

1 Linear Algebra

1.1 Column wise decomposition.

Any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ can be decomposed into the sum of its columns:

$$\mathbf{A} = \sum_{j=1}^n \mathbf{A}_{:j} e_j^\top, \quad (1)$$

where e_j are standard basis vectors of \mathbb{R}^n . Notice that this is a rank 1 decomposition.

1.2 Row wise decomposition.

Any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ can be decomposed into the sum of its rows:

$$\mathbf{A} = \sum_{i=1}^m e_i \mathbf{A}_{i:}^\top, \quad (2)$$

where e_i are standard basis vectors of \mathbb{R}^m . Notice that this is a rank 1 decomposition.

2 Multithreading

Per-thread scratch buffers (aka thread-local workspaces) are a very common pattern to cut allocation overhead and lock contention in multi-threaded code. In this way we do not create lots of short-lived temporaries, and have better cache locality and to avoid false sharing on shared buffers.

3 LLM Training

3.1 Scaling the logits after LLM head.

We usually apply RMS norm to normalize (along the last dimension E , where it stands for model dimension, B means batch size and L means sequence length) the tensor $\mathbf{X} \in \mathbb{R}^{B \times L \times E}$ we feed into LLM head, and obtain the corresponding logits l . After RMS normalization, each tensor corresponding to the token $x_i \in \mathbb{R}^E$ will then have $\text{RMS}(x_t) = 1$. Now notice that for each coordinate $x_{i,t}, t \in [E]$, treating as a random variable, its variance is given by

$$\text{Var}(x_{i,t}) = \mathbb{E}[x_{i,t}^2] - (\mathbb{E}[x_{i,t}])^2, \quad (3)$$

and if it is zero-mean (or small), then $\text{Var}(x_{i,t}) \simeq \mathbb{E}[x_{i,t}^2]$, which is to say that second moment reflects the variance.

The next step is to use the empirical observation that for linear layers, hidden vectors tend to be approximatedly rotation-invariant (isotropic), i.e., each coordinate behaves like the others, so we can use the second moment over the coordinate in a token to replace the actual second moment. And the former, is given by

$$\text{Var}(x_{i,t}) \simeq \frac{1}{E} \sum_{t=1}^E x_{i,t}^2 = 1. \quad (4)$$

Now we start to consider the logits, which is generated by

$$l_{j,i} = w_j^\top x_i = \sum_{t=1}^E w_{j,t} x_{i,t}.$$

If we assume each weight entry $w_{j,t}$ are i.i.d. with variance σ^2 the logits variance is give by

$$\text{Var}(l_{j,i}) = \sum_{t=1}^E \sigma^2 \text{Var}(x_{t,i}) \simeq E\sigma^2.$$

So the standard deviation $\sim \sqrt{E}$. To ensure that logits do not scale with the model dimension, we scale it by \sqrt{E} .

4 Attention

4.1 Multihead Self Attention (MHA)

Consider an input tensor $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$ to an attention layer, where B is the batch size, L is the sequence length, and d_{model} is the model dimension.

1. The first step involves computing queries, keys, and values. We have three matrices, \mathbf{W}_q , \mathbf{W}_k , and $\mathbf{W}_v \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$, and simultaneously perform the following operations

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_q; \quad \mathbf{K} = \mathbf{X}\mathbf{W}_k; \quad \mathbf{V} = \mathbf{X}\mathbf{W}_v. \quad (5)$$

These operations are vectorized, meaning that for each sequence b in the batch of size B , we do

$$\mathbf{Q}_b = \mathbf{X}_b \mathbf{W}_q \quad \forall b \in [B].$$

\mathbf{W}_q , \mathbf{W}_k , and \mathbf{W}_v are trainable parameters shared across the entire batch. The resulting \mathbf{Q} , \mathbf{K} , and \mathbf{V} have the shape $\mathbb{R}^{B \times L \times d_{\text{model}}}$.

2. Next, for multihead attention, we reshape \mathbf{Q} , \mathbf{K} , and \mathbf{V} from shape $\mathbb{R}^{B \times L \times d_{\text{model}}}$ into $\mathbb{R}^{B \times H \times L \times d_{\text{head}}}$, where H is the number of attention heads and d_{head} is the dimension of each head. To achieve this, we first divide d_{model} into H heads, resulting in shapes of $\mathbb{R}^{B \times L \times H \times d_{\text{head}}}$. Then we rearrange into $\mathbb{R}^{B \times H \times L \times d_{\text{head}}}$. Conceptually, each head uses a subset of dimensions from d_{model} to compute scores between queries and keys along the sequence dimension L . We will use the following notations $\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h$ to denote the per head tensor in $\mathbb{R}^{B \times 1 \times L \times d_{\text{head}}}$ for each head $h \in [H]$.
3. In the next step, we perform the attention calculation:

$$\begin{aligned} \mathbf{S}_h &:= \text{Scores}_h(\mathbf{Q}_h, \mathbf{K}_h) = \frac{\mathbf{Q}_h \mathbf{K}_h^\top}{\sqrt{d_{\text{head}}}} + \mathbf{M} \\ \mathbf{A}_h &:= \text{Attention}_h(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{S}_h) \mathbf{V}_h, \quad \forall h \in [H]. \end{aligned}$$

The scaled multiplication of \mathbf{Q}_h and \mathbf{K}_h^\top is vectorized, resulting in $\mathbf{S}_h \in \mathbb{R}^{B \times 1 \times L \times L}$ and $\mathbf{S} \in \mathbb{R}^{B \times H \times L \times L}$.¹ Optionally, we could use a mask matrix to mask out certain tokens, an example would be the causal self attention. To stabilize the gradients, we element-wise divide raw scores by $\sqrt{d_{\text{head}}}$.

¹Here \mathbf{S} is the stack of \mathbf{S}_h along dimension H .

This scaling choice can be justified because each element of $\mathbf{Q}_h \mathbf{K}_h^\top$ represents a dot product between vectors of dimension d_{model} . The variance of this dot product scales as $\text{Var}(\langle q_h, k_h \rangle) \sim d_{\text{model}} \sigma_q^2 \sigma_k^2$. Since variance scales quadratically, we divide by $\sqrt{d_{\text{head}}}$. The softmax operation turns the scores after masking into probabilities, along the last dimension.² Imagine $z = [z_1, \dots, z_L] \in \mathbb{R}^L$ is a row vector, then essentially, softmax defines the operation:

$$\sigma(z)_i := \frac{e^{z_i}}{\sum_{j=1}^L e^{z_j}}. \quad (6)$$

Sometimes we use a numerically stable version to replace it

$$\tilde{\sigma}(z) := \frac{e^{z_i - \max(z)}}{\sum_{j=1}^L e^{z_j - \max(z)}}. \quad (7)$$

It is worth mentioning that in single head attention (scaled dot product), the complexity of computation is $\mathcal{O}(BL^2 d_{\text{model}})$, while for multihead attention, it is the same since we do $\mathcal{O}(H \times BL^2 d_{\text{head}}) = \mathcal{O}(BL^2 d_{\text{model}})$.

4. Finally, we concatenate and mix attention outputs from all heads. Concatenation involves first transposing \mathbf{A}_h to $\mathbb{R}^{B \times L \times H \times d_{\text{head}}}$ and then merging the last two dimensions into $\mathbf{A} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$. This concatenated result is projected using a matrix \mathbf{W}_O , as follows:

$$\text{MHA}(\mathbf{X}) = \mathbf{A} \mathbf{W}_O. \quad (8)$$

The final output retains the shape $\mathbb{R}^{B \times L \times d_{\text{model}}}$.

There a bunch of reasons why we are using multi heads instead of scaled dot product attention.

- **Diversity of learned attention patterns:** Each head learns different attention patterns in parallel. A single attention head computes only one set of attention scores.
- **Subspace specialization:** Instead of operating in d_{model} , each head projects to a lower dimension subspace d_{head} . This suggests that each head operates in a distinct feature subspace.
- **Improved gradient flow and representation mixing:** Independent paths improve gradient flow and richness of learned representations.

Notice that the computational cost are the **SAME**!

