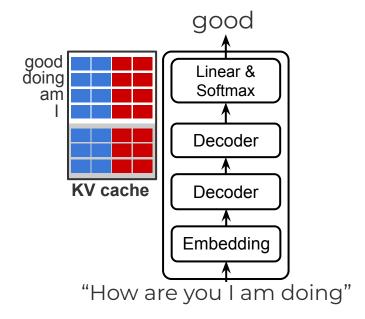
GPT-2: Decoding





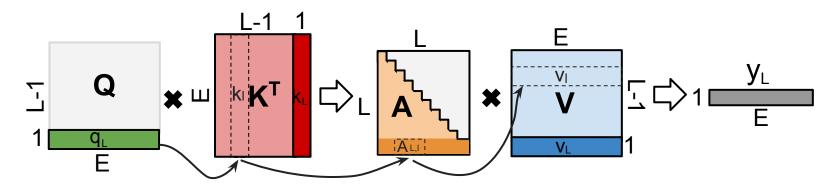
GPT-2: Decoding





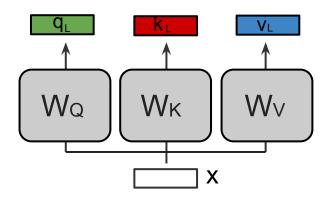
• We can simply select the token with the highest score. But better results are achieved if the model considers other words as well. So a better strategy is to sample a word from the entire list using the score as the probability of selecting that word.





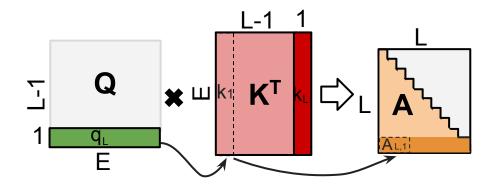
- During the decoding phase, new tokens are continuously generated and must be processed using the buffered K and V vectors to generate subsequent tokens.
- Without a KV cache, all previous K and V vectors must be recomputed, resulting in significant computational overhead.





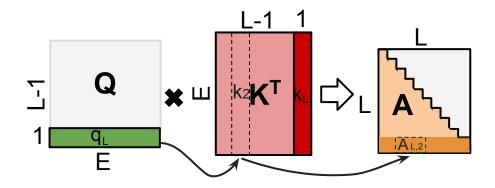
- Given the input, the q_L , k_L , v_L are first computed by passing through the linear layers.
- After that, the k_I, v_I vectors (I=1,...,L-1) are also loaded from the memory.





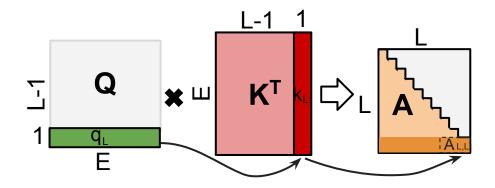
• K and V are loaded from the memory, the q vector of the current token (q⊥) is multiplied with the each of the key vector k_i, (i=1...L) to produce the result A_{L,i}





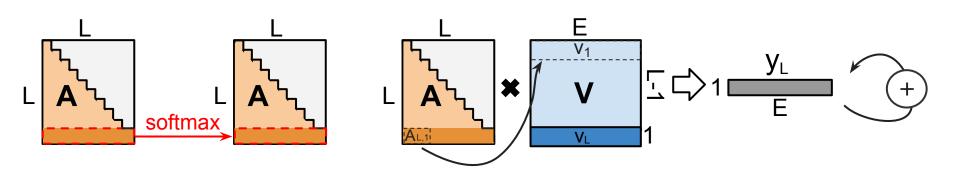
• K and V are loaded from the memory, the q vector of the current token (q⊥) is multiplied with the each of the key vector k_i, (i=1...L) to produce the result A_{L,i}





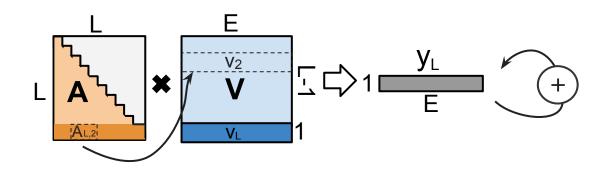
 K and V are loaded from the memory, the q vector of the current token (q⊥) is multiplied with the each of the key vector k_i, (i=1...L) to produce the result A_{L,i}





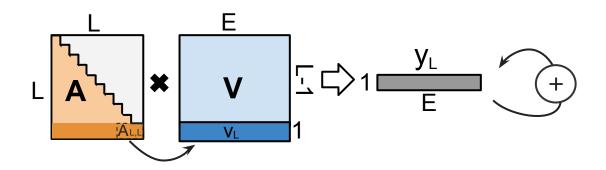
- Afterwards, the computed AL,i (i=1...L) will then passed to softmax function.
- Then each element of A_L will then multiplied with vi (i=1...L) and elementwise sum together.





- Afterwards, the computed AL,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.

