

Transformer Architecture Documentation

A Comprehensive Guide to the Implementation in `Transformer_from_scratch.py`

Based on “Attention Is All You Need”

February 8, 2026

Abstract

This document provides a detailed explanation of the Transformer architecture implementation found in `Transformer_from_scratch.py`. The implementation follows the original “Attention Is All You Need” paper with an innovative addition: *Multi-Head Latent Attention* which compresses key-value pairs for improved efficiency. Each component is explained mathematically and conceptually, with code snippets and dimensional analysis.

Contents

1	Introduction	2
1.1	Key Innovations in This Implementation	2
2	Core Components	2
2.1	Input Embeddings	2
2.2	Positional Encoding	2
2.3	Layer Normalization	3
2.4	Feed-Forward Network	4
2.5	Residual Connections	4
3	Attention Mechanisms	5
3.1	Standard Multi-Head Attention	5
3.2	Multi-Head Latent Attention (Innovation)	5
4	Encoder and Decoder	7
4.1	Encoder Block	7
4.2	Encoder Stack	7
4.3	Decoder Block	7
4.4	Decoder Stack	8
5	Complete Transformer	9
5.1	Projection Layer	9
5.2	Transformer Class	9
5.3	Model Construction	10
6	Data Flow Summary	12
6.1	Training Forward Pass	12
6.2	Inference (Greedy Decoding)	12
7	Key Differences from Original Paper	12
8	Conclusion	12

1 Introduction

The Transformer architecture revolutionized natural language processing by replacing recurrent layers with *attention mechanisms*. Unlike RNNs which process sequences sequentially, Transformers process entire sequences in parallel through self-attention.

1.1 Key Innovations in This Implementation

This implementation includes two attention mechanisms:

1. **Standard Multi-Head Attention** – The classic attention from the original paper
2. **Multi-Head Latent Attention (MLA)** – A memory-efficient variant that compresses keys and values

The `build_transformer()` function uses MLA for all encoder and decoder blocks.

2 Core Components

2.1 Input Embeddings

```
1 class InputEmbeddings(nn.Module):  
2     def __init__(self, d_model: int, vocab_size: int) -> None:  
3         super().__init__()  
4         self.d_model = d_model  
5         self.vocab_size = vocab_size  
6         self.embedding = nn.Embedding(vocab_size, d_model)  
7  
8     def forward(self, x):  
9         return self.embedding(x) * math.sqrt(self.d_model)
```

Listing 1: InputEmbeddings class

Purpose: Convert integer token IDs into dense vector representations.

Mathematical Operation:

$$\text{Embedding}(x) = \text{lookup}(x) \times \sqrt{d_{\text{model}}} \quad (1)$$

Dimensions:

- Input: `(batch, seq_len)` – integer token indices
- Output: `(batch, seq_len, d_model)` – embedding vectors

Why multiply by $\sqrt{d_{\text{model}}}$? The original paper scales embeddings to match the magnitude of positional encodings, preventing the positional information from drowning out semantic information.

2.2 Positional Encoding

```
1 class PositionalEncoding(nn.Module):  
2     def __init__(self, d_model: int, seq_len: int, dropout: float) -> None:  
3         super().__init__()  
4         pe = torch.zeros(seq_len, d_model)  
5         position = torch.arange(0, seq_len).unsqueeze(1)  
6         div_term = torch.exp(torch.arange(0, d_model, 2).float()  
7                             * (-math.log(10000.0) / d_model))  
8         pe[:, 0::2] = torch.sin(position * div_term)
```

```

9     pe[:, 1::2] = torch.cos(position * div_term)
10    self.register_buffer('pe', pe.unsqueeze(0))
11
12    def forward(self, x):
13        x = x + (self.pe[:, :x.shape[1], :]).requires_grad_(False)
14        return self.dropout(x)

```

Listing 2: PositionalEncoding class

Purpose: Inject information about token position since attention has no inherent sense of sequence order.

