

ISC 2018

DLU™ : Deep Learning Unit

Fujitsu's HPC Development Timeline

FUJITSU

K computer

The K computer is still competitive in various fields; from advanced research to manufacturing.



Gordon Bell
Prize Finalist
(2016)

HPCG
No.1
(2016)

Graph500
No.1
(2016)

Deep Learning Unit (DLU™)

DLU is a processor designed for deep learning that has the ability to handle large-scale neural networks.



Post-K Computer

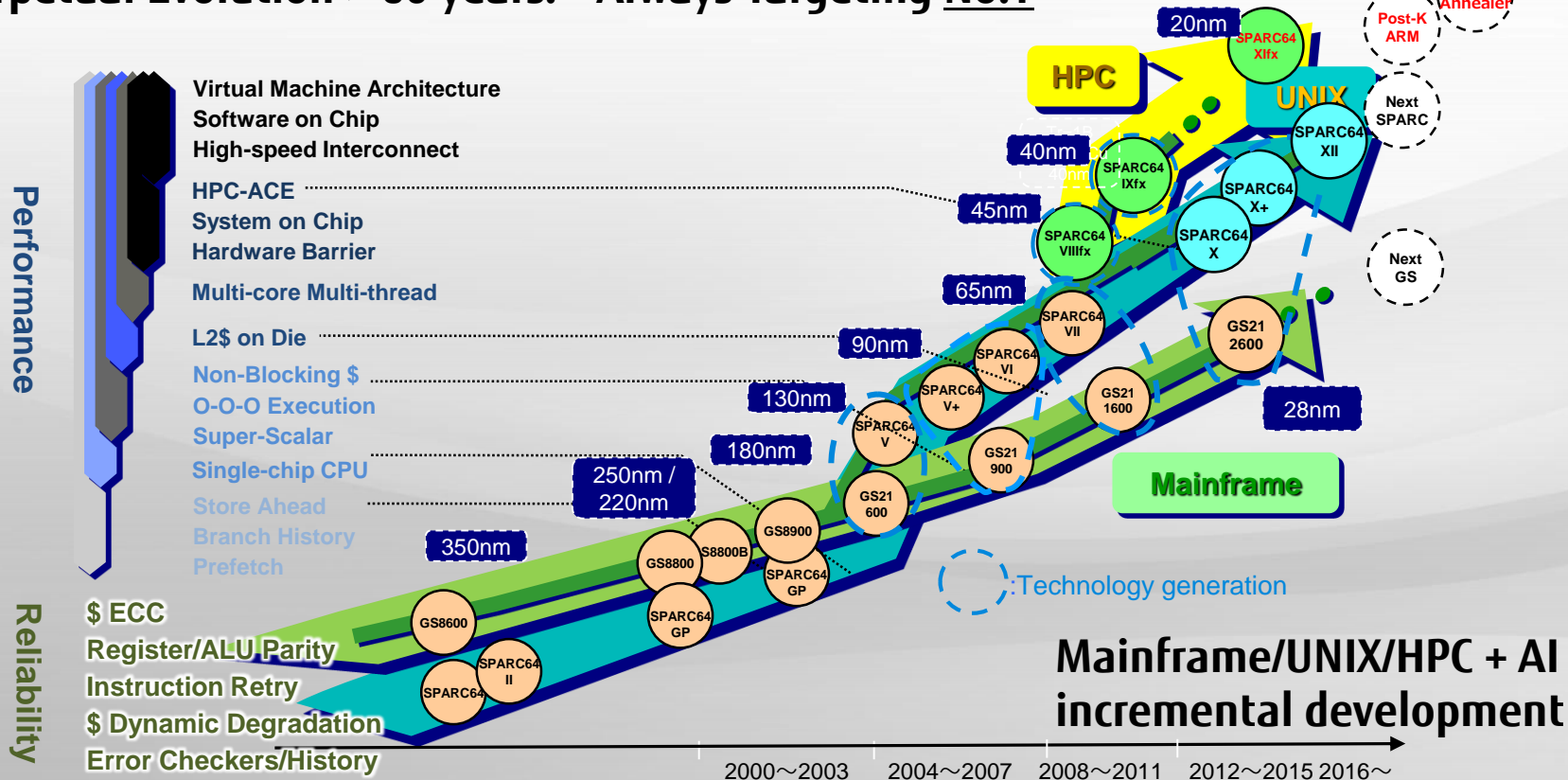
The post-K is under development to achieve superior application performance.



