

# Edge-Ready GenAI: Engineering Low-Latency Solutions for Resource- Constrained Environments

Welcome to our comprehensive exploration of generative AI deployment in resource-constrained environments. As enterprises increasingly demand real-time AI capabilities, understanding how to optimize GenAI for edge computing becomes critical for developers and engineers.

Today we'll examine systematic approaches to maintain performance while significantly reducing computational demands, energy consumption, and latency—opening new possibilities for intelligent edge applications.

**By: Sai Kalyan Reddy Pentaparthi**





About me:

**Sai Kalyan Reddy Pentaparthi**

Principal Engineer Software

**ST Engineering iDirect, Inc.**

Working on enhancing connectivity through Global IP-Satellite Network Infrastructure

[www.linkedin.com/in/saikalyanrp](https://www.linkedin.com/in/saikalyanrp)



# The Edge Computing Frontier



## Limited Computational Resources

Edge devices generally have limited processing power, memory, and storage compared to centralized cloud systems.



## Energy Constraints

Many edge devices run on batteries or operate under strict power limits, restricting AI model complexity.



## Connectivity Challenges

Edge AI solutions must operate reliably even with limited or no connection to cloud resources.



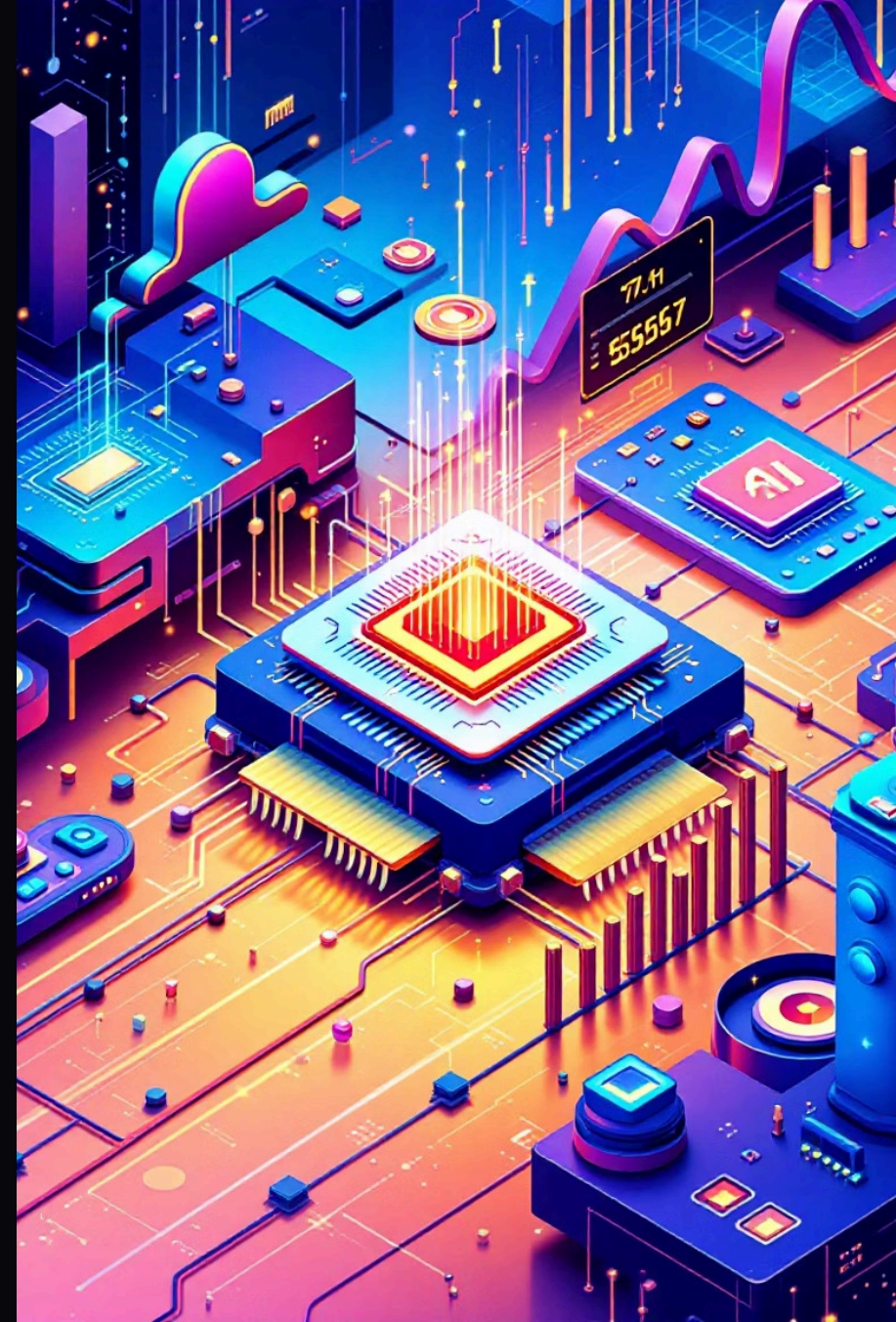
## Real-time Requirements

Applications often require low-latency responses, increasing performance demands on edge models.



# Edge AI Market Trends

- **Strong Market Expansion:**
  - The Edge AI accelerator market is growing rapidly with a **38.9% CAGR**.
  - It is projected to reach **\$7.68 billion by 2027**, showing strong demand for on-device AI.
- **Key Adoption Drivers:**
  - The need for low **latency (78.3% of organizations)** to ensure real-time responses is a primary driver.
  - **Data privacy concerns (64.7%)** also encourage local AI processing.
- **Industry Adapting to Challenges:**
  - **Model optimization (used in 82% of deployments)** is widely applied to run complex AI efficiently on constrained edge hardware.



# The Edge AI Market Evolution

1

## Current State

Limited GenAI deployment at edge due to resource constraints

2

## Next 12 Months

Increasing adoption of optimized edge AI solutions

3

## 18-24 Months

Most enterprises requiring real-time edge AI capabilities

4

## Beyond 2025

Ubiquitous edge AI with specialized hardware acceleration

Market forecasts indicate we're at an inflection point in edge AI adoption. The next 18-24 months will witness a dramatic shift as enterprises increasingly require real-time AI capabilities without the latency penalties of cloud-dependent architectures.

This transition is being accelerated by advances in hardware acceleration, model optimization techniques, and the growing necessity for privacy-preserving local computation.