Mathematical Formula:

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (2)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (3)$$

Key Properties:

- Unique encoding for each position
- Relative positions can be derived via linear transformations
- Values bounded between $[-1, 1]$
- `requires_grad_=False` – positions are fixed, not learned

2.3 Layer Normalization

```

1 class LayerNormalization(nn.Module):
2     def __init__(self, features: int, eps: float = 10**-6) -> None:
3         super().__init__()
4         self.eps = eps
5         self.alpha = nn.Parameter(torch.ones(features))
6         self.bias = nn.Parameter(torch.zeros(features))
7
8     def forward(self, x):
9         mean = x.mean(dim=-1, keepdim=True)
10        std = x.std(dim=-1, keepdim=True)
11        return self.alpha * (x - mean) / (std + self.eps) + self.bias

```

Listing 3: LayerNormalization class

Purpose: Stabilize training by normalizing activations within each layer.

Mathematical Operation:

$$\text{LayerNorm}(x) = \alpha \odot \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta \quad (4)$$

Where:

- μ = mean across last dimension
- σ = standard deviation across last dimension
- α = learnable scale parameter
- β = learnable shift parameter
- ϵ = small constant for numerical stability

Note: This implementation uses *pre-normalization* (normalization before sublayers) via `ResidualConnection`.

2.4 Feed-Forward Network

```
1 class FeedForwardBlock(nn.Module):
2     def __init__(self, d_model: int, d_ff: int, dropout: float) -> None:
3         super().__init__()
4         self.linear_1 = nn.Linear(d_model, d_ff)
5         self.dropout = nn.Dropout(dropout)
6         self.linear_2 = nn.Linear(d_ff, d_model)
7
8     def forward(self, x):
9         return self.linear_2(self.dropout(torch.relu(self.linear_1(x))))
```

Listing 4: FeedForwardBlock class

Purpose: Apply non-linear transformations to each position independently.

Architecture:

$$\text{FFN}(x) = W_2 \cdot \text{Dropout}(\text{ReLU}(W_1 \cdot x + b_1)) + b_2 \quad (5)$$

Dimensions:

- W_1 : $(d_{\text{model}} \times d_{ff})$ – expands dimension
- W_2 : $(d_{ff} \times d_{\text{model}})$ – projects back
- Typical: $d_{ff} = 4 \times d_{\text{model}}$ (e.g., 2048 for $d_{\text{model}} = 512$)

2.5 Residual Connections

```
1 class ResidualConnection(nn.Module):
2     def __init__(self, features: int, dropout: float) -> None:
3         super().__init__()
4         self.dropout = nn.Dropout(dropout)
5         self.norm = LayerNormalization(features)
6
7     def forward(self, x, sublayer):
8         return x + self.dropout(sublayer(self.norm(x)))
```

Listing 5: ResidualConnection class

Purpose: Enable gradient flow through deep networks and stabilize training.

Operation:

$$\text{Output} = x + \text{Dropout}(\text{Sublayer}(\text{LayerNorm}(x))) \quad (6)$$

This is **pre-normalization** (norm before sublayer), which has been shown to train more stably than post-normalization for deep Transformers.

3 Attention Mechanisms

3.1 Standard Multi-Head Attention

```

1 class MultiHeadAttentionBlock(nn.Module):
2     def __init__(self, d_model: int, h: int, dropout: float) -> None:
3         super().__init__()
4         self.d_model = d_model
5         self.h = h
6         assert d_model % h == 0, "d_model is not divisible by h"
7         self.d_k = d_model // h
8
9         self.w_q = nn.Linear(d_model, d_model, bias=False)
10        self.w_k = nn.Linear(d_model, d_model, bias=False)
11        self.w_v = nn.Linear(d_model, d_model, bias=False)
12        self.w_o = nn.Linear(d_model, d_model, bias=False)
13        self.dropout = nn.Dropout(dropout)

```

Listing 6: MultiHeadAttentionBlock class

Core Idea: Split the model dimension into h heads, allowing the model to attend to information from different representation subspaces.

Scaled Dot-Product Attention:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (7)$$

Multi-Head Attention:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (8)$$

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (9)$$

Dimensions:

Tensor	Shape Before	Shape After
Query	(batch, seq, d_{model})	(batch, h , seq, d_k)
Key	(batch, seq, d_{model})	(batch, h , seq, d_k)
Value	(batch, seq, d_{model})	(batch, h , seq, d_k)
Attention Scores	—	(batch, h , seq, seq)
Output	(batch, h , seq, d_k)	(batch, seq, d_{model})

Why $\sqrt{d_k}$? Prevents dot products from growing too large, which would push softmax into regions with extremely small gradients.