Fujitsu Processor Development

Perpetual Evolution > 60 years: Always Targeting No.1

FUJITSU



DLU: Processor Designed for Deep Learning

FUJITSU

DLU

Deep Learning Unit



Utilizing technologies derived
from the K computer



Features

- Architecture designed for deep learning
- Low-power consumption design
- » **Goal: 10x Performance / Watt compared to competitors**
- Scalable design with Tofu interconnect technology
- » **Ability to handle large-scale neural networks**

What's the New Architecture for the DLU?

Domain specific, Optimal precision, and Massively parallel.

Conventional Architecture

General Use

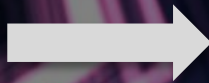
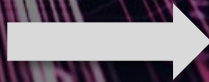
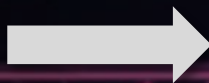
Complicated 0-0-0 cores
w/ cache memory

High Precision

Double/Single precision FP

Sequential + Parallel

Multiple strong cores



The New Architecture

1. Domain Specific

Domain specific cores
w/ large register file

2. Optimal Precision

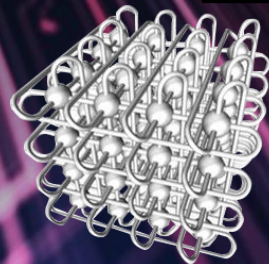
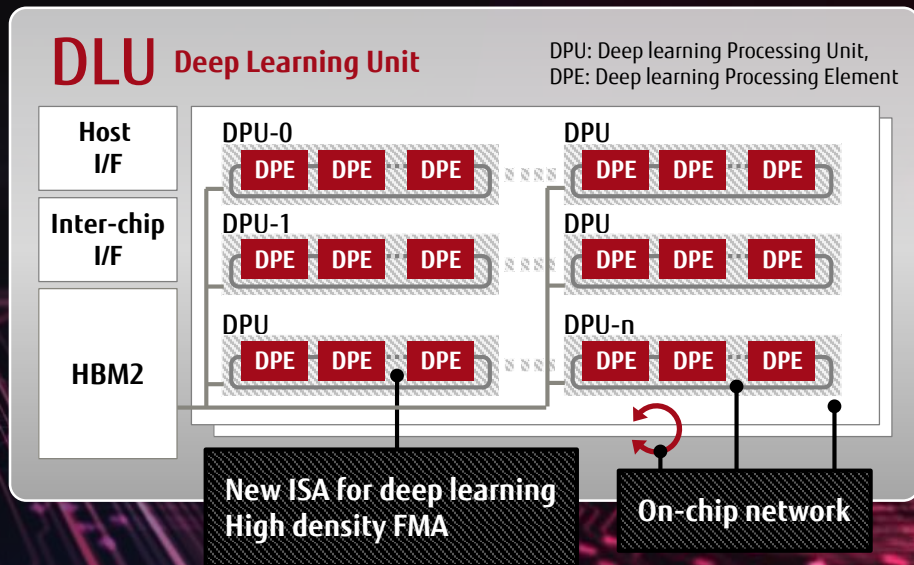
Deep Learning Integer

3. Massively Parallel

Many cores w/ on-chip network

DLU Architecture

- ISA: Newly developed for deep learning
- Micro-Architecture
 - Simple pipeline to remove HW complexity
 - On-chip network to share data between DPUs
- Utilizes Fujitsu's HPC experience, such as high density FMAs and high speed interconnect
- Maximizes performance / watt



