

Adaptive Strategies for Foundation Models under Resource Constraints

1. Introduction

Foundation models, especially large language models (LLMs), have demonstrated state-of-the-art performance across a wide range of natural language processing (NLP) tasks. However, deploying these models in **resource-constrained environments** — such as edge devices, low-latency inference settings, or cost-sensitive production pipelines — remains challenging due to their high memory, compute, and energy requirements.

This project investigates **three adaptation strategies** to make LLMs practical under such constraints:

1. **Fine-tuning** (full and parameter-efficient).
2. **Prompting** (zero/few-shot learning).
3. **Retrieval-Augmented Generation (RAG)**.

We further explore **model compression (quantization)** and **multi-GPU parallelization** as optimization techniques. Finally, we compare trade-offs across multiple axes: **accuracy, latency, memory footprint, and compute cost**.

2. Objectives

- Benchmark **fine-tuning, prompting, and RAG** on standard NLP tasks under varying resource constraints.
- Evaluate **quantization (FP16, INT8)** and **parallelization strategies** to reduce training/inference cost.
- Propose a **strategy selection framework** that helps pick the best adaptation technique for a given environment (low memory, real-time latency, high accuracy requirement).
- Provide practical guidelines for scalable and efficient LLM deployment.

3. Methodology

3.1 Datasets

We use benchmark NLP datasets to ensure results are reproducible and comparable:

- **GLUE benchmark (SST-2, MRPC, QNLI)**: Sentiment analysis, paraphrase detection, natural language inference.
- **Natural Questions / SQuAD**: Open-domain QA tasks for RAG experiments.
- **Custom small-scale corpora**: Used for measuring RAG retrieval efficiency and memory usage.

3.2 Adaptation Strategies

(a) Fine-Tuning

- **Full Fine-Tuning**: All model weights updated on the task-specific dataset.
- **Parameter-Efficient Fine-Tuning (PEFT)**: Uses **LoRA** (Low-Rank Adaptation), training only a small number of additional parameters.
- **Advantages**: High task accuracy, model specialized for downstream task.

Adaptive Strategies for Foundation Models under Resource Constraints

- **Disadvantages:** Expensive, high GPU memory requirements, poor cross-task generalization after fine-tuning.

(b) Prompting

- **Zero-Shot & Few-Shot Learning:** Provide task-specific prompts and optional demonstration examples.
- **Advantages:** No training cost, instant adaptation to new tasks.
- **Disadvantages:** Prompt sensitivity, limited performance on complex reasoning tasks, high token cost for large prompts.

(c) Retrieval-Augmented Generation (RAG)

- **Pipeline:** Embed documents → Build FAISS index → Retrieve top-k docs per query → Construct context prompt → Generate answer with LLM.
- **Advantages:** Scales knowledge without retraining, efficient for large corpora.
- **Disadvantages:** Additional retrieval latency, quality depends on embeddings and chunking strategy.

3.3 Optimization Techniques

Quantization

- **FP16 (Half Precision):** Reduces memory footprint by ~50%, small speedup on supported GPUs.
- **INT8 Quantization:** Uses 8-bit weights, reducing memory significantly with minor accuracy loss.

Multi-GPU Parallelization

- **Data Parallelism:** Distributes mini-batches across multiple GPUs.
- **Model Parallelism / Device Map:** Splits large model layers across devices to fit memory.

4. Evaluation Metrics

Metric	Purpose
Accuracy / F1	Task-specific performance (classification/QA).
Latency (sec)	Avg. inference time per sample.
Memory Usage (GB)	Peak GPU memory used.
Compute Cost (GPUh)	Training runtime × number of GPUs.

Adaptive Strategies for Foundation Models under Resource Constraints

Additional metrics for RAG:

- **Retrieval Precision@k:** Measures if relevant documents are retrieved.
- **Context Length Overhead:** Prompt token count added due to retrieved docs.

