

Advanced Network Architecture research group

λコンピューティング環境構築のための Globus Toolkit を用いた MPI ライブラリの実装と評価

大阪大学 大学院情報科学研究科
井本 舞

発表内容

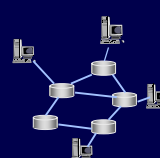
- 4. 研究の背景
- 4. 研究の目的
- 4. 共有メモリシステムを利用した MPI ライブラリの実装方法
- 4. 性能評価
- 4. まとめと今後の課題

2005/4/20 PN研究会 2

研究の背景

- 4. グリッドコンピューティング
 - ネットワークを介して複数の計算機を接続し、計算資源、ストレージを共有
 - ↳ 広域で大規模な計算
 - ↳ 大容量データの転送
- 4. 通信オーバーヘッドが問題
 - TCP/IP が通信に使われる
 - ↳ パケット処理によるオーバーヘッド
 - ↳ パケットロスによる再送遅延

高速かつ、高信頼な通信パイプをエンドユーザに提供する新たな技術が必要



3

λコンピューティング環境の提案

- 4. 計算機、ルータを光ファイバで接続する
- 4. 波長パスを張り、波長パスを通信の最小粒度とする

パケット処理のオーバーヘッドを削除



2005/4/20 4

λコンピューティング環境の提案

- 4. 波長パスを利用して仮想リングを構成する

仮想光リングを構成することで
光リングを専用の高速通信路として利用可能



2005/4/20 5

Globus Toolkit による グリッド環境構築

- 4. Globus Toolkit によってグリッド環境を構築
 - グリッド環境構築のためのミドルウェア
 - 通信、認証、ジョブ管理などを行う
 - ユーザには、各計算機の実装に依存しないインターフェースを提供する
- 4. Globus Toolkit の下位層である通信には、λコンピューティング環境を使用
 - 高速な通信が可能となる

λコンピューティング環境に Globus Toolkit を導入することにより、ユーザはλコンピューティング環境を意識することなく高速な分散計算環境を使うことができる

6

研究の目的

λコンピューティング環境において並列計算を行うことを対象とする

- λコンピューティング環境に Globus Toolkit を導入
- λコンピューティング環境における MPI ライブラリの実装と評価
 - NTT フォトニクス研究所が開発した AWG-STAR システムを利用

2005/4/20 P.N.研究会 7

AWG-STAR システム

- 各ノードは波長ルータ (AWG) に接続し光リングネットワークを構成
- 各ノードは共有メモリボードを搭載
 - 各共有メモリボードは同じデータを保持

2005/4/20 P.N.研究会 8

AWG-STARにおけるデータ共有手法

- 共有メモリボードへのアクセスによりデータ共有が図られる
 - 自ノードの共有メモリへ書き込むことで、他ノードの共有メモリへ反映される
 - 共有データの取得は自ノードの共有メモリボードから読み込むだけで実現できる

2005/4/20 P.N.研究会 9

MPI (Message Passing Interface)

- 並列計算ではメッセージ交換を行いながら計算を進める
- MPI はプロセス間でメッセージを交換するための仕様

```

if(my_rank == 0){
  sprintf(message, "%s", Hello);
  MPI_Send (送信元ランク: 1
            メッセージ: message);
}
else if (my_rank == 1){
  MPI_Recv (送信元ランク: 0
           メッセージ: message);
}

```

2005/4/20 P.N.研究会 10

MPI ライブラリの実装方法

- すべてのライブラリをスクラッチから作成する方法
- 既にあるライブラリを基にして作成する方法
 - Globus Toolkit 上で動作する MPICH-G2 の手法を利用
 - 認証、ジョブ命令送信などは MPICH-G2 と同じ
 - 通信を TCP/IP から AWG-STAR を用いたものに変更

2005/4/20 P.N.研究会 11

λコンピューティング環境における並列計算

MPI アプリケーション	MPI アプリケーション
MPI ライブラリ	共有メモリを用いた MPI ライブラリ
Globus Toolkit	Globus Toolkit
TCP/IP	TCP/IP
イーサネット	イーサネット
	AWG-STAR

