

Notices and Disclaimers

Performance varies by use, configuration and other factors. Learn more at www.Intel.com/PerformanceIndex.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Intel contributes to the development of benchmarks by participating in, sponsoring, and/or contributing technical support to various benchmarking groups, including the BenchmarkXPRT Development Community administered by Principled Technologies.

Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

Some results may have been estimated or simulated.

Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

All product plans and roadmaps are subject to change without notice.

Statements in this document that refer to future plans or expectations are forward-looking statements. These statements are based on current expectations and involve many risks and uncertainties that could cause actual results to differ materially from those expressed or implied in such statements. For more information on the factors that could cause actual results to differ materially, see our most recent earnings release and SEC filings at www.intc.com.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

IXPUG 2023

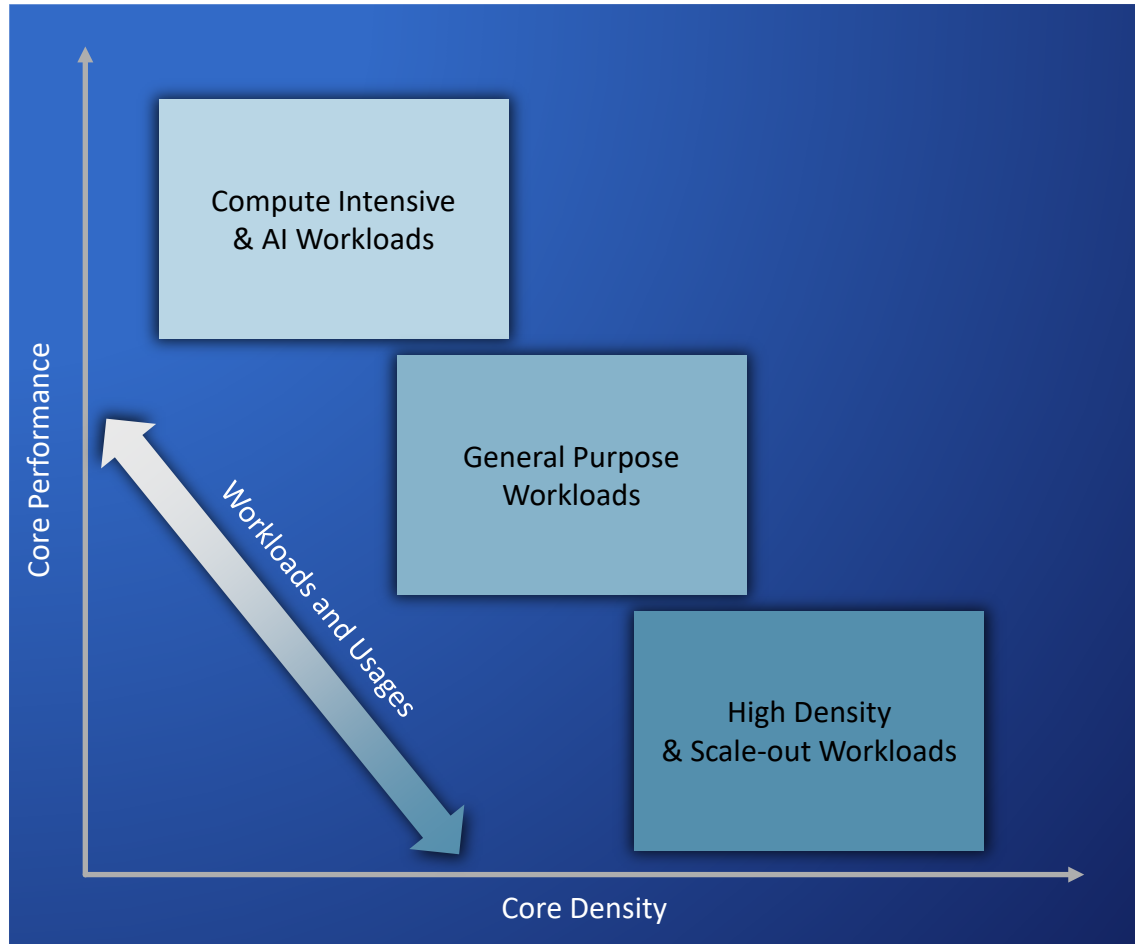
Intel®'s Next Generation Data Center CPU Architecture and the Performance-Core/Efficient-Core Processors Built on this Common Foundation

