

# 事業仕分けに負けない スパコンの作り方

實田 健

日本ヒューレット・パカード株式会社

2010年7月10日



# 会社概要



# 日本ヒューレット・パッカーード株式会社概要



設立	1963年 横河ヒューレット・パッカーード設立 －1999年7月 計測事業と分社化
代表取締役 社長執行役員	小出 伸一
事業	コンピューター、コンピューターシステム、 コンピューター周辺機器、 ソフトウェア製品の開発・製造・輸入・ 販売・リース・レンタルおよびサポート
本社	東京都千代田区五番町7番地
資本金	100億円
売上高	3,630億円(2009年10月期)
社員数	約5,400名(2010年2月現在)
セールス/ サポート拠点	全国39カ所

スパコンは何ですか？



# スパコンとは

スーパーコンピュータ(略称:スパコン)とは、内部の演算処理速度がその時代の一般的なコンピュータより極めて高速な計算機(コンピュータ)のこと。HPCサーバ(High Performance Computing Server)とも呼ばれる。

スーパーコンピュータの定義は時代によって大きく変化するが、一般的にはその時代の最新技術が投入された最高性能の計算機を指す。現時点では一般的に使用されるサーバ機よりも浮動小数点演算が1,000倍以上速いコンピュータを「スーパーコンピュータ」と呼ぶことが多い。

日本の文部科学省の科学技術・学術審議会では2005年現在、1.5TFLOPS以上の演算性能を持つコンピュータを政府調達における「スーパーコンピュータ」と位置付けている。



出展: wikipedia

# スパコンとは(続)

文部科学省の科学技術・学術審議会 のリンクをたどってみた

[http://www.mext.go.jp/b\\_menu/shingi/gijyutu/gijyutu4/002-2/gijiroku/05112901.htm](http://www.mext.go.jp/b_menu/shingi/gijyutu/gijyutu4/002-2/gijiroku/05112901.htm)

## 研究環境基盤部会 学術情報基盤作業部会 コンピュータ・ネットワークワーキンググループ(第9回)議事内容

日時: 平成17年10月17日(月曜日)16時～18時

場所: 文部科学省F1会議室 (古河ビル6階)

平成12年に100GFLOPS以上の演算性能を持つ計算機を政府調達上のスーパーコンピュータと定義したところであるが、急速な計算機開発技術の進展に伴い、100GFLOPS以上の演算性能を持つ計算機の導入事例が増えたことなどから、米国と協議した結果、本年5月より1.5TFLOPS以上の演算性能を持つ計算機を政府調達上のスーパーコンピュータとすることになった。本資料は1.5をTFLOPS基準としてまとめたものである。



# FLOPSって？

出展: wikipedia

FLOPS (フロップス、Floating point number Operations Per Second) はコンピュータの性能指標の一つ。1秒間に浮動小数点数演算が何回できるかという能力を理論的/実際の(実験的)に表したもののこと。

科学技術計算やシミュレーションを行うスーパーコンピュータ等の性能を表す際に用いられることが多い。

## FLOPSの理論値計算方法

CPUクロック × コア数 × 1コア1クロックあたりの浮動小数点数演算回数

Q: 1.5TFLOPSって今のサーバーでどのくらい！？

HP ProLiant DL980G7で計算すると… **3台でおつりができます！**



1台あたり  $2.26\text{GHz} \times 64\text{core} \times 4 = 578.56 \text{ GFLOPS}$

**3台だと約1.73TFLOPS**

\*あくまでも理論値です。



# 文部科学省もわかっていた

研究環境基盤部会 学術情報基盤作業部会  
コンピュータ・ネットワークワーキンググループ(第9回)議事内容より

スーパーコンピュータの定義は変わっていくと思うが、機種更新を行ったら、大体どれくらいの期間スーパーコンピュータと言えるのか。

2年も経つと世界のトップランクのスーパーコンピュータが入れ替わるくらい性能が向上し続けているため、5、6年もスーパーコンピュータという位置づけを維持することはできない。

メーカーは大体5年かけて新しいスーパーコンピュータを開発するため、かつての大型計算機センターでは5年に1回スーパーコンピュータを更新していた。ただ、メーカーごとに開発のサイクルが2～3年程度ずれているため、スーパーコンピュータと言えるのは2～3年ではないか。

国立大学法人化前の計算機に係るレンタル料の予算は、6年間のレンタル期間を前提に配分されていたが、いまや6年のレンタル契約を結ぶと、レンタル期間中にスーパーコンピュータとしての価値がなくなってしまう。



# ちなみに最新スパコン性能ランキング1位は？

世界で最も高速なコンピュータシステムの上位500位までを定期的にランク付けし評価するプロジェクト



2010年 6月のTOP500ランキングNo1は

Rank	Site	Computer	RMax	RPeak	Processor	Cores
1	Oak Ridge National Laboratory	Cray XT5-HE Opteron Six Core 2.6 GHz	1759000	2331000	AMD x86_64 Opteron Six Core	224162

1759000 GFLOPS

すなわち **1.759 PFLOPS**

ちなみに2009年11月 TOP500 HPシェア1位！  
- HP BladeSystemを利用したシステム = **203件**

- HPシステムでは、209件のランクイン

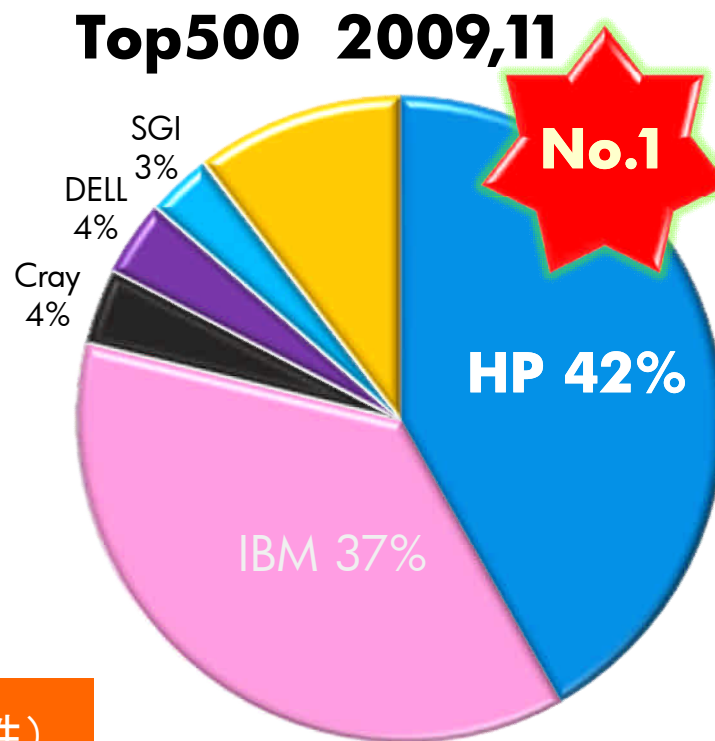
**→42%のシェア No.1**



高いコンピューティング性能



卓越した環境性能(省電力、省スペース性)

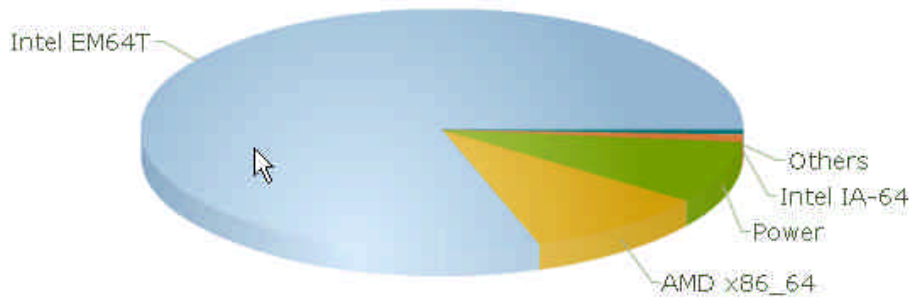
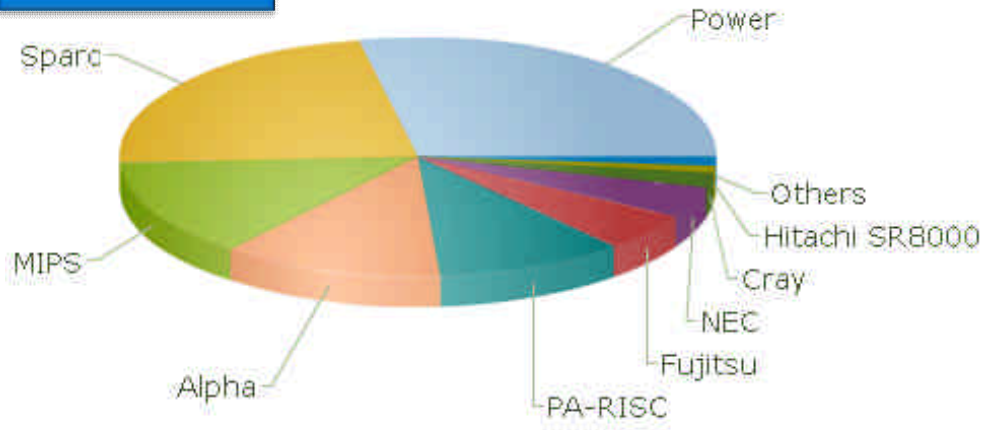


2009年11月 ベンダー別シェア

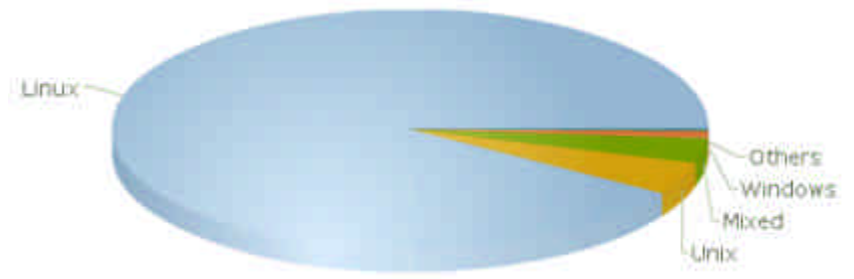
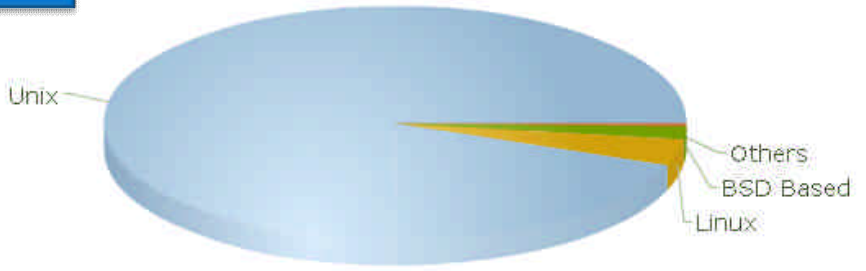


# TOP500の今昔

**プロセッサ** 1999年11月 → 2010年6月



**OS** 1999年11月 → 2010年6月



近年はx86アーキテクチャでOSはLinuxがほとんど



# スパコンのトレンドは

TOP500の推移が象徴するように、スパコンはスパコン専用機から  
いわゆるオープン系へトレンドが移動  
時代の流れとともに



高い専用機を買わなくても、安いサーバーをいっぱい並べて  
並列計算すれば高い性能が出せる！

今やスパコンは**High Performance Computing Cluster**が主流

通称：**PCクラスター**

導入に対する敷居は大きく下がり、様々な分野で活用されている

# スーパーコンピュータの市場は？

## 従来からの“HPC”

### 製造

CAE：衝突解析、流体解析

EDA：半導体設計

### 公共・学術

科学技術計算：汎用計算用途

## “データセンター”

通信：xSP/iDC

Web 2.0

Cloud Computing

Hostingインフラ

通信：R&D

汎用計算/解析

## “金融グリッド”

### 金融

Grid：リスク/デリバティブ分析

Low Latencyシステム

様々な分野でスパコンは  
使われている！！

スパコン市場で共通するキーワードは  
たくさんのノード

CPU性能/大容量メモリ

高速ネットワーク

省電力/省スペース

コストパフォーマンス

安くて!!!

速くて!!!

地球にやさしいスパコンを!

