Empowered by Innovation **NEC**

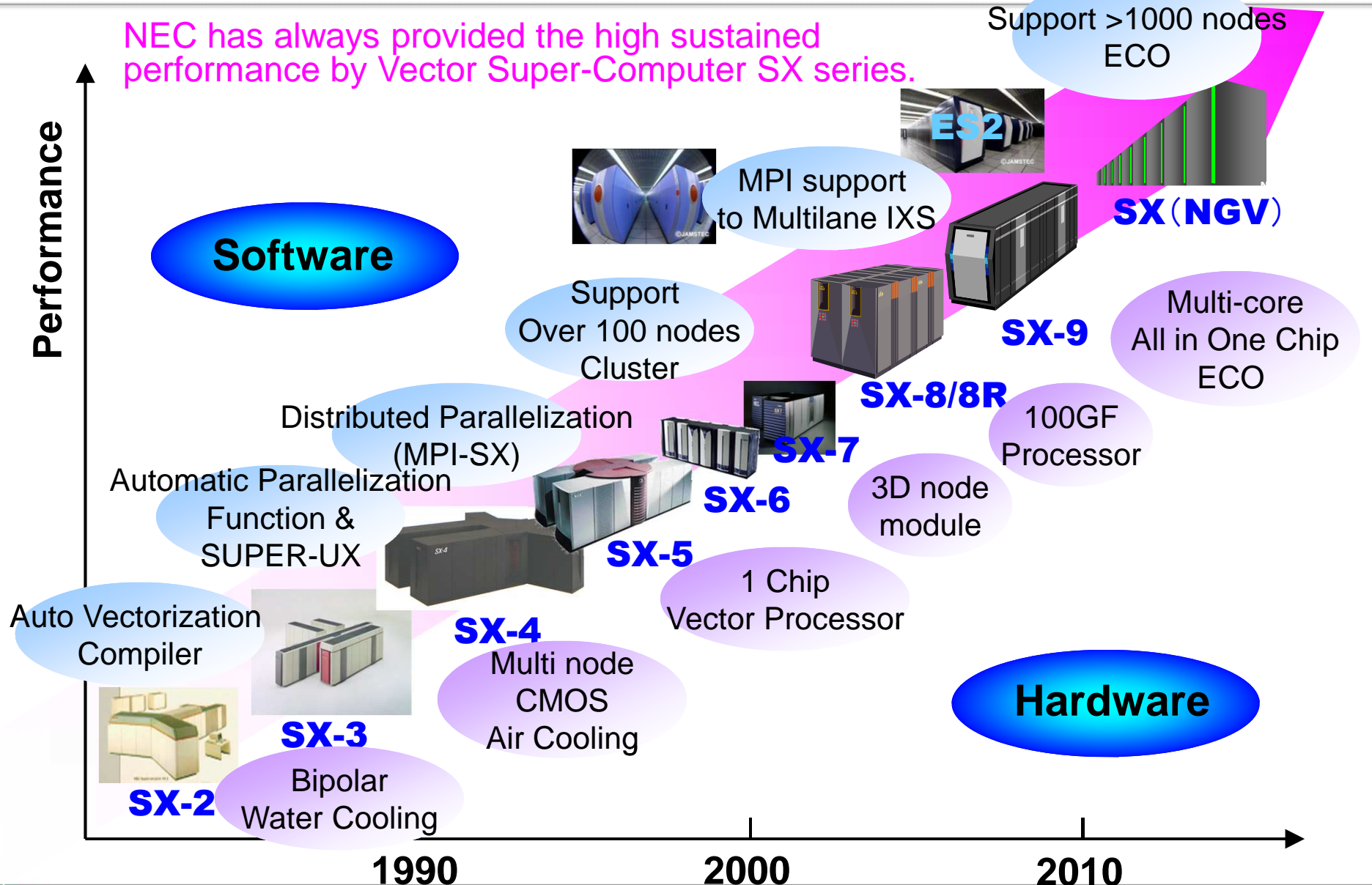# NEC Hybrid Solutions for Meteo Sites

**October 2nd, 2012**

**NEC Corporation**

NEC

# Hybrid Concept, our strategy

▌ COTS(Commercial Off-The-Shelf) are adequate for quite some applications.

▌ But they are not the answer to every HPC-challenge.

▌ Consequently NEC will continue a proprietary vector architecture.

▌ The seamless integration of the vector-system with one build from standard components is the key of NEC's strategy.

▌ In particular when complicated workflows need to be mapped on the best, i.e. most efficient hardware platform, as it is the case in production environments in the weather forecast business.
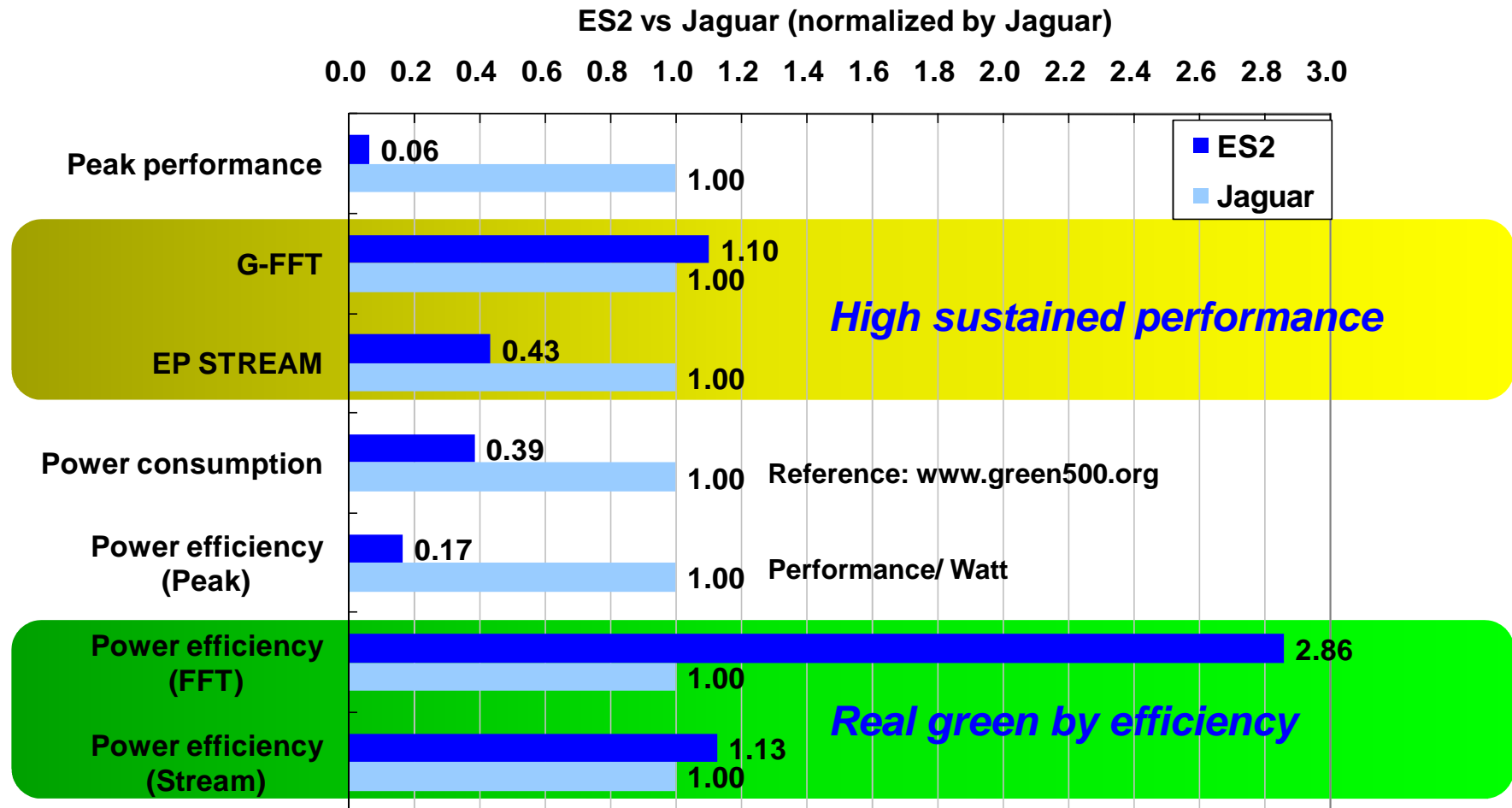
Empowered by Innovation  NEC

