

# 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:

南里豪志、高見利也、本田宏明、薄田竜太郎、  
森江善之、小林泰三

富士通株式会社:

住元真司、安島雄一郎、志田直之、佐賀一繁、野瀬貴史

公益財団法人九州先端科学技術研究所:

柴村英智、曾我武史

京都大学:

深沢圭一郎

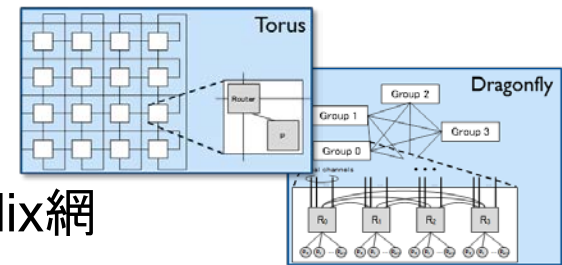
2015年10月15日

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

# 研究の狙い

- エクサスケール計算環境での利用に耐える  
「スケーラブルな通信ライブラリ」  
およびそれを活用したアプリケーションの開発
- 想定するエクサスケール計算環境：
  - 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 high-radix 網
  - プロセス数：数百万
  - メモリ量：10GB程度 / プロセス
- 通信ライブラリ、およびアプリケーションへの要求：  
通信性能と省メモリの両立による高スケーラビリティの達成

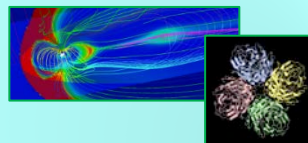


# 通信ライブラリ ACP

## (Advanced Communication Primitives)

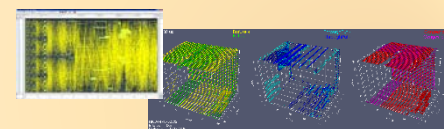
### Applications

ACPライブラリを活用した  
スケーラブルなアプリケーション



### Performance Estimation Tool

ACP向け通信性能  
予測ツール



### ACP Library

#### Middle Layer

コミュニケーションライブラリ  
通信チャンネルを明示的に生成、  
破棄する通信関数群

データライブラリ  
複数プロセスに分散配置  
したデータ構造を操作す  
る関数群

#### Basic Layer

#### 基本層

PGAS(Partitioned Global Address Space)モデルに基づく、通信抽象化層

インターコネクトネットワーク (InfiniBand, Tofu, Ethernet)

# Basic Layer

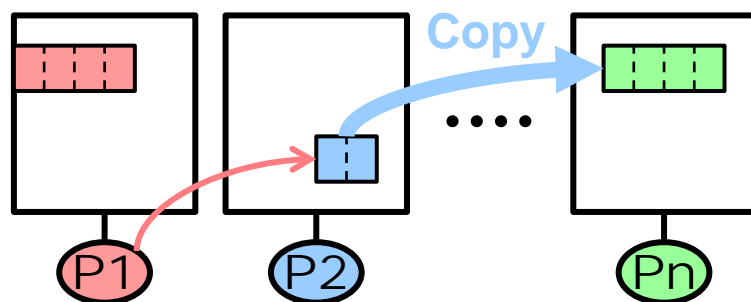
- グローバルメモリ管理

- 各プロセスが登録したメモリ空間を 64bit グローバルアドレスで参照
- 静的領域: 初期化時に登録され、全プロセスから参照可能な領域
- 動的領域: 各プロセスで実行時に登録された領域

- グローバルメモリアクセス

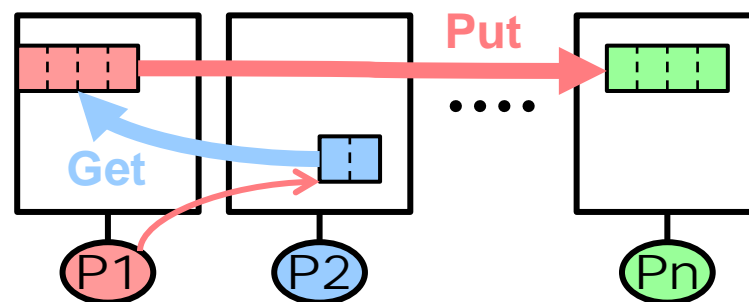
- 任意のグローバルアドレスに対するメモリアクセス (copy, atomic)
- メモリアクセス間の順序関係を記述可能

## GMA



Copy Request

## RMA



Get Request

# サンプルコード

```
#include <stdlib.h>
#include <acp.h>
#define BUF_SIZE 16777216

int main(int argc, char** argv)
{
    volatile acp_ga_t* buf_ga;
    acp_init(&argc, &argv);

    acp_ga_t top = acp_query_starter_ga(acp_rank());
    buf_ga = (volatile acp_ga_t*)acp_query_address(top);

    void* buf = malloc(BUF_SIZE);
    acp_atkey_t key = acp_register_memory(buf, BUF_SIZE, 0);

    *buf_ga = acp_query_ga(key, buf);
    acp_sync();
    /* ... */
    acp_finalize();
    return 0;
}
```

