# LFM2: Liquid Foundation Model 2
## Step-by-Step Implementation Guide

Complete Tutorial

September 6, 2025

## Contents

# 1 Introduction

LFM2 (Liquid Foundation Model 2) represents a hybrid neural architecture combining convolution operations for local processing with attention mechanisms for global dependencies. This tutorial walks through each component step-by-step, showing the architecture diagram first, followed immediately by the corresponding implementation.

# 2 Foundation Components

## 2.1 Configuration and Setup

Let's start with the basic configuration that defines our model:

```python
import math
import torch
import torch.nn as nn
import torch.nn.functional as F
from typing import Optional, Tuple, Dict, Any, Union
from dataclasses import dataclass

@dataclass
class LFM2Config:
    """Configuration class for LFM2 model."""
    vocab_size: int = 32768            # Size of the vocabulary
    hidden_size: int = 1024            # Hidden dimension size
    intermediate_size: int = 2816      # Feedforward intermediate dimension
    num_conv_blocks: int = 10          # Number of LIV convolution blocks
    num_attention_blocks: int = 6      # Number of GQA blocks
    num_attention_heads: int = 16      # Number of attention heads
    num_key_value_heads: int = 4       # Number of key-value heads for GQA
    conv_kernel_size: int = 3          # Kernel size for short convolutions
    max_position_embeddings: int = 32768  # Maximum sequence length
    rms_norm_eps: float = 1e-6         # RMS normalization epsilon
    rope_theta: float = 10000.0        # RoPE base frequency
    attention_dropout: float = 0.0     # Dropout rate for attention
    hidden_dropout: float = 0.0        # Dropout rate for hidden layers
    initializer_range: float = 0.02    # Weight initialization std
    use_cache: bool = True             # Enable KV caching
    tie_word_embeddings: bool = True   # Tie input/output embeddings
```

This configuration defines the hybrid architecture: 10 convolution blocks for local processing, followed by 6 attention blocks for global dependencies. The grouped query attention uses 16 query heads but only 4 key-value heads (4:1 ratio) for memory efficiency.

## 2.2 RMSNorm Implementation

```python
class RMSNorm(nn.Module):
    """Root Mean Square Layer Normalization.

    RMSNorm normalizes using only the root mean square, without centering:
    RMSNorm(x) = (x / RMS(x)) * learnable_scale
    where RMS(x) = sqrt(mean(x^2))
    """

```

```python
 9      def __init__(self, hidden_size: int, eps: float = 1e-6):
10          super().__init__()
11          # Learnable scale parameter (no bias unlike LayerNorm)
12          self.weight = nn.Parameter(torch.ones(hidden_size))
13          self.variance_epsilon = eps
14
15      def forward(self, hidden_states: torch.Tensor) -> torch.Tensor:
16          # Store original dtype for final output
17          input_dtype = hidden_states.dtype
18
19          # Convert to float32 for stable computation
20          hidden_states = hidden_states.to(torch.float32)
21
22          # Compute variance: mean of squared values
23          variance = hidden_states.pow(2).mean(-1, keepdim=True)
24
25          # Normalize by RMS: x / sqrt(variance + eps)
26          hidden_states = hidden_states * torch.rsqrt(variance + self.
                variance_epsilon)
27
28          # Apply learnable scale and convert back to original dtype
29          return self.weight * hidden_states.to(input_dtype)
```

RMSNorm is more stable than LayerNorm because it doesn't center the data (no mean subtraction), only scales by the RMS. This reduces computation and often improves training stability.

## 2.3 Rotary Positional Embedding (RoPE)

```python
 1  class RotaryPositionalEmbedding(nn.Module):
 2      """Rotary Positional Embedding implementation.
 3
 4      RoPE encodes position by rotating query and key vectors by position-dependent
            angles.
 5      This allows the model to understand relative positions naturally.
 6      """
 7
 8      def __init__(self, dim: int, max_position_embeddings: int = 32768,
 9                   base: float = 10000.0):
10          super().__init__()
11          self.dim = dim
12          self.max_position_embeddings = max_position_embeddings
13          self.base = base
14
15          # Precompute frequency inverse: 1 / (base^(2i/d))
16          # This creates different rotation frequencies for each dimension pair
17          inv_freq = 1.0 / (self.base **
18                            (torch.arange(0, self.dim, 2).float() / self.dim))
19          self.register_buffer("inv_freq", inv_freq, persistent=False)
20
21      def forward(self, x: torch.Tensor, seq_len: int) -> Tuple[torch.Tensor, torch.
            Tensor]:
22          # Generate position indices [0, 1, 2, ..., seq_len-1]
23          t = torch.arange(seq_len, device=x.device).type_as(self.inv_freq)
24
25          # Compute frequencies for each position: pos * inv_freq
26          freqs = torch.outer(t, self.inv_freq)
27
```

```
28          # Duplicate frequencies for sine and cosine (each dim pair needs both)
29          emb = torch.cat((freqs, freqs), dim=-1)
30
31          # Return cosine and sine components
32          cos = emb.cos()
33          sin = emb.sin()
34
35          return cos, sin
36
37 def rotate_half(x: torch.Tensor) -> torch.Tensor:
38      """Rotates half the hidden dims of the input.
39
40      For RoPE, we need to rotate pairs of dimensions. This function
41      swaps the first and second half and negates the second half.
42      """
43      x1 = x[..., : x.shape[-1] // 2]  # First half
44      x2 = x[..., x.shape[-1] // 2 :]  # Second half
45      return torch.cat((-x2, x1), dim=-1)  # Rotate: [-x2, x1]
46
47 def apply_rotary_pos_emb(q: torch.Tensor, k: torch.Tensor,
48                          cos: torch.Tensor, sin: torch.Tensor) -> Tuple[torch.
                                Tensor, torch.Tensor]:
49      """Apply rotary positional embedding to query and key tensors.
50
51      The rotation is: x_rotated = x * cos + rotate_half(x) * sin
52      This applies the rotation matrix in complex number form.
53      """
54      q_embed = (q * cos) + (rotate_half(q) * sin)
55      k_embed = (k * cos) + (rotate_half(k) * sin)
56      return q_embed, k_embed
```

RoPE works by treating pairs of dimensions as complex numbers and rotating them by position-dependent angles. This creates a natural way for the model to understand relative positions between tokens.

