

# 省メモリ技術と動的最適化技術による スケーラブル通信ライブラリの開発

九州大学情報基盤研究開発センター

南里 豪志

2011年10月11日

CREST「ポストペタスケール高性能計算に資するシステムソフトウェア技術の創出」領域会議

# 研究のねらい

- ▶ Exa FLOPS級計算機で、スケーラブルな並列アプリケーションの実行を可能とする

スケーラブルな通信ライブラリ実装技術

スケーラブルな並列アプリケーション作成のための  
通信ライブラリ応用技術

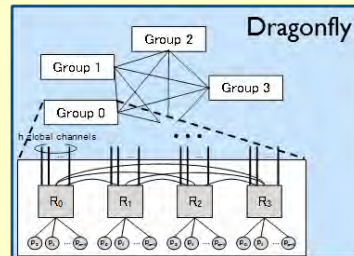
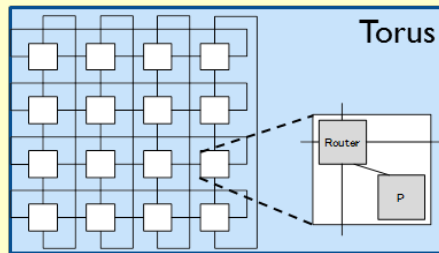
- ▶ 目標：  
通信バッファ領域を総メモリ容量の10%以内に抑え  
ながら、実アプリケーションで数千万～数億プロセス  
までの性能向上を維持する

# 本プロジェクトで想定するExa FLOPS環境

## システム全体

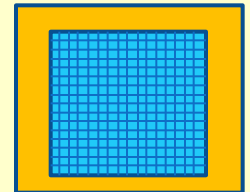
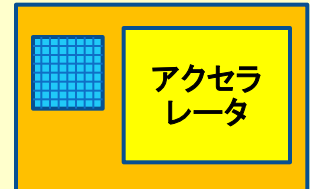
- ▶ ノード数: 数十万ノード
- ▶ インターコネクト:  
多次元トーラス もしくは  
直接網+間接網(Dragon Fly等)

プロセスの配置が性能に大きく影響する



## ノード内

- ▶ コア数:  
数十コア+アクセラレータ  
(heterogeneous)  
  
もしくは  
数百コア~千コア  
(homogeneous)
- ▶ メモリ量: ~数百GB



コア当たりメモリ量の低下

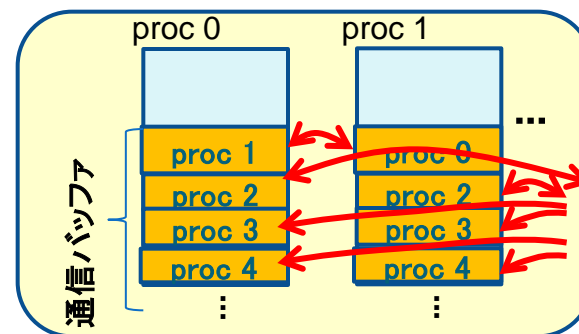
## 並列プログラミング環境

- |                 |                                    |
|-----------------|------------------------------------|
| プロセス数:          | 数千万~数億プロセス                         |
| プロセス当たりメモリ量:    | 1~10GB/プロセス程度                      |
| プログラミングインタフェース: | 高級並列言語(PGAS等)<br>もしくは MPI + OpenMP |

# Exa FLOPSに向けた通信ライブラリの課題

## 1. 数億プロセスに耐えるメモリ管理機構

- ▶ プロセスあたりメモリ容量: 1～10GB程度
- ▶ 従来の通信ライブラリにおける使用メモリ量の問題
  - i) 各プロセスで全プロセスの情報を管理  
= 1プロセス 4Byte でも **400MB/プロセス**
  - ii) 通信相手プロセス毎にバッファを用意  
= 1プロセス 1KBでも、最悪 **100GB/プロセス**



