

## 汎用超並列計算機の上での大規模数値計算

— QCD\_TARO<sup>(+)</sup> プロジェクトの経験から —

山形大学・教育学部 中村純 (Atsushi NAKAMURA)

### 要旨

ベクトル型のスーパー・コンピュータしか使ったことのなかった素人のユーザーが経験した超並列計算の報告をする。単純な幾何学的分割とメッセージ・パッシングを使ってプログラミングを行なった。並列計算の概要を簡単に与えたのち、我々が遭遇したいくつかの問題（ゲージ系に固有の問題、繰り込み変換に伴う問題等）とその解決の試みを紹介する。現在の並列計算のユーザーから見た問題点を報告する。

### 1. 序

並列計算機が新しいスーパー・コンピュータとして注目を集めている。計算機の専門家の間では、並列コンパイラ等の研究が意欲的に進められているが、実用化までの道はまだ遠いらしい。

しかし、我々のような大規模科学計算を進めているユーザーにとって、魅力的なピーク性能を持ったハードが既に存在しているのに、手をこまねいているのは辛い。ベクトル型スーパーコンピュータの出現によって、シミュレーションが科学を進める道具となることは確立しつつあるが、系のサイズ、統計精度が限られた結果が出始めると、これらを克服して、決定的な結果を出したくなる。<sup>(\*)</sup> 比較的廉価で高いパフォーマンスが得られる並列計算機が素朴なプログラムで使えるものなら、どんどん計算を始めたい。

---

<sup>(+)</sup> QCD on Thousand cell ARray processor for Omnipurposes. 千のセルを持った多目的プロセッサの上でのQCD。QCDはQuantum Chromo-Dynamics (量子色力学) の略で、物質の究極の構成子と現在考えられているクォークとグルオンを支配する力学。

<sup>(\*)</sup> われわれの取り扱うモデルは4次元空間  $(x, y, z, t)$  の離散化に基づいている。今一辺に (わずか) 30点取ったとしても、 $10^6$  の格子点が必要になる。各点にあるグルオン場は、4成分を持つベクトルで、各成分は  $3 \times 3$  複素行列で表される。モンテカルロ計算で一点の更新に100回の浮動少数点演算が現われるとして (実際はもっと多い)、(わずか)  $10^4$  回のモンテカルロ・ステップを実行すると、1000ギガの計算が必要となる。一辺の点の数を倍にすれば、格子点は16倍になる。またモンテカルロ計算の統計誤差を半分にしようとするれば、4倍のステップ数が必要になる。

ここで報告するアプリケーション（格子QCD）は、1つの計算に数千時間かかるので、ピーク性能に少しでも近い値を得ることが重要となる。また、これまでベクトル計算機で行なわれてきた計算が、超並列計算機を使っているからという理由で制限をうけては面白くない。

以上のような背景をもって進んでいる我々のプロジェクトでの経験を報告する。

## 2. 幾何学的分割とメッセージ・パッシング

時空間の関数の微分方程式で表されるモデルなら、時空間を離散化することによって、差分方程式になる。対象の物体をいくつかの部分に分けて、それぞれを並列計算機を構成しているプロセッシング・エレメント（PE）に割り当てる。差分方程式なので、1つの点を計算するのに、その近傍の点が必要になる。境界の点に対しては、近傍の点が同じPEの上に無いことがある（図1）。その時は、必要な点の上での値を send, receive で他のPEから持ってくる。これで並列プログラムができる。専門家にはうけないが、ユーザーにはこれが一番わかりやすい。

