**NEC**

2nd Aurora Deep Dive Workshop at RWTH Aachen University
November 28, 2019

# SX-Aurora TSUBASA

Passion for Sustained Performance

**Shintaro MOMOSE, Ph.D.**
shintaro.momose@emea.nec.com
NEC Deutschland GmbH

1. Aurora Over View

2. Vector Engine

3. VE Partitioning Mode

4. VE10E, VE20

5. Performance

6. High Density Product, A412-8

7. Aurora3

Orchestrating a brighter world  NEC

# SX-Aurora TSUBASA

**POINT 1** — **Memory Bandwidth**

1.22TB/s / processor, 150GB/s / core
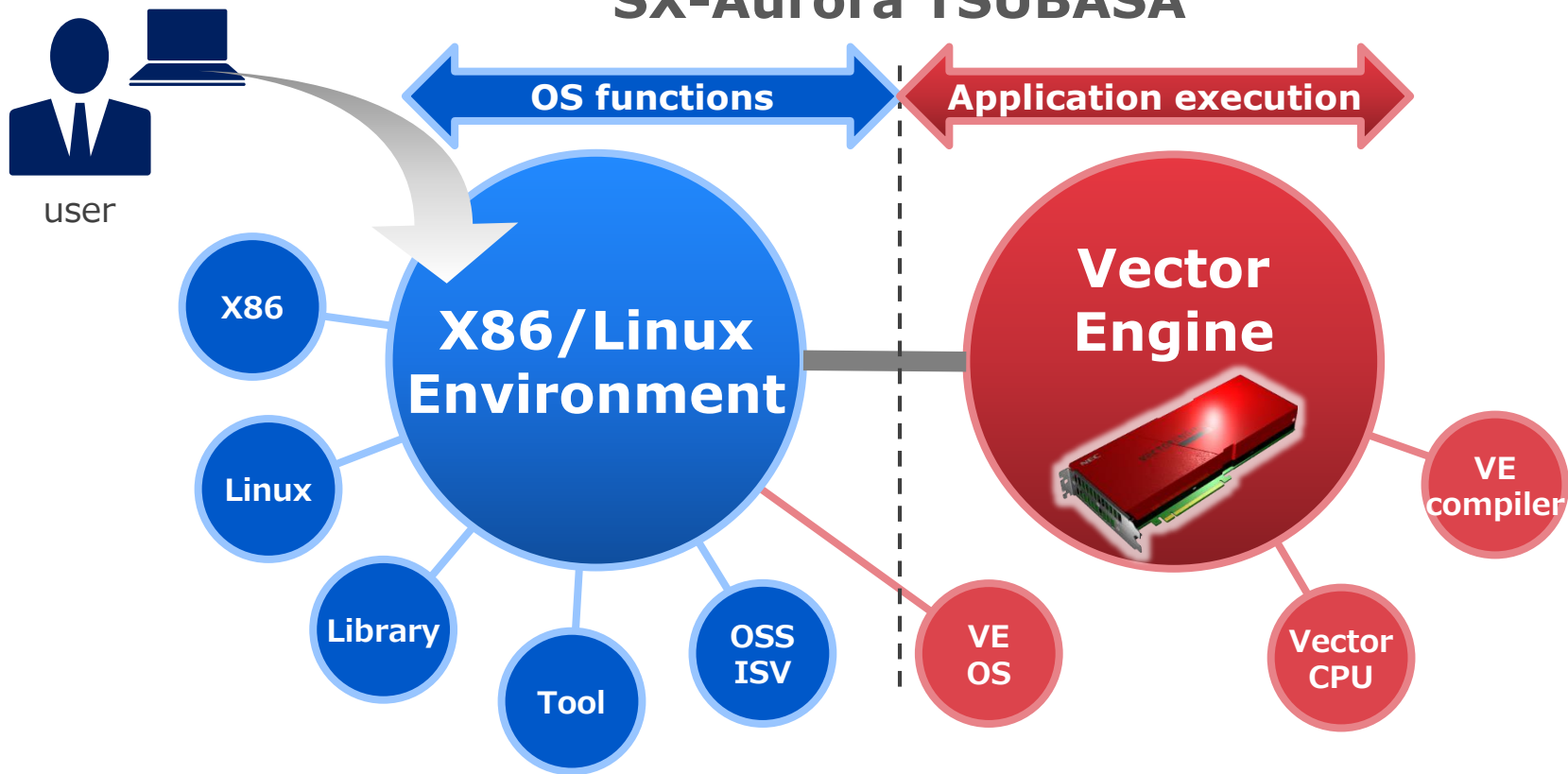
**POINT 2** — **Easy to Use**