# 3   LFM2 Convolution Block

## 3.1   Architecture Diagram

## 3.2   Implementation

```
1 class LFM2ConvBlock(nn.Module):
2      """LFM2 double-gated short-range convolution block.
3
4      This implements Linear Input-Varying (LIV) convolution:
5      1. Project input to 3 components: gates B, C, and values x'
6      2. First gating: B * x' (input-dependent gating)
7      3. Depthwise convolution for local pattern extraction
8      4. Second gating: C * conv_output (input-dependent gating)
9      5. Final projection back to hidden size
10     """
11
12     def __init__(self, config: LFM2Config):
13         super().__init__()
14         self.config = config
```
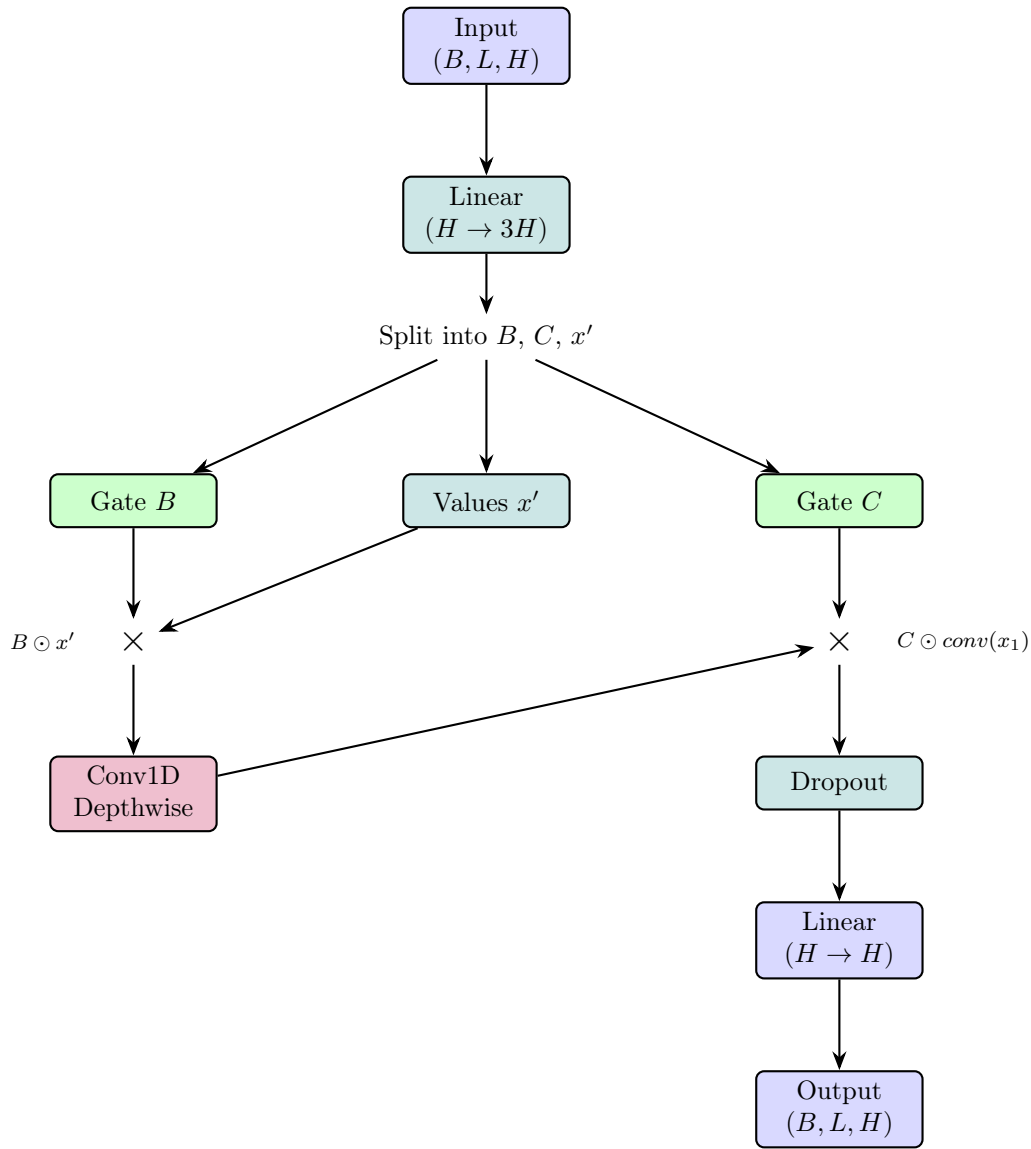
# LFM2 Convolution Block



Figure 1: Double-gated Linear Input-Varying convolution with input-dependent gates B and C.

```
15          self.hidden_size = config.hidden_size
16
17          # Input projection to gates and values (splits into B, C, x)
18          self.input_projection = nn.Linear(
19              self.hidden_size,
20              3 * self.hidden_size,   # Output: 3 * hidden_size for B, C, x
21              bias=False,
22          )
23
24          # Depthwise convolution: each channel processed independently
25          # This is memory efficient and captures local patterns
```

```python
26          self.conv = nn.Conv1d(
27              self.hidden_size,              # Input channels
28              self.hidden_size,              # Output channels (same as input)
29              kernel_size=config.conv_kernel_size,   # Usually 3
30              padding=config.conv_kernel_size // 2,   # Same padding
31              groups=self.hidden_size,       # Depthwise: groups = channels
32              bias=False,
33          )
34
35          # Output projection back to original hidden size
36          self.output_projection = nn.Linear(
37              self.hidden_size, self.hidden_size, bias=False
38          )
39
40          # Dropout for regularization
41          self.dropout = nn.Dropout(config.hidden_dropout)
42
43      def forward(self, x: torch.Tensor) -> torch.Tensor:
44          """Forward pass implementing the LIV convolution.
45
46          Args:
47              x: Input tensor of shape (batch_size, seq_len, hidden_size)
48
49          Returns:
50              Output tensor of same shape as input
51          """
52          batch_size, seq_len, hidden_size = x.shape
53
54          # Step 1: Project input to gates B, C and values x'
55          # Input: (B, L, H) -> Output: (B, L, 3H)
56          projected = self.input_projection(x)
57
58          # Split into three equal parts: B, C, x'
59          # Each has shape (B, L, H)
60          B, C, x_proj = projected.chunk(3, dim=-1)
61
62          # Step 2: First gating - multiply values by gate B
63          # This allows input-dependent filtering before convolution
64          x_gated = B * x_proj   # Element-wise multiplication
65
66          # Step 3: Apply depthwise convolution
67          # Conv1D expects (batch, channels, sequence), so transpose
68          x_conv_input = x_gated.transpose(1, 2)   # (B, L, H) -> (B, H, L)
69          x_conv = self.conv(x_conv_input)         # Apply convolution
70          x_conv = x_conv.transpose(1, 2)          # Back to (B, L, H)
71
72          # Step 4: Second gating - multiply conv output by gate C
73          # This provides input-dependent filtering after convolution
74          x_gated_2 = C * x_conv
75
76          # Step 5: Apply dropout for regularization
77          x_gated_2 = self.dropout(x_gated_2)
78
79          # Step 6: Final output projection
80          output = self.output_projection(x_gated_2)
81
82          return output
```