4.2 Multi Query Attention (MQA)

In MQA, different heads have its own query, but share the same key and value. Specifically, for a head h , we have

$$\begin{aligned} \mathbf{S}_h &:= \text{Scores}(\mathbf{Q}_h, \mathbf{K}) = \frac{\mathbf{Q}_h \mathbf{K}^\top}{\sqrt{d_{\text{head}}}} + \mathbf{M} \\ \mathbf{A}_h &:= \text{Attention}_h(\mathbf{Q}_h, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{S}_h) \mathbf{V}, \quad \forall h \in [H]. \end{aligned}$$

This means that for each head h , we have a separate $\mathbf{Q}_h \in \mathbb{R}^{B \times 1 \times L \times d_{\text{head}}}$ and shared $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{B \times 1 \times L \times d_{\text{head}}}$. Compared to standrad MHA, MQA has the following features:

²This is to say that for each $L \times L$ matrix, we softmax every row.

- **Reduced parameter count:** Each head has its own query only, shared key and value.
- **Smaller activation size (memory usage):** Now $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{B \times 1 \times L \times d_{\text{head}}}$, so the activation size is smaller.
- **Reduced KV cache (fast inferencing):** For transformer based models such as GPT, we generate text one token at a time. To avoid recomputing attention over all previous tokens on every step, we cache k, v (key and value vectors) for all previously seen tokens. Specifically, in standard MHA, for each layer and token generated, we need $2BHLd_{\text{head}}$ for cached \mathbf{K}, \mathbf{V} . In MQA, we share \mathbf{K} and \mathbf{V} so that the cost becomes $2BLd_{\text{head}}$.
- **Minimal accuracy loss.** Used in GPT-3.5, PaLM, LLaMA, etc.

4.3 KV cache

The motivation for KV caching is to enable efficient inference — both in terms of compute time and memory bandwidth. At inference time only, autoregressive models input a sequence of tokens $\{x_t, \dots, x_{t+L-1}\}$ to generate the next token x_{t+L} . To avoid recomputing key and value vectors for all previous tokens every time, we cache the k, v pairs corresponding to the tokens $x_t \dots, x_{t+L-1}$ in the forward pass. Then, when generating x_{t+L+1} , we can reuse the vectors for cached $x_{t+1}, \dots, x_{t+L-1}$ and we only need to compute k, v for x_{t+L} . This mechanism is known as the **KV cache**.

4.4 Grouped Query Attention (GQA)

GQA is like an interpolation between MQA and MHA, where we ask groups of heads to share \mathbf{K}, \mathbf{V} . Specifically, let $g(h)$ be a function that maps a head h to its corresponding group index, then we have

$$\begin{aligned} \mathbf{S}_h &:= \text{Scores}(\mathbf{Q}_h, \mathbf{K}_{g(h)}) = \frac{\mathbf{Q}_h \mathbf{K}_{g(h)}^\top}{\sqrt{d_{\text{head}}}} + \mathbf{M} \\ \mathbf{A}_h &:= \text{Attention}_h(\mathbf{Q}_h, \mathbf{K}_{g(h)}, \mathbf{V}_{g(h)}) = \text{softmax}(\mathbf{S}_h) \mathbf{V}_{g(h)}, \quad \forall h \in [H]. \end{aligned}$$

Benefits:

1. Less memory than MHA.
2. Flexible compute/memory tradeoff by controlling the number of k, v heads.

It is used in LLaMA 2 and Mistral.

4.5 Multihead Latent Attention (MLA)

Before we go into details, we need to first differentiate between self attention and cross attention.

- Self attention is the case when $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ comes from the same input sequence. Its typically used in encoder blocks of BERT, GPT, LLaMA, etc, and decoder blocks in GPT, T5, etc.
- Cross attention refers to the case when \mathbf{Q} comes from one sequence but \mathbf{K}, \mathbf{V} comes from another sequence. It is typically used in the case that decoder attends to encoder outputs (T5, BART), and the case of vision language models where text attends to image, and perceiver-style latent attention.

We can formulate cross attention in the following way: Let $\mathbf{Z} \in \mathbb{R}^{B \times M \times d_{\text{model}}}$ (expanded from $\mathbb{R}^{1 \times M \times d_{\text{model}}}$.) be a target sequence (queries) and $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$ be a source sequence (keys and values),

$$\mathbf{Q} = \mathbf{Z}\mathbf{W}_q; \quad \mathbf{K} = \mathbf{X}\mathbf{W}_k; \quad \mathbf{V} = \mathbf{X}\mathbf{W}_v. \quad (9)$$

Notice that M could be different than L . In the case of cross-attention with a latent array, we often have, $M \ll L$, which significantly reduces computational cost by avoiding full self-attention over the entire input sequence. We then compute the attention scores and softmax:

$$\mathbf{S}_h := \text{Scores}(\mathbf{Q}_h, \mathbf{K}_h) = \frac{\mathbf{Q}_h \mathbf{K}_h^\top}{\sqrt{d_{\text{head}}}} \quad (10)$$

$$\mathbf{A}_h := \text{Attention}_h(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h) = \text{softmax}(\mathbf{S}_h) \mathbf{V}_h, \quad \forall h \in [H]. \quad (11)$$

Notice that in this case $\mathbf{S}_h \in \mathbb{R}^{B \times 1 \times M \times L}$ and $\mathbf{A}_h \in \mathbb{R}^{B \times 1 \times M \times d_{\text{head}}}$ for head $h \in [H]$. After concatenation, we result in $\mathbf{A} \in \mathbb{R}^{B, M \times d_{\text{model}}}$, which is like we are focusing on a smaller sequence. In MLA, latent vector it self is a learnable parameter and shared accross a batch. These latents act like information bottleneck that extract useful features from the long input $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$.

4.6 Latent Transformer Block

This is a key design of Perceiver (2021, DeepMind), Set Transformer and efficient transformers for long inputs (e.g., audio, video, documents). Essentially it can be viewed as cross attention followed by latent self-attention. Mathematically speaking, we are giving a vector $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$. First we are using the latent vector $\mathbf{Z} \in \mathbb{R}^{1 \times M \times d_{\text{model}}}$ (expanded accross batch dimension.) and the input $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$, we have

$$\mathbf{A}_1 = \text{Attention}(\mathbf{X}_q = \mathbf{Z}, \mathbf{X}_k = \mathbf{X}, \mathbf{X}_v = \mathbf{X}),$$

which essentially asks "What should I learn from all of you tokens?". After this step: each latent now contains information extracted from the input. Notice that now $\mathbf{A}_1 \in \mathbb{R}^{B \times M \times d_{\text{model}}}$ is a compressed representation of \mathbf{X} , extracted by the latent array. Then we do normal self attention on the latent variable:

$$\mathbf{A}_2 = \text{Attention}(\mathbf{X}_q = \mathbf{A}_1, \mathbf{X}_k = \mathbf{A}_1, \mathbf{X}_v = \mathbf{A}_1)$$

where each latent vector is allowed to look at other latents, share what they learned and refine itself.

```
for each block:
    z = z + CrossAttention(q ← z, k ← x, v ← x)
    z = z + SelfAttention(q ← z, k ← z, v ← z)
```

Before we actually feed the \mathbf{Z} and \mathbf{X} into the attention block and the final feed forward layer, we first do normlization (LayerNorm in the case of my code).