Krishna Vinod

Principal Engineer, Datacenter Processor Architecture, Intel



Data Center Requirements Are Expanding



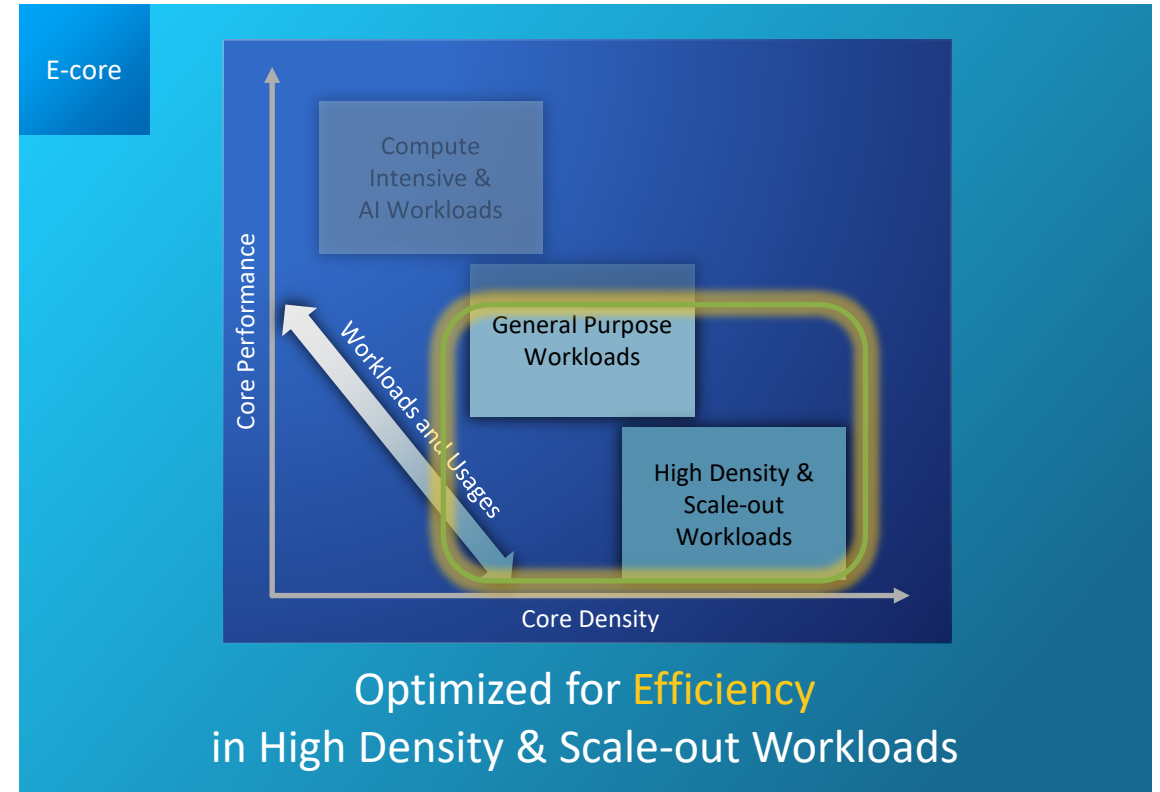
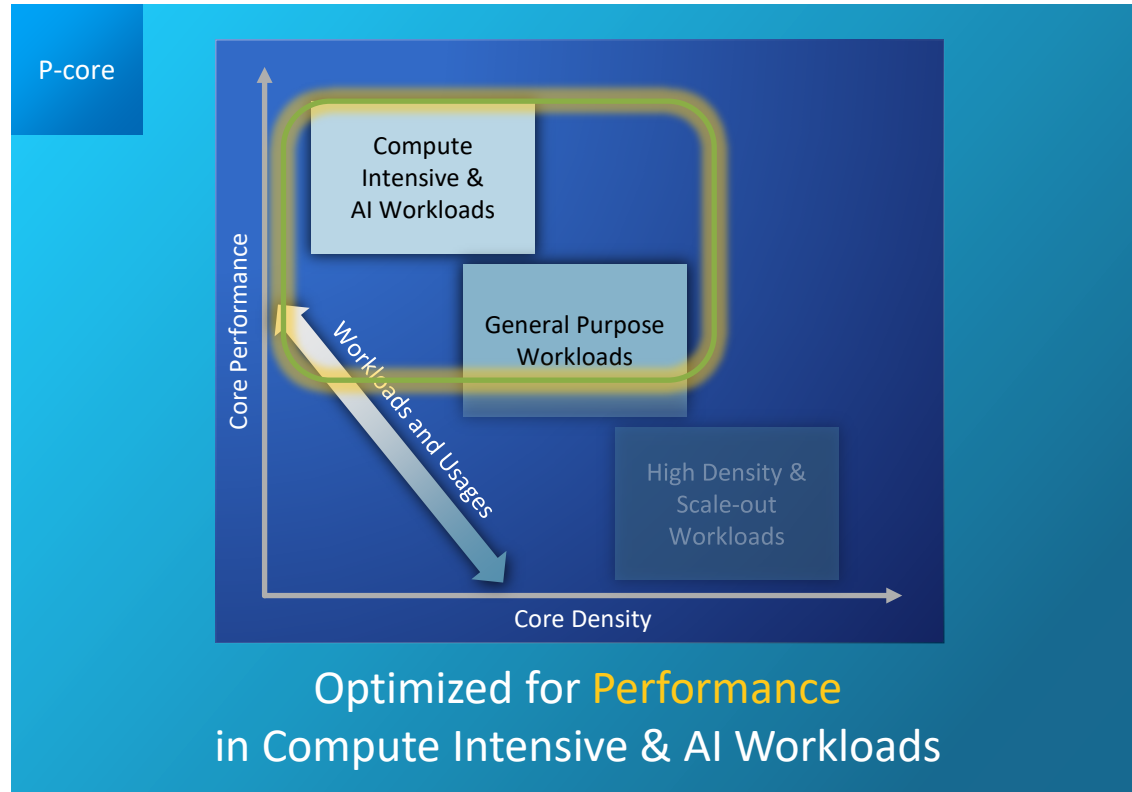
Continuing demand for CPU cores that deliver high performance per core

Growing demand for cores that deliver better density and the best performance per watt at specific performance levels

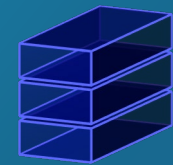
Optimizing for specific types of workloads requires trade-offs between core performance & core density

Expanding deployment models demanding increased power, IO & memory bandwidth

Expanded Xeon Portfolio with Optimized Processors



Common Platform Foundation
& Shared Software Stack



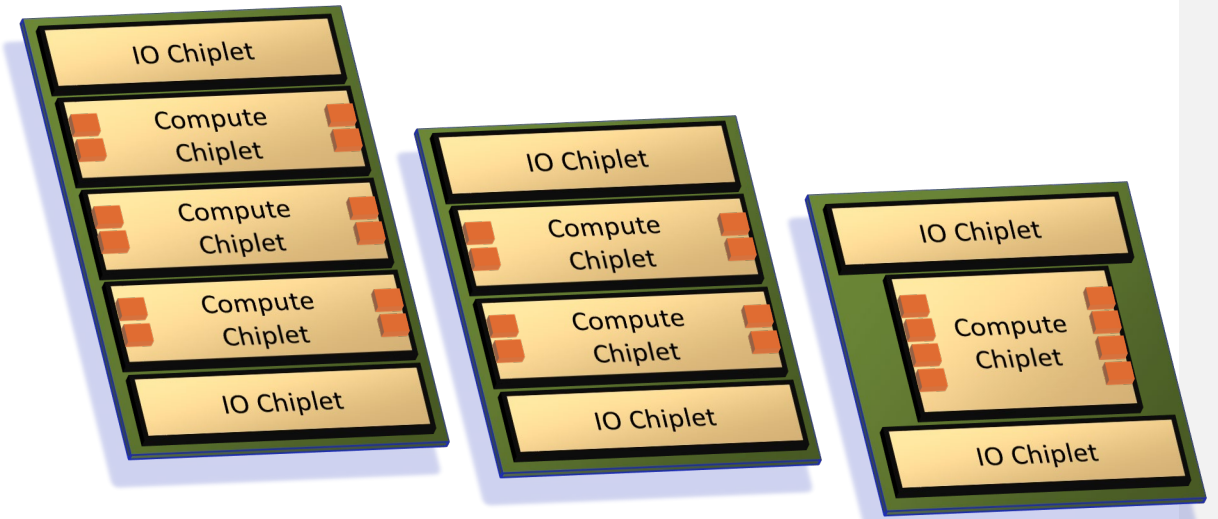
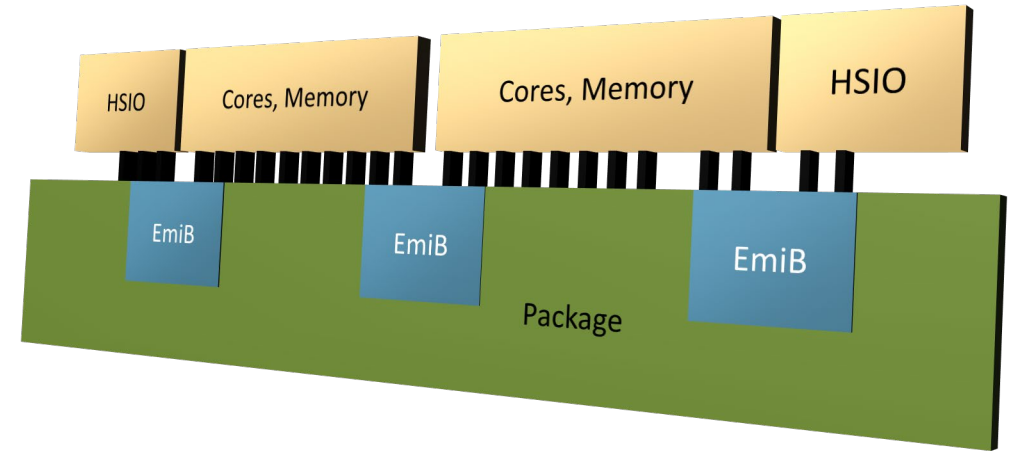
Modular SoC Architecture