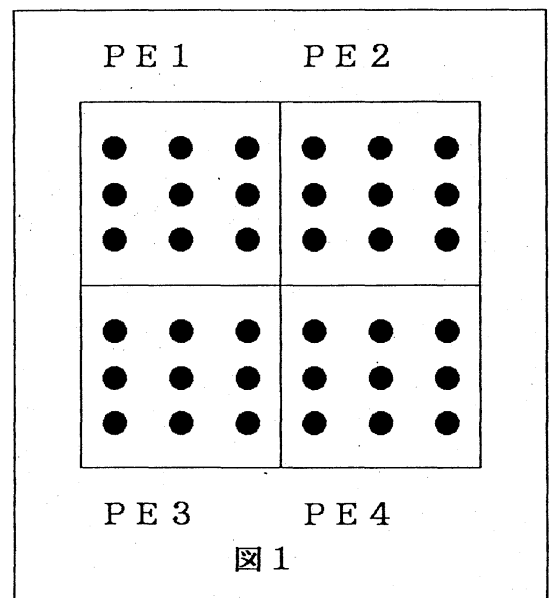


図1

## 3. ゲージ理論

ゲージ理論の特徴は、基本的な変数が点の上ではなく、点を結ぶリンクの上で

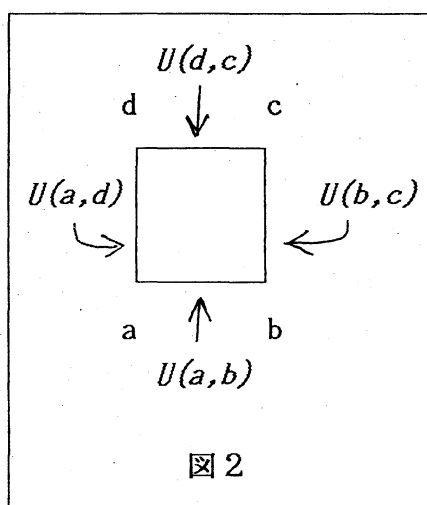


図2

定義されており、

相互作用がリンクをつないでできるループで与えられていることである。図2において、 $a, b, c, d$  は2次元平面上の格子点を表し、 $U(a, b)$ 等が2つの格子点を結ぶリンクの上で定義されたSU(3)行列（ $3 \times 3$ ユニタリ行列で行列式が1のもの）であるとき、相互作用は

$$\text{Tr } U(a, b) U(b, c) U(c, d) + U(d, a)$$

という形をとる。

我々の問題は空間3次元・時間1次元の4次元空間  $(x, y, z, t)$  を離散化したものである。計

算を進めている超並列計算機、富士通AP1000の要素は2次元配列なので、<sup>(+)</sup> 4次元空間を2次元面の上で分割しなければならない。一番簡単なのはxとyの部分分割することである(図3)。尤もらしく「xyパーティション」となす。このときゲージ系の特徴のため、前後左右だけでなく、斜め方向のPE間の通信も必要になる。<sup>(\*)</sup> 現在の世界のトップレベルのスケールの問題

