# Understanding the Internal Mechanics of LLMs

### Muskula Rahul

## Understanding the Internal Mechanics of LLMs

Large Language Models (LLMs) have transformed natural language processing and artificial intelligence over the past few years. While their impressive capabilities are readily apparent, the underlying mechanics that drive these systems remain mysterious to many. This article delves into the technical intricacies that power today's cutting-edge LLMs.

## 1 The Transformer Architecture: Foundation of Modern LLMs

Modern LLMs are built upon the Transformer architecture, first introduced in the seminal paper "Attention Is All You Need" (Vaswani et al., 2017). Unlike previous recurrent neural network approaches, Transformers process entire sequences simultaneously through self-attention mechanisms.

At its core, the Transformer relies on the multi-head attention mechanism:

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

Where $Q$ (queries), $K$ (keys), and $V$ (values) are linear projections of the input embeddings, and $d_k$ is the dimension of the keys, serving as a scaling factor to stabilize gradients during training.

Multi-head attention extends this further by computing attention multiple times in parallel:

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

Where each head is:

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

## 2 The Scaling Laws: Why Size Matters

One of the most significant discoveries in LLM research is the emergence of scaling laws. These empirical relationships show how model performance improves predictably with increases in model size, dataset size, and compute resources.

Kaplan et al. (2020) identified that loss $L$ scales with model parameters $N$ approximately as:

$$L(N) \approx \left(\frac{N_c}{N}\right)^{\alpha_N}$$

Where $N_c$ is a constant and $\alpha_N \approx 0.076$, indicating that each doubling of model size reduces loss by approximately 5%.

Similarly, loss scales with dataset size $D$ as:

$$L(D) \approx \left(\frac{D_c}{D}\right)^{\alpha_D}$$

Where $\alpha_D \approx 0.095$, suggesting that dataset size is slightly more important than model size.

# 3 Training Dynamics: Navigating Loss Landscapes

Training LLMs involves navigating an extraordinarily high-dimensional loss landscape. The objective function in language modeling is typically next-token prediction, formulated as:

$$\mathcal{L}(\theta) = -\frac{1}{|\mathcal{D}|} \sum_{(x,y)\in\mathcal{D}} \log p_\theta(y|x)$$

Where $\theta$ represents model parameters, $\mathcal{D}$ is the training dataset, and $p_\theta(y|x)$ is the probability assigned by the model to the correct next token $y$ given context $x$.

Most LLMs use variants of adaptive optimization algorithms, with AdamW being particularly popular:

$$m_t = \beta_1 m_{t-1} + (1-\beta_1)g_t$$

$$v_t = \beta_2 v_{t-1} + (1-\beta_2)g_t^2$$

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

$$\theta_t = \theta_{t-1} - \eta\frac{\hat{m}_t}{\sqrt{\hat{v}_t}+\epsilon} - \eta\lambda\theta_{t-1}$$

Where $g_t$ is the gradient, $m_t$ and $v_t$ are the first and second moment estimates, $\beta_1$ and $\beta_2$ are decay rates, $\eta$ is the learning rate, $\lambda$ is the weight decay coefficient, and $\epsilon$ is a small constant for numerical stability.

# 4 Tokenization: The Critical First Step

Before processing text, LLMs must convert it into numerical representations through tokenization. Modern LLMs typically use subword tokenization algorithms like Byte-Pair Encoding (BPE) or SentencePiece.

BPE starts with individual characters and iteratively merges the most frequent adjacent pairs to form a vocabulary of specified size. This process can be formalized as:

1. Initialize vocabulary $V$ with all unique characters in corpus.

2. Compute frequencies of all adjacent token pairs $(a, b)$ where $a, b \in V$.

3. Find the most frequent pair $(a, b)$ and add merged token $ab$ to $V$.

4. Replace all occurrences of $(a, b)$ with $ab$ in the corpus.

5. Repeat steps 2-4 until vocabulary reaches desired size or merge frequency falls below threshold.

# 5 Positional Encodings: Capturing Sequential Information

Unlike RNNs, Transformers have no inherent understanding of sequential order. To address this, position encodings are added to input embeddings. The original Transformer used sinusoidal encodings:

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

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

Where $pos$ is the position index, $i$ is the dimension index, and $d_{model}$ is the model dimension.

