

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 LLM Training

2.1 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

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

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

3 Attention

3.1 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 $\sqrt{d_{\text{head}}}$. 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 $\mathcal{O}(H \times BL^2 d_{\text{head}}) = \mathcal{O}(BL^2 d_{\text{model}})$.

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

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

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**!

3.2 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:

- **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.

3.3 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**.

3.4 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.

3.5 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 $\mathbf{X} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$.

3.6 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).

3.7 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 → 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

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

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.

3.8 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)
```

3.9 DCA

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

3.10 Linear Attention

3.11 Sparse Attention

4 Normalizations

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

4.1 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 normlized 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.

4.2 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.

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

- **(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 $x \in \mathbb{R}^d$)

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

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.

5 Positional Embedding

6 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 1: 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 2: Very low-precision formats used for efficient training research (FP8) or deployment quantisation (INT).

7 Computing the Number of Parameters