ノード間の認証
ジョブ実行命令送信
アプリケーション間の認証
ジョブ実行命令送信
アプリケーションにおけるメッセージパッシング

2005/4/20 P.N.研究会 12

共有メモリを用いた MPI ライブラリの実装

- 共有メモリ上で動的にメモリを割り当てる事ができない
 - プロセスの個数を n とすると、共有メモリを $n \times n$ に分割
 - 一つの領域を一組の送信/受信プロセス間でのデータ交換をする領域として用いる

1組の送信・受信プロセス間でデータ交換するための領域

共有メモリ

2005/4/20 PN研究会 13

共有メモリを用いたメッセージパッシング方法

- 共有メモリ上でデータを交換する
 - データをまるごと共有メモリに書き込むことが可能

送信側 ローカルメモリ

共有メモリ

受信側 ローカルメモリ

パケット分割ヘッダ付加

TCP / IP

パケット組立てヘッダ削除

2005/4/20 14

メッセージ送信側の動作

- 送信関数が呼ばれると、送信データを共有メモリに書きこむ
- 送信データを書き込んだ後、受信プロセスにシグナルを送る
 - シグナルは AWG-STAR が提供する機能
 - 任意のプロセスに送信できる

送信プロセス

データ CCC

データ AAA

データ BBB

データ CCC

受信プロセス

共有メモリ

2005/4/20 15

メッセージ受信側の動作

- 受信関数が呼ばれるタイミングと、データを受信するタイミングが異なることを考慮
 - 受信データバッファと受信要求バッファをローカルメモリ上に設ける

共有メモリ

データ AAA

データ

共有メモリからデータの読み込み

受信データバッファ

データ CCC

データ BBB

受信関数

受信要求バッファ

要求 AAA

共有メモリ

PN研究会 16

メッセージ受信側の動作

