

High Performance and Productivity

並列プログラミング課題と挑戦

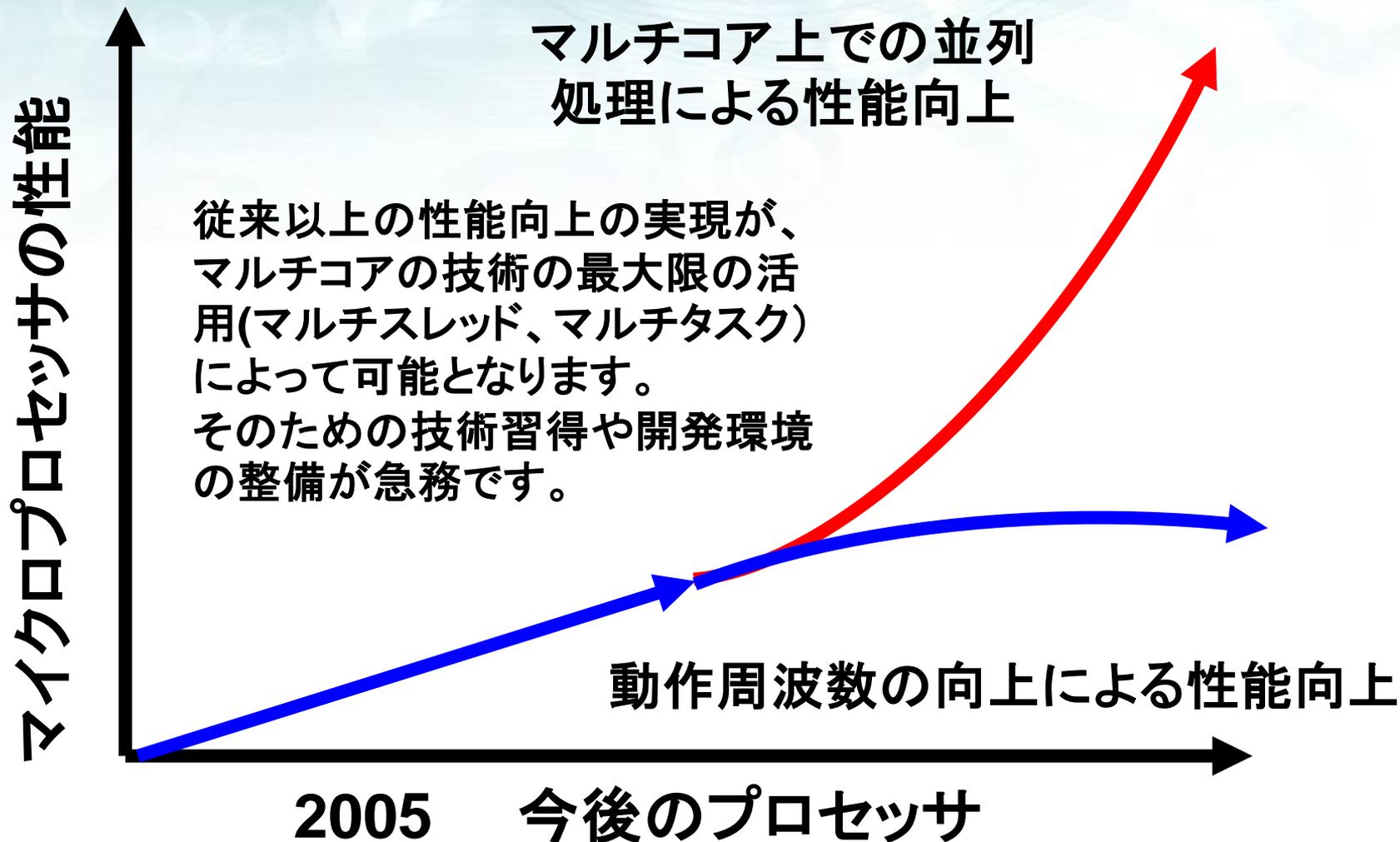
スケーラブルシステムズ株式会社



HPCシステムの利用の拡大の背景

- シミュレーションへの要求
 - より複雑な問題をより精度良くシミュレーションすることが求められている
- HPCシステムでの並列処理の要求の拡大
 1. モデル、アルゴリズム、解析対象は何れもより複雑で、規模の大きなものになっている
 2. マイクロプロセッサのマルチコア化
 3. クラスタに代表される並列計算機の一般化
HPCシステムの一般化
の低コスト化

