

Azure Maia 200

The Inference Efficiency Frontier

A full-stack inference play across silicon, memory hierarchy, and scale-up networking.



- FP4 TENSOR COMPUTE
- INFERENCE-FIRST ON-CHIP FABRIC
- BALANCED MEMORY HIERARCHY
- INTEGRATED NIC TWO-TIER SCALING
- VARIABLE POWER BALANCING

01 THE CORE THESIS

Redefining the Inference Efficiency Frontier

AI inference is an efficiency frontier defined by the intersection of accuracy, latency, throughput, cost, and energy. Microsoft positions Maia 200 as part of a portfolio approach — a heterogeneous Azure fleet that matches the best silicon to the customer workload.

Different applications land at different points on this curve.

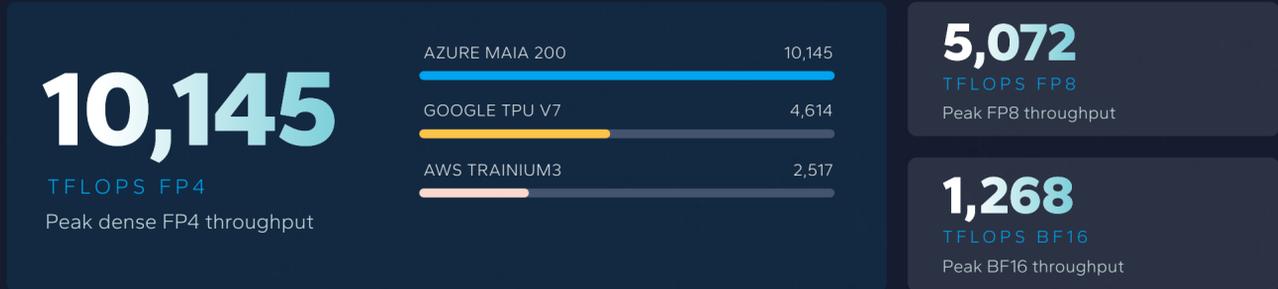
Interactive workloads tend to be latency sensitive, while batch and throughput-oriented workloads often prioritize cost and sustained token output.

- ACCURACY
- LATENCY
- THROUGHPUT
- COST
- ENERGY

02 THE SILICON BET

FP4 Throughput Leadership, Built for Inference Economics

Maia 200 makes a deliberate bet on narrow-precision formats as the path to scaling token generation efficiently. Peak FP4 throughput positions it as a serious competitor to modern AI accelerators in the narrow-precision inference era.



Narrow-precision formats, increasingly including FP4, are becoming central to scaling token generation efficiently.

- Narrow-precision tensor compute
- Inference-focused on-chip fabric
- Integrated NIC designed to enable a two-tier scale-up approach

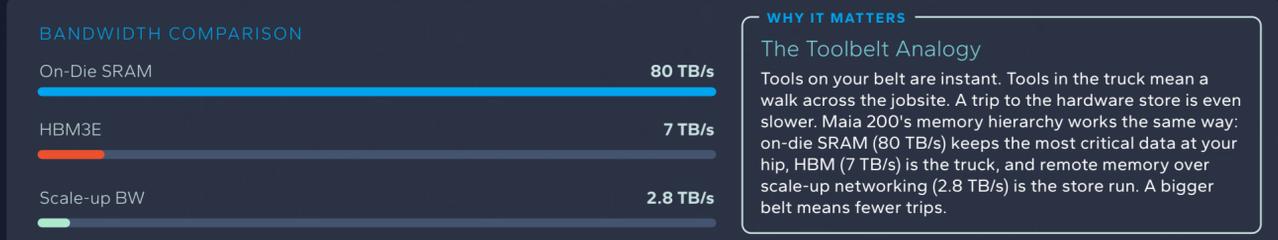
Peak specs are directional indicators; sustained throughput and token economics depend on models, quantization, and scaling efficiency.

03 THE MEMORY WALL

A Tiered Hierarchy for Effective Bandwidth

Large SRAM and explicit locality controls help keep high-value data on die, reduce repeated HBM accesses, and improve sustained throughput. The practical goal: keep data local when possible, fall back to HBM, and use remote HBM through scale-up paths when needed.

<p>TIER 1 • ON-DIE SRAM</p> <p>Ultimate Locality</p> <p>272</p> <p>MB On-Die SRAM</p> <p>80</p> <p>TB/s SRAM Bandwidth</p> <p><i>Keep high-value data on die, reduce repeated HBM accesses</i></p>	<p>TIER 2 • HBM3E</p> <p>Primary Model Storage</p> <p>216</p> <p>GB HBM3E</p> <p>7</p> <p>TB/s HBM Bandwidth</p> <p><i>Primary working set and model weight storage</i></p>
---	--



Inference is frequently memory bound rather than compute bound.

04 NETWORKING

Reliable Chip-to-Chip Communication at Scale

Maia 200 pursues a scale-up approach based on observed inference patterns, pairing an integrated NIC with commodity Ethernet to reduce cost and power while preserving communication performance across larger clusters.

<p>2.8</p> <p>TB/s Bidirectional</p> <p>Scale-up bandwidth per accelerator</p>	<p>6,144</p> <p>Accelerators</p> <p>Maximum scale-up domain size</p>
---	---

- Integrated NIC reduces cost & power
- Two-tier scale-up on commodity Ethernet
- Predictable collective operations
- Improved scheduling & bin packing

Larger scale-up domains can improve scheduling, bin packing, and utilization across diverse workloads — supporting common parallelism patterns in modern inference with strong tensor-parallel performance.