

# How Large Language Models (LLMs) Are Made: A Complete Guide

A comprehensive guide to building GPT-like models from scratch and understanding advanced innovations like DeepSeek

---

## Table of Contents

- [1. Introduction](#)
  - [2. Data Preprocessing Pipeline](#)
  - [3. Attention Mechanisms](#)
  - [4. Transformer Architecture](#)
  - [5. Complete GPT Architecture](#)
  - [6. Training Process](#)
  - [7. DeepSeek Innovations](#)
  - [8. Advanced Optimization Techniques](#)
  - [9. Mixture of Experts \(MoE\)](#)
  - [10. Mixture of Depths \(MoD\)](#)
  - [11. Conclusion](#)
- 

## Introduction

Large Language Models (LLMs) like GPT have revolutionized natural language processing. This guide provides a comprehensive overview of how these models are built from scratch, covering both the foundational GPT architecture and cutting-edge innovations from DeepSeek.

LLMs are transformer-based neural networks trained on vast amounts of text data to predict the next token in a sequence. They learn language patterns, context, and relationships between words to generate human-like text.

---

## Data Preprocessing Pipeline

The journey of building an LLM starts with preparing the data. Here's the step-by-step process:

### 1. Data Loading and Tokenization

```
import tiktoken

# Load raw text data
with open('data.txt', 'r') as f:
    raw_text = f.read()

# Get GPT-2 vocabulary (size ~50,257)
vocab = tiktoken.get_encoding("gpt2")
tokens = vocab.encode(raw_text)
```

### 2. Dataset Creation

Create a dataset class that segments the data into input-target pairs:

```
class LLMDataset:
    def __init__(self, tokens, max_length, stride):
        self.input_ids = []
        self.target_ids = []

        for i in range(0, len(tokens) - max_length, stride):
            input_chunk = tokens[i:i + max_length]
            target_chunk = tokens[i + 1:i + max_length + 1]
            self.input_ids.append(input_chunk)
            self.target_ids.append(target_chunk)
```

### 3. Embedding Layers

Token Embeddings:

- Acts as a lookup table initialized randomly
- Parameters: vocab\_size × embedding\_dim
- Embeddings are learned during training

### Positional Embeddings:

- Helps the model understand token positions
- Parameters: context\_length × embedding\_dim
- Added to token embeddings to create input embeddings

```
import torch.nn as nn
```

```
token_embedding = nn.Embedding(vocab_size, embedding_dim)
pos_embedding = nn.Embedding(context_length, embedding_dim)
```

```
# Final input embeddings
input_embeddings = token_embedding(input_ids) + pos_embedding(positions)
```

---

## Attention Mechanisms

Attention is the core mechanism that allows LLMs to understand relationships between tokens, especially in long sequences where RNNs fail.

### 1. Simplified Attention (Base Level)

The most basic form of attention:

```
# Calculate attention scores by dot product
attention_scores = torch.matmul(embeddings, embeddings.transpose(-2, -1))
```

```
# Apply softmax to normalize scores
attention_weights = torch.softmax(attention_scores, dim=-1)
```

```
# Multiply attention weights by embeddings
context_embeddings = torch.matmul(attention_weights, embeddings)
```

### 2. Self-Attention

Self-attention learns separate weight matrices for queries (Q), keys (K), and values (V):

```
class SelfAttention(nn.Module):
    def __init__(self, embed_dim):
        super().__init__()
        self.W_q = nn.Linear(embed_dim, embed_dim)
        self.W_k = nn.Linear(embed_dim, embed_dim)
        self.W_v = nn.Linear(embed_dim, embed_dim)

    def forward(self, x):
        Q = self.W_q(x)
        K = self.W_k(x)
        V = self.W_v(x)

        # Compute attention scores
        scores = torch.matmul(Q, K.transpose(-2, -1)) / math.sqrt(embed_dim)
        attention_weights = torch.softmax(scores, dim=-1)

        # Compute output
        output = torch.matmul(attention_weights, V)
        return output
```

### 3. Causal Attention

Prevents the model from "peeking" at future tokens during training:

```
def apply_causal_mask(attention_scores):
    seq_len = attention_scores.size(-1)
    mask = torch.triu(torch.ones(seq_len, seq_len), diagonal=1).bool()
    attention_scores.masked_fill_(mask, float('-inf'))
    return attention_scores

# In the attention computation
scores = torch.matmul(Q, K.transpose(-2, -1)) / math.sqrt(d_k)
scores = apply_causal_mask(scores)
attention_weights = torch.softmax(scores, dim=-1)
```