The key innovation here is the double gating mechanism. Unlike standard convolution, both

gates B and C depend on the input, making this a "Linear Input-Varying" convolution. The depth-wise convolution (groups=hidden$_{size}$)$is computationally efficient while still capturing local patterns.$

# 4 Grouped Query Attention

## 4.1 Architecture Diagram

## 4.2 Implementation

```python
class GroupedQueryAttention(nn.Module):
    """Grouped Query Attention implementation.

    GQA reduces memory usage by using fewer key-value heads than query heads.
    For example: 16 query heads, 4 key-value heads means each KV head
    is shared across 4 query heads (16/4 = 4 groups).
    """

    def __init__(self, config: LFM2Config):
        super().__init__()
        self.config = config
        self.hidden_size = config.hidden_size
        self.num_heads = config.num_attention_heads          # e.g., 16
        self.num_key_value_heads = config.num_key_value_heads # e.g., 4
        self.head_dim = self.hidden_size // self.num_heads    # e.g., 64

        # Number of query heads per key-value head
        self.num_key_value_groups = self.num_heads // self.num_key_value_heads

        # Validate configuration
        if self.hidden_size % self.num_heads != 0:
            raise ValueError(f"hidden_size must be divisible by num_heads")

        # Linear projections - note different output sizes
        self.q_proj = nn.Linear(
            self.hidden_size,
            self.num_heads * self.head_dim,         # Full size: 16 * 64
            bias=False,
        )
        self.k_proj = nn.Linear(
            self.hidden_size,
            self.num_key_value_heads * self.head_dim, # Reduced: 4 * 64
            bias=False,
        )
        self.v_proj = nn.Linear(
            self.hidden_size,
            self.num_key_value_heads * self.head_dim, # Reduced: 4 * 64
            bias=False,
        )
        self.o_proj = nn.Linear(
            self.num_heads * self.head_dim,         # Back to full hidden size
            self.hidden_size,
            bias=False,
        )

        # Rotary positional embedding
        self.rotary_emb = RotaryPositionalEmbedding(
```
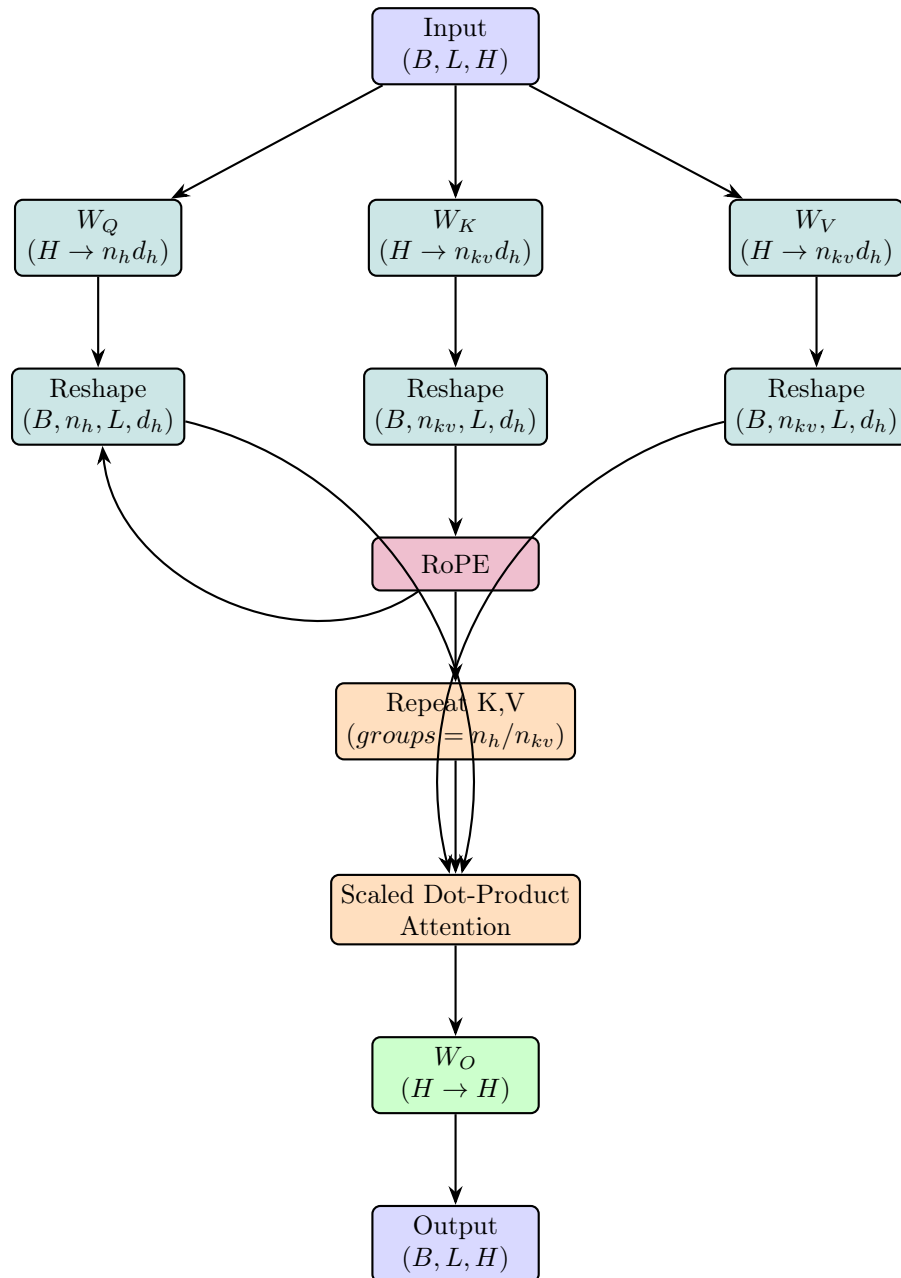
# Grouped Query Attention