# SX History and Technical evolutions



NEC has always provided the high sustained performance by Vector Super-Computer SX series.

Performance

**Software**

**Hardware**

Support >1000 nodes
ECO

ES2

MPI support
to Multilane IXS

**SX(NGV)**

Support
Over 100 nodes
Cluster

Multi-core
All in One Chip
ECO

**SX-9**

Distributed Parallelization
(MPI-SX)

**SX-8/8R**

100GF
Processor

Automatic Parallelization
Function &
SUPER-UX

**SX-7**

3D node
module

**SX-6**

**SX-5**

Auto Vectorization
Compiler

1 Chip
Vector Processor

**SX-4**

Multi node
CMOS
Air Cooling

**SX-3**

Bipolar
Water Cooling

**SX-2**

1990          2000          2010

Empowered by Innovation  **NEC**

# NGV HARDWARE

Empowered by Innovation

# Which is smarter ?

- **Break the POWER WALL by "High computational efficiency"**
- **Higher sustained performance and efficiency are "SX DNA"**

**ES2 vs Jaguar (normalized by Jaguar)**

| | | | |
|---|---|---|---|
| Peak performance | ES2: 0.06 | Jaguar: 1.00 | |
| G-FFT | ES2: 1.10 | Jaguar: 1.00 | *High sustained performance* |
| EP STREAM | ES2: 0.43 | Jaguar: 1.00 | |
| Power consumption | ES2: 0.39 | Jaguar: 1.00 | Reference: www.green500.org |
| Power efficiency (Peak) | ES2: 0.17 | Jaguar: 1.00 | Performance/ Watt |
| Power efficiency (FFT) | ES2: 2.86 | Jaguar: 1.00 | *Real green by efficiency* |
| Power efficiency (Stream) | ES2: 1.13 | Jaguar: 1.00 | |

Legend: ■ ES2  ■ Jaguar

Axis: 0.0 0.2 0.4 0.6 0.8 1.0 1.2 1.4 1.6 1.8 2.0 2.2 2.4 2.6 2.8 3.0

Empowered by Innovation **NEC**

# Targets of Next Generation Vector

## High sustained performance
World's top-class core performance (64GF)
World's top-class memory bandwidth (64GB/s)

**Inherit SX-DNA**

## Low power consumption
World's top-class energy-saving supercomputer
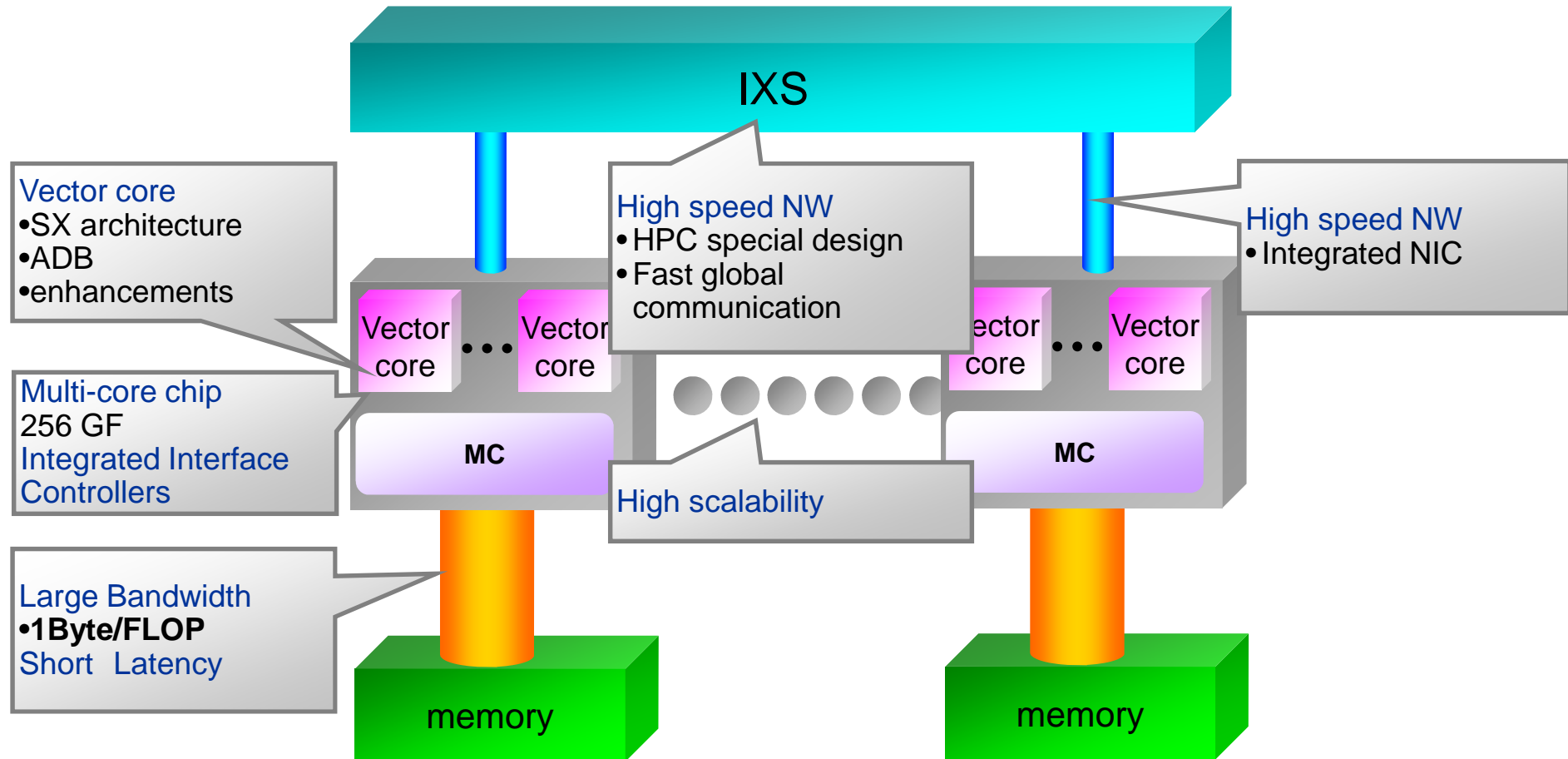
$\dfrac{1}{10}$ compared to SX-9