### 4. Multi-Head Attention

Runs multiple attention heads in parallel to capture different types of relationships:

```
class MultiHeadAttention(nn.Module):
    def __init__(self, embed_dim, num_heads):
        super().__init__()
        self.num_heads = num_heads
        self.head_dim = embed_dim // num_heads

        self.W_q = nn.Linear(embed_dim, embed_dim)
        self.W_k = nn.Linear(embed_dim, embed_dim)
        self.W_v = nn.Linear(embed_dim, embed_dim)
        self.W_o = nn.Linear(embed_dim, embed_dim)

    def forward(self, x):
        batch_size, seq_len, embed_dim = x.size()

        Q = self.W_q(x).view(batch_size, seq_len, self.num_heads, self.head_dim)
        K = self.W_k(x).view(batch_size, seq_len, self.num_heads, self.head_dim)
        V = self.W_v(x).view(batch_size, seq_len, self.num_heads, self.head_dim)

        # Transpose for attention computation
        Q = Q.transpose(1, 2) # (batch, num_heads, seq_len, head_dim)
        K = K.transpose(1, 2)
        V = V.transpose(1, 2)

        # Compute attention
        scores = torch.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.head_dim)
        attention_weights = torch.softmax(scores, dim=-1)
        output = torch.matmul(attention_weights, V)

        # Concatenate heads
        output = output.transpose(1, 2).contiguous().view(batch_size, seq_len, embed_dim)
        return self.W_o(output)
```

---

## Transformer Architecture

The transformer block is the fundamental building unit of LLMs, consisting of several key components:

### 1. Layer Normalization

Prevents exploding/vanishing gradients and stabilizes training:

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

Where:

- $\mu$  is the mean
- $\sigma^2$  is the variance
- $\gamma$  (scale) and  $\beta$  (shift) are trainable parameters

```
class LayerNorm(nn.Module):
    def __init__(self, embed_dim, eps=1e-5):
        super().__init__()
        self.eps = eps
        self.scale = nn.Parameter(torch.ones(embed_dim))
        self.shift = nn.Parameter(torch.zeros(embed_dim))

    def forward(self, x):
        mean = x.mean(dim=-1, keepdim=True)
        var = x.var(dim=-1, keepdim=True, unbiased=False)
        norm_x = (x - mean) / torch.sqrt(var + self.eps)
        return self.scale * norm_x + self.shift
```

### 2. Feed Forward Network (FFN)

A two-layer neural network with GeLU activation:

```
class FeedForward(nn.Module):
    def __init__(self, embed_dim):
        super().__init__()
        self.linear1 = nn.Linear(embed_dim, 4 * embed_dim) # Hidden size is 4x
        self.linear2 = nn.Linear(4 * embed_dim, embed_dim)
        self.gelu = nn.GELU()

    def forward(self, x):
        return self.linear2(self.gelu(self.linear1(x)))
```

**GeLU Activation:** Smoother than ReLU and better for language modeling:

```
def gelu(x):
    return 0.5 * x * (1 + torch.tanh(math.sqrt(2/math.pi) * (x + 0.044715 * x**3)))
```

### 3. Residual Connections

Add the input of each sub-layer to its output to help gradients flow:

```
class TransformerBlock(nn.Module):
    def __init__(self, embed_dim, num_heads):
        super().__init__()
        self.attention = MultiHeadAttention(embed_dim, num_heads)
        self.ffn = FeedForward(embed_dim)
        self.norm1 = LayerNorm(embed_dim)
        self.norm2 = LayerNorm(embed_dim)

    def forward(self, x):
        # Attention with residual connection
        attn_output = self.attention(self.norm1(x))
        x = x + attn_output

        # FFN with residual connection
        ffn_output = self.ffn(self.norm2(x))
        x = x + ffn_output

    return x
```

---

## Complete GPT Architecture

Combining all components into a complete GPT model:

```
class GPTModel(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.token_embedding = nn.Embedding(config.vocab_size, config.embed_dim)
        self.pos_embedding = nn.Embedding(config.context_length, config.embed_dim)

        self.transformer_blocks = nn.ModuleList([
            TransformerBlock(config.embed_dim, config.num_heads)
            for _ in range(config.num_layers)
        ])

        self.final_norm = LayerNorm(config.embed_dim)
        self.output_head = nn.Linear(config.embed_dim, config.vocab_size, bias=False)

    def forward(self, input_ids):
        batch_size, seq_len = input_ids.shape

        # Create embeddings
        token_embeds = self.token_embedding(input_ids)
        pos_embeds = self.pos_embedding(torch.arange(seq_len, device=input_ids.device))
        x = token_embeds + pos_embeds

        # Pass through transformer blocks
        for block in self.transformer_blocks:
            x = block(x)

        # Final normalization and output projection
        x = self.final_norm(x)
        logits = self.output_head(x)

    return logits
```

---

## Training Process

### 1. Loss Function

Cross-entropy loss measures how well predicted logits match target tokens:

```
def compute_loss(logits, targets):
    # Flatten for cross-entropy computation
    logits_flat = logits.view(-1, logits.size(-1))
    targets_flat = targets.view(-1)

    loss = torch.nn.functional.cross_entropy(logits_flat, targets_flat)
    return loss
```

### 2. Training Loop

```

def train_model(model, train_loader, val_loader, num_epochs):
    optimizer = torch.optim.AdamW(model.parameters(), lr=1e-4)

    for epoch in range(num_epochs):
        model.train()
        total_loss = 0

        for batch in train_loader:
            input_ids, targets = batch

            # Forward pass
            logits = model(input_ids)
            loss = compute_loss(logits, targets)

            # Backward pass
            optimizer.zero_grad()
            loss.backward()
            optimizer.step()

            total_loss += loss.item()

        # Validation
        val_loss = evaluate_model(model, val_loader)
        print(f"Epoch {epoch}: Train Loss = {total_loss/len(train_loader):.4f}, Val Loss = {val_loss:.4f}")