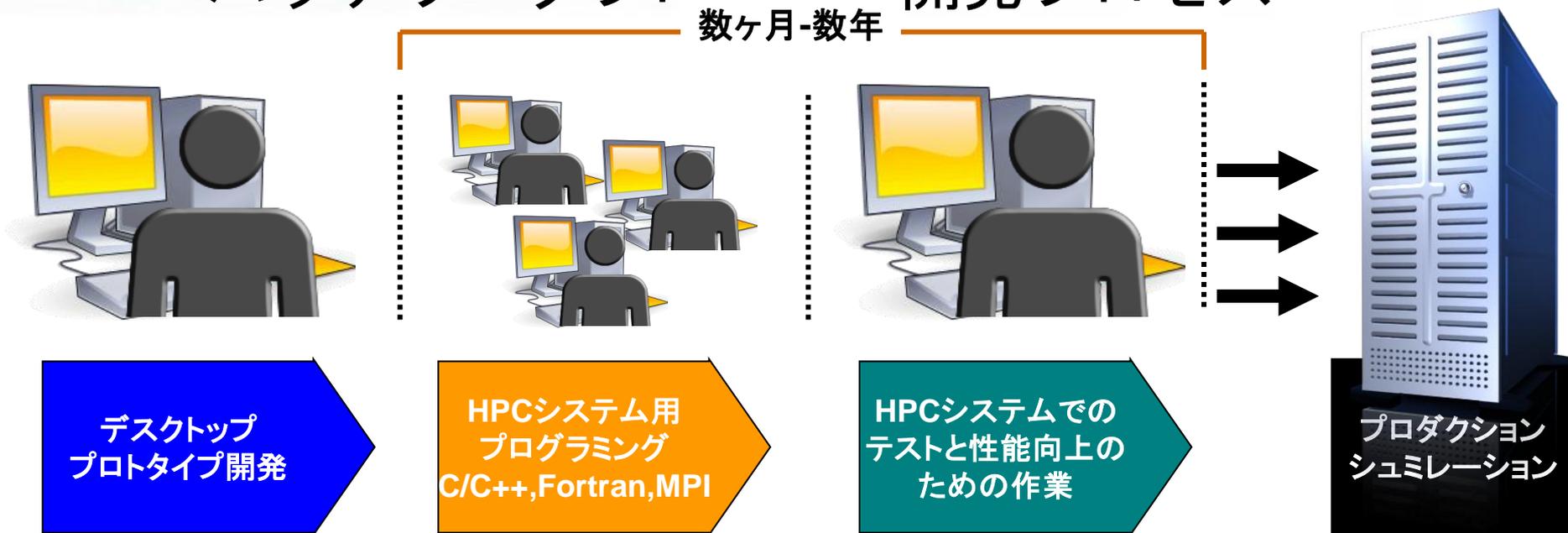
マイクロプロセッサの性能向上 動作周波数からマルチコアへ



並列アプリケーション

- 従来型の並列アプリケーションの開発プロセス

– バッチワークフローでの開発プロセス
数ヶ月-数年

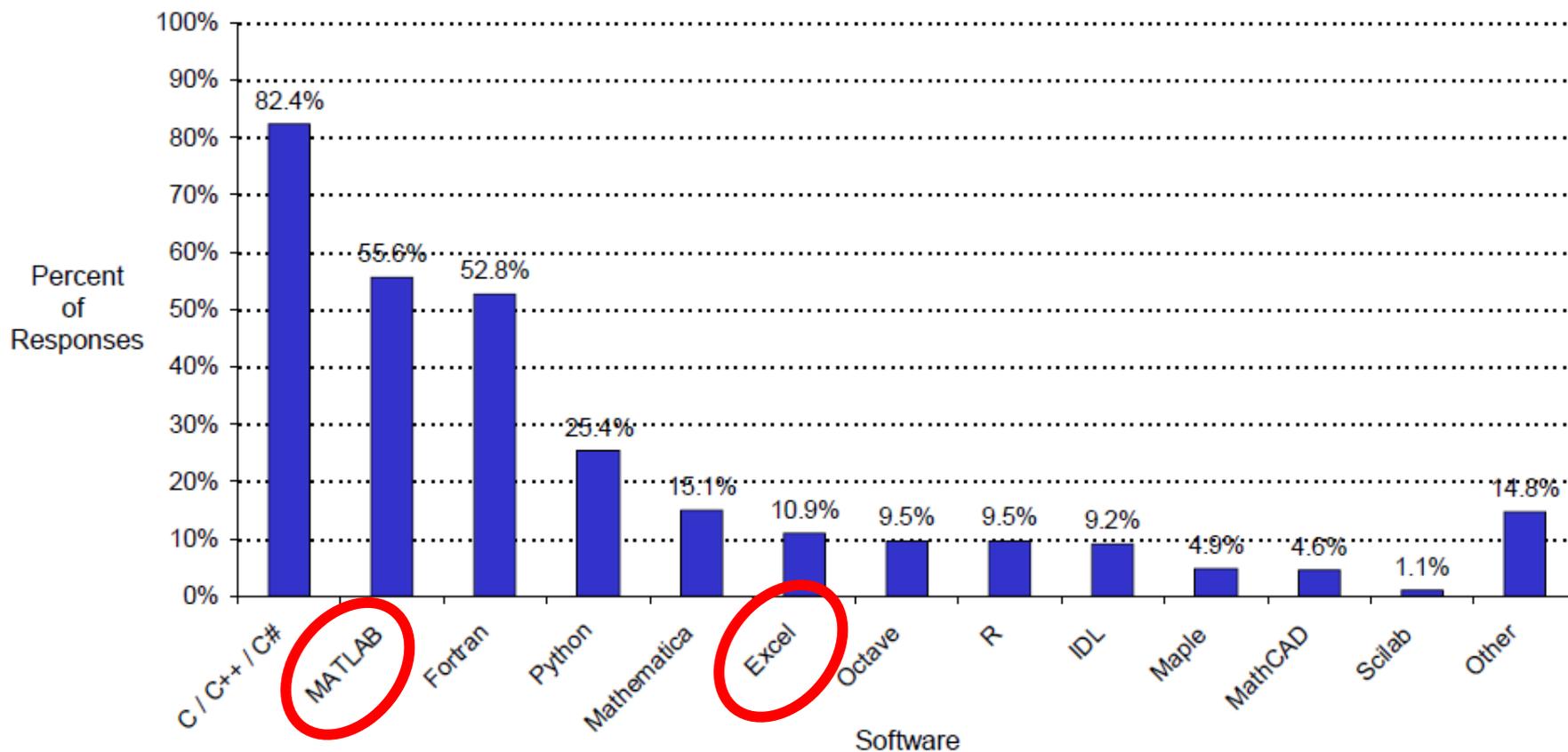


デスクトップ

- デスクトップでのプロトタイプ計算
 - ハイレベル言語（例えば、MATLABなど）
 - マイクロソフト エクセルなどでの処理
- ユーザ
 - システムアーキテクチャやマイクロプロセッサのマイクロアーキテクチャを意識することなくプログラミングを行う
 - GUIを利用した開発環境

並列プログラミングでの課題

PC上でのプログラミングAPI



The Development of Custom Parallel Computing Applications
Simon Management Group
September 2006

複数選択可

スケーラブルシステムズ株式会社

並列アプリケーション

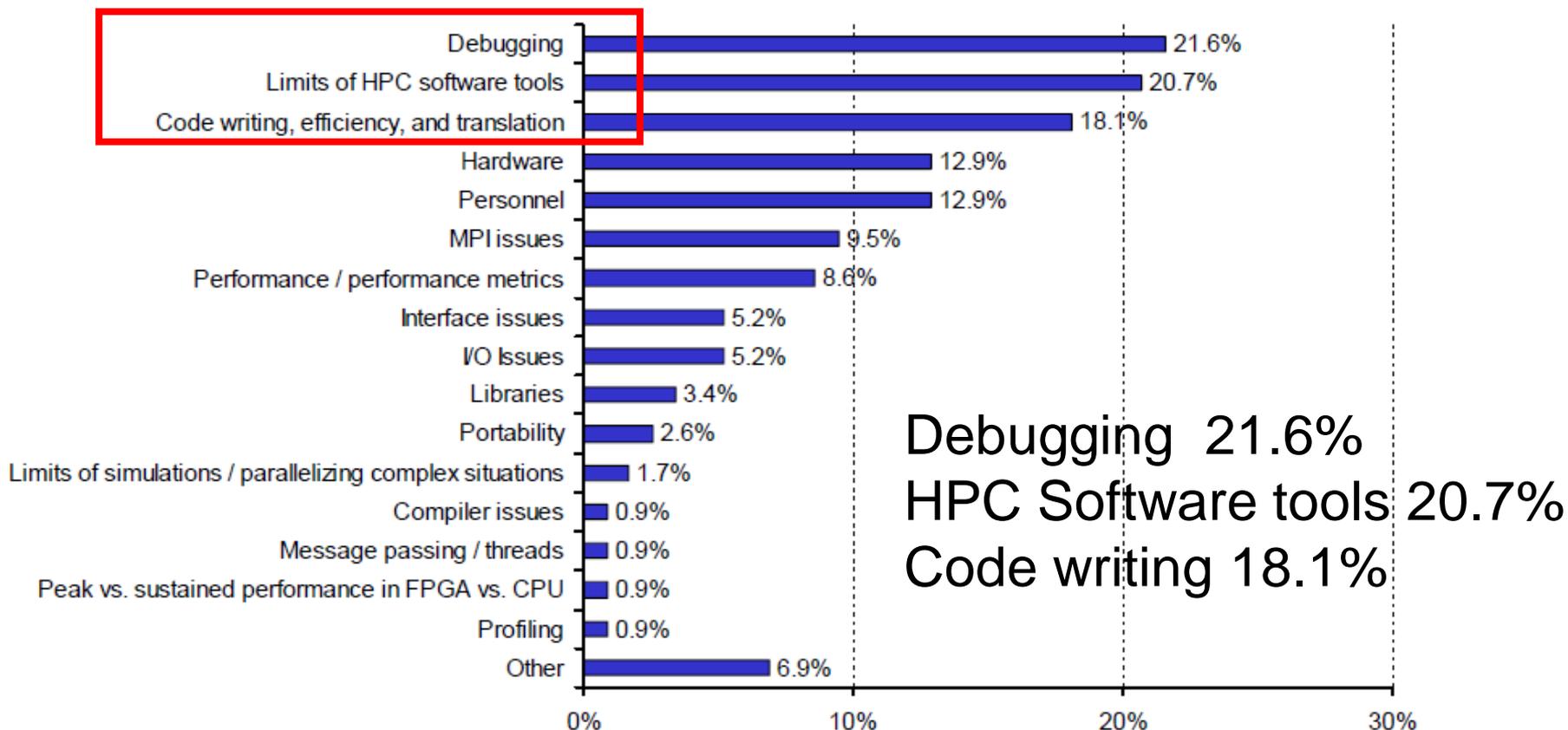
- プログラミング
 - 専門的な知識と並列APIに関する学習
 - C, Fortranなどのコンパイラを利用し、並列処理には、MPI (Message Passing Interface)などのコミュニケーションを明示的に記述
- ユーザ
 - 利用するHPCシステムのアーキテクチャを意識したプログラミング
 - プログラムの開発時にHPCシステムを利用する場合、バッチなどにジョブを投入し、プログラミングの確認を行うことが必要
 - 実際のモデル化やアルゴリズムを実際の解析対象で確認するのはプログラムの完成時まで困難

プログラミングの生産性

- 開発サイクル
 - プログラムのコーディング以外にも様々な作業が必要
- プログラムのデバッグ
 - 実際の解析モデルと入力データが必要
 - 逐次処理や対話処理が必要
- スケーラビリティの実現
 - HPCシステムでの高いスケーラビリティの実現には、高度なプログラミングが必要
 - アルゴリズムの選択やデータのモデル化の検討

並列プログラミングでの課題

- 並列アプリケーション開発でのボトルネック
– デバッグ環境や開発ツールに関する不満



並列プログラミング

並列コンパイラ

並列デバッガ

並列数学ライブラリ

並列コード最適化ツール

スレッド解析ツール
最適化ツール

MPIタスク解析ツール
最適化ツール

統合インターフェイス



ソフトウェアのギャップの解決

デスクトップ

Windows環境
スレッドベースの並列処理
対話処理
豊富なデバッグツールと
開発環境

クラスタシステム

バッチ環境での利用
複雑なデバッグ
MPIなどのメッセージ交換
方式でのプログラミング
Linux (Unix)