Component Details

- Separate Compute and IO silicon chiplets
- EmiB packaging -> high bandwidth / low latency
- Modular die fabric enables flexible construction
- Common IP, Firmware, OS, platform ingredients

Platform Details

- Scalability: 1S-8S (P-core), 1S-2S (E-core)
- Supports range of core counts and thermals
- Memory: up to 12-channels DDR/MCR, 1-2DPC
- Advanced I/O up to 136 lanes PCIe 5.0/CXL 2.0, up to 6 UPI links (144 lanes)
- Self boot

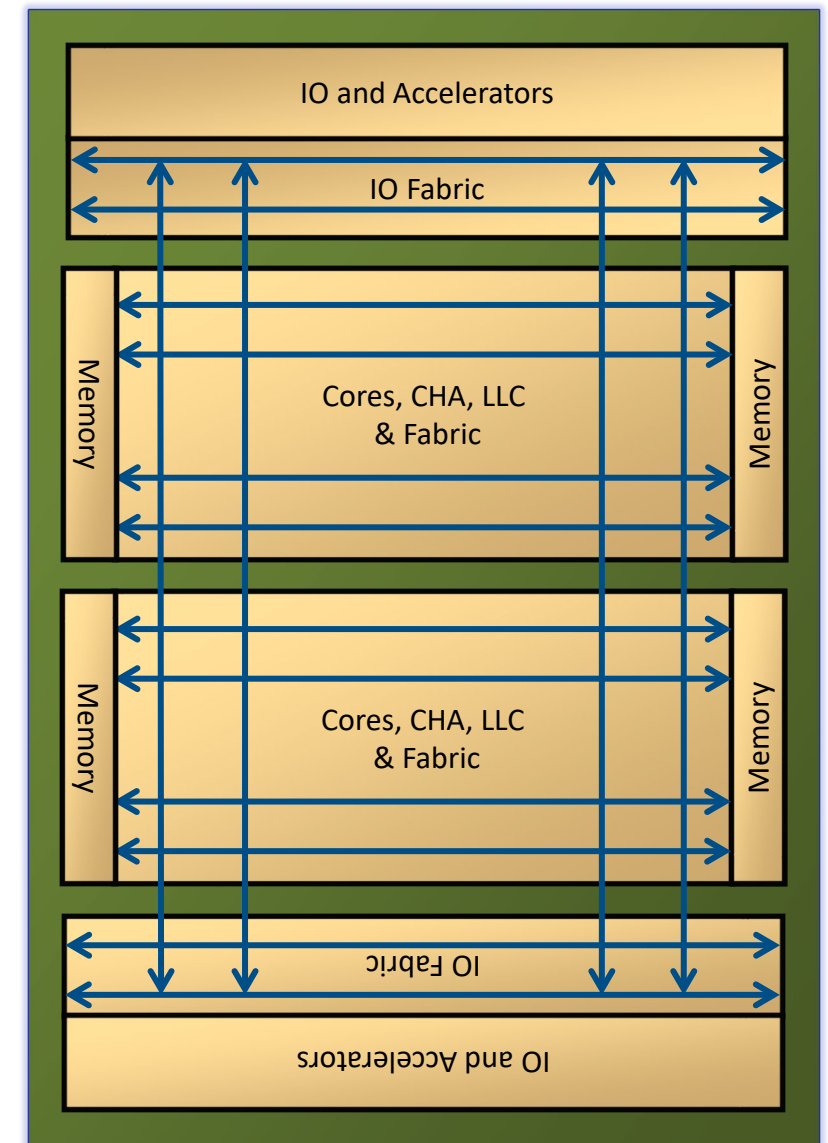


intel.
XEON

Scalability and Flexibility to Deliver a Wide Range of Optimized Products

Modular Mesh Fabric

- Logically Monolithic Mesh enables direct access between agents within the socket
- Last level cache is shared amongst all cores and can be partitioned into per-die sub-NUMA clusters
- EmiB technology extends the high-speed fabric across all die in the package
- Modularity and flexible routing allows per-die definition of rows and columns
- Fabric distributes IO traffic across multiple columns to ease congestion
- Global infrastructure is modular/hierarchical



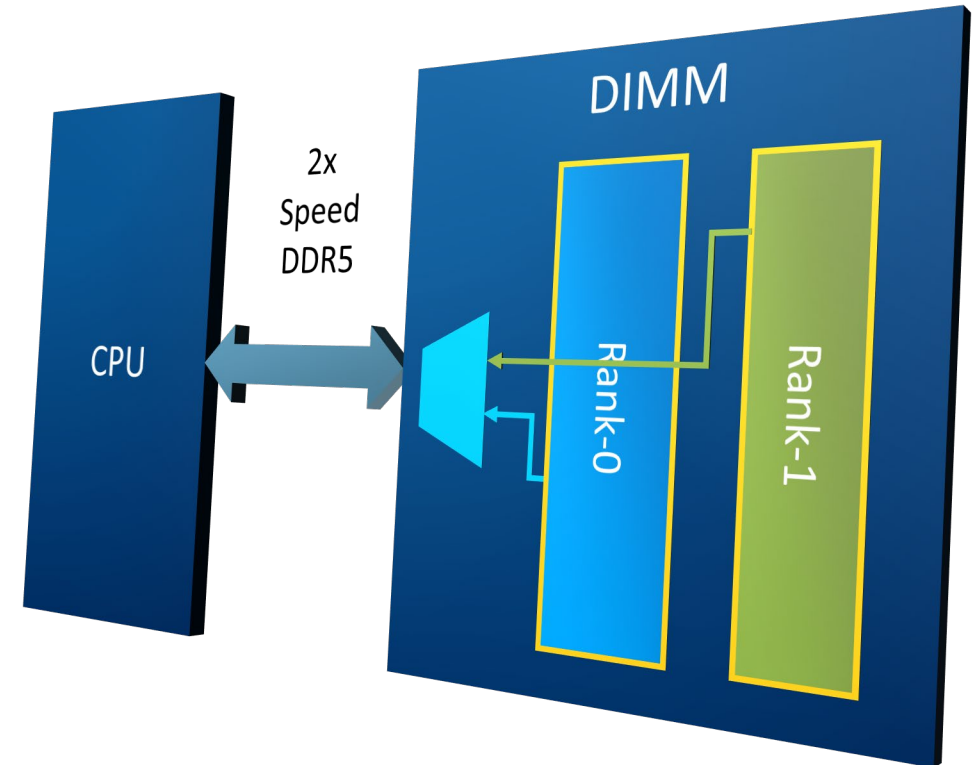
Multiplexed Combined Rank DIMMs

Problem

- New workloads benefiting from more BW/thread
- Adding DIMMs to deliver BW is expensive
- Prefer to increase memory BW @ iso capacity/core