4.7 Pre- and Post- Normlization

In general, there are two ways of doing layer normlization, Post-LN and Pre-LN. In the original implementation of transformer, post-LN is used. However, pre-LN has become the modern default, which is used in GPT-2/3/4, T5, LLaMA, PaLM, Perceiver, etc. The benefits of using pre-LN includes the follows:

- **(Help gradient flow & increasing training stability):** In a deep stack, residual paths carry the untouched signal forward. With Pre-LN, those residual paths also carry unit-variance, zero-mean activations (because they are already normalized). That keeps gradients well-scaled and prevents the exploding / vanishing issues that appeared when stacking 24 - 100+ layers with Post-LN. Empirically, Pre-LN lets you train hundreds (even thousands) of layers with a stable learning rate schedule, whereas Post-LN often needed warm-up tricks or gradient clipping.
- **(Easier optimization of very long sequences):** Cross-entropy loss is applied after the final LayerNorm. With Post-LN every sub-layer's output is renormalized, the network must constantly "undo" those shifts. Pre-LN leaves the residual branch untouched, so the model can accumulate information across time steps or tokens without repeatedly rescaling it.
- **(Faster convergence):** Many ablations show 1.3 - 1.5x faster convergence for GPT/T5 style models when switching from Post-LN \rightarrow Pre-LN. This is because every tensor that flows straight down the stack (both forward activations and backward gradients through the residual skip) has mean 0 and variance 1, which helps stabilize second-moment estimate quickly for Adam.
- **(Safer with half-precision / mixed-precision):** Normalizing before the high-variance matrix multiplications keeps activations in a narrower numeric range, reducing overflow/underflow risk in FP16/BF16 training.

LayerNorm: Mathematically speaking, consider an input \mathbf{X} in the space $\mathbb{R}^{B \times L \times d_{\text{model}}}$, for the l -th token in the b -th batch

$$\mathbf{L}[b, l, :] = \text{LayerNorm}(\mathbf{X}[b, l, :]) = \gamma \cdot \frac{\mathbf{X}[b, l, :] - \mu_{b,l}}{\sqrt{\sigma_{b,l}^2 + \epsilon}} + \beta,$$

where

$$\mu_{b,l} = \frac{1}{d_{\text{model}}} \sum_{i=1}^{d_{\text{model}}} \mathbf{X}[b, l, i], \quad \sigma_{b,l}^2 = \frac{1}{d_{\text{model}}} \sum_{i=1}^{d_{\text{model}}} (\mathbf{X}[b, l, i] - \mu_{b,l})^2,$$

$\gamma, \beta \in \mathbb{R}^{d_{\text{model}}}$ are learned shift and scale vectors. Notice that statistics are computed per sample, per position, no batch coupling, so the network behaves the same in training and inference and is robust to batch-size 1. Its benefits includes:

- **(Zero-mean, unit-var features).** Keeps dot-products in a predictable range, resulting in stable softmax gradients.
- **(Identical behaviour in training / inferencing).** Important for autoregressive generation where batch size changes.
- **(Works with any sequence length).** No running-average statistics needed.

In a simpler form, layernorm can be written as

$$\mathbf{L} = \gamma \odot \frac{\mathbf{X} - \mu}{\sqrt{\sigma + \epsilon}} + \beta,$$

where all operations are elementwise. Notice that $\gamma, \beta \in \mathbb{R}^{d_{\text{model}}}$ are learnable parameters.

4.8 Dropout

During Training: This refers to During training we randomly set a fraction of the sub-layer output activations to zero and scale the survivors. Specifically, for $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$ and the drop out rate $0 < p < 1$, we sample a binary training mask $\mathbf{M} \in \{0, 1\}^{B \times L \times d_{\text{model}}}$, using the Bernoulli(q) distribution where $q = 1 - p$, we then apply

$$\hat{\mathbf{X}} = \frac{1}{q} \mathbf{X} \odot \mathbf{M}.$$

q is referred to as the keep probability, and we are trying to make $\mathbb{E}[\hat{\mathbf{X}}] = \mathbf{X}$.

During Inferencing: At inference time the mask is removed and \mathbf{X} passes through unchanged. Code:

```
torch.nn.Dropout(p=p)
```

4.9 DCA

This one seems to be less well known, on medical image tasks.

4.10 Linear Attention

Let us first look at the asymptotic time complexity and memory complexity of attention.

Operation	Complexity
Projections $\mathbf{X} \mapsto (\mathbf{Q}, \mathbf{K}, \mathbf{V})$	$\mathcal{O}(BLd_{\text{model}}^2)$
Score computation \mathbf{QK}^\top (per head)	$\mathcal{O}(BHL^2d_{\text{head}})$
Softmax over scores	$\mathcal{O}(BHL^2)$
Weighted sum softmax(\mathbf{QK}^\top) \mathbf{V}	$\mathcal{O}(BHL^2d_{\text{head}})$
Output projection (concat heads $\rightarrow d_{\text{model}}$)	$\mathcal{O}(BLd_{\text{model}}^2)$

Table 1: Computation complexity of the main steps in a standard self-attention block. Here B = batch size, L = sequence length, H = number of heads, d_{head} = head dimension, and $d_{\text{model}} = H \cdot d_{\text{head}}$.

For naive self-attention, we can see that as sequence length becomes longer, the dominating term would be $\mathcal{O}(BHL^2d_{\text{head}})$, which becomes $\mathcal{O}(BHL_q L_k d_{\text{head}})$, which is quadratic in L .

Memory consumption:

- Store \mathbf{Q} , \mathbf{K} and \mathbf{V} : $\mathcal{O}(BLd_{\text{model}})$
- Store attention scores/probs: $\mathcal{O}(BHL^2)$
- Output activations for backprob add similar terms.
- By FlashAttention / PyTorch SDPA (streaming): we keep compute the same $\mathcal{O}(BHL^2d_{\text{head}})$ but reduce memory from $\mathcal{O}(BHL^2)$ to roughly $\mathcal{O}(BLd_{\text{model}})$ by not materializing the full $L \times L$ matrix.

Scenario	Time (leading term)	Memory (leading term)
Self-attention (train)	$O(BHL^2d_{\text{head}})$	$O(BHL^2)$
Cross-attention	$O(BHL_qL_kd_{\text{head}})$	$O(BHL_qL_k)$
With FlashAttention / SDPA	same time	$O(BLd_{\text{model}})$
Autoregressive decoding (per token at step t)	$O(Htd_{\text{head}})$	cache $O(Ld_{\text{model}})$

Table 2: Asymptotic compute and memory complexity for different attention scenarios. Here B = batch size, L = sequence length (train), H = number of heads, d_{head} = per-head dimension, $d_{\text{model}} = H \cdot d_{\text{head}}$, L_q/L_k = query/key lengths for cross-attention, and t = current decoding step in autoregressive inference. Autoregressive decoding is with KV cache.

Linear attention: The linear attention is introduced to avoid the quadratic dependence using attention kernels which allow us to obtain linear dependence. Note that the attention of token i is given by

$$\text{attn}_i = \frac{\sum_{j=1}^L \exp\left(\frac{q_i^\top k_j}{\sqrt{d_{\text{head}}}}\right) v_j}{\sum_{j=1}^L \exp\left(\frac{q_i^\top k_j}{\sqrt{d_{\text{head}}}}\right)}.$$

To compute it, we need to perform $\mathbf{Q}\mathbf{K}^\top$ and apply softmax rowwise. We introduce the kernel trick just as we did in SVM, imagine we can write

$$\exp\left(\frac{q_i^\top k_j}{\sqrt{d_{\text{head}}}}\right) = \phi(q_i)^\top \phi(k_j),$$

for some feature map $\phi : \mathbb{R}^{d_{\text{head}}} \mapsto \mathbb{R}^{d_\phi}$ with non-negative outputs, then

$$\text{Numerator: } \sum_{j=1}^L \phi(q_i)^\top \phi(k_j) v_j = \sum_{j=1}^{\max}$$

4.11 Sparse Attention

5 Normalizations

Besides the layer norm mentioned in the previous chapter, there are other norms used quite often.

5.1 BatchNorm

Batch norm is rarely used in language modeling especially in modern architectures. It is not sequence or position aware, and it requires consistent statistics over a batch, with variable-length sequences in language modeling, the batch statistics can be unstable and unreliable. Furthermore, it mixes information from each sequences, but ideally we want to keep it separate. However, it is great for vision tasks

Mathematically, given an input $\mathbf{X} \in \mathbb{R}^{B \times d}$, for each feature j , we have

$$\mu_j = \frac{1}{B} \sum_{i=1}^B \mathbf{X}[i, j], \quad \sigma_j^2 = \frac{1}{B} \sum_{i=1}^B (\mathbf{X}[i, j] - \mu_j)^2.$$

Then each value is normalized in a way that

$$\hat{X}[i, j] = \gamma[j] \frac{X[i, j] - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} + \beta[j],$$

where $\gamma, \beta \in \mathbb{R}^d$ are learnable parameters per feature / channel.