Fortran/C/C++ programing, OpenMP
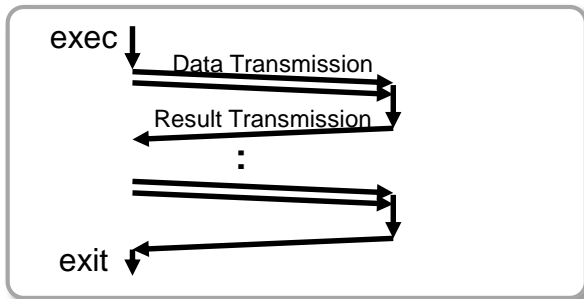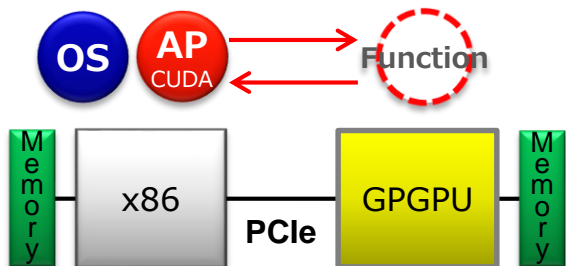Automatic vectorization/parallelization

**POINT 3** — **x86/Linux**

High sustained performance on
x86/Linux environment

Orchestrating a brighter world **NEC**

SX-Aurora TSUBASA

user

OS functions

Application execution

X86

Linux

Library

Tool

X86/Linux Environment

OSS ISV

Vector Engine

VE OS

VE compiler

Vector CPU

Orchestrating a brighter world   NEC

# What is Different from GPGPU?

## GPGPU Architecture

**OS** **AP** CUDA → Function

Memory — x86 — **PCIe** — GPGPU — Memory

exec
Data Transmission
Result Transmission
:
exit

**Frequent PCIe transmission**

## Aurora Architecture

**OS** **AP**

Memory — x86 — **PCIe** — VE — Memory

exec
Start Processing
OS Function ← I/O,etc
exit
End Processing

**Whole AP is executed on VE**

**disadvantage**
- PCIe bottleneck
- Small memory
- Programming difficulty

→ **Advantage**
- Avoiding PCIe bottleneck
- Larger memory
- Standard language

Orchestrating a brighter world  **NEC**

## Programing Environment

**Vector Cross Compiler**

| automatic vectorization | automatic parallelization |

| | |
|---|---|
| Fortran: | F2003, F2008 |
| C/C++: | C11/C++14 |
| OpenMP: | OpenMP4.5 |
| Library: | MPI3.1, libc, BLAS, Lapack, etc |
| Debugger: | gdb, Eclipse parallel tools platform |
| Tools: | PROGINF, Ftrace Viewer |

```
$ vi sample.c
$ ncc sample.c
```

## Execution Environment

```
$ a.out
```

**Vector Host**

**execution**

\Orchestrating a brighter world **NEC**

Native Mode | Accelerator Mode | Scalar Acceleration Mode

Application

Native Mode: VE AP

Accelerator Mode: x86 AP → VEO → VE

Scalar Acceleration Mode: VE AP ← VH call, x86 →

OS: VEOS, Linux

Hardware: VH — VE

Orchestrating a brighter world    NEC

# Vector Engine

\Orchestrating a brighter world  **NEC**

32e

256e

8e

vector register
256e x 64
(128kB)

64e

A
B
C

D = A x B + C

32e

D

FMA x3

Vector Length = 256e (32e x 8 cycle)
307.2GF = 32Flops/cycle x 2(FMA) x 3 x 1.6GHz

Orchestrating a brighter world

NEC

# VE Processor

| VE10E Specification | |
|---|---|
| cores/CPU | 8 |
| core performance | ~307GF(DP) ~614GF(SP) |
| CPU performance | ~2.45TF(DP) ~4.91TF(SP) |
| cache capacity | 16MB shared |
| memory bandwidth | 1.35TB/s |
| memory capacity | 24, 48GB |



**2.45TF**

**307GF** core core core

core core core core

**0.4TB/s**

**3TB/s**

Software controllable cache **16MB**

**1.35TB/s**

HBM2 memory x 6

Orchestrating a brighter world **NEC**

**SPU**
- I cache:             32kB
- O cache:            32kB
- L2 cache:          256kB

**VPU**
- Vector Register (VR)
  128kB = 256e x 64
  376kB physical

**LLC**
- 16MB
- Write back
- ADB/MSHR functions

**Memory**
- 48GB
- HBM2 x6

Orchestrating a brighter world    **NEC**

# Memory Architecture



Single core

Memory (48GB)

Cache (16MB)

Vector register, 256e x 64 (128kB)

VPU (Vector Processing Unit)

Vector Pipeline x 32

VFMA0
VFMA1
VFMA2
ALU0
ALU1
DIV

**SPU**
Scalar Processing Unit

1.38TB/s / processor
(Ave. 170GB/s / core)