静的変数のグローバルアドレスを取得

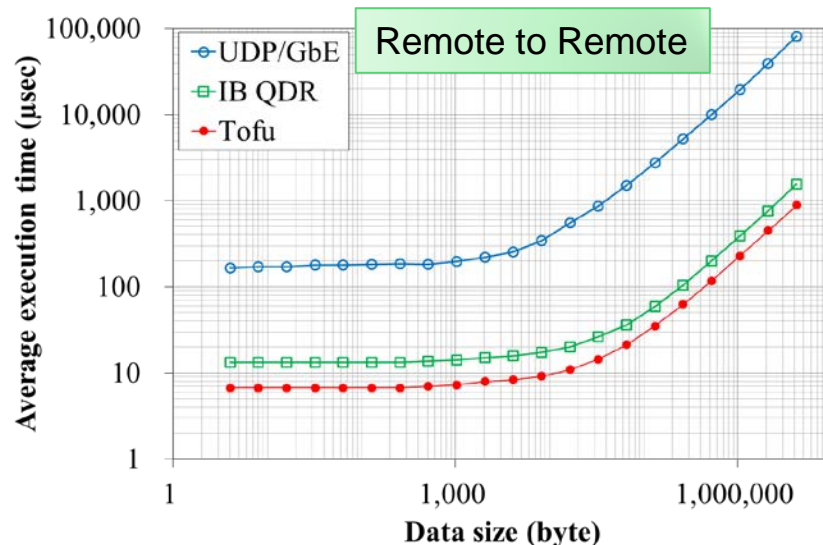
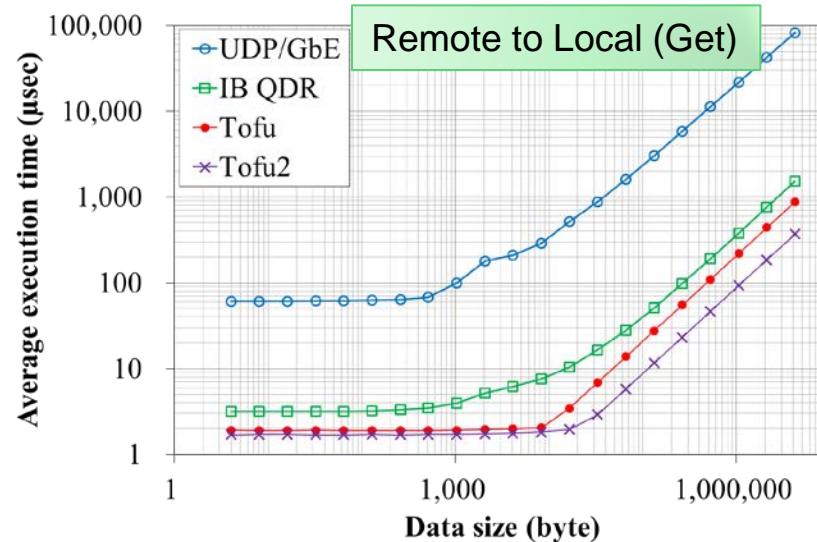
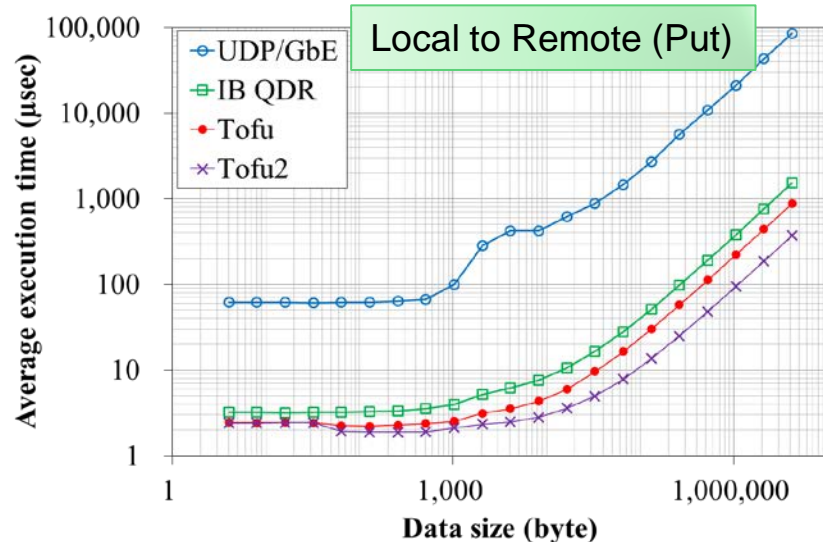
```
acp_ga_t top = acp_query_starter_ga(acp_rank());
buf_ga = (volatile acp_ga_t*)acp_query_address(top);
```

```
void* buf = malloc(BUF_SIZE);
acp_atkey_t key = acp_register_memory(buf, BUF_SIZE, 0);
```