MCR Details

- Fully platform compatible with DDR5 RDIMMs
- Mux merges two ranks at 1x speed into channel running at 2x speed
- Each Rank operates at half the channel speed
- JEDEC standardization in progress today
- Achieving 30-40% more BW than RDIMM



intel.
XEON

MCR DIMMs add bandwidth without the cost of excess capacity

CXL Attached Memory

SW tiering with independent regions

- Interleaving within regions

HW tiering through heterogenous interleaving

- Single native DDR + CXL memory region

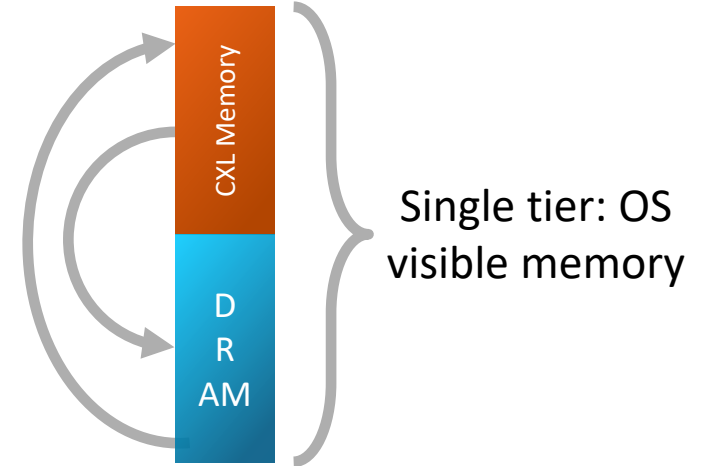
HW tiering with Intel Flat Memory Mode

- Cacheline granular movement
- All memory is addressable (vs. cached)
- Hot lines stay in lower latency memory

HW tiering is fully SW transparent

Flat Memory Mode (1:1 ratio)

50% CXL
memory
+
50% DRAM



intel.
XEON

CXL Memory is a cost efficient, flexible, first-class solution

I/O Die Architecture

Universal IO stacks

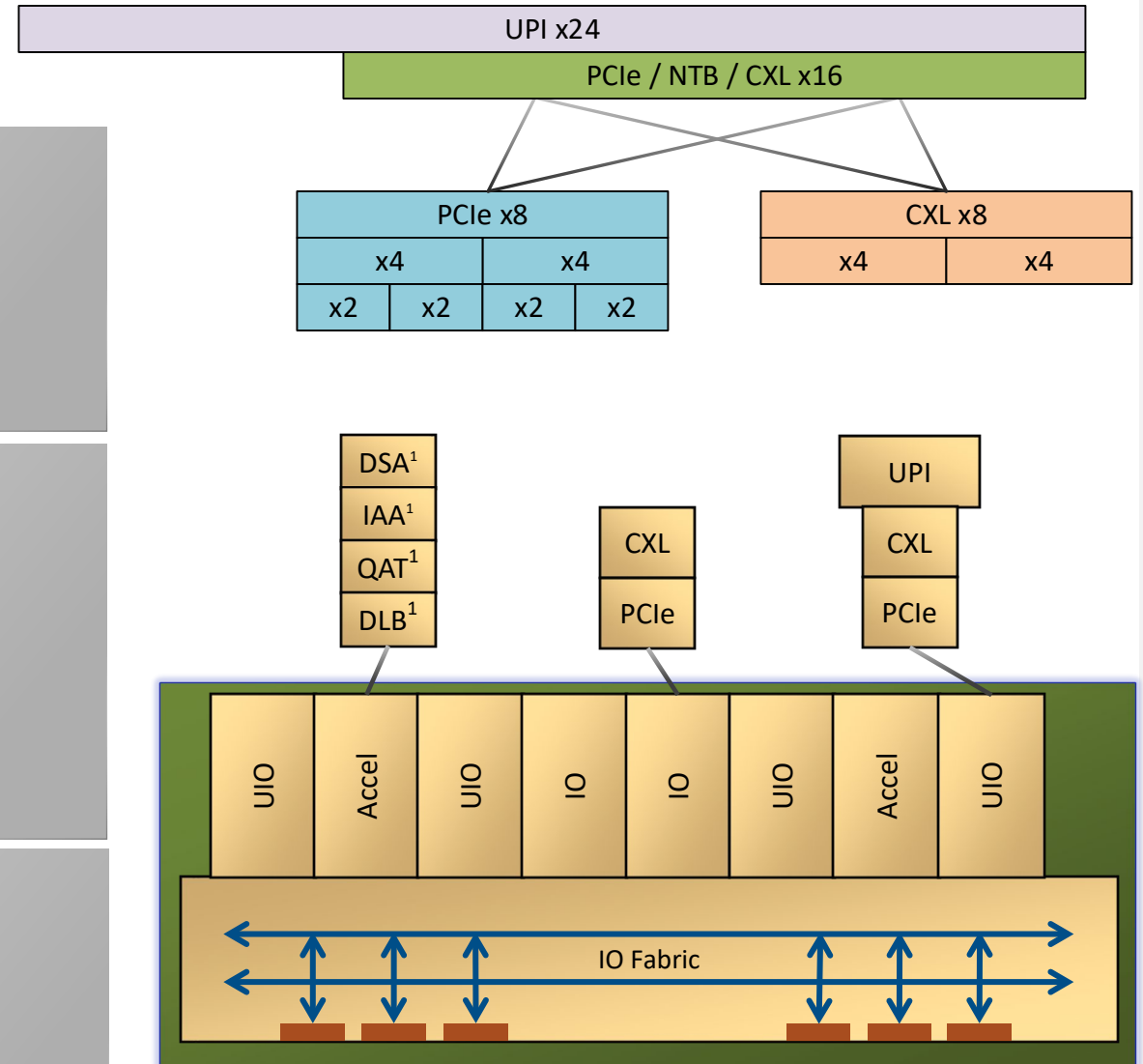
- Each port can operate as UPI (x24), PCIe(x16), CXL(x16)
- PCIe and CXL can be intermixed at x8 granularity
- Integrated workload accelerator engines