If we are given an input 2D/Convolutional data $\mathbf{X} \in \mathbb{R}^{B \times C \times H \times W}$, we have

$$\mu_c = \frac{1}{BHW} \sum_{b,h,w} \mathbf{X}[b, c, h, w], \quad \sigma_c^2 = \frac{1}{BHW} \sum_{b,h,w} (\mathbf{X}[b, c, h, w] - \mu_c)^2,$$

basically, we are computing mean and variance per channel. Notice that **LayerNorm is still the default in most Vision Transformer (ViT-style) backbones**. In CNN, sometimes GroupNorm replaces BatchNorm because it's batch-size independent and mixes well with vision features.

5.2 RMSNorm

Its name is Root Mean Square Layer Normalization, which has become popular in some recent LLMs, especially as a lighter and sometimes more numerically stable alternative to LayerNorm. It normalizes only by root mean square of the features³ instead of mean and variance. Given $X \in \mathbb{R}^{B \times L \times d_{\text{model}}}$

$$\text{RMS}_{b,l} = \sqrt{\frac{1}{d_{\text{model}}} \sum_{i=1}^{d_{\text{model}}} \mathbf{X}[b, l, i]^2 + \epsilon}$$

Then, each element in \mathbf{X} is normalized by

$$\hat{X}[b, l, i] = \frac{\mathbf{X}[b, l, i]}{\text{RMS}_{b,l}} \cdot \gamma_i, \quad \forall i \in [d],$$

where γ_i is the i -th component of the learnable scaling factor $\gamma \in \mathbb{R}^{d_{\text{model}}}$.

LLaMA (and most of Meta-Llama 1/2/3 checkpoints) swap the original LayerNorm for pre-norm RMSNorm.

Below is a complete list of RMSNorm's features.

- **(Rescales by RMS only & scaling only).**
- **(Compute & Memory efficiency).** 30 % cheaper per norm op; overall 2-6 % faster end-to-end in large LLMs.
- **(Numerically stable).** Works well with very deep pre-norm Transformers.

When does it shine: (1) Ultra-deep LLMs, (2) Inference-first or edge deployments, (3) Pre-norm architectures.

RMSNorm really does skip the “subtract-the-mean” step, so a single normalization needs one less reduction-operation and a bit less memory traffic (removing the mean saves one vector reduction, one broadcast, and one add.). This skip would not affect training stability, because pre-norm residuals absorb the offset: in modern transformers RMSNorm sits before each residual branch, any mean shift can be compensated by the next linear layer's bias. To see this: consider the following procedure, (assuming $x \in \mathbb{R}^d$)

$$\boxed{x} \xrightarrow{\text{RMSNorm}} \boxed{y = \gamma \cdot \frac{x}{\|x\|_{\text{RMS}}}} \xrightarrow{\mathbf{W}, b} \boxed{u = \mathbf{W}y + b} \rightarrow \boxed{z = x + u}$$

³Features in this context refer to the elements in dimension d_{model}

y indeed does not have mean zero, denote it as μ_y , since the next operation is affine, we can rewrite it as

$$u = \mathbf{W}(y - \mu_y \mathbf{1}) + (b + \mathbf{W}\mu_y \mathbf{1}),$$

which means that it is equivalent to adjusting the bias. During training, back-prop will simply nudge b so the network learns whatever overall shift is optimal, it does not care where that shift originates. After that u is added to the original x , but the drift in μ_y will not accumulate unchecked because: (1) The residual path still carries the original, unshifted activations, (2) The next block starts with another RMSNorm, which rescales its input (including any offset) back to a controlled RMS before new computations begin.

6 Positional Embedding

7 Precisions

Here is a summary provided by Chat-GPT o3.

Format	Bits	Exp/Mant	Dynamic range [†]	ULP@1	Typical use
FP32	32	8/23	1.2×10^{-38} - 3.4×10^{38}	2^{-23}	Master weights, optimiser
FP16 (IEEE)	16	5/10	6.1×10^{-5} - 6.6×10^4	2^{-10}	Fwd/Bwd on Volta & T4
bfloat16	16	8/7	1.2×10^{-38} - 3.4×10^{38}	2^{-7}	Default on A100/H100, TPU
TF32*	19	8/10	same as FP32	2^{-10}	GEMMs on Ampere

Table 3: Mid-/high-precision dtypes widely used during LLM training.[†]Smallest positive *normal* value to largest finite.*TF32 is a compute mode; tensors stored as FP32.

Format	Bits	Exp/Mant	Dynamic range [†]	ULP@1	Typical use
FP8-E4M3	8	4/3	9.2×10^{-8} - 4.5×10^2	2^{-3}	Research train / fast inf. (H100)
FP8-E5M2	8	5/2	3.0×10^{-10} - 5.7×10^4	2^{-2}	Wider range variant (H100)
INT8 / INT4	8/4	—	$\pm 127 / \pm 7$	1 / 1/16	Post-training inference quant.

Table 4: Very low-precision formats used for efficient training research (FP8) or deployment quantisation (INT).

A floating point number in IEEE-style formats is stored as:

Sign bit | Exponent bits | Mantissa (fraction) bits

1. **Sign bit:** (1 bit), 0 positive and 1 negative.
2. **Exponent bits:** control the scale (powers of 2). More exponent bits means wider dynamic range.
3. **Mantissa bits:** control the precision (how many distinct numbers you can represent between powers of 2). More mantissa bits means smaller ULP⁴ (more precise).

⁴Unit in the Last Place: The smallest possible difference between two representable floating-point numbers around a given value

A number is typically constructed in this way:

$$(-1)^{\text{sign}} \times (1.\text{mantissa bits}) \times 2^{\text{exponent}-\text{bias}}.$$

- **FP32** is used as the safe baseline in deep learning, it is almost always used for the master weights⁵ and optimizer state because it is numerically stable.
- **FP16** is used in forward/backward pass to save memory in Volta(Nvidia V100)/T4 mixed-precision training. Prone to overflow/underflow unless scaled (hence loss scaling).
- **Loss scaling** refers to the practice we multiply our loss by a large constant S before back propagation: $L_{\text{scaled}} = S \cdot L$, this scales up gradients which avoids gradient underflow. However, after computing the gradients, we divide them by S before updating the weights to restore correct magnitude. NVIDIA's AMP (Automatic Mixed Precision) does dynamic loss scaling, adjusting S automatically to prevent overflow.
- **bfloat16** is Brain floating point came from Google Brain, whose exp bits are same as **FP32**, but the mantissa bits are much fewer, this avoids the over/under-flow problem in **BF16**. This comes at a price of lower precision, but it is the default choice on A100 and H100.

FP16 on Volta/T4 is fast but fragile, and requires loss scaling, which is why bfloat16 became popular. It enjoys the same range as FP32, but no scaling headaches.

Its advantage includes: (i) no loss scaling needed to handle with over/under flow, (ii) half the memory of **FP32**, (iii) speed up in matmul on A100/H100, (iv) can be used as a drop in for **FP32**, because the dynamic range is the same.

For downsides: (a) precision loss, (b) optimizers often still keep master weights in **FP32** to avoid cumulative rounding error, (c) numerically sensitive operations such as softmax/ normalization, kernel still uses **FP32** internally.

- **TF32**: NVIDIA introduced with Ampere GPUs (A100, RTX 3000 series), default in cuBLAS/cuDNN matmul on Ampere if you pass FP32 inputs. FP32's size and range, FP16's precision.
- **FP8-E4M3, FP8-E5M2**: We can now tell directly from their names that FP8-E4M3 has a narrower dynamic range at a higher precision. For FP8-E4M3, $\text{ULP}@1$ is $\frac{1}{2^3} = 0.125$ and for FP8-E5M2, the $\text{ULP}@1$ becomes $\frac{1}{2^2} = 0.25$ which is very coarse.
E4M3: good for weights and activations that stay within a moderate range.
E5M2: good for gradients or loss-related values that can swing wildly in magnitude.