400GB/s / core

\Orchestrating a brighter world    **NEC**

**World's first implementation of 6 HBM2 memories**

Orchestrating a brighter world    **NEC**

# Card Implementation



- ■ Standard PCIe implementation
- ■ Connector: PCIe Gen.3 x16
- ■ Double height (same form factor as Nvidia)
- ■ <300W (DGEMM ~210W/VE, STREAM ~200W/VE, HPCG ~215W/VE)

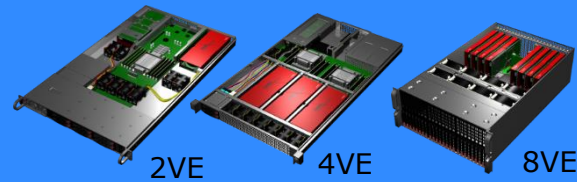\Orchestrating a brighter world   **NEC**

**A500**
**A400**

## Supercomputer Model

- For large scale configuration
- DLC with 40C water

**A300**

## Rack Mount Model

- Flexible configuration
- Air Cooled

2VE     4VE     8VE
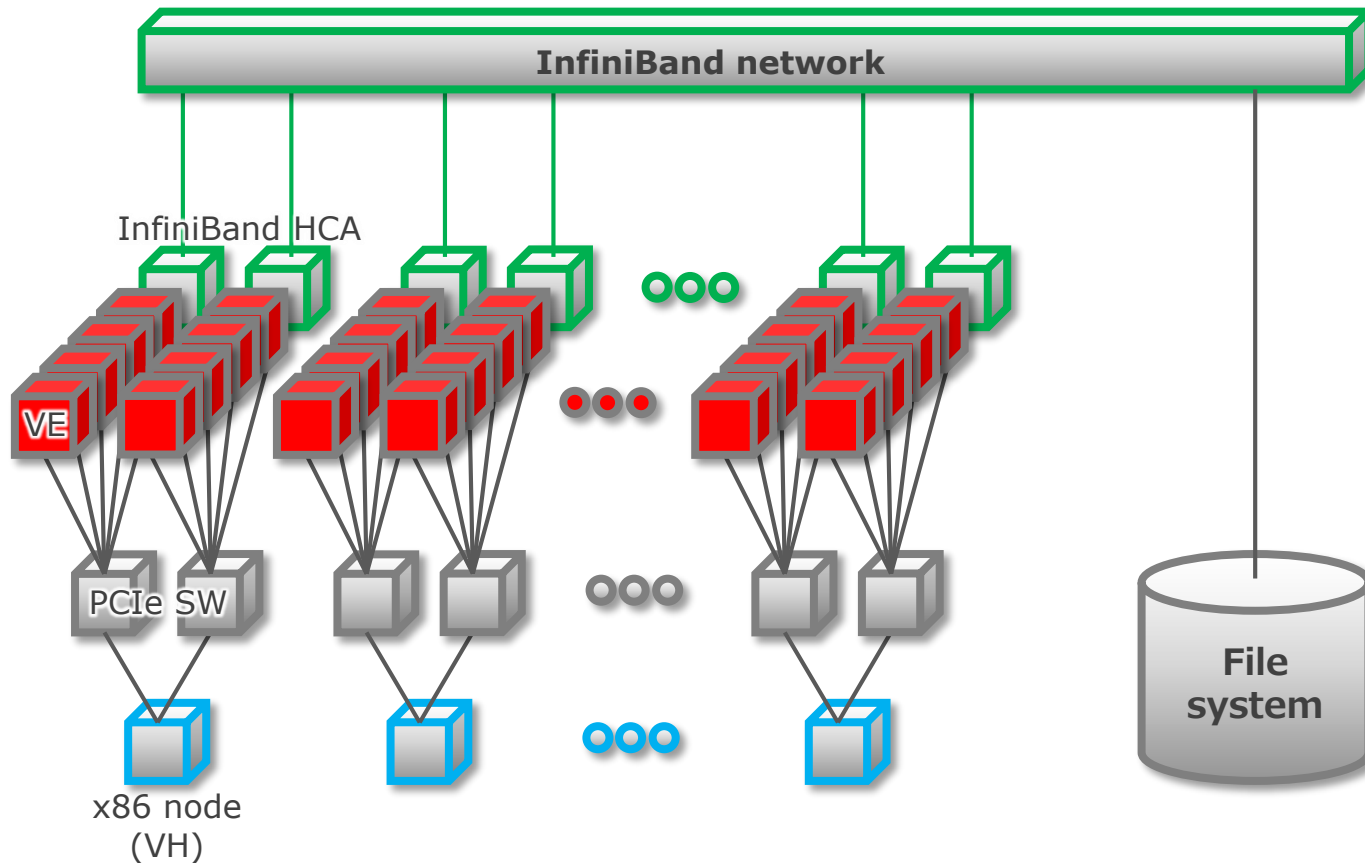
**A100**

## Tower Model

- For developer/programmer
- Tower implementation

1VE

Orchestrating a brighter world    **NEC**

InfiniBand network

InfiniBand HCA

VE

PCIe SW
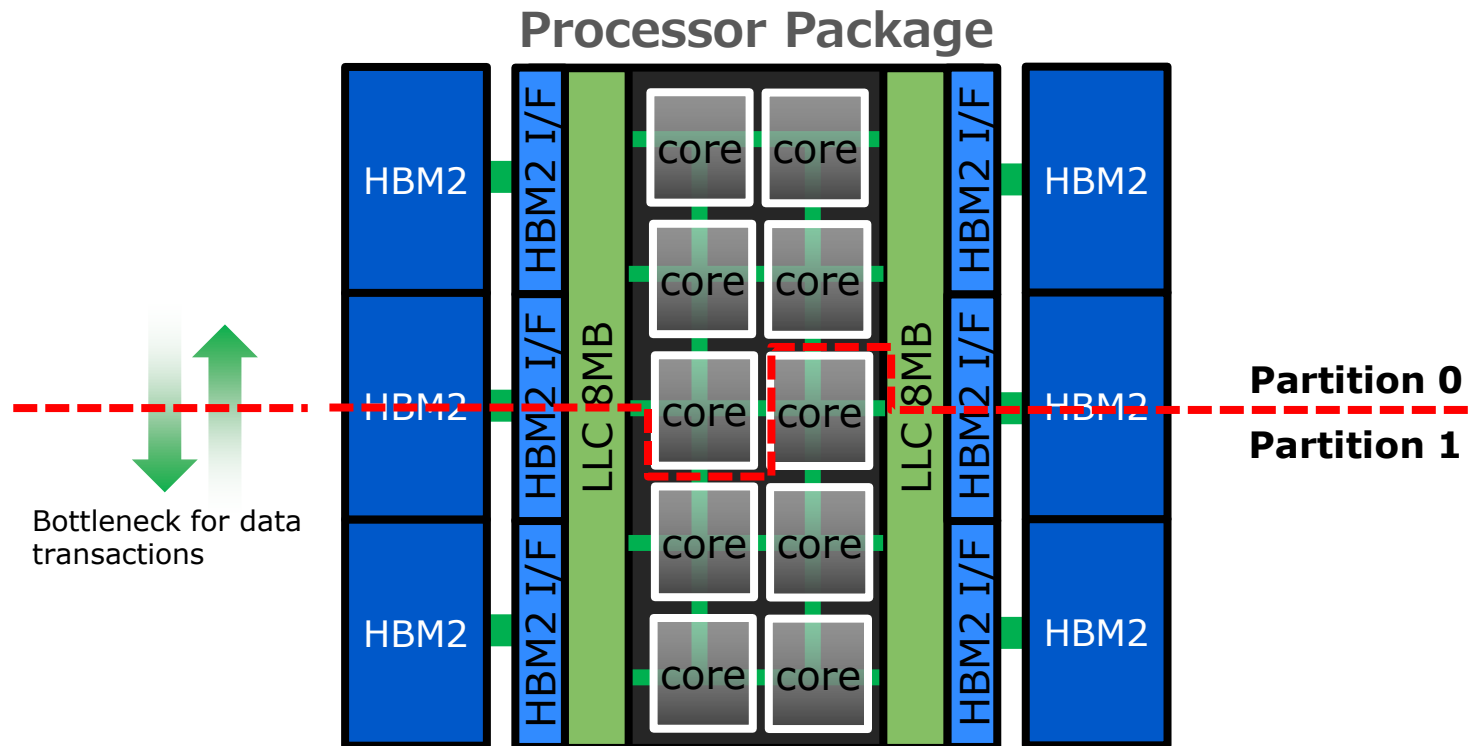
x86 node
(VH)

File system

