

2013年度 実践的並列コンピューティング 第8回

[基礎編最終回]

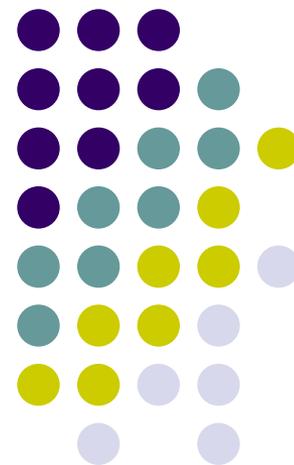
MPIによる

分散メモリ並列プログラミング(3)

遠藤 敏夫

endo@is.titech.ac.jp

2013年6月3日





MPIプログラムの性能を考える

- 前回までは、MPIプログラムの挙動の正しさを議論
- 今回は速度性能に注目

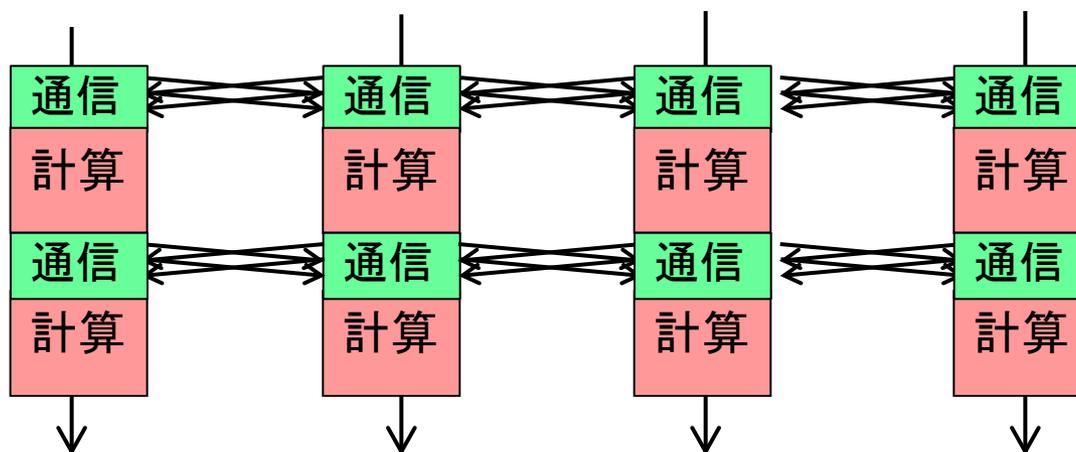
MPIプログラムの実行時間 =

プロセス内計算時間 + プロセス間通信時間

計算量、(プロセス内)ボトルネック有無
メモリアクセス量、計算不均衡...など関係

今回のトピック

MPI版
ステンシル
計算の挙動

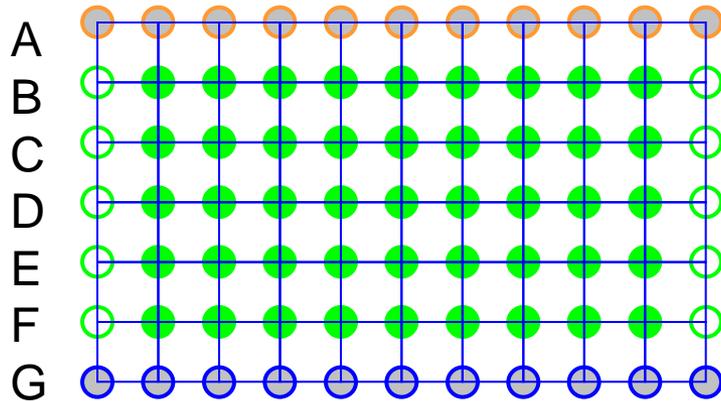


時間
ステップ

ステンシル計算の工夫： 計算と通信のオーバラップ(1)



データの依存関係をより詳細に考えると、他プロセスのデータを必要としない計算が存在する！



行C~Eは通信を待たずに
計算できることに注目
⇒B~Dは通信完了前に
計算してしまってもよい

※ コードがやや複雑になるので、レポート課題[1] (MPIの場合)では必須ではありません



計算と通信のオーバラップ(2)

- オーバラップを組み込んだMPIプログラムの流れ

```
for (it...) {
```

行Bを前のプロセスへ送信開始, Fを次のプロセスへ送信開始

行Aを前のプロセスから受信開始, Gを次のプロセスから受信開始

C~Eの点を計算

上記の全通信終了を待つ(MPI_WaitかMPI_Waitall)