How they (FP8) are used on H100: It can be applied both for training and inferencing.

During **training**:

1. Keep master weights in FP32 (like with FP16/bfloat16 training).
2. Cast activations/gradients to FP8 for GEMMs inside the forward/backward pass.
3. Apply per-tensor or per-channel **scaling** to map values into FP8's limited range.
4. Often mix E4M3 for forward activations, E5M2 for backward gradients.

During **Inferencing**: Post-training quantization to FP8 for ultra-fast inference with minimal memory footprint.

⁵Master weight refers to the full-precision copy of the model's parameter that we keep during mixed-precision training.

Scaling is not optional: FP8’s numeric range is so tiny that, without actively scaling tensors before casting to FP8, we’ll either hit overflow or underflow constantly. The fix is that we multiply the tensor by a scaling factor S before FP8 conversion,

$$x_{\text{scaled}} = x \times S.$$

We may choose S so that $\max x_{\text{scaled}}$ fits nicely into FP8’s max normal value, and we store S as a separate FP32 number. We can later undo the the scaling after computation:

$$y = y_{\text{scaled}} \times S^{-1}.$$

This keeps numbers inside FP8’s safe zone while still preserving the original magnitude relationship.

Per channel vs. per tensor: Often, there are **per channel** scaling and **per tensor** scaling. Let us image that we are doing a fully-connected layer, where we have this weight matrix $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$. Per channel scaling means that we have a scaling factor for each output feature (each row), so there will be d_{out} scalars, while for per tensor scaling, we only have 1 scaling factor. In this sense, per tensor scaling is simpler, but will be suboptimal if some channels are loud and the others are quiet, while per channel scaling has better precisions, especially in convolutional layers or transformers where variance differs accross each head / filter.

Automation on H100: Notice that there is automated scaling on NVIDIA H100, which tracks running max values for each tensor, and picks E4M3 (more precision) for forward activations, E5M2 (more range) for backward gradients. The engine applies scaling transparently, so we mostly see FP8 without manually tuning the scaling factor S .

As an example of scaling, let us consider the previous example of fully connected layer, where we are expected to do $\mathbf{O} = \mathbf{W}\mathbf{X}$, originally, both \mathbf{W} and \mathbf{X} are FP32 (master weights and full precision activations). Before GEMM⁶, we independently scaled them to FP8 (to fit the range) by

$$\mathbf{X}_{\text{scaled}} = \mathbf{X} \cdot S_{\mathbf{X}}; \quad \mathbf{W}_{\text{scaled}} = \mathbf{W} \cdot S_{\mathbf{W}},$$

then we performed the matmul in FP8,

$$\mathbf{O}_{\text{scaled}} = \text{FP8GEMM}(\mathbf{W}_{\text{scaled}}, \mathbf{X}_{\text{scaled}}).$$

To restore the original magnitude, we take advantage of the scaling factor $S_{\mathbf{X}}, S_{\mathbf{W}}$ who are FP32,

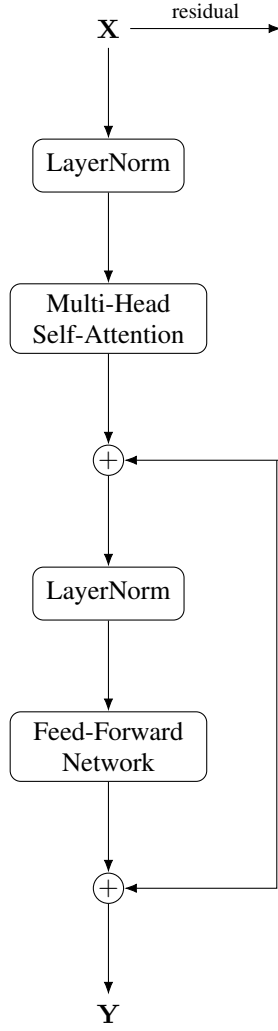
$$\mathbf{O} = \mathbf{O}_{\text{scaled}} \cdot \frac{1}{S_{\mathbf{X}} S_{\mathbf{W}}}.$$

8 Encoder and Decoder Structure

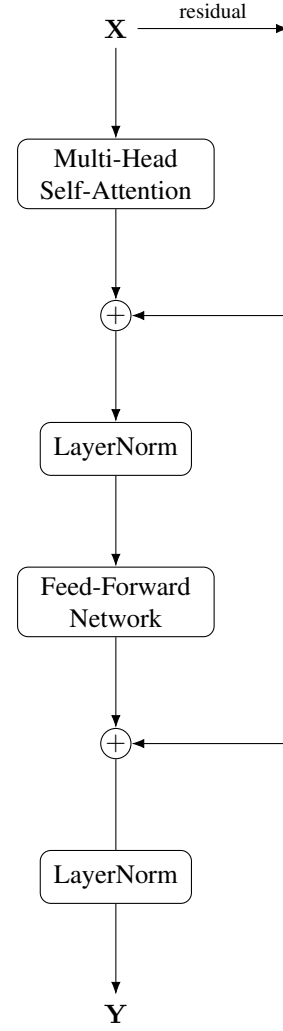
Encoder Typically, encodes looks like:

As we have discussed before, the Pre-LayerNorm fashion is the current state of the art as its training is more stable, as well as its convergence. It helps gradient flow, and deep model stability, does not requires learning rate warm up and works better with long context. It do has downside though, Pre-LN has a weaker normalizing effect so sometimes people add an extra final LayerNorm after the last layer for output stability (“final LN” in GPT models).

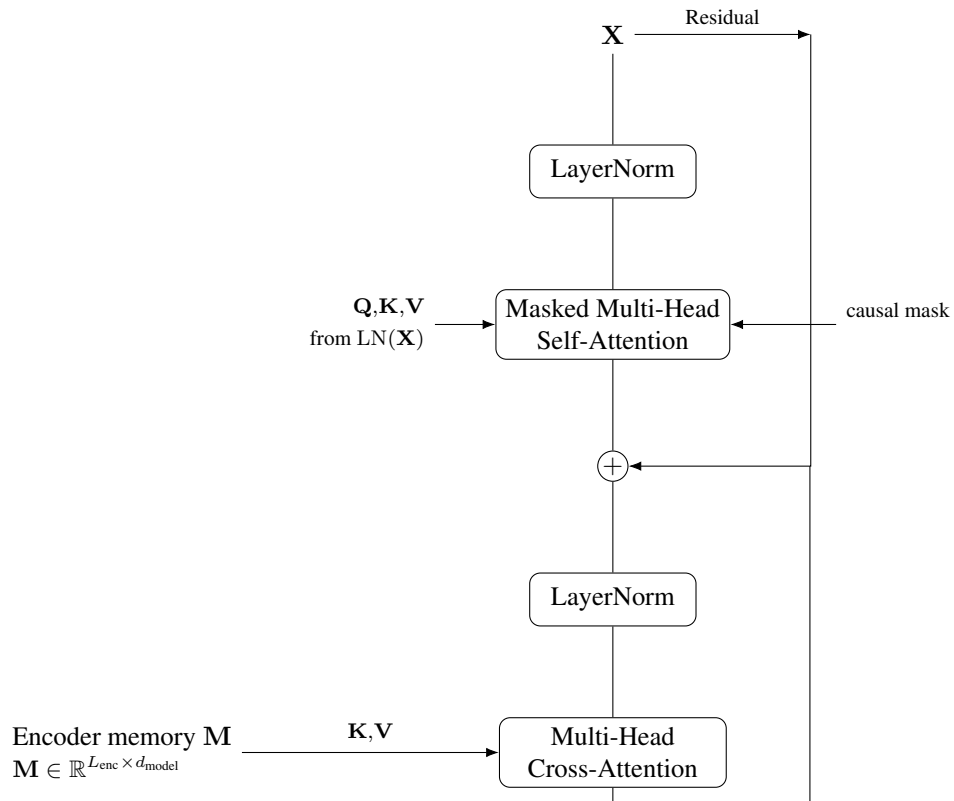
⁶General Matrix-Matrix Multiplication, in the form of $\mathbf{D} = \alpha \cdot \mathbf{A}\mathbf{B} + \beta\mathbf{C}$, where α, β are scalars. This is the workhorse of modern deep learning.