New IO Capabilities

- Full CXL 2.0 support - Memory pooling, interleave, port-bifurcation, hot-plug, etc.
- IO RDT – Extends RDT’s resource monitoring and control to I/O devices and channels
- Secure Interconnect – UPI/PCIe/CXL Link Encryption

Enhanced IO performance

- UPI @ 24GT/s. w/6-links = 1.8x prior gen
- UPI Affinity – keeps traffic closer to the cluster
- Accelerator interface = 64B/cycle = 2x prior gen, 2x DSA
- IO die fabric distributes traffic across all mesh columns



¹Intel® Data Streaming Accelerator (Intel® DSA)
 Intel® In-Memory Analytics Accelerator (Intel® IAA)
 Intel® QuickAssist Technology (Intel® QAT)
 Intel® Dynamic Load Balancer (Intel® DLB)

Compute Die Architecture

Construction

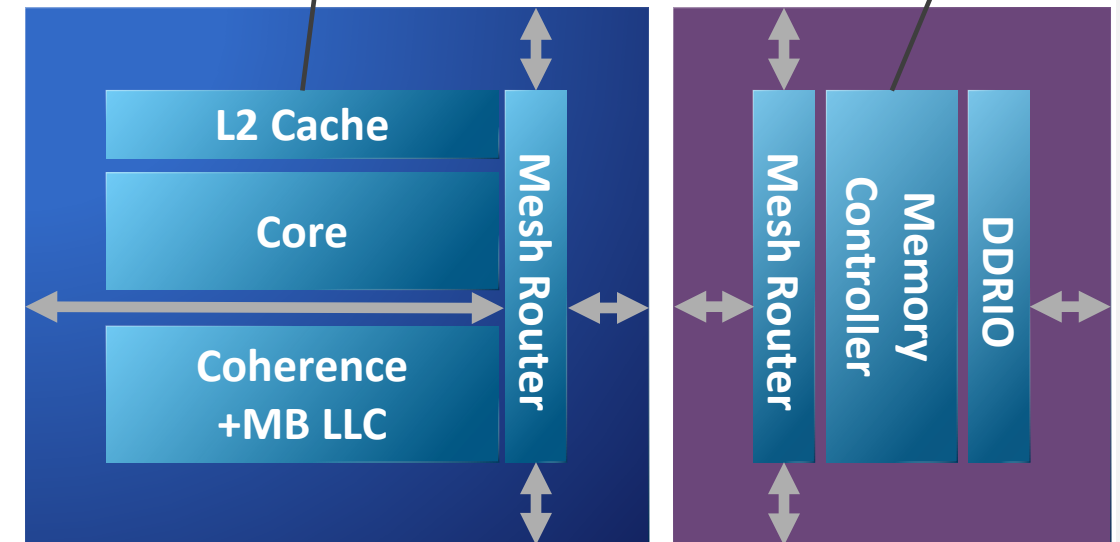
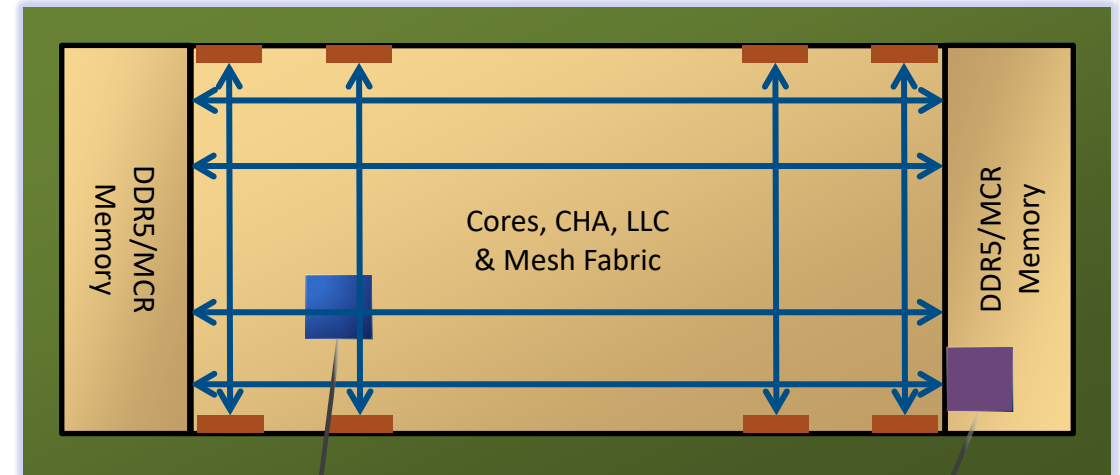
- Higher Performance & Efficiency with Intel 3 process
- Flexible row/column die structure

Core Tile (Mesh building block)

- Contains Core+L2, LLC+SF+CHA slice, mesh fabric interface
- Adaptable to multiple architectures
 - Granite Rapids = P-Core Tiles (Redwood Cove)
 - Sierra Forest = E-Core Tiles (Sierra Glen)