Orchestrating a brighter world   NEC

# The latest version of Aurora Software

| Software | Version | Date |
|---|---|---|
| VEOS | 2.2.0 | 2019/10 |
| MMM | 1.2.17 | 2019/10 |
| VMC Firmware | 1.5.11-1 | 2019/10 |
| InfiniBand for SX-Aurora TSUBASA | 1.3.0 | 2019/10 |
| License Management | 1.3-1 | 2019/10 |
| SDK; NEC C/C++ Compiler, NEC Fortran Compiler | 2.5.1 | 2019/11 |
| SDK; Numeric Library Collection | 2.0.0-2.1 | 2019/10 |
| SDK; binutils | 2.26-2.3 | 2019/11 |
| SDK; Tuning tool PROGINF/FTRACE (veperf) | 2.1.0 | 2019/05 |
| SDK; Tuning tool NEC Ftrace Viewer | 1.0.0 | 2018/02 |
| SDK; NEC Parallel Debugger | 1.0.0 | 2018/02 |
| NEC MPI | 2.3.0 | 2019/10 |
| NEC Network Queuing System V (NQSV) ResourceManager / JobManipulator / JobServer | 1.04 | 2019/10 |
| NEC Scalable Technology File System (ScaTeFS) Server | 3.2 | 2019/10 |
| NEC Scalable Technology File System (ScaTeFS) Client | 3.0.30.5 | 2019/10 |