ワークステーション
サーバ

クラスタ

#Processors

2

4

8

16

32

64

128

プログラミングのギャップ

数ヶ月-数年



プロダクション
シュミレーション

スケーラブルな性能の
アプリケーションの開発

並列プログラミング
C/C++, MPI
OpenMP

テストと性能
向上のための
作業

プロトタイプ
開発

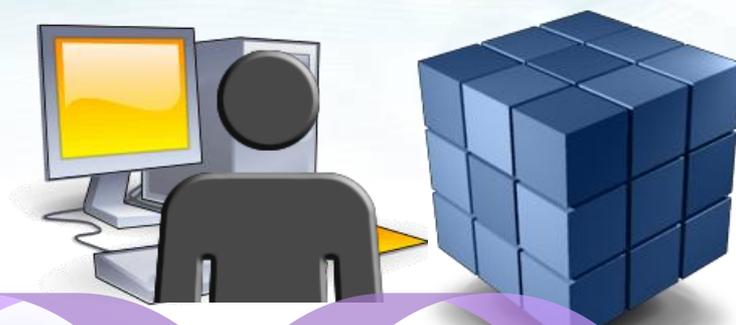
テストと性能
向上の
ための作業

デスクトップ

プロトタイプ開発

プログラミングの生産性の向上

スケーラブルSMPシステム



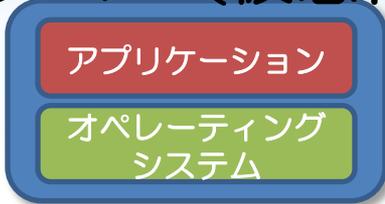
デスクトップ
プロトタイプ開発

HPCシステムでの
テストと性能向上の
ための作業

数ヶ月

ハイエンド仮想化

サーバ（仮想化なし）



サーバ
仮想化

ハイ
エンド
仮想化



仮想化ソフトウェア



一台の仮想マシン



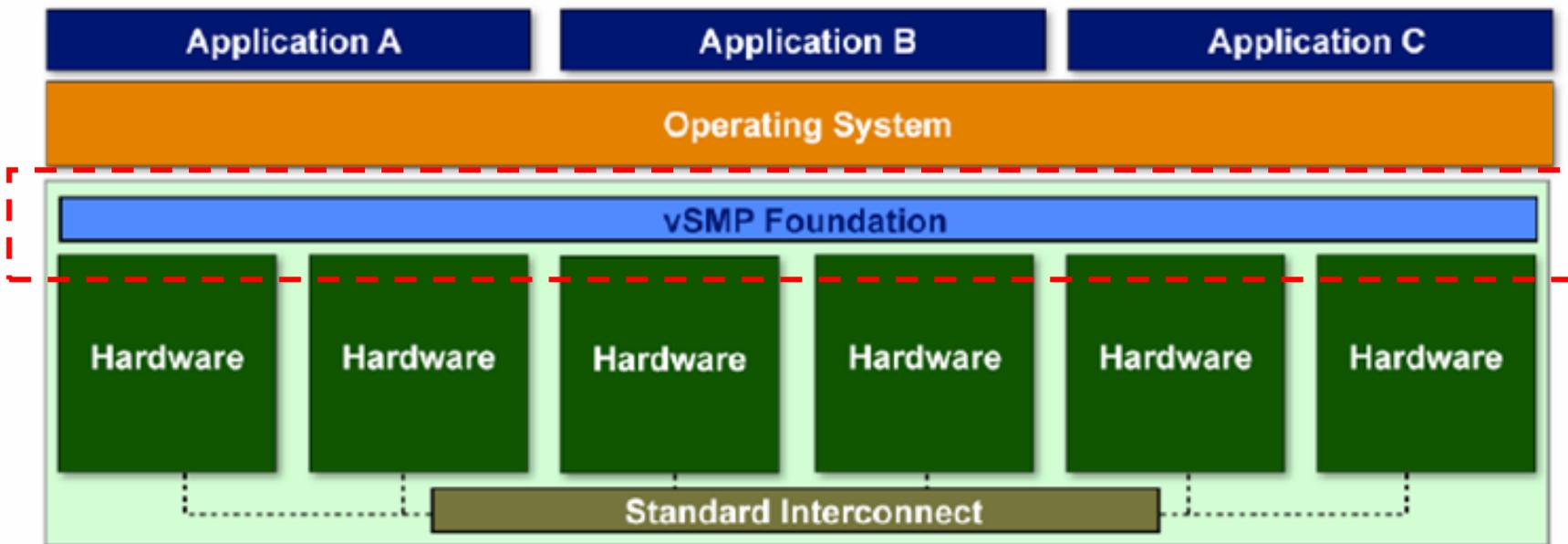
仮想化ソフトウェア



ScaleMP vSMPアーキテクチャ

アプリケーションについては、他のx86システムと100%のバイナリ互換を実現

OSは通常のLinuxディストリビューションが利用可能



Hardwareは一般のx86チップセットと標準インターコネクトでシステムの構築が可能

vSMP Foundation でのシステムのSMP拡張を実現

VXSMP 1440 システム

VXSMP 1440システムとは？

- VXPPO R1440ベースのvSMPシステム
- 1Uの筐体に4台のXeon 55xx/54xx/52xx を搭載し、最大16コアSMPを実現
- 96GBの共有メモリ空間（最大）
- 4台のHDDをOSが共有（最大6TB）



写真はXeon 5400 搭載モデル

VXSMP1400

```
192.168.1.3 - Tera Term VT
File Edit Setup Control Window Help
address sizes : 38 bits physical, 48 bits virtual
power management:

processor      : 14
vendor_id     : GenuineIntel
cpu family    : 6
model         : 23
model name    : Intel(R) Xeon(R) CPU          E5472 @ 3.00GHz
stepping      : 6
cpu MHz       : 3000.110
cache size    : 6144 KB
physical id   : 5
siblings      : 4
core id       : 2
cpu cores     : 4
fpu           : yes
fpu_exception: yes
cpuid level   : 10
wp            : yes
flags         : fpu vme de pse tsc msr pae mce cx8 apic sep mtrr pge
pi mmx fxsr sse sse2 ss ht tm syscall nx lm constant_tsc pni ds_cpl vmmx
bogomips      : 6019.62
clflush size  : 64
cache alignment: 64
address sizes : 38 bits physical, 48 bits virtual
power management:

processor      : 15
vendor_id     : GenuineIntel
cpu family    : 6
model         : 23
model name    : Intel(R) Xeon(R) CPU          E5472 @ 3.00GHz
stepping      : 6
cpu MHz       : 3000.110
cache size    : 6144 KB
physical id   : 5
siblings      : 4
core id       : 3
cpu cores     : 4
fpu           : yes
fpu_exception: yes
cpuid level   : 10
wp            : yes
flags         : fpu vme de pse tsc msr pae mce cx8 apic sep mtrr pge
pi mmx fxsr sse sse2 ss ht tm syscall nx lm constant_tsc pni ds_cpl vmmx
bogomips      : 6022.14
clflush size  : 64
cache alignment: 64
address sizes : 38 bits physical, 48 bits virtual
power management:

[root@localhost ~]#
```