```

### 3. Text Generation

```

def generate_text(model, input_ids, max_new_tokens, temperature=1.0, top_k=50):
    model.eval()

    for _ in range(max_new_tokens):
        with torch.no_grad():
            logits = model(input_ids)
            next_token_logits = logits[0, -1, :] / temperature

            # Apply top-k sampling
            if top_k > 0:
                top_k_logits, top_k_indices = torch.topk(next_token_logits, top_k)
                next_token_logits = torch.full_like(next_token_logits, float('-inf'))
                next_token_logits[top_k_indices] = top_k_logits

            probs = torch.softmax(next_token_logits, dim=-1)
            next_token = torch.multinomial(probs, 1)

            input_ids = torch.cat([input_ids, next_token.unsqueeze(0)], dim=1)

    return input_ids

```

---

## DeepSeek Innovations

DeepSeek introduced several groundbreaking innovations that significantly improve LLM efficiency and performance.

### 1. Key-Value (KV) Caching Problem

Traditional attention has a major inefficiency in inference:

- a. AB -> Inference -> C
- b. ABC -> Inference -> D
- c. ABCD -> Inference -> E

The problem: We recalculate K, V matrices for tokens AB multiple times, wasting computation.

#### Memory requirement for KV cache:

Size =  $l \times b \times n \times h \times s \times 2 \times 2$   
 Where:

- $l$  = number of transformer blocks
- $b$  = batch size
- $n$  = number of heads
- $h$  = head dimension
- $s$  = context length
- 2 = one for Key, one for Value
- 2 = 2 bytes (float16)

This leads to 400GB+ memory requirements for large models!

### 2. Multi-Query Attention (MQA)

**Solution:** Make all attention heads share the same K and V matrices.

- **Advantage:** Reduces KV cache from 400GB to ~3GB
- **Disadvantage:** Loss in performance due to reduced perspectives

```
class MultiQueryAttention(nn.Module):
    def __init__(self, embed_dim, num_heads):
        super().__init__()
        self.num_heads = num_heads
        self.head_dim = embed_dim // num_heads

        self.W_q = nn.Linear(embed_dim, embed_dim) # Multiple queries
        self.W_k = nn.Linear(embed_dim, self.head_dim) # Single key
        self.W_v = nn.Linear(embed_dim, self.head_dim) # Single value
```

### 3. Grouped Query Attention (GQA)

**Balance between MHA and MQA:** Create groups of heads that share K,V matrices.

Example: 64 heads → 8 groups of 8 heads each

[1,1,1,1], [2,2,2,2], [3,3,3,3], [4,4,4,4]

Memory formula:  $l \times b \times g \times h \times s \times 2 \times 2$  where  $g$  = number of groups

### 4. Multi-Head Latent Attention (MLA)

DeepSeek's breakthrough solution that achieves both low memory and high performance:

```
class MultiHeadLatentAttention(nn.Module):
    def __init__(self, embed_dim, num_heads, latent_dim):
        super().__init__()
        self.W_q = nn.Linear(embed_dim, embed_dim)
        self.W_dkv = nn.Linear(embed_dim, latent_dim) # Down projection
        self.W_uk = nn.Linear(latent_dim, embed_dim) # Up projection for K
        self.W_uv = nn.Linear(latent_dim, embed_dim) # Up projection for V

    def forward(self, x):
        Q = self.W_q(x)
        C_kv = self.W_dkv(x) # Cached latent matrix

        K = self.W_uk(C_kv)
        V = self.W_uv(C_kv)

        # Standard attention computation
        scores = torch.matmul(Q, K.transpose(-2, -1))
        attention_weights = torch.softmax(scores, dim=-1)
        output = torch.matmul(attention_weights, V)

        return output
```

**Key insight:** Instead of caching K and V separately, cache a shared latent representation  $C_{kv}$  from which both can be derived.

**Memory formula:**  $l \times b \times n_l \times s \times 2$  where  $n_l$  = latent dimension

### 5. Rotary Position Encoding (RoPE)

Traditional positional encoding adds position information to embeddings, potentially distorting their magnitude. RoPE rotates the vector angle while preserving magnitude.

```
def apply_rope(x, position):
    """Apply rotary position encoding"""
    # Split embedding into pairs
    x1, x2 = x[..., ::2], x[..., 1::2]

    # Compute rotation angles
    theta = 10000 ** (-2 * torch.arange(0, x.size(-1)//2) / x.size(-1))
    angles = position.unsqueeze(-1) * theta

    cos_angles = torch.cos(angles)
    sin_angles = torch.sin(angles)

    # Apply rotation
    rotated_x1 = x1 * cos_angles - x2 * sin_angles
    rotated_x2 = x1 * sin_angles + x2 * cos_angles

    # Recombine
    result = torch.zeros_like(x)
    result[..., ::2] = rotated_x1
```

```
result[..., 1::2] = rotated_x2
```

```
return result
```

## 6. MLA + RoPE Integration

Challenge: Applying RoPE directly to MLA breaks the "absorbed query" optimization.