Figure 2: GQA uses fewer key-value heads than query heads, then repeats K,V to match Q heads.

```
48              self.head_dim,
49              config.max_position_embeddings,
50              config.rope_theta,
51          )
52
53          # Attention dropout
54          self.attention_dropout = nn.Dropout(config.attention_dropout)
```

```python
55
56     def forward(self, hidden_states: torch.Tensor,
57                 attention_mask: Optional[torch.Tensor] = None,
58                 position_ids: Optional[torch.Tensor] = None,
59                 past_key_value: Optional[Tuple[torch.Tensor]] = None,
60                 output_attentions: bool = False,
61                 use_cache: bool = False) -> Tuple[torch.Tensor, Optional[torch.
                       Tensor], Optional[Tuple[torch.Tensor]]]:
62
63         bsz, q_len, _ = hidden_states.size()
64
65         # Step 1: Project input to Q, K, V
66         query_states = self.q_proj(hidden_states)  # (B, L, n_h * d_h)
67         key_states = self.k_proj(hidden_states)    # (B, L, n_kv * d_h)
68         value_states = self.v_proj(hidden_states)  # (B, L, n_kv * d_h)
69
70         # Step 2: Reshape for multi-head attention
71         # Q: (B, L, n_h * d_h) -> (B, n_h, L, d_h)
72         query_states = query_states.view(
73             bsz, q_len, self.num_heads, self.head_dim
74         ).transpose(1, 2)
75
76         # K, V: (B, L, n_kv * d_h) -> (B, n_kv, L, d_h)
77         key_states = key_states.view(
78             bsz, q_len, self.num_key_value_heads, self.head_dim
79         ).transpose(1, 2)
80         value_states = value_states.view(
81             bsz, q_len, self.num_key_value_heads, self.head_dim
82         ).transpose(1, 2)
83
84         # Step 3: Get sequence length for RoPE (including past context)
85         kv_seq_len = key_states.shape[-2]
86         if past_key_value is not None:
87             kv_seq_len += past_key_value[0].shape[-2]
88
89         # Step 4: Apply rotary positional embedding to Q and K
90         cos, sin = self.rotary_emb(value_states, kv_seq_len)
91         query_states, key_states = apply_rotary_pos_emb(
92             query_states, key_states, cos, sin
93         )
94
95         # Step 5: Handle past key-value cache for generation
96         if past_key_value is not None:
97             # Concatenate past and current key-values
98             key_states = torch.cat([past_key_value[0], key_states], dim=2)
99             value_states = torch.cat([past_key_value[1], value_states], dim=2)
100
101        # Prepare cache for next iteration
102        past_key_value = (key_states, value_states) if use_cache else None
103
104        # Step 6: Repeat K,V to match number of query heads (GQA core)
105        # From (B, n_kv, L, d_h) to (B, n_h, L, d_h)
106        key_states = key_states.repeat_interleave(
107            self.num_key_value_groups, dim=1
108        )
109        value_states = value_states.repeat_interleave(
110            self.num_key_value_groups, dim=1
111        )
112
```

```
113        # Step 7: Compute scaled dot-product attention
114        # Q @ K^T / sqrt(d_k)
115        attn_weights = torch.matmul(
116            query_states, key_states.transpose(2, 3)
117        ) / math.sqrt(self.head_dim)
118
119        # Apply attention mask if provided (for padding/causality)
120        if attention_mask is not None:
121            attn_weights = attn_weights + attention_mask
122
123        # Step 8: Apply softmax and dropout
124        attn_weights = F.softmax(attn_weights, dim=-1, dtype=torch.float32)
125        attn_weights = attn_weights.to(query_states.dtype)
126        attn_weights = self.attention_dropout(attn_weights)
127
128        # Step 9: Apply attention to values
129        attn_output = torch.matmul(attn_weights, value_states)
130
131        # Step 10: Reshape back to original format and apply output projection
132        attn_output = attn_output.transpose(1, 2).contiguous()
133        attn_output = attn_output.reshape(bsz, q_len, self.hidden_size)
134        attn_output = self.o_proj(attn_output)
135
136        # Return outputs (optionally include attention weights)
137        if not output_attentions:
138            attn_weights = None
139
140        return attn_output, attn_weights, past_key_value
```

The memory efficiency of GQA comes from step 6: instead of computing separate K,V for each head, we compute fewer K,V heads and repeat them. This reduces the KV cache size significantly during generation.

# 5 SwiGLU Feed-Forward Network

## 5.1 Implementation

```
1  class SwiGLU(nn.Module):
2      """SwiGLU activation function.
3
4      SwiGLU combines Swish activation with Gated Linear Units:
5      SwiGLU(x) = Swish(xW_gate) * (xW_up) * W_down
6
7      Where Swish(x) = x * sigmoid(x) = x *   (x)
8      This has been shown to work better than ReLU-based FFNs.
9      """
10
11     def __init__(self, config: LFM2Config):
12         super().__init__()
13         self.hidden_size = config.hidden_size
14         self.intermediate_size = config.intermediate_size   # Usually ~2.7x
              hidden_size
15
16         # Two parallel transformations to intermediate size
17         self.gate_proj = nn.Linear(
18             self.hidden_size, self.intermediate_size, bias=False
```

```
19          )
20          self.up_proj = nn.Linear(
21              self.hidden_size, self.intermediate_size, bias=False
22          )
23
24          # Down projection back to hidden size
25          self.down_proj = nn.Linear(
26              self.intermediate_size, self.hidden_size, bias=False
27          )
28
29      def forward(self, x: torch.Tensor) -> torch.Tensor:
30          """Forward pass implementing SwiGLU.
31
32          Args:
33              x: Input tensor of shape (..., hidden_size)
34
35          Returns:
36              Output tensor of same shape as input
37          """
38          # Apply both linear projections in parallel
39          gate = self.gate_proj(x)    # "Gate" branch: controls information flow
40          up = self.up_proj(x)        # "Up" branch: carries information
41
42          # Apply SiLU (Swish) activation to gate branch
43          # SiLU(x) = x * sigmoid(x) - smooth, non-monotonic activation
44          swish_gate = F.silu(gate)
45
46          # Multiply gate and up branches (gating mechanism)
47          gated_output = swish_gate * up
48
49          # Project back down to original hidden size
50          return self.down_proj(gated_output)
```

SwiGLU has been empirically shown to outperform ReLU-based feed-forward networks. The gating mechanism allows the network to control information flow, while SiLU provides smooth gradients compared to ReLU.