Pre-LayerNorm Transformer block



Post-LayerNorm Transformer block



Decoder As a further explanation besides the figure: this \mathbf{M} is the matrix of encoder hidden states produced by the last encoder layer. It is what the decoder attends to in cross attention. Let us image that an input $\mathbf{X} \in \mathbb{R}^{B \times L_{\text{enc}} \times d_{\text{model}}}$ is fed into the encoder blocks. After embedding and positional encoding, we have

$$\mathbf{H}^0 \in \mathbb{R}^{B \times L_{\text{enc}} \times d_{\text{model}}}.$$

Each encoder will do

$$\mathbf{H}^{k+1} = \text{Encoder}(\mathbf{H}^k) \in \mathbb{R}^{B \times L_{\text{enc}} \times d_{\text{model}}}, \quad \forall k \in [N_{\text{enc}}],$$

And \mathbf{M} is exactly $\mathbf{M} = \mathbf{H}^{N_{\text{enc}}}$. As a result, for the cross attention layer: $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{B \times L_{\text{enc}} \times d_{\text{model}}}$ and $\mathbf{Q} \in \mathbb{R}^{B \times L_{\text{dec}} \times d_{\text{model}}}$, so the output of decoder block will always be $\mathbb{R}^{B \times L_{\text{dec}} \times d_{\text{model}}}$.

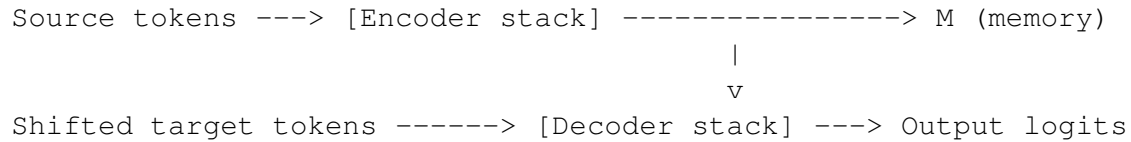
The two types of blocks are different:

- Encoder: Self attention only, no masking, aiming to build a contextual representation of the entire source sequence.
- Decoder: Masked self attention, which prevents seeing future tokens so that the generation is autoregressive. Cross attention, which lets each target position attend to the encoder's output \mathbf{M} (source context) The purpose overall is to generate the target sequence one step at a time, with access to both past target tokens and the entire source sequence.

What happens if we are using one type?

- Encoder: leaking future target tokens during training (no causal mask), so the model wouldn't learn autoregressive generation.
- Decoder: useless masked self-attention in the encoder that blocks half the context for no reason, waste compute on a cross-attention sublayer when no encoder memory exists yet.

The original seq2seq transformer: Things happens in two phases,



Encoder runs once on the full source, decoder runs once (training) or incrementally (inference), always starting with shifted target embeddings as its own input stream. The input of encoder is the source sequence $x = [x_1, x_2, \dots, x_{L_{\text{enc}}}]$, after embedding we get \mathbf{X} , while the input of the decoder is the target sequence y so far, but shifted (\hat{y}) so that the model predicts the next token.

1. We take the original gold target sequence $y = [y_1, y_2, \dots, y_{L_{\text{dec}}}]$.
2. Shift right by one and add the <BOS> special token, and obtain $\hat{y} = [\text{<BOS>}, y_1, \dots, y_{L_{\text{dec}}-1}]$, we then embed to get $\hat{\mathbf{Y}}$ and use it as an input to the decoder blocks.

8.1 Encoder only models

BERT, RoBERTa, DeBERTa, ELECTRA, Sentence-BERT (SBERT).

The architecture mostly are: • Stacks of encoder blocks only, • Full self attention, • Input sequence length stay fixed, • output contextualized embeddings for every token.

Purpose: Understand text: classification, regression, retrieval, token-level labeling (NER, POS tagging, QA span prediction). In those cases the model needs bidirectional context: token sees both left and right neighbors.

For many NLP tasks, we already have the full text and just need to analyze it, not generate it.

Advantages:

- **Better context capture:** every token attends to all others.
- **More efficient training** for non-generative tasks (no need to autoregress).
- **Easier fine-tuning** for classification tasks: just take the [CLS] embedding⁷.

8.2 Decoder only models

GPT family (GPT-2, GPT-3, GPT-4, LLaMA, Mistral, Falcon, etc.); BLOOM, OPT, Pythia; Code generation models (CodeLLaMA, StarCoder)

Architecture: • Stack of decoder blocks only • Causal self-attention (mask future positions) • **No cross-attention** • Input length = current sequence length during generation.

Purpose: Generate text: language modeling, code generation; model learns

$$P(\text{next token} \mid \text{prev tokens})$$

Advantages:

- **Simpler:** do not need an encoder, since we aim to predict the next token given the past. Same architecture works for both pretraining (predict next token) and inference (sample next token)
- **Massive scalability:** can ingest any text corpus, no need for aligned parallel data.
- **Flexible prompts:** can condition on arbitrary text in-context.

8.3 Encoder-decoder models

Often on seq2seq task, requires paired data⁸.

8.4 Prefix-decoder models

Slightly less main stream variant in the transformer family, similar to an interpolation between decoder-only and encoder-decoder models.

Key difference: In prefix decoder blocks, we only have one prefix self attention layer (causal self attention layer with prefix mask), and we do not have two attentions (cross attention + self attention) pattern

⁷In BERT and similar models, we prepend a special token [CLS] (“classification”) to the start of every input sequence before feeding it to the encoder. [CLS] has its own trainable embedding vector in the model’s vocabulary, just like any word. It is treated as position 0 in the sequence and goes through all encoder layers along with the other tokens. After the final encoder layer, [CLS] has a contextualized vector $h \in \mathbb{R}^{d_{\text{model}}}$, which encodes information from the entire sequence (thanks to self-attention).

⁸An example is that machine translation needs paired data: Source: “I like apples”, Target: “J’aime les pommes”.

as in the normal decoder block. We are essentially changing the external information from encoder to prefix region in the mask.

Imagine now we take a sequence of length $L = m + n$

$$\mathbf{x} = [\underbrace{p_1, \dots, p_m}_{\mathbf{p}}, \underbrace{g_1, \dots, g_n}_{\mathbf{g}}],$$

where \mathbf{p} is the prefix tokens which are given in full before text generation starts. This could be tokens from a source sentence (like the “encoder output” flattened into tokens), learnable virtual tokens (prefix tuning), encoded representations from another modality (e.g., image embeddings mapped to token space), or just a long text prompt. \mathbf{g} here is the generation tokens that we will generate causally one by one.

A standard decoder (decoder-only transformers) uses a strict causal mask, which is a lower triangular matrix,

$$\text{Mask}[i, j] = \begin{cases} 1 & \text{if } j \leq i \\ 0 & \text{otherwise} \end{cases}$$

for all $i, j \in [L]$, this says that token i can only attend to positions $j \leq i$. While for a prefix decoder, we relaxes the prefix region

$$\text{Mask}[i, j] = \begin{cases} 1 & \text{if } i \leq m \text{ and } j \leq m & \text{prefix attends to all prefix} \\ 1 & \text{if } i \geq m \text{ and } j \leq m & \text{generation attends to all prefix} \\ 1 & \text{if } i > m, m < j \leq i & \text{generation attends to its own past} \\ 0 & \text{otherwise} \end{cases}$$

An example is that when we want to generate a translation for a sentence \mathbf{t} (of tokens), we will use it as a prefix, and we know that t_0 is translated into s_0 already, in this case, our prefix would be $[t_0, \langle \text{sep} \rangle, s_0, \langle \text{BOS} \rangle]$ and our generation tokens would be \mathbf{t} . We include $\langle \text{BOS} \rangle$ in the prefix because it is a known token which is used to anchor.