各プロセス用に  
通信バッファを用意

# Exa FLOPSに向けた通信ライブラリの課題

## 2. 数億プロセスにおける性能チューニング

▶ 従来の人手による手法や実行前の静的な手法の限界

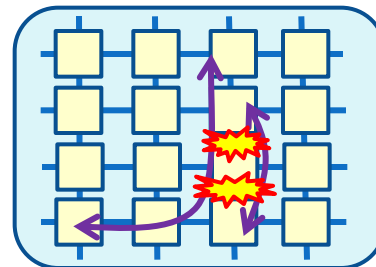
i) 性能最適化の所要時間:

**プロセス数増に伴って最適化の探索空間が爆発的に拡大**

ii) 性能見積もりの複雑化:

**実行時に決定する要素で性能が大きく変動**

例) プロセス配置による通信衝突や通信距離変化



通信衝突

# 本プロジェクトの研究項目

## 課題

使用メモリの削減

実行時の最適化

## 研究項目

(A) スケーラブルな隣接通信及び集団通信のための省メモリアルゴリズム及び動的最適化技術

(九州大学 南里, 森江, 研究員A)

(B) 通信バッファを削減した通信モデルにもとづいた通信プロトコル (富士通 住元, 安島, 三浦, 岡本)

(C) 実行時の状況に応じてパケット送信間隔を動的に制御する通信最適化技術

(九州先端科学技術研究所 柴村, 薄田)