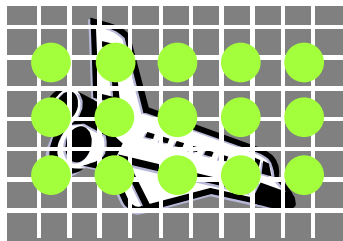
# PCサーバーの使われ方



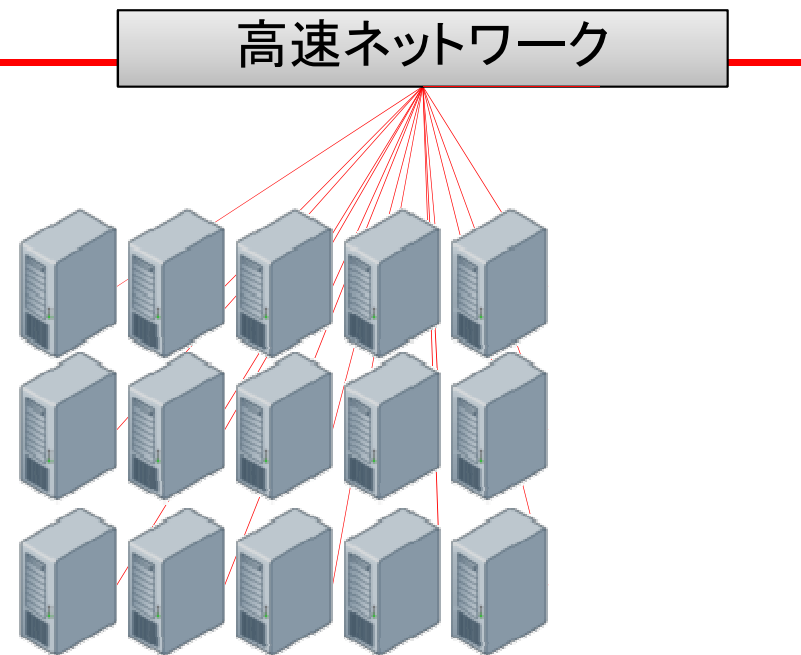
# たくさんのサーバーを使って 並列計算するには？

解析系を例にした場合

①モデル作成 ②メッシュに分割



③並列計算の実行

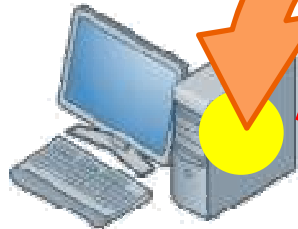
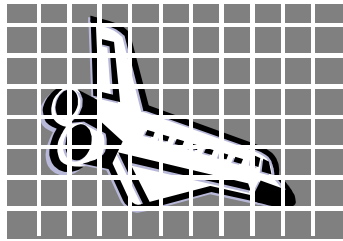




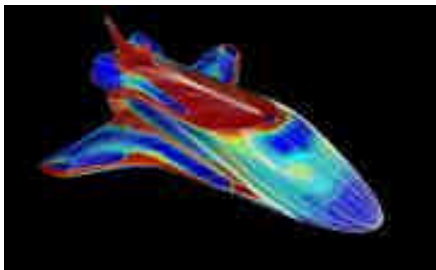
# たくさんのサーバーを使って 並列計算するには？

解析計算を例にした場合

①モデル作成 ②メッシュに分割

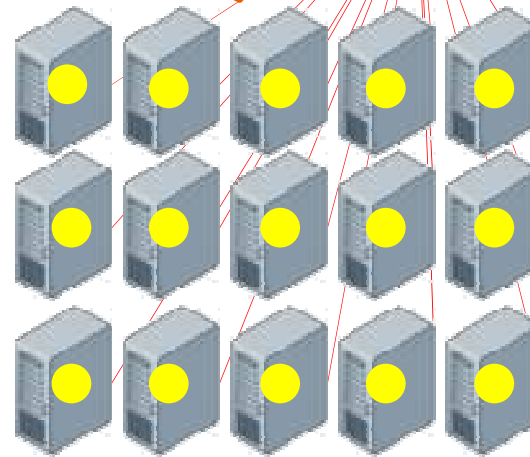


④ポスト処理



③並列計算の実行

高速ネットワーク



メッシュを細かくしたりして、解析精度を上げたい  
→ コンピュータパワーが必要

# HPCクラスタを作る際の検討事項

## ハードウェア選択

CPU性能は？  
メモリの性能・容量は？  
ディスク性能・容量は？  
ネットワークは？

## 環境構築・管理

OSのインストールは？  
環境設定は？  
大量のサーバー管理は？

# クラスターで使われる CPUについて



# X86プロセッサー さらなるマルチコア時代へ

–さらにマルチコア化が進む



	現在	2010年
Intel	DP : Nehalem 4コア MP : Dunnington 6コア	DP : Westmere 6コア MP : Nehalem-EX 8コア
AMD	Istanbul 6コア	Magny-Cours 12コア

–それに伴い、メモリ帯域の拡充も

- Nehalem-EXもNUMAアーキテクチャを採用へ
- AMDもDDR-3メモリを採用

# INTEL “WESTMERE” CPU

## – 業界初の32ナノプロセス採用

- ベースアーキテクチャはNehalemと同一

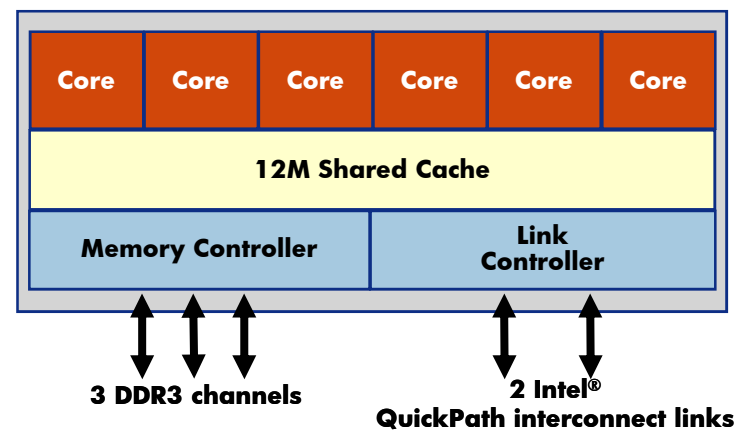
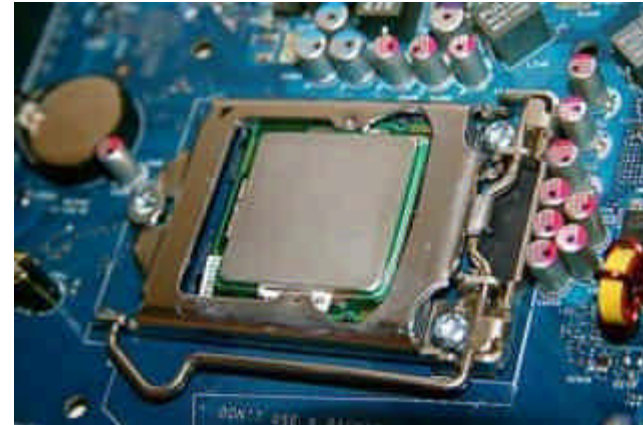
## – CPUの特徴

- 最大6 Core / 12 スレッド
- 12MB shared cache
- CPUあたり 9 メモリスロット。3つのDDR3チャンネル
- インテル ハイパースレッドテクノロジー
- インテル ターボブースト

## – 搭載サーバーの特徴

- 2 CPU / 24 スレッド
- NUMAアーキテクチャ
- 最大18メモリスロット搭載可能

※搭載サーバーモデルにより、仕様は異なる場合があります



# AMD “MAGNY-COURS” CPU

## – 次世代Opteron CPU

- 2way、4way両用

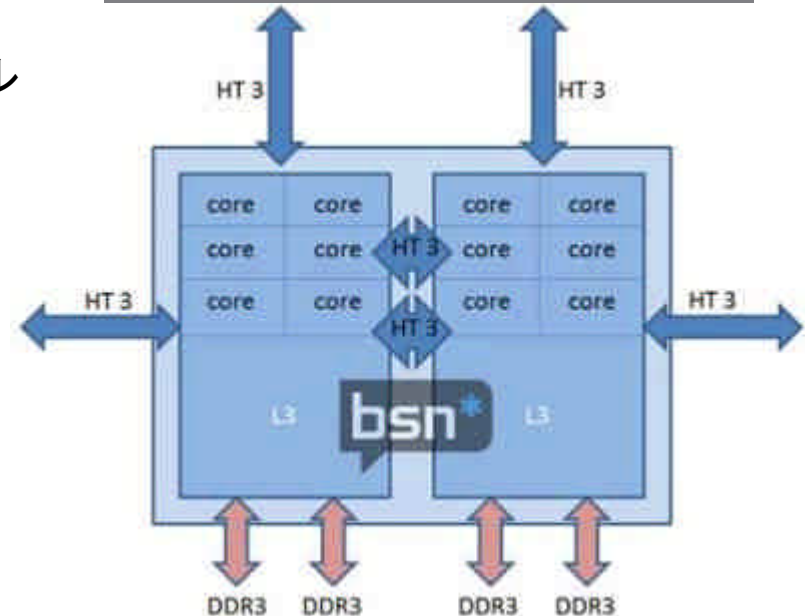
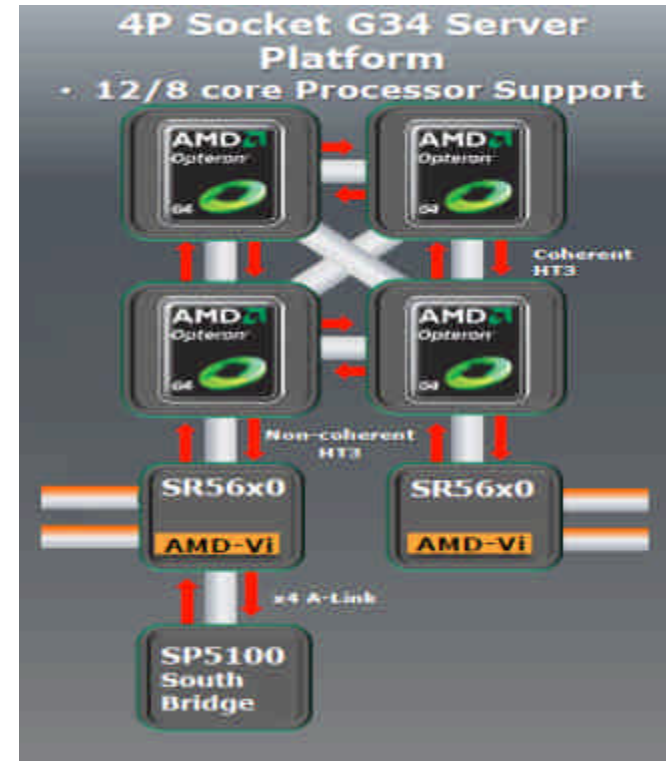
## – CPUの特徴

- 最大12 Core / 12 スレッド
  - 6 Core CPUを2つくっつけた形
- 12MB shared cache
- CPUあたり 12 メモリスロット。4つのDDR3チャンネル

## – 搭載サーバーの特徴

- 2 CPU / 24スレッド。4 CPU / 48スレッド
- NUMAアーキテクチャ
- 2way : 24メモリスロット。4way : 48メモリスロット

※搭載サーバーモデルにより、仕様は異なる場合があります



# クラスターで使われる メモリについて



# 速くて省電力、大容量のDDR3メモリ

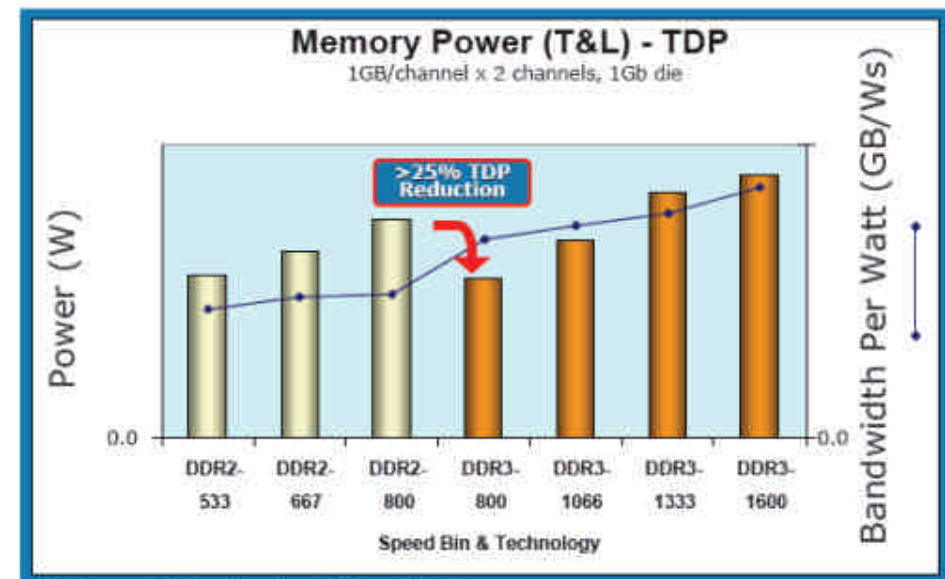
## ➤速い！

### ◆帯域幅が**2倍**で、低レイテンシ

- DDR3:1333MHz vs. DDR2:667MHz

## ➤最大**25%**省電力

- DDR2 1.8Vから、1.5Vへ
- 温度感知センサーがアイドル時の電力を抑制



Data from Intel Developer Forum.



