

Residual Connection (Skip Connection) :

A Residual connection means : instead of only passing the transformed output of a layer, we also add the original input back to it.

$$\text{Output} = F(x) + x$$

Where,

x = input to the layer

$F(x)$ = output of the layer (after transformation)

Why ?

1. Helps gradients flow \rightarrow reduces Vanishing gradient problem.
2. Keeps original information \rightarrow Model doesn't forget the input.
3. Easier training \rightarrow Model can learn "adjustments" (residuals) instead of the whole mapping.

Example :

If input $= x = 5$

and the layer computes $F(x) = 2$

$$\text{Residual output} = F(x) + x = 2 + 5 = 7.$$

In Transformers : Every Attention and Feedforward block has a residual connection around it, followed by Layer Norm.

Normalization in Transformers

- # Keeps activations stable \rightarrow avoids Exploding / Vanishing gradients.
- # Balances features so no token dominates.
- # Works with residual connections to stabilize sums.
- # Uses LayerNorm (not BatchNorm) \rightarrow better for sequence data.

In short : Normalization = stability + smooth training + better convergence.

Batch Normalization (Vs) Layer Normalization

1. Batch Normalization:

\Rightarrow where it normalizes : Across the batch dimension (Example in the mini-batch).

\Rightarrow How it's work :

For each feature (dimension of the hidden layer), it computes mean and variance across all Example in the batch.

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_{i,j}, \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_{i,j} - \mu_j)^2$$

The Normalize Each feature :

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Intuition : Make each feature dimension have

mean 0 and Variance 1 across the batch.

Best for : computer vision (CNNs), where batches are large and features behave consistently.

Example:

$$X = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

Normalize across the batch dimension (rows), per feature (column).

column 1 (features across samples): $[1, 4]$

$$\Rightarrow \text{Mean} = (1 + 4) / 2 = 2.5$$

$$\Rightarrow \text{Variance} = ((1 - 2.5)^2 + (4 - 2.5)^2) / 2$$

$$= (2.25 + 2.25) / 2 = 2.25$$

$$\Rightarrow \text{Variance} =$$

$$\Rightarrow \text{Std} = \sqrt{2.25} = 1.5$$

$$\Rightarrow \text{Normalized} = [(1 - 2.5) / 1.5, (4 - 2.5) / 1.5]$$

$$= [-1.0, 1.0]$$

Column 2 (features across samples): $[2, 5]$

$$\text{Mean} = 3.5, \text{Variance} = 2.25, \text{Std} = 1.5$$

$$\text{Normalized} = [(2 - 3.5) / 1.5, (5 - 3.5) / 1.5]$$

$$= [-1.0, 1.0]$$

Column 3 (features across samples): $[3, 6]$

Mean = 4.5, Variance = 2.25, std = 1.5

$$\text{Normalized} = [(3-4.5)/1.5, (6-4.5)/1.5] \\ = [-1.0, 1.0]$$

BN Result :-

$$\begin{bmatrix} -1 & -1 & -1 \\ 1 & 1 & 1 \end{bmatrix}$$

Layer Normalisation :-

⇒ Where it normalized: Across the features (hidden dimensions) of a single Example.

⇒ How it works:

for each sample, compute mean and Variance across all hidden units in that layer.

$$\mu_i = \frac{1}{H} \sum_{j=1}^H x_{i,j}, \quad \sigma_i^2 = \frac{1}{H} \sum_{j=1}^H (x_{i,j} - \mu_i)^2$$

Then normalize:

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

Intuition: Each token's hidden Vector is normalized independently, across its features.

Best for: MLP models (Transformers, BERT, GPT) where batch sizes can vary and sequence length matters more than batch statistics.

Example 1

$$X = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

Normalize across features (columns) for each sample (row).