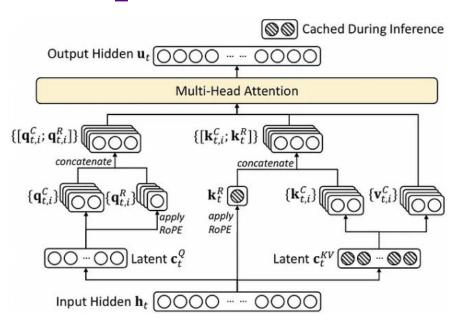




- Afterwards, the computed AL,i (i=1...L) will then passed to softmax function
- Then each element of A_L will then multiplied with vi (i=1...L) and elementwise sum together.



Deepseek V3



$$\mathbf{c}_{t}^{Q} = W^{DQ} \mathbf{h}_{t}, \tag{37}$$

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

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

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

$$\boxed{\mathbf{c}_{t}^{KV}} = W^{DKV}\mathbf{h}_{t},\tag{41}$$

$$[\mathbf{k}_{t1}^{C}; \mathbf{k}_{t2}^{C}; ...; \mathbf{k}_{t,n_b}^{C}] = \mathbf{k}_{t}^{C} = W^{UK} \mathbf{c}_{t}^{KV},$$
 (42)

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

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

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

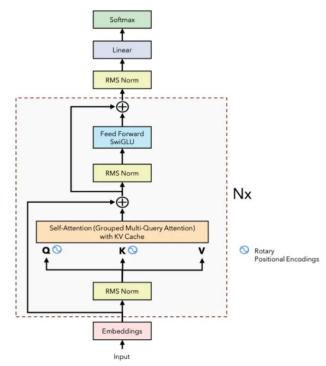
$$\mathbf{o}_{t,i} = \sum_{j=1}^{t} \text{Softmax}_{j} \left(\frac{\mathbf{q}_{t,i}^{T} \mathbf{k}_{j,i}}{\sqrt{d_{h} + d_{h}^{R}}} \right) \mathbf{v}_{j,i}^{C}, \tag{46}$$

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



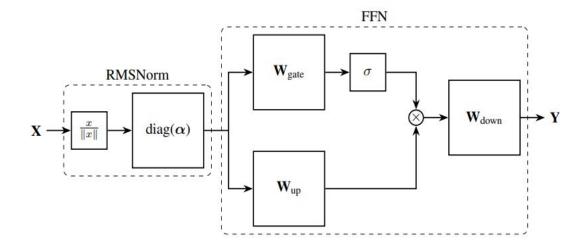
LLaMA

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





MLPs with Gating



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



How LLM is Trained?

The loss function consists of two parts:

$$L_1(\mathcal{U}) = \sum_{i=1}^{n} \log P(u_i|u_{i-k}, \dots, u_{i-1}; \Theta)$$

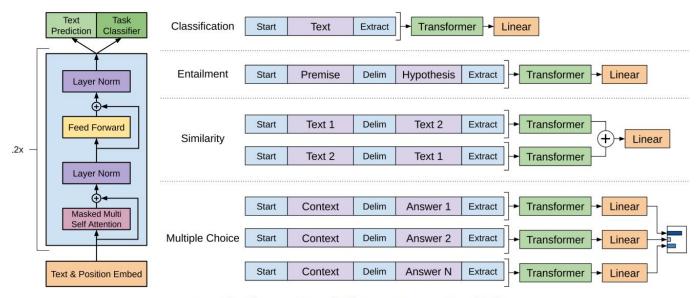
"A newspaper article should contain these five main components: a headline, a byline, a lead/lede paragraph, an explanation, and any other additional information."

A newspaper article should contain these five main **xxx** ——— "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. (BERT)



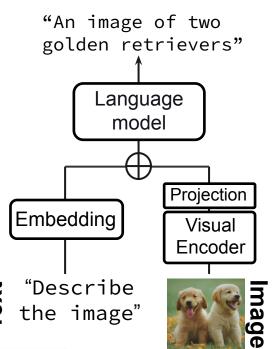
How LLM is Trained?



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



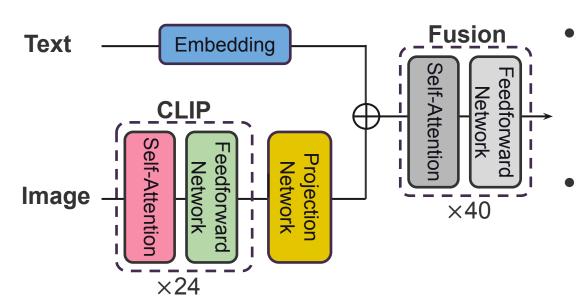
Vision Language Model



 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.



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.

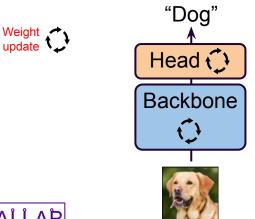
Topics

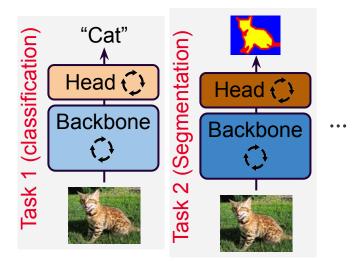
- Transformer basics
- Bert
- Vision transformer
- Large Language Model
- Self-supervised learning for visual model



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.