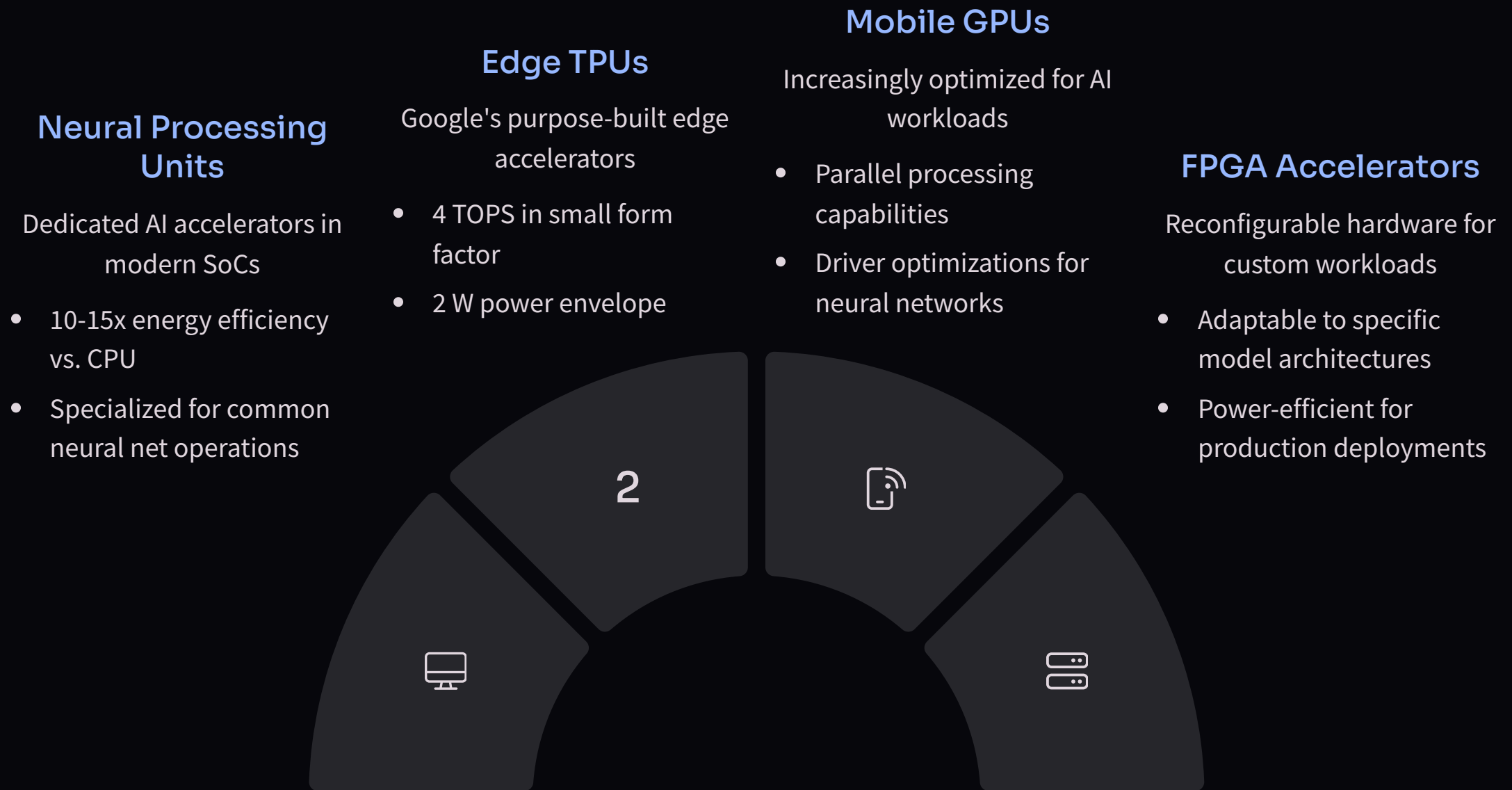


# Edge AI Market Trends

- **Strong Market Expansion:**
  - The Edge AI accelerator market is experiencing rapid growth (**38.9% CAGR**).
  - Projected to reach **\$7.68 billion by 2027**, indicating significant demand for on-device AI.
- **Key Adoption Drivers:**
  - Driven primarily by the need for low **latency (78.3% of organizations)** for real-time responses.
  - **Data privacy concerns (64.7%)** are also a major factor, favoring local processing.
- **Industry Adapting to Challenges:**
  - Widespread **model optimization (used in 82% of deployments)** shows the market actively working to run complex AI efficiently on constrained edge hardware.



# Hardware Acceleration Solutions



Specialized hardware accelerators dramatically improve the performance and energy efficiency of edge AI workloads. Modern system-on-chips (SoCs) increasingly incorporate neural processing units (NPUs) that deliver an order of magnitude better performance-per-watt compared to general-purpose CPUs.

These purpose-built accelerators are specifically designed to handle the computational patterns of neural networks, with optimized datapaths for matrix multiplication and activation functions that dominate AI workloads.

# Edge GPU's: Overview

Table 1: Comparative Specification Overview of Select Edge GPU/Accelerator Families

Feature	NVIDIA Jetson Orin Nano Super 8GB	NVIDIA Jetson Orin NX 16GB	NVIDIA Jetson AGX Orin 64GB	NVIDIA RTX 4000 SFF Ada	AMD Ryzen AI 300 Series (Rep. NPU) <sup>1</sup>	Intel Core Ultra (Rep. NPU) <sup>2</sup>	Intel Arc A380E	Google Coral Dual Edge TPU M.2	Qualcomm Snapdragon X Elite (Rep. NPU) <sup>3</sup>	EdgeCortex SAKURA-I
Manufacturer	NVIDIA	NVIDIA	NVIDIA	NVIDIA	AMD	Intel	Intel	Google	Qualcomm	EdgeCortex
Key Accelerator	GPU + Tensor Cores	GPU + Tensor Cores	GPU + Tensor Cores + DLAs	GPU + Tensor Cores	NPU (XDNA) + iGPU (RDNA)	NPU + iGPU (Arc)	GPU (Arc)	TPU (ASIC)	NPU (Hexagon)	DNA (ASIC/FPGA)
AI Perf. (TOPS)	67 (INT8 Sparse)	100-157 (INT8 Sparse)	275 (INT8 Sparse)	19.2 TFLOPS (FP32)	~45 TOPS (Total System) <sup>4</sup>	~10-11 TOPS (NPU only) <sup>5</sup>	~5 TFLOPS (FP32)	8 (INT8)	45 (NPU only) <sup>6</sup>	40 (INT8 Dense)
Memory Capacity (GB)	8	16	64	20	System RAM	System RAM	6 (GDDR6)	N/A (Uses Host)	System RAM	N/A (Uses Host)
Memory Bandwidth (GB/s)	102	102.4	204.8	160	LPDDR5X (System)	LPDDR5/x (System)	186	N/A	136 (LPDDR5X)	N/A
TDP (W)	7-25	10-25	15-60	70	Laptop TDP	Laptop TDP	75	~4 (Module)	Laptop TDP	5-10
Typical Form Factor	SOM	SOM	SOM	SFF PCIe Card	Integrated SoC	Integrated SoC	PCIe Card	M.2 E-Key	Integrated SoC	PCIe Card / Chip