## Small installation space
Reduce floor space cost
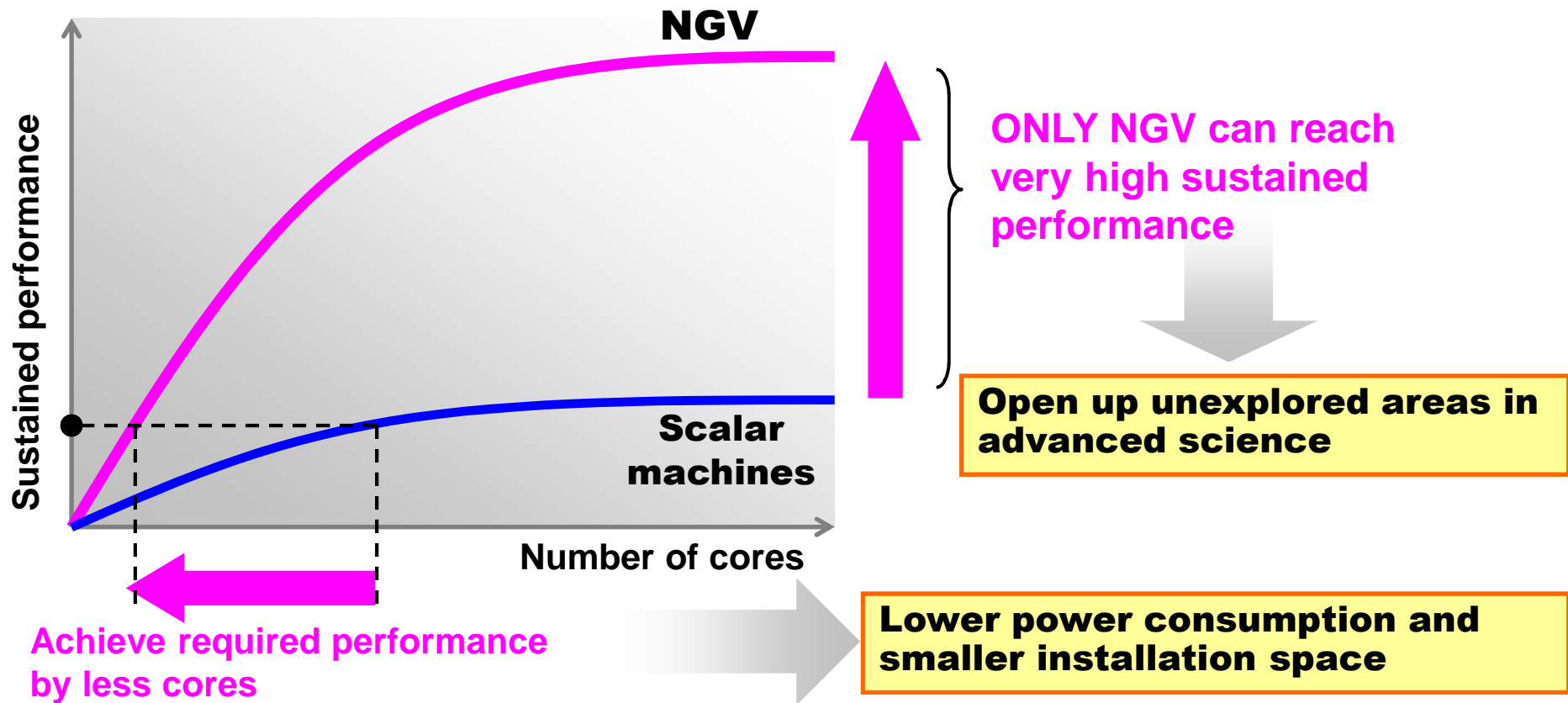
$\dfrac{1}{5}$ Compared to SX-9

Empowered by Innovation  **NEC**

# Next Vector Configuration

The next generation multi-core vector architecture provides
high sustained performance at low power consumption

IXS

**Vector core**
- SX architecture
- ADB
- enhancements

**High speed NW**
- HPC special design
- Fast global communication

**High speed NW**
- Integrated NIC

Vector core ... Vector core

Vector core ... Vector core

**Multi-core chip**
256 GF
Integrated Interface Controllers

MC

MC

High scalability

**Large Bandwidth**
- **1Byte/FLOP**
Short Latency

memory

memory

Empowered by Innovation  **NEC**

# Higher Sustained Performance by Powerful Core

- ■Very high system SUSTAINED performance = SX DNA
  - →required performance with less cores
    - → parallel efficiency
    - → lower power and space
- ■ "Green HPC"

**NGV**

**Sustained performance** (y-axis)

**Number of cores** (x-axis)

**Scalar machines**

**ONLY NGV can reach very high sustained performance**

**Open up unexplored areas in advanced science**

**Achieve required performance by less cores**

**Lower power consumption and smaller installation space**

©NEC Corporation, 2012

Empowered by Innovation **NEC**

# All-in-one Processor

■ **4 powerful cores and each interface controllers are integrated in one-CPU → Power saving**
■ **Compact card design → Space saving**