We prefer prefix decoder sometimes, because

- **Fewer Parameters:** Parameter in the attention part drop roughly half.
- **Lower compute cost:** When prefix length is modest compared to generation length, prefix decoder can be fast overall, especially on hardware where fusing the attention into one pass matters.
- **Simpler architecture:** Do not need a separate encoder-decoder split
- **Easier caching for autoregressive generation:** In standard decoder with cross attention, we need to cache (i) decoder past key/values for masked self attention (ii) the encoder outputs for cross attention. In a prefix decoder, the prefix tokens are just part of the same KV-cache.

Down sides: (1) less flexibility in terms of shared layer stack and parameters. (2) for long prefixes, attention is more expansive.

8.5 Summary

As a summary:

Model type	Pros	Cons
Encoder-only	<ul style="list-style-type: none"> • Strong understanding with bidirectional context • Great for classification, retrieval, and embedding learning 	<ul style="list-style-type: none"> • Not suited for generative tasks
Decoder-only	<ul style="list-style-type: none"> • Simpler architecture • Easy to train generatively • Highly flexible for any prompt-based generation 	<ul style="list-style-type: none"> • Lacks bidirectional context • Weaker for pure understanding tasks without adaptation
Encoder-decoder	<ul style="list-style-type: none"> • Best for sequence-to-sequence tasks (translation, summarization, speech-to-text) 	<ul style="list-style-type: none"> • Requires paired data • Heavier compute at inference (two stacks)

Table 5: Comparison of encoder-only, decoder-only, and encoder-decoder Transformer architectures.

8.6 The Linear Layers

FFNs (Feed Forward Networks) are a part of the encoder, decoder block, there are several design choices here. They are also known as MLPs (Multi-Layer Perceptrons), and they are used to transform the input features into a higher-dimensional space, apply non-linear activation, and then project back to the original dimension. For the 1- d case, imagine we have a vector $x \in \mathbb{R}^{d_{\text{model}}}$, then the FFN does the following:

$$\text{FFN}(x) = \mathbf{W}_2 \cdot \sigma(\mathbf{W}_1 x + b_1) + b_2,$$

where $mW_1 \in \mathbb{R}^{d_{\text{ffn}} \times d_{\text{model}}}$ is the matrix that projects the input to a higher dimension, σ is a non-linear activation function (e.g., ReLU, GELU), and $\mathbf{W}_2 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ffn}}}$ projects it back to the original dimension. d_{ffn} is often $2 \text{ to } 4 \times d_{\text{model}}$. Now if we consider the true input $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$, then basically, for each sequence in the batch, we apply the same FFN to each token, and there will be no interactions between different tokens.

Why do we put FFNs in the encoder/decoder blocks? The FFN is used to introduce non-linearity into the model, allowing it to learn more complex representations. Self-attention is linear in the feature dimension for a fixed set of attention weights. It mixes tokens but does not increase the per-token expressivity much. To see this, consider the following attention formula:

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_{\text{model}}}} \right)}_{:=\mathbf{A}} \mathbf{V},$$

if the attention weight matrix \mathbf{A} is fixed, then the output is a linear combination of the input \mathbf{V} which is linear in \mathbf{X} , i.e., linear in the feature dimension.

If we **remove the FFN**, the model would reduce to mostly linear mixing layers accross tokens, which leads to collapse of expressivity (model underfits complex transformations), and much worse performance.

Design choices for FFNs: There are several design choices for FFNs, which can affect the model's performance and efficiency.

- Expansion ratio: $d_{\text{ffn}}/d_{\text{model}}$.

This is typically around $2.4\times$. In the original Transformer paper, it is $4\times$, while in LLaMA, it is $2\times$. In many large models, FFNs account for 50% – 60% of total parameters.

- Activation function: ReLU, GELU, SiLU, etc.

The original Transformer paper uses ReLU, but GELU is more popular in modern models (BERT, GPT-2/3).

ReLU: piecewise linear, fast, but can lead to dead neurons (zero gradients for negative inputs).

$$\text{ReLU}(x) = \max(0, x).$$

GELU: a smoother, probabilistic activation function that approximates the Gaussian distribution.

$$\text{GELU}(x) = x \cdot \Phi(x) = x \cdot \frac{1}{2} \left(1 + \text{erf} \left(\frac{x}{\sqrt{2}} \right) \right),$$

where $\Phi(x)$ is the cumulative distribution function of the standard normal distribution, and erf is the error function.

- Pros: **Smooth, i.e., differentiable everywhere**⁹, which enables better gradient flows. It often yields better performance in transformers, understanding and generation tasks.
- Cons: **Computationally more expensive**, slightly more than ReLU.

SiLU: another smooth activation function, models such as EfficientNet used it.

$$\text{SiLU}(x) = x \cdot \sigma(x),$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function. It is sometimes called the Swish-1 (Swish with $\beta = 1$.)

8.7 Gated Linear Units

GLU (Gated Linear Unit) takes an input, splits it into two parts (value & gate), and uses one to gate the other via an elementwise product. Mathematically, for $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$, we first apply two independent linear projections using $\mathbf{W}_g, \mathbf{W}_v \in \mathbb{R}^{d_{\text{ffn}} \times d_{\text{model}}}$,

$$\mathbf{G} = \mathbf{W}_g \mathbf{X} + \mathbf{B}_g, \quad \mathbf{V} = \mathbf{W}_v \mathbf{X} + \mathbf{B}_v;$$

We then apply a nonlinear activation to the gate part, after which we do elementwise multiplication.

$$\mathbf{G}' = \sigma(\mathbf{G}), \quad \text{GLU}(\mathbf{X}) = \mathbf{V} \odot \mathbf{G}'.$$

Note that we are applying activation to \mathbf{G} (the gate) instead of \mathbf{V} , because we are applying the non-linearity here to shape its gating behaviour. This separation between gate and value allows us to learn specialized gating pattern, which acts like an learned feature selector.

SwiGLU: we are essentially replacing the original activation function $\sigma(\cdot)$ with $\text{SiLU}(\cdot)$, which is smoother. This is used in LLaMA, PaLM, etc.

GEGLU: we are essentially replacing the original activation with $\text{GELU}(\cdot)$ or its approximations with tanh.

⁹Smoothness in this context refers to infinitely differentiable functions C^∞ , not the smoothness in the optimization theory sense.

8.8 Mixture-of-Experts FFNs

In a standard Transformer block, the FFN is a single feed-forward network (sometimes gated like SwiGLU) applied to every token. In MoE, instead of one FFN, we have E separate FFN “experts”. For each token, a router chooses a subset (often 1 or 2) of experts to run. The outputs are combined according to the router’s weights. Its advantages include

- **Scaling parameters without scaling compute:** We can have lots of experts trained but we only invoke one of them.
- **Specialization:** The model specializes experts for different types of inputs (topics, syntactic structures, etc.).

Mathematically, let $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$, E be the number of experts in this case.

The forward pass:

- (i) Router logits:

$$\mathbf{R} = \mathbf{W}_r \mathbf{X} + \mathbf{B}_r \in \mathbb{R}^{B \times L \times E},$$

where $\mathbf{B}_r = \mathbf{1} b_r^\top$ is the extended bias matrix.

- (ii) Gating probabilities:

$$\mathbf{P} = \text{softmax}(\mathbf{R}) \in \mathbb{R}^{B \times L \times E}.$$

- (iii) Top- k routing: Keep the best- k experts per token according to the probabilities, and renormalize. Specifically, for each token l in each batch b , denote $\mathcal{T}_k(b, l)$ as the indices of top- k experts according to $\mathbf{P}[b, l, :]$, we do

$$\tilde{\mathbf{P}}[b, l, d] = \frac{\mathbf{P}[b, l, d]}{\sum_{j=1}^{|\mathcal{T}_k(b, l)|} \mathbf{P}[b, l, j]}, \quad \forall d \in \mathcal{T}_k(b, l), \quad |\mathcal{T}_k(b, l)| = k,$$

and sets the other coordinates to be zero, which allows us to obtain $\tilde{\mathbf{P}}$. The renormalization here is to avoid suppressing magnitude purely because we dropped low-prob experts.