動的に確保したローカルメモリを登録

# Basic Layerの基本性能

## • UDP、Tofu、InfiniBandで稼働



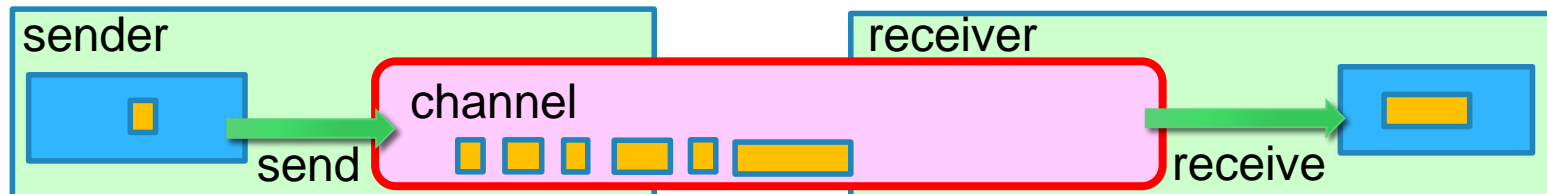
- 各インターコネクトで、ハード性能に近い性能を達成
- 小サイズでの遅延時間については、現在、チューニング中
- Tofu2 の Remote to Remoteは未計測

# Basic Layer メモリ使用量見積もり (100万プロセス)

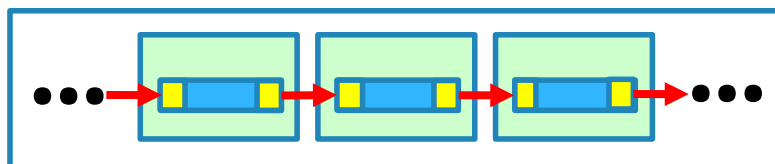
	InfiniBand	Tofu	UDP
リンク数 に比例	176MB @ 100万プロセス <ul style="list-style-type: none"> <li>登録メモリ領域管理 4B</li> <li>QueuePair 160B</li> <li>スターターメモリ情報 12B</li> </ul>	69MB @ 100万プロセス <ul style="list-style-type: none"> <li>コマンド受信バッファ 64B</li> <li>Tofuアドレステーブル 4B</li> <li>Tofu経路テーブル 1B</li> </ul>	18MB @ 100万プロセス <ul style="list-style-type: none"> <li>待ち受けIPアドレス 4B、 待ち受けポート番号 2B</li> <li>送信シーケンス番号 4B、 受信シーケンス番号 4B</li> <li>リンク番号 4B (初期化、リ セット時使用、常時確保)</li> </ul>
メモリ登 録数に 比例		9KB @ 128 登録 <ul style="list-style-type: none"> <li>登録メモリ管理テーブル 40B</li> <li>論理アドレス検索テーブル 16B (比例) + 8 bytes × 256 (固定)</li> </ul>	
その他	2MB <ul style="list-style-type: none"> <li>コマンドキュー/バッファ 960KB</li> <li>登録メモリ領域テーブル 88B × 255</li> <li>登録メモリキャッシュ 10KB × 1024</li> <li>CompleteQueue 128B</li> </ul>	262KB <ul style="list-style-type: none"> <li>コマンドキュー兼コマンド送信 バッファ 64B × 4,096</li> </ul>	647KB <ul style="list-style-type: none"> <li>コマンドキュー 80B × 4,096</li> <li>コマンドステーション 1,504B × 64</li> <li>デリゲートステーション 3,480 × 64</li> </ul>
合計	約178MB	約70MB	約19MB

# 通信チャネルインタフェース

- 2プロセス間のSend/Recv通信



- 特徴:
  - 片方向  $\Rightarrow$  役割 (sender or receiver) に応じて必要最小限のバッファ確保が可能
  - in order  $\Rightarrow$  MPIのタグのような out-of-order のメッセージ管理が不要
- プログラム例
  - 片方向一次元シフト通信



```
ch0 = acp_create_ch(left, myrank);
ch1 = acp_create_ch(myrank, right);
```

チャネル生成

```
for (...){
    req0 = acp_nbseend(ch0, addr0, size);
    req1 = acp_nbrecev(ch1, addr1, size);
    acp_wait_ch(req0);
    acp_wait_ch(req1);
    calc();
}
```

通信

計算

```
req0 = acp_nbfree_ch(ch0);
req1 = acp_nbfree_ch(ch1);
acp_wait_ch(req0);
acp_wait_ch(req1);
```

チャネル解放



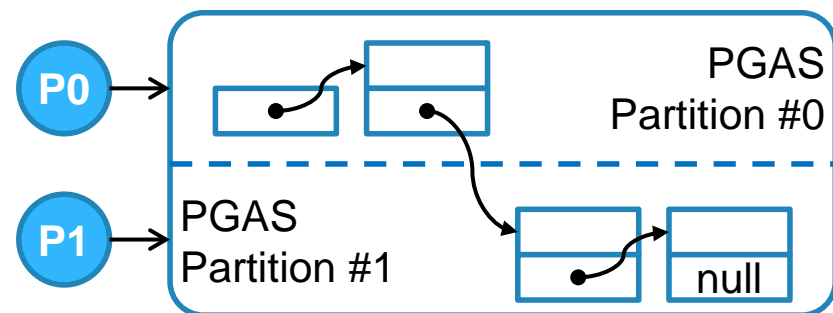
# データライブラリ

- データ構造を複数プロセスに分散
  - データ生成時に配置を明示的に指定
- データの生成, 操作, 破棄は非同期
  - 計算と通信のオーバーラップを促進する
- データ構造の型