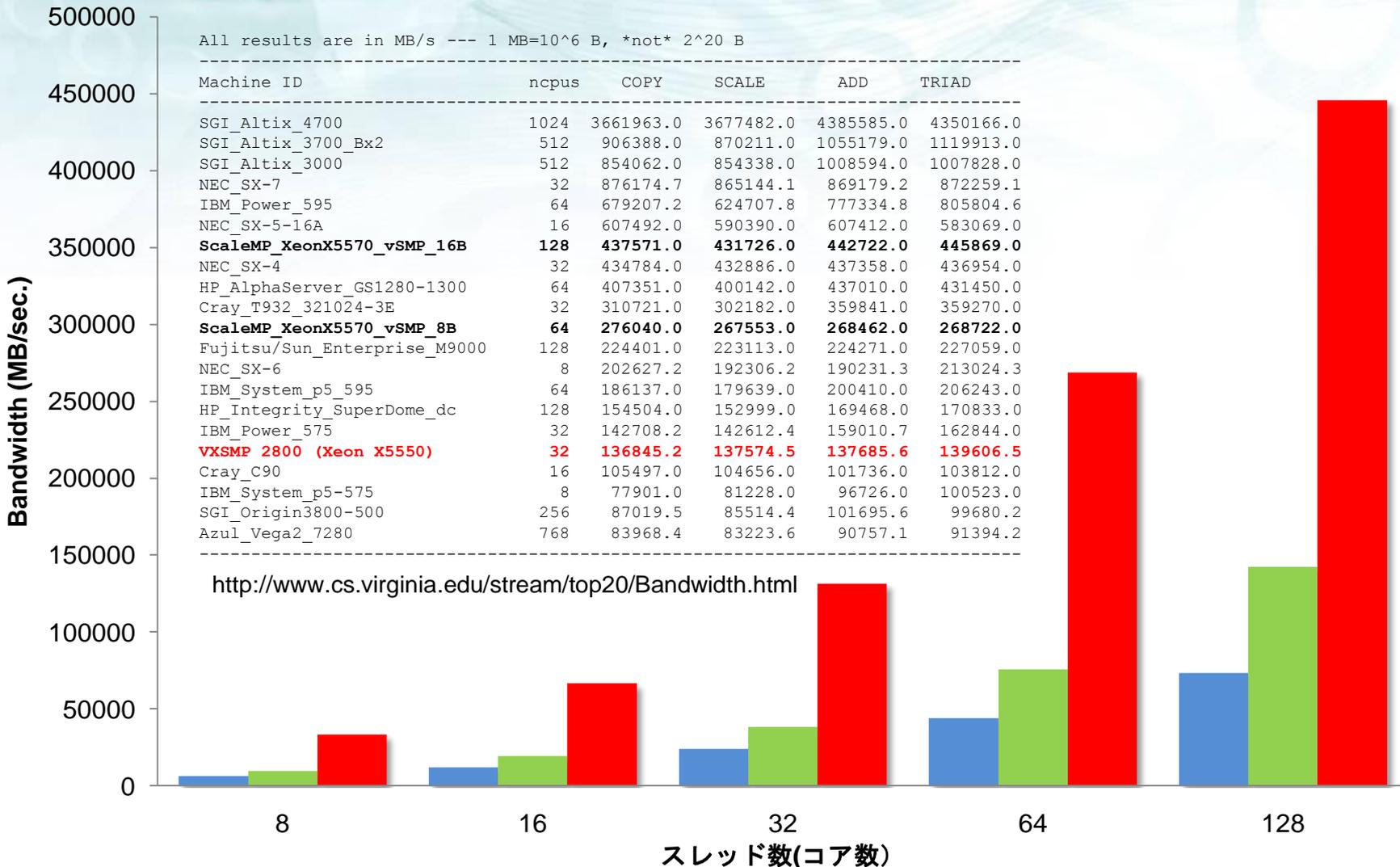
```
192.168.1.3 - Tera Term VT
File Edit Setup Control Window Help
top - 02:25:24 up 37 min, 2 users, load average: 7.38, 2.26, 0.92
Tasks: 239 total, 17 running, 222 sleeping, 0 stopped, 0 zombie
Cpu0 : 98.7%us, 1.3%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu1 : 89.5%us, 1.3%sy, 0.0%ni, 9.2%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu2 : 90.7%us, 1.3%sy, 0.0%ni, 8.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu3 : 88.2%us, 2.6%sy, 0.0%ni, 9.2%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu4 : 90.7%us, 1.3%sy, 0.0%ni, 8.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu5 : 90.7%us, 1.3%sy, 0.0%ni, 8.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu6 : 89.5%us, 2.6%sy, 0.0%ni, 7.9%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu7 : 94.7%us, 0.0%sy, 0.0%ni, 5.3%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu8 : 100.0%us, 0.0%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu9 : 100.0%us, 0.0%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu10 : 98.7%us, 1.3%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu11 : 100.0%us, 0.0%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu12 : 98.7%us, 1.3%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu13 : 100.0%us, 0.0%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu14 : 100.0%us, 0.0%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Cpu15 : 100.0%us, 0.0%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Mem: 63109756k total, 27359056k used, 35750700k free, 121032k buffers
Swap: 4192956k total, 0k used, 4192956k free, 14255632k cached

  PID USER      PR  NI  VIRT  RES  SHR  S %CPU  %MEM    TIME+  COMMAND
 4605 root        25   0 15.3g 12g  640 R   25  20.0   0:11.76 xlinpack_xeon64
 4608 root        25   0 15.3g 12g  640 R   25  20.0   0:12.06 xlinpack_xeon64
 4555 root        25   0 15.3g 12g  640 R   25  20.0   0:25.65 xlinpack_xeon64
 4602 root        25   0 15.3g 12g  640 R   25  20.0   0:11.82 xlinpack_xeon64
 4603 root        25   0 15.3g 12g  640 R   25  20.0   0:11.75 xlinpack_xeon64
 4604 root        25   0 15.3g 12g  640 R   25  20.0   0:11.72 xlinpack_xeon64
 4606 root        25   0 15.3g 12g  640 R   25  20.0   0:12.00 xlinpack_xeon64
 4607 root        25   0 15.3g 12g  640 R   25  20.0   0:12.04 xlinpack_xeon64
 4609 root        25   0 15.3g 12g  640 R   25  20.0   0:12.02 xlinpack_xeon64
 4601 root        25   0 15.3g 12g  640 R   24  20.0   0:06.09 xlinpack_xeon64
 4598 root        25   0 15.3g 12g  640 R   23  20.0   0:05.88 xlinpack_xeon64
 4595 root        25   0 15.3g 12g  640 R   23  20.0   0:05.87 xlinpack_xeon64
 4596 root        25   0 15.3g 12g  640 R   23  20.0   0:05.89 xlinpack_xeon64
 4597 root        25   0 15.3g 12g  640 R   23  20.0   0:05.89 xlinpack_xeon64
 4599 root        25   0 15.3g 12g  640 R   23  20.0   0:06.00 xlinpack_xeon64
 4600 root        25   0 15.3g 12g  640 R   23  20.0   0:05.90 xlinpack_xeon64
 3807 root        25   0 97480 1200  888 S    0   0.0   0:00.12 automount
    1 root        15   0 10304  660  552 S    0   0.0   0:01.98 init
    2 root         RT   0  0  0  0 S    0   0.0   0:00.00 migration/0
    3 root         RT   0  0  0  0 S    0   0.0   0:00.00 ksoftirqd/0
    4 root         RT   0  0  0  0 S    0   0.0   0:00.00 watchdog/0
    5 root         RT   0  0  0  0 S    0   0.0   0:00.00 migration/1
    6 root         RT   0  0  0  0 S    0   0.0   0:00.00 ksoftirqd/1
    7 root         RT   0  0  0  0 S    0   0.0   0:00.00 watchdog/1
    8 root         RT   0  0  0  0 S    0   0.0   0:00.00 migration/2
    9 root         RT   0  0  0  0 S    0   0.0   0:00.00 ksoftirqd/2
   10 root         RT   0  0  0  0 S    0   0.0   0:00.00 watchdog/2
   11 root         RT   0  0  0  0 S    0   0.0   0:00.00 migration/3
   12 root         RT   0  0  0  0 S    0   0.0   0:00.00 ksoftirqd/3
   13 root         RT   0  0  0  0 S    0   0.0   0:00.00 watchdog/3
```

Xeon 5400 搭載モデル
スケラブルシステムズ株式会社

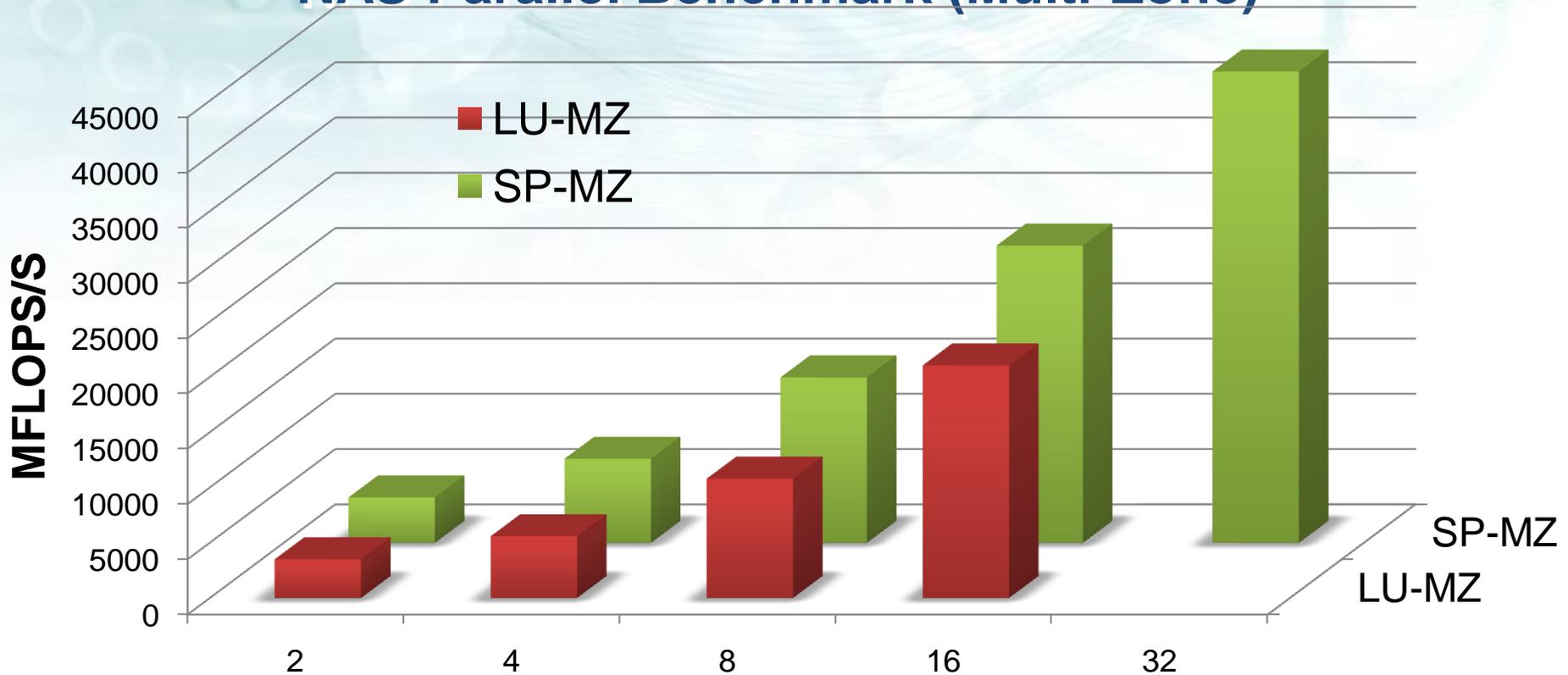
Stream (OMP)ベンチマーク

■ 1333MHz FSB(128cores/16 boards) ■ 1600MHz FSB(128cores/16 boards) ■ 6.4GT/s QPI(128cores/16 boards)