**Solution:** Split Q and K into RoPE and non-RoPE components:

```
class MLA_with_RoPE(nn.Module):
    def __init__(self, embed_dim, latent_dim, rope_dim):
        super().__init__()
        self.rope_dim = rope_dim
        self.non_rope_dim = embed_dim - rope_dim

        # Separate projections for RoPE and non-RoPE parts
        self.W_q_rope = nn.Linear(embed_dim, rope_dim)
        self.W_q_non_rope = nn.Linear(embed_dim, self.non_rope_dim)

        self.W_dkv = nn.Linear(embed_dim, latent_dim)
        self.W_uk_rope = nn.Linear(latent_dim, rope_dim)
        self.W_uk_non_rope = nn.Linear(latent_dim, self.non_rope_dim)

    def forward(self, x, position):
        # Split query
        q_rope = apply_rope(self.W_q_rope(x), position)
        q_non_rope = self.W_q_non_rope(x)

        # Cached latent
        c_kv = self.W_dkv(x)

        # Split key
        k_rope = apply_rope(self.W_uk_rope(c_kv), position)
        k_non_rope = self.W_uk_non_rope(c_kv)

        # Combine and compute attention
        Q = torch.cat([q_rope, q_non_rope], dim=-1)
        K = torch.cat([k_rope, k_non_rope], dim=-1)

        # Rest of attention computation...
```

---

## Advanced Optimization Techniques

### 1. Temperature Scaling

Controls randomness in text generation:

```
def apply_temperature(logits, temperature):
    """Higher temperature = more random, lower = more deterministic"""
    return logits / temperature

# Usage in generation
logits = model(input_ids)
scaled_logits = apply_temperature(logits, temperature=0.8)
probs = torch.softmax(scaled_logits, dim=-1)
```

### 2. Top-k Sampling

Limits sampling to the k most likely tokens:

```
def top_k_sampling(logits, k=50):
    top_k_logits, top_k_indices = torch.topk(logits, k)
    filtered_logits = torch.full_like(logits, float('-inf'))
    filtered_logits[top_k_indices] = top_k_logits
    return filtered_logits
```

---

## Mixture of Experts (MoE)

MoE replaces dense FFN layers with sparse expert networks, dramatically improving efficiency.

### 1. Basic MoE Concept

Instead of one large FFN, use multiple smaller "expert" networks:

```

class MoELayer(nn.Module):
    def __init__(self, embed_dim, num_experts, top_k):
        super().__init__()
        self.num_experts = num_experts
        self.top_k = top_k

        # Router network
        self.router = nn.Linear(embed_dim, num_experts)

        # Expert networks
        self.experts = nn.ModuleList([
            FeedForward(embed_dim) for _ in range(num_experts)
        ])

    def forward(self, x):
        batch_size, seq_len, embed_dim = x.shape
        x_flat = x.view(-1, embed_dim)

        # Route to experts
        router_logits = self.router(x_flat)
        routing_weights = torch.softmax(router_logits, dim=-1)

        # Select top-k experts
        top_k_weights, top_k_indices = torch.topk(routing_weights, self.top_k)
        top_k_weights = top_k_weights / top_k_weights.sum(dim=-1, keepdim=True)

        # Combine expert outputs
        output = torch.zeros_like(x_flat)
        for i in range(self.top_k):
            expert_idx = top_k_indices[:, i]
            weight = top_k_weights[:, i]

            for expert_id in range(self.num_experts):
                mask = (expert_idx == expert_id)
                if mask.any():
                    expert_output = self.experts[expert_id](x_flat[mask])
                    output[mask] += weight[mask].unsqueeze(-1) * expert_output

        return output.view(batch_size, seq_len, embed_dim)

```

## 2. Load Balancing

Ensures experts are used evenly to prevent some from being underutilized:

```

def compute_load_balance_loss(router_logits, selected_experts):
    """Auxiliary loss to encourage balanced expert usage"""
    # Expert importance (sum of routing probabilities)
    expert_importance = torch.softmax(router_logits, dim=-1).sum(dim=0)

    # Fraction of tokens routed to each expert
    num_tokens = selected_experts.numel()
    expert_usage = torch.bincount(selected_experts.flatten(),
                                   minlength=router_logits.size(-1)).float()
    expert_fraction = expert_usage / num_tokens

    # Load balance loss
    num_experts = router_logits.size(-1)
    load_balance_loss = num_experts * torch.sum(expert_fraction * expert_importance)

    return load_balance_loss

```

## 3. DeepSeek MoE Innovations

### A. Auxiliary Loss-Free Load Balancing:

```

class AdaptiveLoadBalancing(nn.Module):
    def __init__(self, num_experts, learning_rate=0.1):
        super().__init__()
        self.expert_bias = nn.Parameter(torch.zeros(num_experts))
        self.lr = learning_rate

    def update_bias(self, expert_usage, target_usage):
        """Update expert bias based on usage statistics"""
        load_violation = expert_usage - target_usage
        self.expert_bias.data += self.lr * torch.sign(load_violation)

    def forward(self, router_logits):
        return router_logits + self.expert_bias