Row 1 (features = $[1, 2, 3]$):

$$\text{Mean} = (1 + 2 + 3) / 3 = 2$$

$$\begin{aligned} \text{Variance} &= ((1-2)^2 + (2-2)^2 + (3-2)^2) / 3 \\ &= (1 + 0 + 1) / 3 = 0.67 \end{aligned}$$

$$\text{Std} \approx 0.82$$

$$\begin{aligned} \text{Normalized} &= [(1-2)/0.82, (2-2)/0.82, (3-2)/0.82] \\ &\approx [-1.22, 0, 1.22] \end{aligned}$$

Row 2 (features = $[4, 5, 6]$):

$$\text{Mean} = (4 + 5 + 6) / 3 = 5$$

$$\begin{aligned} \text{Variance} &= ((4-5)^2 + (5-5)^2 + (6-5)^2) / 3 \\ &= (1 + 0 + 1) / 3 = 0.67 \end{aligned}$$

$$\text{Std} = 0.82$$

$$\begin{aligned} \text{Normalized} &= [(4-5)/0.82, (5-5)/0.82, \\ &\quad (6-5)/0.82] \approx [-1.22, 0, 1.22] \end{aligned}$$

LN Result:

$$\begin{bmatrix} -1.22 & 0 & 1.22 \\ -1.22 & 0 & 1.22 \end{bmatrix}$$

Why Transformers uses Layer Normalization (and not Batch Normalization)

Batch Normalization :

- ⇒ looks at all samples in a batch and normalizes them together.
- ⇒ Work great in images (large batches, pixels are independent)
- ⇒ Problem in NLP / Transformers → batch size can be small and sentence lengths vary → unstable

Layer Normalization :

- ⇒ Looks at one token (word) at a time and normalizes across its features.
- ⇒ Doesn't care about batch size.
- ⇒ Perfect for text → keeps each word's hidden state stable.

Simple Analogy :

Batch Norm : "class average" → depends on how many students are in the class.

Layer Norm : "your own score normalized against your subjects" → doesn't care about class size, always consistent.

In short : Transformers use Layer Norm because :

1. Works even with small batches.
2. Works with variable-length sentences.
3. Same behavior in training & testing.

Feed-Forward Network (FFN) in Transformers:

- # The FFN is a small neural network inside each Transformer block. It processes each token's vector individually (no interaction with other tokens).

FFN in Deep Learning (DL) vs FFN in Transformers:

Aspect	FFN in DL	FFN in Transformers
Purpose	Maps input \rightarrow output (classification, regression, etc)	Refines each token's representation after attention.
Input	Whole dataset features (Eg: image pixels, tabular data, word embeddings)	One token vector at a time (Eg: 512-d vector)
Application	Applied once to entire input	Applied independently to every token (same weight shared)
Structure	Fully connected layers + activation (Eg: ReLU, GELU)	Same: 2 linear layers + activation (usually ReLU / GELU)
Hidden size	Depends on task (user-defined)	usually 4x layer than model dimension (Eg: $512 \rightarrow 2048 \rightarrow 512$)
Role	End-to-end learner for tasks	Adds non-linearity + expressive power inside Transformer block
Analogy	Whole system brain (thinking & deciding)	Individual token's private thought after group discussion (attention)

Masked Multi-head Attention :-

⇒ used in decoder's first attention layer.

Problem it solves :

While generating a sequence (Eg: a translation), the model must not peek at future words.

Example : if we're predicting the 3rd word, the model should only look at words 1 and 2, not 4 and 5, ...

How it works :

Normal self-attention lets each token attend to all tokens.

Masking = block future tokens by setting their attention scores to $-\infty$.

After Softmax → probabilities for those future tokens become 0.

Math :

Self-attention score :

$$\text{Score}(Q, K) = \frac{QK^T}{\sqrt{d_k}}$$

Masked Version :

$$\text{Score}(Q, K) = \begin{cases} \frac{QK^T}{\sqrt{d_k}}, & j \leq i \\ -\infty, & j > i \end{cases}$$

Where,

i = current token position

j = future position

Why we must not peek at future words?