行B, Fの点を計算

二つの配列の切り替え

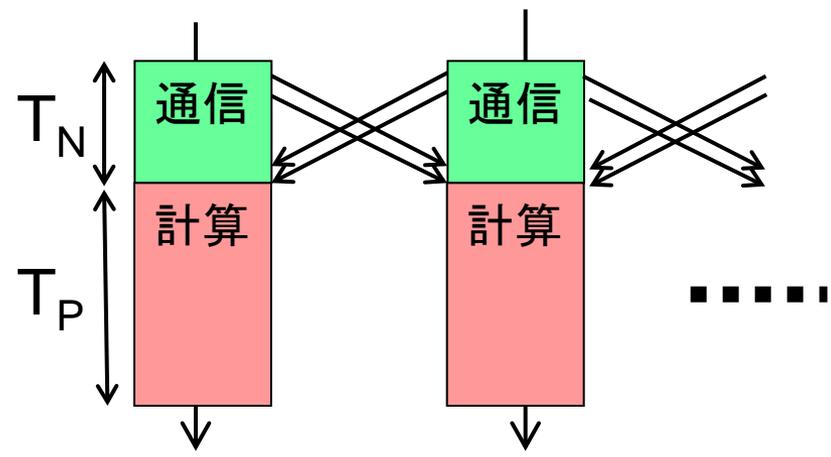
```
}
```



計算と通信のオーバーラップ(3)

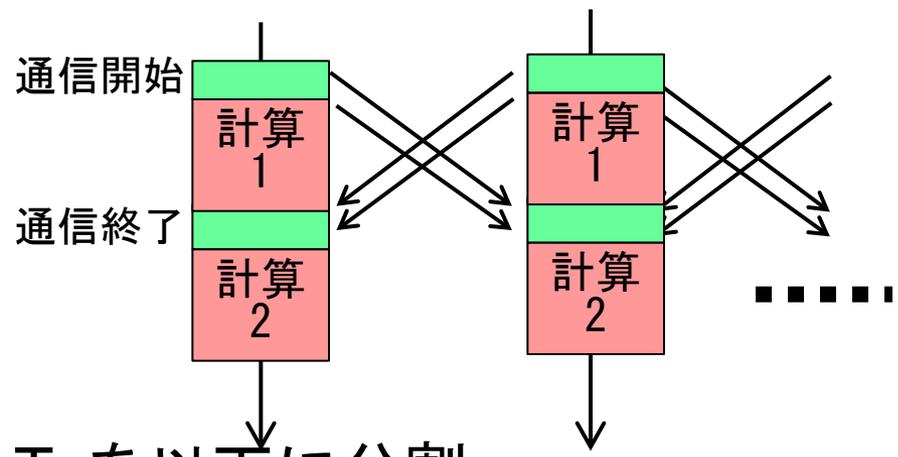
時間ステップあたりの全体時間を T 、通信時間 T_N 、
計算時間 T_P とする

オーバーラップなし



$$T = T_N + T_P$$

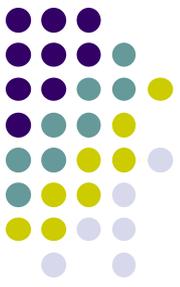
オーバーラップあり



T_P を以下に分割

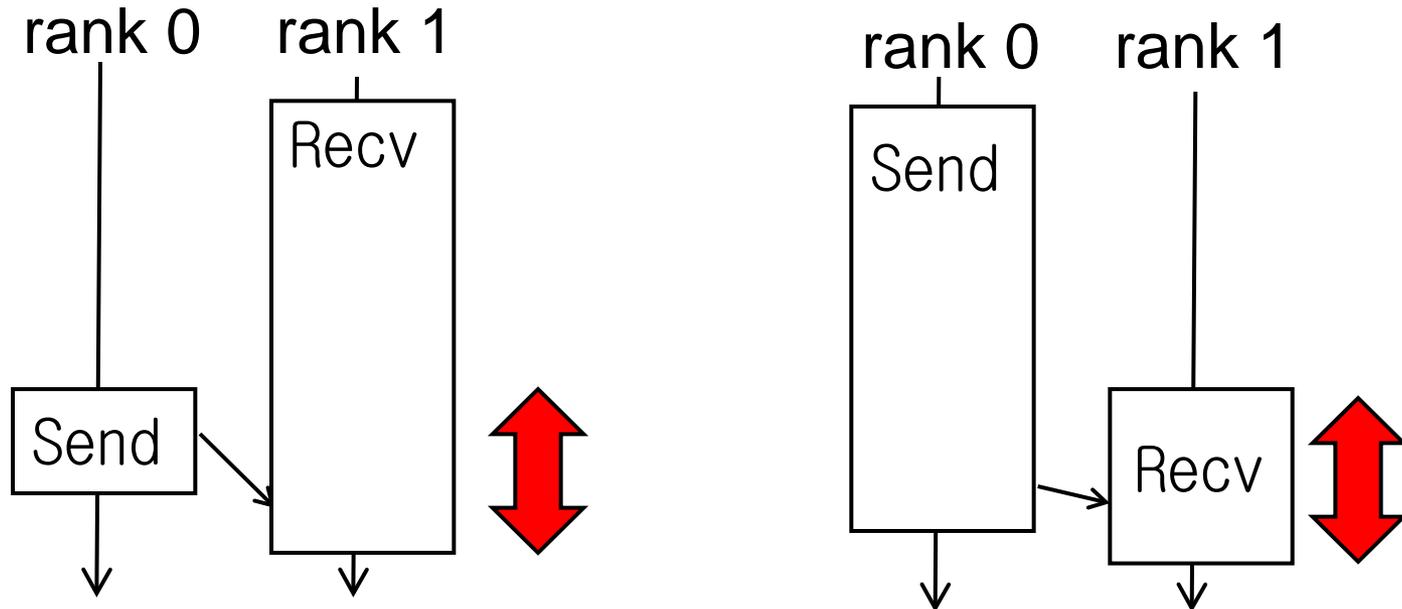
- T_{P1} : オーバーラップ可能部分
- T_{P2} : オーバーラップ不可部分

$$T = \max(T_N, T_{P1}) + T_{P2}$$



通信時間は何によって決まるか

- MPIプログラムの性能を理解するためには、通信時間の理解が必要



- おおまかには、待ち合わせ時間 + **実際の転送時間**



転送時間の性質 (1)