# ハイパフォーマンスの追求 メモリ容量と高速化の両立

## – Nehalemサーバーのメモリ構成ルール

- 最速の1,333MHzで稼働させるには、メモリをチャンネルあたり1枚ずつしかさすことができない
  - チャンネルあたり2枚さした時点で、1,333MHz対応メモリを利用しても、1,066MHzが上限に

## – HPの場合

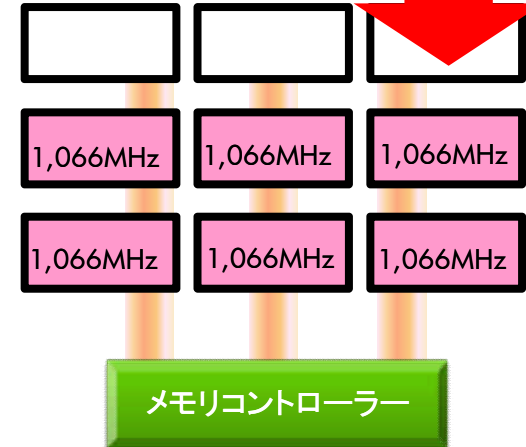
- チャンネルあたり2枚さしても、1,333MHzを実現することが可能
  - 最新BIOSにupdate / BIOSにて機能を有効に

自社販売メモリの綿密なテスト体制と、  
自社開発のBIOSだからこそ本機能を迅速に実現



従来、および他社

性能Down



HPの場合

性能Up



# クラスターで使われる ネットワークについて



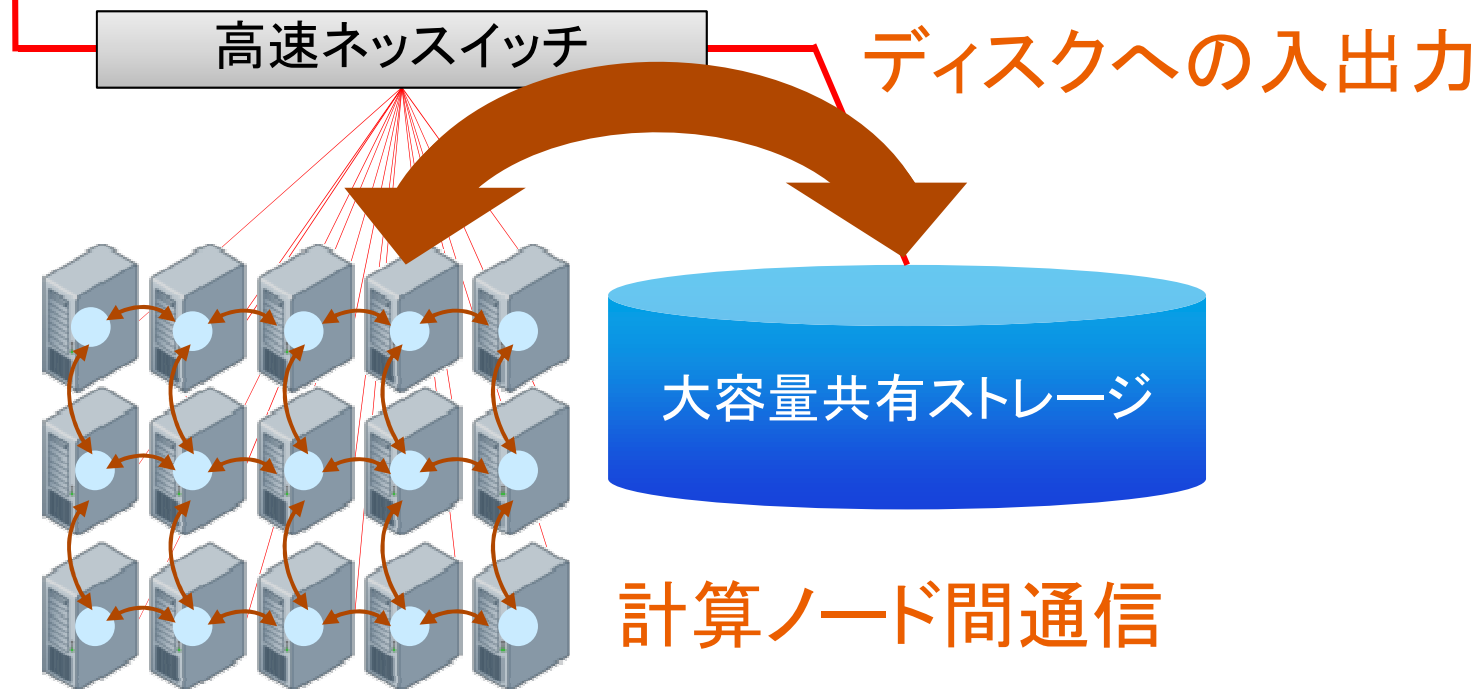
# ネットワークも重要



CPUやメモリが速くても  
通信速度が遅いと、計算時間も遅くなる！

1GbEthernet? 10GbEthernet?

それもあるけど...

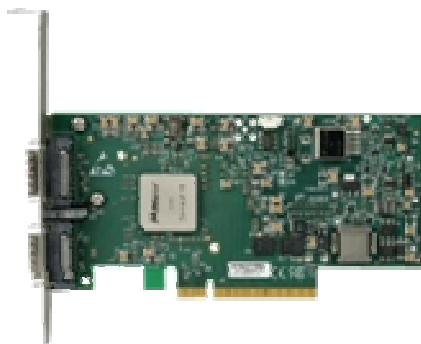


# ハイパフォーマンスネットワークの追求 INFINIBANDはQDRの時代へ

## - Infiniband QDR製品群



c-Class QDR Mezzanine HCA



QDR PCI-Express G2 HCA



c-Class内蔵 QDR Switch



Voltaire OEM 36port QDR Switch

クラスタの全体性能の向上のために必要不可欠  
HP自身が責任を持ってクラスタプラットフォームを認定しています

# スパコン界のトレンド GPUコンピューティングと 高速半導体ストレージ



# アクセラレーターソリューション

- x86システムで標準化しながら、特定用途には局所的な性能を提供

## ▶ GPGPU

NVIDIA Tesla GPU



## ▶ 半導体ストレージ

IOアクセラレーター  
(HP BladeSystem向け)



# GPUとは

## GPU: Graphics Processing Unit

3D グラフィックスの表示に必要な計算処理を行なう半導体チップ。

- ・画像処理を専門とした補助演算装置
- ・単純な処理の反復作業に向いている

IT用語辞典e-Wordsより  
<http://e-words.jp/w/GPU.html>



PCクラスターでやるような計算にも使ってしまうおう！

というわけで

GPGPU (General Purpose GPU) が今  
スパコン界で注目を集めている！

# 何故 GPGPUなのか？

- ・コアが多い！

CPUだとせいぜい8コアとか12コアとか  
GPUなら数百コア！

- ・メモリ帯域幅も広い！

Nvidia Tesla c2070



CPUと比較して何十倍もの演算性能を発揮することが出来て  
かつ、コストパフォーマンスが高い！

当然注意点もある

- ・利用においてはコーディングが重要（GPU用にコーディングする必要あり）  
Cライクな開発環境が用意されている CUDA etc
- ・冷却も考えておかないとダメ



# TOP500の2位に GPU搭載システムがランクイン

Rank	Site	Computer	RMax	RPeak	Processor	Cores
1	Oak Ridge National Laboratory	Cray XT5-HE Opteron Six Core 2.6 GHz	1759000	2331000	AMD x86_64 Opteron Six Core	224162
2	National Supercomputing Centre in Shenzhen (NSCS)	Nebulae - Dawning TC3600 Blade, Intel X5650, NVIDIA Tesla C2050 GPU	1271000	2984300	Intel EM64T Xeon X56xx (Westmere-EP)	120640

Rmaxは1.27PFlops

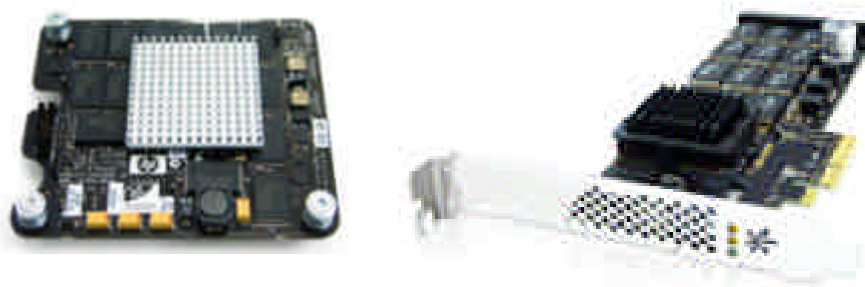
RPeakは2.98PFlops!



# 高速半導体ストレージ I/Oアクセラレーター

## I/O アクセラレーターの特徴

- 不揮発性NAND Flashメモリーを採用
- PCIeスロットへ直接接続し広い転送帯域を確保
- 24+1チャンネルx8メモリーバンク構造によりハードディスク200個以上のパフォーマンスを発揮
- ラックスペースを削減
- 稼働部分の無い設計でハードウェア障害を削減
- 消費電力はハードディスク1台分と同等に削減
- ハイレベルのエラー訂正機能や予備容量の搭載によりNANDフラッシュメモリーの寿命を長期化

