

Model Compression and Deployment Techniques

October 22, 2025

Contents

1 Model Pruning, Distillation, and Quantization

Compression techniques shrink model size, reduce latency, and lower memory footprint while preserving accuracy. Figure ?? compares pruning, distillation, and quantization workflows.

1.1 Pruning

Pruning removes redundant parameters or structures. Given weights \mathbf{W} and pruning mask \mathbf{M} , the effective parameters are $\tilde{\mathbf{W}} = \mathbf{M} \odot \mathbf{W}$. Common schemes include:

- **Unstructured pruning:** Zeroes individual weights with small magnitude; implements via l_0 regularization or iterative magnitude pruning. Sparse linear algebra is required to realize runtime gains.
- **Structured pruning:** Drops channels, filters, or attention heads to match hardware constraints. Channel pruning often optimizes

$$\min_{\mathbf{M}} \mathcal{L}(\mathbf{M} \odot \mathbf{W}) + \lambda \|\mathbf{M}\|_0 \quad \text{s.t.} \quad \sum_c M_c \leq K, \quad (1)$$

where K is the target channel budget.

- **Dynamic pruning:** Chooses masks conditioned on the input (e.g., SkipNet, dynamic token pruning). Reinforcement signals balance accuracy and cost.

Lottery ticket experiments reveal subnetworks that, when retrained, match original performance. In practice, pruning is paired with fine-tuning or knowledge distillation to recover accuracy.

1.2 Knowledge Distillation

Distillation transfers knowledge from a teacher model f_T to a student f_S . The blended objective combines hard labels and soft teacher logits:

$$\mathcal{L}_{\text{KD}} = (1 - \alpha) \mathcal{L}_{\text{CE}}(f_S(\mathbf{x}), \mathbf{y}) + \alpha T^2 \text{KL} \left(\sigma \left(\frac{f_T(\mathbf{x})}{T} \right) \parallel \sigma \left(\frac{f_S(\mathbf{x})}{T} \right) \right), \quad (2)$$

where T is the temperature, α the distillation weight, and σ denotes softmax. Variants include intermediate feature matching, attention transfer, and self-distillation with ensemble teachers.

1.3 Quantization

Quantization maps high-precision weights and activations to lower bit-width representations:

$$q = \text{clip} \left(\text{round} \left(\frac{x}{s} \right) + z, q_{\min}, q_{\max} \right), \quad (3)$$

with scaling factor s , zero-point z , and quantized range $[q_{\min}, q_{\max}]$. Key modes:

- **Post-training quantization (PTQ):** Calibrates scale factors using a small dataset; suitable for INT8 inference when accuracy drop is acceptable.
- **Quantization-aware training (QAT):** Simulates quantization during training via straight-through estimators (STE), enabling INT8 or INT4 deployment with minimal degradation.
- **Mixed-precision quantization:** Assigns bit widths per layer to meet accuracy/latency targets; solved via integer programming or reinforcement learning.

Activation outlier suppression (e.g., SmoothQuant) facilitates INT8 inference for transformers by redistributing scaling factors across weights and activations.

1.4 Workflow Integration

Compression pipelines often interleave techniques: prune the teacher, distill into a compact student, and then quantize. Hardware-aware neural architecture search (NAS) explores architectures that better survive quantization. The following pseudo-code outlines a joint workflow:

Listing 1: Combined pruning, distillation, and quantization-aware training.

```
1 teacher = load_pretrained_model()
2 student = initialize_compact_model()
3
4 # Structured pruning on teacher
5 for step in range(prune_steps):
6     loss = training_step(teacher, data_batch)
7     loss.backward()
8     apply_structured_pruning(teacher, sparsity_schedule(step))
9
10 # Distill knowledge to student
11 for epoch in range(kd_epochs):
12     for batch in dataloader:
13         teacher_logits = teacher(batch.inputs).detach()
14         loss = kd_loss(student(batch.inputs), batch.labels, teacher_logits,
15                        alpha=0.7, temperature=4.0)
16         loss.backward()
17         optimizer.step()
18         optimizer.zero_grad()
19
20 # Quantization-aware fine-tuning
21 quantizer = prepare_qat(student, bitwidth=8)
22 for epoch in range(qat_epochs):
23     for batch in dataloader:
24         output = quantizer(batch.inputs)
25         loss = criterion(output, batch.labels)
26         loss.backward()
27         optimizer.step()
28         optimizer.zero_grad()
29 export_int8(quantizer, path="student_int8.onnx")
```