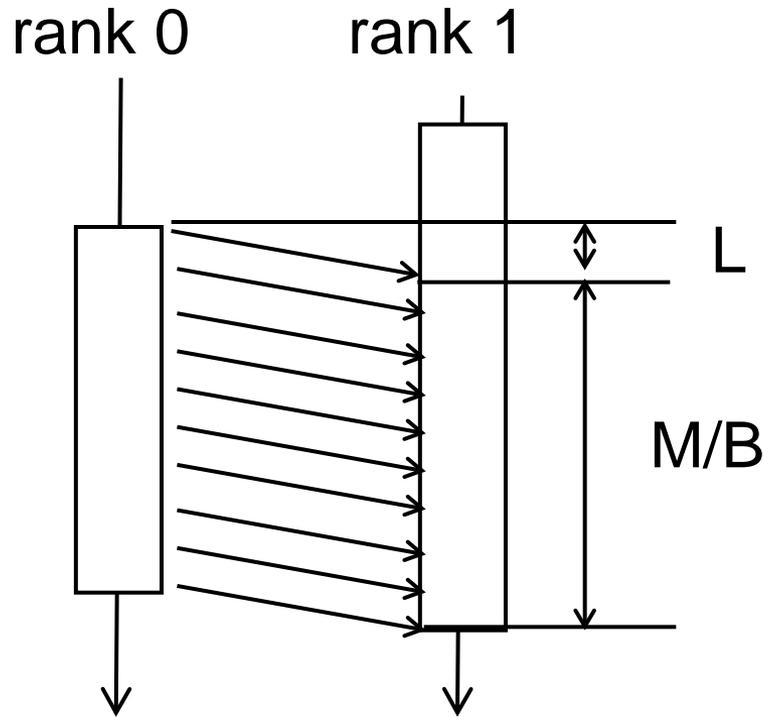
- 転送時間は一般に、
 - メッセージサイズが大きいほど長い
 - メッセージ数が多いほど長い
 - 10KB × 1回の転送時間 < 1KB × 10回の転送時間
- 転送時間の単純なモデル:

$$T = M / B + L$$

転送時間 (sec) メッセージサイズ (Byte) バンド幅 (Byte/sec) 通信レイテンシ (狭義の) (sec)



転送時間の性質 (2)



$$T = M / B + L$$

L: (狭義の)通信レイテンシ

通信の最小単位が、送信者から受信者に届く時間

B: バンド幅

通信路が一秒間に通すことのできるデータ量

※Byte, bitの差に注意。1Byte=8bit

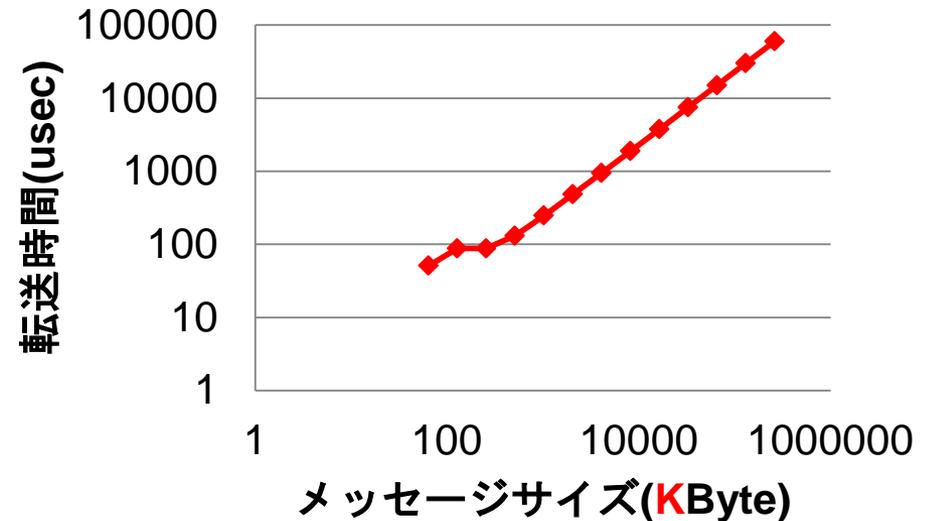
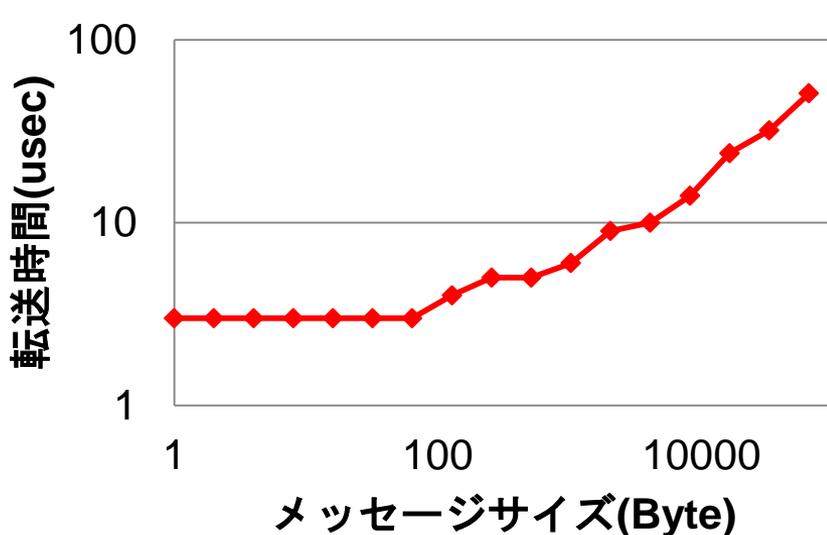
「ギガビットイーサネット」は1Gbps = 125MB/s