- vector      可変長一次元配列
- list        双方向リンクリスト
- deque      双方向キュー
- map        連想配列
- set        集合

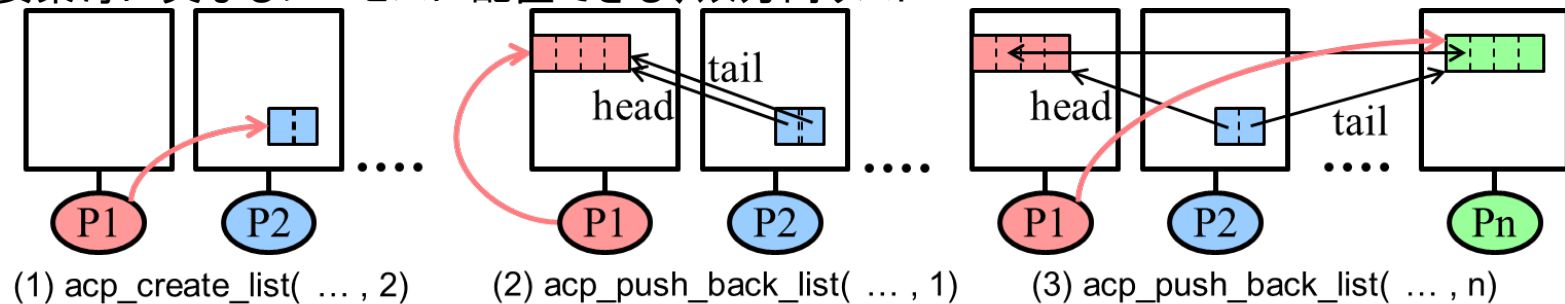
C++言語の標準テンプレートライブラリ(STL)を参考

- グローバルメモリアロケータ
  - データライブラリの基盤技術
  - 指定したプロセスに動的にメモリ割り当て



# API例

- メモリアロケータ関数
  - データの生成、破棄で内部的に使用、ユーザーも使用可能
- ベクタ関数
  - 単一プロセスに配置される、可変長配列
- リスト関数
  - 要素毎に異なるプロセスに配置できる、双方向リスト



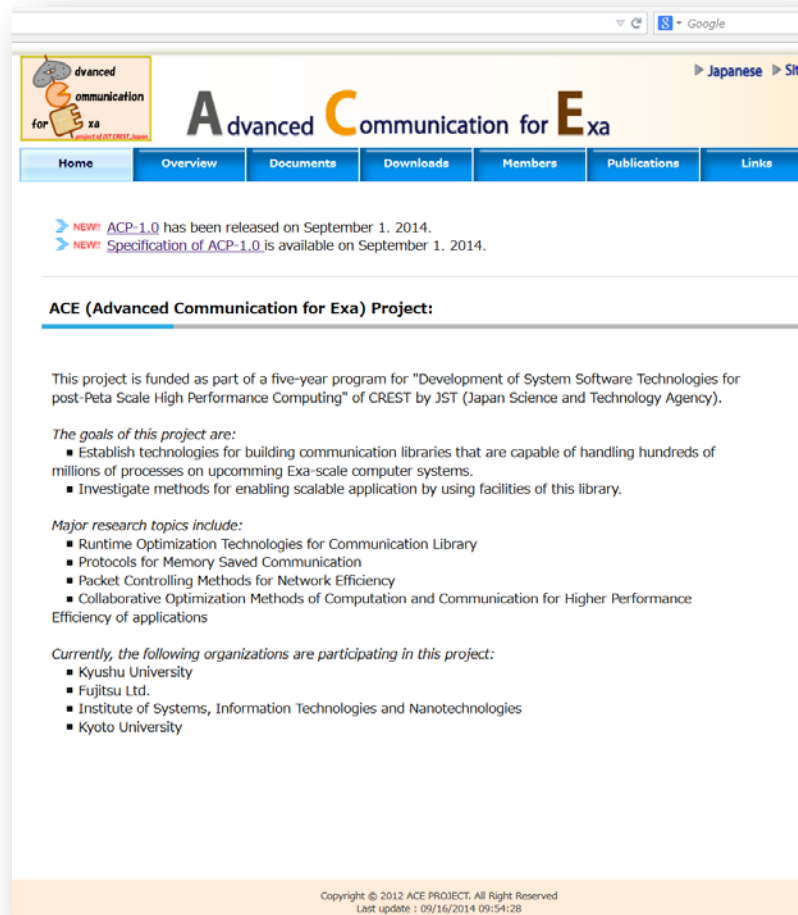
名称	定義
メモリ割当	<code>acp_ga_t <b>acp_malloc</b>(size_t size, int rank);</code>
メモリ解放	<code>void <b>acp_free</b>(acp_ga_t ga);</code>
ベクタ生成	<code>acp_vector_t <b>acp_create_vector</b>(size_t nelem, size_t elsize, int rank);</code>
ベクタ末尾要素追加	<code>void <b>acp_push_back_vector</b>(acp_vector_t vector, acp_ga_t ga);</code>
リスト生成	<code>acp_list_t <b>acp_create_list</b>(size_t elsize, int rank);</code>
リスト先頭要素追加	<code>void <b>acp_push_front_list</b>(acp_list_t list, acp_ga_t ga, int rank);</code>