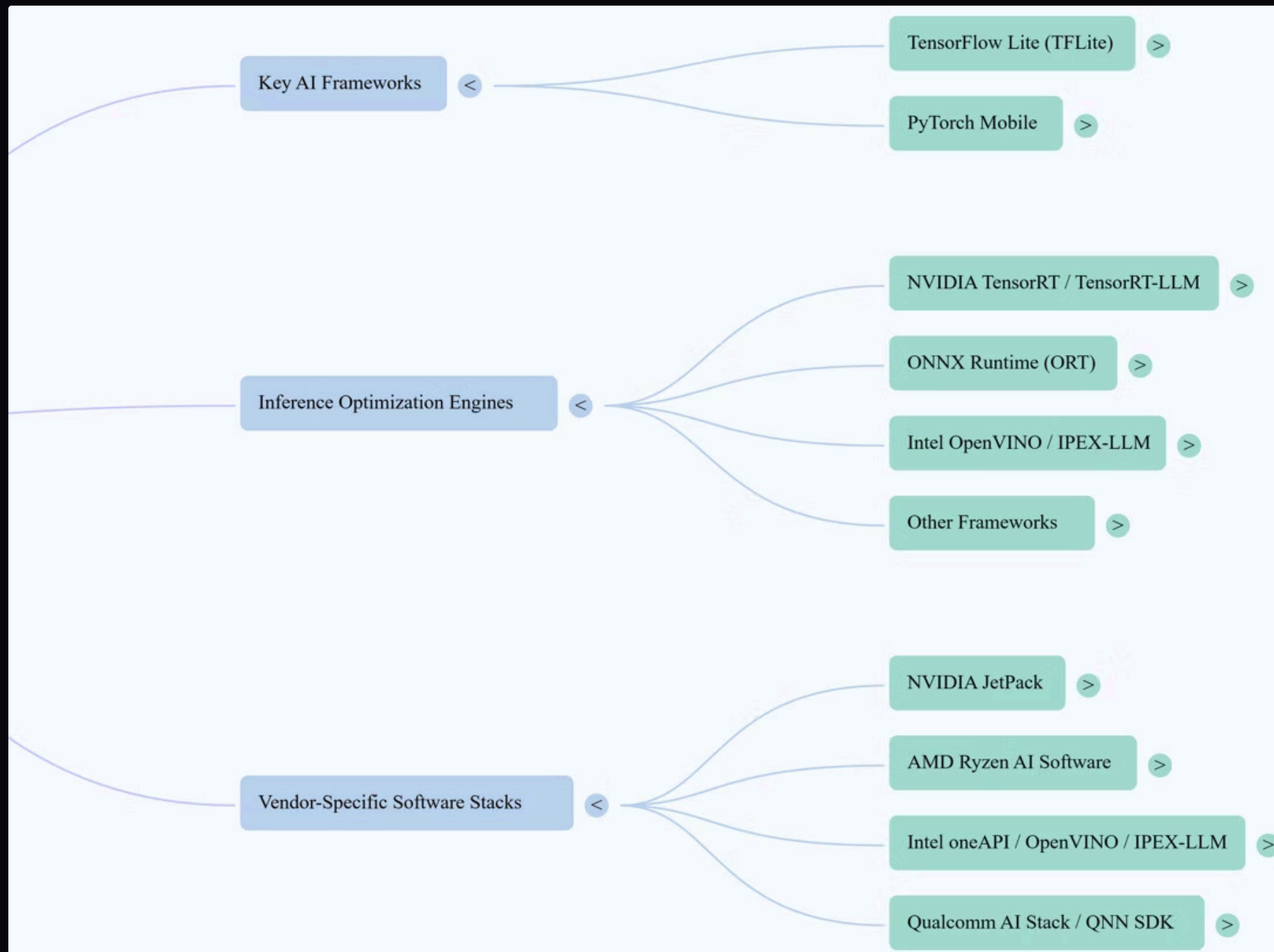


# Edge GPU's: Inference Performance

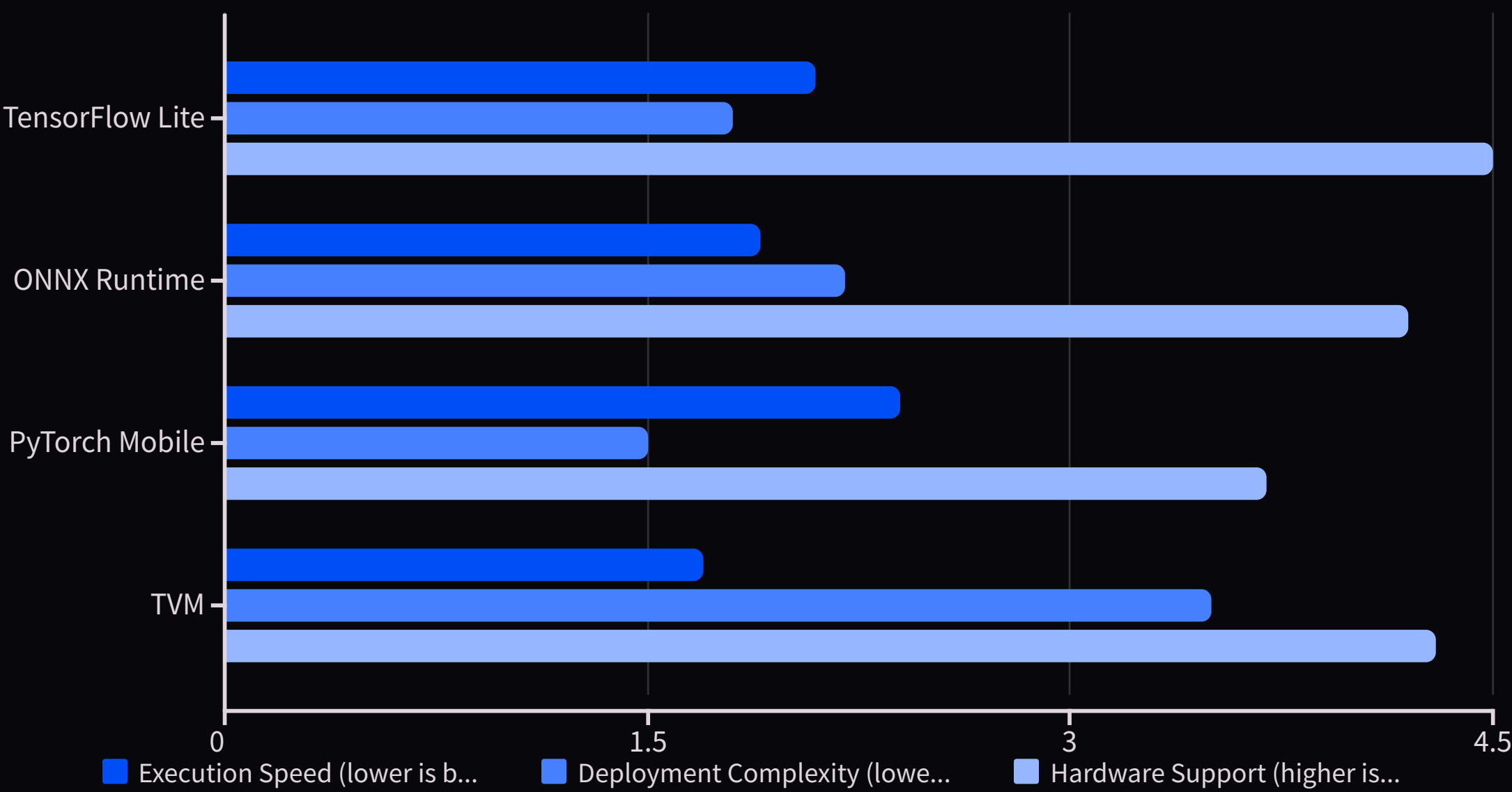
Table 2: LLM Inference Performance Comparison (Illustrative - Based on Fragmented Data)

Hardware	LLM Model	Quantization	Software/Framework	Throughput (Tokens/Sec)	Latency (TTFT ms / ITL ms)	Power (W)	Source/Notes
NVIDIA GPU (RTX 4090)	Mistral 7B	INT4 (AWQ)	TensorRT-LLM v0.7.1	~170	N/A / ~5.9	N/A	Single user benchmark. Significantly faster than llama.cpp (~100 TPS) on same HW.
NVIDIA GPU (RTX 3090)	Mistral 7B	INT4 (AWQ)	TensorRT-LLM v0.7.1	~142	N/A / ~7.0	N/A	Single user benchmark. ~62% faster than llama.cpp (~87 TPS) on same HW.
NVIDIA Jetson Orin Nano Super	Llama 3.1 8B	Optimized	TensorRT-LLM / JetPack 6+	<i>Data Needed</i>	<i>Data Needed</i>	7-25W	Capability confirmed, performance claimed 1.7x higher than Orin Nano. Specific benchmarks pending.
NVIDIA Jetson Orin NX 16GB	Llama 2 7B Chat	GPTQ (INT4?)	TensorRT-LLM v0.12 (Preview)	<i>Data Needed</i>	<i>Data Needed</i>	10-25W	Feasibility shown, requires JetPack 6.1+. Specific benchmarks pending.
AMD Ryzen AI 300 Series	Llama 3.1 8B	INT4 (W4A16)	ONNX Runtime + Vitis AI EP	<i>Data Needed</i>	<i>Data Needed</i>	Laptop TDP	Supported via LLM Hybrid OGA flow (NPU+iGPU). Specific performance data needed.
AMD Ryzen AI 300 Series	Mistral 7B Instruct	INT4 (W4A16)	ONNX Runtime + Vitis AI EP	<i>Data Needed</i>	<i>Data Needed</i>	Laptop TDP	Supported via LLM Hybrid OGA flow (NPU+iGPU). Specific performance data needed.
Intel Core Ultra (NPU)	Llama / Mistral	Optimized	IPEX-LLM / OpenVINO	<i>Data Needed</i>	<i>Data Needed</i>	Laptop TDP	Experimental support exists. Performance highly dependent on software maturity.
Intel Arc A770 (16GB)	DeepSeek-R1 67B	Q4_K_M	llama.cpp (via IPEX-LLM)	<i>Data Needed</i>	<i>Data Needed</i>	~225W TDP	Feasibility shown for very large models on desktop Arc GPU with specific framework.
Mobile SoC (Dimensity 9300)	7B Model	N/A (FP16?)	llama.cpp	~5x vs SD870 (decode)	N/A	Mobile TDP	Relative CPU performance shown, not absolute TPS/Latency on GPU/NPU.

# Beyond Physical Hardware



# Framework Implementation Options



Selecting the appropriate framework for edge deployment is critical. TensorFlow Lite offers excellent hardware support through delegations to NPUs and GPUs, while ONNX Runtime provides superior cross-platform model portability and slightly better raw execution speed for many workloads.

PyTorch Mobile excels in developer experience and deployment simplicity, making it attractive for rapid prototyping. TVM offers the best performance through deep compiler optimizations but requires significantly more implementation complexity and expertise to fully leverage its capabilities.

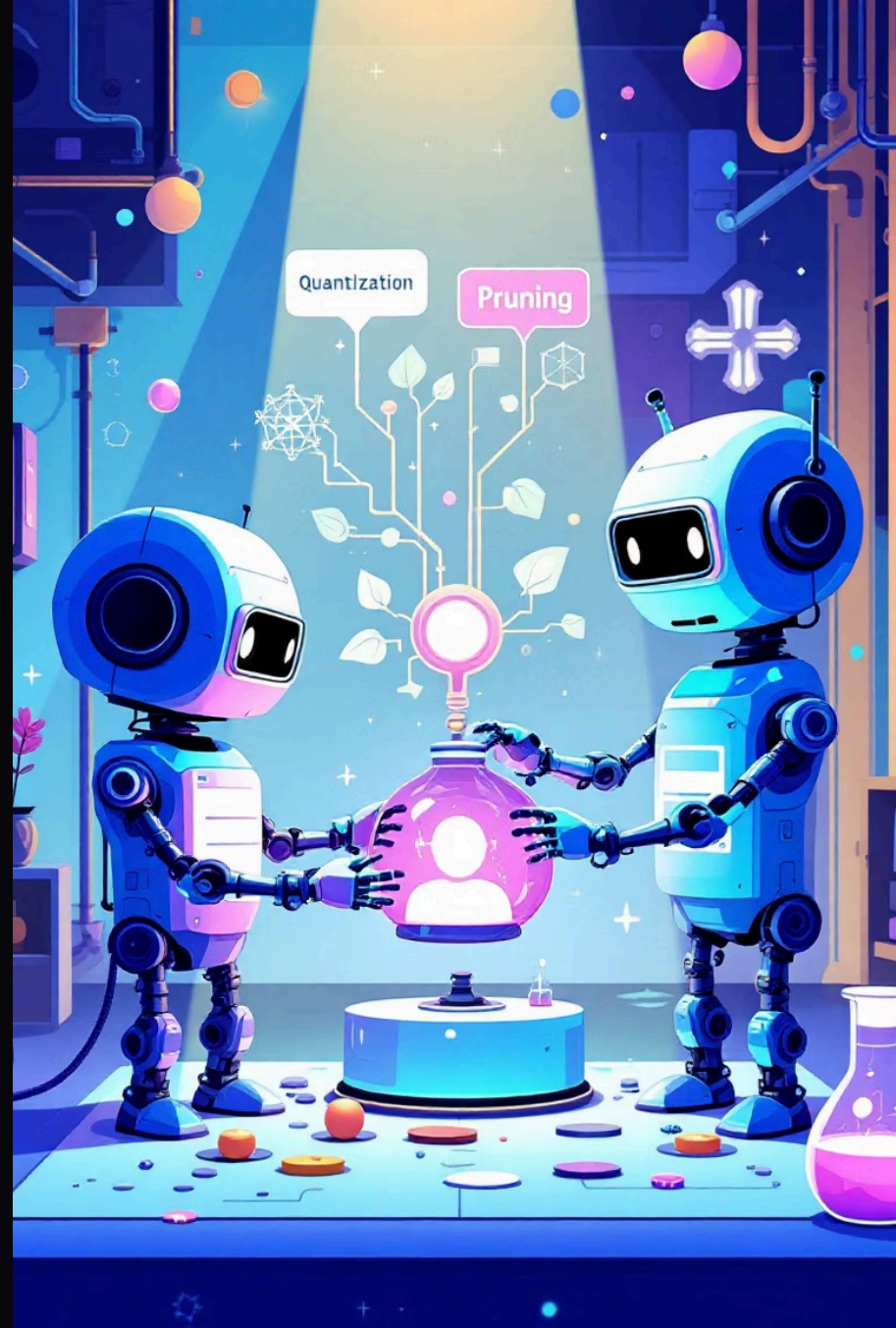


# LLM Optimization Techniques

Key techniques to optimize large language models (LLMs) include:

1. Quantization (Most crucial technique)
2. Network Pruning
3. Knowledge Distillation
4. Low-Rank Approximation / Factorization
5. Memory Optimizations

Our talk will focus on the top 3 techniques in this list.





