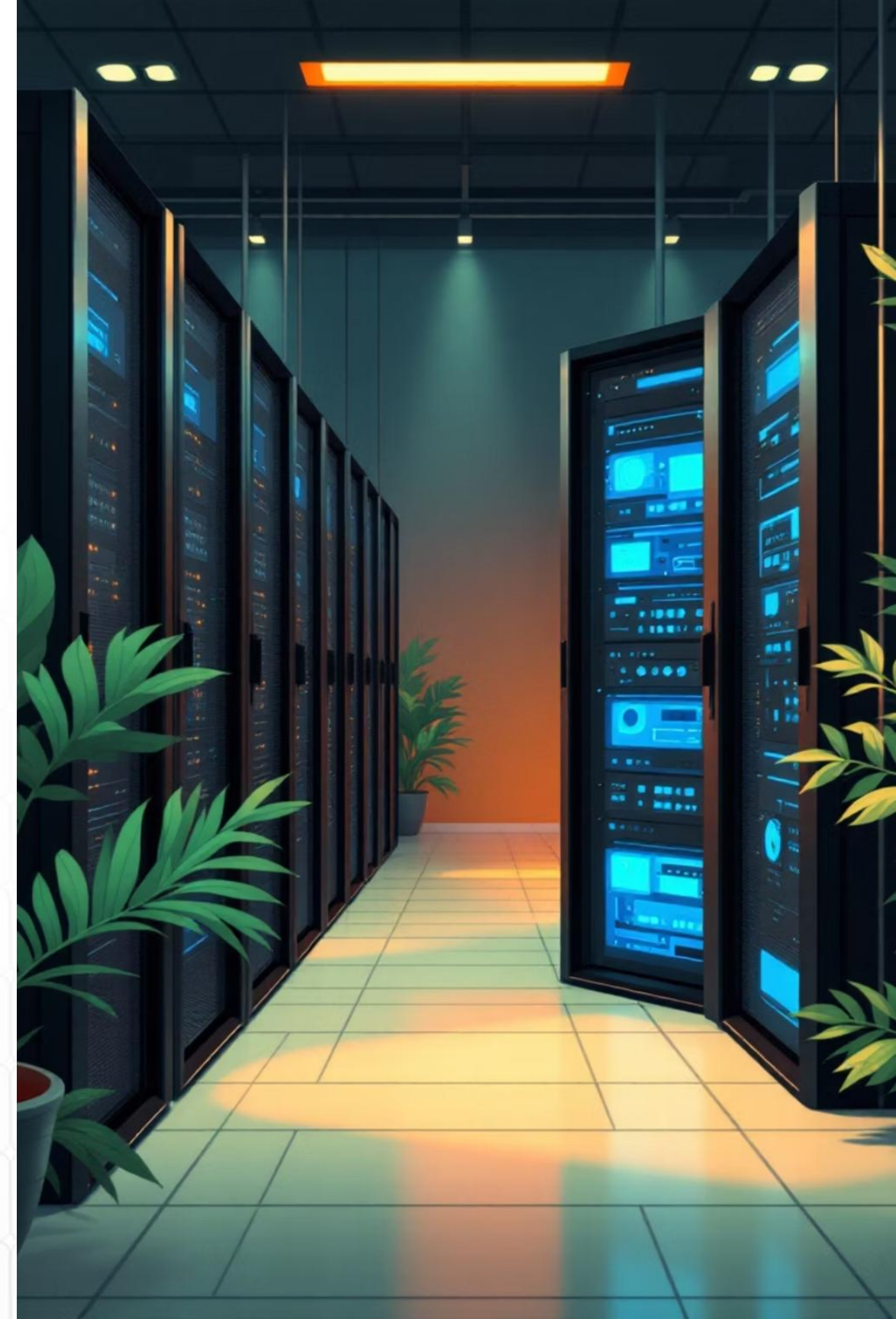


# Quantizing AI: Cutting Model Size Without Cutting Performance

Understanding how software quantization techniques and hardware acceleration work together to deliver high-performance, cost-effective AI deployment on Intel Xeon processors.



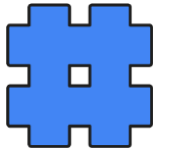
# Monolith to Service-Oriented Architecture (SOA)

Quantization is a software technique that reduces the numerical precision of AI model parameters and activations, transforming high-precision formats like FP32 into more efficient representations such as INT8 or BF16.

## The Core Tradeoff

By accepting a slight loss in model accuracy, quantization delivers significant gains in inference speed, memory efficiency, and computational cost - making AI deployment more practical and scalable.