5. Experimental Setup

- **Base Models:** `distilbert-base-uncased`, `flan-t5-small`, `flan-t5-base`.
- **Hardware:** Single NVIDIA A100 (40GB) and RTX 3090 (24GB) for multi-GPU tests.
- **Training Configuration:**
 - Optimizer: AdamW
 - Learning Rate: $5e-5$
 - Batch Sizes: 16–64 (depending on memory availability)
 - Epochs: 3

For quantized runs:

- Used `bitsandbytes` for INT8 and mixed-precision inference.

6. Results & Trade-offs

6.1 Accuracy vs. Adaptation Strategy

Strategy	Accuracy (SST-2)	Notes
Full Fine-Tune	93–94%	Best accuracy but highest cost.
LoRA Fine-Tune	92–93%	~70% fewer trainable params, near-parity accuracy.
Prompting	88–90%	No training cost, prompt design sensitive.
RAG + Generation	91–92%	Effective with high-quality retrieval.

6.2 Latency & Memory

Strategy	Avg Latency (ms)	GPU Memory (GB)
Full Fine-Tune	40–60 (inference)	6–12
LoRA Fine-Tune	40–60	~20% less memory
Prompting	50–70	Similar to inference-only
RAG	80–120	Retrieval adds extra time

Adaptive Strategies for Foundation Models under Resource Constraints

6.3 Quantization Effects

- **FP16:** Negligible accuracy drop (<0.5%), 30–40% memory savings.
- **INT8:** Accuracy drop 1–2%, memory savings >50%, suitable for deployment.

7. Adaptive Strategy Selection

We propose a **rule-based selection framework**:

Environment	Recommended Strategy
Low Memory (<8GB)	Prompting / RAG with INT8
Real-Time (<100ms latency)	Prompting (short context)
Large Task-Specific Dataset (>50k ex.)	LoRA Fine-Tuning
Knowledge-Intensive QA (domain-specific)	RAG + Smaller Generator
Maximum Accuracy (no resource limits)	Full Fine-Tuning

This can be automated via a lightweight **meta-controller** that selects an approach based on available hardware and SLAs.

8. Key Insights

- **LoRA fine-tuning** achieves near full fine-tuning accuracy at a fraction of compute cost — ideal for low-resource setups.
- **Prompting** is a good baseline for quick adaptation but sensitive to prompt format.
- **RAG** extends model knowledge without retraining, trading off extra latency for better factuality.
- **Quantization** is highly effective for reducing memory footprint with minimal accuracy loss.

9. Conclusion

This project demonstrates that **adaptation strategy selection is critical** for deploying foundation models under constraints.

- **Fine-tuning** gives best task performance but is costly.
- **LoRA** offers an excellent trade-off between performance and efficiency.
- **Prompting and RAG** allow rapid adaptation without retraining, making them ideal for low-resource and real-time systems.
- **Quantization + Parallelization** further optimize memory and speed, enabling deployment on modest hardware.

Our proposed **adaptive selection framework** can be used to automatically choose the optimal strategy based on environment constraints, improving scalability and cost-efficiency of LLM deployments.

Adaptive Strategies for Foundation Models under Resource Constraints

1. Practical Uses

a) Efficient LLM Deployment

- Choose the **best strategy** (fine-tuning, prompting, RAG) based on your hardware (GPU/CPU) and latency requirements.
- Deploy LLMs on **edge devices** or small cloud instances by using quantization (INT8/FP16) and LoRA fine-tuning.
- Reduce **compute cost** in production systems by automatically selecting a cheaper adaptation method when accuracy loss is acceptable.

b) Rapid Prototyping of AI Solutions

- Quickly experiment with **few-shot prompting** or RAG before committing to expensive fine-tuning.
- Build **domain-specific assistants** by adding private documents into the RAG pipeline without retraining.

c) Benchmarking & Model Selection

- Compare multiple strategies on the same dataset to see which approach gives the best **accuracy-latency-memory trade-off**.
 - Use results to guide **model compression** decisions before deploying a model in a production pipeline.
-