# 6  LFM2 Transformer Block

## 6.1  Architecture Diagram

## 6.2  Implementation

```
1  class LFM2Block(nn.Module):
2      """LFM2 transformer block with GQA and SwiGLU.
3
4      Uses pre-normalization pattern:
5      1.  x   = x + Attention(RMSNorm(x))
6      2.  x   =  x   + SwiGLU(RMSNorm( x  ))
7
8      Pre-norm helps with training stability compared to post-norm.
9      """
10
11      def __init__(self, config: LFM2Config):
12          super().__init__()
```

# LFM2 Transformer Block
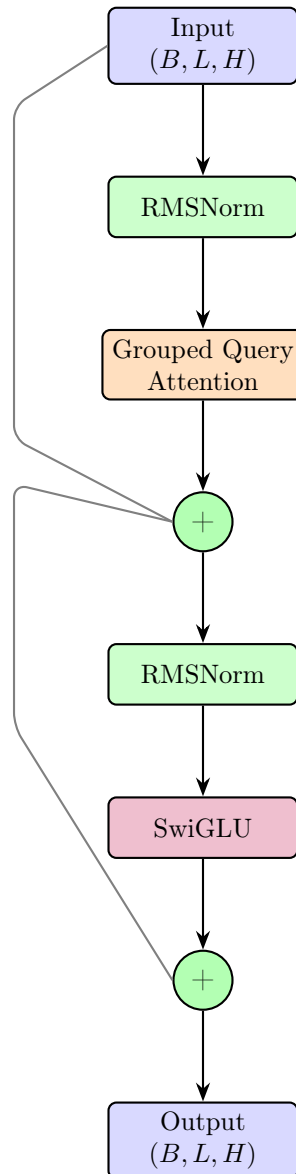


Figure 3: Pre-normalization transformer block with residual connections.

```
13          self.config = config
14          self.hidden_size = config.hidden_size
15
16          # Main components
17          self.self_attn = GroupedQueryAttention(config)
18          self.mlp = SwiGLU(config)
19
20          # Layer normalization (pre-norm pattern)
21          self.input_layernorm = RMSNorm(
22              config.hidden_size, eps=config.rms_norm_eps
```

```python
23        )
24        self.post_attention_layernorm = RMSNorm(
25            config.hidden_size, eps=config.rms_norm_eps
26        )
27
28    def forward(self, hidden_states: torch.Tensor,
29                attention_mask: Optional[torch.Tensor] = None,
30                position_ids: Optional[torch.Tensor] = None,
31                past_key_value: Optional[Tuple[torch.Tensor]] = None,
32                output_attentions: Optional[bool] = False,
33                use_cache: Optional[bool] = False) -> Tuple[torch.Tensor, Optional
                    [torch.Tensor], Optional[Tuple[torch.Tensor]]]:
34        """Forward pass implementing pre-norm transformer block."""
35
36        # First residual block: Self-attention
37        residual = hidden_states  # Store for residual connection
38
39        # Pre-normalization before attention
40        hidden_states = self.input_layernorm(hidden_states)
41
42        # Apply self-attention
43        hidden_states, self_attn_weights, present_key_value = self.self_attn(
44            hidden_states=hidden_states,
45            attention_mask=attention_mask,
46            position_ids=position_ids,
47            past_key_value=past_key_value,
48            output_attentions=output_attentions,
49            use_cache=use_cache,
50        )
51
52        # Add residual connection
53        hidden_states = residual + hidden_states
54
55        # Second residual block: Feed-forward network
56        residual = hidden_states  # Store for next residual connection
57
58        # Pre-normalization before MLP
59        hidden_states = self.post_attention_layernorm(hidden_states)
60
61        # Apply SwiGLU feed-forward network
62        hidden_states = self.mlp(hidden_states)
63
64        # Add residual connection
65        hidden_states = residual + hidden_states
66
67        # Prepare outputs
68        outputs = (hidden_states,)
69
70        if output_attentions:
71            outputs += (self_attn_weights,)
72
73        if use_cache:
74            outputs += (present_key_value,)
75
76        return outputs
```

The pre-normalization pattern (norm before operation) has become standard in modern transformers because it provides more stable gradients during training compared to post-normalization.