1. Casuality in language:

- # If it "peeked ahead", it would be cheating - it wouldn't be true generation.

2. During Inference (generation):

- # At test time, the model doesn't know the future (because it hasn't generated it yet).

- # So training must mimic this condition - that's why we ~~enforce~~ enforce masking during training.

3. Analogy:

- # Imagine predicting the next word in a sentence game.

- # If you secretly saw the whole sentence beforehand, you'd ace it - but you wouldn't have learned how to predict.

- # Masking ensures the model learns to predict step by step.

Cross-Attention:

- # Used in decoder after masked self-attention.

Problem it solves:

- # Decoder needs to use encoder's knowledge (the input sentence) while generating output.

How it works:

Query (Q) = decoder hidden states (target tokens)
Key (K) & Values (V) = encoder outputs (input context)

This allows each decoder token to look at the entire encoded input.

Math:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

But here;

Q = decoder states

K, V = encoder outputs

This way, decoder words "align" with relevant source words.

Transformer Architecture: Encoder & Decoder

Transformers have two main parts:

Encoder \rightarrow understands the input

Decoder \rightarrow Generates the output

This Design is mainly used in sequence to sequence tasks (like translation).

Encoder:

Input: Sequence of tokens (Eg: Sentence in English)

Steps inside Each Encoder block:

1. Embedding + Positional Encoding

2. Multi-Head Attention
3. Residual + Layer Norm
4. Feed-Forward Network (FFN)
5. Residual + Layer Norm (again)

This block repeated N times (Eg: 6 to 12 layers)

Output: Context-rich token Embeddings.

Decoder:

Input: previously generated tokens (Eg: partial translated sentence)

Goal: Generate the next token.

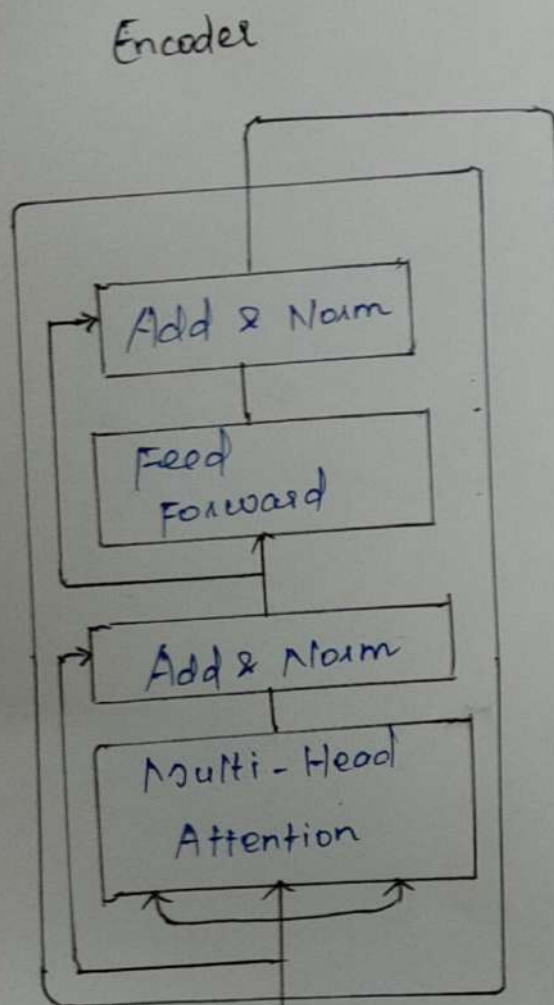
Steps inside each decoder block:

1. Embedding + Positional Encoding.
2. Masked Multi-Head Self-Attention
2. Residual + Layer Norm
4. Cross Attention (Key Difference from Encoder)
5. Residual + Layer Norm
6. Feed-Forward Network (FFN)
7. Residual + Layer Norm.