# Quantization Techniques

## Post-Training Quantization

Converting pre-trained model weights from floating-point (FP32) to lower precision formats (INT8, INT4) after training completion.

- Minimal development effort
- No retraining required
- Moderate accuracy trade-off

## Quantization-Aware Training

Incorporating quantization effects during the training process to minimize accuracy loss when deploying with reduced precision.

- Higher implementation complexity
- Better accuracy preservation
- Model learns to compensate for quantization artifacts

## Dynamic Range Quantization

Adaptively adjusting quantization parameters based on the statistical properties of activations during inference.

- Balances accuracy and performance
- Adapts to input characteristics
- Lower memory bandwidth requirements



# Precision-Aware Quantization

## Post-Training Quantization

Simple conversion from FP32 to INT8/INT4 after training is complete. Minimal accuracy loss for many models.

## Quantization-Aware Training

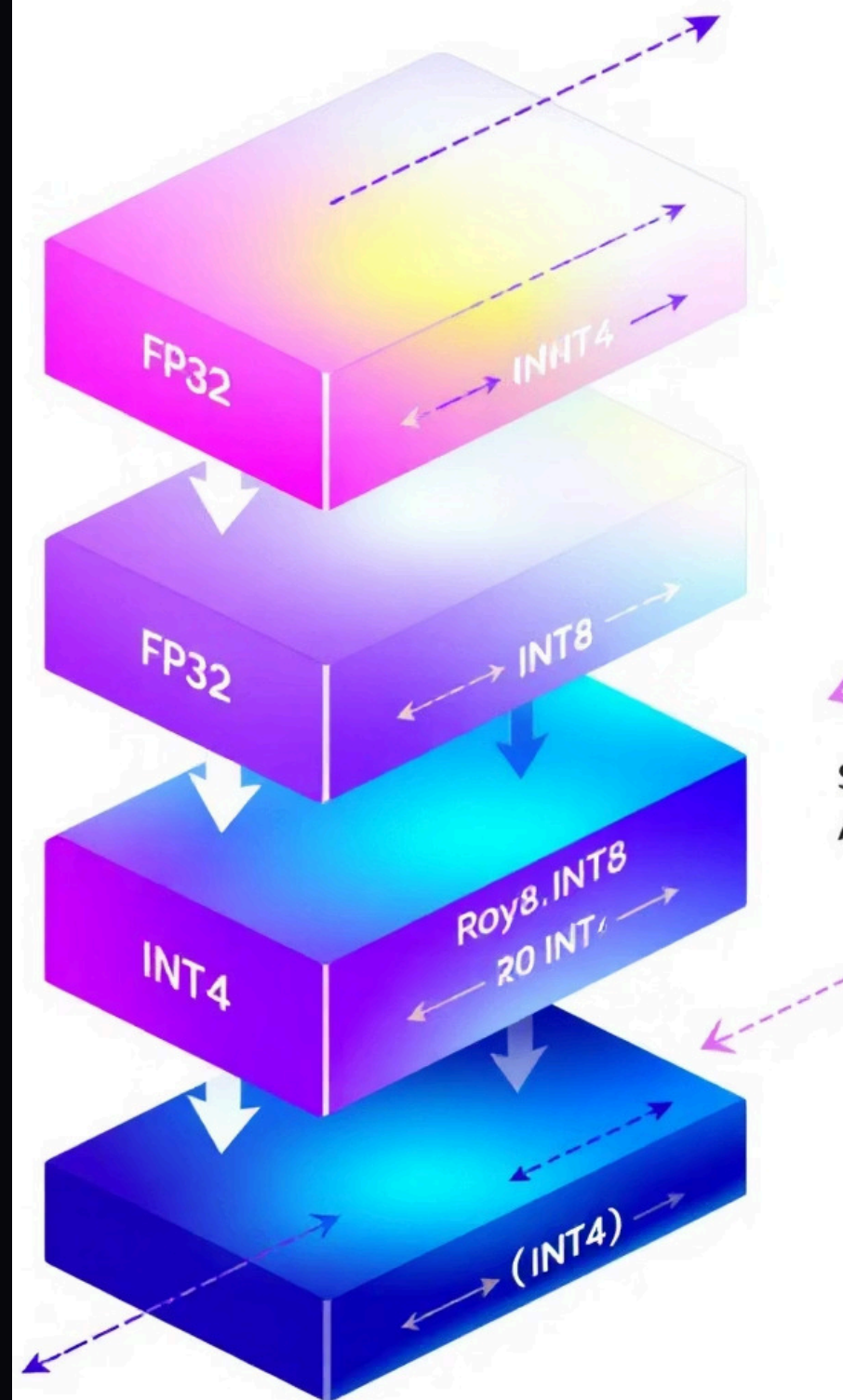
Incorporates quantization effects during training process. Models learn to compensate for reduced precision.

## Mixed-Precision Deployment

Different precision for different layers based on sensitivity analysis. Critical layers may retain higher precision.

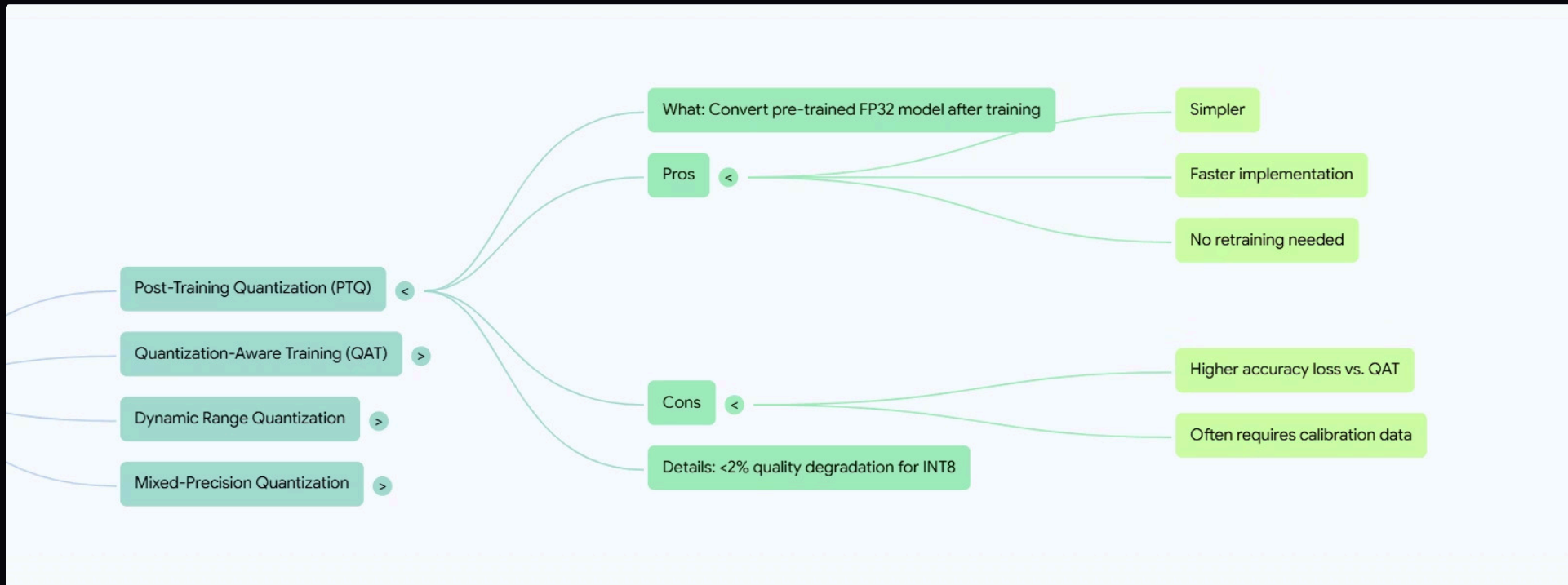
Quantization dramatically reduces the computational and memory requirements of neural networks by representing weights and activations at lower precision. Converting from 32-bit floating point to 8-bit integer can yield 4x smaller models with 3-4x faster inference.