# 7    Complete LFM2 Model

## 7.1    Architecture Diagram
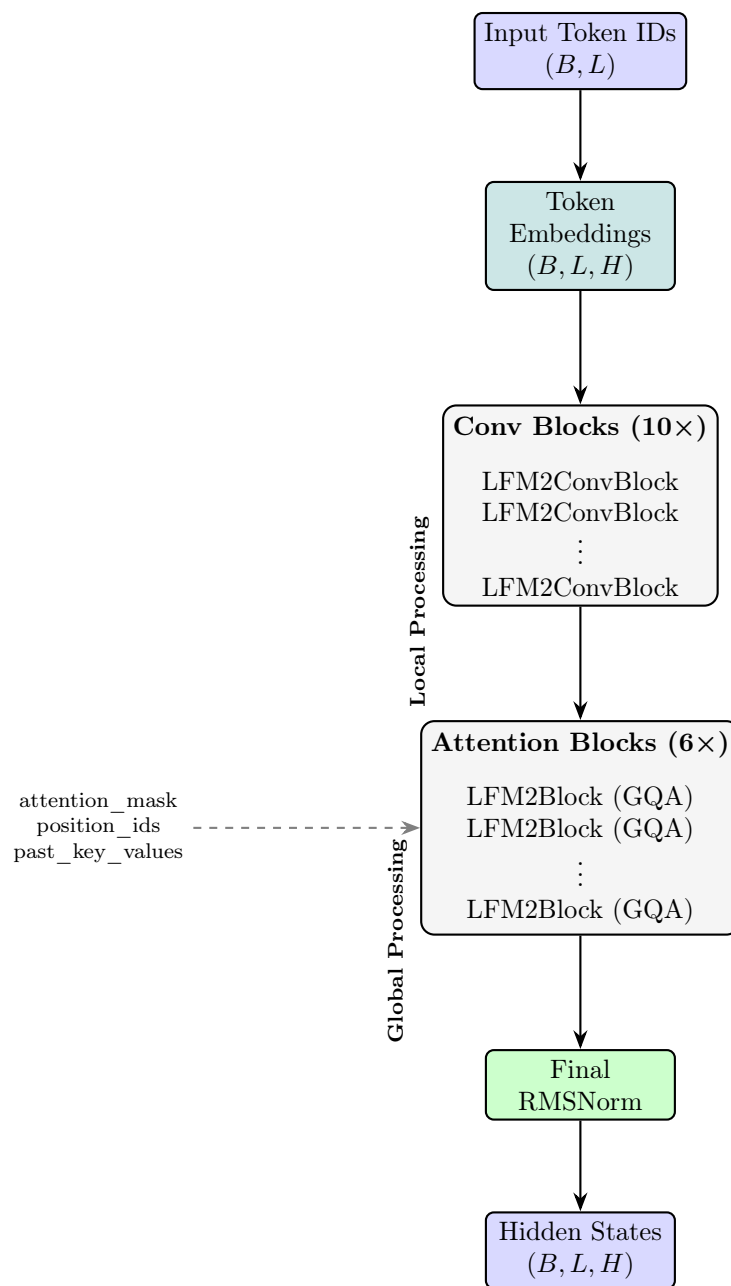
**LFM2 (Liquid Foundation Model 2)**

Figure 4: Complete LFM2 hybrid architecture: convolution for local patterns, attention for global dependencies.

## 7.2   Implementation

```python
class LFM2Model(nn.Module):
    """Complete LFM2 model with hybrid architecture.

    Architecture:
    - Token embeddings
    - 10 LFM2ConvBlocks (local pattern processing)
    - 6 LFM2Blocks with GQA (global dependency modeling)
    - Final RMSNorm

    This hybrid approach allows efficient local processing followed
    by global attention for long-range dependencies.
    """

    def __init__(self, config: LFM2Config):
        super().__init__()
        self.config = config
        self.vocab_size = config.vocab_size
        self.hidden_size = config.hidden_size

        # Token embeddings: convert token IDs to dense vectors
        self.embed_tokens = nn.Embedding(config.vocab_size, config.hidden_size)

        # Model layers: hybrid convolution + attention
        self.layers = nn.ModuleList()

        # Stage 1: Add convolution blocks (local processing)
        for i in range(config.num_conv_blocks):
            self.layers.append(LFM2ConvBlock(config))

        # Stage 2: Add attention blocks (global processing)
        for i in range(config.num_attention_blocks):
            self.layers.append(LFM2Block(config))

        # Final layer normalization
        self.norm = RMSNorm(config.hidden_size, eps=config.rms_norm_eps)

        # Initialize weights
        self.apply(self._init_weights)

    def _init_weights(self, module: nn.Module) -> None:
        """Initialize model weights following standard practices."""
        if isinstance(module, nn.Linear):
            # Normal initialization for linear layers
            torch.nn.init.normal_(
                module.weight, mean=0.0, std=self.config.initializer_range
            )
            if module.bias is not None:
                torch.nn.init.zeros_(module.bias)
        elif isinstance(module, nn.Embedding):
            # Normal initialization for embeddings
            torch.nn.init.normal_(
                module.weight, mean=0.0, std=self.config.initializer_range
            )
        elif isinstance(module, nn.Conv1d):
            # Normal initialization for convolutions
            torch.nn.init.normal_(
```

```python
57                    module.weight, mean=0.0, std=self.config.initializer_range
58                )
59                if module.bias is not None:
60                    torch.nn.init.zeros_(module.bias)
61
62    def forward(self, input_ids: torch.LongTensor,
63                attention_mask: Optional[torch.Tensor] = None,
64                position_ids: Optional[torch.LongTensor] = None,
65                past_key_values: Optional[Tuple[Tuple[torch.Tensor]]] = None,
66                use_cache: Optional[bool] = None,
67                output_hidden_states: Optional[bool] = None,
68                return_dict: Optional[bool] = None) -> Union[Tuple, Dict[str,
                    torch.Tensor]]:
69        """Forward pass through the complete LFM2 model."""
70
71        # Set default values
72        use_cache = use_cache if use_cache is not None else self.config.use_cache
73        output_hidden_states = (output_hidden_states if output_hidden_states is
            not None
74                                else self.config.output_hidden_states)
75        return_dict = return_dict if return_dict is not None else self.config.
            use_return_dict
76
77        # Step 1: Convert token IDs to embeddings
78        hidden_states = self.embed_tokens(input_ids)
79        batch_size, seq_length = hidden_states.shape[:2]
80
81        # Step 2: Initialize position IDs if not provided
82        if position_ids is None:
83            device = input_ids.device if input_ids is not None else hidden_states.
                device
84            position_ids = torch.arange(seq_length, dtype=torch.long, device=
                device)
85            position_ids = position_ids.unsqueeze(0).expand(batch_size, -1)
86
87        # Step 3: Initialize past key values for caching
88        if past_key_values is None:
89            past_key_values = [None] * len(self.layers)
90
91        # Step 4: Prepare attention mask for causal modeling
92        if attention_mask is not None:
93            attention_mask = self._prepare_attention_mask(
94                attention_mask, seq_length, batch_size
95            )
96
97        # Step 5: Initialize output containers
98        all_hidden_states = () if output_hidden_states else None
99        next_decoder_cache = () if use_cache else None
100
101        # Step 6: Process through all layers
102        for idx, decoder_layer in enumerate(self.layers):
103            if output_hidden_states:
104                all_hidden_states += (hidden_states,)
105
106            past_key_value = past_key_values[idx] if past_key_values is not None
                else None
107
108            if isinstance(decoder_layer, LFM2ConvBlock):
109                # Convolution blocks: simple forward pass (no attention cache)
```

```python
110                        hidden_states = decoder_layer(hidden_states)
111                        layer_outputs = (hidden_states,)
112                    else:
113                        # Attention blocks: full forward pass with caching
114                        layer_outputs = decoder_layer(
115                            hidden_states,
116                            attention_mask=attention_mask,
117                            position_ids=position_ids,
118                            past_key_value=past_key_value,
119                            use_cache=use_cache,
120                        )
121                        hidden_states = layer_outputs[0]
122
123                    # Collect cache from attention layers only
124                    if use_cache and isinstance(decoder_layer, LFM2Block):
125                        next_decoder_cache += (layer_outputs[2],)
126
127            # Step 7: Apply final normalization
128            hidden_states = self.norm(hidden_states)
129
130            # Add final hidden state to output
131            if output_hidden_states:
132                all_hidden_states += (hidden_states,)
133
134            # Prepare final cache
135            next_cache = next_decoder_cache if use_cache else None
136
137            # Return in requested format
138            if not return_dict:
139                return tuple(v for v in [hidden_states, next_cache, all_hidden_states]
140                        if v is not None)
141
142            return {
143                "last_hidden_state": hidden_states,
144                "past_key_values": next_cache,
145                "hidden_states": all_hidden_states,
146            }
147
148        def _prepare_attention_mask(self, attention_mask: torch.Tensor,
149                            seq_length: int, batch_size: int) -> torch.Tensor:
150            """Prepare causal attention mask combining padding and causality."""
151
152            # Create causal mask (lower triangular)
153            causal_mask = torch.triu(
154                torch.ones((seq_length, seq_length), dtype=torch.bool),
155                diagonal=1,  # Upper triangle (including diagonal) = True (masked)
156            ).to(attention_mask.device)
157
158            # Expand attention mask to 4D: (batch, 1, seq, seq)
159            expanded_mask = (attention_mask[:, None, None, :]
160                        .expand(batch_size, 1, seq_length, seq_length)
161                        .to(torch.float32))
162
163            # Convert to mask format: 1.0 = masked, 0.0 = not masked
164            expanded_mask = 1.0 - expanded_mask
165            causal_float_mask = causal_mask.to(torch.float32)
166
167            # Combine padding and causal masks
168            combined_mask = expanded_mask + causal_float_mask[None, None, :, :]
```

```
168
169            # Convert to attention weights: 0 = attend, -inf = mask
170            return combined_mask.masked_fill(combined_mask > 0, float("-inf"))
```