**NGV CPU**

**CPU card**

**I/O Controller**
Connection to storage device, Ethernet

**Network controller**
8GB/s/direction, Fat-tree

**Powerful core**
World's fastest CPU core
64GF x 4cores
1MB cache/core (target)

memory

Performance:        256GF
Memory bandwidth: 256GB/s

**Memory Controller**
256GB/s BW control

**Very large memory BW**
World's largest BW 256GB/s

37cm

11cm

Empowered by Innovation    **NEC**

# Packaging

## A prototype of 2 nodes card module



**Water Cooling**
CPU is cooled by water

**CPU**
4 cores
256GF, 256GB/s, 1B/F

**Memory**
16 DIMMs / CPU
64GB

**Node card**
1 CPU
11cm x 37cm

Empowered by Innovation **NEC**

# Rack Implementation

- 64 nodes are implemented in one rack = 16TF
- hybrid cooling by air and liquid
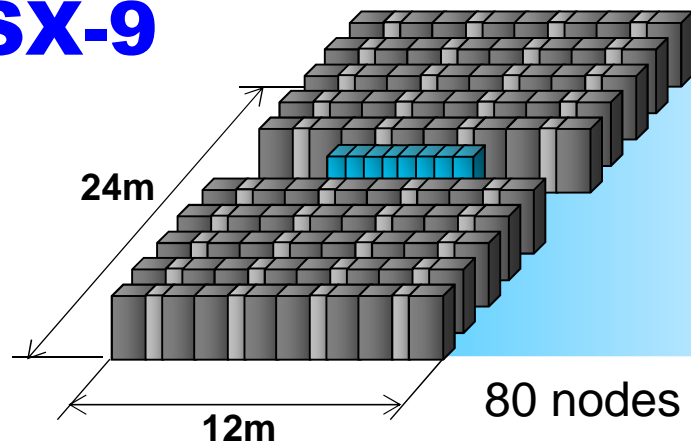- 10x higher performance and over 2x smaller rack vs. SX-9

**Node card 256GF**

110

37cm

330

11cm

**Rack 64 nodes 16TF**

64 cards

2.0m

1.2m

0.7m

**SX-9 rack 1 node 1.6TF**

1.8m

1.8m

1.1m

Empowered by Innovation

**NEC**

# Downsizing and Power Saving

Providing 5x smaller space and 10x lower power consumption compared to SX9 by GREEN design and compact implementation.

**Comparison with same performance (131TF)**

**SX-9**

24m

12m

80 nodes

**25m pool size**

**SX (NGV)**

7m

8m

512 nodes

**Meeting room size**

| | | |
|---|---|---|
| 131TF | | 131TF |
| 288m$^2$ | space 1/5 | 56m$^2$ |
| 2.4MW | power 1/10 | 0.24MW |

Empowered by Innovation **NEC**

# NGV SOFTWARE

Empowered by Innovation    **NEC**

# NGV Software Overview

## Provide hybrid cluster solution
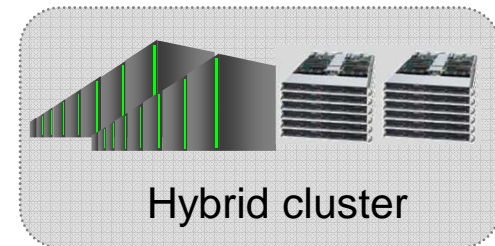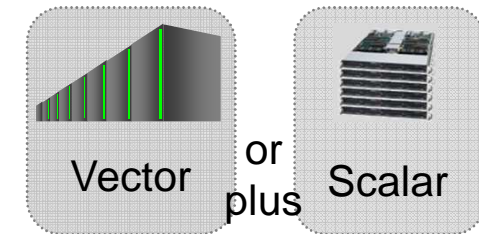### - Integrate vector and scalar cluster as a single system -

### Background

- Demanding computation power for HPC applications
- Not one kind of architecture will fulfill all requirements

Vector   or plus   Scalar

### Solution

- Provide hybrid solution
  - Job collaboration using workflow tools
  - Integrated scheduling (assign right node to right job)
  - New shared file system
- Provide large cluster solution
  - Integrated single system management of vector and scalar cluster
  - Enhanced scalability and reliability

Hybrid cluster

- And much more
  - Sophisticated OS and compiler compliant with standards
  - MPI-3 support, enhanced performance (memory and interconnect)
  - User-friendly tools, easy-to-use debugging environment, etc.