3.2 Multi-Head Latent Attention (Innovation)

```

1 class MultiHeadLatentAttentionBlock(nn.Module):
2     def __init__(self, d_model: int, h: int, dropout: float):
3         super().__init__()
4         self.d_model = d_model
5         self.h = h
6         assert d_model % h == 0, "d_model is not divisible by h"
7         self.d_k = d_model // h
8
9         assert d_model % 4 == 0, "d_model is not divisible by 4"
10        self.latent_dim = d_model // 4 # Compression factor of 4

```

```

12     self.w_q = nn.Linear(d_model, d_model, bias=False)
13     # KEY INNOVATION: Compression + Decompression
14     self.w_kv_compress = nn.Linear(d_model, self.latent_dim, bias=False)
15     self.w_k_decompress = nn.Linear(self.latent_dim, d_model, bias=False)
16     self.w_v_decompress = nn.Linear(self.latent_dim, d_model, bias=False)
17     self.w_o = nn.Linear(d_model, d_model, bias=False)
18     self.dropout = nn.Dropout(dropout)

```

Listing 7: MultiHeadLatentAttentionBlock class

Key Innovation: Instead of computing keys and values at full dimension, MLA:

1. **Compresses** input to a latent representation: $c = W_{\text{kv_compress}} \cdot x$
2. **Decompresses** to get keys and values: $K = W_{\text{k_decompress}} \cdot c$, $V = W_{\text{v_decompress}} \cdot c$

Memory Savings:

- Standard: Store $K, V \in \mathbb{R}^{d_{\text{model}}}$
- MLA: Store latent $c \in \mathbb{R}^{d_{\text{model}}/4}$
- **Reduction:** 75% memory savings for KV cache!

Mathematical Flow:

$$c = W_{\text{compress}} \cdot x \in \mathbb{R}^{d_{\text{model}}/4} \quad (10)$$

$$K = W_{\text{k_decompress}} \cdot c \in \mathbb{R}^{d_{\text{model}}} \quad (11)$$

$$V = W_{\text{v_decompress}} \cdot c \in \mathbb{R}^{d_{\text{model}}} \quad (12)$$

$$Q = W_q \cdot x \in \mathbb{R}^{d_{\text{model}}} \quad (13)$$

Trade-offs:

- ✓ Reduced memory footprint
- ✓ Faster inference (smaller KV cache)
- ✗ Slight capacity reduction (information bottleneck)
- ✗ Additional computation for compression/decompression

4 Encoder and Decoder

4.1 Encoder Block

```
1 class EncoderBlock(nn.Module):
2     def __init__(self, features: int,
3                  self_attention_block: MultiHeadLatentAttentionBlock,
4                  feed_forward_block: FeedForwardBlock, dropout: float) -> None:
5         super().__init__()
6         self.self_attention_block = self_attention_block
7         self.feed_forward_block = feed_forward_block
8         self.residual_connections = nn.ModuleList([
9             ResidualConnection(features, dropout) for _ in range(2)
10        ])
11
12     def forward(self, x, src_mask):
13         x = self.residual_connections[0](x,
14             lambda x: self.self_attention_block(x, x, x, src_mask))
15         x = self.residual_connections[1](x, self.feed_forward_block)
16
17         return x
```

Listing 8: EncoderBlock class

Architecture:

1. Self-Attention with residual: $x = x + \text{Attn}(\text{LN}(x))$
2. Feed-Forward with residual: $x = x + \text{FFN}(\text{LN}(x))$

Self-Attention: Each position attends to all positions in the same sequence (uses `src_mask` for padding).

4.2 Encoder Stack

```
1 class Encoder(nn.Module):
2     def __init__(self, features: int, layers: nn.ModuleList) -> None:
3         super().__init__()
4         self.layers = layers
5         self.norm = LayerNormalization(features)
6
7     def forward(self, x, mask):
8         for layer in self.layers:
9             x = layer(x, mask)
10
11    return self.norm(x)
```

Listing 9: Encoder class

Stacks N (default 6) `EncoderBlocks` with a final layer normalization.