\Orchestrating a brighter world    NEC

# VE Partitioning Mode

\Orchestrating a brighter world **NEC**

# Partitioning Mode for Cache Acceleration

**Processor Package**



**Partition 0**
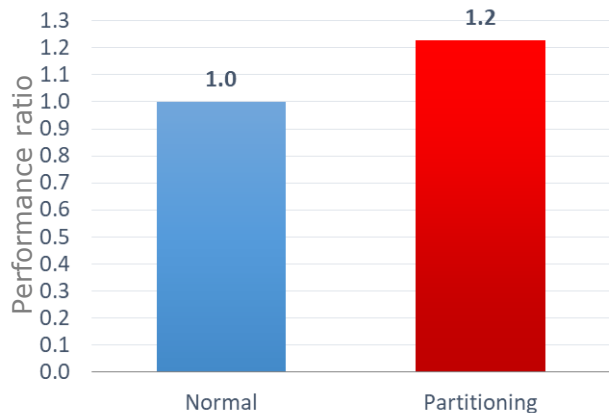
**Partition 1**

Bottleneck for data transactions

**The partitioning mode:**
One VE processor works as two partitions in order to avoid data transaction bottleneck. Due to this function, the sustained cache bandwidth is improved

\Orchestrating a brighter world  **NEC**

# Performance of the Partitioning Mode

## Front Flow Blue (CFD)



| | Normal | Partitioning |
|---|---|---|
| Performance ratio | 1.0 | 1.2 |

## HPCG

| | Normal | Partitioning |
|---|---|---|
| Performance ratio | 1.0 | 1.1 |

## HIMENO Benchmark (CFD, Poisson Equation)

| | Normal | Partitioning |
|---|---|---|
| Performance ratio | 1.0 | 1.0 |

**Cache bandwidth dependency**

Orchestrating a brighter world  NEC

# VE10E, VE20

\Orchestrating a brighter world    **NEC**

# Roadmap



Memory bandwidth / VE

**VE30**

**VE20**

3+TF
1.5TB/s memory bandwidth

**VE10E**

2.45TF
1.35TB/s memory bandwidth

**VE10**

2.45TF
1.22TB/s memory bandwidth

2019    2020    2021    2022

\Orchestrating a brighter world    **NEC**

Memory capacity / Memory bandwidth

| Memory capacity | Memory bandwidth |
| --- | --- |
| | 1.53TB/s |
| 48GB | 1.35TB/s |
| | 1.22TB/s |
| 24GB | 0.75TB/s |

20B  20A

10BE*  10AE*

*10AE is 1.584GHz
*10BE is 1.408GHz

10B  10A

10C

| | 2.15TF | 2.45TF | 3.07TF |
| --- | --- | --- | --- |
| Frequency | 1.4GHz | 1.6GHz | |
| Cores | 8core | | 10core |

Orchestrating a brighter world **NEC**

# Performance

\Orchestrating a brighter world **NEC**

# DGEMM (Calculation Capability Evaluation)



**HPL/DGEMM like application is not target of Aurora**

Chart: Sustained performance [TF]

| Processor | Performance |
|-----------|-------------|
| Xeon Gold 6142 (2CPU) 2017 | 2.10 |
| Tesla V100 (1GPU) 2018 | 6.63 |
| ARM A64FX (1CPU) 2021 | 2.50 |
| Aurora VE10A (1CPU) 2018 | 2.43 |
| Aurora VE10AE (1CPU) 2019 | 2.40 |
| Aurora VE20A (1CPU) 2020 | 3.03 |