The LFM2Model orchestrates the entire hybrid architecture. The key insight is processing through convolution blocks first for local patterns, then attention blocks for global dependencies.

## 7.3   Language Modeling Head

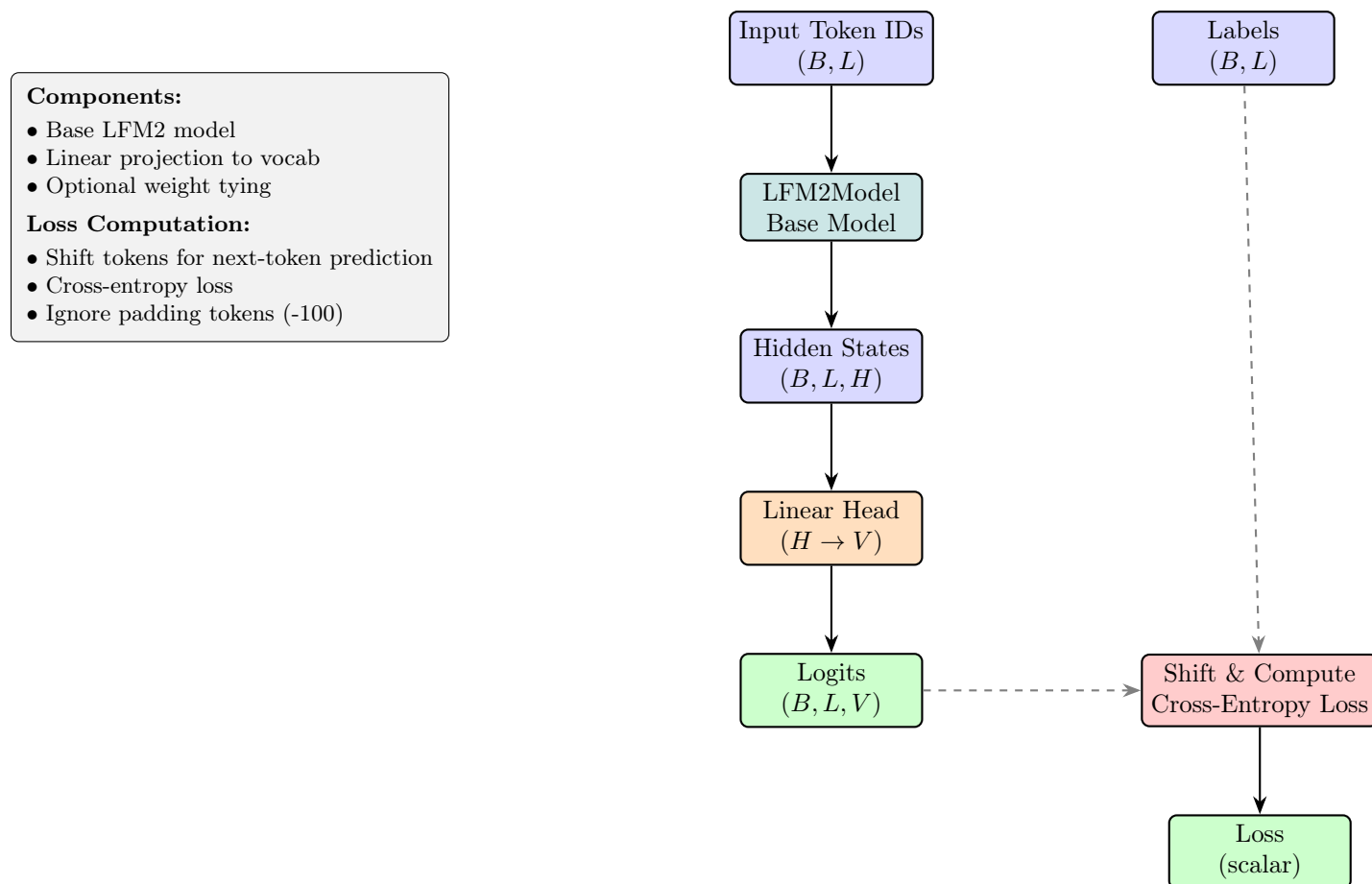**LFM2ForCausalLM - Language Modeling Head**



Figure 5: Language modeling head converts hidden states to vocabulary logits for text generation.