※Tを通信レイテンシと呼ぶこともある



TSUBAMEの通信性能は？

- TSUBAME2の2ノード、OpenMPI 1.6.3で測定

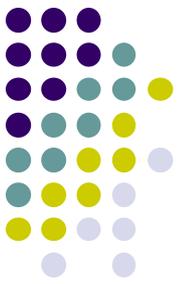


- 切片がL、傾きの逆数がBと考えると

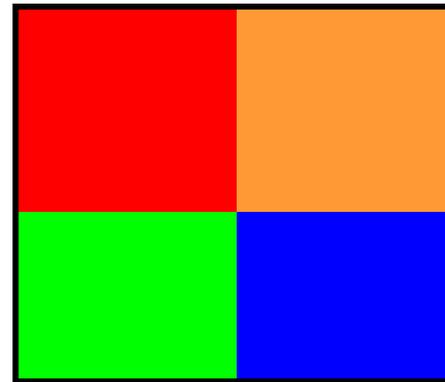
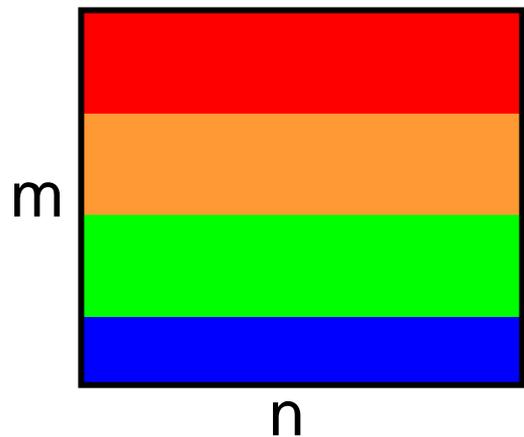
$$L \doteq 3(\text{us}), B \doteq 4.4 (\text{GB/s})$$

理論値が4GB/s x 2本 = 8GBなので、55%くらい

ステンシル計算のさらなる工夫: 多次元分割による通信量削減



$m \times n$ サイズの二次元領域の場合



各プロセスは、上下左右の隣接プロセスと境界交換

(Advectionの場合、斜め隣りも)

ープロセスあたりの

- 計算量: $O(mn/p)$
- 通信量: $O(n)$

ープロセスあたりの

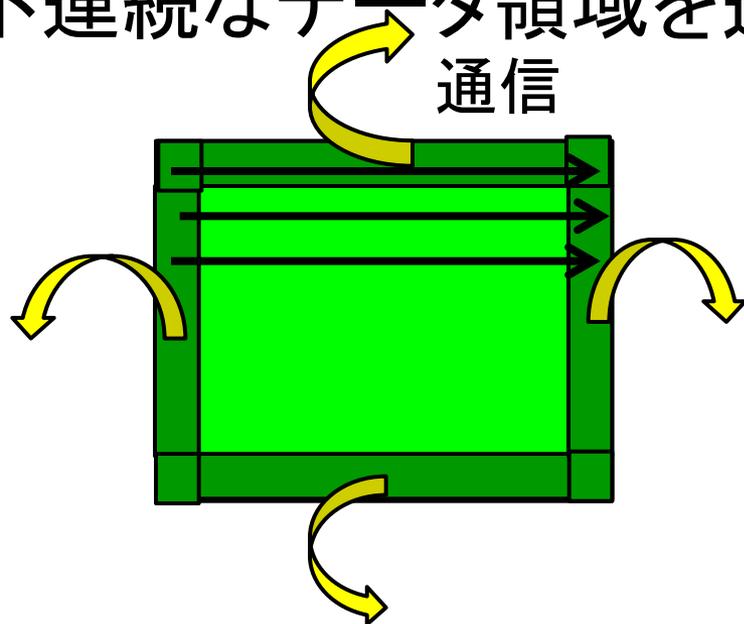
- 計算量: $O(mn/p)$
- 通信量: $O((m+n)/p^{1/2})$