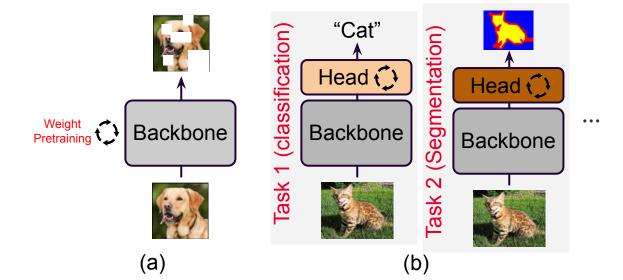






Self-Supervised Learning

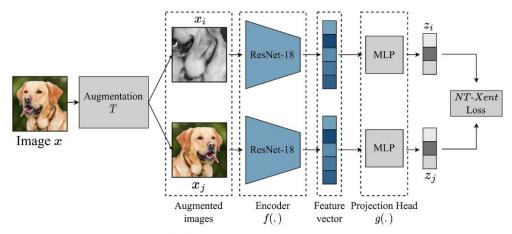
• Self-Supervised Learning (SSL) is a paradigm that leverages intrinsic structures within unlabeled data to create pretext tasks, enabling models to learn meaningful representations that can be fine-tuned for downstream applications.





Popular SSL Methods: Contrastive Learning

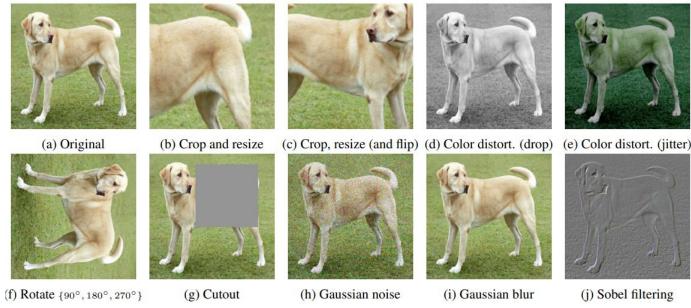
 Contrastive Learning is a framework in which models learn meaningful representations by contrasting positive pairs (similar data points) with negative pairs (dissimilar data points), encouraging the embedding space to capture semantic similarities and differences.







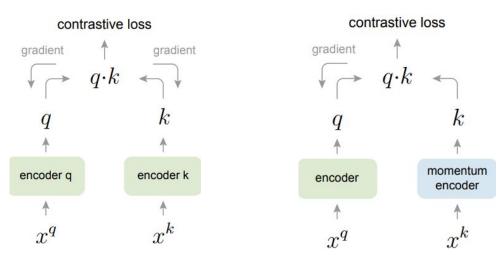
Popular SSL Methods: Contrastive Learning



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



Momentum Contrast SSL (MoCo)



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

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

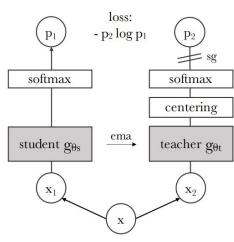
Only the encoder is updated.



Knowledge Distillation with No Labels (DINO)

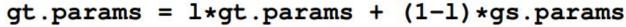
Algorithm 1 DINO PyTorch pseudocode w/o multi-crop.

```
# gs, gt: student and teacher networks
# C: center (K)
# tps, tpt: student and teacher temperatures
# 1, 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 = 1*gt.params + (1-1)*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()
```



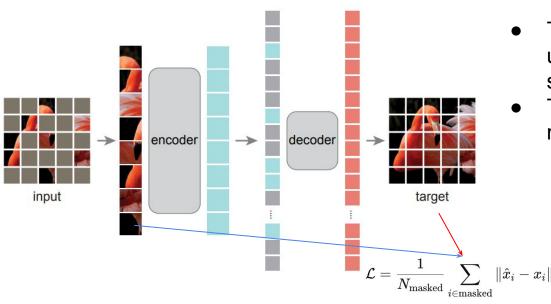
DINO (2021)

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



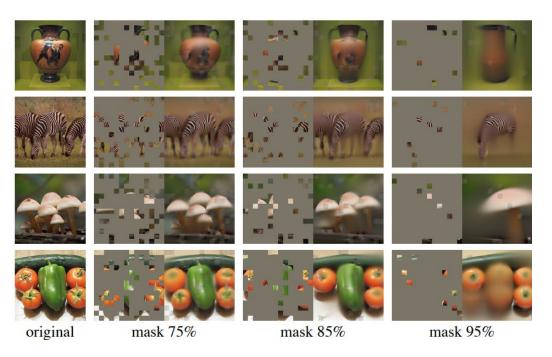


Masked AutoEncoder (MAE)



- 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.

Self-Supervised Learning: Masked Autoencoder





JEPA