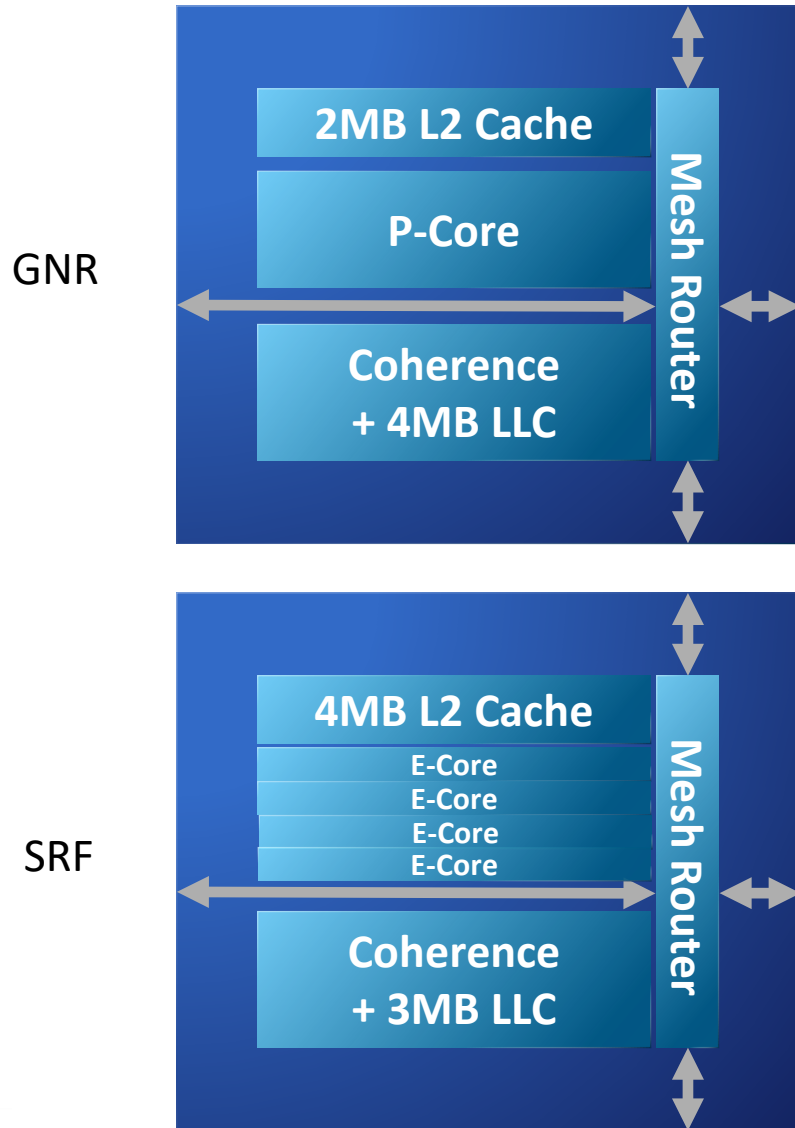
Advanced Memory System

- Common controller/IO supports DDR or MCR memory
- Full support for CXL attached memory



Core Tile

Core Tile Architecture



1 E-Core per module

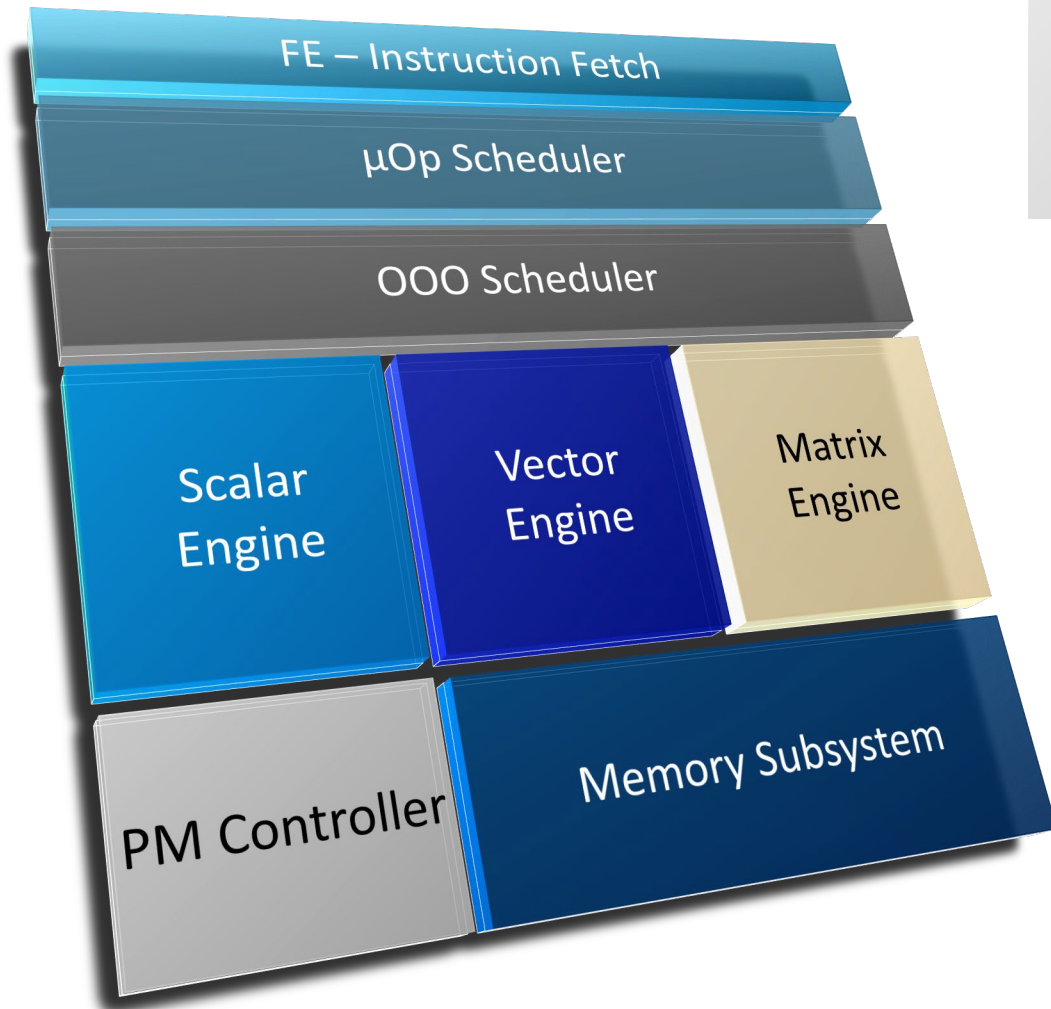
P-Core is Multi-threaded
2 Threads/Core

2 or 4 cores per module
Shared L2 cache
Shared frequency and voltage domain
Shared mesh fabric interface

E-Core is single threaded

LLC slice shared among all cores in socket
High bandwidth pipeline per slice

P-Core: Performance Optimized Core



Proven Intel Xeon Architecture

- Optimized for high performance per core
- Built on the latest Intel 3 process technology
- Improved power efficiency