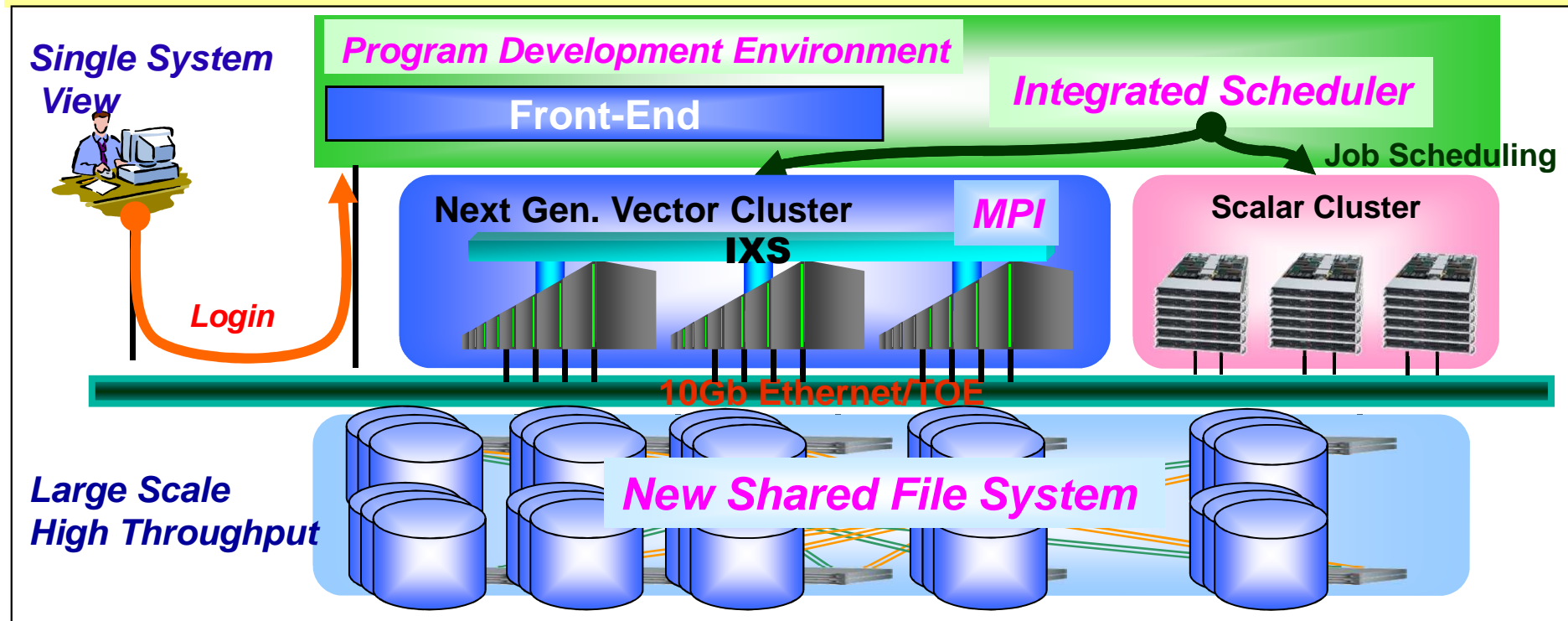**Integrated scheduler (Batch system)**

**Shared File System**

**OS & Compiler**

**MPI**

**Tools**

Empowered by Innovation   NEC

# System Overview

## ■Single system solution – Integrated scheduler –

- ●Supports vector clusters and scalar clusters together as a single system
- ●Easy to manage a system with more than 1000 nodes

*Single System View*

*Program Development Environment*

**Front-End**

*Integrated Scheduler*

Job Scheduling

**Next Gen. Vector Cluster** *MPI*

**IXS**

**Scalar Cluster**

*Login*

10Gb Ethernet/TOE

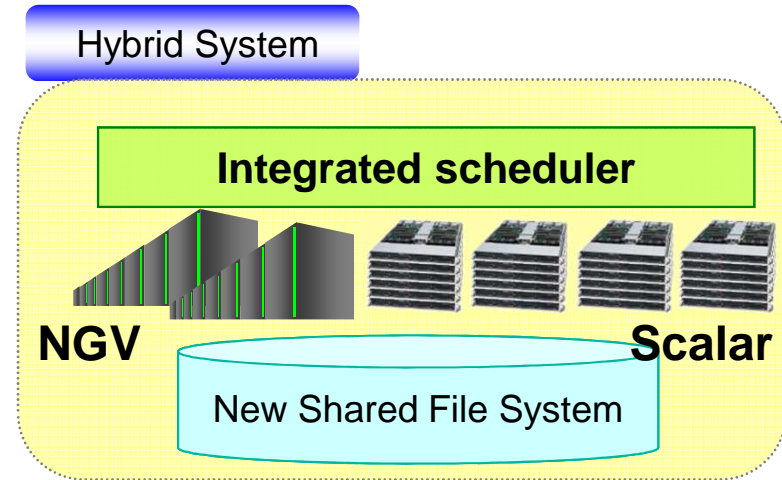*Large Scale High Throughput*

*New Shared File System*

## ■New I/O solution

- ●Realizes new shared file system with huge capacity and large scale using multiple IO servers
- ●Provides high speed IO using proprietary protocol lighter than NFS

Empowered by Innovation **NEC**

# Integrated Scheduler

**Integrated scheduler realizes enhanced hybrid system running real workflow!**
- Vector cluster and scalar cluster managed as a single system
- Collaboration scheduling of vector jobs and scalar jobs using a **workflow script**

Hybrid System

Integrated scheduler

**NGV**　　　　　　　　　　　　**Scalar**

New Shared File System

Vector and scalar system closely coupled

**Easy operation of large scale cluster system**
- Enhanced scalability
- Inter cluster scheduling
- Ensemble job (parameter sweep)

Empowered by Innovation　**NEC**

# Job collaboration realizes seamless usage of hybrid system

## Efficient execution of Job collaboration

- Job execution order can be specified by workflow tools
  - Serial/Parallel execution
  - Conditional branch by exit code

- Collaboration scheduling of vector and scalar jobs
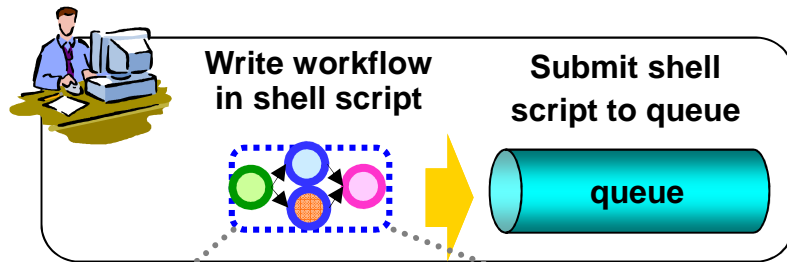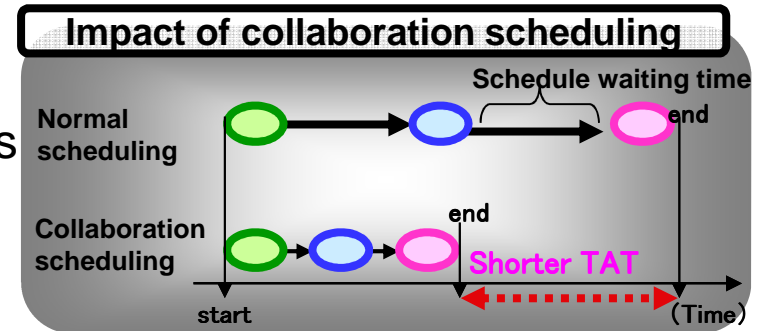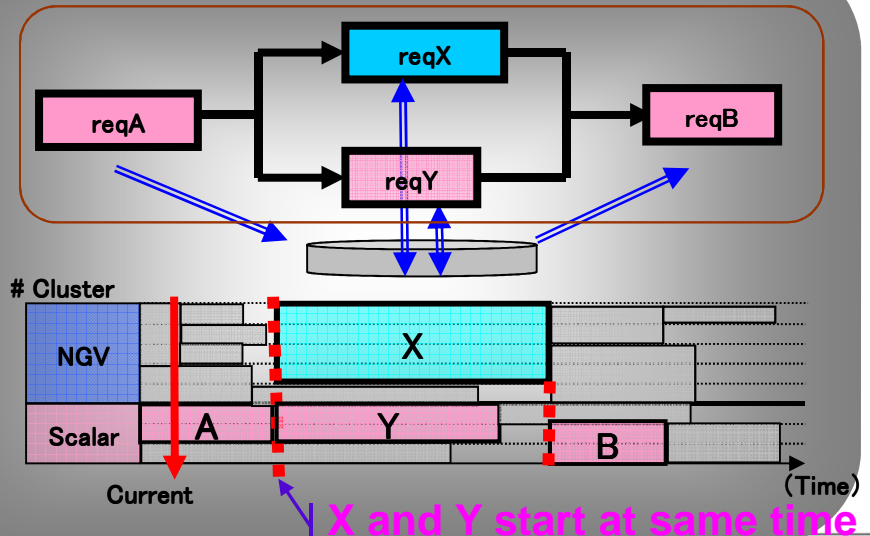  - Collaboration jobs to be executed consecutively to shorten TAT