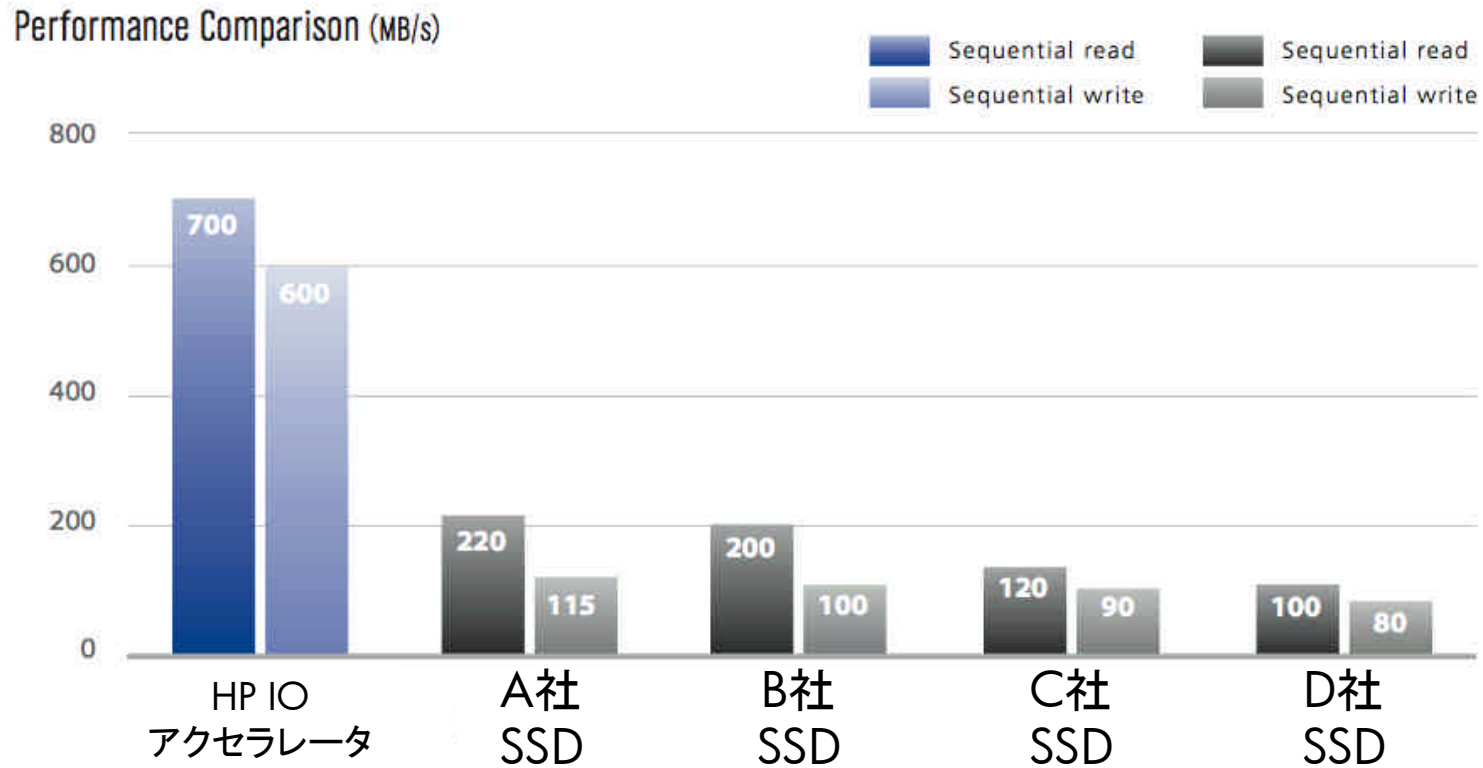


I/O Accelerator

あ、SSDのことでしょ？

**否** SSDとは違います

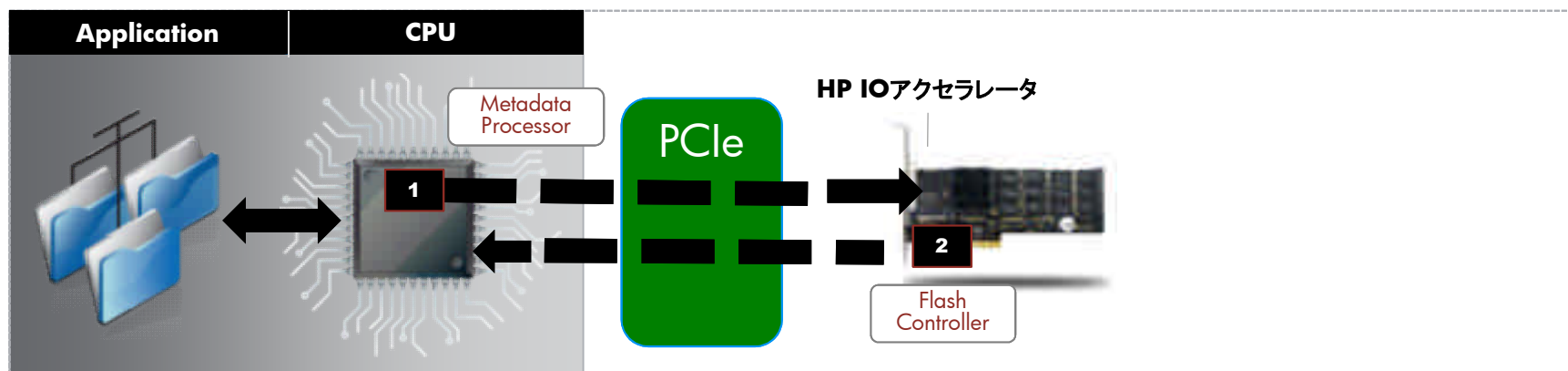
各社SSDとのパフォーマンス比較



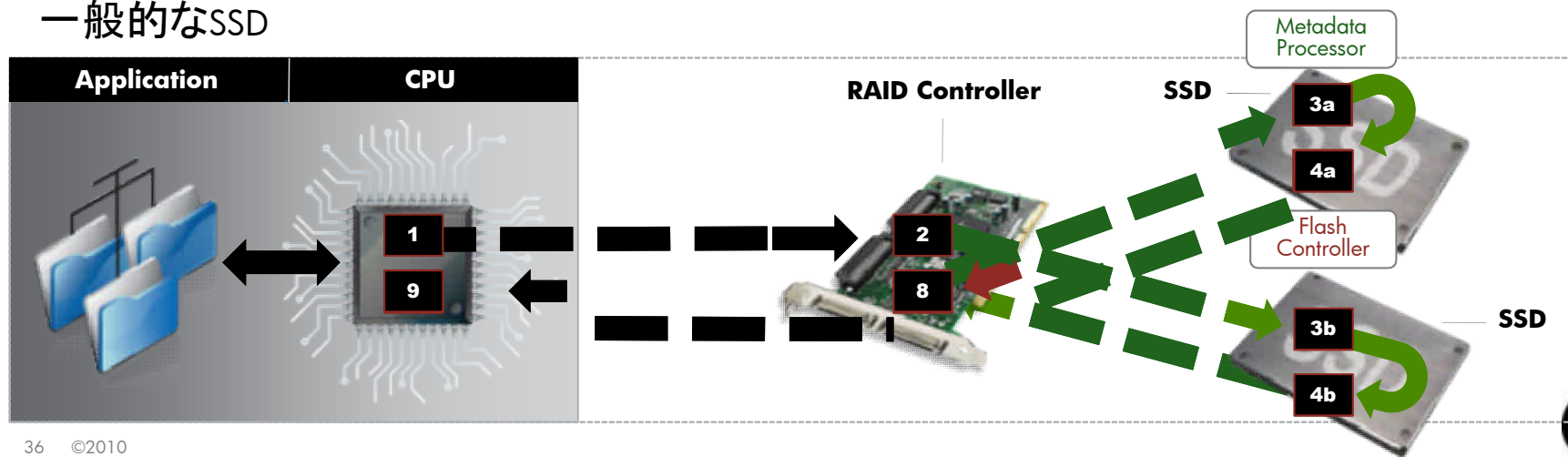
# 速さの秘訣①

## コンテキストスイッチングによる高速化

HP IOアクセラレータ

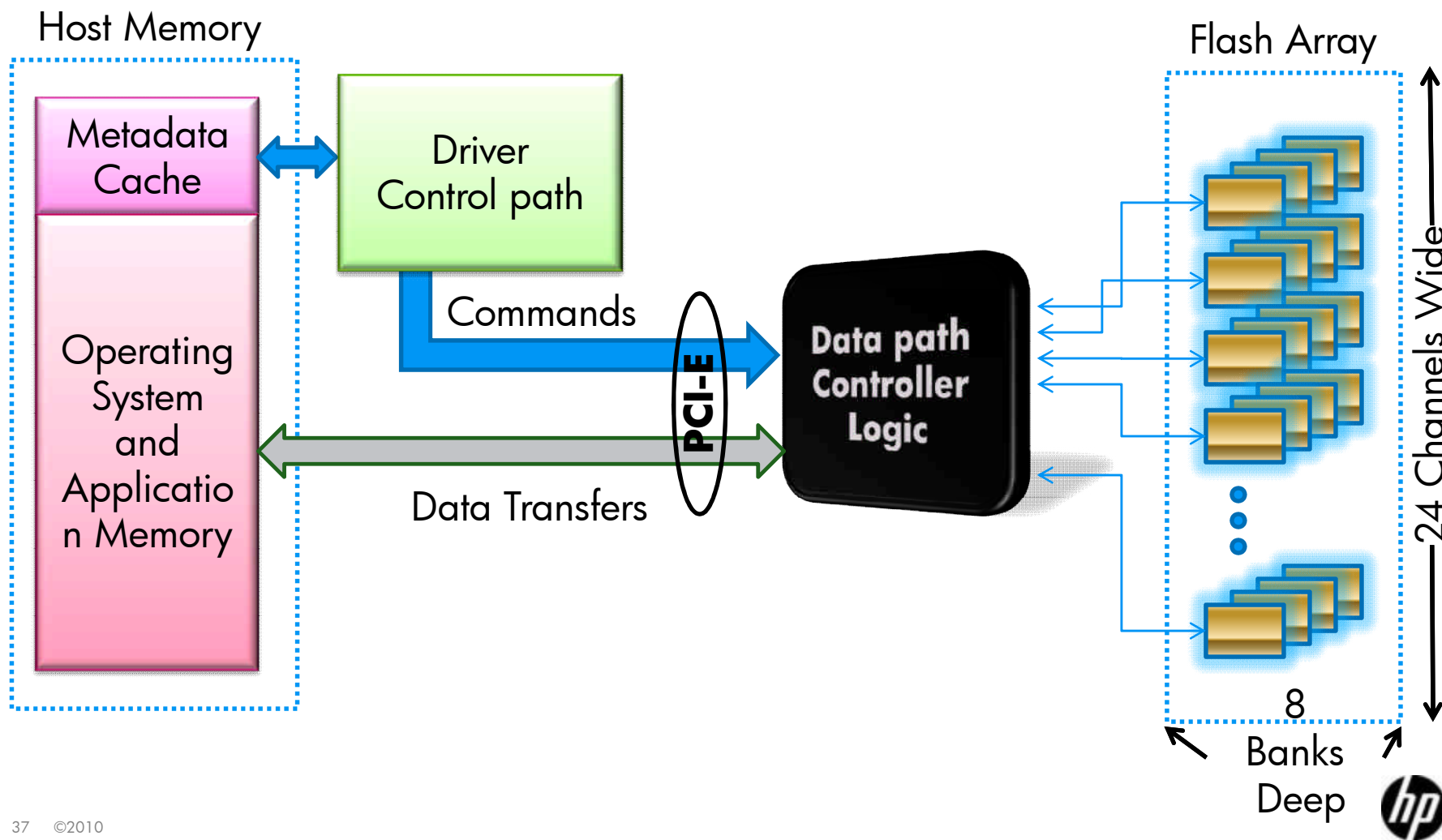


一般的なSSD



# 速さの秘訣②

## IOアクセラレータ アーキテクチャ概要



# スパコンでトレンドと言えば

## TSUBAME2.0

クラウドwatch

東工大、最高性能**2.4**ペタフロップスを実現するグリーンスパコンを開発開始

[http://cloud.watch.impress.co.jp/docs/news/20100617\\_374856.html](http://cloud.watch.impress.co.jp/docs/news/20100617_374856.html)

東京工業大学は6月16日、学術国際情報センターが中心となり、日本電気株式会社(以下、NEC)と米国ヒューレット・パッカード(以下、HP)などの企業連合と合同で、今年11月に日本初のペタコンとして稼働予定のクラウド型グリーンスーパーコンピュータ「TSUBAME2.0」の開発を開始したと発表した。





# TSUBAME2.0の特徴

## 1. 世界一クラスのペタコン: 倍精度2.4ペタフロップス

- 最新型GPU・CPUによるベクトル・スカラー混合アーキテクチャ: 圧倒的計算性能・バンド幅
  - **2.4 Petaflops**, メモリバンド幅**0.72**ペタバイト/s (地球シミュレータの**4.3**倍)
- 世界最高速の200テラビット級のバイセクションバンド幅を実現した光ネットワーク
- 最新のSSDなどを活用した多階層ストレージによる15PBの大規模化と高速化(0.66TB/s)

## 2. 世界一環境「グリーンスパコン」

- TSUBAME1同等のエネルギー消費・30倍の電力性能比・PUE=1.28・「Green500」世界一?
  - GPU+マルチコアCPUの大幅活用による高効率化
  - 最先端の冷却法: 密閉型の水冷ラック, 負荷集中での高熱勾配, 夏季の負荷キャップ
  - JST-CREST Ultra Low Power-HPC等の研究成果の応用
- => **PUE 1.28**以下(他の国内のスパコンセンターは**1.7~2.0**程度)

## 3. 「クラウド型スパコン」: 総合的学内ITホスティング

- Windows HPC/Linuxなど複数OS, 複数環境のサポート
- 仮想化による種々のデータセンターホスティング機能のサポート
- 教育用/Kioskシステムのバックエンド化、全学アカウント・総合的学内ITの集中化・費用削減

## 4. 東工大GSICでの種々の基礎研究・メーカー共同開発の成果

- **JST-CREST “Ultra Low Power HPC”**, 科研特定領域「情報爆発」, 文科省-国立情報研 **NAREGI / e-Science**等、多くの基礎研究
- 海外企業とも: **NVIDIA CUDA CoE** (日本初), **Microsoft TCI** 包括共同研究契約
- **NEC, HP, NVIDIA, Microsoft, Voltaire, DDN** 等との共同開発体制



# 東工大TSUBAME2.0 2010年11月稼働予定



 東京工業大学  
Tokyo Institute of Technology

日本全土のスパコン全て



我国のスパコン(2010年合算1ペタフロップス以下)



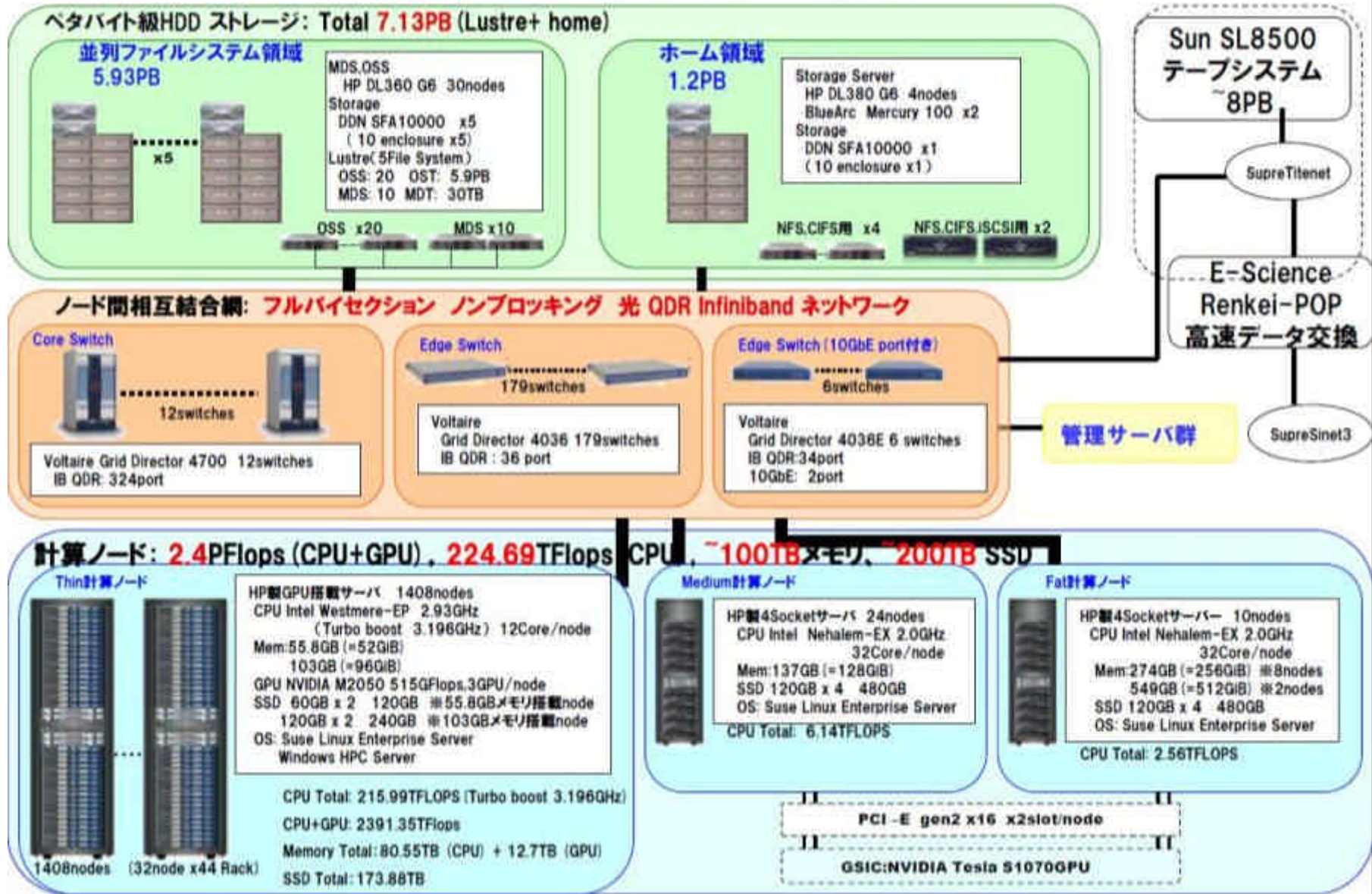
2.4 ペタフロップス (単精度4.8ペタ)  
0.72 ペタバイト/秒メモリバンド幅  
世界トップ超高速光結合網-200テラビット/秒  
総合15ペタバイト階層ストレージ  
60ラック (200m<sup>2</sup>, TSUBAME1以下)  
グリーン (1MW程度、PUE=1.28高効率水冷)  
クラウド (WindowsHPC+Linux+VM)

