# Building Blocks for largest system

Tech Foundation
Today

Architecture

Power & Thermals

Systems

Process & Packaging

Silicon

Memory & I/O

Software

Silicon   Systems   Software

Silicon    Systems    Software

Hybrid
Compute
Cluster
in a Package

# UCIe
## Universal Chiplet
## Interconnect Express

Reduced Time-to-Market

Reduced IP Porting Costs

Smaller, Higher Yielding Components

Optimal, per-Component Tech

Leaders in semiconductors, packaging, IP, cloud service providers joining forces

Alibaba Group

AMD

arm

ASE GROUP

Google Cloud

intel

Meta

Microsoft

nVIDIA

Qualcomm

SAMSUNG

tsmc

High-Speed Standardized Chip-to-Chip Interface

Tuned to Workloads ... Enabling Differentiation

Advanced Packaging

# Integrated Optical I/O

Key innovation for growing data rates, energy efficiency and channel loss minimization needs

XPU

EMIB

Optical I/O Tile

2023 target

| Bandwidth | Reach | Shoreline Density | Energy Efficiency | Latency |
|---|---|---|---|---|
| ~1 Tbps per fiber | >100m | >4x PCIe6 | Trending 3pJ/b | <10ns + TOF |

Silicon   **Systems**   Software

# The Open Standards—CXL™ is The Future

## Solve these challenges with CXL

- Scaling challenges: latency, bandwidth, capacity
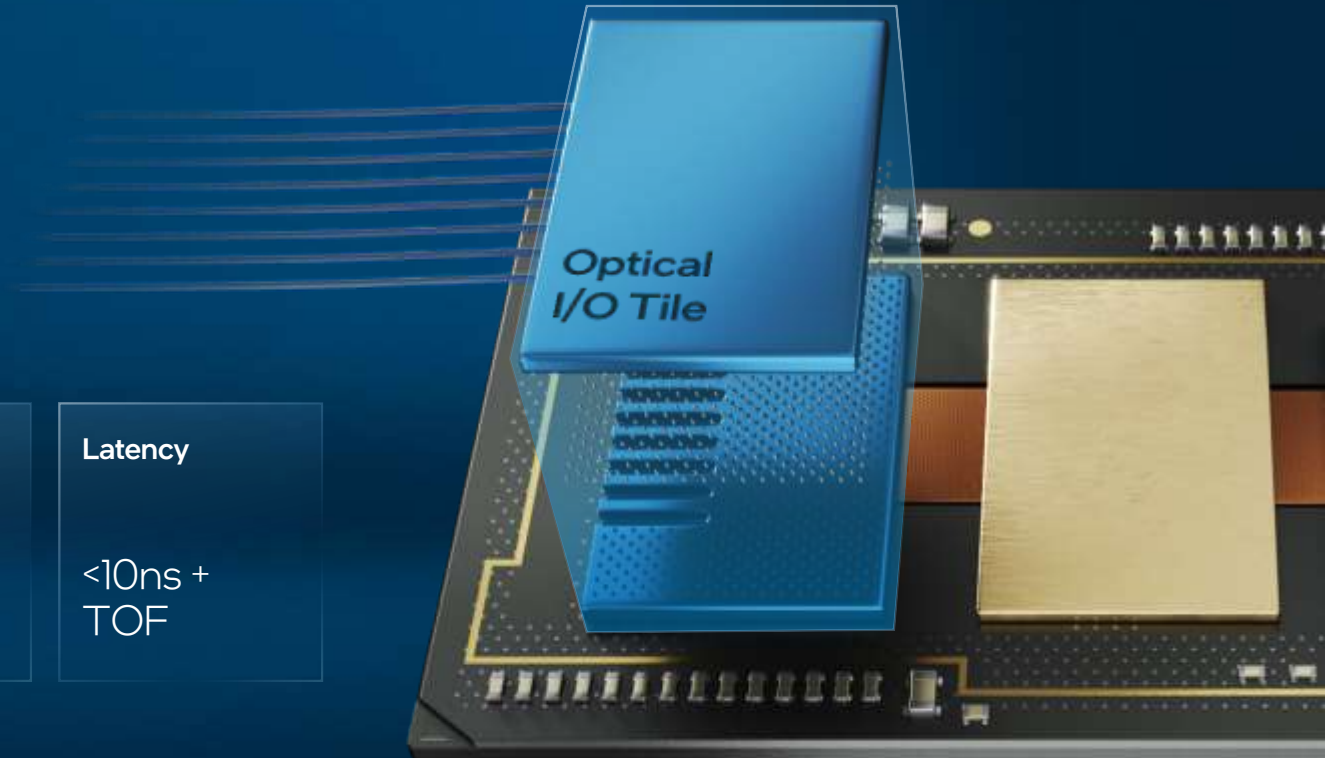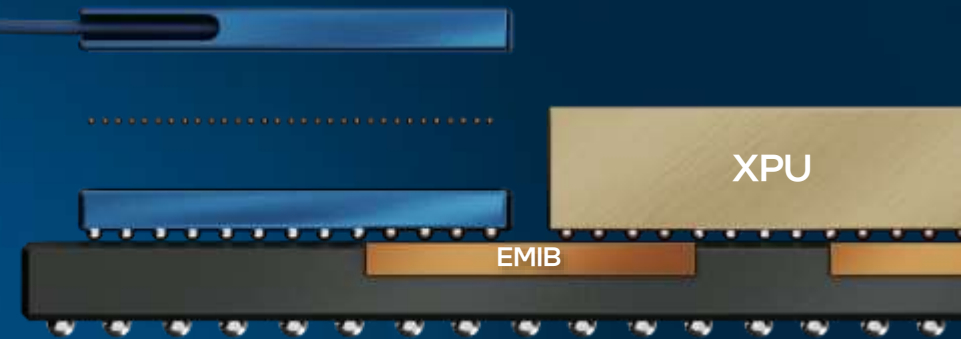- Increase of heterogenous computing + massive datasets + demanding workloads = need for more accessible data, faster
- These computing demands require a more efficient interconnect between CPU and traditional I/O interface