通信量を削減可能



多次元分割と不連続データ

- 多次元分割で通信量合計を削減できるが、不連続なデータ領域を通信する必要がおこる



Row-majorならばの場合、左右プロセスと通信する領域は不連続になってしまう

※ この多次元分割もコードがやや複雑になるので、レポート課題[1] (MPIの場合)では必須ではありません



不連続データを通信するには

考えられる実現方法

- 一要素ずつ通信を繰り返す

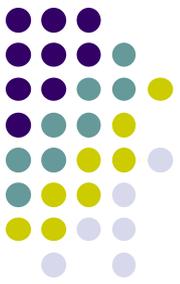
⇒ 性能が大きく低下！メッセージ数が莫大になり、レイテンシを何度も食らうので

- 連続領域に一度コピーしてから通信

⇒ 実装が面倒

- 派生データ型の利用(次ページ以降)

派生データ型(1)



- ユーザが新たにデータ型(MPI_Datatype)を作ることができる
 - MPI_Type_vector関数
 - 他に、MPI_Type_struct, MPI_Type_indexedなど

```
MPI_Datatype mytype;  
MPI_Type_vector(lm, 1, ln, MPI_DOUBLE, &mytype); ←型作成  
MPI_Type_commit(&mytype); ← 利用前に必要な決まり  
:  
(MPI_Send/MPI_Recvなどにmytypeを利用可能)  
MPI_Send(mybuf, 1, mytype, dst, 100, MPI_COMM_WORLD);  
:  
MPI_Type_free(&mytype);
```



派生データ型(2)

`MPI_Type_vector(count, blocklen, stride, oldtype, &newtype)`

等間隔で飛び飛びとなるデータ領域を表す派生データ型を作る

。

count: ブロックの個数

blocklen: 一ブロックの要素数

stride: ブロック間の幅

oldtype: 元となるデータ型

newtype: 新たにつくられる派生データ型

多次元分割はいつも有利なわけではない

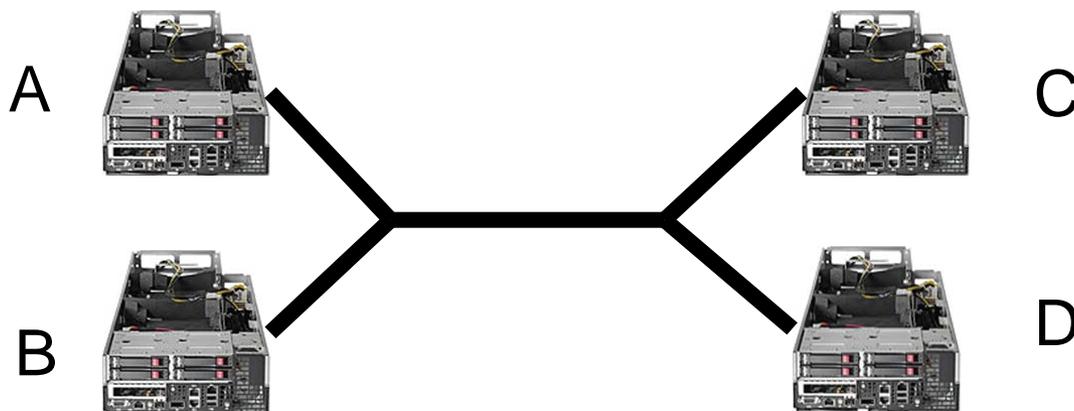


- 派生データ型を使ってさえ、連続領域よりも不連続領域の通信オーバーヘッドは大きくなる
 - 結局MPI関数内部でキャッシュミスやメモリコピーのオーバーヘッド
- 一般的には、プロセス数が大きくなるほど有利
- 三次元領域の計算で、あえて二次元分割にした方が有利なケースも



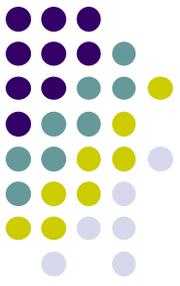
大規模実行で性能を遅くする要素

- 通信路の混雑 (congestion)

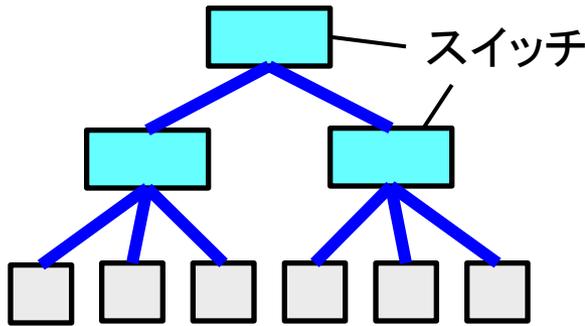


- どのリンクも1GB/sと仮定する。
- A-C間、B-D間も混雑がなければ1GB/sで通信可能
- しかし、A-C, B-D間の同時通信が起こると理論的には半分の0.5GB/sずつになってしまう

