

Multi-Scale Trade-off Analysis

Metric	Cloud Server	Edge Device	Microcontroller
Model Size	1.24 MB	0.41 MB	438.2 KB
Accuracy	83.5%	68.1%	70.8%
Latency	12.4 ms	1.17 ms	1.13 ms
Memory Usage	1.25 MB	0.42 MB	448.2 KB
Power Estimate	X.X W	X.X mW	X.X mW
Development Complexity	Medium	Medium	High

Optimization Strategy Effectiveness

- Mixed Precision
 - o Training Time: 116.2 seconds for 10 epoch
 - o Accuracy: 82.5%

Model Start training using the best performance model in part1, maybe this is why it overpasses the baseline model accuracy (81.85%). However, it also gives a much longer training time, which is unexpected.

- Distributed training scaling efficiency
 - o OneDeviceStrategy: 139.1 seconds
 - o MirroredStrategy: 174.6 seconds
 - o Scaling Efficiency: 0.796
 - o Maybe mirrored strategy is not optimized for single GPU, resulting in lower scaling efficiency.

- Batch Processing optimization benefits

```
=== GRADIENT ACCUMULATION THROUGHPUT ANALYSIS ===
Testing throughput for batch size: 64
  Accumulation steps: 1
  Results: 2613.29 samples/sec, 40.83 batches/sec
Testing throughput for batch size: 128
  Accumulation steps: 2
  Results: 2338.89 samples/sec, 18.27 batches/sec
Testing throughput for batch size: 256
  Accumulation steps: 4
  Results: 2837.72 samples/sec, 11.08 batches/sec
Testing throughput for batch size: 512
  Accumulation steps: 8
  Results: 3738.11 samples/sec, 7.30 batches/sec
Testing throughput for batch size: 1024
  Accumulation steps: 16
  Results: 4186.59 samples/sec, 4.09 batches/sec
Optimal batch size for throughput: 1024
Best throughput: 4186.59 samples/sec
Testing throughput for batch size: 1024
  Accumulation steps: 16
  Results: 4006.68 samples/sec, 3.91 batches/sec
```

With greater batch, the batch processing time becomes much longer, but the inference time between samples becomes shorter.

- Knowledge distillation effectiveness

```
STREAMLINED ANALYSIS: part2_cloud_optimization/student_model.keras

ARCHITECTURE & MEMORY ANALYSIS
-----
Total Parameters: 82,586
Trainable Parameters: 82,138
Layer conv2d: 448 params, dtype=float32, policy=<Policy "float32">, size=1.75 KB
Layer batch_normalization: 64 params, dtype=float32, policy=<Policy "float32">, size=0.25 KB
Layer conv2d_1: 2,320 params, dtype=float32, policy=<Policy "float32">, size=9.06 KB
Layer batch_normalization_1: 64 params, dtype=float32, policy=<Policy "float32">, size=0.25 KB
Layer conv2d_2: 4,640 params, dtype=float32, policy=<Policy "float32">, size=18.12 KB
Layer batch_normalization_2: 128 params, dtype=float32, policy=<Policy "float32">, size=0.50 KB
Layer conv2d_3: 9,248 params, dtype=float32, policy=<Policy "float32">, size=36.12 KB
Layer batch_normalization_3: 128 params, dtype=float32, policy=<Policy "float32">, size=0.50 KB
Layer conv2d_4: 18,496 params, dtype=float32, policy=<Policy "float32">, size=72.25 KB
Layer batch_normalization_4: 256 params, dtype=float32, policy=<Policy "float32">, size=1.00 KB
Layer conv2d_5: 36,928 params, dtype=float32, policy=<Policy "float32">, size=144.25 KB
Layer batch_normalization_5: 256 params, dtype=float32, policy=<Policy "float32">, size=1.00 KB
Layer dense_1: 1,290 params, dtype=float32, policy=<Policy "float32">, size=5.04 KB
Model File Size: 0.39 MB
Theoretical Memory: 0.32 MB
Training Memory: 2.01 MB
Inference Memory: 0.33 MB
Memory footprint: 2.34 MB

PERFORMANCE & SPEED ANALYSIS
-----
Test Accuracy: 73.87%
Test Loss: 0.9342
Single Sample Time: 2.37 ± 19.09 ms
GPU Available: Yes
GPU Memory Growth: True
```

```
STREAMLINED ANALYSIS: part2_cloud_optimization/teacher_model.keras

ARCHITECTURE & MEMORY ANALYSIS
-----
Total Parameters: 324,394
Trainable Parameters: 323,498
Layer conv2d: 896 params, dtype=float32, policy=<Policy "float32">, size=3.50 KB
Layer batch_normalization: 128 params, dtype=float32, policy=<Policy "float32">, size=0.50 KB
Layer conv2d_1: 9,248 params, dtype=float32, policy=<Policy "float32">, size=36.12 KB
Layer batch_normalization_1: 128 params, dtype=float32, policy=<Policy "float32">, size=0.50 KB
Layer conv2d_2: 18,496 params, dtype=float32, policy=<Policy "float32">, size=72.25 KB
Layer batch_normalization_2: 256 params, dtype=float32, policy=<Policy "float32">, size=1.00 KB
Layer conv2d_3: 36,928 params, dtype=float32, policy=<Policy "float32">, size=144.25 KB
Layer batch_normalization_3: 256 params, dtype=float32, policy=<Policy "float32">, size=1.00 KB
Layer conv2d_4: 73,856 params, dtype=float32, policy=<Policy "float32">, size=288.50 KB
Layer batch_normalization_4: 512 params, dtype=float32, policy=<Policy "float32">, size=2.00 KB
Layer conv2d_5: 147,584 params, dtype=float32, policy=<Policy "float32">, size=576.50 KB
Layer batch_normalization_5: 512 params, dtype=float32, policy=<Policy "float32">, size=2.00 KB
Layer dense_1: 33,024 params, dtype=float32, policy=<Policy "float32">, size=129.00 KB
Layer dense_1: 2,570 params, dtype=float32, policy=<Policy "float32">, size=10.04 KB
Model File Size: 1.32 MB
Theoretical Memory: 1.24 MB
Training Memory: 5.70 MB
Inference Memory: 1.25 MB
Memory footprint: 6.95 MB

PERFORMANCE & SPEED ANALYSIS
-----
Test Accuracy: 82.31%
Test Loss: 1.2696
Single Sample Time: 1.58 ± 2.86 ms
GPU Available: Yes
GPU Memory Growth: True
```