## Go Faster
- Improves data handling and reduces I/O bottlenecks

## Do More
- Memory bandwidth and capacity expansion, efficient access across shared memory
- Enables CPU and accelerators to share memory resources for higher performance

## Save More
- Improves TCO

# CXL™ Consortium – Scope and Feature
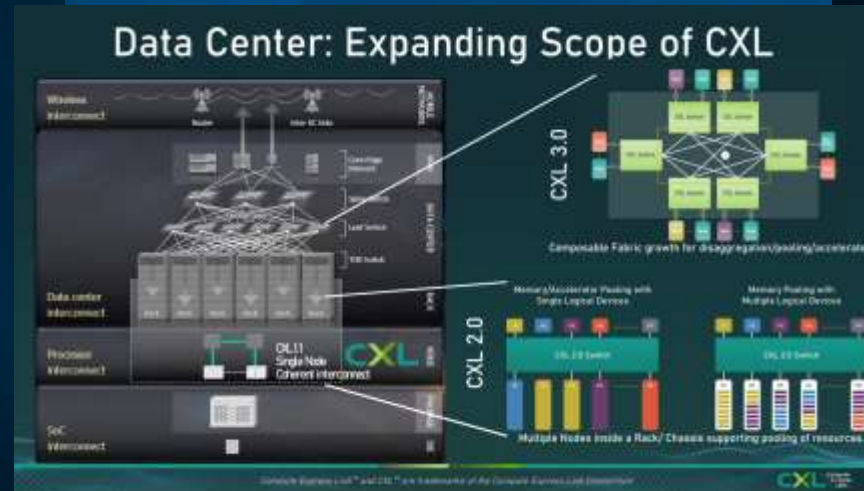
## Coherent Interface
Leverage PCIe with three multiplexed protocols
Built on top of PCIe® infrastructure

## Low Latency
CXL.Cache/CXL.Memory targets near CPU cache coherent latency (<200ns load to use)

## Asymmetric
Complexity
Eases burden of cache coherence interface designs for devices



Data Center: Expanding Scope of CXL

## Heterogeneity
Enable 3 type of devices based on 3 protocols:
memory, cache, IO

## Modularity & disaggregation
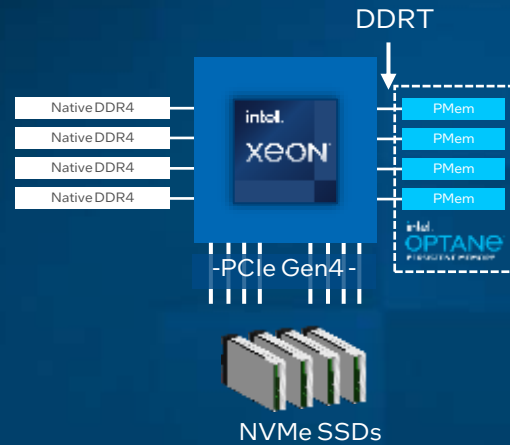Allocate and deallocate resources on demand,
Enhance memory and cache coherency
Peer to peer memory access