Advanced techniques like mixed-precision quantization allow for targeted optimization, preserving full precision only where absolutely necessary while aggressively quantizing less sensitive network components.

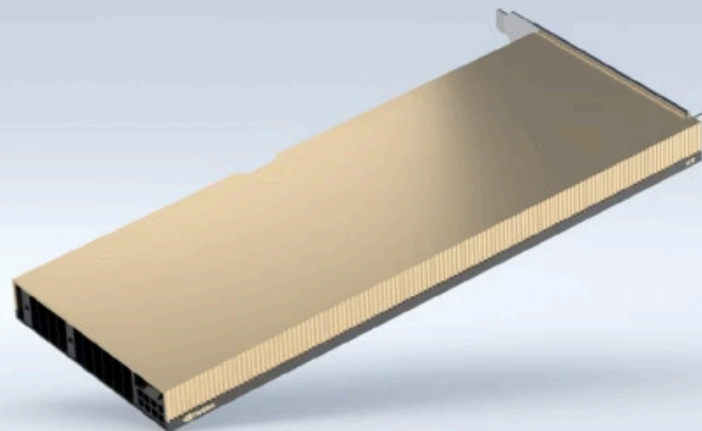




# Quantization Technique: PTQ



# GPU requirements for Float32 vs Int8 Weights



## NVIDIA A10

Accelerated Graphics and Video with AI  
for Mainstream Enterprise Servers

### Enrich Graphics and Video Applications with Powerful AI

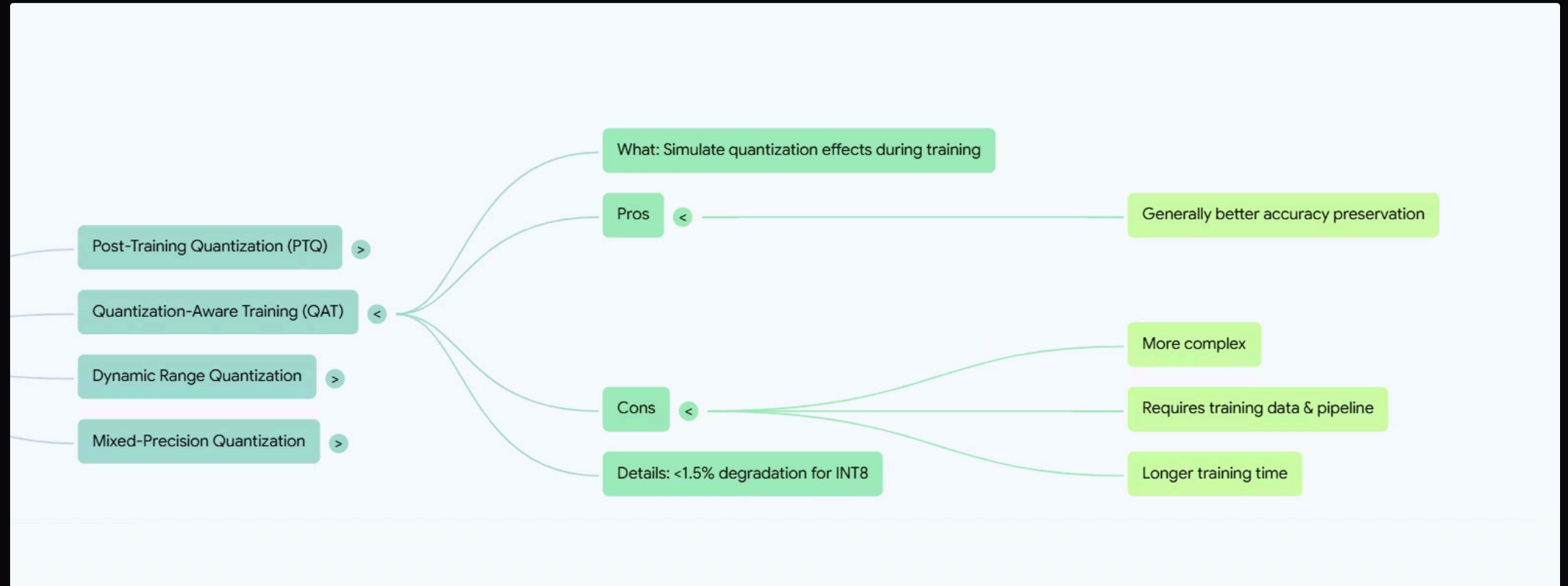
The NVIDIA A10 Tensor Core GPU combines with NVIDIA RTX Virtual Workstation (vWS) software to bring mainstream graphics and video with AI services to mainstream enterprise servers, delivering the solutions that designers, engineers, artists, and scientists need to meet today's challenges. Built on the latest NVIDIA Ampere architecture, the A10 combines second-generation RT Cores, third-generation Tensor Cores, and new streaming microprocessors with 24 gigabytes (GB) of GDDR6 memory—all in a 150W power envelope—for versatile graphics, rendering, AI, and compute performance. From virtual workstations, accessible anywhere in the world, to render nodes to the data centers running a variety of workloads, A10 is built to deliver optimal performance in a single-wide, full-height, full-length PCIe form factor.

NVIDIA A10 is supported as part of NVIDIA-Certified Systems™, in the on-prem data center, in the cloud, and at the edge. NVIDIA A10 builds on the rich ecosystem of AI frameworks from the NVIDIA NGC™ catalog, CUDA-X™ libraries, over 2.3 million developers, and over 1,800 GPU-optimized applications to help enterprises solve the most critical challenges in their business.

### SPECIFICATIONS

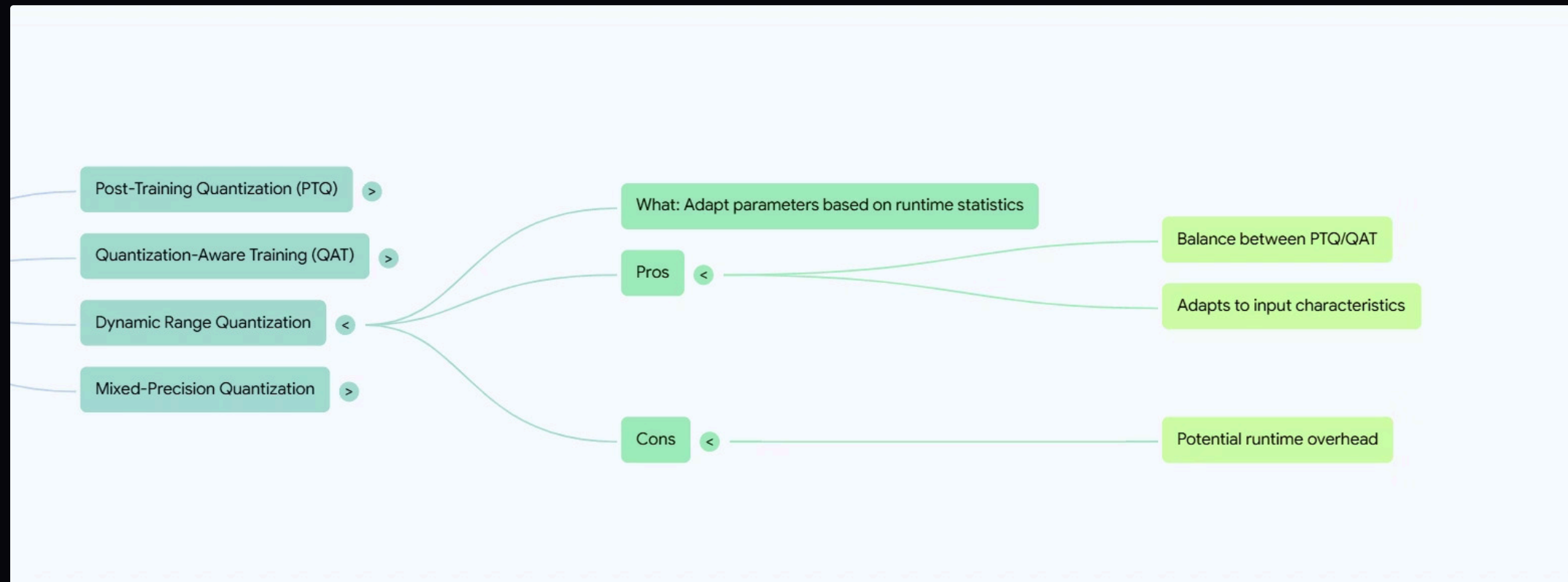
FP32	31.2 TF
TF32 Tensor Core	62.5 TF   125 TF*
BFLOAT16 Tensor Core	125 TF   250 TF*
FP16 Tensor Core	125 TF   250 TF*
INT8 Tensor Core	250 TOPS   500 TOPS*
INT4 Tensor Core	500 TOPS   1000 TOPS*
RT Cores	72
Encode / Decode	1 encoder 2 decoders (+AV1 decode)
GPU Memory	24 GB GDDR6
GPU Memory Bandwidth	600 GB/s
Interconnect	PCIe Gen4: 64 GB/s
Form Factor	1-slot FHFL
Max TDP Power	150W
vGPU Software Support	NVIDIA vPC/vApps, NVIDIA RTX™ vWS, NVIDIA AI Enterprise
Secure and Measured	Yes (optional)