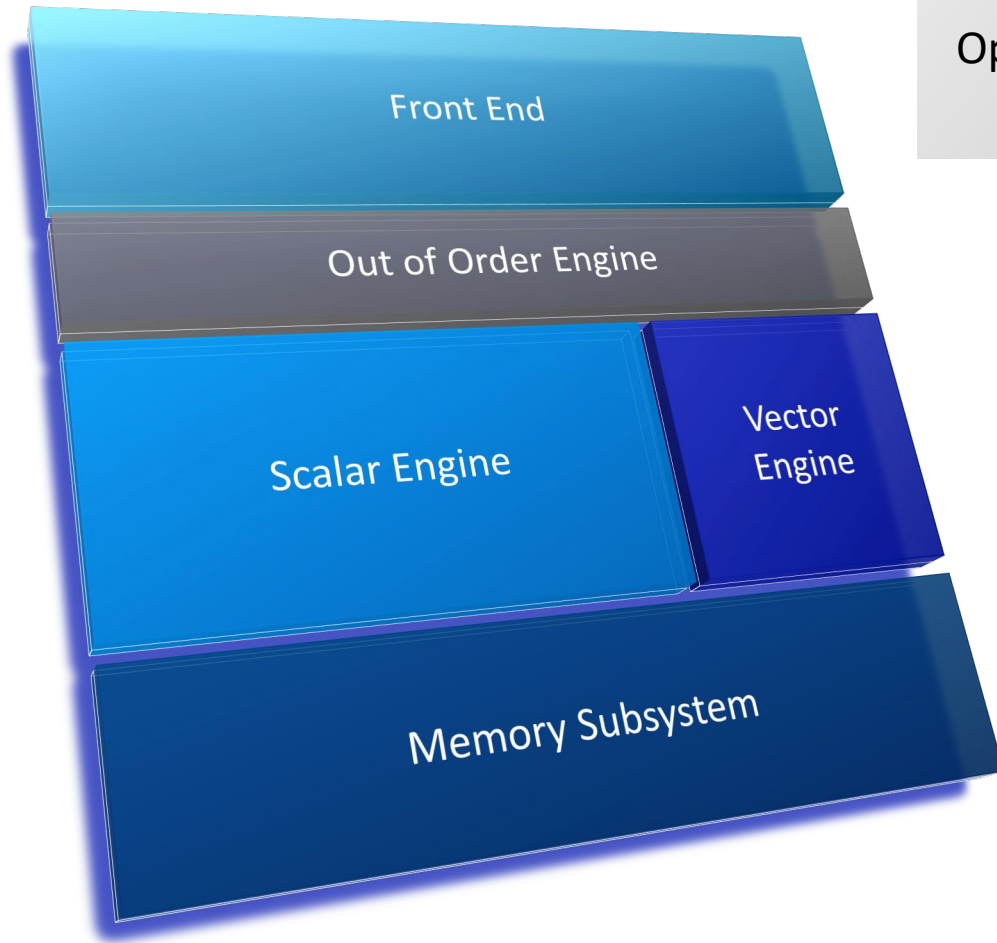
New Software Capabilities

- Matrix Engine supports FP16 for AI/ML
- More memory encryption keys w/ 256b strength
- Code SW prefetch and taken branch hints
- Per-Thread memory bandwidth allocation
- L2 cache allocation and code/data prioritization

Enhanced μArch

- 64KB, 16-way I-cache
- Improved branch predictor and mis-recovery
- 3-cycle FP multiplication
- More outstanding memory requests and prefetch capabilities

E-Core: Efficiency Optimized Core



Designed for scalable throughput performance
Optimized for power and density efficient throughput with:

