

Building a Low-Latency Mixture-of-Experts Orchestrator for Domain-specific data

Anubhav 2022TT12161
Arnav Gupta 2022TT11352

Technical Report : [Sample Video](#) , [Repo](#)

1. Executive Summary

This project implements an intelligent multi-agent mental health support system designed to run entirely on mobile devices without internet connectivity. The system features a lightweight orchestration layer that intelligently routes user queries to specialized fine-tuned language models, each trained on specific mental health conditions (Anxiety, Bipolar Disorder, Depression, OCD, and Schizophrenia). Through a comprehensive pipeline of data generation, model fine-tuning, quantization, and mobile optimization, we achieved a fully functional on-device AI system that maintains privacy while providing specialized mental health support.

2. Intelligent Orchestration Architecture

2.1 Problem Statement

A single general-purpose language model often lacks depth in specialized domains. For mental health applications requiring nuanced understanding of distinct conditions, a multi-agent approach with intelligent routing provides superior performance compared to monolithic models.

2.2 Orchestrator Design

The orchestration system employs a classical machine learning approach optimized for efficiency and interpretability:

Architecture Components:

- **Feature Extraction:** TF-IDF (Term Frequency-Inverse Document Frequency) vectorization with 5,000 maximum features
- **N-gram Range:** Unigrams and bigrams (1,2) to capture both individual terms and contextual phrases
- **Classification:** Multinomial Logistic Regression with 1,000 training iterations
- **Target Classes:** Five mental health categories mapping to specialized models

Technical Specifications

```
python

vectorizer = TfidfVectorizer(
    max_features=5000,
    ngram_range=(1, 2),
    min_df=2,
    max_df=0.95
)

classifier = LogisticRegression(
    max_iter=1000,
    random_state=42,
    multi_class='multinomial'
)
```

2.3 Training Process

The orchestrator learns from the **23,200 Q&A pairs** used to train specialized models, creating a coherent routing mechanism:

1. **Data Loading:** Aggregates questions from all five mental health domain JSON files
2. **Stratified Splitting:** 80/20 train-test split maintaining class balance
3. **Vectorization:** Converts text into sparse TF-IDF feature vectors
4. **Training:** Fits multinomial logistic regression to learn classification boundaries
5. **Evaluation:** Achieved 87% accuracy

2.4 Inference Pipeline

When a user submits a query:

1. **Input Processing:** Text is transformed using the pre-fitted TF-IDF vectorizer
2. **Classification:** Logistic regression predicts the most relevant mental health category
3. **Probability Distribution:** Returns confidence scores for all five categories
4. **Routing Decision:** Directs query to the corresponding specialized model
5. **Response Generation:** Selected model generates contextually appropriate response

Performance Characteristics:

```
=====
TRAINING RESULTS
=====

Accuracy: 0.8688

Classification Report:
precision    recall    f1-score   support
anxiety      0.93     0.85      0.89      638
bipolar      0.85     0.90      0.87     1322
depression   0.82     0.83      0.82      939
ocd          0.89     0.92      0.91     1306
schiz        0.88     0.74      0.81      435

accuracy           0.87
macro avg       0.87     0.85      0.86     4640
weighted avg    0.87     0.87      0.87     4640

Model saved to mental_health_orchestrator.pkl
```

- Inference Latency: <10ms on CPU
- Model Size: ~50MB (serialized pickle file)
- Memory Footprint: Minimal during inference
- Accuracy: 87% routing accuracy

2.5 Advantages of Orchestration

This architecture provides several critical benefits:

- **Specialization:** Each model develops deep expertise in its domain
- **Efficiency:** Lightweight classifier routes queries without loading all models
- **Scalability:** New mental health domains can be added by training additional specialized models
- **Interpretability:** Probability distributions provide transparency in routing decisions
- **Resource Optimization:** Only one specialized model needs to be loaded at inference time