OpenMPベンチマーク

NAS Parallel Benchmark (Multi-Zone)



OpenMPスレッド数/N プロセッサコア

著名な公開ベンチマークツールであるNAS Parallel Benchmark (NPB) の一つであるNPB-MZ (NPB Multi-Zone)はより粒度の大きな並列化の提供を行っています。NPB-MZでは、ハイブリッド型の並列処理やネストしたOpenMPのテストが可能です。ここでの結果は、OpenMPだけでの並列処理の性能を評価しています。

OpenMP プログラム コンパイルと実行例

```
$ cat -n pi.c
 1  #include <omp.h>                                // OpenMP実行時関数呼び出し
 2  #include <stdio.h>                                // のためのヘッダファイルの指定
 3  #include <time.h>
 4  static int num_steps = 1000000000;
 5  double step;
 6  int main ()
 7  {
 8      int i, nthreads;
 9      double start_time, stop_time;
10      double x, pi, sum = 0.0;
11      step = 1.0/((double) num_steps);
12      #pragma omp parallel private(x)              // OpenMPサンプルプログラム:
13      {                                             // 並列実行領域の設定
14          nthreads = omp_get_num_threads();       // 実行時関数によるスレッド数の
15          #pragma omp for reduction(+:sum)        // "for" ワークシェア構文
16          for (i=0;i< num_steps; i++){          // privateとreduction指示
17              x = (i+0.5)*step;                  // の指定
18              sum = sum + 4.0/(1.0+x*x);
19          }
20      }
21      pi = step * sum;
22      printf("%5d Threads : The value of PI is %10.7f¥n",nthreads,pi);

```

取得

句

```
$ icc -O -openmp pi.c
pi.c(14) : (col. 3) remark: OpenMP DEFINED LOOP WAS PARALLELIZED.
pi.c(12) : (col. 2) remark: OpenMP DEFINED REGION WAS PARALLELIZED.
$ setenv OMP_NUM_THREADS 2
$ a.out
 2 Threads : The value of PI is  3.1415927
```

OpenMP指示行

OpenMP実行時関数

コンパイルとメッセージ

環境変数の設定

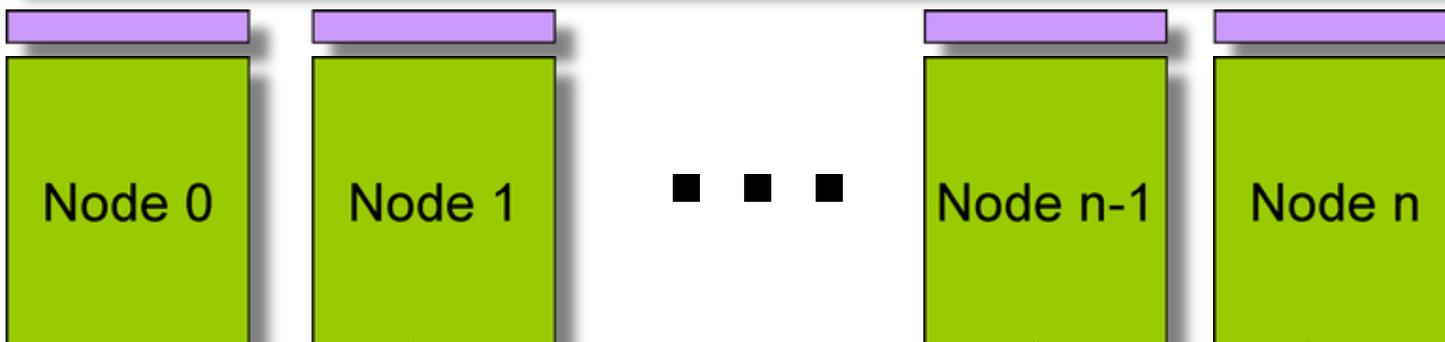
分散仮想共有メモリ (DVSM) インテル クラスタ OpenMP

共有データ



マルチスレッド化されたプログラム

DVSM



Node 0

Node 1

■ ■ ■

Node n-1

Node n



ネットワーク、スイッチ等

Cluster OpenMP プログラム コンパイルと実行例

クラスタ間共有データの定義

```
$ cat -n cpi.c
 1 #include <omp.h>
 2 #include <stdio.h>
 3 #include <time.h>
 4 static int num_steps = 1000000;
 5 double step;
 6 #pragma intel omp sharable(num_steps)
 7 #pragma intel omp sharable(step)
 8 int main ()
 9 {
10 int i, nthreads;
11 double start_time, stop_time;
12 double x, pi, sum = 0.0;
13 #pragma intel omp sharable(sum)
14 step = 1.0/(double) num_steps;
15 #pragma omp parallel private(x)
16 {
17     nthreads = omp_get_num_threads();
18 #pragma omp for reduction(+:sum)
19     for (i=0;i< num_steps; i++){
20         x = (i+0.5)*step; // の指定
21         sum = sum + 4.0/(1.0+x*x);
22     }
23 }
24 pi = step * sum;
25 printf("%5d Threads : The value of PI is %10.7f¥n",nthreads,pi);
26 }
27
```

// OpenMP実行時間関数呼び出し
// のためのヘッダファイルの指定

OpenMP実行時間関数

// OpenMPサンプルプログラム:
// 並列実行領域の設定

// 実行時間関数によるスレッド数の取得
// "for" ワークシェア構文
// privateとreduction指示句

コンパイルとメッセージ

```
$ icc -cluster-openmp -O -xT cpi.c
cpi.c(18) : (col. 1) remark: OpenMP DEFINED LOOP WAS PARALLELIZED.
cpi.c(15) : (col. 1) remark: OpenMP DEFINED REGION WAS PARALLELIZED.
$ cat kmp_cluster.ini
--hostlist=node0,node1 --processes=2 --process_threads=2 --no_heartbeat --startup_timeout=500
$ ./a.out
4 Threads : The value of PI is 3.1415927
```

並列実行処理環境の設定

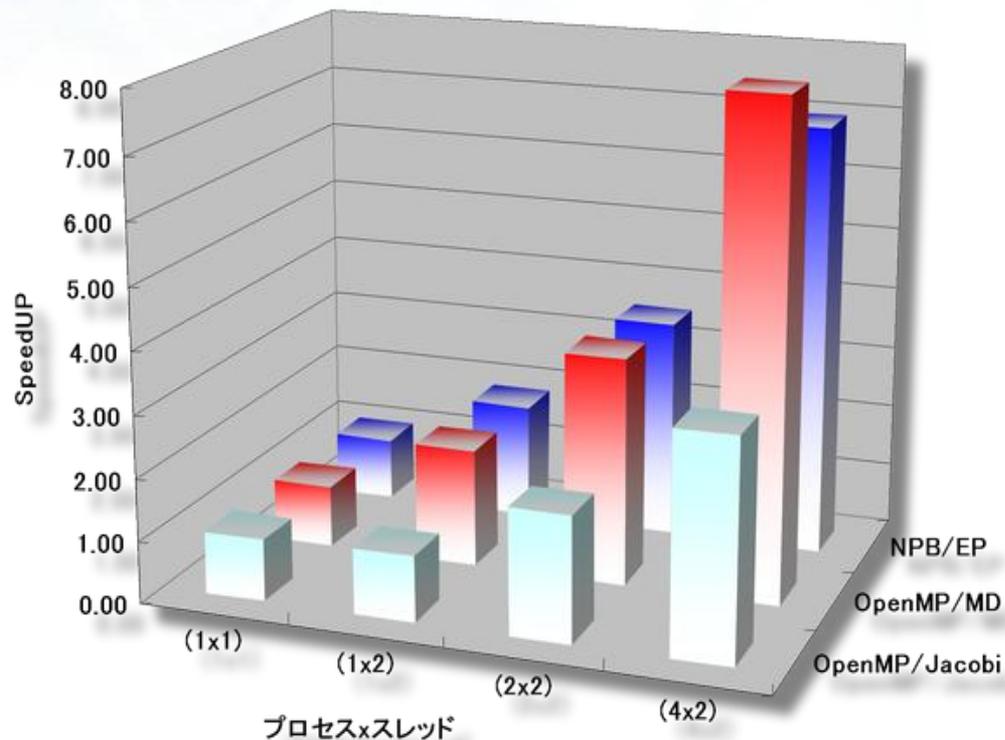
Cluster OpenMP プログラム

ベンチマークシステムラビリティサンプル

- NEXXUS 4820-PT
- 2.66GHz/1066MHz FSB/16GB Memory/InfiniBand

プログラムサンプル

- NAS Parallel Benchmark / EP ベンチマーク
- OpenMPサンプルプログラム（分子動力学サンプル、nparts=10000で実行）
- OpenMPサンプル（Jacobi法サンプル、5000x5000）