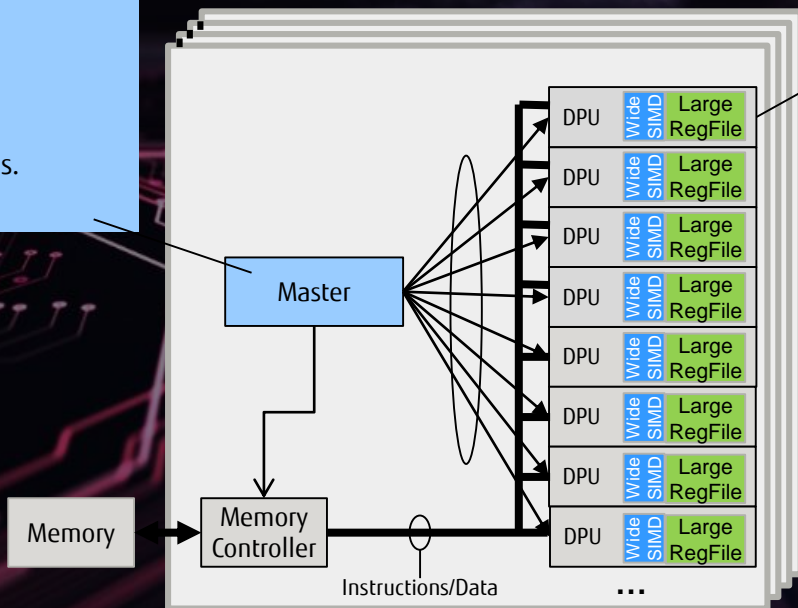
Fujitsu's interconnect technology
Large scale DLU interconnect through off-chip network

Heterogeneous Cores and Large Register File

The combination of few large core (Master) and many small execution cores (DPU) results in more performance with less power consumption, compared to a conventional homogeneous structure

Master Core: Memory Access and DPU control

- Push & Pull instructions and data for DLUs.
- Start/stop execution of DLUs

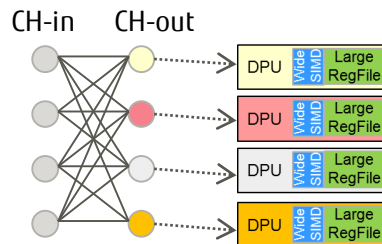


DPU: Execution

- Wide SIMD execution and Large Register File.
- Execute DL operations based on master core's control