```

### B. Shared Experts:

```

class DeepSeekMoE(nn.Module):
    def __init__(self, embed_dim, num_routed_experts, num_shared_experts, top_k):

```



```

super().__init__()

# Always-active shared experts
self.shared_experts = nn.ModuleList([
    FeedForward(embed_dim) for _ in range(num_shared_experts)
])

# Routed experts
self.routed_moe = MoELayer(embed_dim, num_routed_experts, top_k)

def forward(self, x):
    # Shared expert outputs (always computed)
    shared_output = sum(expert(x) for expert in self.shared_experts)

    # Routed expert outputs
    routed_output = self.routed_moe(x)

    return shared_output + routed_output

```

### C. Fine-Grained Expert Segmentation:

Instead of having large experts, break them into smaller, more specialized ones:

```

# Traditional: 1 expert with 4096 neurons
traditional_expert = FeedForward(embed_dim=1024, hidden_dim=4096)

# DeepSeek: 64 experts with 64 neurons each (same total parameters)
segmented_experts = nn.ModuleList([
    FeedForward(embed_dim=1024, hidden_dim=64)
    for _ in range(64)
])

```

---

## Mixture of Depths (MoD)

Mixture of Depths is an advanced technique that dynamically routes tokens through different depths of transformer layers, optimizing computation and model capacity.

### 1. Concept

Instead of passing all tokens through every transformer layer, MoD selectively routes tokens to different depths based on their complexity or importance.

### 2. Benefits

- **Efficiency:** Reduces computation by skipping layers for simpler tokens.
- **Flexibility:** Allows the model to allocate more resources to complex tokens.
- **Scalability:** Enables deeper models without proportional increase in computation.

### 3. Implementation Sketch

## Mixture of Depths (MoD)

Mixture of Depths (MoD) is an innovative approach that makes transformer layers adaptive, allowing tokens to "skip" certain layers when full processing isn't necessary. This technique significantly reduces computational cost while maintaining model performance.

### 1. Core Concept

Traditional transformers process every token through every layer uniformly. MoD introduces the idea that not all tokens need the same amount of processing - some tokens can be understood with fewer layers while others require deeper processing.

**Key Insight:** Simple tokens (like common words, punctuation) often don't benefit from deep processing, while complex tokens (rare words, context-dependent terms) need more computational attention.

### 2. MoD Architecture

```

class MixtureOfDepthsBlock(nn.Module):
    def __init__(self, embed_dim, num_heads, capacity_factor=1.0):
        super().__init__()
        self.embed_dim = embed_dim
        self.capacity = int(capacity_factor * embed_dim) # How many tokens to process

        # Router to decide which tokens need processing

```

```

self.router = nn.Linear(embed_dim, 1)

# Standard transformer components
self.attention = MultiHeadAttention(embed_dim, num_heads)
self.ffn = FeedForward(embed_dim)
self.norm1 = LayerNorm(embed_dim)
self.norm2 = LayerNorm(embed_dim)

def forward(self, x):
    batch_size, seq_len, embed_dim = x.shape

    # Compute routing scores
    router_logits = self.router(x).squeeze(-1) # (batch_size, seq_len)

    # Select top-k tokens for processing based on capacity
    top_k = min(self.capacity, seq_len)
    routing_weights, selected_indices = torch.topk(
        torch.softmax(router_logits, dim=-1),
        top_k,
        dim=-1
    )

    # Create output tensor (start with input)
    output = x.clone()

    # Process only selected tokens
    for batch_idx in range(batch_size):
        selected_tokens = x[batch_idx, selected_indices[batch_idx]] # (top_k, embed_dim)
        weights = routing_weights[batch_idx] # (top_k,)

        # Apply transformer layers to selected tokens
        # Attention
        normed_tokens = self.norm1(selected_tokens)
        attn_output = self.attention(normed_tokens.unsqueeze(0)).squeeze(0)
        selected_tokens = selected_tokens + attn_output

        # FFN
        normed_tokens = self.norm2(selected_tokens)
        ffn_output = self.ffn(normed_tokens)
        processed_tokens = selected_tokens + ffn_output

        # Weight the processed tokens
        weighted_tokens = processed_tokens * weights.unsqueeze(-1)

        # Update output
        output[batch_idx, selected_indices[batch_idx]] = weighted_tokens

    return output

```

### 3. Routing Strategies

#### A. Learned Routing:

The router learns which tokens need processing through gradient descent:

```

class LearnedMoDRouter(nn.Module):
    def __init__(self, embed_dim, hidden_dim=128):
        super().__init__()
        self.router = nn.Sequential(
            nn.Linear(embed_dim, hidden_dim),
            nn.ReLU(),
            nn.Linear(hidden_dim, 1),
            nn.Sigmoid() # Output probability of needing processing
        )

    def forward(self, x):
        return self.router(x)

def select_tokens_probabilistic(router_probs, capacity, training=True):
    """Select tokens probabilistically during training, deterministically during inference"""
    if training:
        # Gumbel-Softmax for differentiable sampling
        gumbel_noise = -torch.log(-torch.log(torch.rand_like(router_probs)))
        logits = torch.log(router_probs + 1e-8) + gumbel_noise
        selected = torch.topk(logits, capacity, dim=-1)[1]
    else:
        # Deterministic selection during inference
        selected = torch.topk(router_probs, capacity, dim=-1)[1]

    return selected