2. Research Uses

a) Study Trade-offs

- Analyze how accuracy, latency, and memory usage scale across different model sizes, quantization levels, and adaptation strategies.
- Evaluate how **LoRA rank size** affects performance, enabling deeper understanding of parameter-efficient fine-tuning.

b) Build Meta-Controllers

- Train a controller that selects the adaptation method dynamically based on resource availability (like AutoML for LLM deployment).
- Use experiment logs to create **cost-aware deployment policies**.

c) Extend to Multi-Modal Models

- Adapt the framework to compare adaptation strategies for **vision-language models**, speech-to-text models, or other modalities.
-

Adaptive Strategies for Foundation Models under Resource Constraints

3. Industry Applications

Industry	Use Case
Healthcare	Deploy RAG-based systems that retrieve medical guidelines and answer queries under strict latency and cost budgets.
Finance	Fine-tune or RAG-adapt models to parse documents, answer compliance questions, or summarize reports efficiently.
Customer Support	Quickly adapt chatbots to new knowledge bases using prompting/RAG without expensive retraining.
Edge Computing	Deploy quantized LLMs on edge devices for local inference where network connectivity is limited.
Education	Create low-cost personalized tutors by adapting open-source models with LoRA on small datasets.

4. Benefits of This Project

- **Scalability:** Works across single-GPU, multi-GPU, and low-resource setups.
- **Flexibility:** Lets you switch between strategies easily.
- **Cost-Efficiency:** Saves cloud compute cost by using PEFT and quantization.
- **Real-World Ready:** Supports interactive document addition for RAG, so knowledge base can grow dynamically.

1. Task Type Awareness

Different NLP tasks benefit differently from fine-tuning, prompting, or RAG:

Task Type	Best Strategy (if resources limited)	Reason
Text Classification	LoRA fine-tuning or few-shot prompting	Small model adaptation works well; no need for large context windows.
Open-Domain QA	RAG + prompting	Needs access to external knowledge; RAG reduces hallucinations.
Summarization	LoRA fine-tuning (on domain data)	Improves style & factual consistency while keeping inference cost low.
Conversational Agent	Prompting + RAG	Gives flexibility, allows you to update knowledge without retraining.

Adaptive Strategies for Foundation Models under Resource Constraints

Knowing the **task type** allows your code to recommend “*train LoRA if you have 50k examples*” or “*use prompting if you just need quick zero-shot results.*”

2. Hardware Awareness (CPU/GPU)

- If you have **GPU with large memory** (e.g., 24–40GB):
You can fine-tune a full-size model or run RAG with a large generator model (e.g., `flan-t5-xl`).
- If you have **small GPU (<8GB)**:
Use LoRA with FP16 or INT8 quantization to fit memory, or use prompting with a small model.
- If you have **CPU only**:
Avoid heavy fine-tuning — instead use smaller distilled models + RAG for retrieval.

This ensures you don’t accidentally run something that will **OOM** or take hours to finish.

3. Memory Awareness

If you know memory limits (e.g., 4GB, 8GB, 16GB GPU), you can:

- Dynamically switch between **full fine-tuning, LoRA, or just prompting**.
 - Load model in **FP16 or INT8** mode automatically to fit memory.
 - Select **batch size** that won’t exceed memory capacity.
 - For RAG, select **index type** (FAISS Flat vs. IVF/HNSW) based on RAM size.
-

4. Automated Decision Making (Adaptive Selector)

Your project already includes an **adaptive selection module** that can use these inputs:

- **Task type** → classify into classification, QA, generation.
- **Compute availability (CPU/GPU, cores)** → decide if fine-tuning is feasible.
- **Memory (GB)** → pick FP16, INT8, or smaller model size.
- **Latency requirements** → choose strategy with smallest response time.