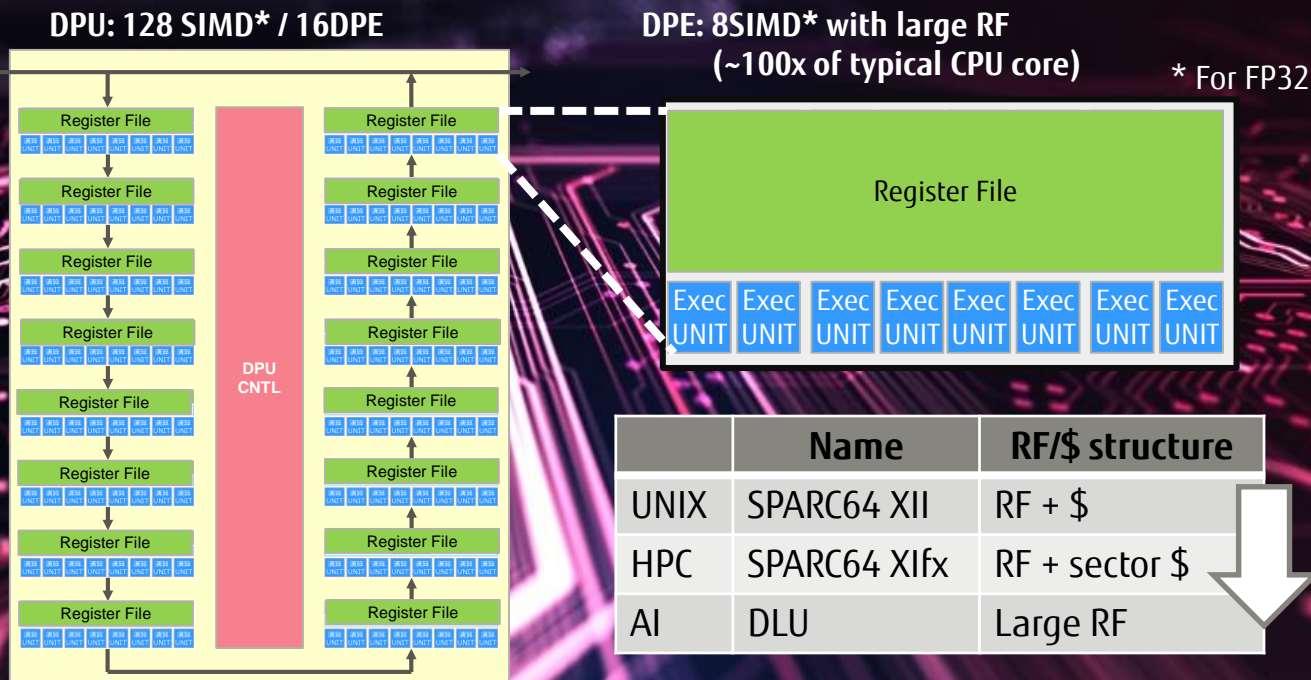
How to utilize many DPUs

- e.g. convolution layer
one CH-out / DPU
multiple batch / DPU



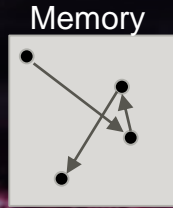
DPE & Large RF (Register File)

- DPU consists of 16 DPEs connected with on-chip network
- DPE includes large RF and wide SIMD execution units to realize an efficient Deep Learning engine.
 - RF is fully SW controllable unlike cache to extract full HW potential



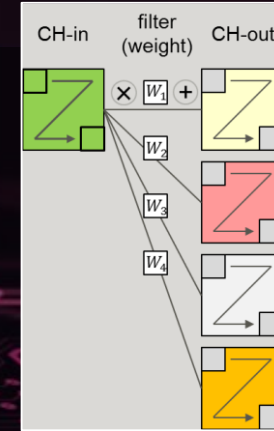
Domain specific architecture - Why cache memory removed -

General Processors



- Complex hardware to achieve high performance for any applications with various data access patterns
- E.g.
 - Large cache memory with cache tags and LRU replacement Unit
 - Hardware Prefetch Engine