```

#### B. Content-Based Routing:

Route based on token characteristics:

```

class ContentBasedMoRouter(nn.Module):
    def __init__(self, embed_dim, vocab_size):
        super().__init__()
        self.embed_dim = embed_dim

        # Attention-based importance scoring
        self.attention_scorer = nn.MultiheadAttention(embed_dim, num_heads=8)

        # Token frequency bias (common tokens need less processing)
        self.frequency_bias = nn.Embedding(vocab_size, 1)

    def forward(self, x, token_ids):
        batch_size, seq_len, embed_dim = x.shape

        # Compute self-attention to measure token importance
        attn_output, attn_weights = self.attention_scorer(x, x, x)
        importance_scores = torch.mean(attn_weights, dim=1) # Average attention received

        # Apply frequency bias
        freq_bias = self.frequency_bias(token_ids).squeeze(-1)

        # Combine scores
        routing_scores = importance_scores.mean(dim=1) - 0.1 * freq_bias

        return torch.sigmoid(routing_scores)

```

## 4. Dynamic Capacity Allocation

Adjust the number of processed tokens based on layer depth and computational budget:

```

class DynamicCapacityMoD(nn.Module):
    def __init__(self, embed_dim, num_layers, base_capacity=0.5):
        super().__init__()
        self.num_layers = num_layers

        # Capacity decreases with depth (early layers process more tokens)
        self.layer_capacities = [
            base_capacity * (1.0 - 0.1 * i) for i in range(num_layers)
        ]

        self.layers = nn.ModuleList([
            MixtureOfDepthsBlock(embed_dim, capacity_factor=cap)
            for cap in self.layer_capacities
        ])

    def forward(self, x):
        processing_stats = []

        for i, layer in enumerate(self.layers):
            # Track which tokens are processed at each layer
            tokens_processed = min(
                int(self.layer_capacities[i] * x.size(1)),
                x.size(1)
            )
            processing_stats.append(tokens_processed)

            x = layer(x)

        return x, processing_stats

```

## 5. Load Balancing for MoD

Ensure computational load is distributed evenly across layers:

```

def compute_mod_load_balance_loss(processing_counts, target_budget):
    """
    Encourage even distribution of computational load across layers

    Args:
        processing_counts: List of tokens processed per layer
        target_budget: Target average tokens per layer
    """
    processing_tensor = torch.tensor(processing_counts, dtype=torch.float32)

    # Variance penalty - encourage even distribution
    variance_penalty = torch.var(processing_tensor)

    # Budget penalty - stay within computational budget
    total_processed = torch.sum(processing_tensor)
    budget_penalty = torch.relu(total_processed - target_budget)

    return variance_penalty + 0.1 * budget_penalty

```

# Usage in training loop

```

class MoDLoss(nn.Module):
    def __init__(self, alpha=0.01):
        super().__init__()
        self.alpha = alpha

    def forward(self, logits, targets, processing_stats, target_budget):
        # Standard language modeling loss
        lm_loss = F.cross_entropy(
            logits.view(-1, logits.size(-1)),
            targets.view(-1)
        )

        # Load balancing loss
        lb_loss = compute_mod_load_balance_loss(processing_stats, target_budget)

        return lm_loss + self.alpha * lb_loss

```

## 6. MoD Training Strategies

### A. Curriculum Learning:

Start with high capacity and gradually reduce it during training:

```

class MoDCurriculumScheduler:
    def __init__(self, initial_capacity=1.0, final_capacity=0.3, total_steps=100000):
        self.initial_capacity = initial_capacity
        self.final_capacity = final_capacity
        self.total_steps = total_steps

    def get_capacity(self, step):
        """Linearly decrease capacity during training"""
        progress = min(step / self.total_steps, 1.0)
        capacity = self.initial_capacity - progress * (
            self.initial_capacity - self.final_capacity
        )
        return capacity

    def update_model_capacity(self, model, step):
        current_capacity = self.get_capacity(step)
        for layer in model.mod_layers:
            layer.capacity = int(current_capacity * layer.max_capacity)

```

### B. Auxiliary Losses:

Guide the router to make better decisions:

```

def compute_router_auxiliary_losses(router_outputs, targets, token_ids):
    """Multiple auxiliary losses to train the router effectively"""

    # 1. Entropy regularization - encourage diverse routing decisions
    routing_probs = torch.sigmoid(router_outputs)
    entropy_loss = -torch.mean(
        routing_probs * torch.log(routing_probs + 1e-8) +
        (1 - routing_probs) * torch.log(1 - routing_probs + 1e-8)
    )

    # 2. Prediction-based routing - tokens that are harder to predict need more processing
    with torch.no_grad():
        prediction_difficulty = F.cross_entropy(
            targets.view(-1, targets.size(-1)),
            token_ids.view(-1),
            reduction='none'
        ).view(token_ids.shape)

    difficulty_correlation = F.mse_loss(
        routing_probs.squeeze(-1),
        prediction_difficulty
    )

    # 3. Sparsity regularization - encourage processing only necessary tokens
    sparsity_loss = torch.mean(routing_probs)

    return {
        'entropy': 0.01 * entropy_loss,
        'difficulty': 0.1 * difficulty_correlation,
        'sparsity': 0.05 * sparsity_loss
    }