- (iv) Dispatch: For each expert e , define the packed input by selecting tokens routed to e as \mathbf{X}_e . Let N_e be the number of tokens whose top- k includes expert e , $\pi_e(\cdot)$ be the corresponding map of expert e (from $\{1, \dots, N_e\} \mapsto \{[b, l]\}$) to the original indices pair (b, l) and let $\eta_e(\cdot)$ be the map from $\{[b, l]\}$ to the index in $[N_e]$. Then the input for an expert is

$$\mathbf{X}_e \in \mathbb{R}^{N_e \times d_{\text{model}}}, \quad \text{where} \quad \mathbf{X}_e[j, :] = \mathbf{X}[\pi_e(j), :], \quad \forall j \in [N_e]$$

with the renormalized gate weight

$$g_e \in \mathbb{R}^{N_e}, \quad \text{where} \quad g_e[j] = \tilde{\mathbf{P}}[\pi_e(j), e].$$

- (v) Expert maps: Depending on the expert type (dense FFN / Gated FFN), we create

$$\mathbf{Y}_e = \text{FFN}_e(\mathbf{X}_e) \in \mathbb{R}^{N_e \times d_{\text{model}}}.$$

(vi) Gate scaling and combine: Apply per token gate (weight) to the output

$$\tilde{\mathbf{Y}}_e = \text{diag}(g_e) \cdot \mathbf{Y}_e, \quad (\text{in ML literature } g_e \odot \mathbf{Y}_e \text{ using broadcast.})$$

Then the output $\mathbf{Z} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$ would be filled such that

$$\mathbf{Z}[b, l, :] = \sum_{e \in \mathcal{T}_k(b, l)} \tilde{\mathbf{Y}}_e[\eta_e(b, l), :]$$

before outputting.

Capacity constraint: In practical implementation, we often set **capacity constraint** in its formalization, which is the limit on how many tokens each expert is allowed to process in a single forward pass. The reason to have it is

- **Load imbalance is common:** one “hot” expert could get far more tokens than the average, causing GPU memory overflow for that expert, and slower steps since we need to wait for it.
- **Hardware needs fixed allocation:** In distributed training, each expert lives on one or more GPUs, memory buffers for expert input/output must be pre-allocated. If the number of tokens per expert varies wildly, we cannot pre-allocate efficiently without wasting huge amounts of memory.

We set a capacity per expert:

$$C_e = \left\lfloor \alpha \cdot \frac{BL}{E} \right\rfloor,$$

where E is the number of experts and $\alpha \geq 1$ is the capacity factor to allow slackness. If an expert get more than C_e tokens:

1. **Dropping:** (common in Switch Transformers) Excess tokens are simply dropped for that expert: they don’t get processed there, and their contribution from that expert is zeroed.
2. **Rerouting** (less common) Excess tokens are sent to a backup expert.

If an expert get fewer than C_e tokens: one may use **padding** we pad its buffer to C_e entries for vectorization but those padded entries are ignored.

Load balancing auxiliary loss: We need an extra loss term during training so that the router doesn’t collapse onto just a few experts. The reason is that, without any extra incentive, the model may prefer some experts a lot more than others, leaving many experts underused or unused. We have

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{task}} + \lambda \cdot \mathcal{L}_{\text{aux}},$$

where $\mathcal{L}_{\text{task}}$ is the task loss, and in the common GShard / Switch Transformer¹⁰, we have

$$\mathcal{L}_{\text{aux}} = E \cdot \sum_{e=1}^E f_e \cdot m_e.$$

¹⁰Both of these are specific MoE implementations from Google that popularized large-scale sparse FFNs. Switch is basically a simplified, more efficient GShard.

f_e is defined as the fraction of tokens assigned to expert e ,

$$f_e = \frac{N_e}{BL} = \frac{1}{BL} \sum_{b=1}^B \sum_{l=1}^L \mathbf{1}^\top \cdot \{e \in \mathcal{T}_k(b, l)\},$$

where $\{e \in \mathcal{T}_k(b, l)\}$ is the indicator vector (sum of corresponding one hot vector) in $\mathbb{R}^{B \times L}$ suggesting if expert e is selected by this token, while m_e is defined as

$$m_e = \frac{1}{BL} \sum_{b=1}^B \sum_{l=1}^L P[b, l, e],$$

which is the average probability of assigning to expert e . One can show that the loss is smaller when f_e and m_e are more evenly distributed across experts. \mathcal{L}_{aux} can be thought of as an load balancing regularizer.

The all to all trick: In distributed training, experts are sharded across devices (e.g., 1 expert per GPU or multiple experts per GPU.) Naively, we must

1. Send tokens to the device hosting their assigned expert.
2. Process them locally in the expert's feed-forward network.
3. Send back the processed outputs to the original device to continue the model pipeline.

This is a typical many-to-many communication pattern, where every GPU may need to send tokens to every other GPU. The most straightforward way is to let all tokens sent to all devices before filtering using a series of `all_gather` operations, but this can waste bandwidth.

We can do it with the **all to all** trick, which uses a collective communication primitive where each GPU directly sends only the tokens that the other GPUs need, in one coordinated call. In NCCL and similar libraries

```
torch.distributed.all_to_all_single(output, input, ...)
```

In practice, MoE frameworks do two all-to-all per MoE layer:

1. **Dispatch:** tokens \rightarrow owning expert GPU.
2. **Combine:** expert outputs \rightarrow original token order.

In this way, redundant communication is reduced. Real world examples include GShard / Switch Transformer / DeepSpeed-MoE / Megatron-MoE.

Processes and threads: The legacy `nn.DataParallel` (**DP**) is based on one process multi-threading, which is generally slower and less scalable (single optimizer state, host-side bottlenecks, GIL contention, extra device = host hops). This is why **DDP** (multiprocess, one rank per GPU) is recommended.

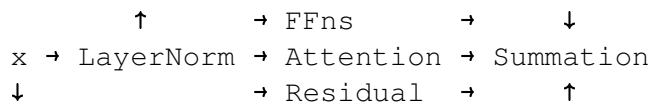
Reasons:

1. **CUDA context is per-process:** Each process gets a clean, single-GPU context, which leads to simpler, safer memory management and fewer heisenbugs.
2. **Communication stacks (NCCL)** are designed around ranks = processes.

3. **Python GIL:** Python threads can't execute Python bytecode truly in parallel. Many CUDA ops release the GIL(Global Interpreter Lock), but coordination/launch logic still contends. One process per GPU avoids that bottleneck.
4. **Failure isolation:** If one GPU OOMs or a kernel errors, you don't take down every replica in the same process.
5. **Deterministic performance:** Independent schedulers per process avoid thread scheduling contention and reduce cross-device interference within a single interpreter.

8.9 Parallel Attention + FFN

The parallel design changes the flow so that attention and FFN happen at the same time on the same input (after normalization), instead of one after the other. Their outputs are summed together with the residual connection in a single step.



Why do we do this:

- **Parallelism:** Attention and FFN can be computed concurrently.
- **Fewer LayerNorms:** Saves computation.
- **Information Flow:** In sequential design, the FFN only sees the post-attention representation. In parallel, FFN sees the original representation directly, possibly preserving more raw token information.
- **Better gradient flow:** Shorter path between input and output in the computation graph.

This architecture is used in some newer LLMs, e.g. PaLM and GPT-JT, to improve speed without hurting performance.

Some practical notes:

- **Scaling:** Some designs scale the attention and FFN outputs to balance their contributions.
- **DropPath / Stochastic depth:** In residual networks, instead of always computing the residual branch, we may randomly drop it during training according to some probabilities, so that the model learns to rely on multiple paths, not just one. In parallel attention, we may consider dropping attention branch and FFN branch independently. In this way we reduces overfitting by introducing randomness, and encourages robustness.
- **Variance stability:** When we sum two residual branches instead of one, we risk blowing up the variance of activations early in training. This is another reason why we do scaling:

$$x' = x + \frac{1}{\sqrt{2}} (\text{attn_out} + \text{ffn_out})$$

We can also control by careful weight initializations, or considering set a learnable scaling parameters like α_{attn} and α_{ffn} which is small intially and train up.

9 Computing the Number of Parameters

10 Fintuning