## 7.4   Attention Mask Preparation

```
1  class LFM2ForCausalLM(nn.Module):
2      """LFM2 model for causal language modeling (text generation)."""
3
4      def __init__(self, config: LFM2Config):
```

## Attention Mask Preparation
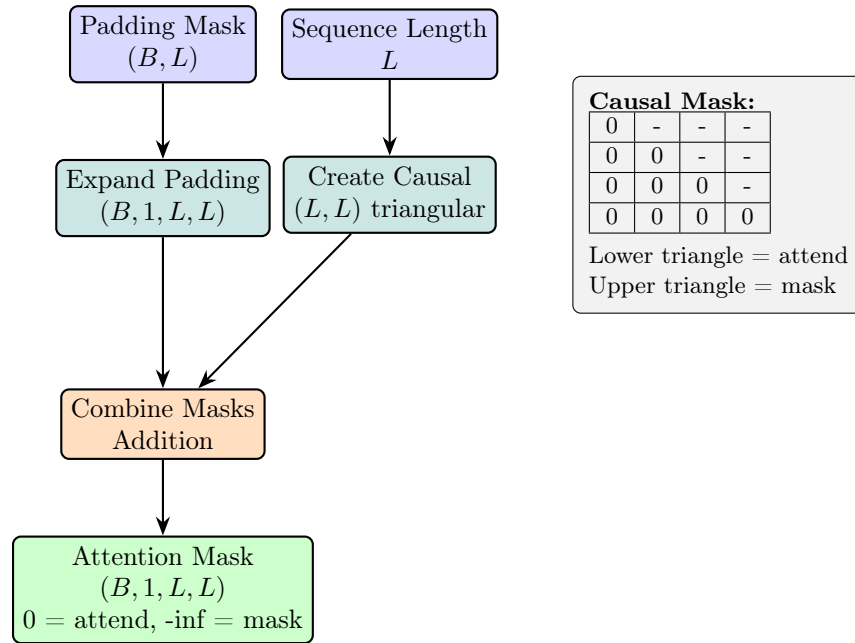


Figure 6: Attention mask combines padding mask with causal mask for autoregressive generation.

```
 5          super().__init__()
 6          self.config = config
 7          self.model = LFM2Model(config)
 8          self.vocab_size = config.vocab_size
 9
10          # Language modeling head: hidden states -> vocabulary logits
11          self.lm_head = nn.Linear(config.hidden_size, config.vocab_size, bias=False
                )
12
13          # Optionally tie input and output embeddings (parameter sharing)
14          if config.tie_word_embeddings:
15              self.lm_head.weight = self.model.embed_tokens.weight
16
17          self.apply(self._init_weights)
18
19      def forward(self, input_ids: torch.LongTensor,
20                  attention_mask: Optional[torch.Tensor] = None,
21                  labels: Optional[torch.LongTensor] = None,
22                  **kwargs) -> Union[Tuple, Dict[str, torch.Tensor]]:
23          """Forward pass for causal language modeling."""
24
25          # Forward through base model
26          outputs = self.model(
27              input_ids=input_ids,
28              attention_mask=attention_mask,
29              **kwargs
30          )
31
```

```
32            hidden_states = outputs["last_hidden_state"] if isinstance(outputs, dict)
                  else outputs[0]
33
34            # Compute logits for each position and vocabulary item
35            logits = self.lm_head(hidden_states)
36
37            # Compute loss if labels are provided (training)
38            loss = None
39            if labels is not None:
40                # Shift for next-token prediction: predict token i+1 from tokens 0..i
41                shift_logits = logits[..., :-1, :].contiguous()  # Remove last
                      position
42                shift_labels = labels[..., 1:].contiguous()      # Remove first
                      position
43
44                # Cross-entropy loss
45                loss_fct = nn.CrossEntropyLoss(ignore_index=-100)
46                shift_logits = shift_logits.view(-1, self.vocab_size)
47                shift_labels = shift_labels.view(-1)
48
49                # Move to same device for model parallelism
50                shift_labels = shift_labels.to(shift_logits.device)
51                loss = loss_fct(shift_logits, shift_labels)
52
53            # Return outputs
54            if isinstance(outputs, dict):
55                return {
56                    "loss": loss,
57                    "logits": logits,
58                    "past_key_values": outputs.get("past_key_values"),
59                    "hidden_states": outputs.get("hidden_states"),
60                }
61            else:
62                output = (logits,) + outputs[1:]
63                return (loss,) + output if loss is not None else output
```

## 7.5   Model Creation and Testing

```
1  def create_lfm2_model(model_size: str = "700M") -> LFM2ForCausalLM:
2      """Create LFM2 model with predefined configurations."""
3
4      # Predefined size configurations
5      size_configs = {
6          "350M": {
7              "hidden_size": 768,
8              "intermediate_size": 2048,
9              "num_attention_heads": 12,
10             "num_key_value_heads": 3,
11         },
12         "700M": {
13             "hidden_size": 1024,
14             "intermediate_size": 2816,
15             "num_attention_heads": 16,
16             "num_key_value_heads": 4,
17         },
18         "1.2B": {
```

```python
19              "hidden_size": 1536,
20              "intermediate_size": 4096,
21              "num_attention_heads": 24,
22              "num_key_value_heads": 6,
23          },
24      }
25
26      # Create and configure model
27      config = LFM2Config()
28      config.__dict__.update(size_configs[model_size])
29
30      return LFM2ForCausalLM(config)
31
32  def test_lfm2_model():
33      """Test LFM2 model with dummy inputs."""
34      print("Creating LFM2 model...")
35      model = create_lfm2_model("350M")
36      model.eval()
37
38      # Create dummy inputs
39      batch_size, seq_len = 2, 128
40      input_ids = torch.randint(0, model.config.vocab_size, (batch_size, seq_len))
41      attention_mask = torch.ones(batch_size, seq_len)
42
43      print(f"Input shape: {input_ids.shape}")
44      print(f"Model parameters: {sum(p.numel() for p in model.parameters()):,}")
45
46      # Forward pass
47      with torch.no_grad():
48          outputs = model(input_ids=input_ids, attention_mask=attention_mask)
49
50      print(f"Output logits shape: {outputs['logits'].shape}")
51      print("Test completed successfully!")
52
53  # Run the test
54  test_lfm2_model()
```

# 8 Conclusion

This step-by-step implementation of LFM2 demonstrates several key innovations:

- **Hybrid Architecture**: Convolution blocks for local patterns + attention blocks for global dependencies

- **Grouped Query Attention**: Reduces memory usage by sharing key-value pairs across query heads

- **Modern Components**: RMSNorm, RoPE, SwiGLU for improved training and performance

- **Efficient Design**: Double-gated convolutions and depthwise operations for computational efficiency

The implementation provides a complete, working LFM2 model suitable for language modeling tasks, with careful attention to memory efficiency and modern best practices.