Later models adopted learned positional encodings, while more recent architectures like RoPE (Rotary Position Embedding) incorporate position information directly into attention calculations:

$$RoPE(q,k,m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} q_0 \\ q_1 \end{pmatrix} \cdot \begin{pmatrix} \cos n\theta & -\sin n\theta \\ \sin n\theta & \cos n\theta \end{pmatrix} \begin{pmatrix} k_0 \\ k_1 \end{pmatrix}^T$$

# 6 Computational Efficiency Innovations

Training and running LLMs requires enormous computational resources. Several innovations have made this more manageable:

## 6.1 Mixed Precision Training

Computing operations in lower precision (FP16 or BF16 instead of FP32) reduces memory requirements and computation time. This is formalized as:

$$\text{Loss}_{scaled} = \text{Loss} \times \text{scale\_factor}$$

$$\text{Gradients}_{FP16} = \text{Compute\_Gradients}(\text{Loss}_{scaled})$$

$$\text{Gradients}_{FP32} = \text{Gradients}_{FP16}/\text{scale\_factor}$$

## 6.2 Attention Optimizations

FlashAttention algorithms reduce memory bandwidth bottlenecks by recomputing attention on-the-fly rather than storing the full attention matrix:

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

Traditional implementation: $O(N^2)$ memory for the attention matrix $\hat{P}$.

FlashAttention: $O(N)$ memory by tiling the computation into smaller blocks that fit in faster memory (SRAM).

# 7 Inference Optimization Techniques

Large LLMs present significant inference challenges. Key techniques to address these include:

## 7.1 KV Caching

During autoregressive generation, previous key-value pairs from attention layers are cached to avoid recomputation:

$$\text{KV-Cache}_t = \{\text{KV-Cache}_{t-1}, (K_t, V_t)\}$$

Where $K_t$ and $V_t$ are the key-value pairs generated at time step $t$.

## 7.2 Quantization

Reducing the precision of model weights and activations after training:

$$W_q = \text{round}\left(\frac{W - \min(W)}{\max(W) - \min(W)} \times (2^b - 1)\right) \times \frac{\max(W) - \min(W)}{2^b - 1} + \min(W)$$

Where $W$ is the original weight tensor, $W_q$ is the quantized tensor, and $b$ is the bit width (e.g., 4 or 8 bits).

# 8 Understanding LLM Capabilities Through Scaling

Emergent abilities—capabilities not present in smaller models but appearing in larger ones—remain one of the most fascinating aspects of LLMs. These can be visualized through step function-like improvements:

$$P(correct) = \sigma\left(\alpha \log\left(\frac{N}{N_0}\right)\right)$$

Where $P(correct)$ is the probability of correctly performing a task, $N$ is the number of parameters, $N_0$ is a threshold parameter count, and $\sigma$ is the logistic function.

# 9 The Challenge of Context Windows

LLMs' ability to process long contexts has improved dramatically. Extending context windows introduces quadratic computational complexity in standard attention:

$$\text{Complexity}_{\text{attention}} = O(L^2 \cdot d)$$

Where $L$ is sequence length and $d$ is hidden dimension.

Techniques like sparse attention patterns reduce this to:

$$\text{Complexity}_{\text{sparse}} = O(L \cdot \log(L) \cdot d)$$

# 10 Conclusion

The internal mechanics of LLMs represent one of the most sophisticated achievements in AI engineering. From the elegant mathematics of attention to the empirical scaling laws that guide development, these systems combine theoretical insights with practical engineering solutions. As research continues, we can expect further refinements to these architectures that will expand their capabilities while addressing current limitations in reasoning, factuality, and computational efficiency.

Understanding these mechanics not only satisfies scientific curiosity but provides the foundation for developing more capable, efficient, and reliable AI systems in the future.