

# 富岳上での **FDPS** の性能デバッグ

牧野淳一郎  
神戸大学

**2021/10/26** 富岳計算宇宙惑星内部ワークショップ

# 話の構成

1. 問題はなにか:**FDPS** を使ったプログラムの富岳上での性能と安定性
2. 測定結果: 何が遅いか?
3. 問題と対応
4. **GPLUM** における問題と対応
5. まとめ

# 問題はなにか:**FDPS** を使ったプログラムの富岳上での性能と安定性

**FDPS** を使って富岳で動いている/動いて欲しいプログラムは沢山ある。

**ASURA/FDPS**、**GPLUM**、**FDPS-SPH**、**PETAR**、あと他にも。あと神戸大学/理研の大石先生のところで防災シミュレーションで**SPH**。多分**DEM**とかも。

どうも色々問題があってあんまり使えてない。

- 遅い。
- しばらく回すと落ちる。

というわけで、「ある程度直した」という話を今日は。

# 測定

- 測定に使ったコード: **nbody-with-center**  
(<https://github.com/jmakino/nbody-with-center>) 牧野が書いた、惑星リング用のコード。書いたというか **N**体サンプルプログラムをちょっといじって中心天体重力と粒子同士の非弾性衝突をいれただけ。
- **FDPS** の領域分割とツリーは円筒座標で作るモードで動作。
- とりあえず簡単なコードで様子を見るため。
- ノード数: **1k-4k**
- 粒子数: **2-200M**
- 見込み角  $\theta = 0.5$ 、リング半径 **1**、幅=**0.1**、**4**ステップに**1**回領域分割(いらないけど時間測定のため)

# 最初の測定結果

## 1 ステップの計算時間

項目	時間(ミリ秒)
全体	<b>105</b>
領域分割	<b>1.5</b>
粒子交換	<b>10.5</b>
相互作用計算	<b>54</b>
「その他」	<b>38</b>

- とにかくものすごく遅い。ハードウェアの性能がちゃんとでたら**1**ミリ秒のオーダー、まあ通信が××だとしても **10**ミリ秒くらいになって欲しい。
- この測定はアプリケーションプログラム側で、あんまりバリアとかちゃんといれてないので例えば「その他」はなにかとか怪しい。

# とりあえずいれた修正

- **LET**(相互作用計算に必要な他のノードの情報。 **Local Essential Tree**) の交換に **MPI\_ALLTOALLV** を使う標準モードの代わりに **P2P\_EXACT** に切り替え。これは、領域のトップレベルの情報は **ALLGATHER** で共有して、それより細かい情報が必要なところだけ **P2P** で交換する
- 相互作用計算を **SIMD** 化。あんまり気合いはいってなくて最内側ループが連続アクセスでベクトル化されるように成分毎の配列にコピーしてから計算するくらい。

# 書換え後の性能

## External times

項目	時間 (ミリ秒)
<b>Total</b>	<b>59</b>
領域分割	<b>1.7</b>
粒子交換	<b>8.4</b>
相互作用計算	<b>44</b>
その他	<b>4.6</b>

## Internal times

項目	時間 (ミリ秒)
領域分割	<b>3.6</b>
粒子交換	<b>22</b>
プロセス内/全体ツリー構成	<b>9.4/9.5</b>
相互作用計算	<b>3.2</b>
<b>LET</b> 構築/通信	<b>6.5/15</b>

**SIMD** 化で相互作用計算カーネル「だけは」すごく速くなった。

# 問題と対応

- 粒子交換
- **LET** 構築/通信
- ツリー構築
- 領域分割



# 粒子交換

これは、粒子が動いた/領域分割が変わったことで、属するべきところに送るところ。

- リング用円筒座標コードで、粒子はそんなに遠くまでいかないはず。
- 実際に通信している相手を見たら何故かほぼ全部になっていた。角度方向で両隣以外は、実際にはそこまで動く粒子がないがチェックしていた。
- これは **FDPS** の性能バグで、(この円筒座標用の) 周期境界での粒子交換アルゴリズムが間違っていた。修正した。

修正後の性能:

**22ms → 2.2ms**

# LET 構築/通信

- P2P\_EXACT はリング向けではあんまり、、、
- P2P\_FAST に変更。

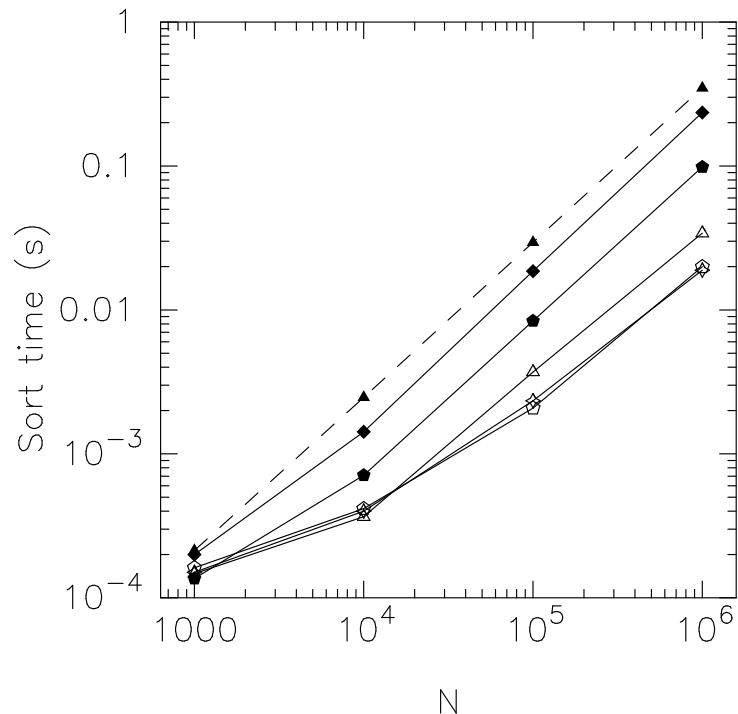
**LET 構築時間:** 6.5ms → 3.5ms

**LET 通信時間:** 15ms → 0.3ms

# ツリー構築

- **FDPS** では、ツリーはモートンキー作ってからソートする。**OpenMP** で性能だすために並列マージソートを使っていたが、これが富岳では非常に遅くなっていた。
- とりあえず **1** コアで **std::sort** にしてみた  
プロセス内ツリー構築時間: **9.4ms** → **0.56ms**
- あと、並列サンプルソートも実装した。  
プロセス内ツリー構築時間: **0.56ms** → **0.44ms**  
(**1** ノードあたりの粒子数が多いとちゃんと **20** 倍くらい速くなる)

# 並列サンプルソートの性能



**160** バイトのデータを **64** ビットキーでソートする時間 (経過時間)

破線+黒三角は **1** コア **std::sort**。黒四角、五角、中抜き三角、四角、五角は **2**、**4**、**12**、**24**、**48** スレッドでの並列サンプルソート。

**24** スレッドまではいいが **48** スレッドは、、、どこが遅いかとかもうちょっと調べて性能あげられるかもしれない。

サンプルソートライブラリとして <https://github.com/jmakino/sortlib> で公開

# 領域分割

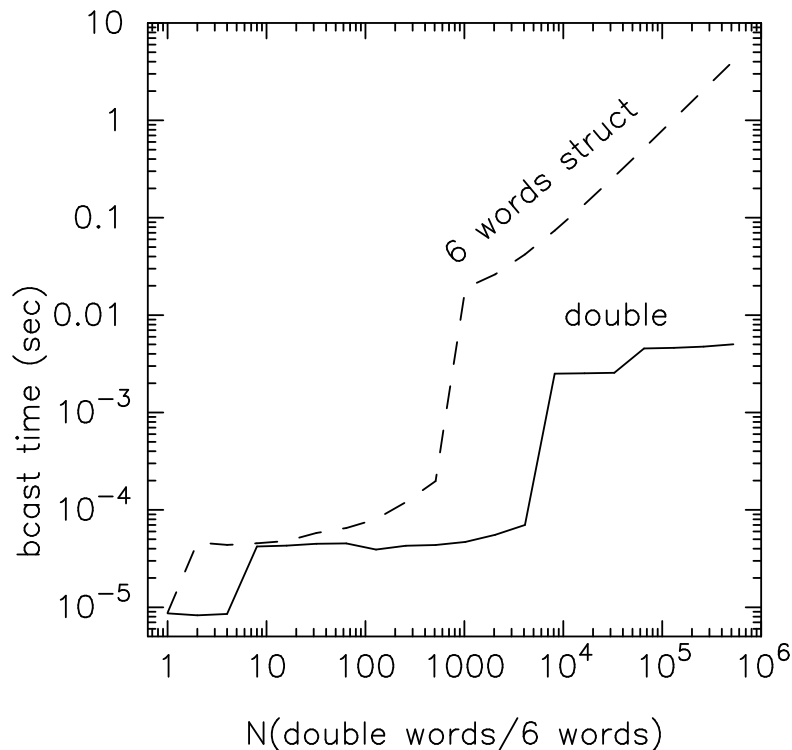
領域分割の操作毎にバリアいれて時間測った結果わかったこと:

- サンプル粒子あつめるところ (何故か **MPI\_AllgatherV**) は十分速い。
- ランク **0** での計算はそこそこ。
- 作った領域分割の放送が滅茶苦茶に遅い。単に **MPI\_Bcast 1** 度呼ぶだけでデータ量も上のサンプル粒子もらうところより桁で少ないが時間は桁で遅い。 **MPI** のバグとしか思えない (詳しくは細野君の発表で)

