

# Fujitsu AI-HPC Platform - Powering the Next Era of Simulation



## Trends in HPC

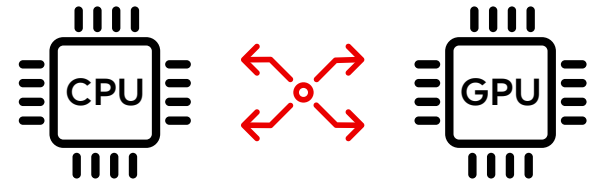
Ever-Growing data volumes and needs for computational power

Needs for higher computational density and power efficiency

Rapid AI evolution and expansion of applicable areas

## AI-HPC platform powering the Next Era of Simulation

- Platform for fusion of AI and simulation
- CPU + GPU tightly coupled architecture
- Leverages both CPU and GPU strength
- Targeted around 2030



# Revolutionizing Simulation with AI-HPC Platform

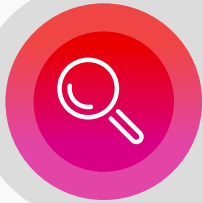


## AI-enhanced Simulation

Significantly improve accuracy and speed of simulation using AI

Faster Computation

Higher Prediction Accuracy



## Simulated Data-Driven Science

Discover unknown behaviors hidden within the complexity of real-world physics

Massive Data Exploration

Hidden Pattern Detection



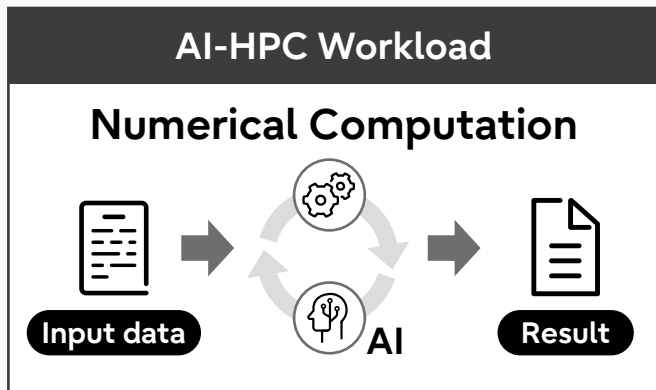
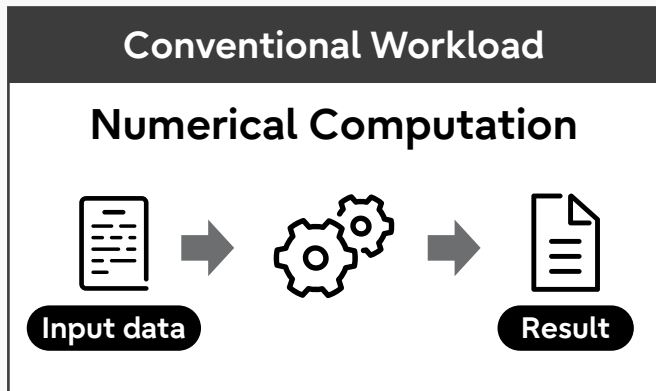
## Agentic Simulation

AI agents autonomously plan, execute, and refine multi-step simulation loops

Autonomous Trial & Error

Beyond Human Intuition

# AI-HPC Workload and CPU Requirements



## Demands on CPU



### Fine-grained Simulation ↔ AI Interaction

Simulation and AI exchange data at every step, creating a tightly coupled feedback loop



### Agent Reasoning & Planning

Agent planning is a sequential decision chain — single-thread performance directly determines reasoning speed



### Low-latency Real-time Response

Real-time inference in the control loop demands both low latency and high throughput

## CPU Requirements



### On-CPU AI Processing

SVE2 + SME2 and optional NPU enable AI inference directly on CPU — eliminating data transfer overhead

**SME2**



### Tight GPU Coupling

High-bandwidth, low-latency CPU-GPU interconnect for seamless data flow in the simulation-AI loop