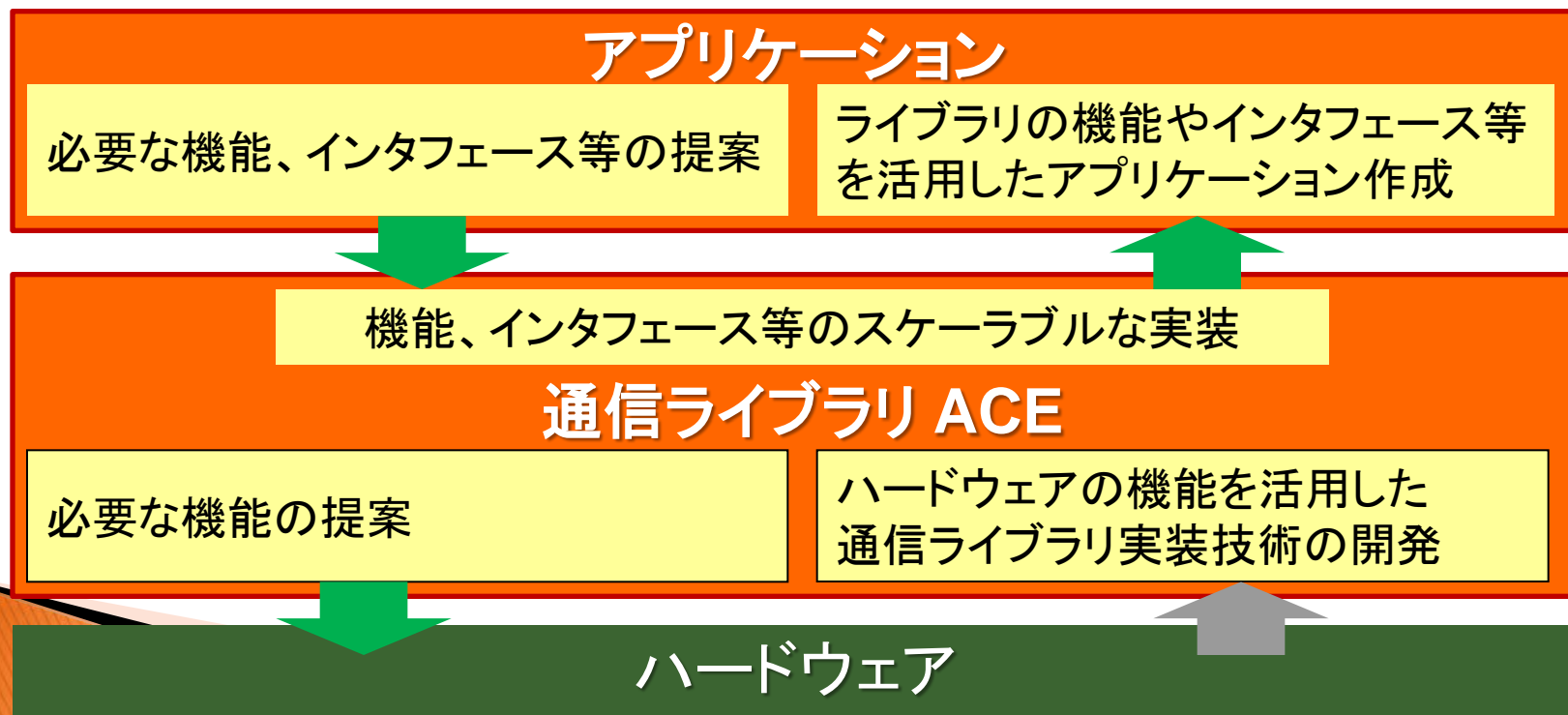
(D) スケーラブルな並列アプリケーション作成のための通信ライブラリ応用技術

(九州大学 高見, 稲富, 深沢, 本田)

スケーラブル通信ライブラリ  
ACE(Advanced Comm.  
lib for Exa)

# 研究のポイント

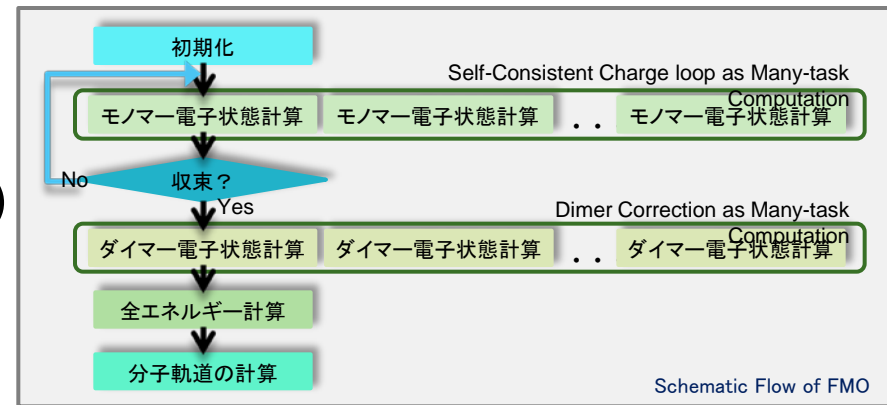
- ▶ 通信ライブラリの低遅延化、高スループット化、および省メモリ化
  - アルゴリズム、通信路制御、プロトコルの開発
- ▶ ハードウェア、通信ライブラリ、アプリケーションの co-design



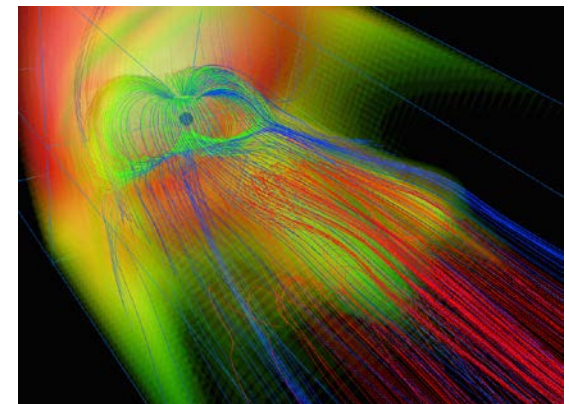


# Exa FLOPSに向けたアプリケーションからの要請

- ▶ 通信を削減、隠蔽、効率化する技術
- ▶ 例) 量子化学計算プログラム OpenFMO
  - 集団通信(Allreduce等)を隠蔽したい
  - 負荷均等化(master-worker)のための排他制御用通信を効率よく行いたい



- ▶ 例) 電磁流体プログラム
  - 領域間の隣接通信を効率よく行いたい



電磁流体シミュレーション結果 8



# ACE

## Advanced Communication library for Exa

既存のアプリケーション:

- ACEの省メモリ技術と動的最適化技術によりスケーラビリティ向上

ACEのインタフェースを活用するアプリケーション

- さらに非ブロッキング集団通信, 隣接通信等の活用で通信を削減・隠蔽・効率化

### スケーラブル通信ライブラリ ACE

通信インタフェース

Put/Get, Send/Recv

隣接通信(ブロッキング, 非ブロッキング)