シングルAPIでの並列処理

MPI
OpenMP

Cluster OpenMPは、ノード内(SMP)とノード間で同一の並列APIでのプログラミングを可能とします。

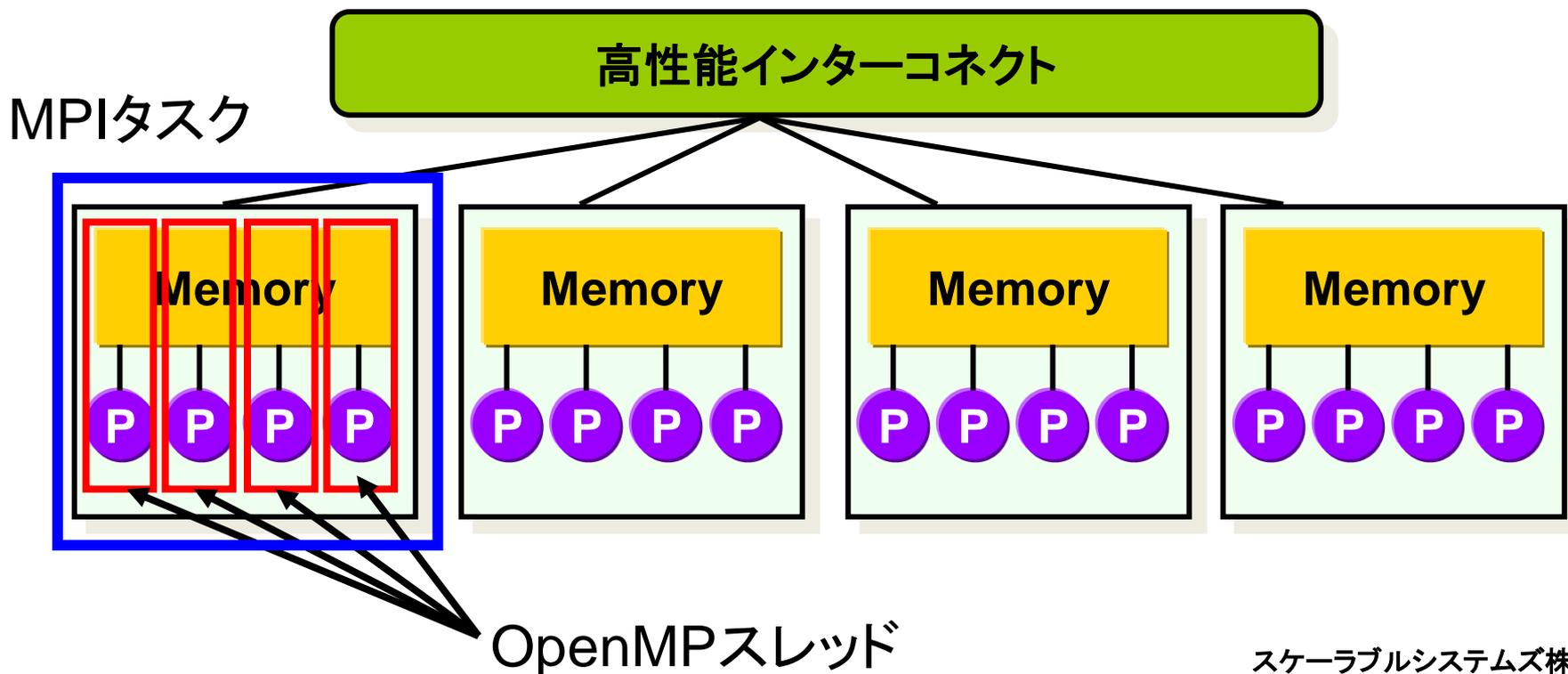
Vertical Scaling



Horizontal Scaling MPI

MPI/OpenMPハイブリッドモデル

- MPIでは領域分割などの疎粒度での並列処理を行う
- OpenMPは、各MPIタスク内で、ループの並列化などにより細粒度での並列化を担う
- 計算は、タスク-スレッドの階層構造を持つ



MPI/OpenMPハイブリッドコード

- MPIで並列化されたアプリケーションにOpenMPでの並列化を追加
- MPI通信とOpenMPでのワークシェアを利用して効率良い並列処理の実現

Fortran

```
include 'mpif.h'
program hybsimp

call MPI_Init(ierr)
call MPI_Comm_rank (... , irank, ierr)
call MPI_Comm_size (... , isize, ierr)
! Setup shared mem, comp. & Comm
!$OMP parallel do
  do i=1,n
    <work>
  enddo
! compute & communicate
call MPI_Finalize(ierr)
end
```

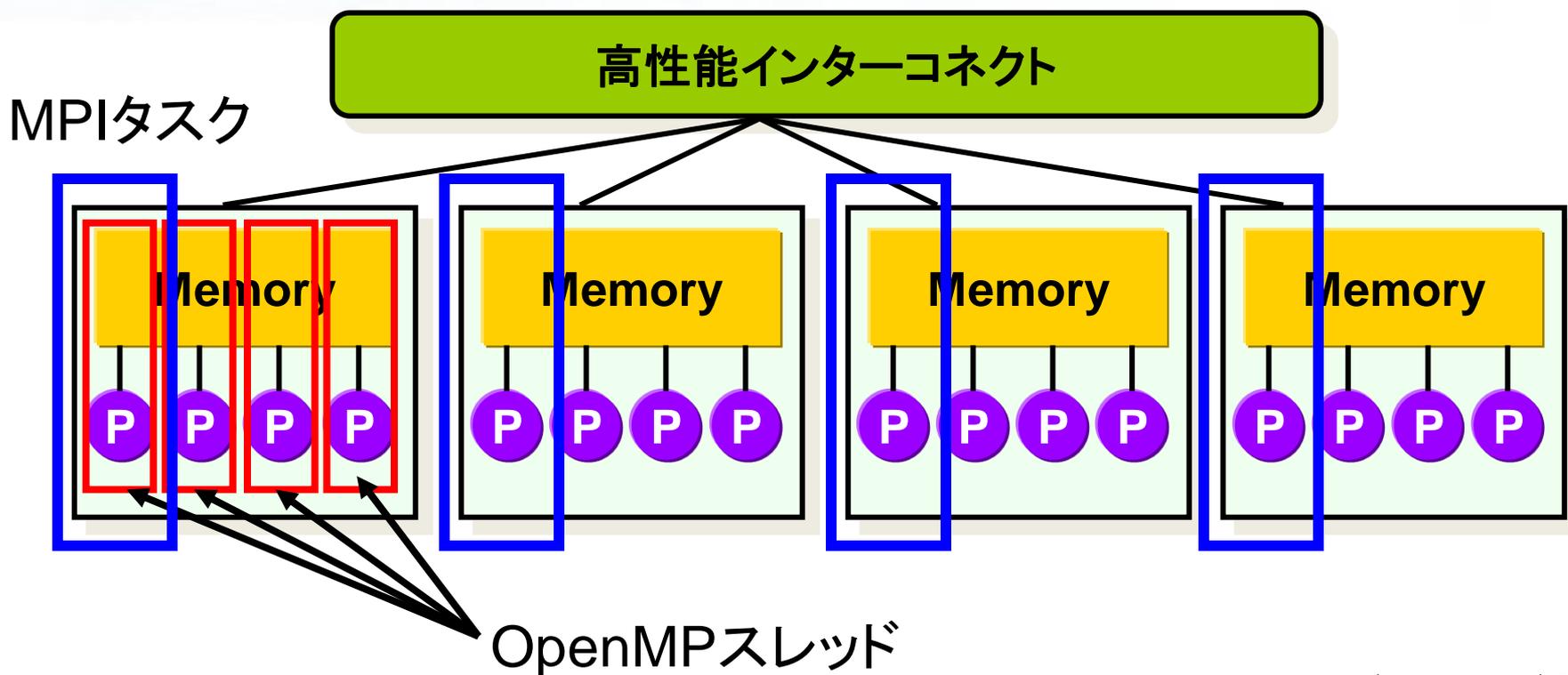
C/C++

```
#include <mpi.h>
int main(int argc, char **argv){
int rank, size, ierr, i;

ierr= MPI_Init(&argc,&argv[]);
ierr= MPI_Comm_rank (... ,&rank);
ierr= MPI_Comm_size (... ,&size);
//Setup shared mem, compute & Comm
#pragma omp parallel for
  for(i=0; i<n; i++){
    <work>
  }
// compute & communicate
ierr= MPI_Finalize();
```