Compression Workflow Overview

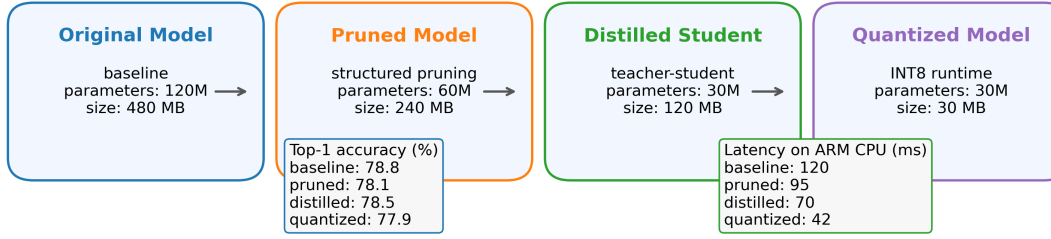


Figure 1: Comparison of pruning, distillation, and quantization pipelines. The bottom row shows accuracy vs. model size trade-offs.

2 Deployment to Mobile and Edge (TensorRT, ONNX, TFLite)

Edge deployment requires toolchains that convert trained models into device-specific runtimes. Figure ?? summarizes the major ecosystems.

2.1 ONNX as an Interchange Format

The Open Neural Network Exchange (ONNX) defines an intermediate representation (IR) with operator sets (opsets). Exporting from PyTorch or TensorFlow yields a portable graph:

$$\text{Graph} = (\mathcal{V}, \mathcal{E}, \mathcal{O}), \quad (4)$$

where \mathcal{V} are tensors, \mathcal{E} edges, and \mathcal{O} operator nodes. Version alignment between exporter and runtime is critical. Shape inference and constant folding reduce redundancy before deployment.

2.2 TensorRT Optimization

TensorRT compiles ONNX graphs into CUDA kernels. Major passes include layer fusion, precision calibration (FP16/INT8), and kernel auto-tuning. The optimization profile specifies input ranges for dynamic shapes. Execution providers (TensorRT EP) integrate into ONNX Runtime to fall back to CPU/GPU operators when unsupported.

2.3 TensorFlow Lite (TFLite)

TFLite converts TensorFlow SavedModels into a flatbuffer format with selective operator kernels for mobile CPUs, GPUs, and NPUs. Quantization-aware training can export INT8 kernels compatible with Edge TPU. Delegate mechanisms (e.g., NNAPI, Core ML) offload computation to vendor accelerators.

2.4 Deployment Checklist

- Validate numerical parity between source framework and exported model via golden tests.
- Profiles memory usage and latency under realistic batch sizes and sequence lengths.
- Integrate fallback paths: e.g., if TensorRT fails to build an engine, fall back to ONNX Runtime GPU.
- Monitor operator coverage; custom ops require plugin development or graph rewriting.

The following script demonstrates exporting a PyTorch model to ONNX and building a TensorRT engine using the Python API:

Listing 2: PyTorch to ONNX export and TensorRT engine building.

```
1 import torch
2 import onnx
3 import tensorrt as trt
4
5 model = build_model().eval().cuda()
6 dummy = torch.randn(1, 3, 224, 224, device="cuda")
7 torch.onnx.export(model, dummy, "model.onnx",
8                   input_names=["input"], output_names=["logits"],
9                   opset_version=17, do_constant_folding=True,
10                  dynamic_axes={"input": {0: "batch"}, "logits": {0: "batch"}})
11
12 onnx_model = onnx.load("model.onnx")
13 onnx.checker.check_model(onnx_model)
14
15 logger = trt.Logger(trt.Logger.INFO)
16 builder = trt.Builder(logger)
17 network = builder.create_network(1 << int(trt.NetworkDefinitionCreationFlag.
18    EXPLICIT_BATCH))
19 parser = trt.OnnxParser(network, logger)
20 with open("model.onnx", "rb") as f:
21     parser.parse(f.read())
22
23 config = builder.create_builder_config()
24 config.set_memory_pool_limit(trt.MemoryPoolType.WORKSPACE, 1 << 30)
25 if builder.platform_has_fast_int8:
26     config.set_flag(trt.BuilderFlag.INT8)
27     config.int8_calibrator = gather_calibration_data()
28
29 engine = builder.build_engine(network, config)
30 with open("model.plan", "wb") as f:
31     f.write(engine.serialize())
```