4.3 Decoder Block

```
1 class DecoderBlock(nn.Module):
2     def __init__(self, features: int,
3                  self_attention_block: MultiHeadLatentAttentionBlock,
4                  cross_attention_block: MultiHeadLatentAttentionBlock,
5                  feed_forward_block: FeedForwardBlock, dropout: float) -> None:
6         super().__init__()
7         self.self_attention_block = self_attention_block
8         self.cross_attention_block = cross_attention_block
9         self.feed_forward_block = feed_forward_block
10        self.residual_connections = nn.ModuleList([
```

```

11         ResidualConnection(features, dropout) for _ in range(3)
12     ])
13
14     def forward(self, x, encoder_output, src_mask, tgt_mask):
15         x = self.residual_connections[0](x,
16             lambda x: self.self_attention_block(x, x, x, tgt_mask))
17         x = self.residual_connections[1](x,
18             lambda x: self.cross_attention_block(x, encoder_output,
19                                                 encoder_output, src_mask))
20         x = self.residual_connections[2](x, self.feed_forward_block)
21     return x

```

Listing 10: DecoderBlock class

Architecture (3 sublayers):

1. **Masked Self-Attention:** Each position attends to previous positions only (causal masking)
2. **Cross-Attention:** Queries from decoder attend to Keys/Values from encoder output
3. **Feed-Forward:** Position-wise transformation

Cross-Attention:

$$\text{CrossAttn}(Q_{\text{decoder}}, K_{\text{encoder}}, V_{\text{encoder}}) \quad (14)$$

This is where the decoder “looks at” the source sequence representation.

4.4 Decoder Stack

```

1 class Decoder(nn.Module):
2     def __init__(self, features: int, layers: nn.ModuleList) -> None:
3         super().__init__()
4         self.layers = layers
5         self.norm = LayerNormalization(features)
6
7     def forward(self, x, encoder_output, src_mask, tgt_mask):
8         for layer in self.layers:
9             x = layer(x, encoder_output, src_mask, tgt_mask)
10    return self.norm(x)

```

Listing 11: Decoder class

Stacks N DecoderBlocks with a final layer normalization.

5 Complete Transformer

5.1 Projection Layer

```
1 class ProjectionLayer(nn.Module):
2     def __init__(self, d_model, vocab_size) -> None:
3         super().__init__()
4         self.proj = nn.Linear(d_model, vocab_size)
5
6     def forward(self, x) -> None:
7         return self.proj(x)
```

Listing 12: ProjectionLayer class

Converts decoder output from d_{model} dimensions to vocabulary size for softmax prediction.

Output: (batch, seq_len, vocab_size) – logits for each token.

5.2 Transformer Class

```
1 class Transformer(nn.Module):
2     def __init__(self, encoder: Encoder, decoder: Decoder,
3                  src_embed: InputEmbeddings, tgt_embed: InputEmbeddings,
4                  src_pos: PositionalEncoding, tgt_pos: PositionalEncoding,
5                  projection_layer: ProjectionLayer) -> None:
6         super().__init__()
7         self.encoder = encoder
8         self.decoder = decoder
9         self.src_embed = src_embed
10        self.tgt_embed = tgt_embed
11        self.src_pos = src_pos
12        self.tgt_pos = tgt_pos
13        self.projection_layer = projection_layer
14
15    def encode(self, src, src_mask):
16        src = self.src_embed(src)
17        src = self.src_pos(src)
18        return self.encoder(src, src_mask)
19
20    def decode(self, encoder_output, src_mask, tgt, tgt_mask):
21        tgt = self.tgt_embed(tgt)
22        tgt = self.tgt_pos(tgt)
23        return self.decoder(tgt, encoder_output, src_mask, tgt_mask)
24
25    def project(self, x):
26        return self.projection_layer(x)
```

Listing 13: Transformer class

Three Main Operations:

1. `encode(src, src_mask):`

- Embed + positional encode source
- Pass through encoder stack
- Returns: encoder representation

2. `decode(encoder_output, src_mask, tgt, tgt_mask):`

- Embed + positional encode target
- Pass through decoder stack (with cross-attention)

- Returns: decoder representation

3. project(x):