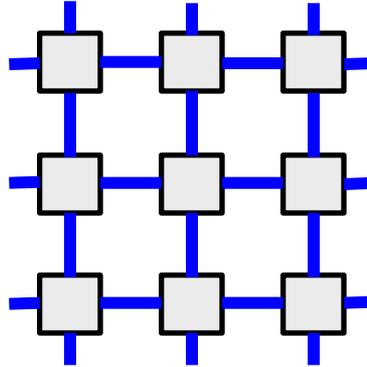
様々なネットワークポロジ



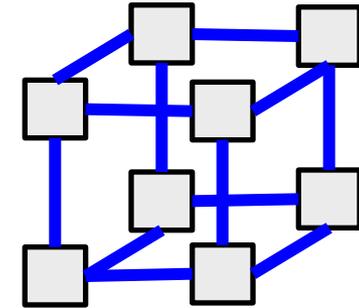
ツリー構造



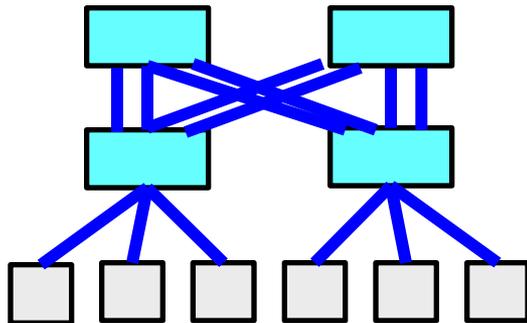
メッシュ/トーラス構造



ハイパーキューブ構造



ファットツリー構造



TSUBAME2など

京、BlueGeneなど

- トーラス: ノード数増加時のパーツ数が少ない
 - TSUBAME2ではそれほど混雑が顕在化することはない
- ⇒ それでも約500ノード以上になると問題が見えてくる

基礎編の終わりに:なぜ並列プログラムは遅くなるのか?(1)



下記のように様々な性能低下原因が考えられる!!

●アルゴリズム上の問題

- 排他制御などによるボトルネック
- 計算量の不均衡...

●システム内での混雑

- メモリへのアクセスの集中
- ノード間ネットワークの混雑
- ストレージへのアクセスの集中・経路でのネットワーク混雑...
 - TSUBAMEではMPIもストレージ(/home, /work)も同じInfiniBandを利用
 - 他プログラムの影響もありうる

基礎編の終わりに:なぜ並列プログラムは遅くなるのか?(2)



- 性能向上のためには、今、なぜ性能が低下しているかの原因を特定する必要
 - 原因特定だけでも決して容易ではない!
 - 別のスパコンに移ると性能が全く変わる可能性も...
- 問題を切り分けるにはどのような実験をすればよいか考える必要がある
 - 通信時間のみの測定により、原因がノード内なのかノード間通信なのか...
 - 1ノードと複数ノードの比較も使える場合も
 - ノード内であれば、(明示的な)ボトルネックがあるか否か、メモリアクセス量はどうか...
- その議論の過程で、スパコンの構造の知識も必要に

再掲:一般的なスパコンの構造

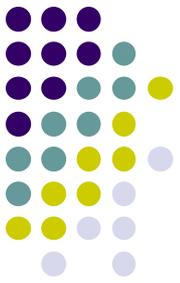


- システム = 多数の**計算ノード** + **外部ストレージ**
 - パーツ間は**ネットワーク**で接続
- 計算ノード = 1以上の**プロセッサ** + **メモリ** + **ローカルストレージ**
 - パーツ間は**PCI-e, QPI**などの**通信路**で接続
- プロセッサ = 1以上の**コア** (+ **L3 キャッシュ**等)



TSUBAME2のネットワーク構成

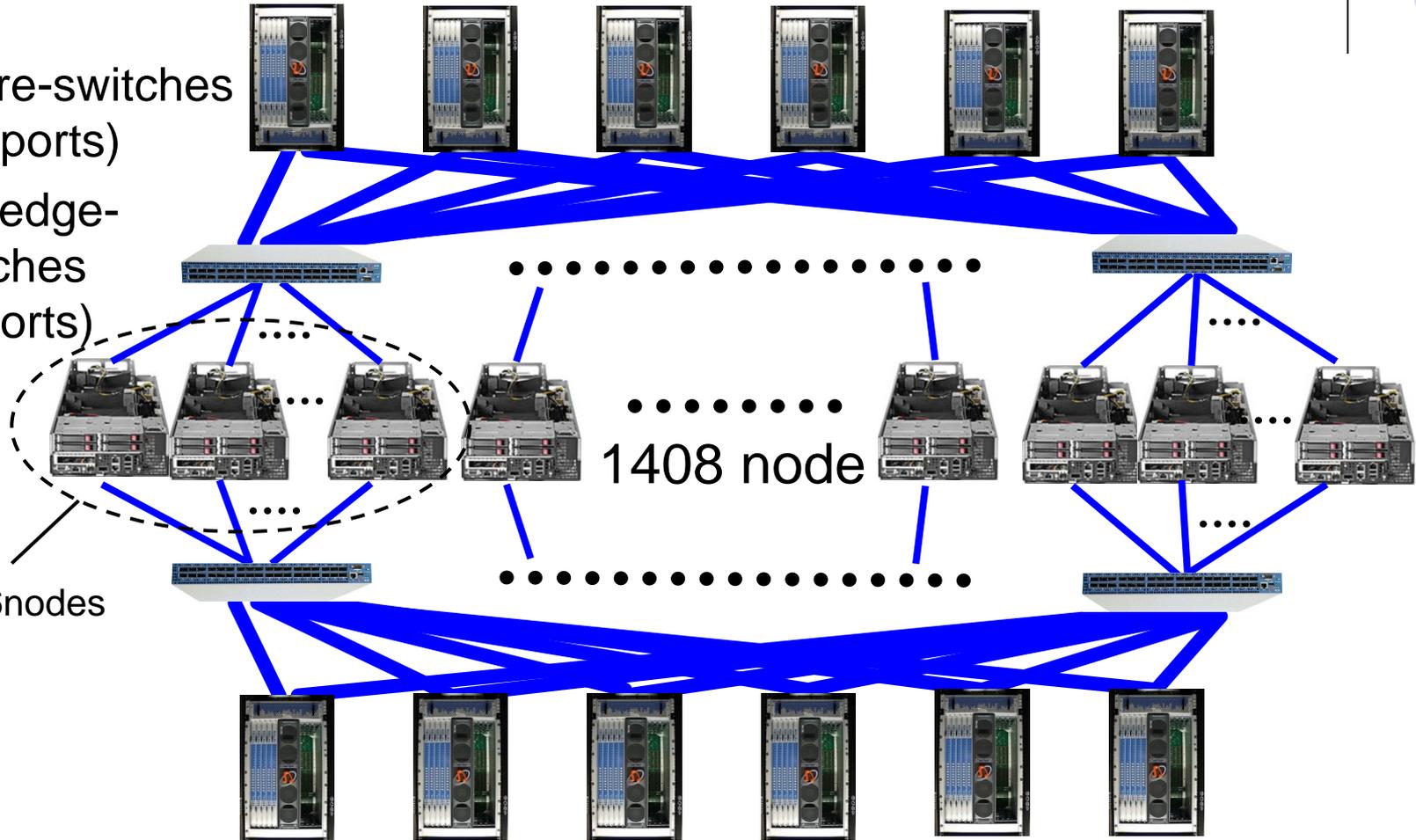
Dual-rail Full-bisection Fat-tree topology