**NVLink Fusion**



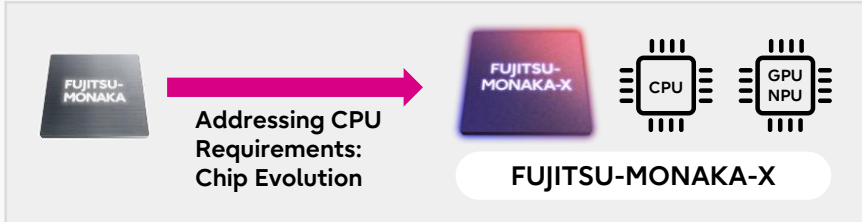
### Many High-Capability Cores

Massive core count with strong single-thread performance for orchestration, agent reasoning, and broad-spectrum simulation

**>200 Cores**

# FUJITSU-MONAKA-X Processor

## Arm-based Processor Enhanced for AI



### Scale and future Enhanced

- ArmV9-A Architecture
- >200 cores / socket
- Ultra low voltage
- 3D chiplet
  - Core die 1.4nm
  - SRAM /IO die
- RAS/ Enhanced security CCA
- Low latency cluster cache

### Use Case Enhanced

- SVE2
  - 512bit SME2
- NPU (AI)
- NVLink Fusion
- Borderless Coherent UArch

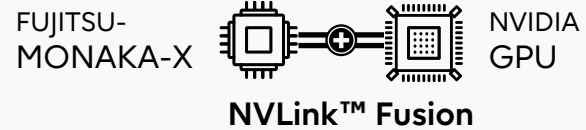
**Upscale Application thrupt / performance**

**Upscale Application coverage**

- Legacy HPC support
- Enhanced GPU collaboration for AI for science further

**CPU-GPU Tight Integration technology**

## KEY Tech: CPU-GPU Tight Integration



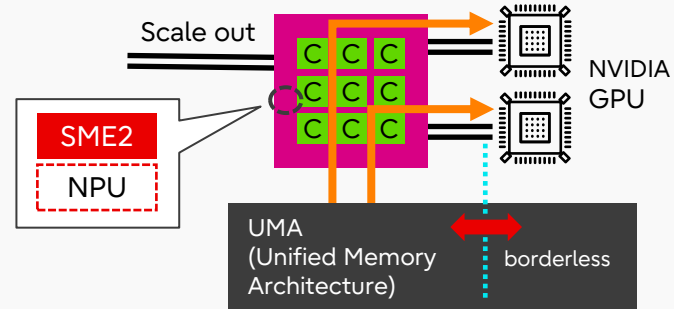
### CPU Strengths -Low Latency SME2

- Complex control flow, latency-sensitive tasks, and irregular memory access
  - Simulation, Realtime AI, Multimodal processing

### GPU Strengths -High thrupt NVLink

- Massively parallel data-parallel execution with regular memory access patterns
  - Large scale DL/ML training, Image/signal recognition

### Bringing together the best of CPU-GPU

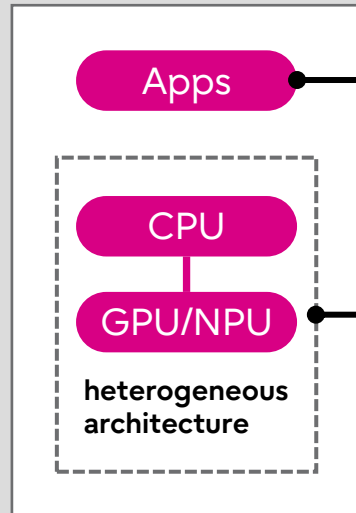


# Advancing AI/HPC Through Software Evolution

- Future AI/HPC platforms require both efficient hardware utilization and software-driven innovation

## Co-Design-Driven Platform Optimization

Next-generation AI/HPC platforms are optimized through hardware–software co-design across system architectures and node configurations for diverse workloads and efficient execution