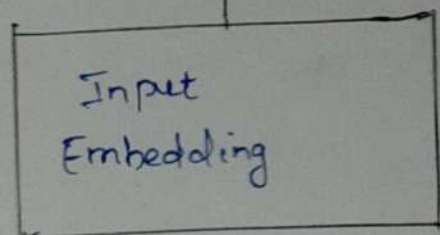
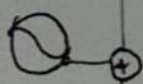
Repeated N times.

Final Layer: Linear + Softmax \rightarrow

predicts probability of next word.



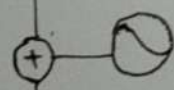
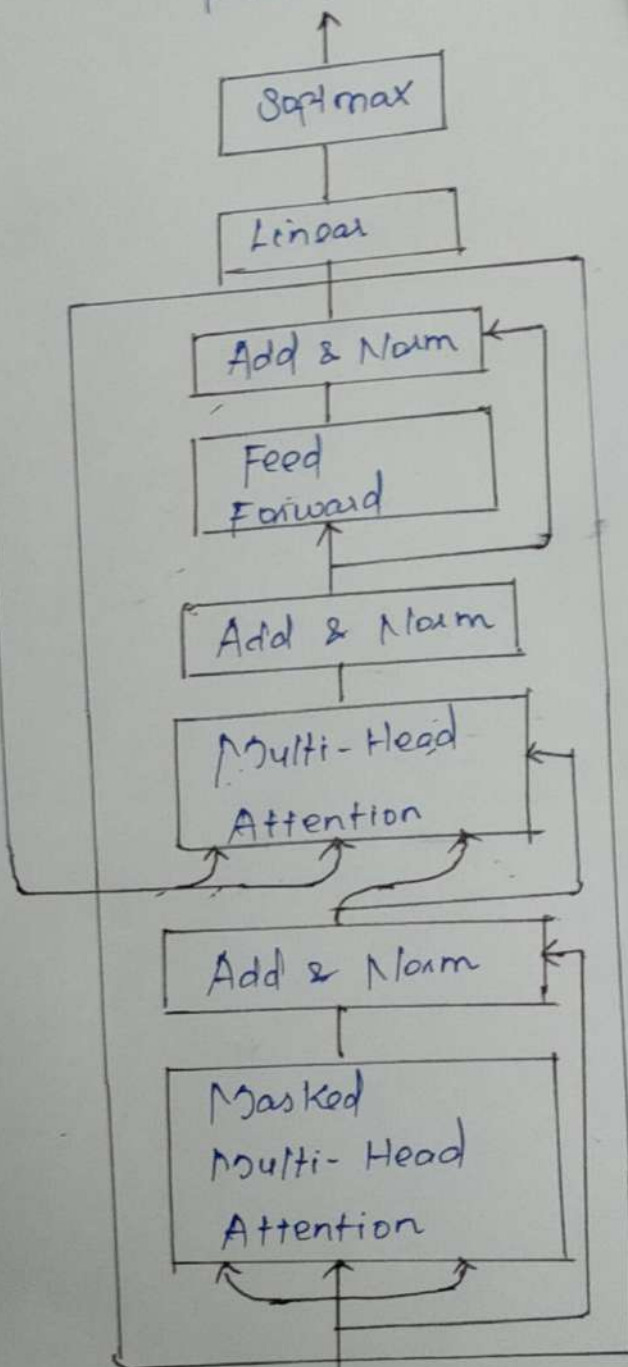
Positional
Encoding



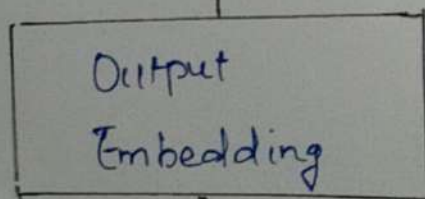
inputs

Decoder

Output
probabilities



Positional
Encoding



outputs (shifted right)