# Meeting the compute needs of the future



SOLVE THE WORLD MOST CHALLENGING PROBLEMS BY PROVIDING ZETTA-SCALE COMPUTING CAPABILITIES BEFORE THE END OF THE DECADE
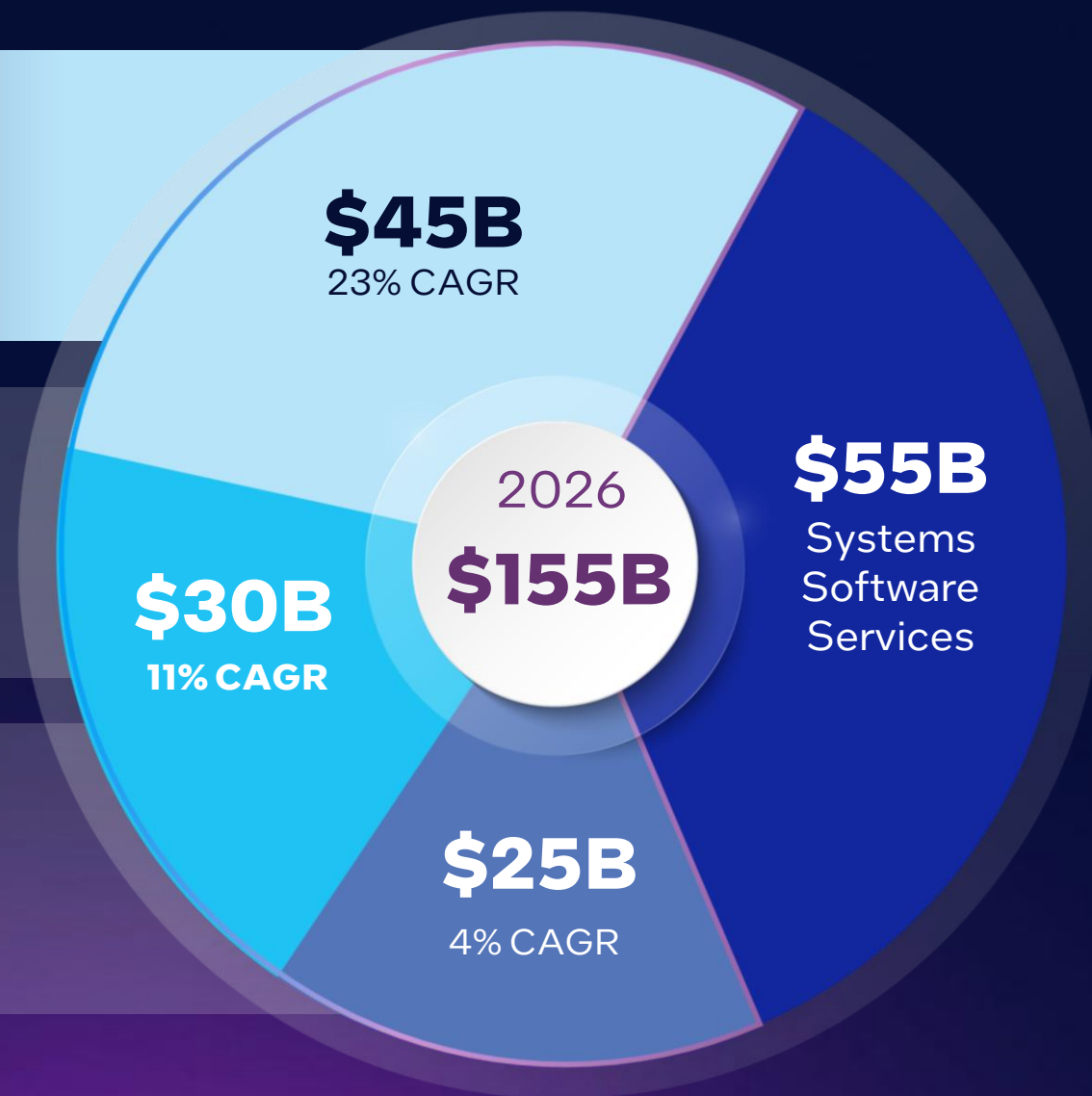
# HPC & Data Center GPU Segments