### Accelerating Simulations with AI

Leveraging technologies such as hybrid surrogate modeling & simulation with accuracy-preserving quantization  
Continuously advancing Scientific AI, Autonomous Scientific Discovery, and Autonomous Research Agents

### Enhancing AI/HPC Software on FUJITSU-MONAKA-X

Continued OSS community collaboration for efficient utilization of FUJITSU-MONAKA-X CPU, GPU, and future heterogeneous accelerator architectures

1. SVE2/SME2 and microarchitecture-optimized HPC compilers and libraries
2. High-performance optimization for major AI libraries, including low-precision computing
3. AI/HPC acceleration across heterogeneous compute — CPUs, GPUs, and NPUs

# FugakuNEXT – Japan’s Flagship Platform

## FugakuNEXT development Goals

### FugakuNEXT, AI-HPC platform Ecosystem

#### Made with Japan

Co-developed by RIKEN, NVIDIA, Fujitsu



#### Technological Breakthroughs

- Tight CPU–GPU integration
- Modernization of applications

#### Sustainability and Continuity

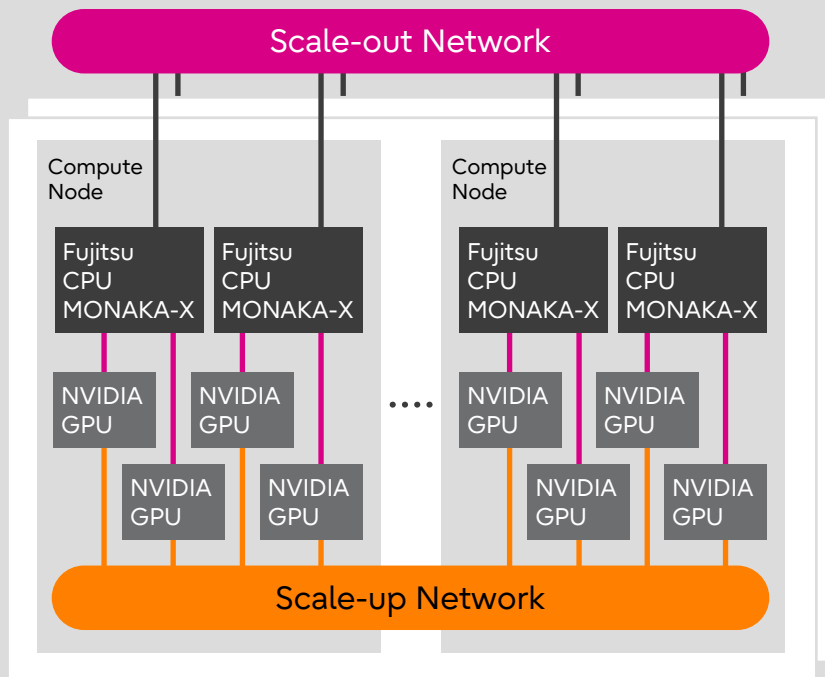
- Sustainable software ecosystem
- Energy-efficient system

## Schedule



# The world's top-level AI-HPC platform

- FugakuNEXT is not merely a larger supercomputer; it is a sovereign AI-HPC infrastructure designed to support both large-scale AI training and traditional scientific computing.



Specification	FugakuNEXT	
	CPU	GPU
Number of Nodes	$\geq 3,400$	
FP64 Vector perf.	$\geq 48$ PFLOPS	$\geq 2.6$ EFLOPS
FP16/BF16 Matrix perf.	$\geq 1.5$ EFLOPS	$\geq 150$ EFLOPS
FP8 Matrix perf.	$\geq 3.0$ EFLOPS	$\geq 300$ EFLOPS
FP8 Matrix sparsity perf.	-	$\geq 600$ EFLOPS
Memory Capacity	$\geq 10$ PiB	$\geq 10$ PiB
Memory BW	$\geq 7$ PB/s	$\geq 800$ PB/s

**Thank you**