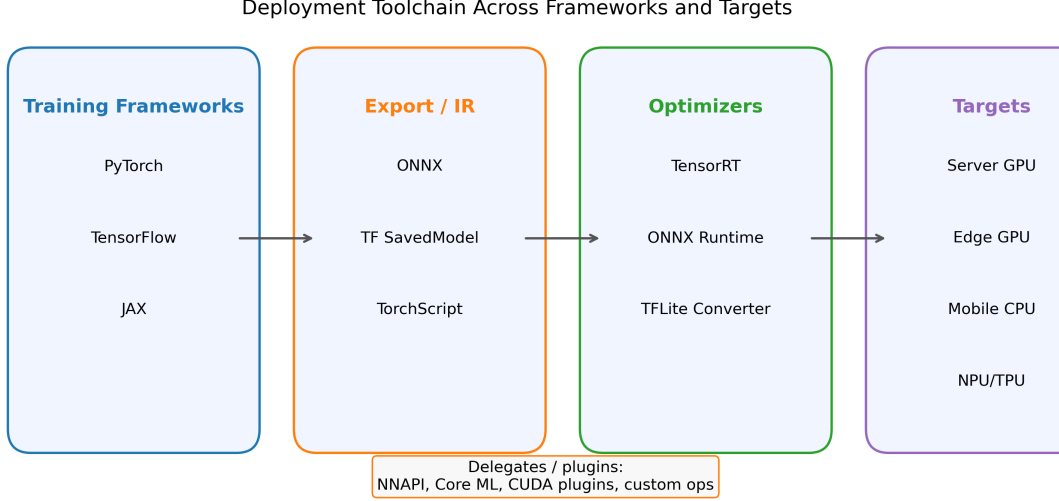


Figure 2: Deployment toolchains across ONNX, TensorRT, and TensorFlow Lite. Optional delegates target vendor accelerators.

3 Inference Acceleration

Inference optimization targets latency, throughput, and energy efficiency. Figure ?? shows a latency decomposition, while Figure ?? outlines common acceleration strategies.

3.1 Kernel and Graph Optimizations

Operator fusion merges sequences like Conv-BN-ReLU into a single kernel, reducing memory traffic. Graph compilers (TVM, XLA, TorchInductor) apply loop tiling, vectorization, and layout transformations. The latency L can be approximated as

$$L = \sum_{i=1}^N \left(\frac{C_i}{\text{FLOP/s}} + \frac{M_i}{\text{BW}} \right), \quad (5)$$

where C_i is compute cost and M_i memory traffic. Tuning seeks to minimize both components via scheduling.

3.2 Batching and Dynamic Serving

Batching amortizes overhead across requests. For online systems with arrival rate λ and service rate μ , queueing delay is governed by the Erlang C formula. Dynamic batching (Triton Inference Server) accumulates requests up to latency budgets. For autoregressive models, speculative decoding and scheduling partial beams reduce token latency.

3.3 Hardware Acceleration

Specialized accelerators (Edge TPU, NVIDIA Tensor Cores) offer mixed-precision support. Memory-bound models benefit from high-bandwidth memory (HBM) and sparsity-aware hardware. For edge deployments, CPU vector engines (NEON, AVX512) are leveraged via libraries like XNNPACK and oneDNN.

3.4 Monitoring and A/B Testing

Production systems require continuous measurement of latency percentiles (P50/P95/P99), throughput, and energy consumption. Canary releases test optimized models on a subset of traffic to ensure stability. Rollback procedures and feature flags guard against degraded user experience.