3. Training Data Generation with Llama 3.2

```
=====
1. HIGH-LEVEL MODEL STRUCTURE
=====

LlamaForCausalLM(
    (model): LlamaModel(
        (embed_tokens): Embedding(128256, 3072)
        (layers): ModuleList(
            (0-27): 28 x LlamaDecoderLayer(
                (self_attn): LlamaAttention(
                    (q_proj): Linear(in_features=3072, out_features=3072, bias=False)
                    (k_proj): Linear(in_features=3072, out_features=1024, bias=False)
                    (v_proj): Linear(in_features=3072, out_features=1024, bias=False)
                    (o_proj): Linear(in_features=3072, out_features=3072, bias=False)
                )
                (mlp): LlamaMLP(
                    (gate_proj): Linear(in_features=3072, out_features=8192, bias=False)
                    (up_proj): Linear(in_features=3072, out_features=8192, bias=False)
                    (down_proj): Linear(in_features=8192, out_features=3072, bias=False)
                    (act_fn): SiLU()
                )
                (input_layernorm): LlamaRMSNorm((3072,), eps=1e-05)
                (post_attention_layernorm): LlamaRMSNorm((3072,), eps=1e-05)
            )
        )
        (norm): LlamaRMSNorm((3072,), eps=1e-05)
        (rotary_emb): LlamaRotaryEmbedding()
    )
    (lm_head): Linear(in_features=3072, out_features=128256, bias=False)
)
```

3.1 Automated Q&A Pair Generation

High-quality training data is crucial for fine-tuning specialized models. We developed an automated pipeline using Meta's Llama-3.2-1B-Instruct to extract knowledge from mental health literature and generate conversational Q&A pairs.

3.2 Pipeline Architecture

Input Processing:

1. **PDF Extraction:** PyPDF2 library extracts text from mental health books
2. **Text Cleaning:** Regular expressions normalize whitespace and formatting
3. **Chunking Strategy:** Text divided into 400-word segments to maintain context while enabling parallel processing

Generation Configuration:

```
python

pairs_per_chunk = 20      # Q&A pairs generated per text segment
max_chunks = 100          # Number of chunks processed
chunk_size = 400           # Words per chunk
```

3.3 Llama 3.2 Generation Parameters

The generation process leverages Llama-3.2-1B-Instruct with carefully tuned parameters:

```
model = AutoModelForCausalLM.from_pretrained(
    "meta-llama/Llama-3.2-1B-Instruct",
    torch_dtype=torch.float16,
    device_map="auto"
)

outputs = model.generate(
    inputs,
    max_new_tokens=1200,
    temperature=0.7,
    top_p=0.9,
    do_sample=True
)
```

3.5 Output Statistics

For each mental health domain:

- **Target Generation:** ~2,000 Q&A pairs per book (100 chunks × 20 pairs)
- **Success Rate:** 80-90% successful JSON parsing
- **Quality Control:** Iterative generation with progress checkpoints every 5 chunks
- **Final Datasets:** 23K conversational Q&A pairs

Example Generated Pair:

```
{
  "question": "How can anxiety affect your relationships?",
  "answer": "Anxiety can make it difficult to connect with others, leading to feelings of isolation and social withdrawal."
},
{
  "question": "Can anxiety be a symptom of a underlying medical condition?",
  "answer": "Yes, anxiety can be a symptom of conditions such as thyroid disorders, chronic pain, or autoimmune diseases."
},
{
  "question": "What are some common triggers for anxiety?",
  "answer": "Common triggers include social situations, public speaking, and certain medications."
},
```

4. Fine-Tuning Gemma 3 with LoRA

```
=====
1. HIGH-LEVEL MODEL STRUCTURE
=====

Gemma3ForCausalLM(
    (model): Gemma3TextModel(
        (embed_tokens): Gemma3TextScaledWordEmbedding(262144, 1152, padding_idx=0)
        (layers): ModuleList(
            (0-25): 26 x Gemma3DecoderLayer(
                (self_attn): Gemma3Attention(
                    (q_proj): Linear(in_features=1152, out_features=1024, bias=False)
                    (k_proj): Linear(in_features=1152, out_features=256, bias=False)
                    (v_proj): Linear(in_features=1152, out_features=256, bias=False)
                    (o_proj): Linear(in_features=1024, out_features=1152, bias=False)
                    (q_norm): Gemma3RMSNorm((256,), eps=1e-06)
                    (k_norm): Gemma3RMSNorm((256,), eps=1e-06)
                )
                (mlp): Gemma3MLP(
                    (gate_proj): Linear(in_features=1152, out_features=6912, bias=False)
                    (up_proj): Linear(in_features=1152, out_features=6912, bias=False)
                    (down_proj): Linear(in_features=6912, out_features=1152, bias=False)
                    (act_fn): PytorchGELUTanh()
                )
                (input_layernorm): Gemma3RMSNorm((1152,), eps=1e-06)
                (post_attention_layernorm): Gemma3RMSNorm((1152,), eps=1e-06)
                (pre_feedforward_layernorm): Gemma3RMSNorm((1152,), eps=1e-06)
                (post_feedforward_layernorm): Gemma3RMSNorm((1152,), eps=1e-06)
            )
        )
        (norm): Gemma3RMSNorm((1152,), eps=1e-06)
        (rotary_emb): Gemma3RotaryEmbedding()
        (rotary_emb_local): Gemma3RotaryEmbedding()
    )
    (lm_head): Linear(in_features=1152, out_features=262144, bias=False)
)
```