# Quantization Technique: QAT

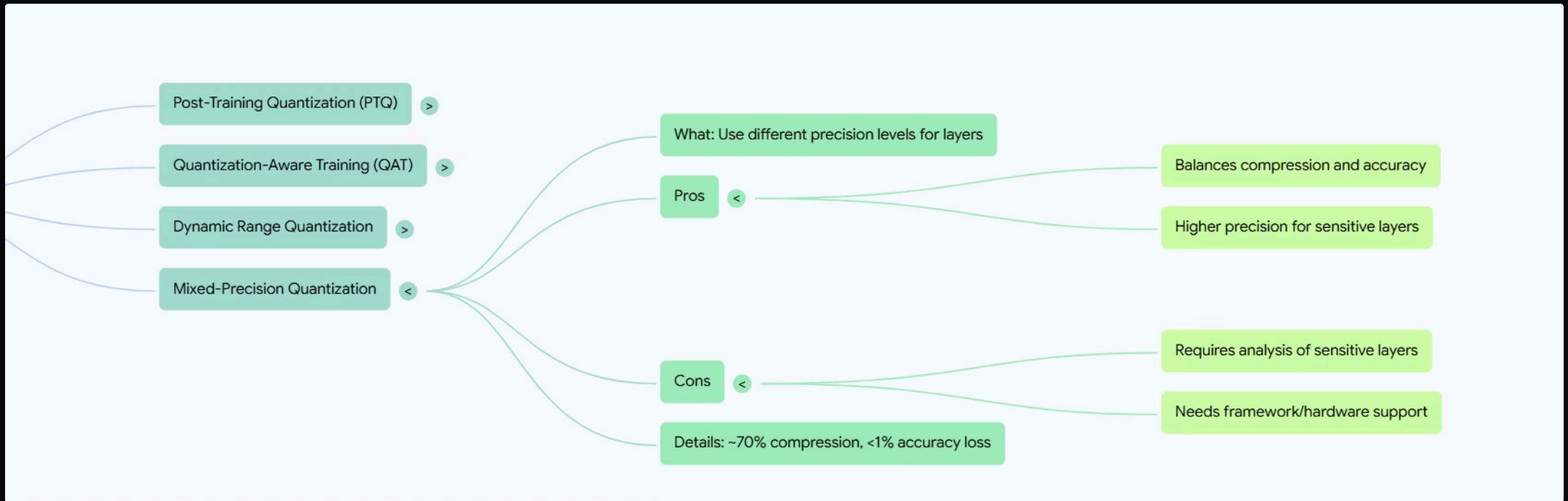




# Quantization Technique: Dynamic Range



# Quantization Technique: Mixed Precision



# Neural Network Pruning Techniques



## Identify Redundancy

Analyze weight distributions and activation patterns to locate non-essential parameters



## Structured Pruning

Remove entire filters, channels, or neurons for hardware compatibility



## Iterative Magnitude Pruning

Gradually remove small weights while retraining to maintain accuracy



## Sensitivity Analysis

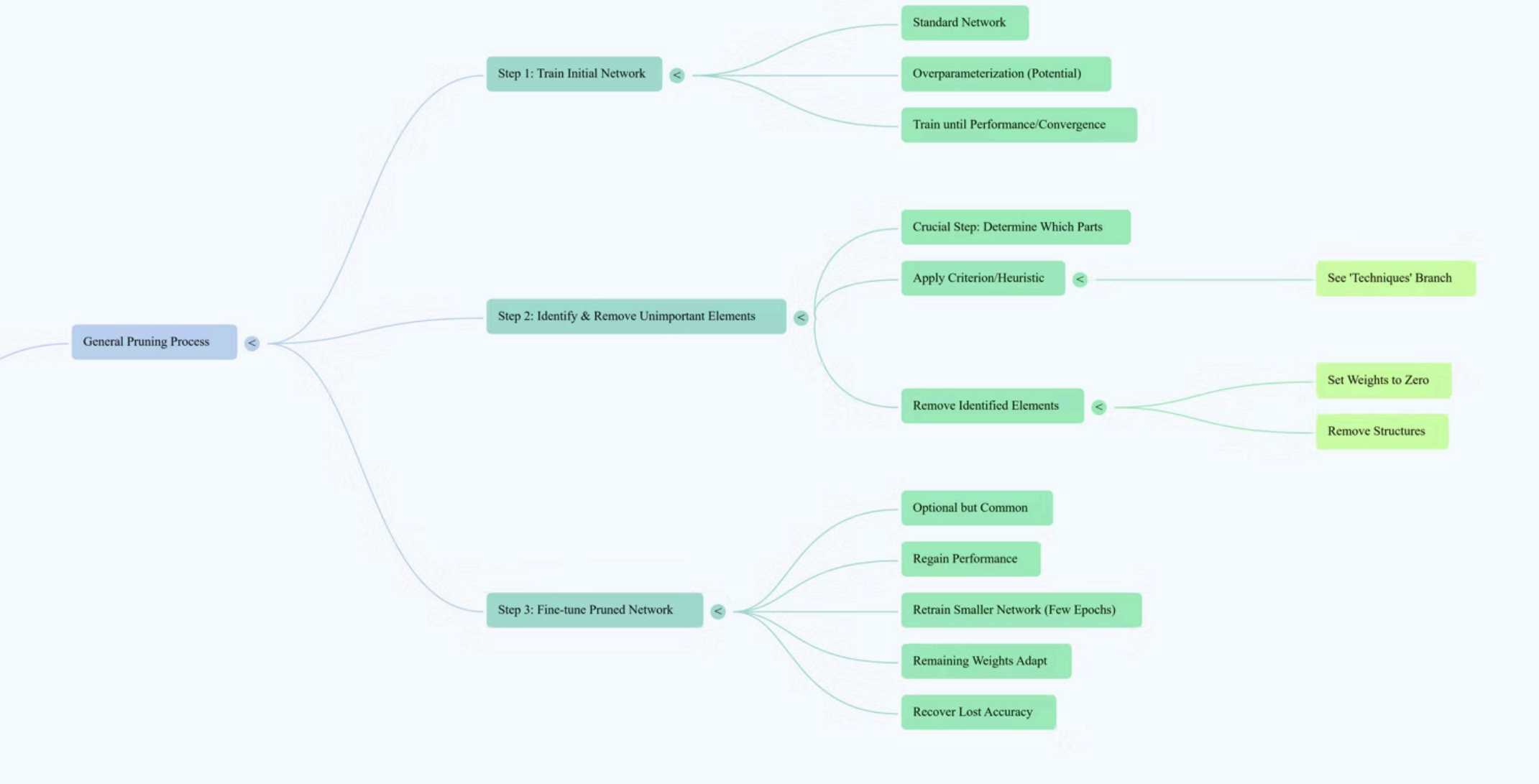
Evaluate impact of pruning on overall model performance

Targeted pruning can reduce model footprints by 50-90% with minimal accuracy impact. By systematically eliminating redundant parameters, we create leaner networks that maintain essential performance characteristics while requiring significantly fewer computational resources.

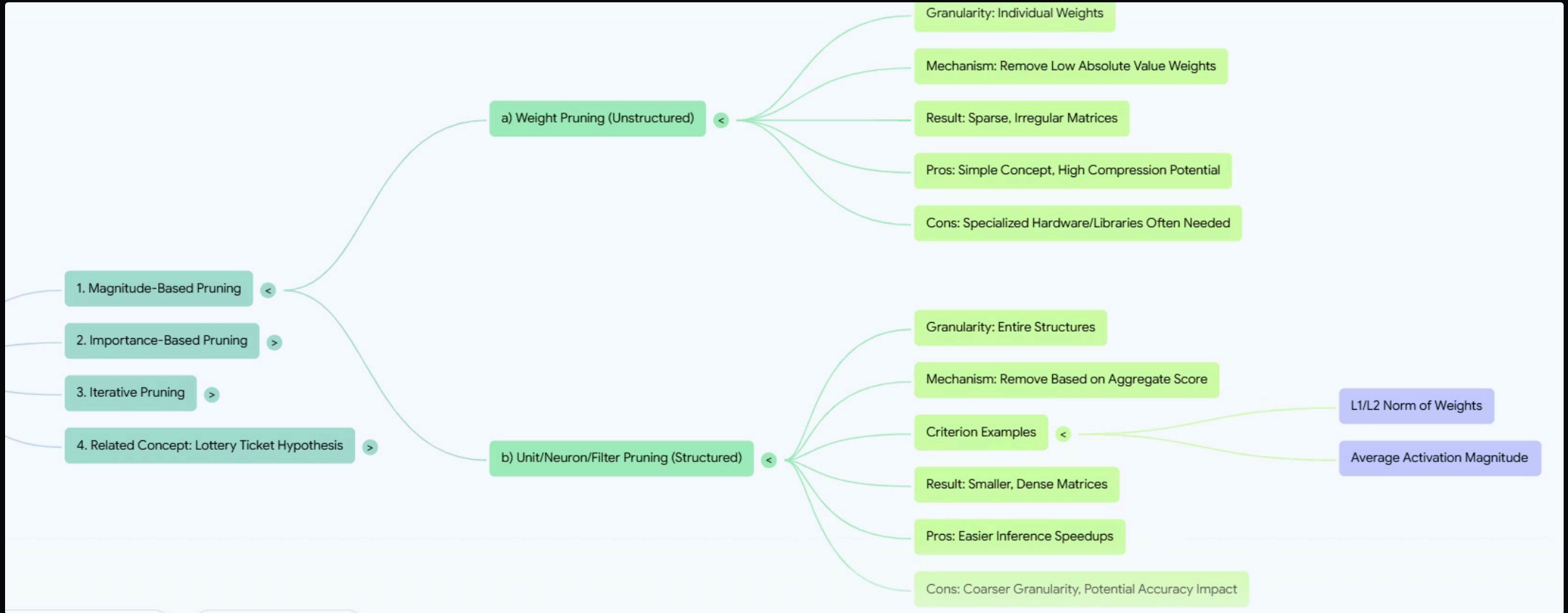
Our research shows that iterative pruning combined with fine-tuning yields the best results, allowing models to adapt to their reduced parameter space while preserving critical functional pathways.