# ACPライブラリ公開

- Webサイト

<http://ace-project.kyushu-u.ac.jp>

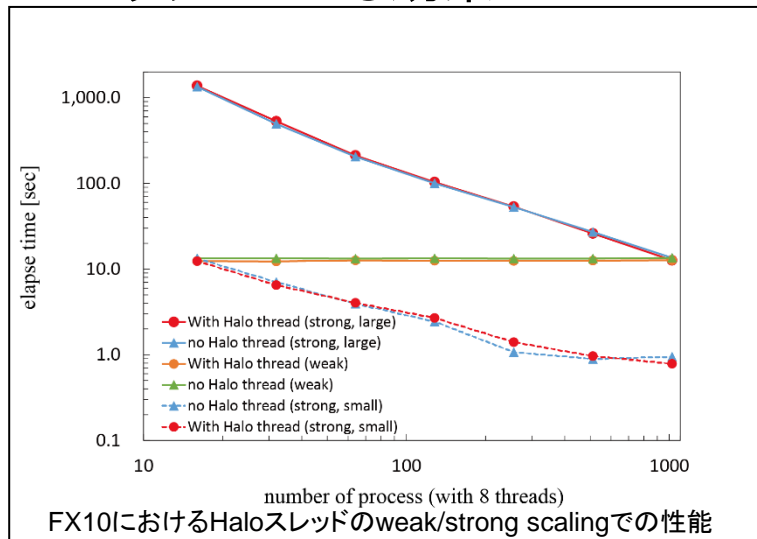
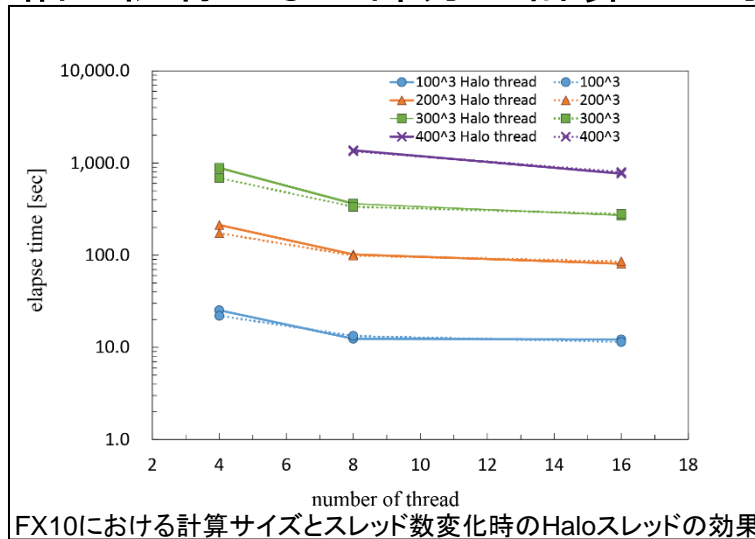
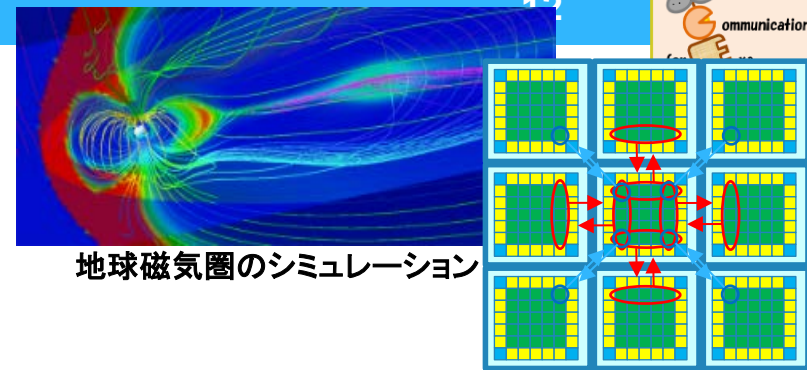
- V1.1.1 公開中