集団通信(ブロッキング, 非ブロッキング)

遠隔Atomic通信

実装技術

省メモリと動的最適化による集団/隣接通信実装

パケット送信間隔の動的制御

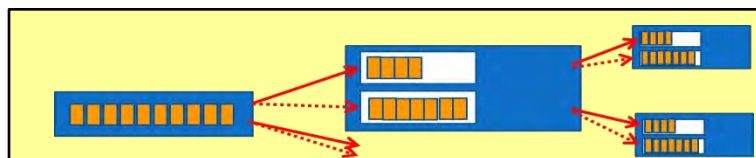
省メモリ通信プロトコル

ネットワークハードウェア

(RDMA, 遠隔Atomic操作, パケット送信間隔調整等)

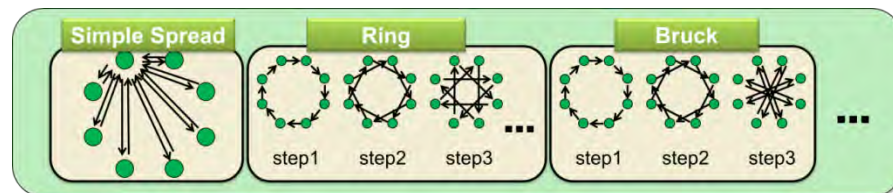
# 動的最適化技術と省メモリ技術による 通信インタフェースの実装

- ▶ ランク配置に応じた隣接通信
- ▶ 省メモリ集団通信アルゴリズム



共有バッファを用いた木構造  
による省メモリアルゴリズム

- ▶ スケーラブルな動的最適化技術
  - 集団通信、隣接通信のアルゴリズム選択、経路選択、セグメントサイズ調整

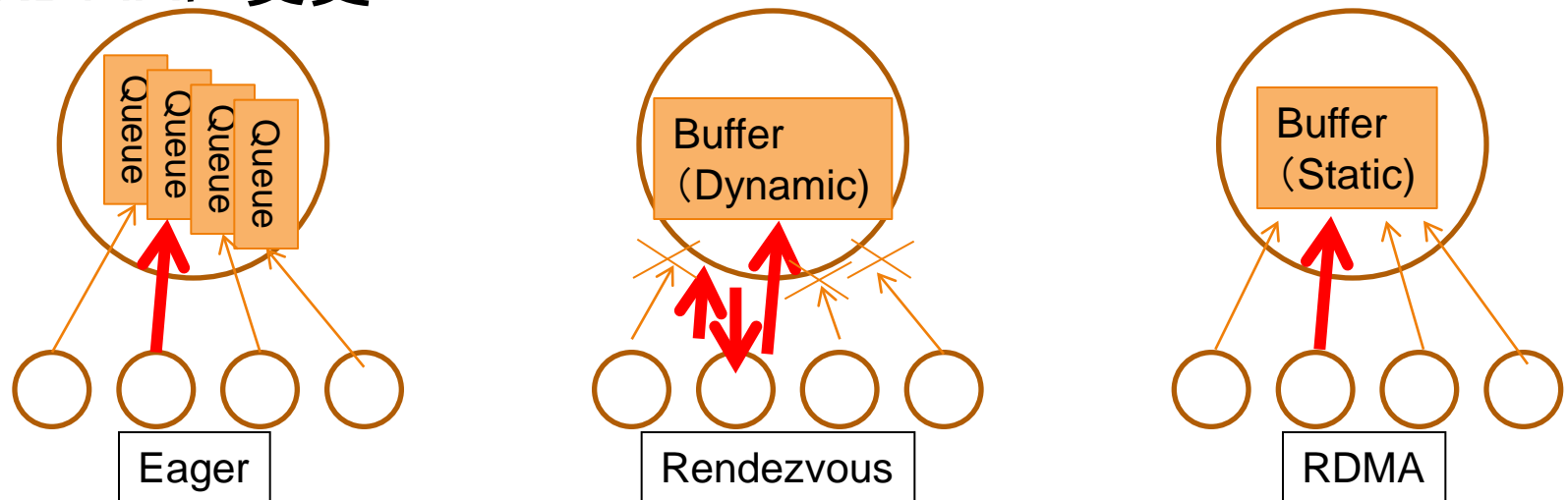


適切な集団通信アルゴリズムの選択

- 非ブロッキング通信の動的最適化技術
- 動的最適化機構の低オーバーヘッド化、省メモリ化
  - ・ 静的な情報の活用、分散型の管理機構

# 省メモリ通信プロトコル

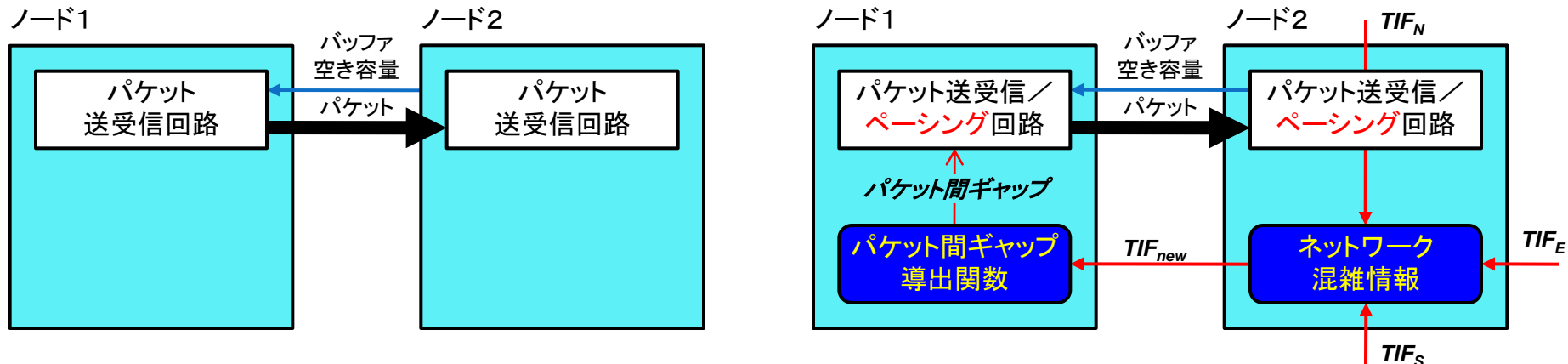
- ▶ 基本通信プロトコル: メモリ使用量の多い Eager、遅延の大きい Rendezvous から、省メモリ・低遅延の RDMAに変更