6 core-switches
(324ports)

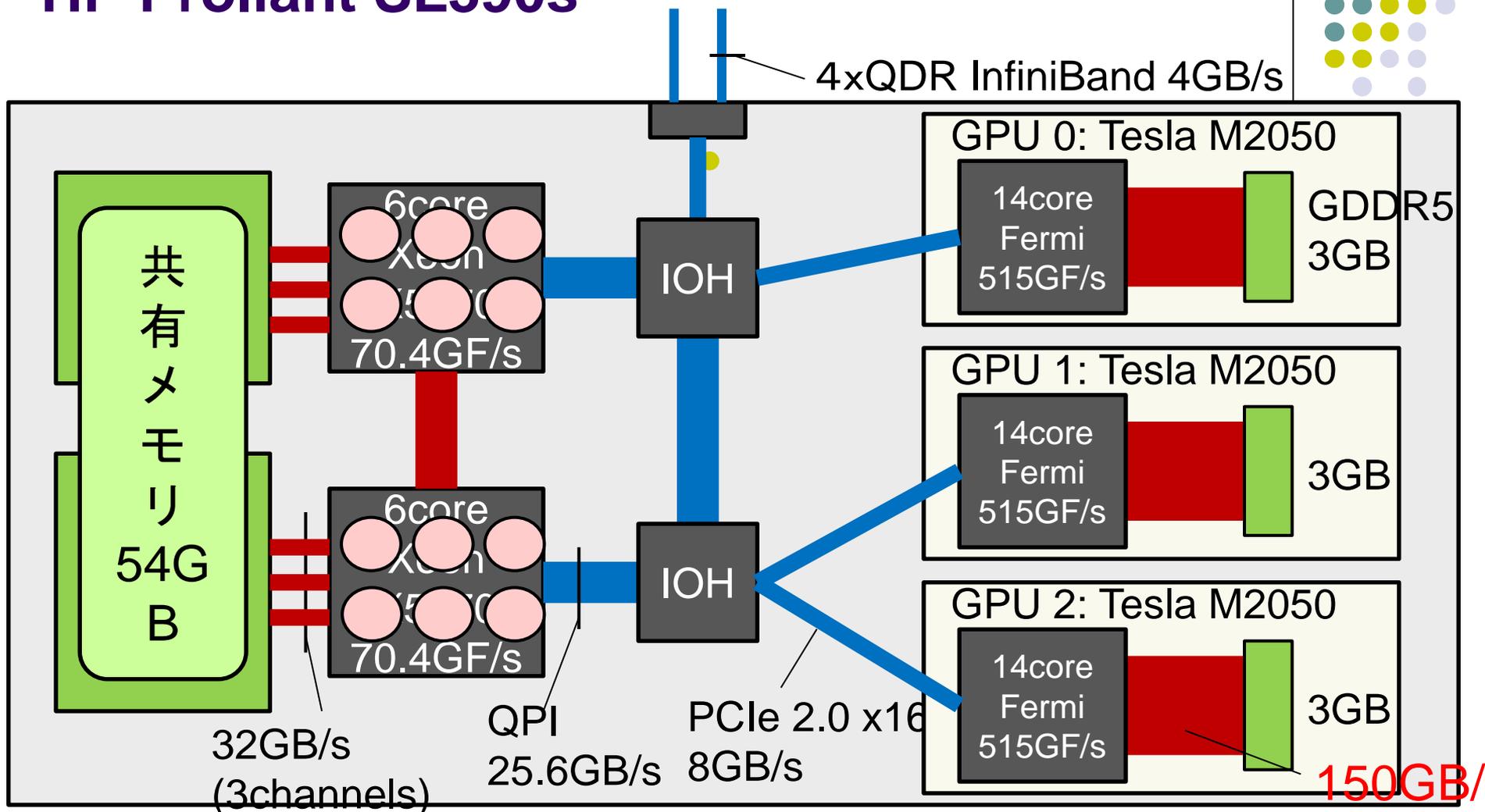
~90 edge-
switches
(36ports)