# 電磁流体シミュレーション のStencil計算

- Haloスレッド

- Stencil計算における袖通信、およびその通信に依存する計算を担当するスレッド
- 袖に依存しない部分の計算とのオーバーラップによる効果

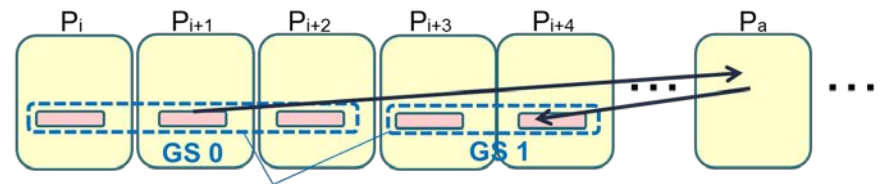
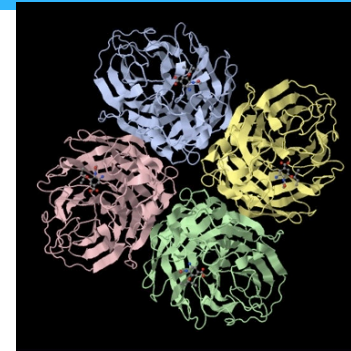


- Halo通信のフレームワーク化

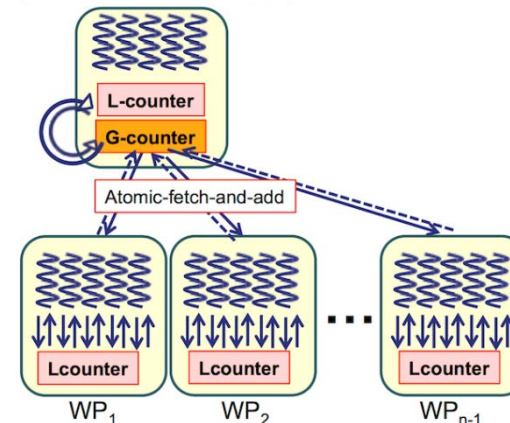
- 通信と計算のオーバーラップを簡単に記述
- ACPによる省メモリ実装

# FMO計算

- FMO計算: master / worker モデル
  - 大規模分子の量子科学計算をタスク並列処理
- ACPによる実装
  - 共有ワークスペース
    - 全プロセスが非同期アクセス可能な連続メモリ領域
    - 1ノードに納まらない大規模連続領域を用いた計算が可能
  - グローバルカウンタ
    - 動的負荷分散に使用
    - Remote Atomic Fetch and Add による実装

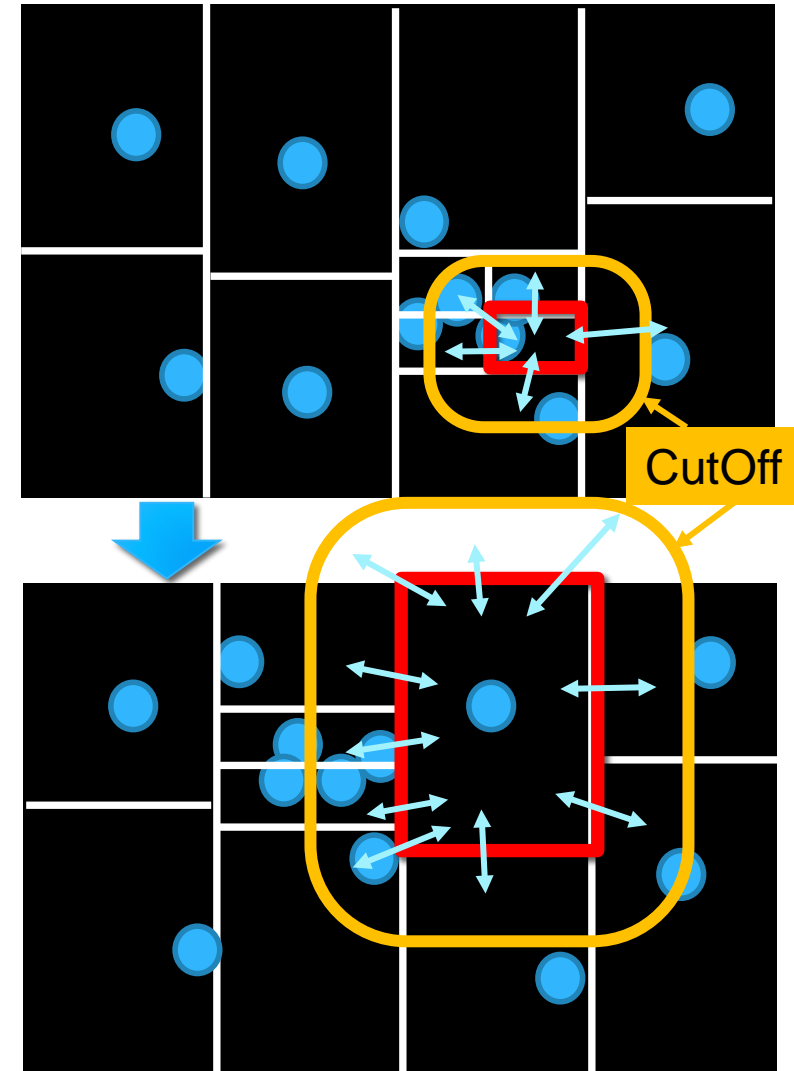


Worker-Proc 0 (WP 0)  
(グローバルカウンタを保持)