**Super Compute**
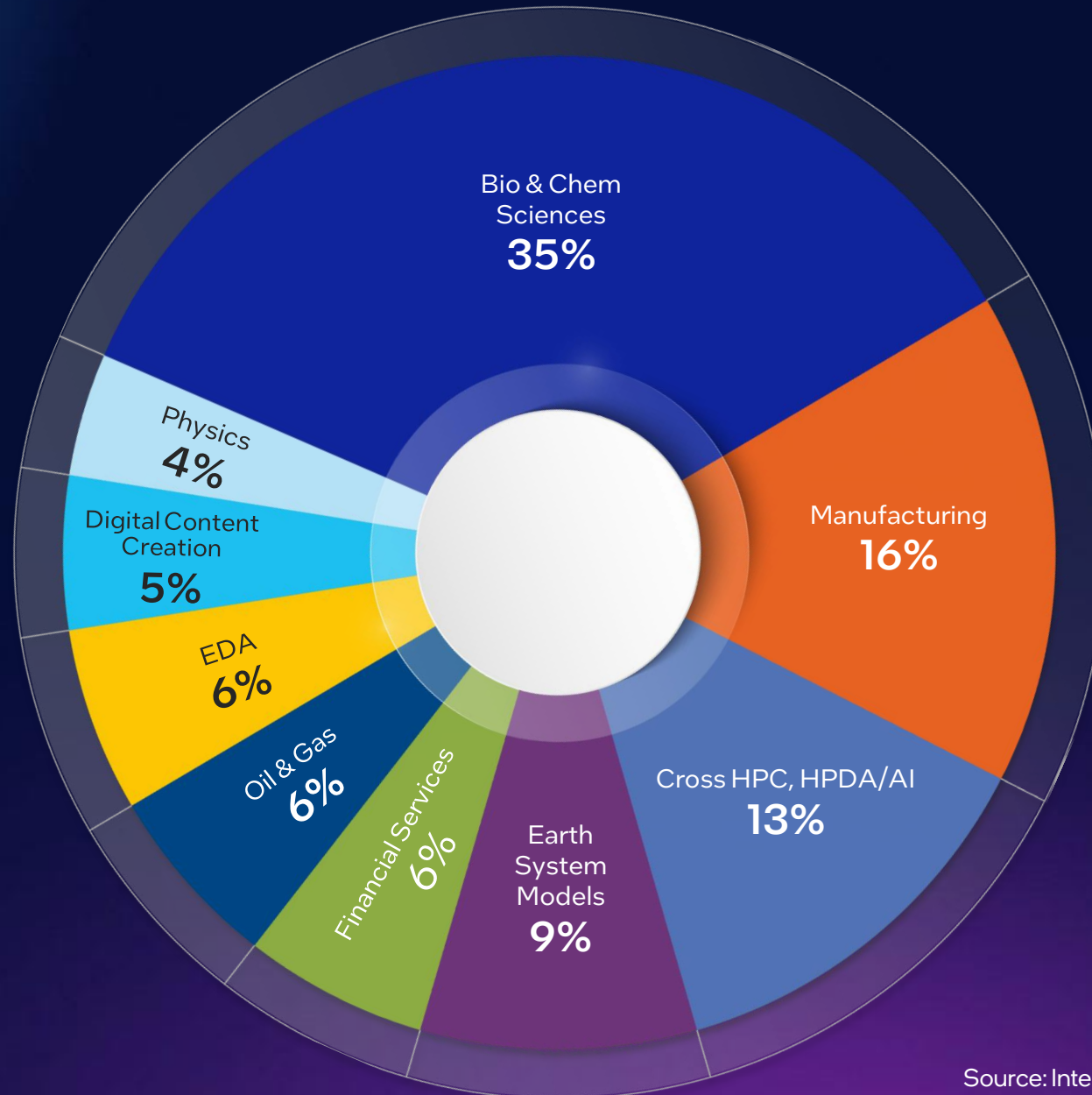
HPC - AI
Media & Visual Cloud

**Custom Compute**

Blockchain
Supercomputing at the
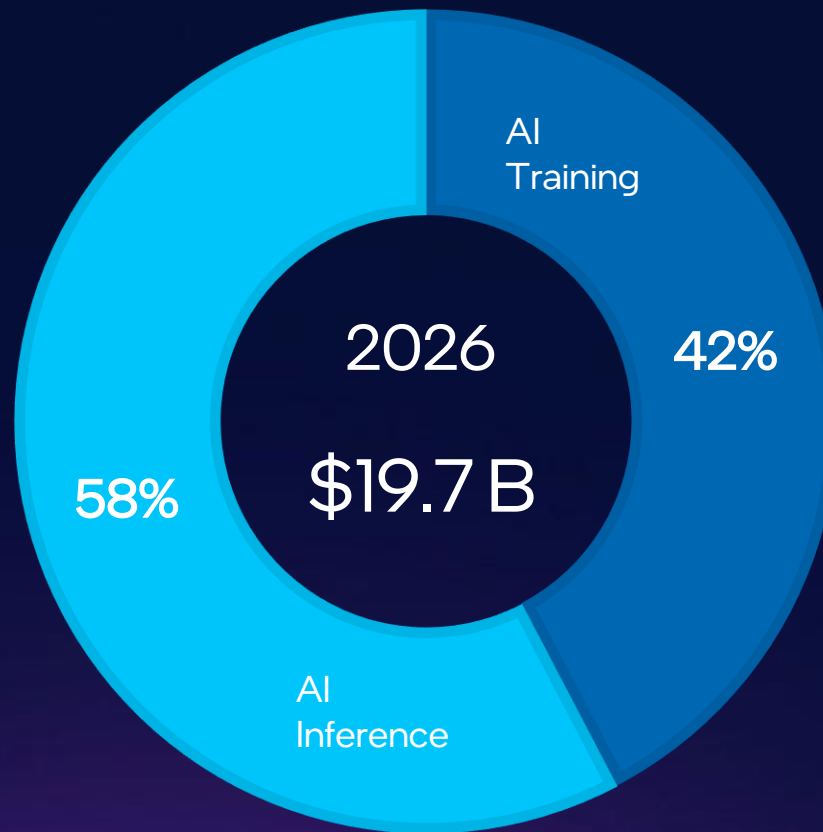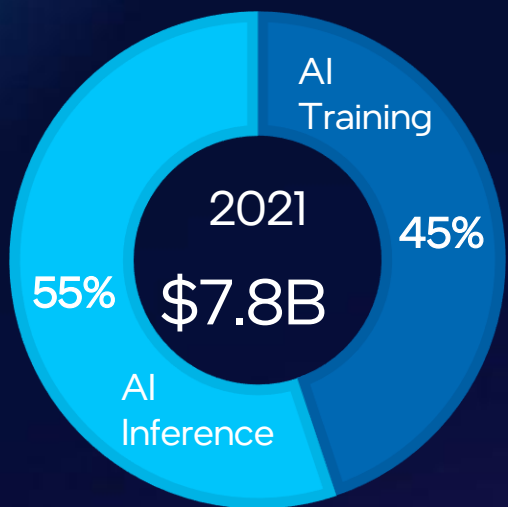Edge