- Linear projection to vocabulary
- Returns: logits for next-token prediction

5.3 Model Construction

```

1 def build_transformer(src_vocab_size: int, tgt_vocab_size: int,
2                     src_seq_len: int, tgt_seq_len: int,
3                     d_model: int = 512, N: int = 6, h: int = 8,
4                     dropout: float = 0.1, d_ff: int = 2048):
5     # Embeddings
6     src_embed = InputEmbeddings(d_model, src_vocab_size)
7     tgt_embed = InputEmbeddings(d_model, tgt_vocab_size)
8
9     # Positional encodings
10    src_pos = PositionalEncoding(d_model, src_seq_len, dropout)
11    tgt_pos = PositionalEncoding(d_model, tgt_seq_len, dropout)
12
13    # Encoder blocks with MultiHeadLatentAttentionBlock
14    encoder_blocks = []
15    for _ in range(N):
16        encoder_self_attention = MultiHeadLatentAttentionBlock(d_model, h,
17 dropout)
18        feed_forward = FeedForwardBlock(d_model, d_ff, dropout)
19        encoder_block = EncoderBlock(d_model, encoder_self_attention,
20                                      feed_forward, dropout)
21        encoder_blocks.append(encoder_block)
22
23    # Decoder blocks with MultiHeadLatentAttentionBlock
24    decoder_blocks = []
25    for _ in range(N):
26        decoder_self_attention = MultiHeadLatentAttentionBlock(d_model, h,
27 dropout)
28        decoder_cross_attention = MultiHeadLatentAttentionBlock(d_model, h,
29 dropout)
30        feed_forward = FeedForwardBlock(d_model, d_ff, dropout)
31        decoder_block = DecoderBlock(d_model, decoder_self_attention,
32                                     decoder_cross_attention, feed_forward,
33                                     dropout)
34        decoder_blocks.append(decoder_block)
35
36    # Assemble
37    encoder = Encoder(d_model, nn.ModuleList(encoder_blocks))
38    decoder = Decoder(d_model, nn.ModuleList(decoder_blocks))
39    projection_layer = ProjectionLayer(d_model, tgt_vocab_size)
40
41    transformer = Transformer(encoder, decoder, src_embed,
42                             tgt_embed,
43                             src_pos, tgt_pos, projection_layer)
44
45    # Xavier initialization
46    for p in transformer.parameters():
47        if p.dim() > 1:
48            nn.init.xavier_uniform_(p)
49
50    return transformer

```

Listing 14: build_transformer function

Default Hyperparameters:

Parameter	Value	Description
d_{model}	512	Model dimension
N	6	Number of layers
h	8	Number of attention heads
d_{ff}	2048	Feed-forward hidden dimension
dropout	0.1	Dropout rate

6 Data Flow Summary

6.1 Training Forward Pass

1. Source Processing:

$$\text{src} \rightarrow \text{Embedding} \rightarrow \text{Positional Encoding} \quad (15)$$

$$\rightarrow \text{Encoder} \rightarrow \text{encoder_output} \quad (16)$$

2. Target Processing:

$$\text{tgt} \rightarrow \text{Embedding} \rightarrow \text{Positional Encoding} \quad (17)$$

$$\rightarrow \text{Decoder (with cross-attention)} \quad (18)$$

$$\rightarrow \text{decoder_output} \quad (19)$$

3. Prediction:

$$\text{decoder_output} \rightarrow \text{Projection} \rightarrow \text{Softmax} \rightarrow \text{Loss} \quad (20)$$

6.2 Inference (Greedy Decoding)

1. Encode source sequence once: $\text{memory} = \text{encode}(\text{src})$
2. Start with [SOS] token
3. For each position i :
 - Decode: $\text{output} = \text{decode}(\text{memory}, \text{generated_tokens})$
 - Project: $\text{logits} = \text{project}(\text{output})$
 - Take argmax: $\text{token}_i = \arg \max(\text{logits})$
 - Append to generated sequence
4. Stop when [EOS] token generated

7 Key Differences from Original Paper

1. **Multi-Head Latent Attention:** This implementation uses KV compression, reducing memory by 75%.
2. **Pre-Normalization:** Layer norm happens before attention/FFN (not after).
3. **Xavier Initialization:** Uses `xavier_uniform_` instead of the paper's scheme.

8 Conclusion

This implementation provides a clean, modular Transformer architecture with the innovative addition of Multi-Head Latent Attention. The code is well-structured with clear separation of concerns:

- **Embeddings:** Handle input representation
- **Attention:** Captures relationships (with memory-efficient MLA)
- **Encoder/Decoder:** Process sequences bidirectionally/unidirectionally
- **Residual Connections:** Enable deep network training

The use of `MultiHeadLatentAttentionBlock` throughout makes this implementation particularly suitable for resource-constrained environments where KV cache memory is a bottleneck.