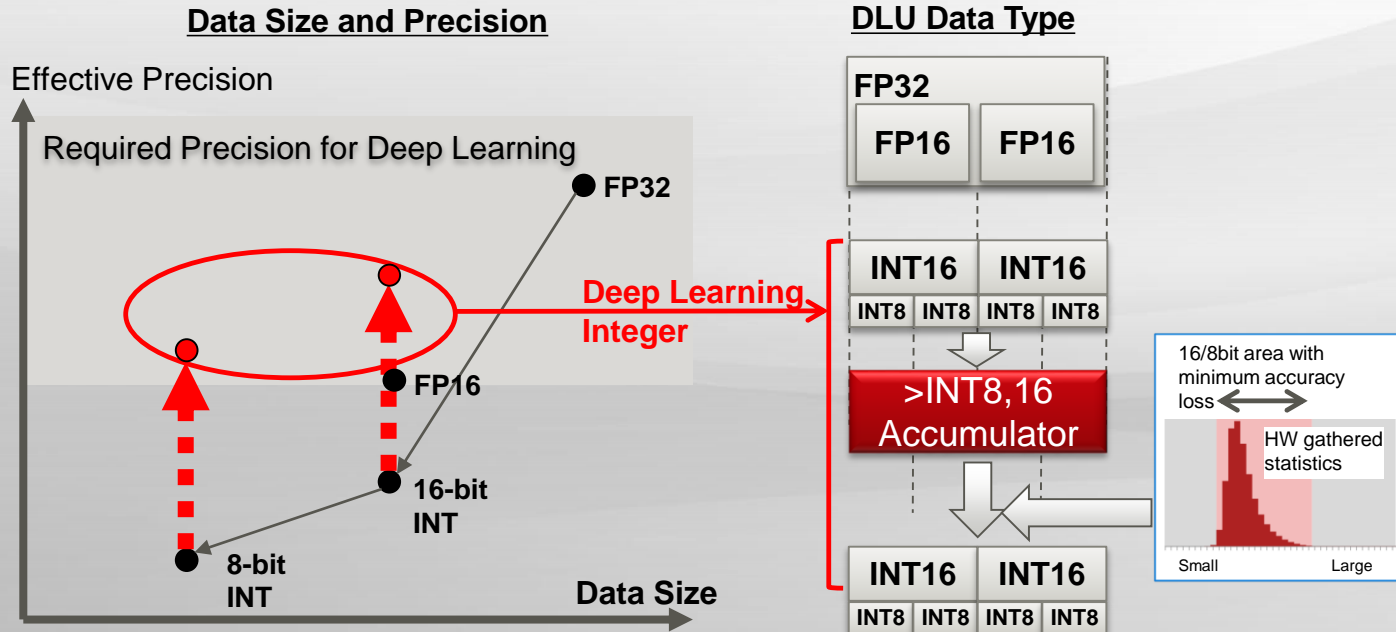
DLU (Domain Specific)



- Simple hardware focusing on simple memory access patterns
- E.g. Convolution Layer
 - CH-in data can be shared among CH-out calculation at all DPUs
 - Memory access patterns are continuously and predictable (software controllable)

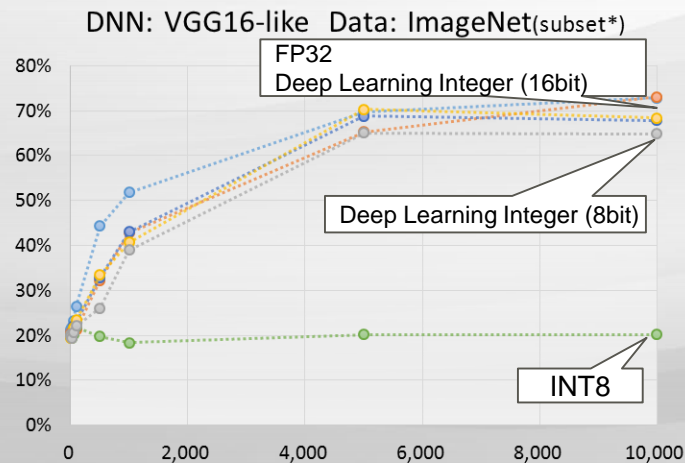
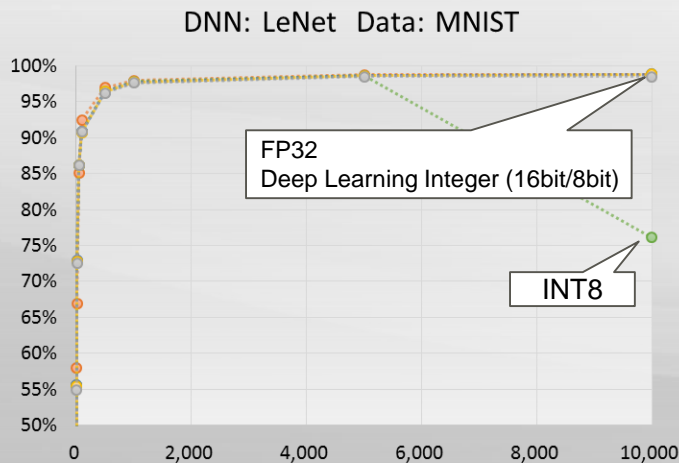
DLINT : Deep Learning Integer

- Fujitsu's "DLINT" realizes necessary accuracy for Deep Learning with only a 16 or 8 bits data size (i.e. less power consumption compared with FP32)
- Training results with DLINT8/16 can be converted to the conventional 8/16-bit INT for inference.