**Impact of collaboration scheduling**

Schedule waiting time

Normal scheduling

Collaboration scheduling

end

Shorter TAT

start (Time)

**Write workflow in shell script**

**Submit shell script to queue**

queue

**Image of shell script**

```
#!/bin/bash

REQA=`qsub -q Scalar reqA`
REQX=`qsub -q NGV --after $REQA reqX`
REQY=`qsub -q Scalar --same $REQX reqY`
REQB=`qsub -q Scalar --after $REQX,$REQY reqB`

qwait $REQB
```
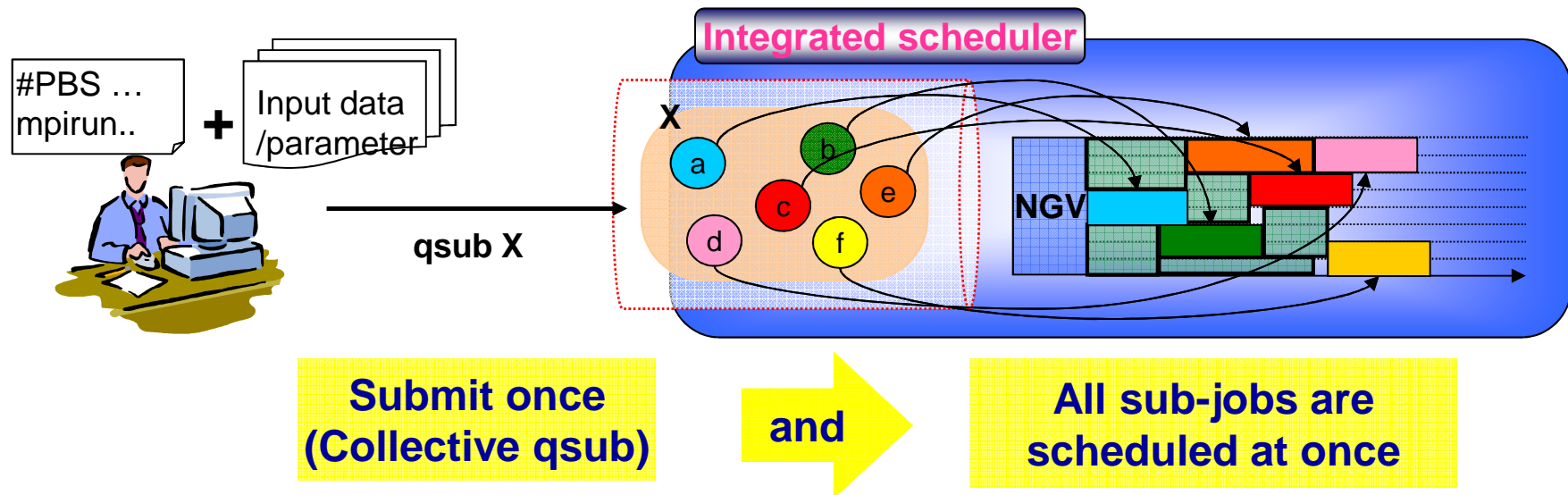
**Workflow execution and collaboration scheduling**

reqX

reqA

reqY

reqB

\# Cluster

NGV

X

Scalar

A

Y

B

Current

(Time)

X and Y start at same time

©NEC Corporation, 2012

Empowered by Innovation **NEC**

# Ensemble job supported

## Run same job with many different parameters (parameter sweep)

- Submit once and thousands of jobs are scheduled immediately
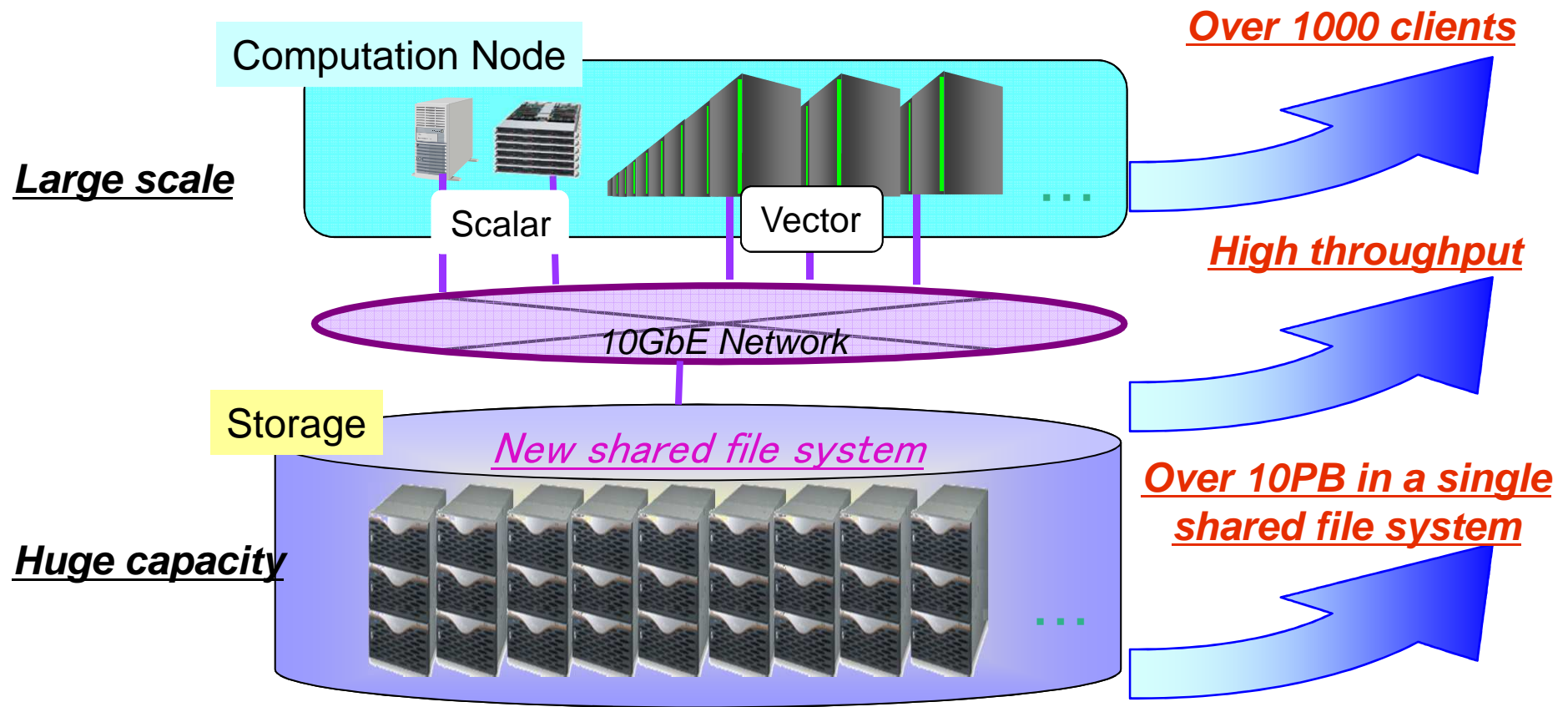- Sub-job for each input file is created automatically