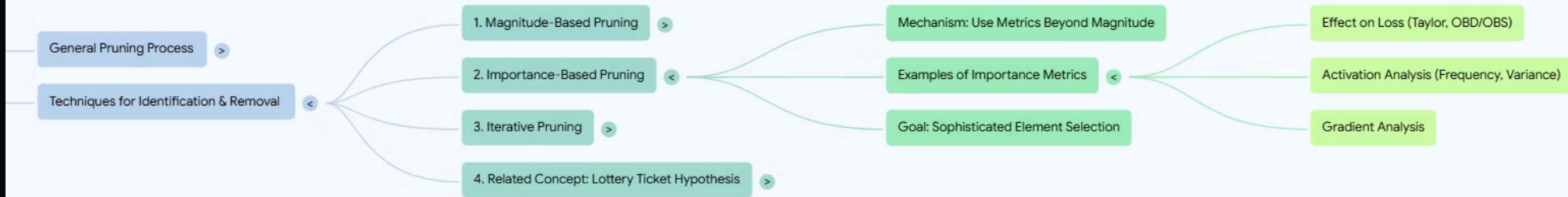
# Pruning Process



# Pruning Technique: Magnitude-Based



# Pruning Technique: Importance-Based





# Pruning Technique: Iterative-Based

1. Magnitude-Based Pruning

>

2. Importance-Based Pruning

>

3. Iterative Pruning

<

4. Related Concept: Lottery Ticket Hypothesis

>

Process Schedule/Approach

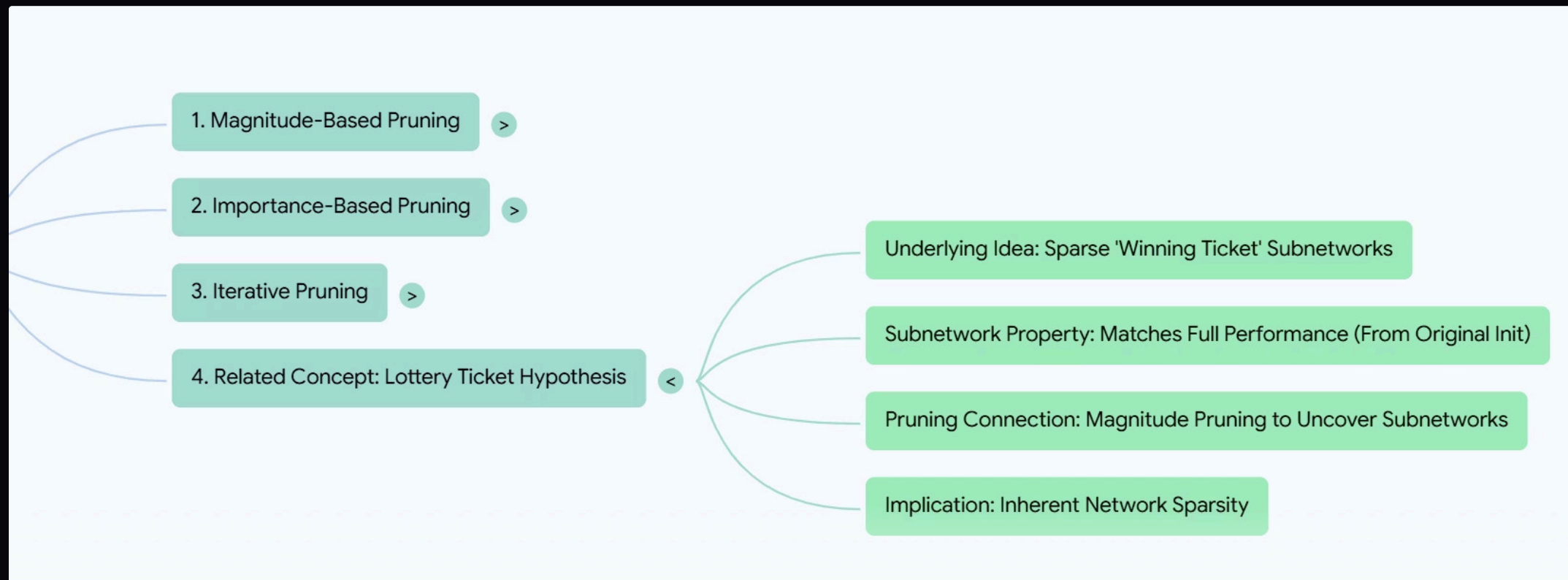
Mechanism: Repeat [Prune -> Fine-tune]

Goal: Reach Target Sparsity Gradually

Benefit: Better Accuracy for Sparsity

Trade-off: More Time-Consuming

# Pruning Technique: Lottery Ticket Hypothesis



# Knowledge Distillation for Edge



## Teacher Model

Large, high-accuracy model with complex capabilities



## Knowledge Transfer

Soft targets encoding learned relationships



## Student Model

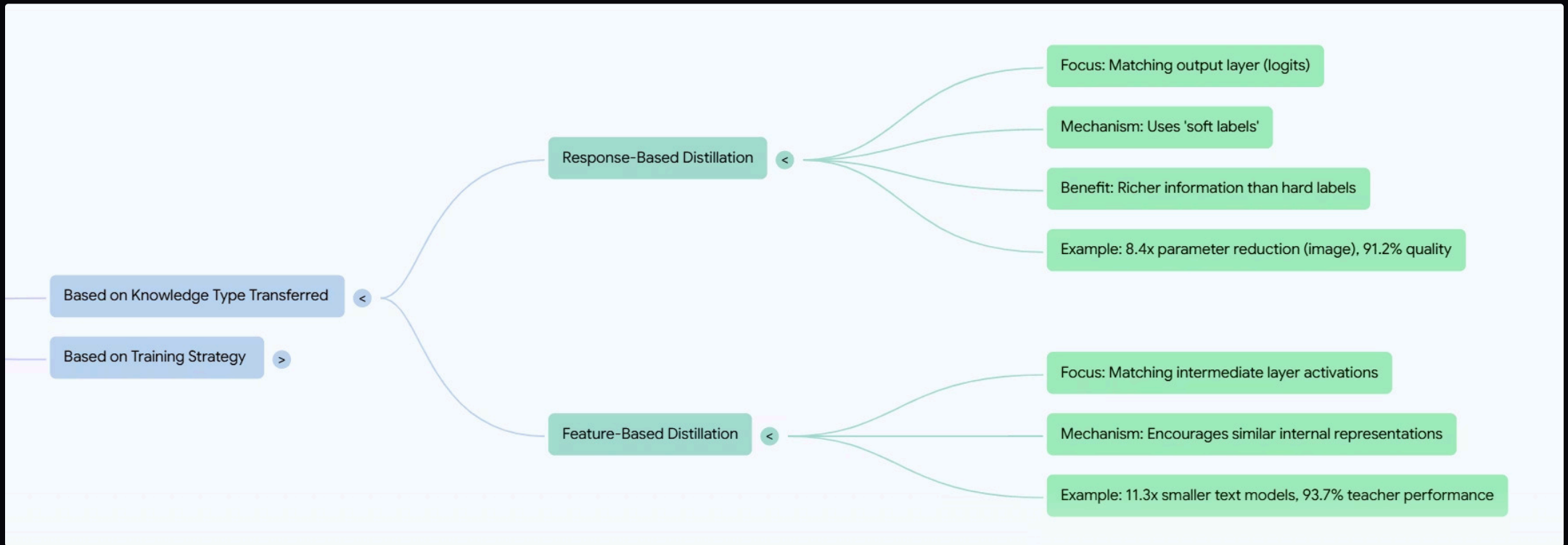
Lightweight model for edge deployment

Knowledge distillation accelerates inference by transferring capabilities from comprehensive teacher models to lightweight student models engineered specifically for edge environments. This approach allows small models to benefit from the learning of much larger networks.