14~16nodes

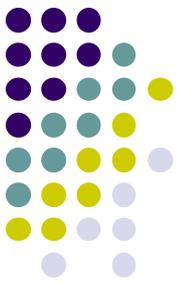


TSUBAME2のノードアーキテクチャ

HP Proliant SL390s



HyperThreadingのため、2ソケット×6コア×2HT=24プロセッサに見える。/proc/cpuinfoファイルの内容を参照



本授業のレポートについて

- 基礎編から一問＋応用編から一問、計二問のレポート提出を必須とします
- 基礎編
 - OpenMP+MPIの、選択問題の中から一問以上
- 応用編
 - CUDA
 - PGAS (予定) } この中から一問以上



基礎編課題説明/Report (1)

以下の[1]—[3]のどれか一つについてレポートを提出してください。二つ以上でも良い。

[1] Advectionサンプルプログラムを、以下のいずれかの方法で並列化してください。

a) OpenMP

b) MPI

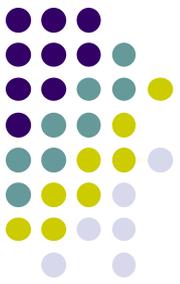
- MPIで一般のサイズに対応するには端数処理が必要である。本レポートではそれは必須ではない
- より良いアルゴリズムにしてもよい。ブロック化・計算順序変更でキャッシュミスが減らせないか？

基礎編課題説明/Report (2)



[2] MPIで並列化され，メモリ利用量を抑えた行列積プログラムを実装してください

- mm-mpiサンプルの改造でよい
- データ分割は本授業の通りでもそれ以外でもよい
- 今回のスライドのアルゴリズムよりも進化した、SUMMA (Scalable Universal Matrix Multiplication Algorithm)[Van de Geijn 1997] もok
- 端数処理はあった方が望ましいが，必須ではない



基礎編課題説明/Report (3)

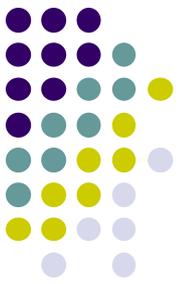
[3] 自由課題: 任意のプログラムを, OpenMPまたはMPI(MPI-2も可)を用いて並列化してください.

- 単純な並列化で済む問題ではないことが望ましい
 - スレッド・プロセス間に依存関係がある
 - 均等分割ではうまくいかない、など
- たとえば, 過去のSuperConの本選問題
<http://www.gsic.titech.ac.jp/supercon/>
たんぱく質類似度(2003), N体問題(2001)・・・
入力データは自分で作る必要あり
- たとえば, 自分が研究している問題



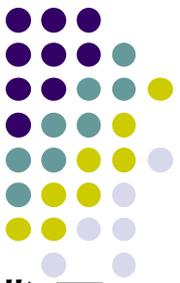
課題の注意

- いずれの課題の場合も、レポートに以下を含むこと
 - 計算・データの割り当て手法の説明
 - TSUBAME2などで実行したときの性能
 - プロセッサ(コア)数を様々に変化させたとき. 大規模のほうがよい. XXコア以上で発生する問題に触れているとなお良い
 - 問題サイズを様々に変化させたとき(可能な問題なら)
 - 高性能化のための工夫が含まれているとなお良い
 - 「XXXのためにXXXを試みたが高速にならなかった」のような失敗でも可
 - プログラムについては, zipなどで圧縮して添付
 - 困難な場合, TSUBAME2の自分のホームディレクトリに置き, 置き場所を連絡. 分かりにくいディレクトリ名推奨



基礎編の課題の提出について

- 提出期限
 - 7/1 (月)
- レポート形式
 - PDF, Word, テキストファイルのいずれか (その他は相談)
 - プログラムも添付 (前ページ参照)
- 送り先: endo@is.titech.ac.jp
- メール題名: ppcomp report



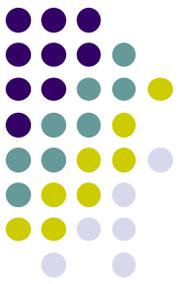
About Account

- TSUBAME2のアカウントができたなら、連絡してください。授業用のTSUBAMEグループへ登録します。

Subject: TSUBAME2 ppcomp account

To: endo@is.titech.ac.jp

- 専攻・研究室
- 学年
- 氏名
- アカウント名



次回/Next Lecture

- 6/10(月)
 - 応用編: CUDAによるGPUプログラミング (1)
- 6/17(月)は休講予定

- スケジュールについてはOCW pageも参照
 - <http://www.el.gsic.titech.ac.jp/~endo/>
→ 2013年度前期情報(OCW) → 講義ノート