2010年合算 1ペタフロップス以下  
全国60か所  
年額300 億円以上

東京工業大学学術国際情報センター  
<http://www.gsic.titech.ac.jp/tsubame2>  
配布資料(2)『TSUBAME2.0の概要』(東工大)



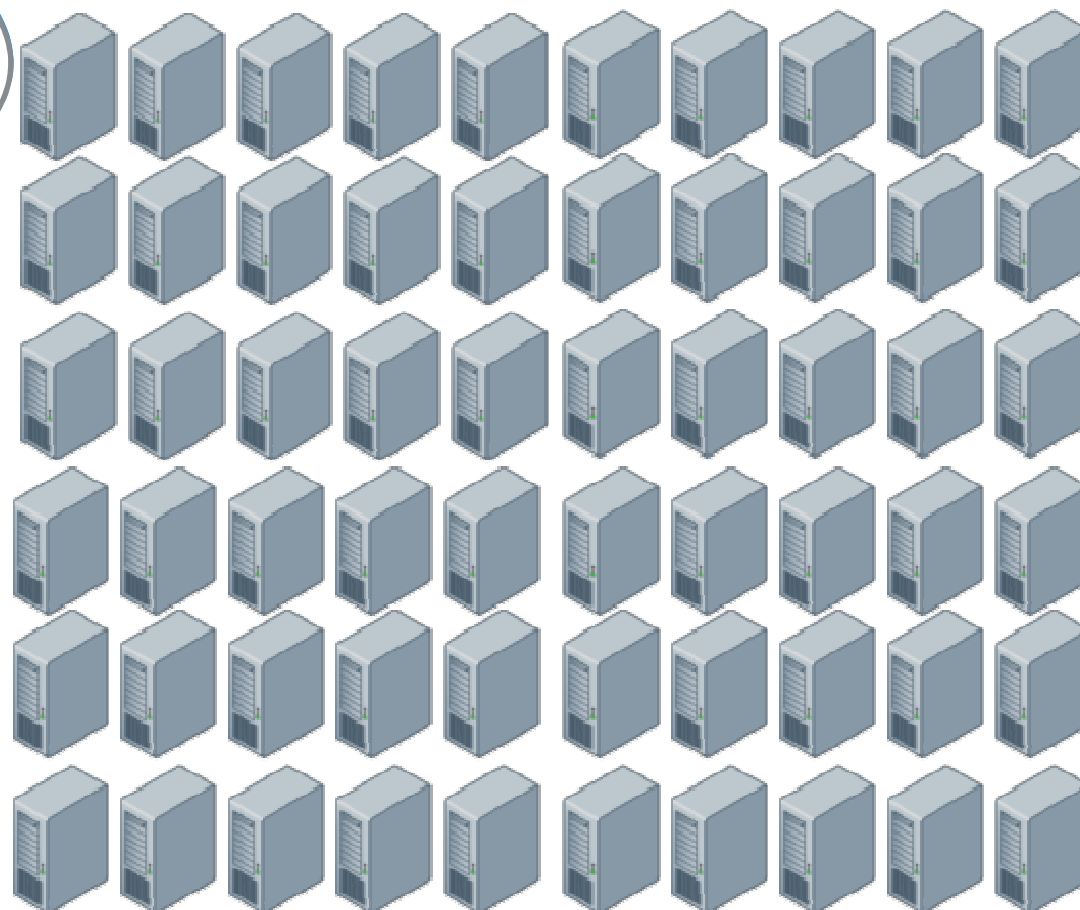
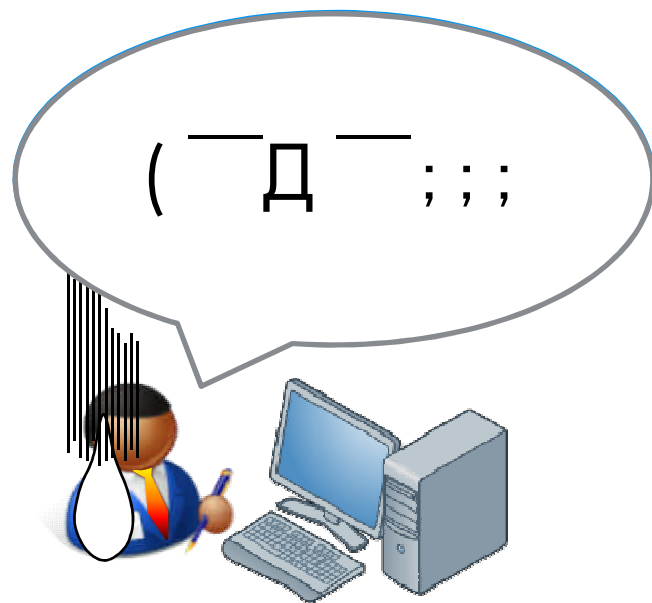
# TSUBAME2.0 システム概念図



# クラスターの管理方法



# 大量のサーバーを構築・管理するのは大変



# クラスター一元管理ツール

## CMU (CLUSTER MANAGEMENT UTILITY)

多くのHPC環境で導入実績

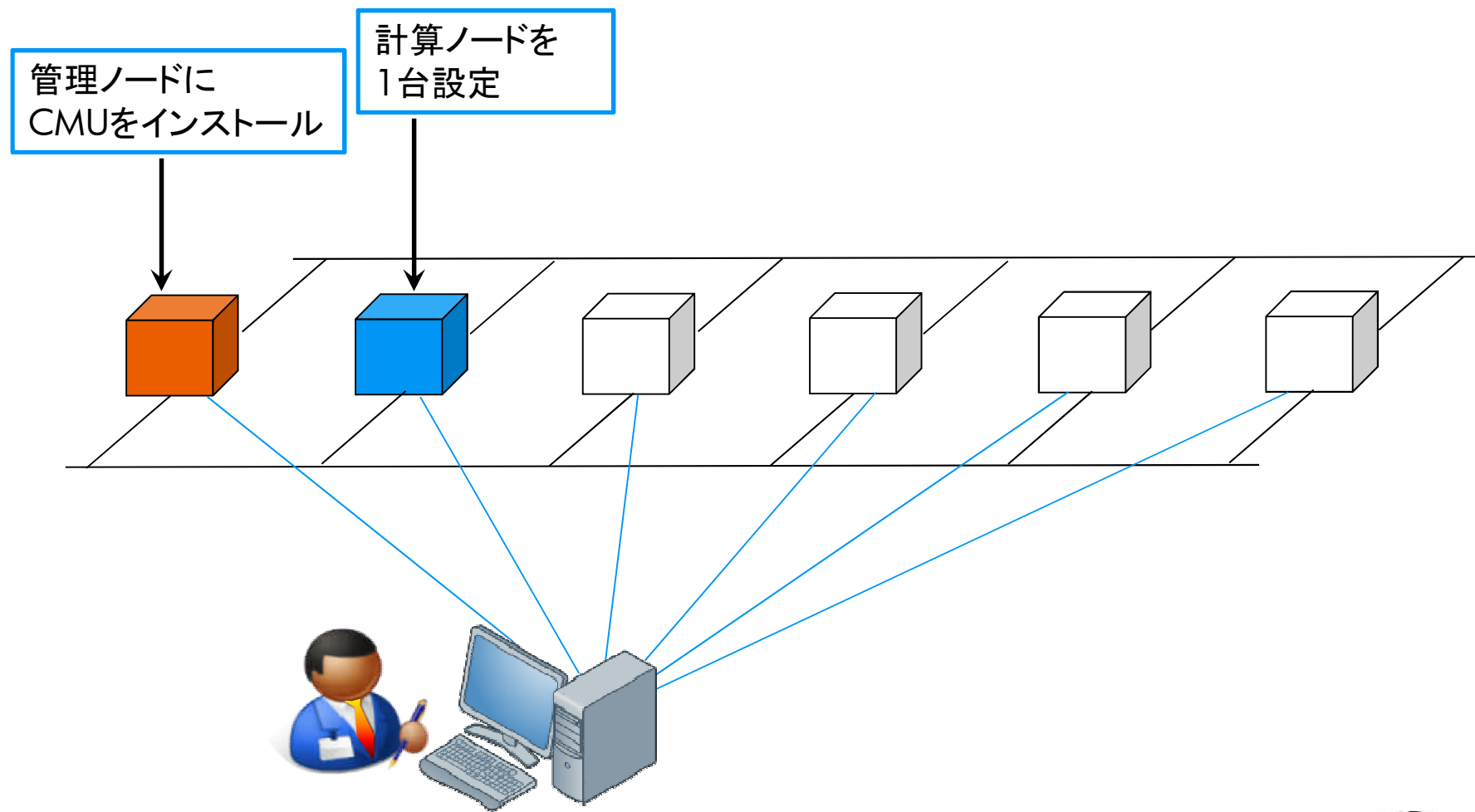
Linux環境一元管理ツールCMU

- 容易に一元管理
  - 一元リソース管理
  - 同時コマンド投入
  - プロビジョニング
- 多くの実績
  - 解析用Linuxクラスタで  
ほぼデフォルトで利用
- ハードと一括したサポート