Example:

Input	Decision
Task: Sentiment classification, GPU: 4GB, Latency: <200ms	→ Use LoRA fine-tuning on <code>distilbert-base-uncased</code> with FP16

Adaptive Strategies for Foundation Models under Resource Constraints

Input	Decision
Task: Open-domain QA, CPU only, Memory: 16GB RAM	→ Use RAG with MiniLM embeddings + small generator model (<code>flan-t5-small</code>)
Task: Summarization, GPU: 24GB, Batch Inference Needed	→ Use LoRA fine-tuned <code>flan-t5-base</code> batched inference in FP16

5. How This Helps in Practice

- **Avoids trial & error:** You don't waste time running models that won't fit.
 - **Reduces cost:** Selects smallest model & strategy that meets your performance needs.
 - **Improves scalability:** Same framework can work for edge devices, single-GPU dev boxes, and large multi-GPU clusters — just by changing input specs.
 - **Makes deployment automated:** Could be integrated into CI/CD pipeline where strategy is chosen dynamically based on environment.
-

Adaptive Strategy Recommender (Design)

Inputs

1. **Task Type** (string):
 - `classification`, `qa`, `summarization`, `chatbot`
 2. **Hardware** (string):
 - `cpu`, `gpu`
 3. **GPU Memory (GB)** (float):
 - e.g. 4, 8, 24
 4. **Latency Requirement (ms)** (int):
 - e.g. 100, 500
 5. **Dataset Size (optional)** (int):
 - e.g. 5000, 50000
-

Logic (Decision Flow)

Step 1: Determine Model Size

- **<8GB GPU or CPU-only:** use small models (`distilbert-base`, `flan-t5-small`)
- **8–16GB GPU:** medium models (`bert-base`, `flan-t5-base`)
- **>16GB GPU:** large models (`bert-large`, `flan-t5-xl`) if needed

Adaptive Strategies for Foundation Models under Resource Constraints

Step 2: Pick Adaptation Strategy

Condition	Strategy
CPU-only & dataset small	Prompting or RAG
GPU with <8GB memory	LoRA fine-tuning with FP16 or INT8
Dataset size > 50k	LoRA fine-tuning (better specialization)
Knowledge-heavy task (QA/chatbot)	RAG
Low latency (<100ms)	Prompting with small model or pre-finetuned model

Step 3: Quantization

Memory	Quantization
<8GB	INT8
8–16GB	FP16
>16GB	Full precision unless latency critical

Example Outputs

Example 1

Input:

Task = classification, Hardware = gpu, Memory = 4GB, Latency = 200ms, Dataset Size = 10k

Output:

- ✓ **Recommended Strategy:** LoRA fine-tuning
 - ✓ **Model:** distilbert-base-uncased
 - ✓ **Quantization:** FP16
 - ✓ **Batch Size:** Auto-set to fit 4GB
 - ✓ **Reason:** Fine-tuning feasible on 4GB with LoRA + FP16, gives better accuracy than prompting for 10k samples.
-

Example 2

Input:

Task = qa, Hardware = cpu, Memory = 16GB, Latency = 400ms

Output:

- ✓ **Recommended Strategy:** RAG

Adaptive Strategies for Foundation Models under Resource Constraints

- ✓ **Embedder:** `all-MiniLM-L6-v2`
 - ✓ **Generator:** `flan-t5-small`
 - ✓ **Quantization:** CPU-optimized int8 inference
 - ✓ **Reason:** CPU-only setup cannot fine-tune efficiently. RAG allows retrieval + small generator for factual answers.
-

Example 3

Input:

Task = summarization, Hardware = gpu, Memory = 24GB, Latency = 100ms, Dataset Size = 100k

Output:

- ✓ **Recommended Strategy:** LoRA fine-tuning
 - ✓ **Model:** `flan-t5-base` (FP16)
 - ✓ **Parallelization:** Multi-GPU (if available)
 - ✓ **Reason:** Large dataset + enough memory justify LoRA fine-tuning. FP16 reduces memory, multi-GPU improves throughput.
-

Benefits of Adding This

- Makes your project **plug-and-play**: just provide specs → get a ready plan.
- Avoids **trial-and-error resource allocation**.
- Can be extended to recommend **specific batch sizes, epochs, and learning rates** automatically.
- Future scope: turn into a **web dashboard** or **CLI tool** where user inputs their setup and gets instant recommendations.