

Model Quantization and Fine-tuning Platform

Summary: Develop a platform that enables efficient model quantization and fine-tuning for deploying large language models on resource-constrained environments. **Problem Statement:** Large language models require significant computational resources, limiting their deployment in edge environments. Your task is to create a platform that automates model quantization, fine-tuning, and optimization for specific use cases while maintaining performance quality. The system should support various quantization techniques, provide performance benchmarking, and enable easy deployment to different hardware configurations. **Steps:** • Design automated quantization pipeline supporting multiple techniques (INT8, INT4, dynamic) • Implement fine-tuning workflows with parameter-efficient methods (LoRA, QLoRA) • Create performance benchmarking suite measuring accuracy, speed, and memory usage • Build deployment optimization for different hardware targets (CPU, GPU, mobile) • Develop model comparison and selection tools based on constraints • Include monitoring and quality assessment for quantized models **Suggested Data Requirements:** • Pre-trained model checkpoints and configuration files • Domain-specific fine-tuning datasets • Hardware performance benchmarks and constraints • Quality evaluation datasets for model comparison **Themes:** GenAI & its techniques, Quantization, Fine-tuning The steps and data requirements outlined above are intended solely as reference points to assist you in conceptualising your solution. **PRD (Product Requirements Document)** Product Vision and Goals To democratize LLM deployment on edge devices by automating optimization, reducing model size by 4x-8x while retaining >95% accuracy. **Goals:** Support 10+ quantization methods, integrate with 5 hardware types, and provide one-click deployment. **Target Audience and Stakeholders**

Primary Users: ML engineers, mobile app developers, IoT specialists. **Stakeholders:** Hardware vendors for benchmarks, end-users for inference. **Personas:** An edge AI developer optimizing GPT-J for Raspberry Pi.

Key Features and Functionality

Auto-quantization with technique selection. PEFT (Parameter-Efficient Fine-Tuning) workflows. Multi-metric benchmarking dashboard. Hardware-specific exporters (e.g., TFLite for mobile). Model selector with constraint-based ranking. Post-deployment monitoring for drift.

Business Requirements

Open-source with enterprise edition for cloud integration. Integration with Hugging Face Hub for model loading.

Success Metrics

Efficiency: >2x speed-up on target hardware. **User Adoption:** 1000+ downloads in first year. **Quality:** Perplexity <5% increase post-quantization.

Assumptions, Risks, and Dependencies

Assumptions: Users have basic PyTorch knowledge. **Risks:** Accuracy loss in quantization; mitigate with calibration datasets. **Dependencies:** Libraries like bitsandbytes for QLoRA, public models from HF.

Out of Scope

Custom hardware acceleration (e.g., FPGA design). Online learning during inference.

FRD (Functional Requirements Document) System Modules and Requirements

Quantization Pipeline (FR-001):

Input: Model checkpoint, calibration data. **Functionality:** Support PTQ (Post-Training Quant), QAT; techniques: static INT8, dynamic, FP16. **Output:** Quantized model with config.

Fine-Tuning Workflow (FR-002):

Input: Quantized model, dataset. **Functionality:** Apply LoRA/QLoRA; trainers with PEFT library. **Output:** Adapted model adapters.

Benchmarking Suite (FR-003):

High-Level Flow:

Config input -> Run pipeline stages sequentially or parallel. Store artifacts in registry.

LLD (Low Level Design)

Quantization LLD:

Static INT8: `model = torch.quantization.quantize(model, qconfig_spec, inplace=False)` Calibration: Run forward passes on 1000 samples.

Fine-Tuning LLD:

LoRA Config: `from peft import LoraConfig; config = LoraConfig(r=16, lora_alpha=32)` Trainer: `from transformers import Trainer; trainer.train()`

Benchmark LLD:

Accuracy: `from evaluate import load; acc = load("accuracy").compute(preds, refs)` Latency: with `torch.profiler.profile(): model(input); print(profile.key_averages())`

Comparison LLD:

Pareto: Use `scipy.optimize` for multi-objective ranking.

Pseudocode textclass QuantFinePlatform: `def init(self): self.hf_hub = HFHub()`

```
def quantize(self, model_name, tech='int8', calib_data):
    model = self.hf_hub.load_model(model_name)
    if tech == 'int8':
        q_model = torch.quantization.quantize_dynamic(model, {nn.Linear: torch.qint8})
    elif tech == 'int4':
        q_model = bitsandbytes.quantize(model, 4)
    q_model.calibrate(calib_data)
    return q_model
```

```
def fine_tune(self, q_model, dataset, method='qlora'):
    config = LoraConfig(...) if method == 'lora' else QLoRAConfig(...)
    peft_model = get_peft_model(q_model, config)
    trainer = Trainer(peft_model, train_dataset=dataset, eval_dataset=val)
    trainer.train()
    return peft_model
```

```
def benchmark(self, models, eval_data, hardware='cpu'):
    results = []
    for m in models:
        acc = evaluate_model(m, eval_data)
        lat, mem = profile_inference(m, hardware)
        results.append({'acc': acc, 'lat': lat, 'mem': mem})
    return results
```

```
def deploy(self, model, target='mobile'):
    if target == 'mobile':
        converted = convert_to_tflite(model)
    return converted
```

```
def compare(self, benchmarks, constraints):
    filtered = [b for b in benchmarks if b['mem'] < constraints['max_mem']]
    ranked = sort_by_pareto(filtered, keys=['acc', '-lat'])
    return ranked[0]
```