V100 result: AMD NEXT HORIZON
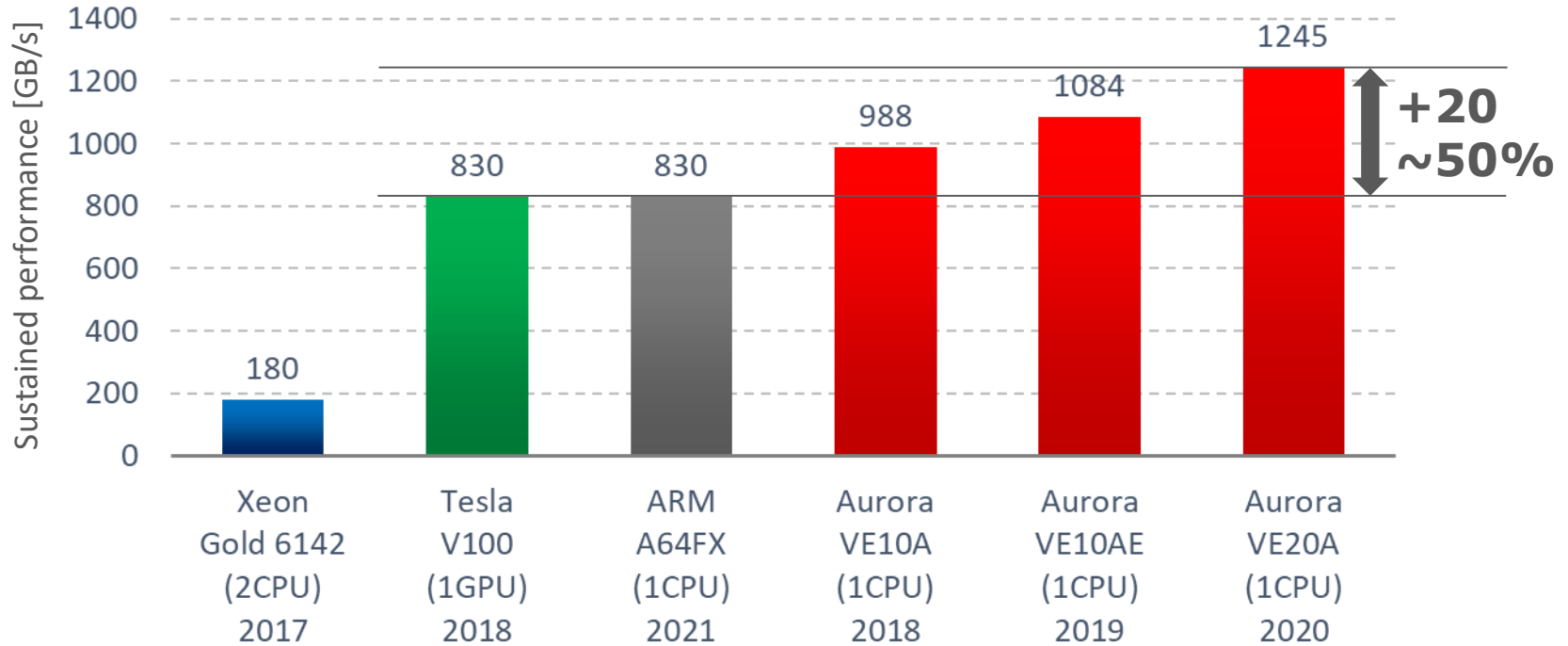http://ir.amd.com/static-files/ef99f84b-e1ad-4e12-8058-f3488f4c47b7

ARM A64FX result: The post-K project and Fujitsu ARM-SVE enabled A64FX processor
https://indico.math.cnrs.fr/event/4705/attachments/2362/2942/CEA-RIKEN-school-19013.pdf

**VE10A: 210W/card**

© NEC Corporation 2019

Orchestrating a brighter world  NEC

# STREM TRIAD (Memory Bandwidth Evaluation)



STREAM TRIAD (Memory Bandwidth Evaluation)

Sustained performance [GB/s]

| | Xeon Gold 6142 (2CPU) 2017 | Tesla V100 (1GPU) 2018 | ARM A64FX (1CPU) 2021 | Aurora VE10A (1CPU) 2018 | Aurora VE10AE (1CPU) 2019 | Aurora VE20A (1CPU) 2020 |
|---|---|---|---|---|---|---|
| Value | 180 | 830 | 830 | 988 | 1084 | 1245 |