4.1 Base Model Selection

Google's Gemma-3-1B-Instruct was selected for its optimal balance of:

- **Size** : 1 billion parameters suitable for mobile deployment
- **Performance** : Competitive with larger models on specialized tasks
- **Optimization** : Native bfloat16 support for efficient training
- **Licensing**: Permissive terms for commercial applications

4.2 LoRA (Low-Rank Adaptation) Methodology

Instead of full fine-tuning, we employ LoRA to reduce computational requirements while maintaining performance

$$W' = W_0 + BA$$

Where:

- W_0 = frozen pre-trained weights
- B, A = trainable low-rank matrices
- r = rank (dimensionality bottleneck)

Configuration:

```
# LoRA configuration optimized for Gemma
peft_config = LoraConfig(
    r=16,
    lora_alpha=32,
    lora_dropout=0.05,
    bias="none",
    task_type="CAUSAL_LM",
    target_modules=["q_proj", "k_proj", "v_proj", "o_proj", "gate_proj", "up_proj", "down_proj"]
)

model = get_peft_model(model, peft_config)
```

Target Modules: LoRA adapters applied to all attention projections (q, k, v, o) and feed-forward network layers (gate, up, down) for comprehensive adaptation.

4.3 Memory-Efficient Training

4-bit Quantization During Training:

```
# Quantization config for memory efficiency
bnb_config = BitsAndBytesConfig(
    load_in_4bit=True,
    bnb_4bit_quant_type="nf4",
    bnb_4bit_compute_dtype=torch.bfloat16, # Gemma uses bfloat16
    bnb_4bit_use_double_quant=True,
)

# Load model and tokenizer
print("Loading Gemma 3 1B model and tokenizer...")
model = AutoModelForCausalLM.from_pretrained(
    model_name,
    quantization_config=bnb_config,
    device_map="auto",
    trust_remote_code=True
)
```

This reduces memory footprint from ~8GB to ~2GB, enabling training on consumer GPUs.

4.4 Training Configuration

```
# Training arguments
training_args = TrainingArguments(
    output_dir=output_dir,
    num_train_epochs=3, # Reduced for Gemma (learns faster)
    per_device_train_batch_size=4, # Smaller batch for stability
    gradient_accumulation_steps=8, # Compensate for smaller batch
    learning_rate=1e-4, # Lower LR for Gemma
    bf16=True, # Use bfloat16 for Gemma
    save_steps=100,
    logging_steps=10,
    save_total_limit=2,
    warmup_steps=50,
    optim="paged_adamw_8bit",
    remove_unused_columns=False,
    report_to="none",
)
```

Training Duration: 2-3 hours per model on **Kaggle 2 X T4 GPU**

4.5 Model Merging for Deployment

LoRA produces lightweight adapter files (~50-100MB), but deployment requires a standalone model in task format.

5. Challenges of Large Language Models

5.1 Computational Constraints

Large language models present significant deployment challenges:

- **Memory Requirements:** Modern LLMs require substantial RAM for inference. A 7B parameter model in FP16 requires ~14GB of memory just for weights, making deployment on resource-constrained devices impractical.
- **Inference Latency:** Autoregressive generation requires sequential token production, with latency proportional to model size. On mobile devices, unoptimized models generate <1 token/second, creating poor user experience.
- **Energy Consumption:** Each inference passes through billions of parameters drains battery rapidly, limiting practical mobile usage to minutes rather than hours.

5.2 Storage Limitations

- **On-Device Storage:** Mobile devices have finite storage, and users are reluctant to allocate multiple gigabytes to single applications. A 7B parameter model in FP16 format occupies ~14GB, exceeding available storage on many devices.

- **Model Distribution:** Downloading multi-gigabyte models over cellular networks is impractical, limiting accessibility and increasing user friction during installation.

5.3 Privacy and Connectivity

While large models offer impressive capabilities, cloud-based inference raises privacy concerns for sensitive mental health conversations. Edge deployment requires dramatic model compression to maintain functionality while respecting user privacy.