## Scalable
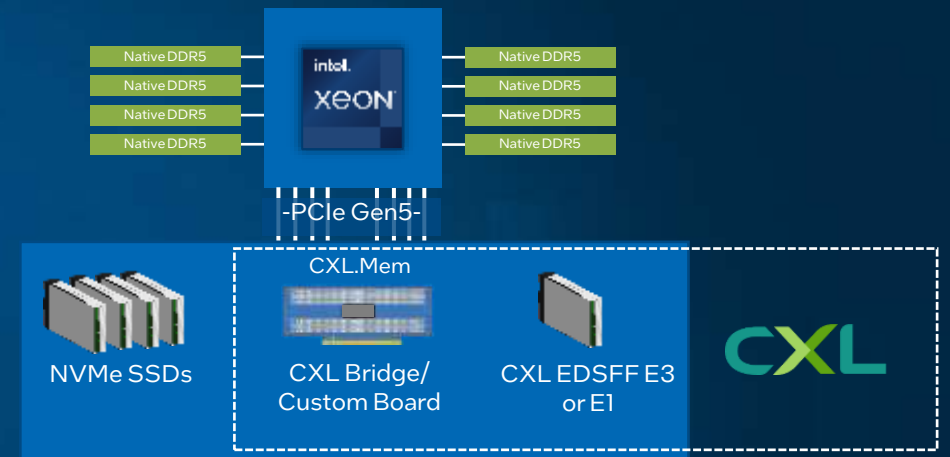High speed low latency fabric
Composable fabric

enable efficient resources sharing and modularity

# Possible Usage Model Transition Examples



DDRT

Native DDR4
Native DDR4
Native DDR4
Native DDR4

intel. XEON

PMem
PMem
PMem
PMem

intel. OPTANE PERSISTENT MEMORY

-PCIe Gen4 -

NVMe SSDs

Native DDR5 | Native DDR5
Native DDR5 | Native DDR5
Native DDR5 | Native DDR5
Native DDR5 | Native DDR5

intel. XEON

-PCIe Gen5-

CXL.Mem

NVMe SSDs | CXL Bridge/ Custom Board | CXL EDSFF E3 or E1

CXL

**Today: Intel® Optane™ Technology**
Intel® Xeon® CPU + DDR4 +
Intel® Optane™ Technology on DDRT

**Future: CXL 2.0**
Intel® Xeon® CPU + DDR5 +
Third-Party CXL Memory Products

| | | |
|---|---|---|
| **Memory Mode** | Cost-effective memory to meet workload needs and expanded memory capacity for workload scale-up. | Adds industry standard protocol for software to cache or tier memory. |
| **Memory Expansion, Augmentation** | Path for higher capacity than DRAM for workloads to scale-up. | Adds additional memory capacity and memory bandwidth to existing DDR attached DRAM. |
| **Persistent Memory** | Fast storage for meta data storage, fast write logs, and caching/tiering acceleration of other storage. | Adds persistence support (non-volatile memory over CXL). |
| **High-Endurance, Performant SSDs** | Very fast storage "replacement" for meta data storage, fast write logs, caching/tiering acceleration. | Adds persistence support (memory-semantic storage over CXL). |

# Coherent Discrete Accelerator for 4th Gen. Intel® Xeon® Scalable Processors

## Compute Express Link (CXL) Interface

intel. **XEON**®

Cache & Memory Coherent

- Low Latency
- Increased Memory Space
- Custom Acceleration

intel. **AGILEX**™ 7

Processor Memory

## Accelerates Diverse Workloads
(Including data analytics, database acceleration and function-as-a-service)

FPGA Memory

*Available on selected Intel Agilex 7 I-series and M-series FPGAs which contain at least one R-Tile.

https://www.computeexpresslink.org/

intel.