- ▶ RDMA, Atomic操作によるリモート並列アルゴリズム・ライブラリを整備
  - 片方向通信で、リモートノードの資源を動的利用
    - Mutex, Lock-Free List, Lock-Free Queue

# 動的パケット送信間隔制御技術

- ▶ ネットワークの混雑情報を収集し, パケット送信間隔を調整 ⇒ 通信スループット最適化
  - TIF(Traffic Information Flit: 経路上の混雑情報通知用フリット)による情報収集技術
  - 最適なパケット送信間隔の導出関数
    - ・ 特に集団通信, 隣接通信等のパターン通信



# スケーラブルな アプリケーション作成技術

- ▶ 量子化学計算プログラム OpenFMO
  - 非ブロッキング集団通信による通信時間隠蔽
  - 遠隔Atomic通信による効率的な排他制御
- ▶ 電磁流体プログラム
  - ランク配置に応じた動的最適化を行う非ブロッキング隣接通信による通信の隠蔽と効率化
  - 複数の領域をまとめることによる通信回数削減
- ▶ その他、幅広いアプリケーションに応用
- ▶ 通信ライブラリへのフィードバック

# 全体スケジュール



# 本研究の波及効果

- ▶ 高級並列言語処理系との連携
  - 省メモリ、動的最適化技術： 下位通信レイヤのスケラブルな実装
  - アプリ作成技術： コンパイラの自動最適化機能への活用
- ▶ MPI規格への寄与
  - アプリケーションによる評価にもとづいた、実用的なインタフェースの提案