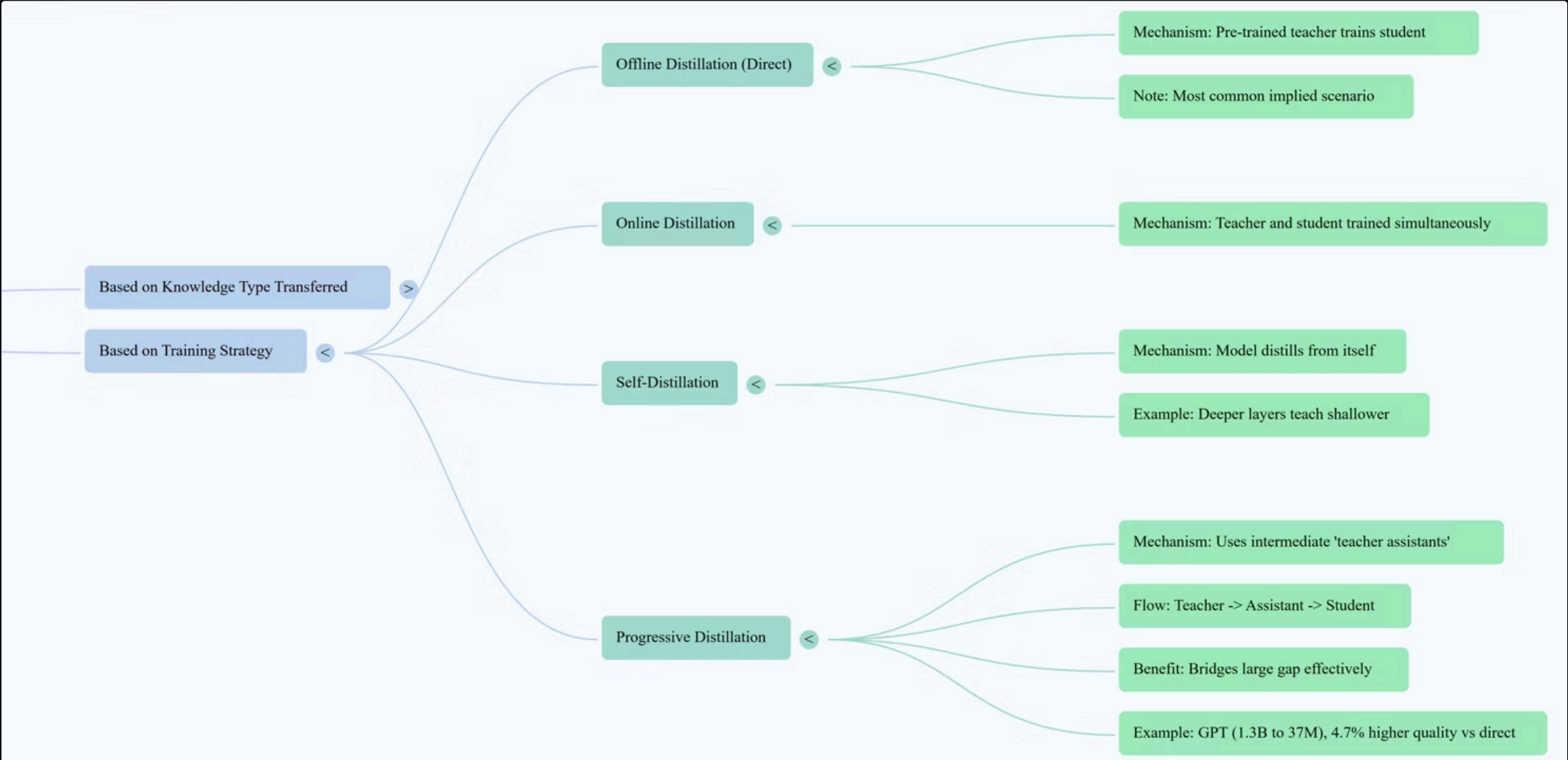
Rather than simply training on hard labels, student models learn from the probability distributions of teacher models, capturing nuanced relationships between classes and enabling better generalization despite their reduced capacity.



# Knowledge Distillation: Based on Type Transferred



# Knowledge Distillation: Based on Training Strategy



# Edge-Optimized Model Architectures



## MobileNets

Depthwise separable convolutions reduce computation by 8-9x while maintaining reasonable accuracy for vision tasks.



## ShuffleNet

Channel shuffling and grouped convolutions optimize for devices with limited computing resources.



## SqueezeNet

Fire modules with squeeze layers reduce parameter count while preserving performance on classification tasks.



## EfficientNet

Compound scaling method optimizes width, depth, and resolution for superior efficiency-accuracy tradeoff.

Purpose-built architectures designed specifically for resource-constrained environments deliver substantially better performance than simply shrinking standard models. These specialized architectures incorporate operations like depthwise separable convolutions that dramatically reduce computational requirements.

The latest generation of edge-optimized models achieves performance approaching that of models 5-10x their size through architectural innovations rather than simply scaling down existing designs.



# Case Study: Edge Voice Assistant

98.2%

Wake Word Accuracy

False activation rate below 0.5%

87.3%

Command Recognition

Across various acoustic environments

76MB

Model Size

Full wake+command system footprint

85ms

Response Latency

End-to-end processing time

Our edge-deployed voice assistant demonstrates how optimization techniques deliver compelling real-world performance. Using a combination of neural network pruning and 8-bit quantization, we reduced the model footprint by 73% compared to the baseline while maintaining excellent accuracy.

Knowledge distillation from a larger teacher model enabled our compact architecture to achieve wake word detection accuracy of 98.2% with a false activation rate below 0.5%. The entire system operates with an end-to-end latency of 85ms, well below the 100ms threshold for perceived real-time interaction.



# Key Takeaways & Implementation Roadmap

## Establish Performance Targets

Define clear latency, accuracy, and resource budgets for your edge AI application. Consider both average and worst-case scenarios.

## Apply Integrated Optimization

Combine pruning, quantization, and distillation techniques rather than relying on a single approach. The compound effect delivers significantly better results.

## Test on Target Hardware

Benchmark on actual deployment devices, as emulators often miss critical performance characteristics. Profile to identify bottlenecks.

## Iterate with Real-World Data

Continuously refine your models based on edge-collected data to adapt to deployment conditions not seen during initial training.

Successfully deploying GenAI at the edge requires a systematic approach that considers the entire optimization pipeline. By leveraging the techniques presented today, you can achieve breakthrough performance metrics on resource-constrained devices.

Remember that edge AI optimization is an iterative process. Start with establishing clear performance targets, apply integrated optimization techniques, thoroughly test on target hardware, and continuously refine your models with real-world data. This methodical approach will help you navigate the unique challenges of edge AI development.

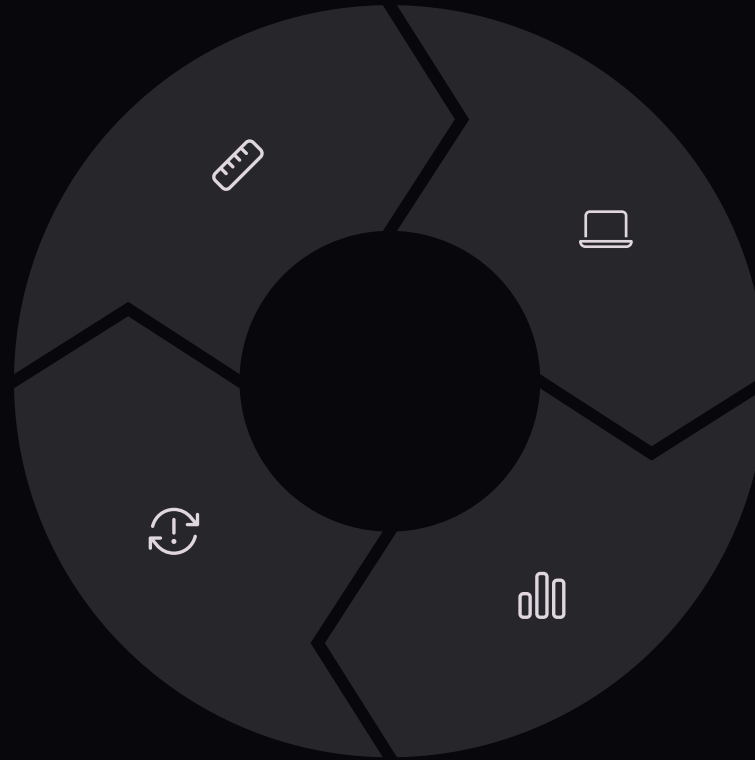
# Benchmarking for Consistent Performance

## Define Metrics

Establish key performance indicators including latency, throughput, accuracy, and energy consumption

## Iterate Optimization

Use benchmark results to guide further refinements of model architecture and parameters



## Test Environment

Create standardized testing procedures that simulate real-world deployment conditions

## Measure Performance

Collect comprehensive metrics across different hardware targets and workloads

Systematic benchmarking is essential for validating optimization results and ensuring consistent performance across diverse deployment scenarios. Our framework includes tools for automated testing across multiple hardware targets and operating conditions.

# Implementation Roadmap



## Audit Model Requirements

Identify target hardware constraints and performance needs



## Prototype Optimization Pipeline

Test optimization techniques individually to assess impact



## Implement Combined Approach

Integrate pruning, quantization, and distillation in sequence



## Deploy and Monitor

Release optimized models with telemetry for continuous improvement

Start your edge AI optimization journey with a comprehensive audit of your model requirements and hardware constraints. Develop a systematic pipeline that combines multiple optimization techniques, with careful validation at each stage. Implement a monitoring system after deployment to gather real-world performance data and guide future improvements.

Thank you