OpenMP/MPIハイブリッドモデル

- MPIは実績のある高性能な通信ライブラリ
- 計算と通信を非同期に実行することも可能
- 通信はマスタースレッド、シングルスレッド、全スレッドで実行することが可能



OpenMP/MPIハイブリッドコード

- OpenMPのプログラムにMPI通信を追加
- 既存のOpenMPプログラムの拡張やスレッドプログラムの新規開発時のオプションとして選択
- MPIは非常に高速また最適化されたデータ通信ライブラリ

Fortran

```
include 'mpif.h'  
program hybmas  
  
!$OMP parallel  
  
!$OMP barrier  
!$OMP master  
call MPI_<Whatever>(...,ierr)  
!$OMP end master  
!$OMP barrier  
  
!$OMP end parallel  
end
```

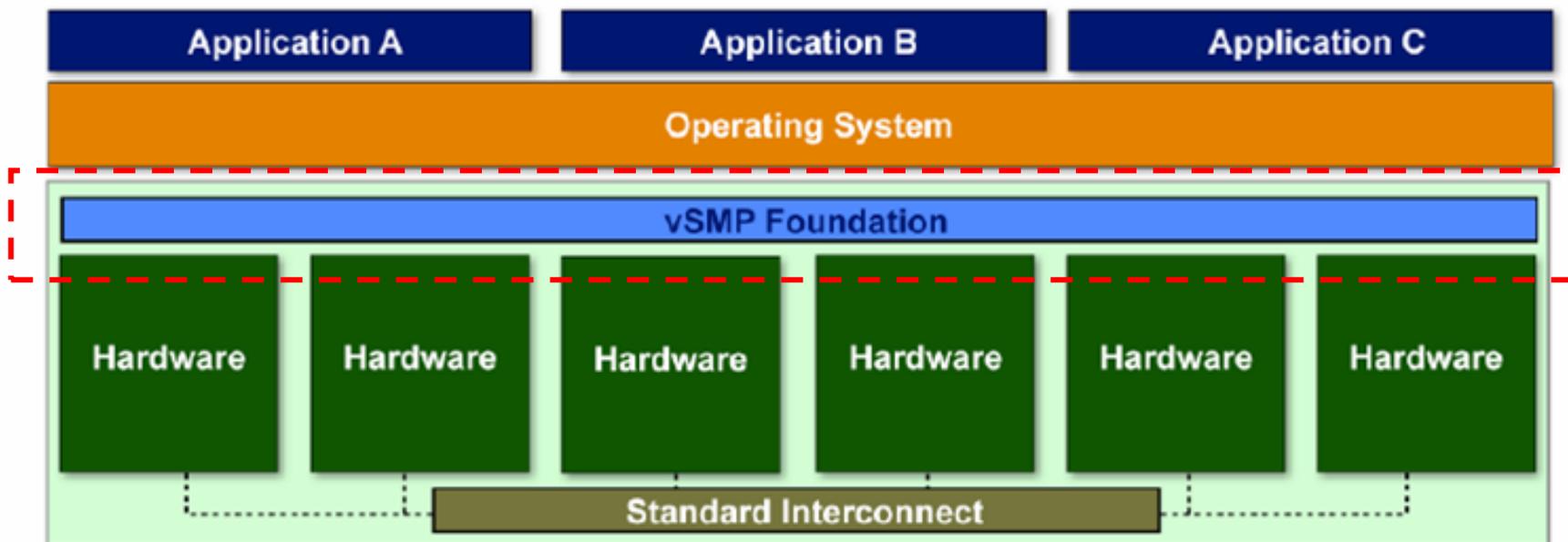
C/C++

```
#include <mpi.h>  
int main(int argc, char **argv){  
int rank, size, ierr, i;  
  
#pragma omp parallel  
{  
#pragma omp barrier  
#pragma omp master  
{  
ierr=MPI_<Whatever>(...)  
}  
#pragma omp barrier  
  
}
```

ScaleMP vSMPアーキテクチャ

アプリケーションについては、他のx86システムと100%のバイナリ互換を実現

OSは通常のLinuxディストリビューションが利用可能



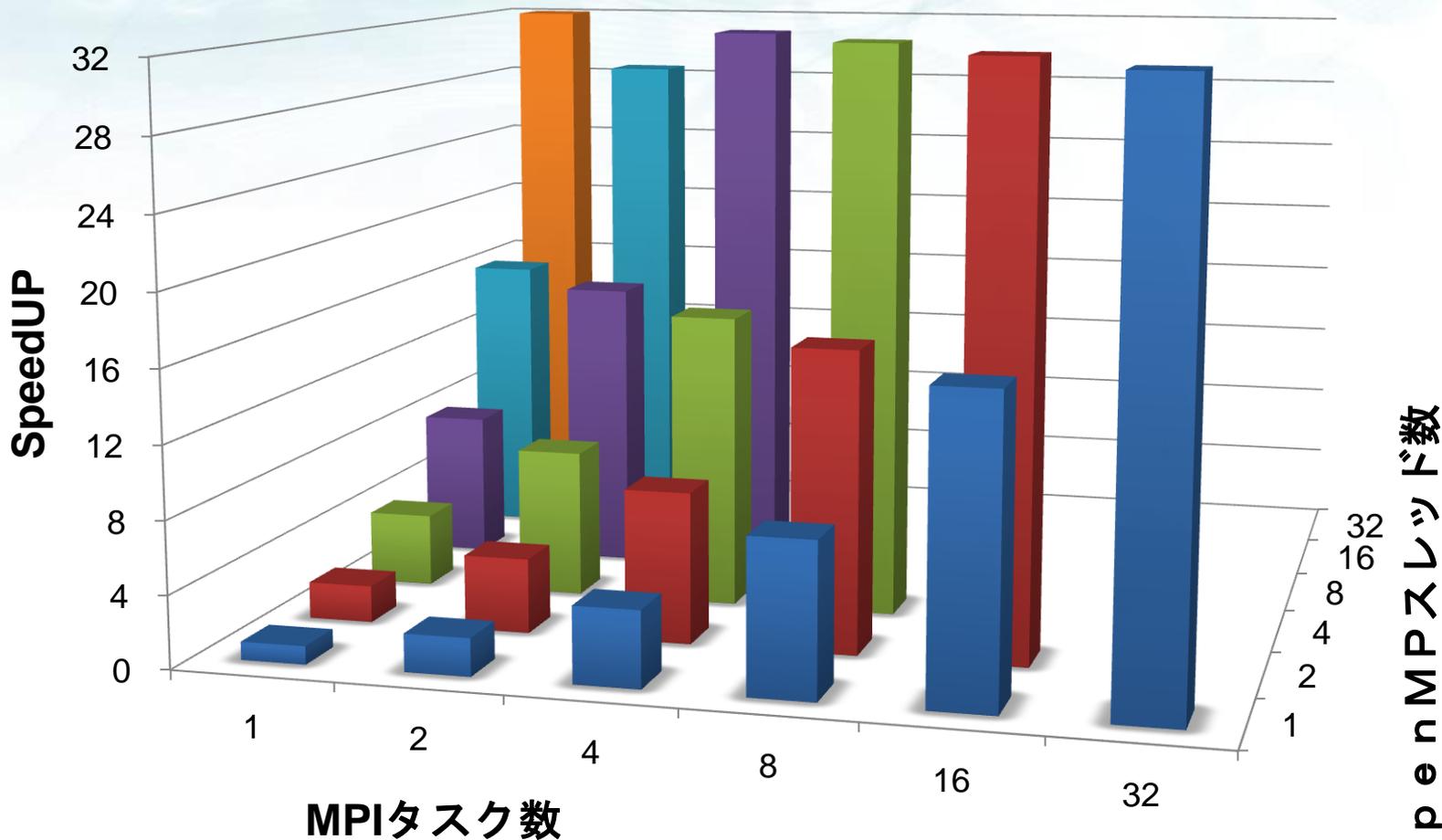
Hardwareは一般のx86チップセットと標準インターコネクでシステムの構築が可能