# 重力N体シミュレーション

- Adaptive P<sup>3</sup>M法
  - 粒子分布に応じた領域分割  
⇒ 非構造格子
  - 不規則な通信
    - 粒子の移動
    - 領域分割の変更
- ACPによる実装
  - 通信用領域の動的確保、解放
  - グローバルメモリアクセス (copy, atomic) による効率的な通信



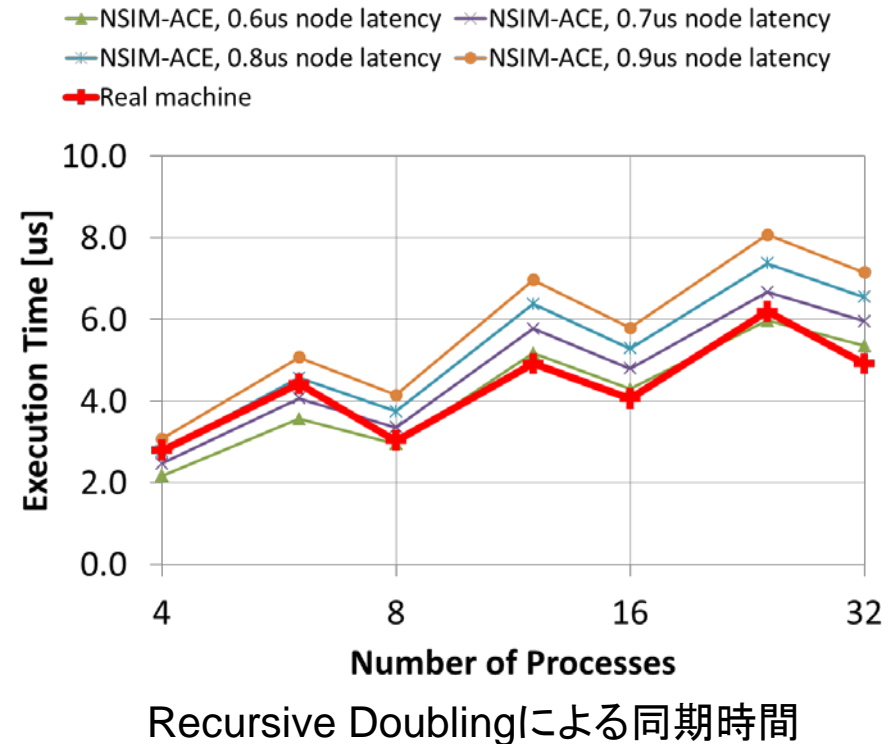
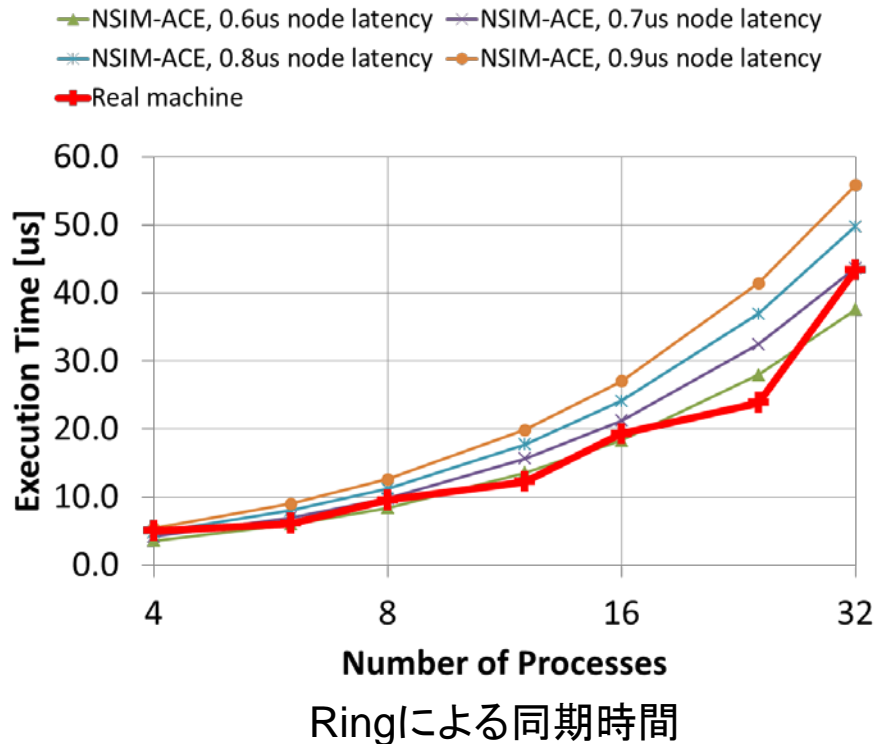
# 性能予測・解析環境 NSIM-ACE