Listing 3: Triton Inference Server dynamic batching configuration (YAML).

```

1 name: "resnet_triton"
2 platform: "tensorrt_plan"
3 max_batch_size: 32
4 input [
5   { name: "input", data_type: TYPE_FP16, dims: [3, 224, 224] }
6 ]
7 output [
8   { name: "logits", data_type: TYPE_FP16, dims: [1000] }
9 ]
10 dynamic_batching {
11   preferred_batch_size: [4, 8, 16, 32]
12   max_queue_delay_microseconds: 2000
13 }
14 instance_group [
15   { kind: KIND_GPU, count: 2, gpus: [0, 1] }
16 ]

```

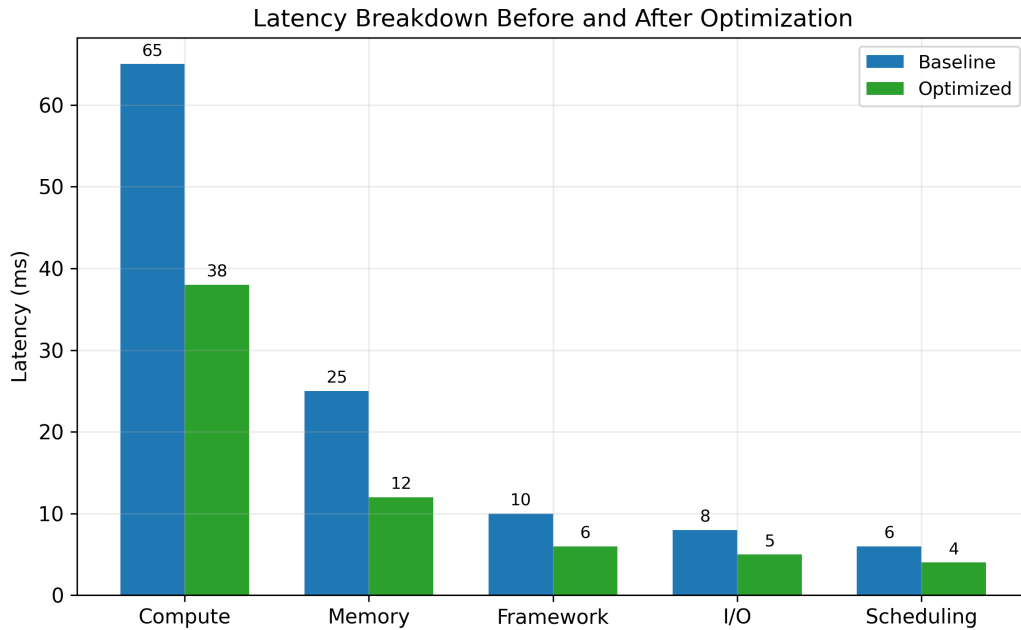


Figure 3: Latency breakdown across compute, memory, and I/O components. Profiling highlights the dominant bottleneck.

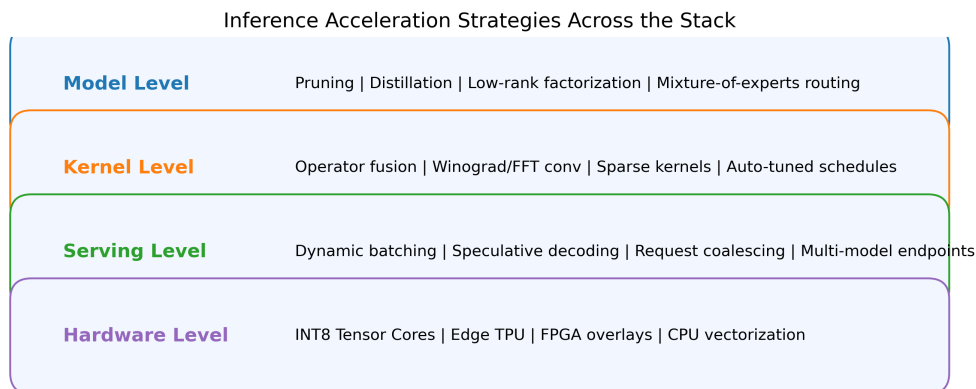


Figure 4: Inference acceleration strategies spanning model, kernel, serving, and hardware layers.

Further Reading

- Song Han et al. “Learning both Weights and Connections for Efficient Neural Networks.” NIPS 2015.
- Geoffrey Hinton et al. “Distilling the Knowledge in a Neural Network.” NIPS 2015 Workshop.
- Jacob et al. “Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference.” CVPR 2018.
- NVIDIA. “TensorRT Developer Guide.” 2023.
- Jared Casper et al. “Amazon SageMaker Inference Recommender.” 2022.