**+20 ~50%**

**VE10A: 195W/card**

ARM A64FX result: The post-K project and Fujitsu ARM-SVE enabled A64FX processor
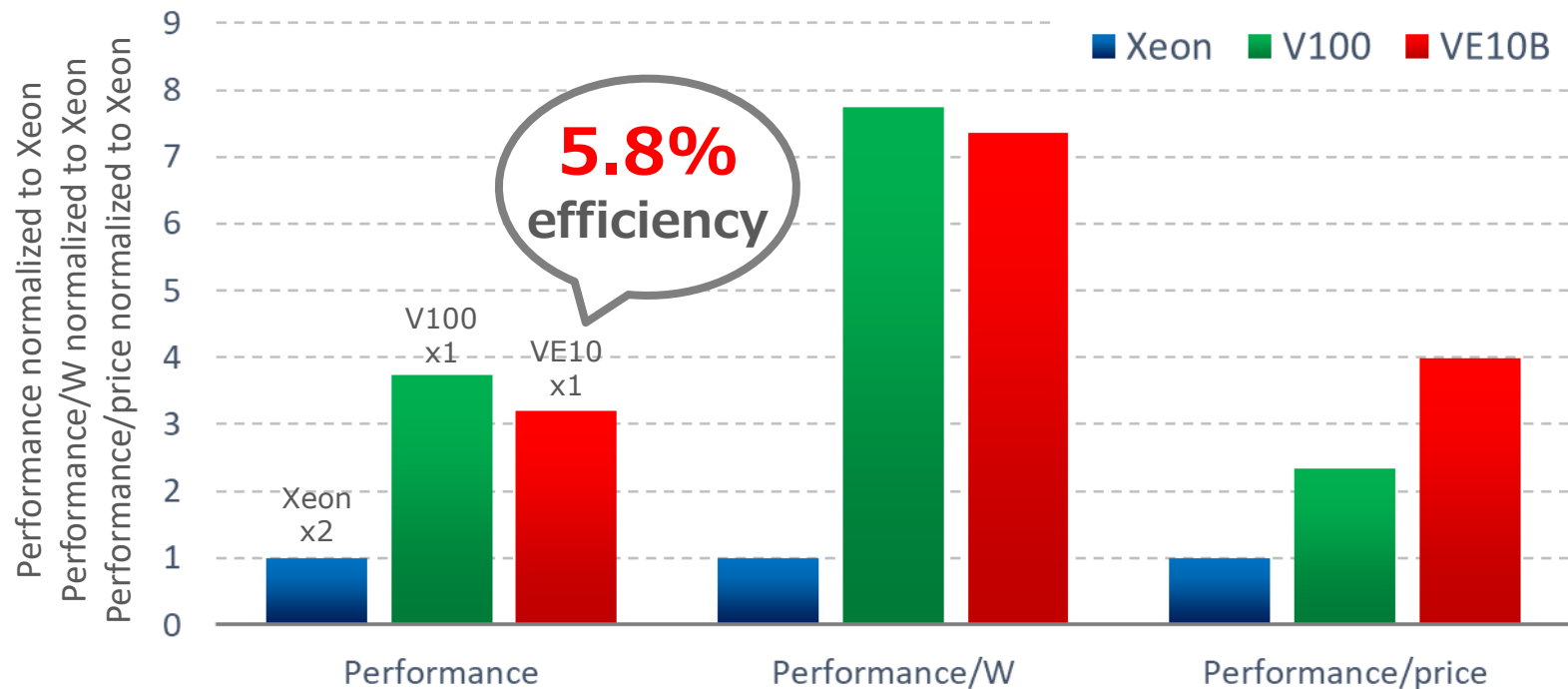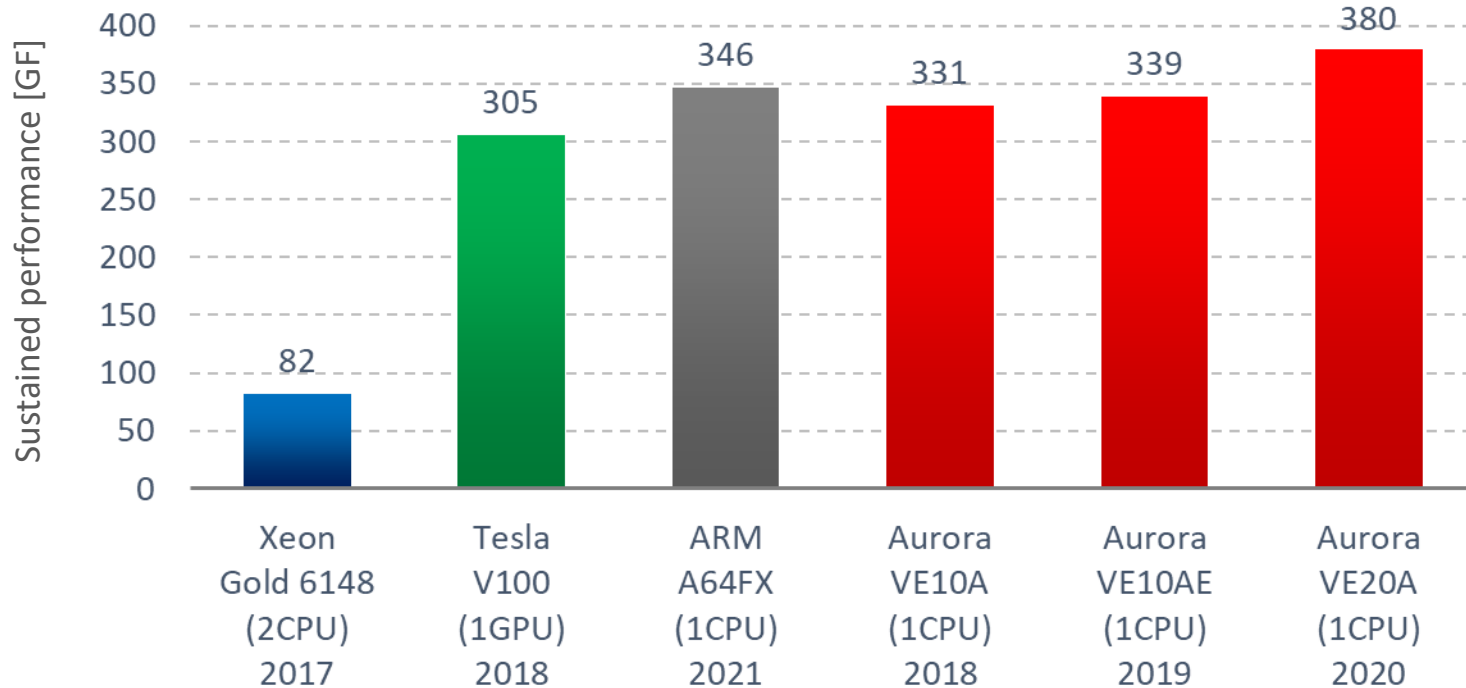https://indico.math.cnrs.fr/event/4705/attachments/2362/2942/CEA-RIKEN-school-19013.pdf