- NSIM-ACE = NSIM + RDMAシミュレーション機能
  - メッセージパッシング／RDMAの両通信モデルをサポート
  - トポロジや通信パターンによる通信衝突を踏まえた通信時間推定
- エクサスケール計算環境をターゲットとした性能予測
  - ACPアプリケーションやACPライブラリの開発・評価に有効
    - 通信モデル比較(MPI vs RDMA)
    - ライブラリ内で用いるアルゴリズム選定
    - スケーラビリティ評価
- NSIM-ACEの開発
  - RDMAモデルに基づいた性能予測機構の実装(完了)
  - ACP準拠のAPI関数の実装(一部完了、進行中)
  - シミュレーション精度の検証(実施中)

# NSIM-ACEのシミュレーション精度検証

## 同期通信の精度検証

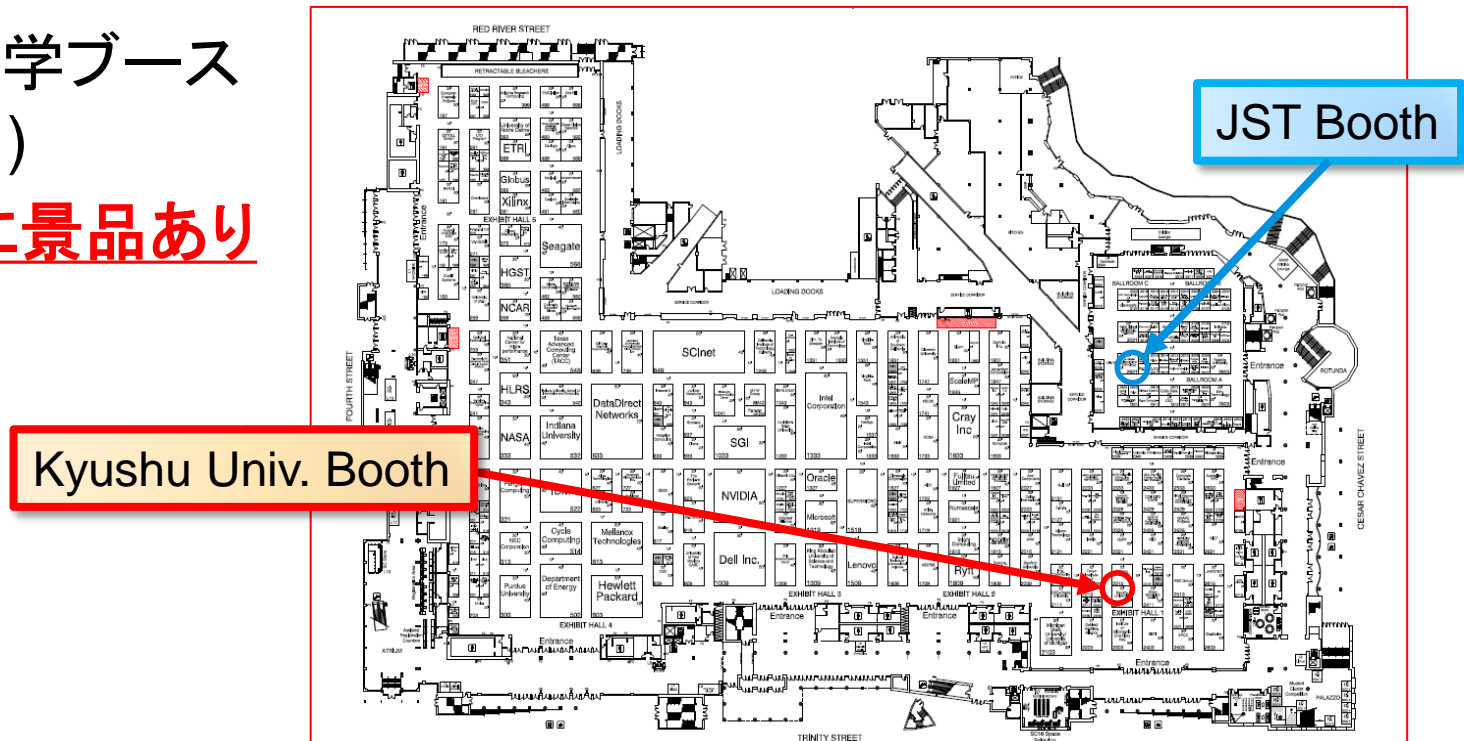
- 実機のネットワークノード仕様をNSIM-ACEに設定
- 実機で観測困難なパラメータは較正により補正
- 実用的な精度でのシミュレーション結果(予測時間、振る舞い)を確認





# ACPチュートリアル @ SC15

- 内容
  - ACPライブラリのインタフェース紹介
  - 実機を使ったプログラミング実習
- 場所:
  - 九州大学ブース  
(#2311)
- 参加者に景品あり



# まとめ

- プロジェクトの目的:  
省メモリと通信性能を両立するスケーラブルな通信ライブラリ開発
- 進捗状況:
  - ACPライブラリ開発、公開
  - アプリケーション開発
  - シミュレータ NSIM-ACE開発
- 今後
  - インタフェースの充実化
  - アプリケーションへの応用、評価、およびフィードバック
  - ACPを用いた上位レイヤの構築
    - DSL、MPI