**Visual Compute**

Gaming
Content Creation
Metaverse

**$45B**
23% CAGR

**$55B**
Systems
Software
Services

2026
**$155B**

**$30B**
11% CAGR

**$25B**
4% CAGR

# Top HPC Applications by Vertical



Bio & Chem Sciences 35%
Manufacturing 16%
Cross HPC, HPDA/AI 13%
Earth System Models 9%
Financial Services 6%
Oil & Gas 6%
EDA 6%
Digital Content Creation 5%
Physics 4%

intel VISION

# AI Training and Inference Opportunity

CPU and GPU



2021
$7.8B

AI Training 45%
AI Inference 55%

2026
$19.7 B

AI Training 42%
AI Inference 58%

intel
VISION

# HPC - AI Super Compute Strategy

Install Base

**4ᵗʰ Gen Intel® Xeon® Scalable Processors**
codenamed Sapphire Rapids

**Intel® Xeon® Processors**
codenamed Sapphire Rapids **HBM**

**Intel® oneAPI**

Open Ecosystem

**Ponte Vecchio**

HPC Workloads

AI ML/DL Training

Rendering Visualization

Scientific Research

Compute Workloads

Performance

intel VISION

# Annual Refresh Cycle

**oneAPI**

**oneAPI**

**oneAPI**

## Falcon Shores

**XPU**

**>5x**
Memory Capacity & B/W
Compute density in x86 socket
Performance/Watt

New Tile Based Flexible & Scalable Architecture

Scalable Architecture for all Super Computing Workloads

x86 + Xe

Ponte Vecchio

4th Gen Xeon HBM

Arctic Sound-M

**4th Gen Intel Xeon® processors**

Ponte Vecchio Next

Xeon HBM

Next

Arctic Sound Next

**Xeon** Emerald Rapids

2022

2023

2024

**intel VISION**

| | | | |
|---|---|---|---|
| **Compute** | Up to **128** Ray Tracing Units | **Highest Compute Density** socket & node | Up to **128 $X^e$ Cores** |
| **Memory** | Up to **64MB** L1 cache | Up to **408MB** L2 Cache | Up to **128GB HBM2e** |
| **I/O** | Up to **8 Fully Connected GPUs** | PCIe **Gen 5** | **$X^e$ Link** High-Speed Coherent Unified GPU Fabric |
| **Technology** | **EMIB** | **Foveros** | Intel 7 TSMC N5 TSMC N7 |

# Ponte Vecchio
## $X^e$ HPC based GPU

## Up to 2.6x[1] Perf
over best in market today

On Track
for Aurora 2 Exaflop Supercomputer[2]

[1]Based on pre-production measurements vs A100.
Learn more at www.intel.com/PerformanceIndex. Results may vary.
[2]>2 exaflop peak performance

intel VISION

# Ponte Vecchio GPU Boards & Systems

| GPU BOARDS Intel Branded | OAM SUBSYSTEMS | FULL SYSTEMS |
|---|---|---|
| PCIe Add In Cards | x4 GPU Subsystem | 1U 4 GPU Server |
| OAM Modules | x8 GPU Subsystem | 4U-7U GPU Server |

# Compute Accelerator Market Segmentation by GPU

| | Segment | Number of GPUs |
|---|---|---|
|  | Exascale | > 10K |
|  | Exascale Follow On | Up to 10K |
|  | Hyperscalers CSP | Up to 10K |
|  | Large Enterprise and Next Wave CSPs | Up to 1K |
| SME  | Enterprise AI and HPC | Up to 100 |

intel
VISION

# Xe Link for Scalability

Enabling a high number of coherent and unified accelerators



Flexibility for Scale Up and Scale Out across GPUs and Nodes

Note:
Xe Links connect gluelessly between PVCs
The Xe Link Logo above is not representing an additional System Device

intel
VISION

Monte Carlo

Black Scholes

Binomial

# Leadership Performance on $X^e$-HPC

intel.