Deep Front-End

- With on-demand length decode

Wide Back-End

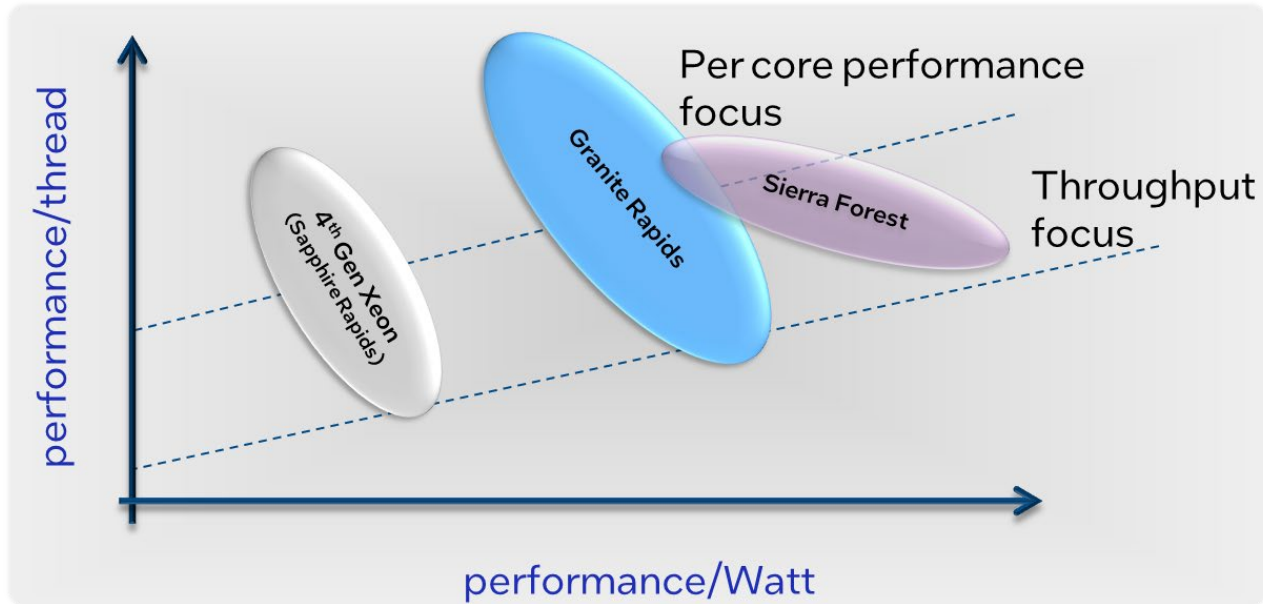
- With many execution ports

Optimized Design

- For latest transistor technologies

Data Center Ready

Deployment View



Granite Rapids

A more balanced machine for Peak performance and AI capabilities

Major improvements for broad datacenter WLS, 2-3x better than 4th gen Xeon

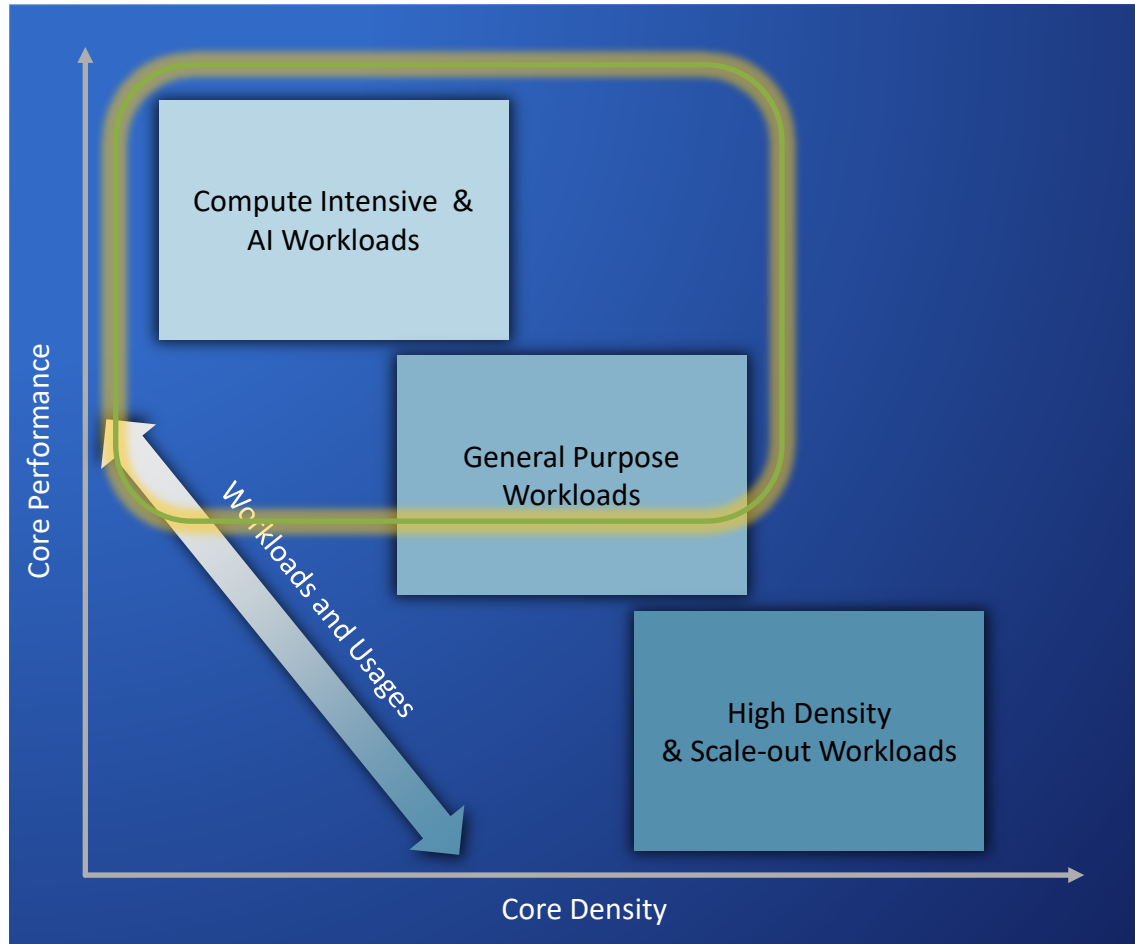
Sierra Forest

Throughput efficiency

2.4x performance per watt over 4th gen Xeon

2.5x rack density over 4th gen Xeon

Harness the Power of Xeon with P-cores



Optimized to deliver significant per-core performance to meet the growing needs of AI & compute intensive workloads

A common HW, FW and software foundation provides compatibility and consistency across a wide range of products

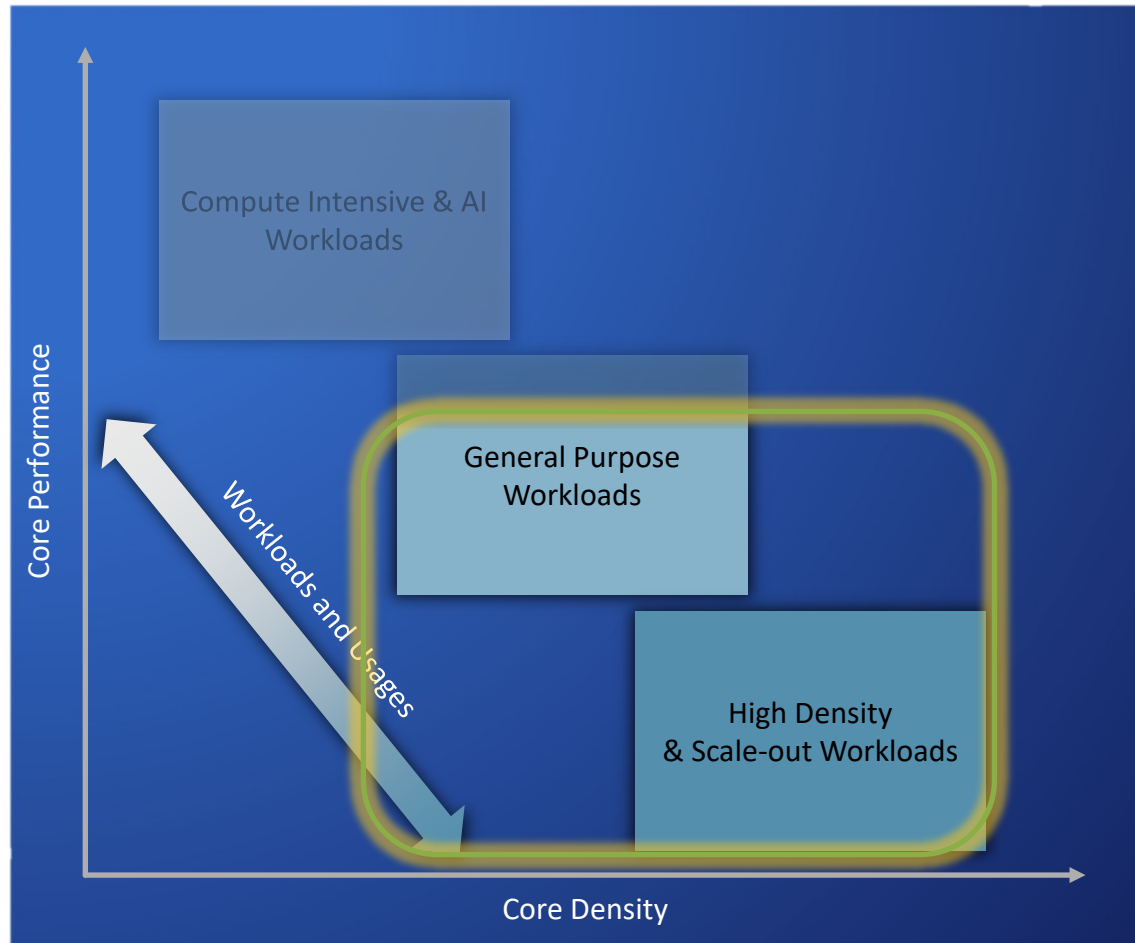
Platform level (IO, memory and coherent) performance scales to match expanding workload requirements

MCR DIMM and CXL attached memory enable efficient and flexible platform performance options



Architected for scalability and flexibility for expanding optimization points

Harness the Efficiency of Xeon with E-cores



Improve OpEx and CapEx
Increased performance per Watt at increased core density

Augments Xeon breadth of coverage
Same hardware, software and firmware

Vector and AI Instruction Support
Including support for FP16, BF16 and INT8 data types

Focus on throughput, density and efficiency
Optimize across full range of utilization

intel.
XEON

Architected For Throughput Efficiency Needs of the Next Decade of Compute

intel®

References

- **TDX**

<https://www.intel.com/content/www/us/en/developer/articles/technical/intel-trust-domain-extensions.html>

- **VT-rp**

<https://www.intel.com/content/dam/www/central-libraries/us/en/documents/intel-virtualization-technologies-white-paper.pdf>

- **CET**

<https://www.intel.com/content/www/us/en/developer/articles/technical/technical-look-control-flow-enforcement-technology.html>

- **VNNI**

<https://community.intel.com/t5/Blogs/Tech-Innovation/Artificial-Intelligence-AI/Deep-Learning-Performance-Boost-by-Intel-VNNI/post/1335670>

- **CXL**

<https://www.computeexpresslink.org/>

- **New instruction set extensions**

<https://www.intel.com/content/www/us/en/developer/articles/technical/intel-sdm.html#inpage-nav-4>