# Types of Quantization Techniques



## Downcasting

Converts model weights to lower precision formats like BF16, with computation and activations remaining in BF16 throughout inference.

- Approximately 50% memory savings vs FP32
- Simpler implementation with minimal accuracy impact



## Linear Quantization

Stores weights as INT8 integers, then dequantizes back to FP32 during computation using scale and zero-point adjustments.

- Up to 75% memory reduction
- Preserves accuracy closer to original model





# Why Quantization Matters for AI Deployment

## Accelerated Performance

Dramatically speeds up model inference and training operations, reducing latency for real-time AI applications.

## Resource Efficiency

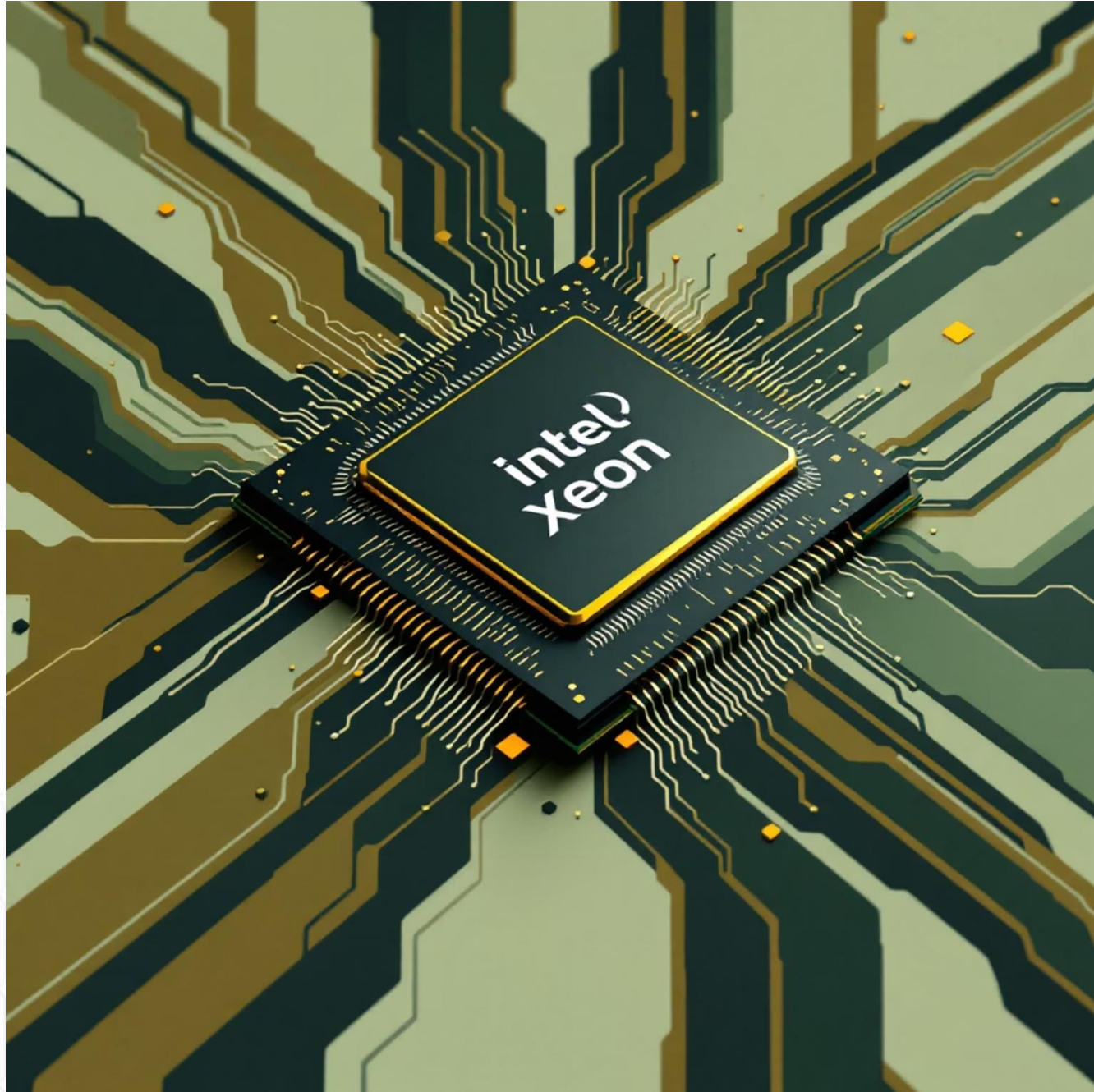
Enables deployment on resource-constrained hardware like edge devices, mobile platforms, and cost-optimized cloud instances.

## Optimized Power & Memory

Reduces memory bandwidth requirements and power consumption, lowering operational costs across data centers.