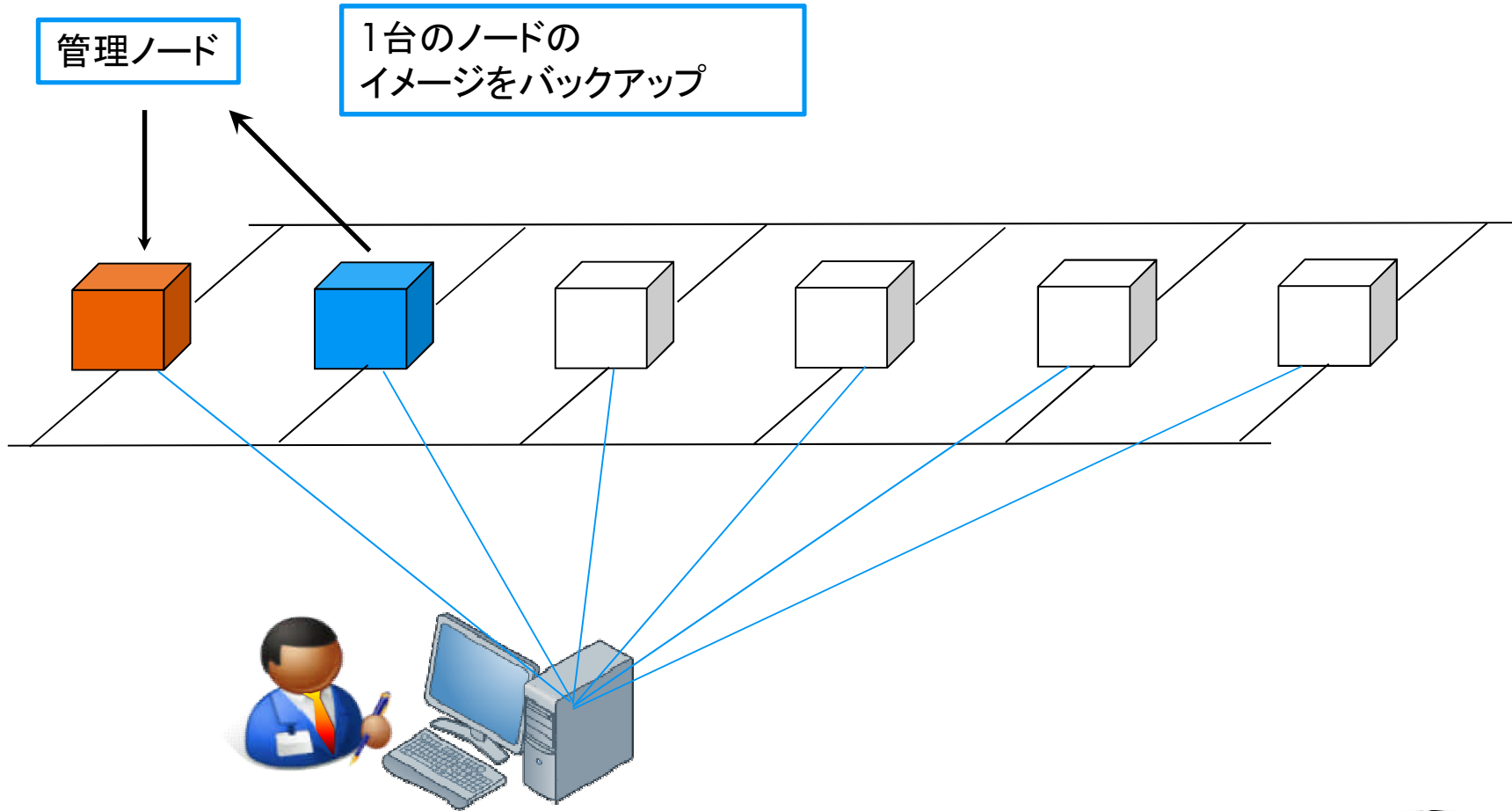
# CMUによるクラスタ環境セットアップ

CMU Cloning



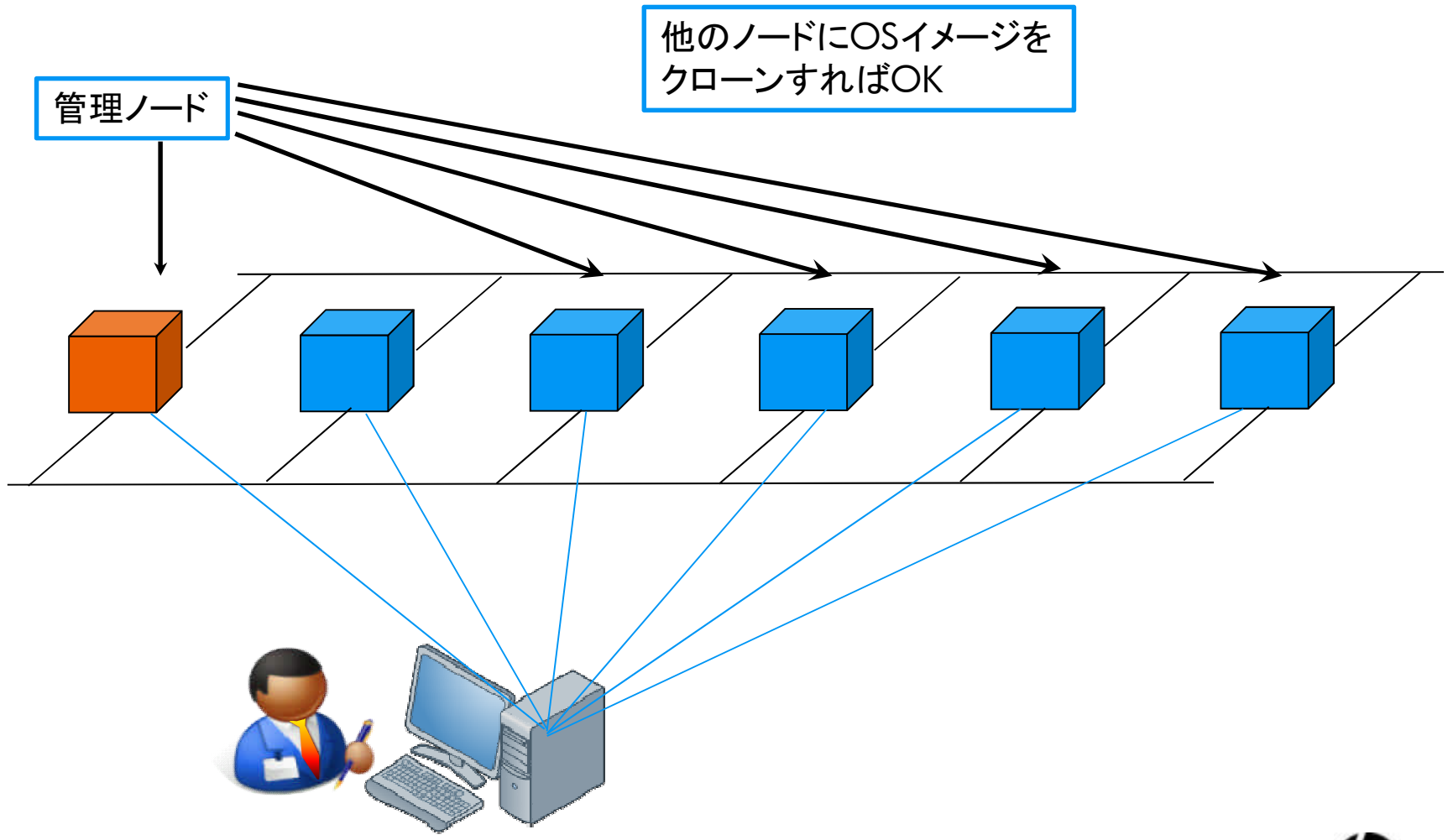
# CMUによるクラスタ環境セットアップ

CMU Cloning



# CMUによるクラスタ環境セットアップ

CMU Cloning



# クローンイメージ配布方式

- 2段階のツリーアルゴリズム

1. 管理ノード



セカンダリサーバ(計算ノード)

2. セカンダリサーバ



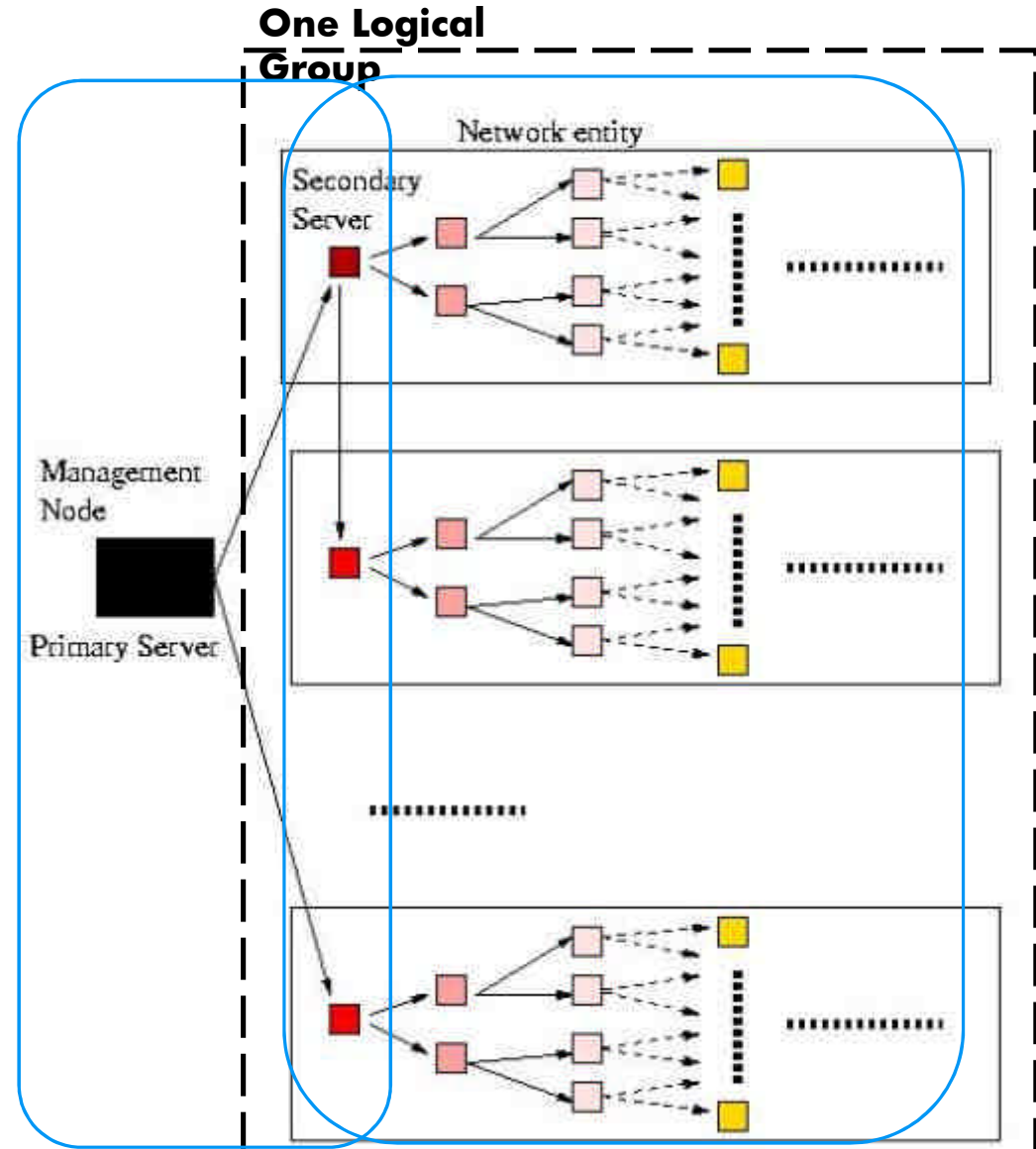
ネットワークエンティティ内の  
計算ノード

利点:

高速なイメージ配布

管理ノードの負荷の低減

スイッチ間トラフィックの低減





# CMUモニタ

The screenshot displays the CMU V4.1.9 monitoring interface. The main window shows a 'GROUP OVERVIEW' with two circular gauges: 'CPU load' (101.7) and 'CPU frequency' (2499.0 MHz). A context menu is open over the CPU load gauge, listing various monitoring sensors. A 'Logical Group' sidebar on the left shows a tree of nodes with their status icons. A 'Node state' legend is at the bottom left, and an 'Alert message table' is at the bottom right.

グループサマリ

ノードの状態

CPU load

CPU frequency

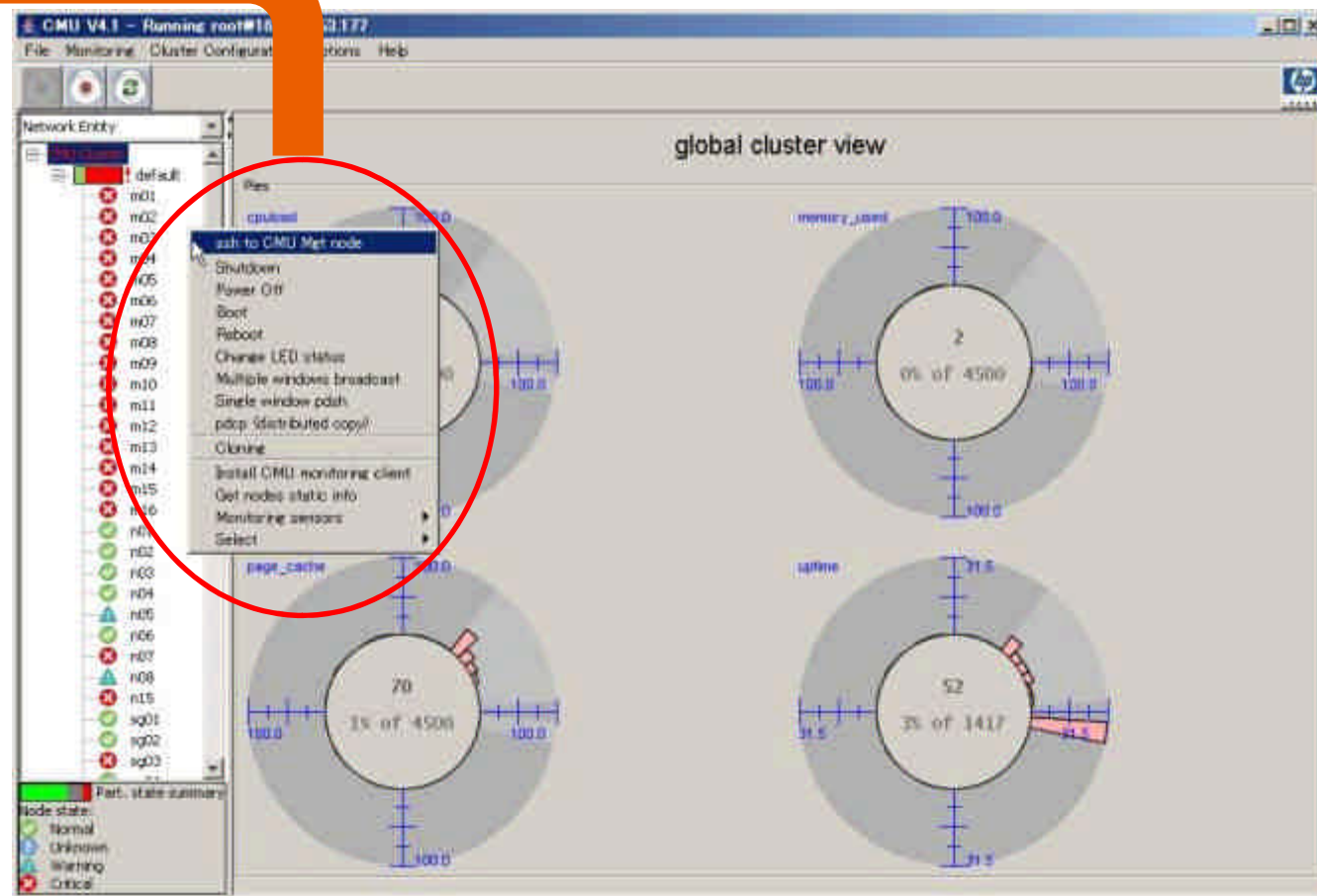
見たい項目を選択

Alert message table



# CMUリモートコマンド

- ssh to CMU Mgt node
- Shutdown
- Power Off
- Boot
- Reboot
- Change LED status
- Multiple windows broadcast
- Single window pdsh
- pdcp (distributed copy)
- Cloning
- Install CMU monitoring client
- Get nodes static info
- Monitoring sensors ▶
- Select ▶



シャットダウンや  
ブート、リブート  
sshのコネクション  
LEDのオンオフなど  
リモートからコマンドを  
実行可能

# ノード管理

## 計算ノードの登録情報の管理

The screenshot displays the CMU V4.1 Node Management interface. On the left, a tree view shows the cluster structure with nodes n01 through n15 and various storage and management nodes. The main area contains a table of node details.