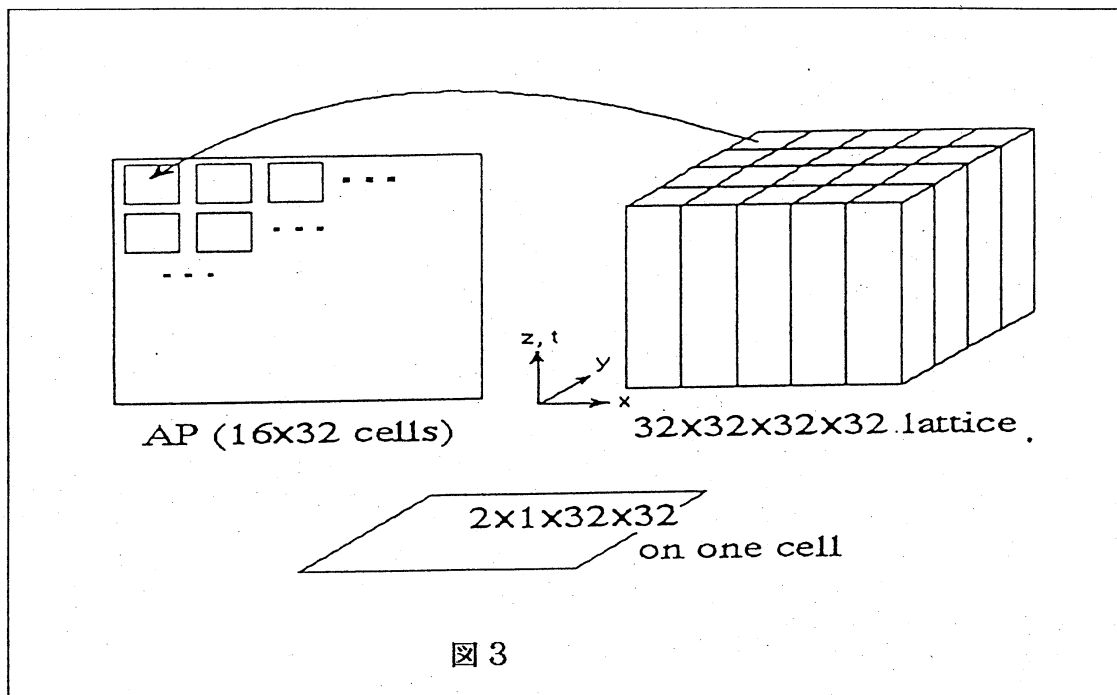


図3

は $(x, y, z, t)$ 方向の点が $32 \times 32 \times 32 \times 32$ のオーダのものである。一方512個のPEの構成は $16 \times 32$ なので、さきほどの「xyパーティション」の分割の方法では、1つのPEの上に行っているのは $2 \times 1 \times 32 \times 32$ となる。

#### 4. 繰り込み変換

繰り込み変換では、近傍の値を平均して新しく疎視化した系をつくっていく(図4)。つまり、 $32 \times 32 \times 32 \times 32$ という系から、 $16 \times 16 \times 16 \times 16$ を作り、更にそれから $8 \times 8 \times 8 \times 8$ という系を作りというように進んでいく。先程の「xyパーティション」では、 $16 \times 16 \times 16 \times 16$ という大きさの系を $16 \times 32$ 個のPEの上に乗せることはできない。そこでz方向

<sup>(+)</sup> 現在N-jigenというライブラリーが用意され、ソフト的に任意次元の構成として扱うことができるようになった。

<sup>(\*)</sup> ループにそって行列の積を計算するさい、まずコの字型の部分から転送を行えば、斜め方向のPE間の通信は必要がなくなる。広島大学・日置慎治氏のご指摘による。

も分割することにして「xyzパーティション」と呼ぶことにした(図5)。  
 まず、全格子をz軸に垂直な方向に8つにスライスする。一方512個の

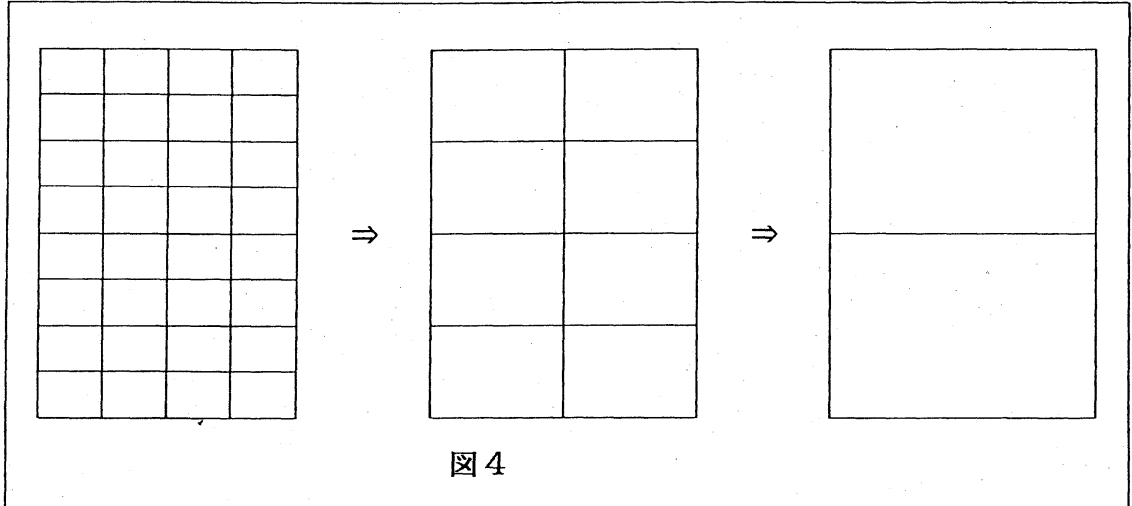


図4

PEを64個ずつのクラスターに分割する。(プログラミングにおける仮想的な分割。ハード的には何も意味しない。)各クラスターをスーパーPEと呼ぶことにする。xyzパーティションでは、PE全体でみると複雑な通信が必要で、xyパーティションの時のように、近傍のPEとの通信だけというようにはいかなくなる。