Accuracy of Deep Learning Integer

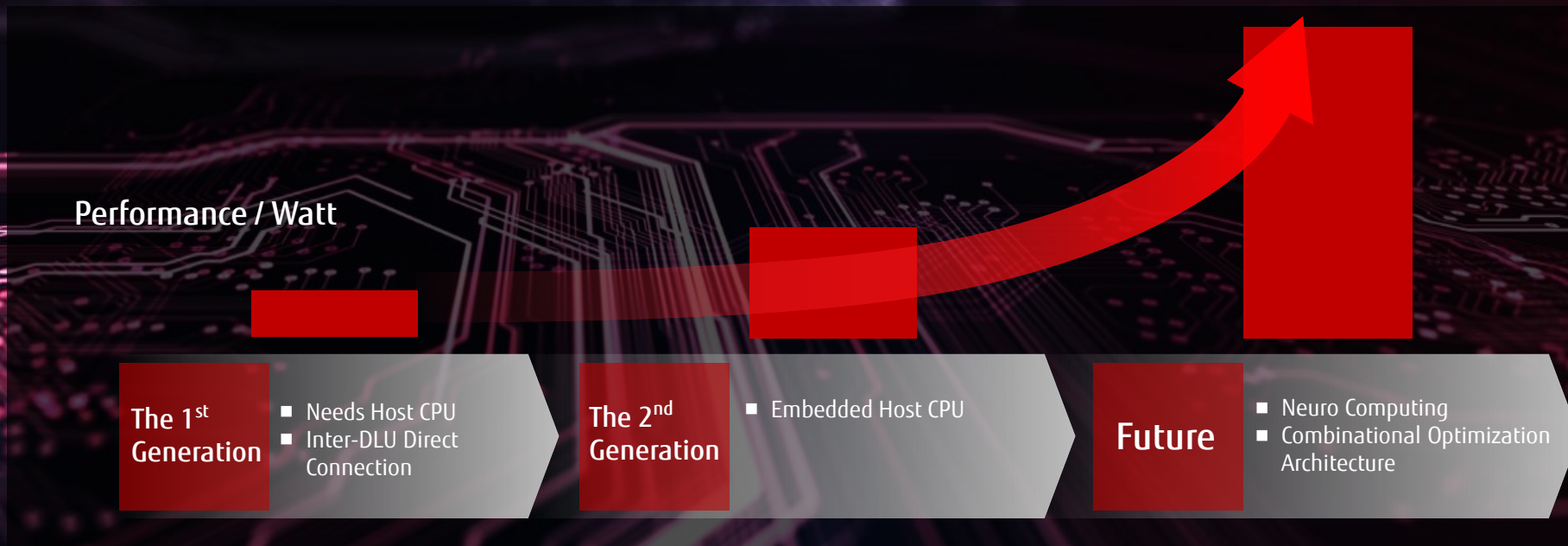
DLINT has shown similar accuracy with FP32 precision