Forum TERATEC 23

Silicon    Systems    **Software**

# Unifying Software Stacks Critical...

Applications
Middleware, Frameworks, Runtimes
Low Level Libraries
OSes & Virtualization
HAL

CPU
GPU
AI

intel Forum TERATEC 23

# Maximizing Impact,
# Minimizing "energy to solution" through...

**Open**
Ecosystem

**Choice**
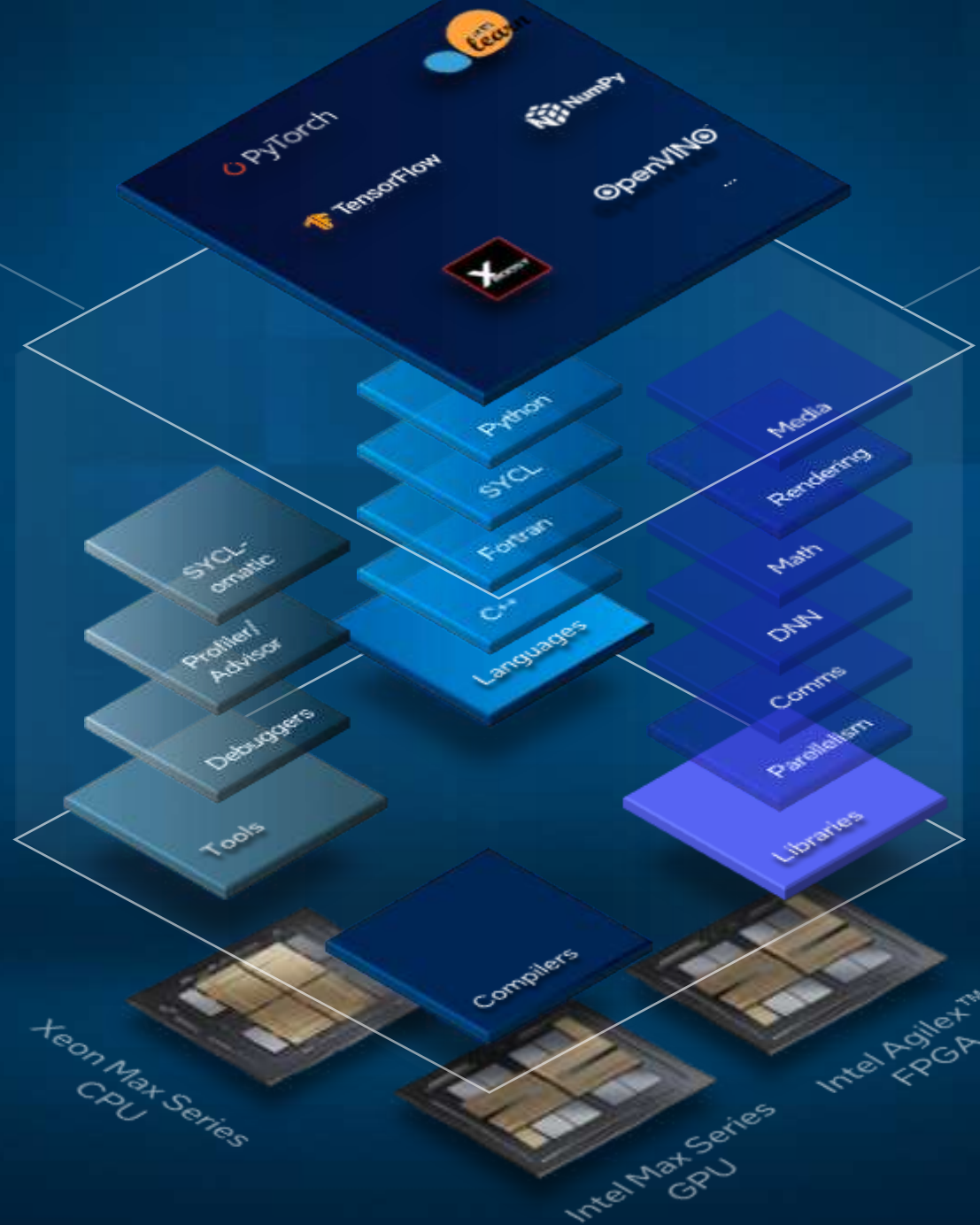Compatibility

**Trust**
Workloads

**Scale**
Delivery & Deployment

intel

Forum
TERATEC 23