```

## 7. MoD Performance Benefits

### Computational Savings:

```

def analyze_mod_efficiency(model, test_data):
    """Analyze computational savings from MoD"""
    total_tokens = 0

```

```

processed_tokens = 0
with torch.no_grad():
    for batch in test_data:
        batch_size, seq_len = batch.shape
        total_tokens += batch_size * seq_len * model.num_layers

        # Count actually processed tokens
        _, processing_stats = model(batch)
        processed_tokens += sum(processing_stats) * batch_size

efficiency_ratio = processed_tokens / total_tokens
flops_reduction = 1.0 - efficiency_ratio

print(f"Efficiency Ratio: {efficiency_ratio:.3f}")
print(f"FLOPs Reduction: {flops_reduction:.1%}")

return efficiency_ratio, flops_reduction

```

### Quality vs. Efficiency Trade-off:

```

def evaluate_mod_tradeoff(model, capacities, test_data):
    """Evaluate performance at different capacity levels"""
    results = []

    for capacity in capacities:
        # Set capacity
        for layer in model.mod_layers:
            layer.capacity = int(capacity * layer.max_capacity)

        # Evaluate
        perplexity = evaluate_perplexity(model, test_data)
        efficiency = capacity # Approximate efficiency

        results.append({
            'capacity': capacity,
            'perplexity': perplexity,
            'efficiency': efficiency,
            'score': perplexity * efficiency # Lower is better
        })

    return results

```

## 8. Advanced MoD Variants

### A. Hierarchical MoD:

Different capacity strategies for different types of layers:

```

class HierarchicalMoD(nn.Module):
    def __init__(self, config):
        super().__init__()

        # Different capacity strategies for different layer types
        self.attention_layers = nn.ModuleList([
            MixtureOfDepthsBlock(config.embed_dim, capacity_factor=0.8)
            for _ in range(config.num_layers)
        ])

        self.ffn_layers = nn.ModuleList([
            MixtureOfDepthsBlock(config.embed_dim, capacity_factor=0.4) # FFN needs less capacity
            for _ in range(config.num_layers)
        ])

    def forward(self, x):
        for attn_layer, ffn_layer in zip(self.attention_layers, self.ffn_layers):
            x = attn_layer(x) # Process attention with higher capacity
            x = ffn_layer(x) # Process FFN with lower capacity
        return x

```

### B. Context-Aware MoD:

Adjust processing based on sequence complexity:

```

class ContextAwareMoD(nn.Module):
    def __init__(self, embed_dim):
        super().__init__()
        self.complexity_estimator = nn.Sequential(
            nn.Linear(embed_dim, 64),
            nn.ReLU(),
            nn.Linear(64, 1),
            nn.Sigmoid()
        )

    def estimate_sequence_complexity(self, x):
        """Estimate how complex the sequence is"""

```

```

# Use attention patterns to measure complexity
attn_entropy = self.compute_attention_entropy(x)

# Use token diversity
token_diversity = self.compute_token_diversity(x)

# Combine measures
complexity = (attn_entropy + token_diversity) / 2
return complexity

def adaptive_capacity(self, x, base_capacity=0.5):
    complexity = self.estimate_sequence_complexity(x)
    # More complex sequences get higher capacity
    adaptive_capacity = base_capacity + 0.3 * complexity
    return min(adaptive_capacity, 1.0)

```

## 9. Combining MoD with MoE

MoD and MoE can be combined for maximum efficiency:

```

class MoD_MoE_Layer(nn.Module):
    def __init__(self, embed_dim, num_experts, mod_capacity=0.5):
        super().__init__()

        # MoD router
        self.mod_router = nn.Linear(embed_dim, 1)
        self.mod_capacity = mod_capacity

        # MoE layer (only applied to selected tokens)
        self.moe_layer = MoELayer(embed_dim, num_experts, top_k=2)

        # Standard layers for non-selected tokens
        self.simple_ffn = nn.Linear(embed_dim, embed_dim)

    def forward(self, x):
        batch_size, seq_len, embed_dim = x.shape

        # MoD selection
        mod_scores = torch.sigmoid(self.mod_router(x)).squeeze(-1)
        num_selected = int(self.mod_capacity * seq_len)

        # Select tokens for deep processing
        selected_mask = torch.topk(mod_scores, num_selected, dim=-1)[1]

        output = x.clone()

        for batch_idx in range(batch_size):
            # Deep processing with MoE for selected tokens
            selected_indices = selected_mask[batch_idx]
            selected_tokens = x[batch_idx, selected_indices]

            processed_tokens = self.moe_layer(selected_tokens.unsqueeze(0)).squeeze(0)
            output[batch_idx, selected_indices] = processed_tokens

            # Simple processing for non-selected tokens
            non_selected_mask = torch.ones(seq_len, dtype=torch.bool)
            non_selected_mask[selected_indices] = False

            if non_selected_mask.any():
                simple_tokens = x[batch_idx, non_selected_mask]
                simple_processed = self.simple_ffn(simple_tokens)
                output[batch_idx, non_selected_mask] = simple_processed

        return output