**Integrated scheduler**

```
#PBS …
mpirun..
```
+
```
Input data
/parameter
```

qsub X

X: a b c d e f

NGV

**Submit once (Collective qsub)**   and   **All sub-jobs are scheduled at once**

- Collective qdel, specific qdel, etc. are also supported

```
Ex) qdel X : delete all sub-jobs
    qdel X.a : delete sub-job X.a
```

- Convenient parameter generation features
  - Sequential data generation (sub-job number, date, time, etc.)
  - Generate parameter from listed filename, etc.

```
Ex) IN-DATA.%(date:0530-0601)
    -> IN-DATA.0530
       IN-DATA.0531
       IN-DATA.0601
```

Empowered by Innovation **NEC**

# New shared file system

New shared file system provides fast I/O to massive data

**Large scale**

Computation Node

Scalar

Vector

**Over 1000 clients**

10GbE Network

**High throughput**

Storage

*New shared file system*

**Huge capacity**

**Over 10PB in a single shared file system**

Empowered by Innovation **NEC**

# File system image

Provide single file system view for large scale cluster

©NEC Corporation, 2012

Empowered by Innovation **NEC**

# Performance improvement  - key points -
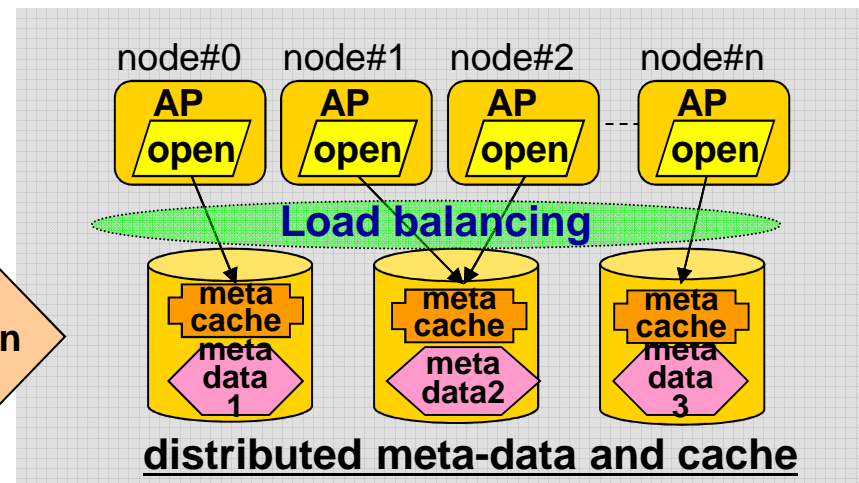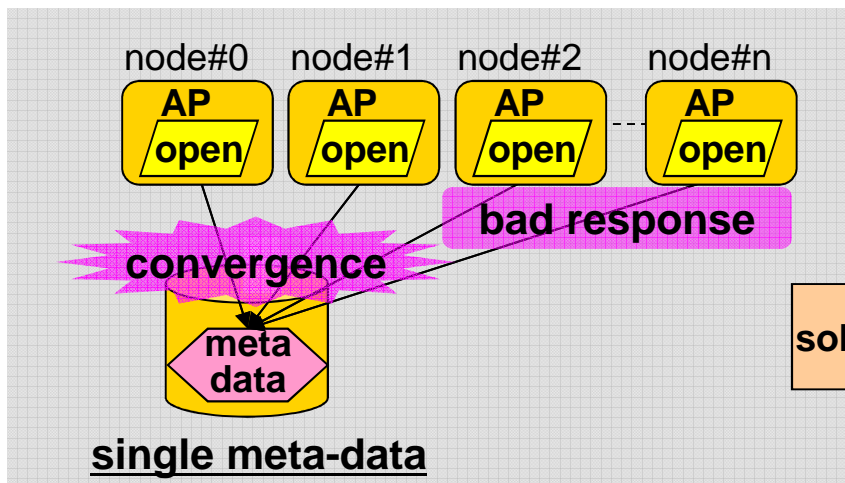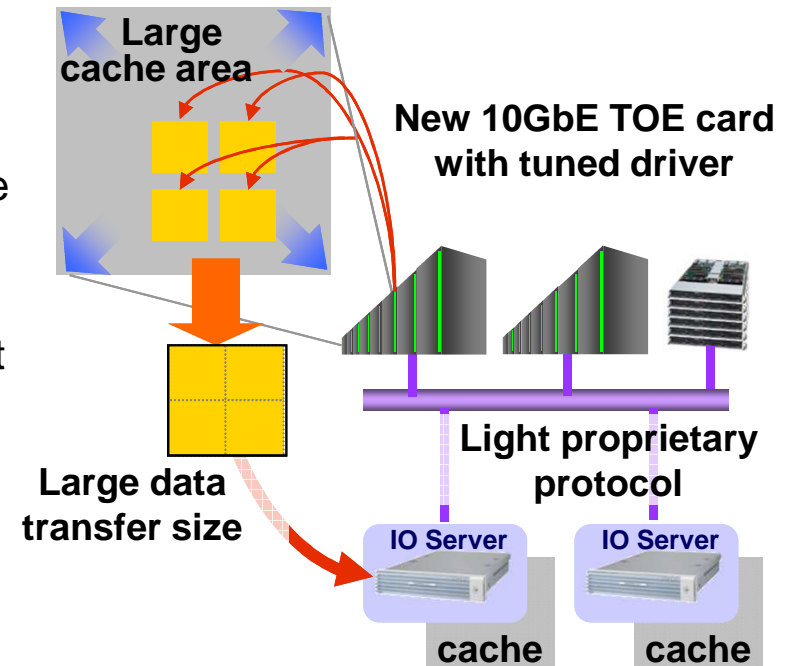
## Light proprietary I/O protocol

- Efficient data transfer between server and client (Gather data into single request and send together)
- Adopt next generation 10GbE TOE card and optimize network driver for the new card

## Data cache

- Efficient IO handling using large data cache on client and IO servers.