Name	IP Address	Network	MAC Address	Lg Group	Mgt IP	Mgt Card
n01	172.31.1.1	255.255.192.0	00-18-78-CD-8C-4C	b465	172.31.17.1	ilo/ilo2
n02	172.31.1.2	255.255.192.0	00-1C-C4-13-0E-EE	b465	172.31.17.2	ilo/ilo2
n03	172.31.1.3	255.255.192.0	00-1C-C4-13-DE-46	b465	172.31.17.3	ilo/ilo2
n04	172.31.1.4	255.255.192.0	00-1C-C4-13-5E-7E	b465	172.31.17.4	ilo/ilo2
n05	172.31.1.5	255.255.192.0	00-1C-C4-8D-9C-86	b465	172.31.17.5	ilo/ilo2
n06	172.31.1.6	255.255.192.0	00-17-64-77-00-34	b465	172.31.17.6	ilo/ilo2
n07	172.31.1.7	255.255.192.0	00-1C-C4-C2-8D-06	b465	172.31.17.7	ilo/ilo2
n08	172.31.1.8	255.255.192.0	00-18-78-E4-32-84	b465	172.31.17.8	ilo/ilo2
n09	172.31.1.9	255.255.192.0	00-1C-C4-C2-E0-B4	b465	172.31.17.9	ilo/ilo2
m10	172.31.1.10	255.255.192.0	00-1E-0B-1C-C3-3C	b465	172.31.17.10	ilo/ilo2
m11	172.31.1.11	255.255.192.0	00-1E-0B-1C-09-7A	b465	172.31.17.11	ilo/ilo2
m12	172.31.1.12	255.255.192.0	00-1C-C4-C3-6D-06	b465	172.31.17.12	ilo/ilo2
m13	172.31.1.13	255.255.192.0	00-1E-0B-1C-D3-E4	b465	172.31.17.13	ilo/ilo2
m14	172.31.1.14	255.255.192.0	00-19-8B-24-DE-3E	b465	172.31.17.14	ilo/ilo2
m15	172.31.1.15	255.255.192.0	00-1E-0B-1C-9C-7A	b465	172.31.17.15	ilo/ilo2
m16	172.31.1.16	255.255.192.0	00-18-78-E2-82-EC	b465	172.31.17.16	ilo/ilo2
n01	172.31.0.1	255.255.192.0	00-25-83-A3-E3-88	rhe53_72GBx2_RAID1	172.31.16.1	ilo/ilo2
n02	172.31.0.2	255.255.192.0	00-25-83-A3-E2-48	rhe53_72GBx2_RAID1	172.31.16.2	ilo/ilo2
n03	172.31.0.3	255.255.192.0	00-24-81-AD-3D-A0	rhe53_72GBx2_RAID1	172.31.16.3	ilo/ilo2
n04	172.31.0.4	255.255.192.0	00-23-7D-F1-82-08	rhe53_72GBx2_RAID1	172.31.16.4	ilo/ilo2
n05	172.31.0.5	255.255.192.0	00-24-81-AD-84-70	rhe53_72GBx2_RAID1	172.31.16.5	ilo/ilo2
n06	172.31.0.6	255.255.192.0	00-24-81-AC-3D-28	rhe53_72GBx2_RAID1	172.31.16.6	ilo/ilo2
n07	172.31.0.7	255.255.192.0	00-24-81-AD-ED-08	rhe53_72GBx2_RAID1	172.31.16.7	ilo/ilo2
n08	172.31.0.8	255.255.192.0	00-24-81-AD-39-E8	rhe53_72GBx2_RAID1	172.31.16.8	ilo/ilo2
n15	172.31.0.15	255.255.192.0	00-1E-0B-ED-61-A0	rhe53_72GBx2_RAID1	172.31.16.14	ilo/ilo2
sq01	172.31.0.51	255.255.192.0	00-12-79-C6-F5-05	segment-sfs20	172.31.16.51	ilo/ilo2
sq02	172.31.0.52	255.255.192.0	00-12-79-C6-11-85	segment-sfs20	172.31.16.52	ilo/ilo2
sq03	172.31.0.53	255.255.192.0	00-12-79-C7-87-DF	segment-sfs20	172.31.16.53	ilo/ilo2
sq04	172.31.0.54	255.255.192.0	00-12-79-C9-8A-5F	segment-sfs20	172.31.16.54	ilo/ilo2



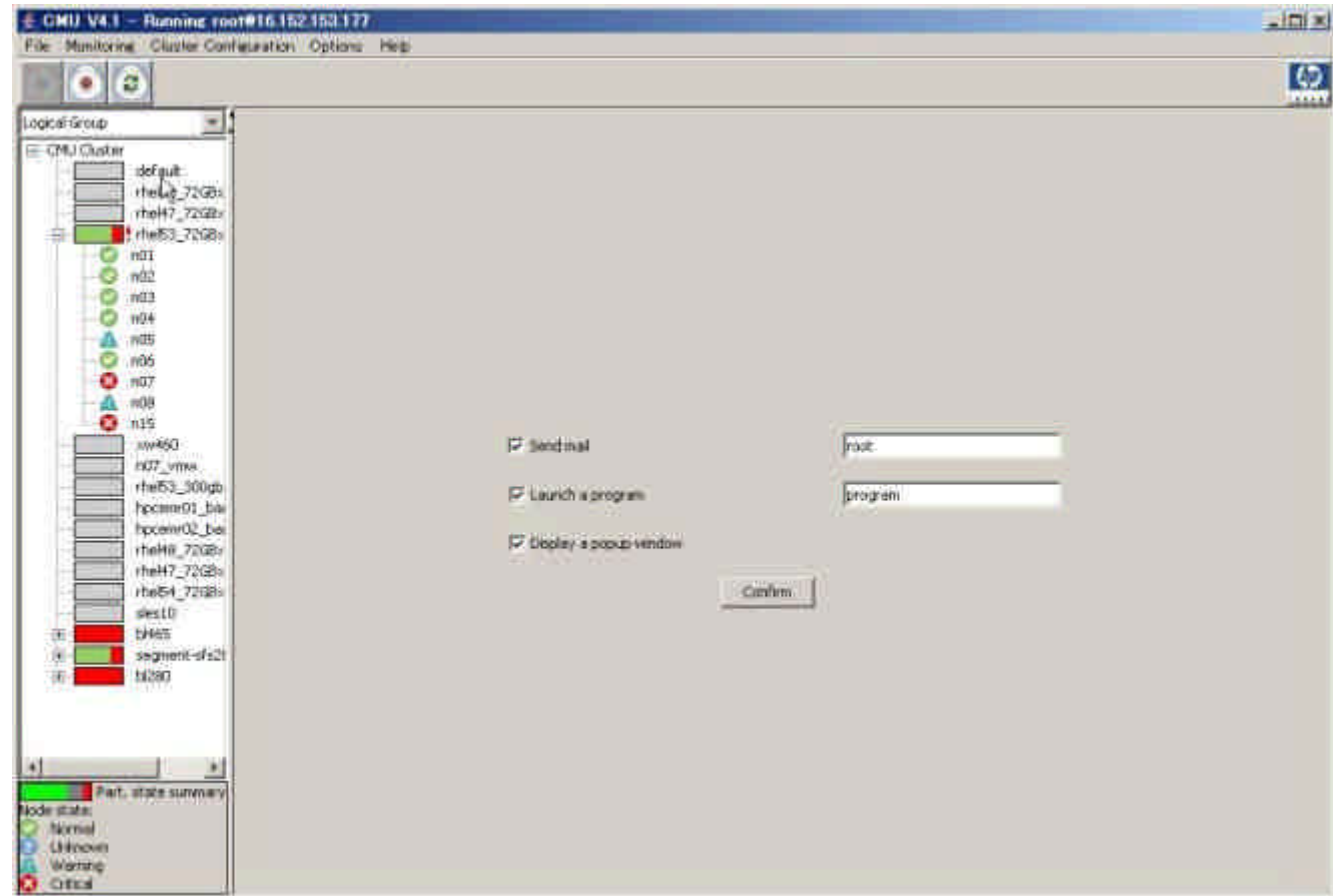
# イベント処理設定

管理ノードから計算ノードへのネットワーク状態を監視し  
状態遷移(切断、接続再開)をトリガとして処理を実行

メール送信

ユーザスクリプトの実行

ポップアップウィンドウ表示



# コンソールブロードキャスト

– 選択した複数ノードにコマンドを一括入力できる

- 接続経路: ssh(telnet)または管理プロセッサ経由

マルチウィンドウ形式

- キー入力や結果を簡単に確認でき、エディタも使用可能



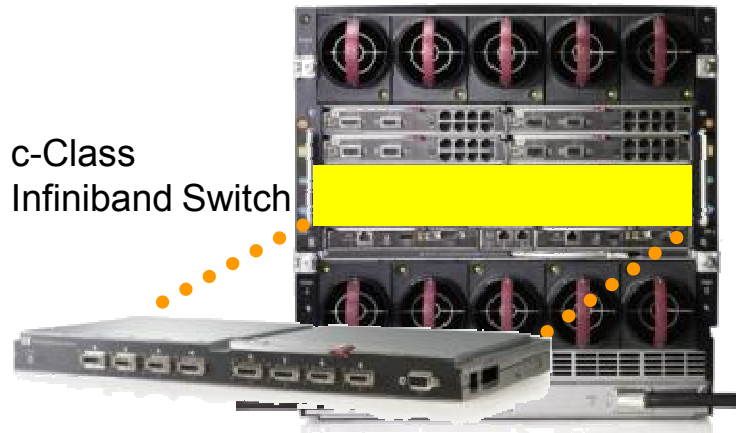
# クラスターに最適な サーバーについて



# HP BLADESYSTEM FOR HPC



様々な用途向けに選べるサーバラインナップ



c-Class  
Infiniband Switch

Infiniband Switchや10Gb Switchも  
エンクロージャに収納



c-Class Enclosure

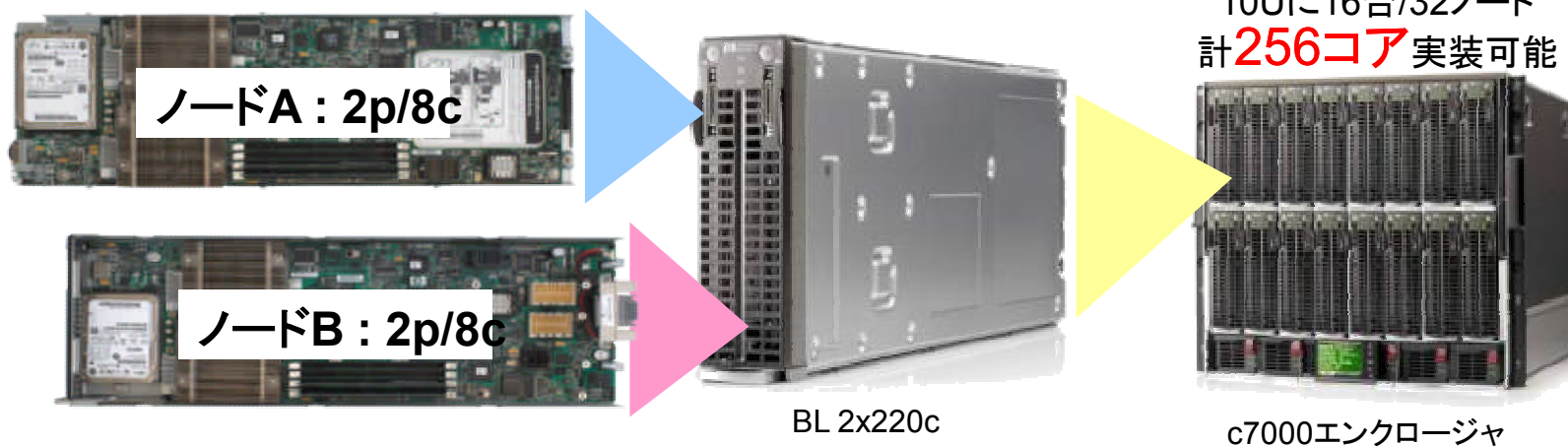
小規模から大規模まで、幅広いシステム規模に対応



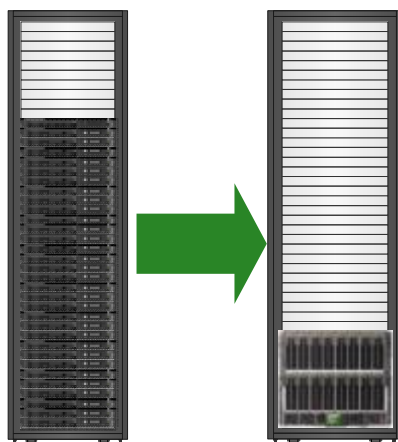
# 超高密度「2 IN 1 ブレード」 HP BLADESYSTEM BL2X220C