(*) ImageNet(subset): image size=96x96, #categories=25

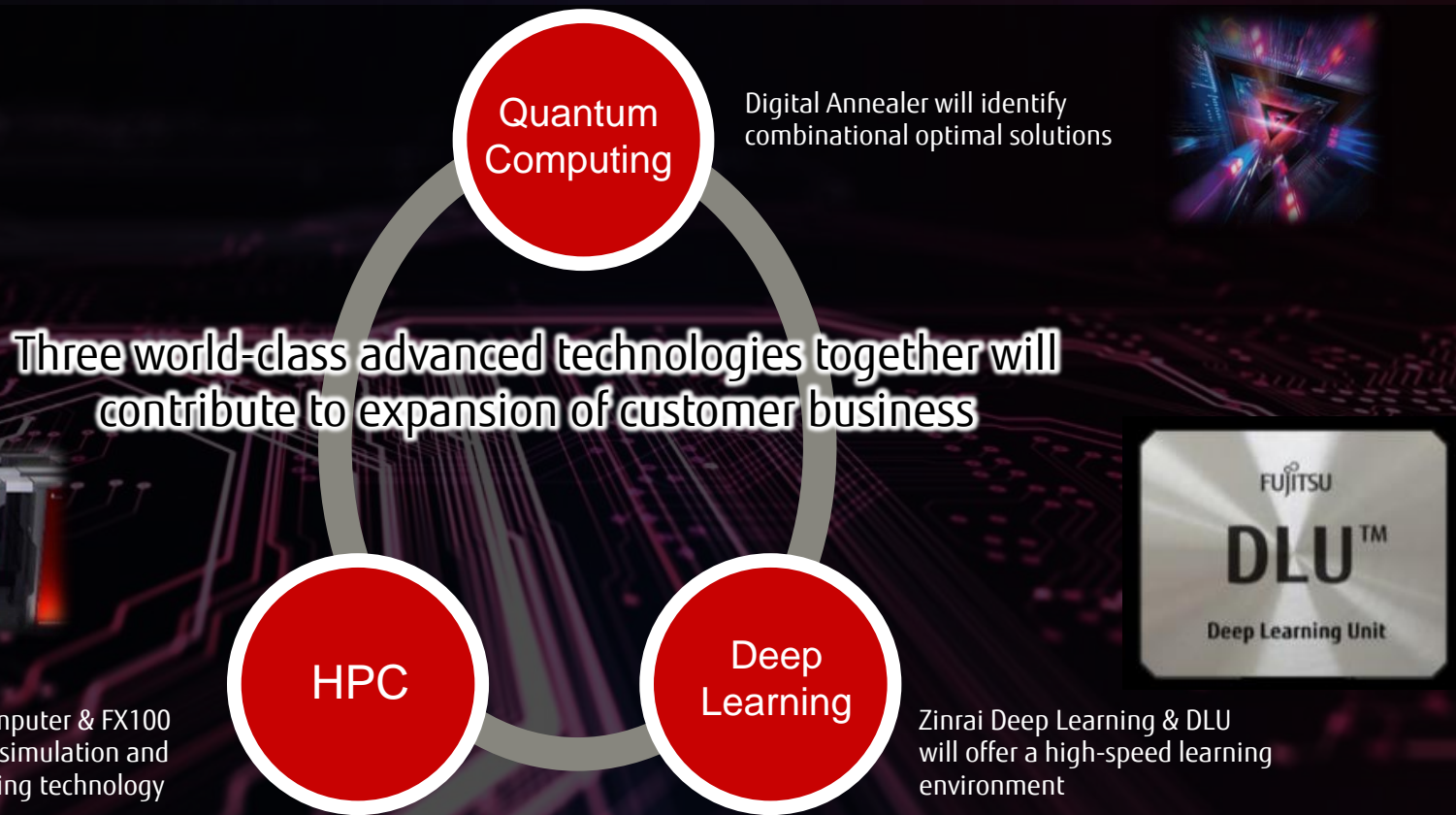
DLU Roadmap

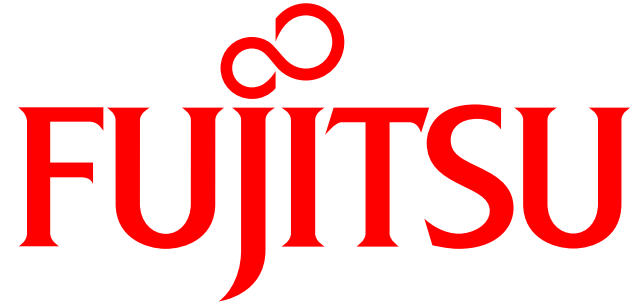
- Multiple generations of DLUs over time, as we currently do for HPC/UNIX/Mainframe processors



* Subject to change without notice

AI will be accelerated by three technologies





shaping tomorrow with you