Orchestrating a brighter world  **NEC**

**VE10A: 200W/card**

# HIMENO Benchmark (CFD, Poisson Equation)



Sustained performance [GF]

| | Xeon Gold 6148 (2CPU) 2017 | Tesla V100 (1GPU) 2018 | ARM A64FX (1CPU) 2021 | Aurora VE10A (1CPU) 2018 | Aurora VE10AE (1CPU) 2019 | Aurora VE20A (1CPU) 2020 |
|---|---|---|---|---|---|---|
| Value | 82 | 305 | 346 | 331 | 339 | 380 |

V100 result: Performance evaluation of a vector supercomputer SX-aurora TSUBASA
https://dl.acm.org/citation.cfm?id=3291728

ARM A64FX result: Supercomputer "Fugaku" Formerly known as Post-K
https://www.fujitsu.com/global/Images/supercomputer-fugaku.pdf

## VE10A: 210W/card

Orchestrating a brighter world   NEC

## Frovedis: NEC's Sparc library fully optimized for Aurora



Bar chart — Speedup normalized to Xeon

Legend:
- Sparc/Xeon Gold 6126 12 cores
- Frovedis/Aurora VE10B 8 cores

| | Logistic Regression | Singular value decomposition | K-means |
|---|---|---|---|
| Sparc/Xeon Gold 6126 12 cores | 1.0 | 1.0 | 1.0 |
| Frovedis/Aurora VE10B 8 cores | 113.2 | 56.8 | 42.8 |

\Orchestrating a brighter world  **NEC**

# Frovedis: NEC's Sparc library fully optimized for Aurora



Q1: group by/aggregate
Q3: filter, join, group by/aggregate
Q5: filter, join, group by/aggregate (larger join)
Q6: filter, group by aggregate

Chart values (Frovedis/VE): Q1: 10.1, Q3: 33.8, Q5: 47.3, Q6: 34.8

Legend: ■ Spark/x86  ■ Frovedis/VE

\Orchestrating a brighter world    NEC

# High Density Product A412-8

A4:         A400 Series
1:           VE10E Generation
2:           Rome processor
-8:         Maximum number of VEs

\Orchestrating a brighter world  **NEC**

Direct Liquid Cooling
## VE DLC
To provide higher density
Hot water cooling

High density VH server



## 8VE/2U

**One VH consists of:**
VE x8
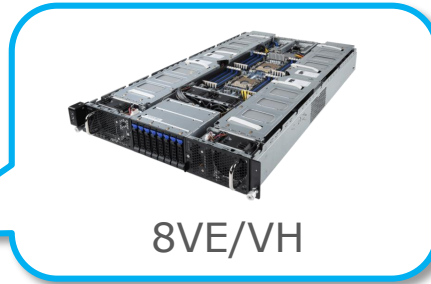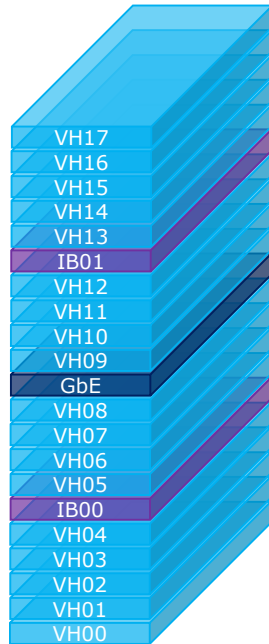AMD Rome processor x1
IB HCA x2

Performance:          19.6TF /VH
Memory bandwidth:  10.8TB/s /VH

# High Density Rack w/ DLC

**144VE/rack can be operated with 35C water**
- High density and hot water cooling model
- VE10AE/10BE are supported



8VE/VH

## x18

144VE/rack
Up to 352TF/rack
194TB/s/rack

| | |
|---|---|
| ■ 194TB/s | = x86 processor x 1000 |
| ■ Power: | 30kW for application run |
| ■ Cooling: | DLC for VE and X86 |
| ■ Dimensions: | W800 x D1400 x H2200 (47U) |

Rack labels (top to bottom): VH17, VH16, VH15, VH14, VH13, IB01, VH12, VH11, VH10, VH09, GbE, VH08, VH07, VH06, VH05, IB00, VH04, VH03, VH02, VH01, VH00

Orchestrating a brighter world  NEC