```

## 10. MoD Implementation Best Practices

### A. Gradient Handling:

```

class StraightThroughMoDRouter(nn.Module):
    """Router with straight-through gradients for discrete selection"""

    def __init__(self, embed_dim):
        super().__init__()
        self.router = nn.Linear(embed_dim, 1)

    def forward(self, x, capacity, training=True):
        router_probs = torch.sigmoid(self.router(x))

        if training:
            # Use Gumbel-Softmax for differentiable discrete sampling
            return self.gumbel_softmax_selection(router_probs, capacity)
        else:
            # Use deterministic selection during inference

```

```
return self.deterministic_selection(router_probs, capacity)
```

```
def gumbel_softmax_selection(self, probs, capacity, temperature=1.0):  
    # Implementation of differentiable discrete selection  
    gumbel_noise = -torch.log(-torch.log(torch.rand_like(probs)))  
    logits = torch.log(probs + 1e-8) + gumbel_noise / temperature  
  
    # Straight-through estimator  
    hard_selection = torch.topk(logits, capacity, dim=-1)[1]  
    soft_selection = torch.softmax(logits / temperature, dim=-1)  
  
    # Straight-through: forward uses hard, backward uses soft  
    return hard_selection + soft_selection - soft_selection.detach()
```

## B. Memory Management:

```
def memory_efficient_mod_forward(mod_layer, x, capacity_ratio=0.5):  
    """Memory-efficient implementation of MoD forward pass"""  
    batch_size, seq_len, embed_dim = x.shape  
    capacity = int(capacity_ratio * seq_len)  
  
    # Process in chunks to save memory  
    chunk_size = min(1000, seq_len) # Process 1000 tokens at a time  
  
    outputs = []  
    for start_idx in range(0, seq_len, chunk_size):  
        end_idx = min(start_idx + chunk_size, seq_len)  
        chunk = x[:, start_idx:end_idx]  
  
        # Apply MoD to chunk  
        chunk_output = mod_layer(chunk)  
        outputs.append(chunk_output)  
  
        # Clear intermediate computations  
        del chunk  
        torch.cuda.empty_cache() if torch.cuda.is_available() else None  
  
    return torch.cat(outputs, dim=1)
```

MoD represents a significant advancement in making transformer models more efficient by allowing adaptive computation. When combined with other techniques like MoE, it enables the creation of highly efficient yet powerful language models that can scale to massive sizes while maintaining reasonable computational costs.

---

```
class MixtureOfDepths(nn.Module):  
    def __init__(self, transformer_blocks, router):  
        super().__init__()  
        self.transformer_blocks = transformer_blocks # List of layers  
        self.router = router # Decides depth per token  
  
    def forward(self, x):  
        batch_size, seq_len, embed_dim = x.size()  
  
        # Get routing decisions: depth per token (shape: batch_size x seq_len)  
        depths = self.router(x)  
  
        outputs = torch.zeros_like(x)  
  
        for depth in range(len(self.transformer_blocks)):  
            # Select tokens that need to go through this depth  
            mask = (depths > depth)  
  
            if mask.any():  
                # Process tokens at this depth  
                x_masked = x[mask]  
                x_processed = self.transformer_blocks[depth](x_masked)  
                outputs[mask] = x_processed  
  
        return outputs
```

## 4. Challenges

- Designing an effective router that balances load and performance.
- Managing gradient flow through variable-depth paths.
- Ensuring stable training with dynamic routing.

---

## Conclusion

This comprehensive guide covered the journey from basic GPT architecture to cutting-edge innovations in DeepSeek. Key takeaways:

## Foundation Concepts:

- **Data preprocessing** with tokenization and embeddings
- **Attention mechanisms** from simple to multi-head causal attention
- **Transformer blocks** with layer normalization, FFN, and residual connections
- **Training process** with cross-entropy loss and generation strategies

## DeepSeek Innovations:

- **Multi-Head Latent Attention (MLA)** - solves KV caching memory issues while maintaining performance
- **Rotary Position Encoding (RoPE)** - preserves embedding magnitude while encoding positions
- **Mixture of Experts (MoE)** - enables sparse computation for massive model scaling
- **Advanced load balancing** - ensures efficient expert utilization

## Key Performance Improvements:

- **Memory efficiency:** From 400GB to 6GB KV cache requirement
- **Computational efficiency:** Sparse activation through MoE
- **Model quality:** Better position encoding and attention mechanisms

These innovations have pushed the boundaries of what's possible with LLMs, enabling larger, more efficient, and more capable language models. The field continues to evolve rapidly, with new architectures and training techniques constantly emerging.

For implementation details and code examples, refer to the original repository's Jupyter notebooks and detailed documentation.