This project addresses these challenges through aggressive quantization and **format optimization, reducing model size by 72.5 % (2GB → 550MB) while maintaining response quality.**

6. Quantization and Mobile Optimization

6.1 Fundamentals of Quantization

Quantization reduces model precision to decrease memory footprint and accelerate inference:

Mathematical Foundation: Convert high-precision floating-point values to lower-precision representations:

$$Q(x) = \text{round}(x / \text{scale}) + \text{zero_point}$$

Where:

- scale: Mapping between float and integer ranges
- zero_point: Offset for asymmetric distributions

Precision Comparison:

- FP32: 32-bit floating point (4 bytes per parameter)
- FP16/BF16: 16-bit formats (2 bytes per parameter)
- INT8: 8-bit integer (1 byte per parameter)
- INT4: 4-bit integer (0.5 bytes per parameter)

6.2 Dynamic INT4 Quantization

We employ dynamic int4 quantization with block-wise scaling for optimal compression:

Block-128 Quantization: Weights divided into 128-element blocks, each with independent scaling factors:

`quantize="dynamic_int4_block128"`

Advantages:

- Granular Precision: Smaller blocks preserve outlier values better than global quantization
- Reduced Error: Per-block scaling minimizes quantization error accumulation
- Hardware Efficiency: 128-element blocks align with mobile GPU architectures

6.3 Conversion Pipeline

Step 1: PyTorch to TensorFlow Lite

Using Google's AI Edge Torch converter:

```
!python3 convert_gemma3_to_tflite.py \
--quantize="dynamic_int4_block128" \
--checkpoint_path="/mnt/c/Users/User/Downloads/results/gemma-anxity-merged/model" \
--output_path="/mnt/c/Users/User/Downloads/results/gemma-anxity-merged/model/gemma3.tflite" \
--prefill_seq_lens=1024 \
--kv_cache_max_len=2048 \
--mask_as_input=True
```

Key Parameters:

- prefill_seq_len: Maximum input token length (1024 for extended context)
- kv_cache_max_len: Key-value cache size for efficient autoregressive generation
- mask_as_input: Attention mask handling for variable-length sequences

Step 2: MediaPipe Task Bundling

TFLite models require tokenizer bundling for mobile deployment:

```
tflite_model="PATH/gemma.tflite" # @param {type:"string"}  
tokenizer_model="PATH/tokenizer.model" # @param {type:"string"}  
start_token=<bos> # @param {type:"string"}  
stop_token=<eos> # @param {type:"string"}  
output_filename="PATH/gemma.task" # @param {type:"string"}  
enable_bytes_to_unicode_mapping=False # @param ["False", "True"] {type:"raw"}  
  
config = bundler.BundleConfig(  
    tflite_model=tflite_model,  
    tokenizer_model=tokenizer_model,  
    start_token=start_token,  
    stop_tokens=[stop_token],  
    output_filename=output_filename,  
    enable_bytes_to_unicode_mapping=enable_bytes_to_unicode_mapping,  
)  
bundler.create_bundle(config)
```

6.4 Size Reduction Results

Compression Analysis:

Stage	Size	Reduction
Merged Model (BF16)	~2 GB	Baseline
TFLite (INT4 block128)	~490 MB	75% ↓
Final .task Bundle	~550MB	72.5% ↓

This dramatic compression makes on-device deployment feasible while maintaining response quality through careful quantization strategy.

9. Conclusion

This project demonstrates a complete pipeline for deploying specialized AI agents on resource-constrained mobile devices. Through intelligent orchestration, automated data generation, efficient fine-tuning with LoRA, and aggressive quantization, we created a privacy-preserving mental health support system that operates entirely offline.

Key Innovations:

1. Lightweight orchestration using classical ML for intelligent routing
2. Automated Q&A generation from domain literature using Llama 3.2
3. Efficient fine-tuning with LoRA reducing training costs by 90%
4. Extreme compression through INT4 quantization achieving 97% size reduction
5. Mobile-optimized deployment using Google AI Edge infrastructure

The resulting system provides specialized mental health support while respecting user privacy and device constraints, making advanced AI assistance accessible without cloud dependency.

10. References

- Meta AI: Llama 3.2 Model Family
- Google: Gemma 3 Language Models
- HuggingFace: PEFT and Transformers Libraries
- Google AI Edge Torch: Mobile Conversion Tools
- MediaPipe: On-Device ML Framework