# Compression Decision Tree: Choose Your Strategy
# Based on Platform Constraints and Requirements



**What is your target platform?**

- Server (80GB)
- Edge (16GB)
- Mobile (4GB)
- MCU (512MB)

**Server (80GB)** → Need max accuracy?
- Yes → FP32 Full model → 700GB 100% acc
- No → INT8 Quantization → 175GB 99% acc

**Edge (16GB)** → Need fast inference?
- Yes → INT8 + Pruning → 50GB 97% acc Fast!
- No → INT4 Quantization → 87GB 97% acc

**Mobile (4GB)** → Size critical?
- Yes → Distillation + INT8 → 10GB 92% acc
- No → INT8 Quantization → 175GB 99% acc

**MCU (512MB)** → Complex task?
- Yes → INT4 + Distillation → 5GB 90% acc
- No → Binary Networks → 500MB 85% acc

Server: Quantize if possible

Edge: INT4 or prune

Mobile: Distill + quantize

MCU: Extreme compression

Decision  Method  Result