とりあえずの対応: **MPI\_Bcast** やめて全プロセスで同じ計算するようにした。あと、計算ではソートは並列サンプルソートに変えて、他にも **OpenMP** 化さぼったところをちゃんといれた。

領域分割時間:     8ms → 0.4ms

# 富岳 MPI\_Bcast の性能 (詳しくは細野君発表で)



**1024** ノード、 **48** スレッド/プロセス

**6 words**= 倍精度 **6** 語の構造体を送る

まともな **MPI** 実装なら、構造体  $N$  個送  
るのは  $6N$  語送るのと同じ時間ですむは  
ず。あと、必要時間がジャンプしないよ  
うに適切なアルゴリズムが選ばれるはず。

はずなんだけどねえ、、、

# 現状

項目	時間(ミリ秒)
全体	<b>7.2</b>
領域分割	<b>0.41</b>
粒子交換	<b>0.78</b>
相互作用計算	<b>5.8</b>
ツリー構築	<b>1</b>
<b>LET</b> 構築	<b>1.6</b>
<b>LET</b> 交換	<b>0.4</b>
ツリー辿る	<b>0.44</b>
相互作用カーネル	<b>1.57</b>
その他	<b>0.1</b>

大分速くなった。

**200** 万粒子に **1024** ノードは実際上は無理で、**72** ノード **9.4ms**、**144** ノード **8.35ms** あたりがまあ使ってもいいくらい。

**144** ノード以下だと **1** ノード **4** プロセスのほうが速い。

**2G** 粒子 **1152** ノードで **0.85** 秒。この辺で相互作用カーネルの効率が低いのがみえてくる。

# FDPS としての改善の余地

- 領域分割には  $O(p^{1/3})$  ( $p$ : プロセス数) のアルゴリズムがあるのでそっちを発掘して切り替える。
- 粒子交換: 粒子を見つけるところがなんか遅いのもうちょっとなんとかする。  
**slow**
- **LET** と 相互作用計算のツリー辿るところが遅いのをなんとかする。これはだいぶ書き直して、**x86** では **1.5** 倍くらい速くなったが富岳ではまだ。



# GPLUM

- **FDPS** 速くなったのでまともに速く動くようになったかということそうでもなく...
- 富岳では少なくとも **12** スレッド/プロセス、大ノード数で性能と安定性を維持するには **48** スレッド/プロセスが望ましいが、その辺想定してなかった結果色々なことが。

# GPLUM—具体的な問題

- 近接粒子リストを **std::vector** を使って作っている。
- さらに、**std::map** を使って粒子 ID からランクやインデックスをひけるようにしている

これらがスレッド数多いとボトルネックになる。あと細かいことが色々。

- **std::vector** は全然別の **vector** をアクセスしてても、サイズが変わると共通のヒープからメモリとってくるので、**push\_back** とかで要素追加する操作はスケールしない
- **std::map** はそもそも要素追加が並列にできない

現行の **GPLUM** では、これらを「なるべく」使わないように修正。**vector** は残ってるかも。100万粒子で1ステップ10ms くらい？

# 色々わかったこと

- 非常に注意深く使えば富岳の通信はかなり速い。
- 注意深く=測定して、理屈に合わない性能のそこはドキュメント読んで理屈にあう性能がでるものに変える。
- 特に **ALLTOALL(V)** は避けるべき。
- あと、**MPI derived data type** は富岳では使ってはいけない。
- 「これ使うと遅いんですが」「では、使わないようにして下さい」的
- 何故かこの作業のあと落ちないで何ステップでも問題なく走るようになった。  
**ALLTOALL(V)** が悪い気がする。

## 色々わかったこと2

- **STL** は便利だがメモリ管理周りは容易に **OpenMP** 環境で性能ボトルネックになるので、そうならないようにする必要あり。使わなければ問題は起きない。
- **OpenMP parallel for** のスレッド起動オーバーヘッドは富岳でもあんまり小さくない。サンプルソートの実行時間に影響があったので、**10** マイクロ秒くらいある？(要調査) 性能だすためにはなるべく **parallel section** の単位を大きくするべき。

# まとめ

- 富岳での **FDPS** の性能の問題点を **1024** ノード **200** 万粒子で詳細測定し、比較的容易に直せるところは直した。
- **SIMD** 化はしたが最適化はできてない相互作用カーネルの実行時間が見えるくらいにやっとなった。まだ実行効率としては低い。
- 惑星リングなら (まだ試してないが) 相互作用リスト再利用モードにすれば相互作用カーネルがメインにできるはず。
- 再利用モード使えない計算だとツリー辿るところをもうちょっとなんとか、、、
- 大雑把にあって **100** ノード超えると **1** ノード **1** プロセスのほうが通信性能的に安全。

- この辺の改良はいった **FDPS 7.1** 公開されました。