## Avoid congestion on meta-data server!

- **Meta-data distribution,** not single MDS like Lustre
- Reduce network traffic using meta-data cache

Large cache area

New 10GbE TOE card with tuned driver

Light proprietary protocol

Large data transfer size

IO Server    IO Server

cache    cache

node#0    node#1    node#2    node#n

AP open    AP open    AP open    AP open

bad response

convergence

meta data

single meta-data

solution

node#0    node#1    node#2    node#n

AP open    AP open    AP open    AP open

Load balancing

meta cache    meta cache    meta cache

meta data 1    meta data2    meta data 3

distributed meta-data and cache

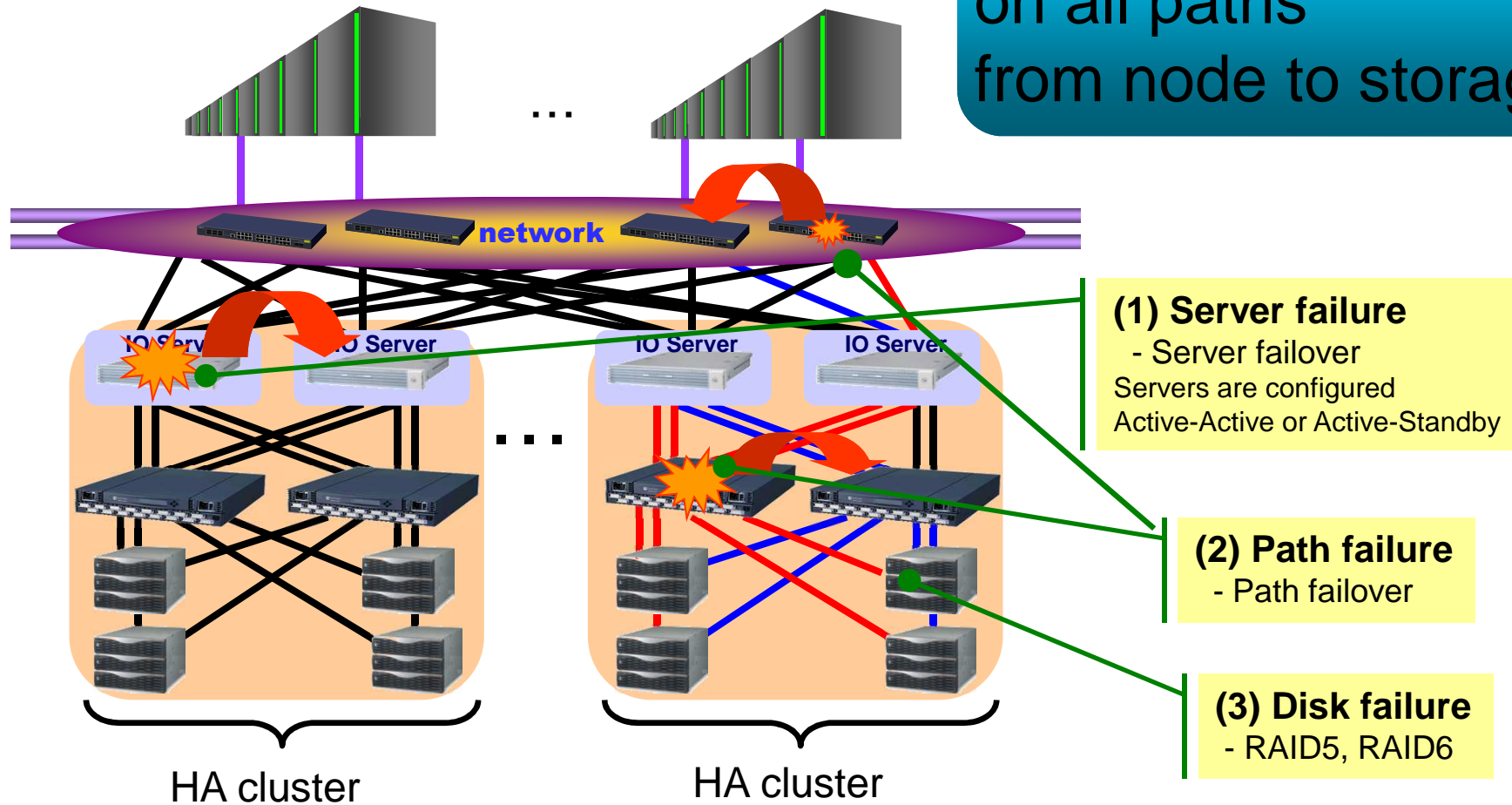Empowered by Innovation    NEC

# Reliability

**Responses to system failure**
(Realization of high availability)
- System continuity and data preservation

**Possible redundancy on all paths from node to storage**



network

IO Server    IO Server    IO Server    IO Server

**(1) Server failure**
 - Server failover
Servers are configured
Active-Active or Active-Standby

**(2) Path failure**
 - Path failover

**(3) Disk failure**
 - RAID5, RAID6

HA cluster          HA cluster
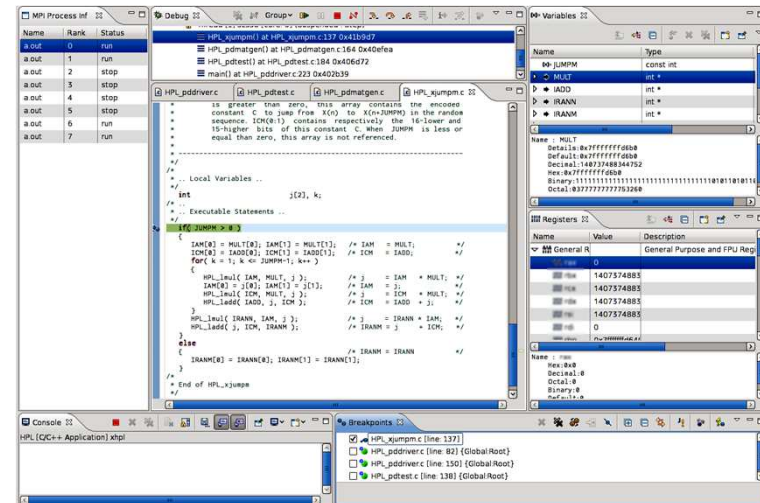
Empowered by Innovation   **NEC**

# And much more ...

**Fortran 2003, C/C++ compilers and MPI**
- ISO standard compliant
- Sophisticated automatic vectorization and parallelization
- Sophisticated usage of ADB (~ vector cache)
- Automatic optimization and optimization by directives
- OpenMP and MPI-3 support

**GUI Tools and debuggers**
- GUI Performance analysis/tuning tool
- GUI debugger

All software are tuned and optimized to
extract best performance out of NGV

Empowered by Innovation    NEC

Empowered by Innovation

**NEC**