## Broad Application Impact

Powers faster AI across NLP, computer vision, recommendation systems, and generative AI workloads at scale.



## Intel Advanced Matrix Extensions (AMX)

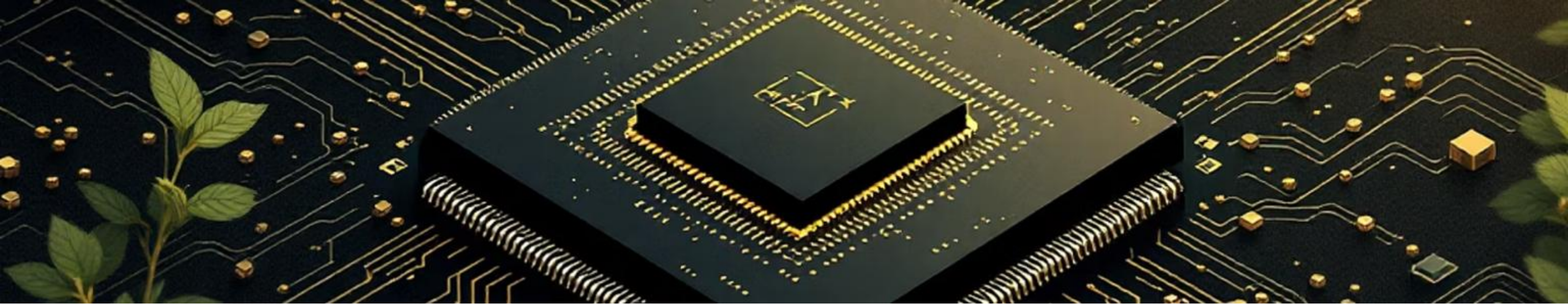
Intel AMX is a **hardware accelerator** integrated directly into Intel Xeon Scalable processors, specifically engineered to accelerate the matrix multiplication operations that are fundamental to AI workloads.

### Supported Precision Formats

- **BF16** for both training and inference
- **INT8** for optimized inference operations

AMX brings GPU-like acceleration capabilities to CPUs, eliminating the need for discrete accelerators in many AI deployment scenarios.





# How Intel AMX Achieves High Performance

01

## Tile Architecture

Eight 2D registers, each storing 1KB of data, enable large matrix chunks to be processed simultaneously rather than element-by-element.

02

## TMUL Engine

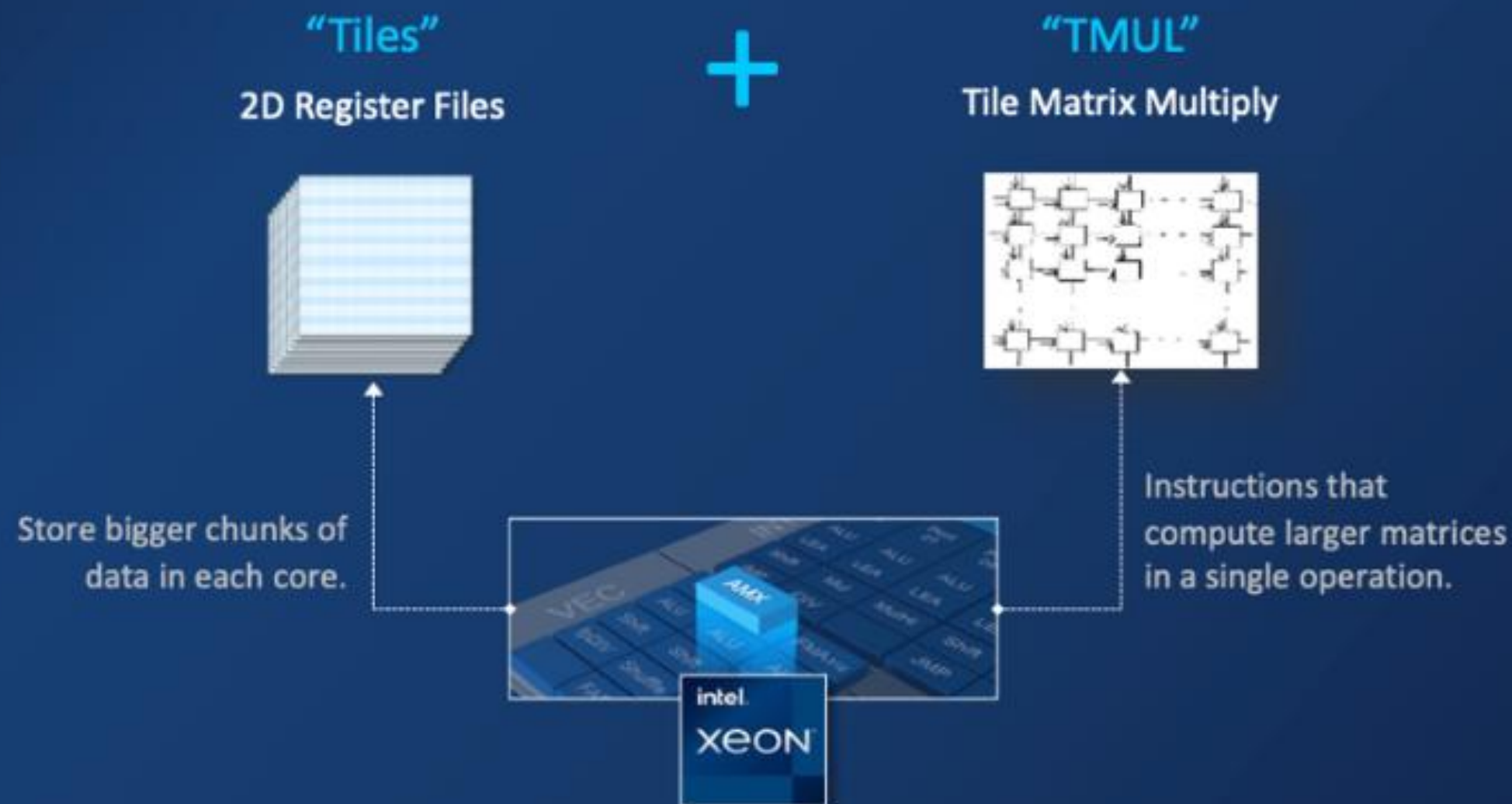
The Tile Matrix Multiplication accelerator engine executes massive matrix operations in single instructions, drastically improving throughput.