コンピューティング性能と環境性能を併せ持つ

HPCプラットフォームに最適なブレードサーバ



同CPUの1Uラックマウントサーバ32台と比較して



- 60% 省スペース
- 52% 省電力
- 64% 軽量

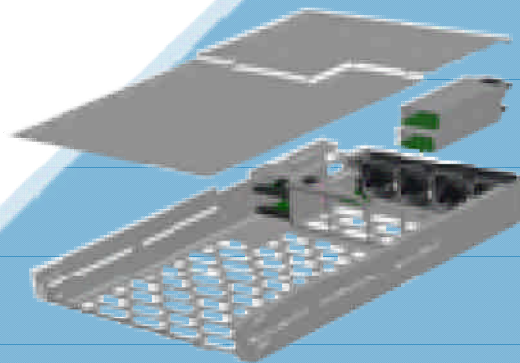


# HP PROLIANT 製品ラインナップ

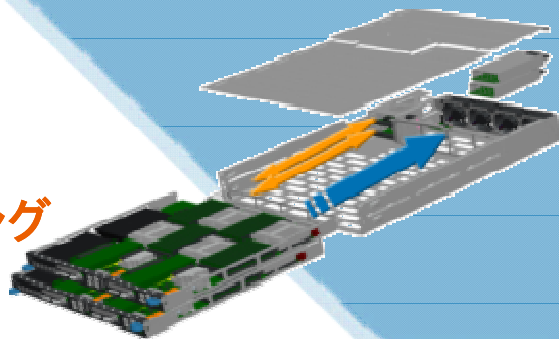
## SL : SCALABLE LINE



前面I/Oケーブルリング



2U シャーシ  
電源効率と冷却効率を追求  
ファンとパワーサプライ(高効率)の共有

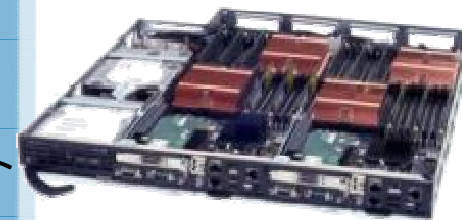


**SL160z G6**  
メモリ、IO重視  
18 DIMM & 2 PCIeスロット

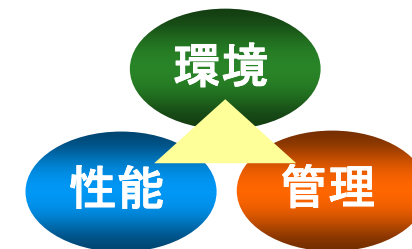
**SL170z G6**  
ディスク容量重視  
6 x 3.5inch ハードディスク

**SL2x170z G6**  
実装密度重視  
1トレイに2ノード

**SL165z G6**  
AMDモデル  
6コアCPU



# スケーラブルサーバー : SL2X170Z G6 SPECPOWERでNO.1を達成



HP ProLiant SL2x170h G6

- 2Uに独立した4ノードを収納可能
- ファンと電源は全ノードで共有し、効率化
- 前面着脱式/ケーブルリングによる抜群の管理性
- 最速CPUも搭載可。高性能HPCシステムにも

## 主なベンダの結果

ベンダ	機種	Result (/watt)
<b>HP</b>	<b>ProLiant SL2x170z G6</b>	<b>2,316</b>
IBM	iDataPlex dx360 M2	<b>2,231</b>
Dell	Power Edge R610	<b>1,930</b>
Fujitsu	PRIMERGY TX100 S1	<b>1,500</b>

<[http://www.spec.org/power\\_ssj2008/results/](http://www.spec.org/power_ssj2008/results/)>

クラスタなどのスケールアウトシステムに  
SL2x170z G6が最適であることを証明!!

# 消費電力の話



サーバーを大量に使うとなると・・・

電気を大量に消費する

大量に熱が発生する

冷やすためにはエアコンがいる

さらに電気を消費する

地球に優しくないですね。



# サーバーコンポーネントから、データセンター全体まで HPの省電力ソリューション

サーバシステムだけではなく、データセンター環境全体の最適化とコスト削減を実現

## データセンター

データセンター環境のリアルタイム視覚化  
データセンター運用の最適化  
エネルギー削減



データセンターレベル

サーバーシステム  
レベル

## サーバーシステム

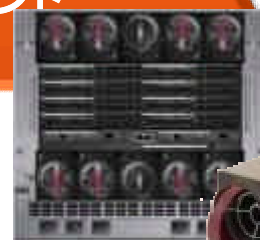
- 大量の温度センサ:**Sea of Sensors**
- 必要な場所だけ冷やすスマート冷却ファン
- 電力上限を制御するパワーキャッピング機能

## サーバーコンポーネント

- 省電力CPUの採用:**Xeon Nehalem**
- 速くて省電力、大容量:**DDR-3** メモリ
- 80 PLUS GOLD取得の共通パワーサプライ



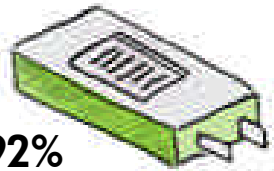
サーバーコンポーネント



# 無駄な電力を使わない工夫 高効率パワーサプライを全機種共通で

**80 PLUS GOLD ...**  
電力効率の最高金賞

460W AC  
最大変換効率92%



750W AC  
最大変換効率92%



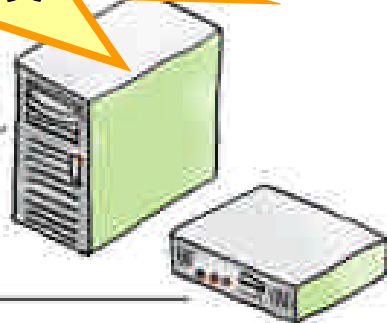
1200W AC  
最大変換効率90%



48Vdc 1200W  
最大変換効率90%



共通の  
スロット



**G6**



HP Power Advisor

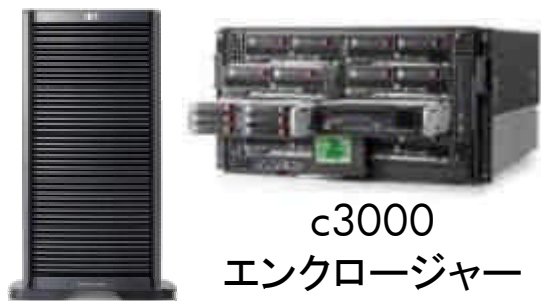
失われた電力は熱になる → 電力と冷却、両方に直結

# 最新の高効率パワーサプライが 幅広いラインナップで利用可能

サーバー、ストレージなどの小中型機を中心に  
電源モジュールスロットを標準化



DLシリーズ



MLシリーズ

c3000  
エンクロージャー

最大460W出力 (AC100~240V): 変換効率 最大92%



最大750W出力 (AC100~240V): 変換効率 最大92%



最大1200W出力 (AC100~240V): 変換効率 最大90%



最大1200W出力 (DC -48V): 変換効率 最大90%



電力効率向上

開発点数の絞込みにより、投資を集中化！  
変換効率**90%以上**のパワーサプライで省電力に

# 電力効率の「最高金賞」パワーサプライ HP ProLiant の採用機種



DL385 G6 : 標準搭載



DL370 G6 : 標準搭載



DL380 G6 : 標準搭載



DL1000 : 標準搭載



c7000エンクロージャ  
標準搭載



DL360 G6 : 標準搭載



SL6000 : 標準搭載



DL320 G6 : オプション (CTO)



ML370 G6  
標準搭載



ML350 G6  
標準搭載



ML330 G6  
オプション



ML150 G6  
オプション



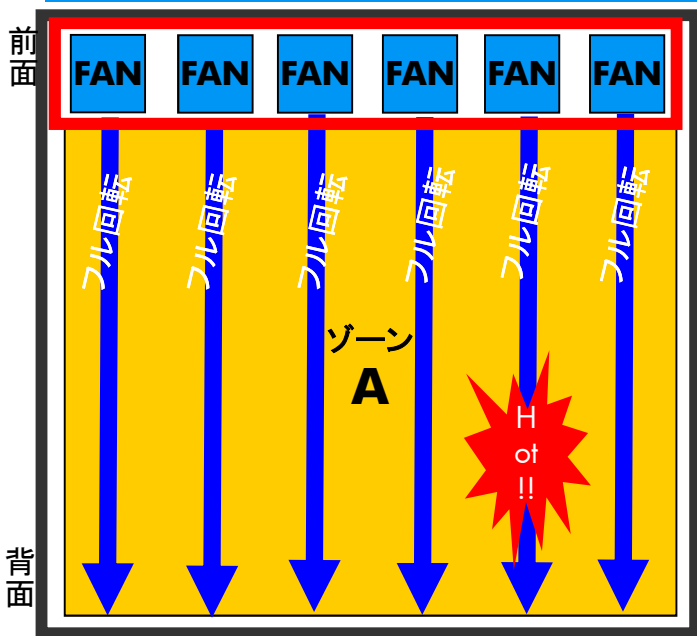
DL160 G6 : オプション





# 無駄な電力を使わない工夫 スマートな冷却

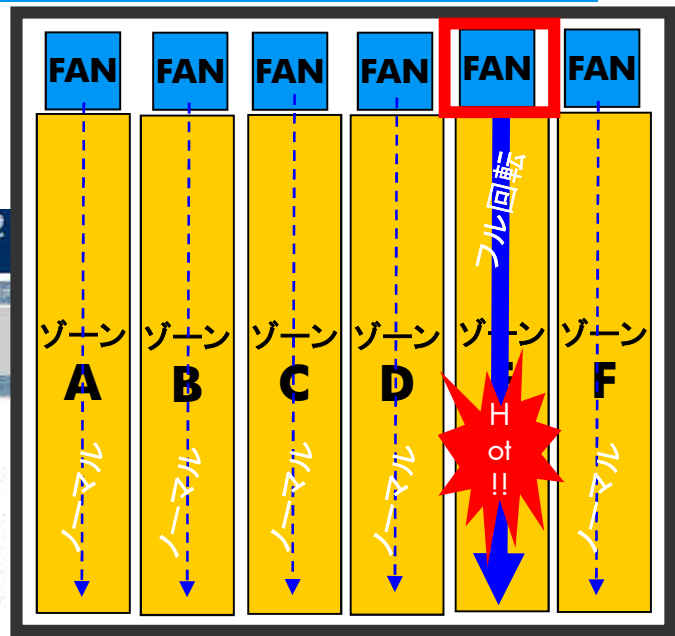
**必要な場所だけ冷やす。**  
HP BladeSystem c-Classエンクロージャの技術をML/DLにも応用



従来製品

**G6**

例: DL380



**【こんなに違う!!】**  
回転速度によるファンの電力消費

ノーマル時: **0.8W**程度

フル回転時: **6-12W**

**10倍以上の差**

**15台で最大840W、  
30台で最大約1600Wにも  
なる電力削減効果!!**

# データセンター環境は、はたして健全か??

「冷やしすぎ」と「熱溜まり」??



冷却の無駄と高温のリスク...

現状はあるべき姿なのか??



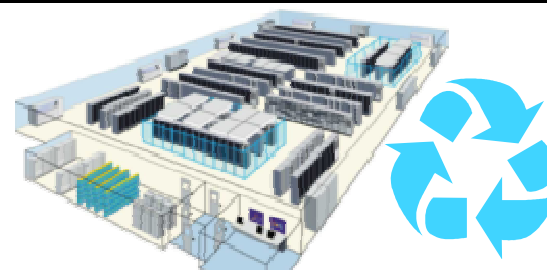
目指すべき指針はなんとなくあるが...

ファシリティ機器は正常稼動??



電源、冷却装置の異常動作は多大な影響に

増強、投資の正しい判断??



既存環境をもっと有効活用できるのでは?

高まるエネルギー、コスト削減要求  
データセンター環境の最適化のためには、情報の可視化が不可欠

# HP Data Center Environmental Edge

データセンター環境をリアルタイムに可視化し、エネルギー使用の効率化とコスト削減を支援

## – ワイヤレスセンサーによる情報の収集と、集計、および可視化を実現

- フロアの温度、湿度、気圧の情報収集
- ワイヤレスのため、手軽に導入可能



## – エネルギー使用状況をリアルタイムに確認し、適正なデータセンター環境の実現を支援

- 熱溜まり、冷やし過ぎの回避
- 安全で効率的な電力提供の実現
- 無駄を排除し、コストを適正化



# Insight Environmental Observerソフトウェア

## 表示項目

- ラックレイアウト
- 温度、湿度、気圧の分布
- 日時での平均/最大温度、湿度、気圧

すべてのセンサーの一覧表示

過去の環境情報のリプレイ

ホットスポット発生領域

気圧

湿度

ラック温度  
上、中、下の  
高さ毎に表示



最後に



# AVATAR – HP PROLIANT

- 10,000 jobs and an estimated 1.3 to 1.4 million tasks per day
- More than 4,000 HP ProLiant Blades power the processing, running 24 hours a day
- 34 racks of HP Servers – each with 4 HP BladeSystem Chassis, 32 servers (16 BL2 × 220c)
- over 30,000 processing cores
- 104 TB RAM

世界一の映画はHPのテクノロジーで作られました