- 受信関数が呼ばれるタイミングと、データを受信するタイミングが異なることを考慮
 - 受信データバッファと受信要求バッファをローカルメモリ上に設ける

共有メモリ

受信データバッファ

データ CCC

データ BBB

受信要求バッファ

要求 DDD

受信関数

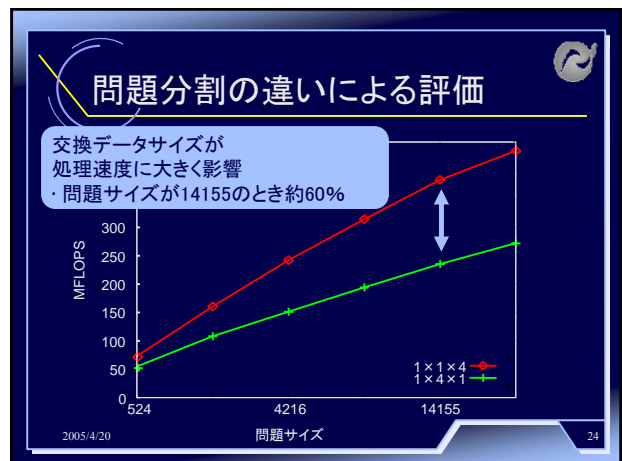
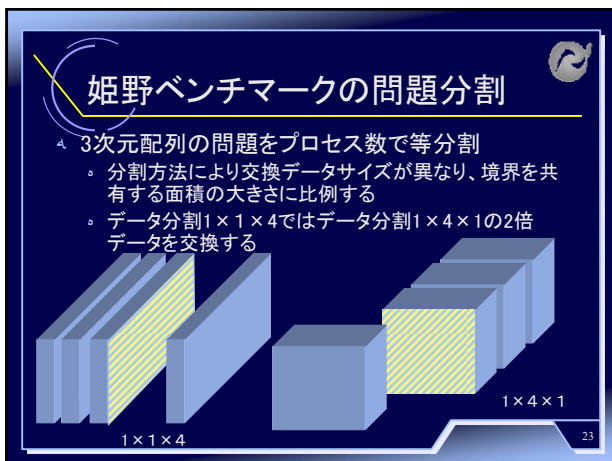
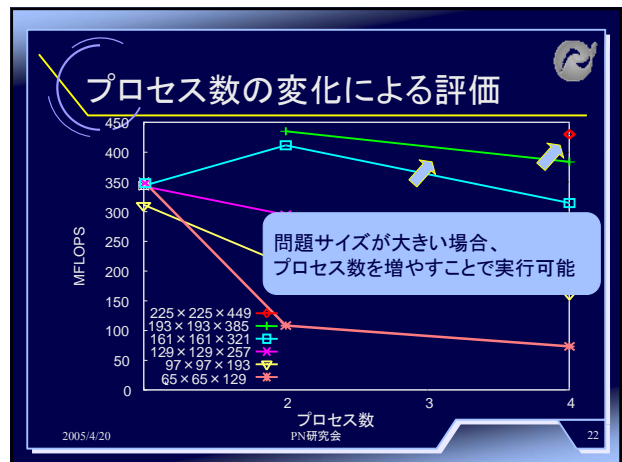
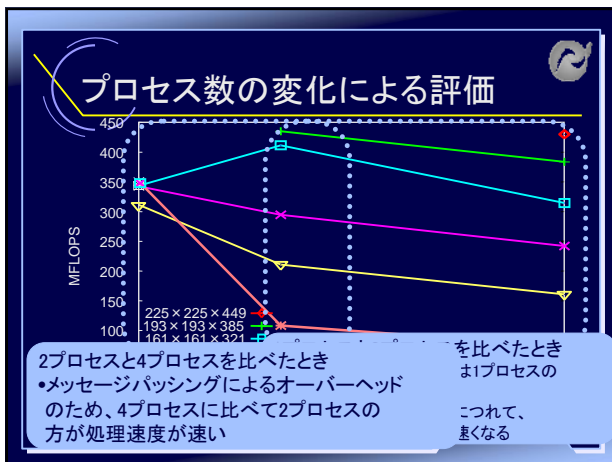
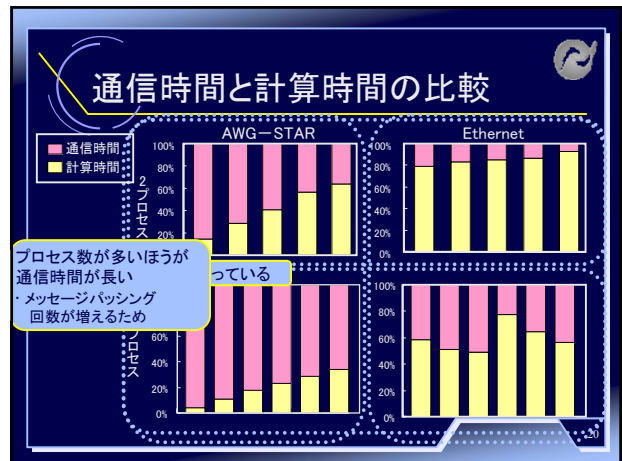
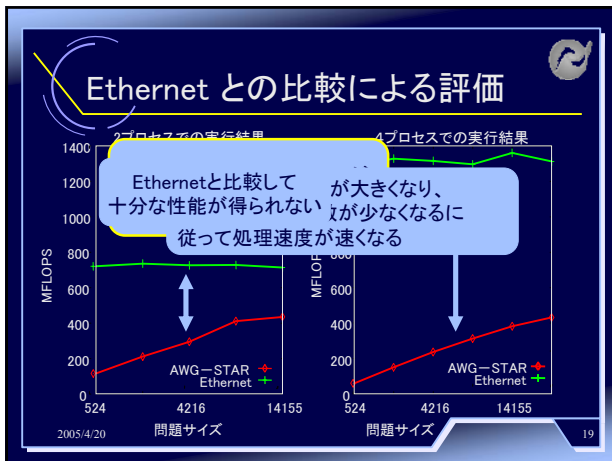
共有メモリ

2005/4/20 PN研究会 17

アプリケーションによる評価

- 実験環境: ノード計算機最大4台で実行
 - 1ノード計算機上で1プロセス
- 評価アプリケーション: 姫野ベンチマーク
 - ポアソン方程式をヤコビ法で解くときの処理速度を計測するベンチマーク
 - 3次元配列の問題をプロセス数で等分割して並列計算
 - データ交換に MPI を利用
 - 様々な問題サイズを設定できる
 - 問題サイズと交換するデータサイズが比例
 - 問題サイズとデータ交換の回数が反比例

2005/4/20 PN研究会 18



まとめと今後の課題

まとめ

- λ コンピューティング環境に Globus Toolkit を導入した
- λ コンピューティング環境における MPI ライブラリ
の設計、開発、実装、評価を行った
- 共有メモリへのアクセス遅延が影響している

今後の課題

- 共有メモリの効率的な利用方法の検討
 - ↳ 動的なメモリ割り当て
 - ↳ メモリマネージャの開発