03

## Optimized Data Flow

By keeping large matrices in tiles, AMX minimizes memory access overhead and maximizes computational efficiency for AI kernels.

# Intel® Advanced Matrix Extensions (AMX) built-in for AI



# Quantization vs Intel AMX: Complementary Technologies

Aspect	Quantization	Intel AMX
Nature	Software technique	Hardware accelerator
Purpose	Reduce precision to shrink model size and speed up computation	Accelerate matrix math operations for AI inference on CPUs
Data Types	Converts from FP32 to INT8 or BF16	Specialized execution for INT8 and BF16 formats
Role	Prepares model for efficient execution	Executes quantized model with hardware acceleration

These technologies work in tandem: quantization optimizes the model, while AMX optimizes execution.



# The Complete Workflow: Quantization + AMX



## Model Training

AI model is trained using high-precision FP32 format to ensure maximum accuracy during learning phase.



## Software Quantization

Quantization toolkit converts trained weights to INT8 or BF16, significantly reducing model footprint.



## Hardware Deployment

Quantized model is deployed on Intel Xeon CPU equipped with AMX accelerator capabilities.



## Accelerated Inference

AMX tiles and TMUL engine execute highly efficient matrix multiplications on quantized operations.



## Optimized Results

Faster inference, reduced memory usage, and lower operational costs with minimal accuracy degradation.



# Performance Benefits at Scale

10x

## Performance Improvement

Up to 10x faster execution in AI tasks including speech recognition, natural language processing, and generative AI workloads.

75%

## Memory Reduction

Reduced memory footprint enables larger models or batch sizes on existing hardware infrastructure.

## Lower Total Cost of Ownership

Minimizes power consumption and hardware requirements, eliminating the need for discrete GPU accelerators in many scenarios.

## Simplified Deployment Pipeline

Leverages existing CPU infrastructure for efficient AI execution across computer vision, recommendation systems, and enterprise applications.



# Key Takeaways

## Software Optimization

Quantization compresses AI models through precision reduction, enabling faster inference with minimal accuracy loss.

## Hardware Acceleration

Intel AMX provides specialized CPU hardware designed for efficient low-precision matrix operations.

## Combined Impact

Together, they deliver **scalable, high-performance AI deployment** with significant cost savings and simplified infrastructure requirements.



Completed

Benchmark completed successfully!

Fastest

INT8

11.56 tokens/sec

Most efficient

INT8

2.0 GB memory

Models tested

4

Completed

Max speedup

1.61x

vs FP32 baseline

Performance comparison

Model		Size (MB)	Time (s)	Speed (tok/s)	Relative
FP32		30,633	35.69	7.17	1.00x
BFloat16		15,317	58.44	4.38	0.61x
BFloat16 + AMX		15,317	23.15	11.06	1.54x
INT8	INT8	2,005	22.14	11.56	1.61x