vSMP Foundation でのシステムのSMP拡張を実現

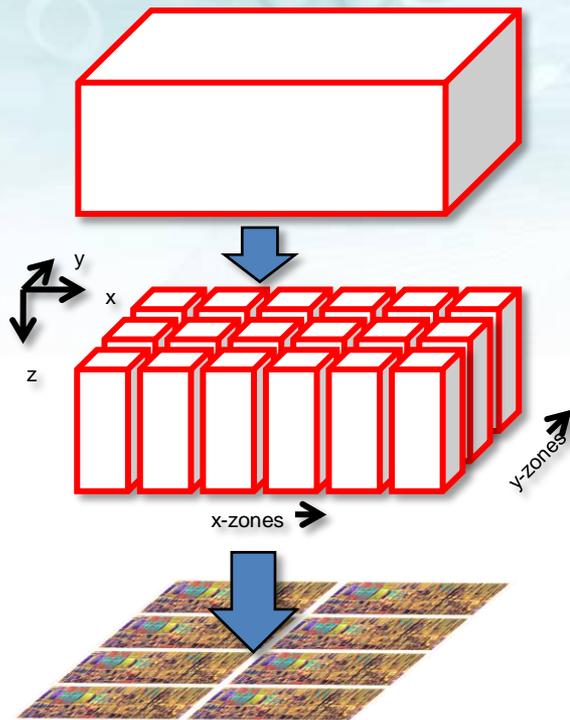
OpenMP/MPI/ハイブリッド

Hybrid OpenMP MPI Benchmarkproject ("homb")

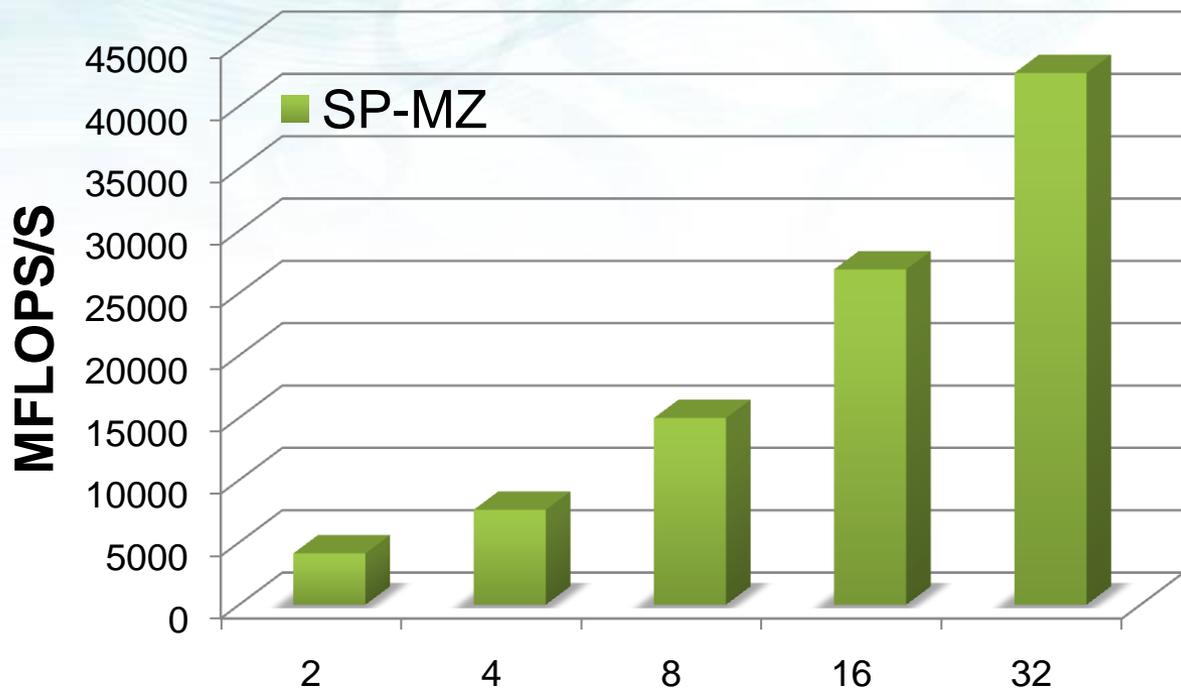
This is the Hybrid OpenMP MPI Benchmarkproject ("homb")
This project was registered on SourceForge.net on May 16, 2009, and is described by the project team as follows:
HOMB is a simple benchmark based on a parallel iterative Laplace solver aimed at comparing the performance of MPI, OpenMP, and hybrid codes on SMP and multi-core based machines.



OpenMPベンチマーク



NAS Parallel Benchmark (Multi-Zone)



OpenMPスレッド数/N プロセッサコア

著名な公開ベンチマークツールであるNAS Parallel Benchmark (NPB) の一つであるNPB-MZ (NPB Multi-Zone)はより粒度の大きな並列化の提供を行っています。NPB-MZでは、ハイブリッド型の並列処理やネストしたOpenMPのテストが可能です。ここでの結果は、OpenMPだけでの並列処理の性能を評価しています。

Xeon 5550 (2.66GHz) vSMP Foundation

ソフトウェアのギャップの解決

デスクトップ

Windows環境
スレッドベースの並列処理
対話処理
豊富なデバッグツールと
開発環境

クラスタシステム

バッチ環境での利用
複雑なデバッグ
MPIなどのメッセージ交換
方式でのプログラミング
Linux (Unix)

ワークステーション
サーバ

vSMP Foundation
プラットフォーム
Cluster OpenMP

クラスタ

#Processors

2

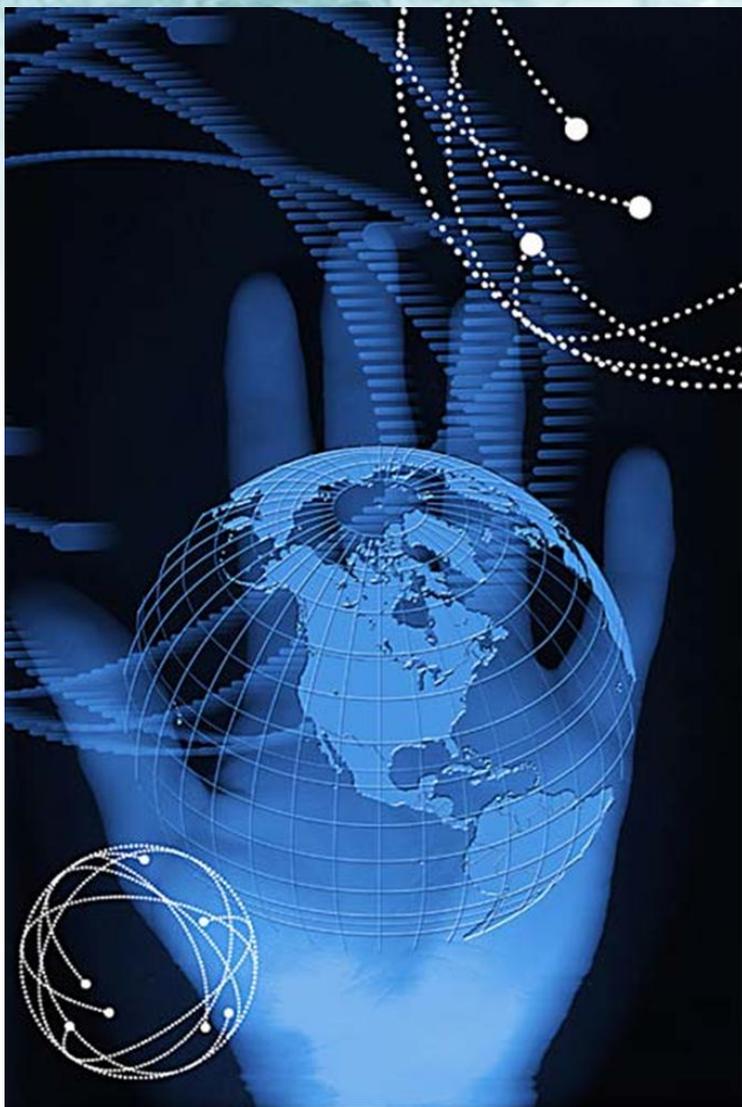
4

8

16

32

64



お問い合わせ

0120-090715 

携帯電話・PHSからは(有料)

03-5875-4718

9:00-18:00 (土日・祝日を除く)

WEBでのお問い合わせ

www.sstc.co.jp/contact

この資料の無断での引用、転載を禁じます。

社名、製品名などは、一般に各社の商標または登録商標です。なお、本文中では、特に®、TMマークは明記しておりません。

In general, the name of the company and the product name, etc. are the trademarks or, registered trademarks of each company.

Copyright Scalable Systems Co., Ltd. , 2009. Unauthorized use is strictly forbidden.

10/16/2009

スケーラブルシステムズ株式会社