- Student model with much smaller size got accuracy >70%

Deployment Strategy Recommendation

- Scenario A: Real-time Video Processing

Since real-time video processing requires extremely low latency, the communication time between cloud servers and devices is not acceptable. So, in this situation, optimizing a high accuracy model to edge device is suitable. Not only will it reduce the communication time between device and server, it can also

support by specialized hardware to provide further speed up. Also, by using cloud backup, the model can be continuously improved. For optimization priority, choosing a lightweight and efficient architecture which designed for speed is the easiest way. Design algorithm suitable for specific hardware and quantization is also a important point.

- Scenario B: IoT Sensor Network

In Scenario B, it requires an extreme power constraint solution. As a result, TinyML optimization is the best choice. TinyML provides models with extremely small size and power consumption. Although it also results in worse accuracy, but with occasional cloud sync, it can be optimized in real life problems.

- Scenario C: Mobile Application

In Scenario C, it needs to balance accuracy and efficiency. Which corresponds to cascade deployment. At most of the time, a core on-device baseline model with high efficiency take place to provide basic functionality. On device models also provide offline ability. If user's device can provide enough computational power and energy, a larger but more accurate model will take place to provide higher accuracy.

Optimization Effectiveness:

1. Cloud Optimization: Cloud optimization strategy focuses on accelerating training and optimizing for large model training. By using mixed precision training, hardware matrix calculation speed can be accelerated. When using distributed training, multiple GPU can be used for training a single model at the same time, which greatly speeds up the training process. With gradient accumulation, when the training data is too large to handle, it can still provide statistical stability and final accuracy.
2. Edge optimization: At the edge, models must run in real-time with limited hardware. The focus changes from training to inference. By using pruning to increase the sparsity of the model, the model size and computational flops will decrease with minimum accuracy trade off. By using post training quantization, the memory requirement of model decreases significantly with slightly accuracy drop.
3. TinyML optimization: The techniques used by TinyML and Edge optimization are similar. But TinyML is used much more aggressively. For example, extreme pruning and reduce precision from 32bit to 8bit. Which often results in significant drop of accuracy. However, it gains a much smaller model size, latency, and energy consumption. Which will be capable of deploying microcontrollers. By using model architecture optimization, it can retain an acceptable accuracy with extremely small model size.

Resource Constraint Impact:

Cloud Server has the most computational resources, memory, and power for AI models. In cloud server, model latency and size are not the major optimization challenges but model accuracy is. So, we focus on optimizing large model training improvements.

Edge devices have limited but still reasonable resources and often require the model to handle real time processing cases. So, the focus will be on reducing the model size and latency while preserving model accuracy.

Microcontrollers have extremely limited resource constraints. The focus point is to first make the model work on this situation. So, we need to aggressively reduce model size and energy consumption.

Development Trade-offs:

In multi scale optimization, the first challenge will be tooling. There are different tools that need to be handled for implementing different optimization techniques. A team must master different usage of tools. It also makes debugging more complex. The second challenge will be the proliferation of model variants. To find the optimal model, a lot of evaluation metrics should be used to measure the different variants of models developed.

Real-world Deployment:

To move towards real-world production deployment, performance monitoring should be added to the pipeline. On production deployment, the real-world situation may be different from what we expected during development. Some elements may not be considered which will affect the result. So, we need a performance monitoring system to clarify the model's performance is the same as we expected. We also need versioning, when we face some unexpected problems, we can quickly return to the version that works just fine.

Future Evaluation:

In the future, there may be some specialized AI accelerators or techniques that make the optimization process less complicated and can easily reach the performance target.