格子には周期境界条件が課してあるので、 $z = 1$ と $z = z_{Max}$  ( $= 32$ )が

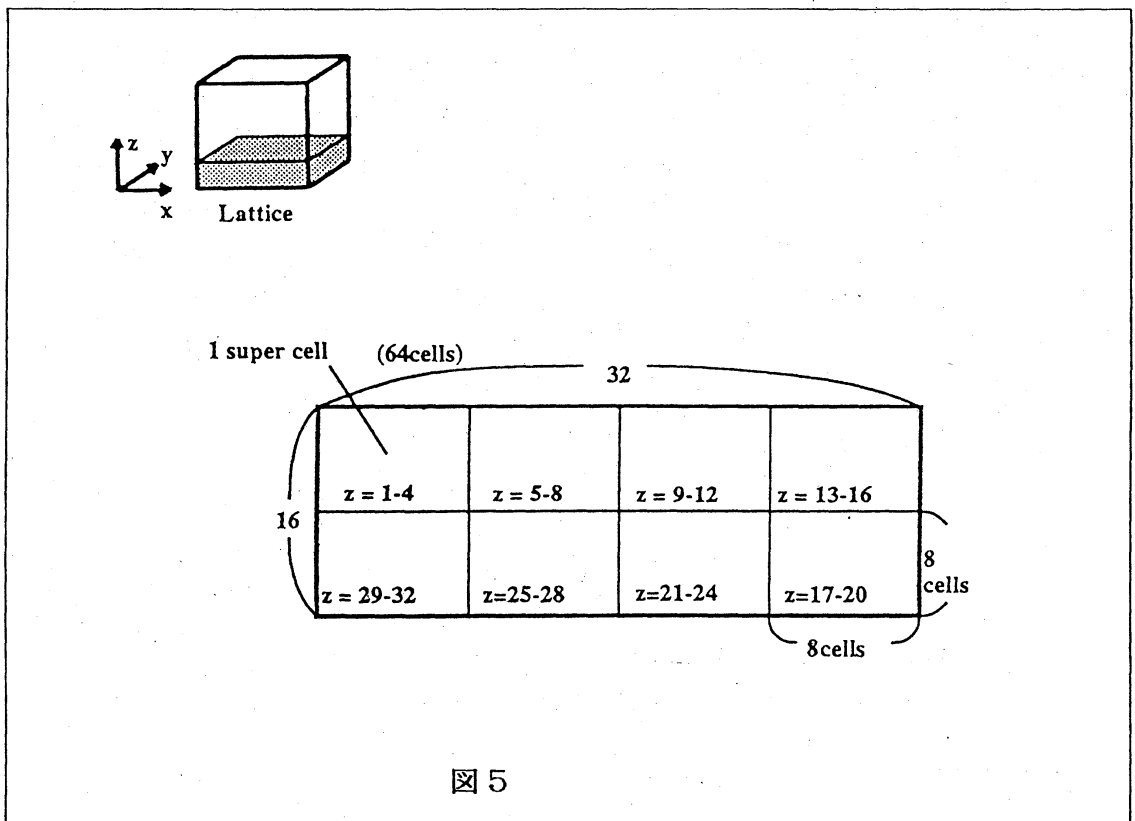


図5

少しでも近くなり、PE間の通信路の距離が小さくなるように、図5のようにスーパーPEを配置した。(\*)

#### 4. 並列計算におけるGATHER/SCATTER

我々のようなタイプの計算における並列の問題というのはいったい何なのかを考えてみる。我々はシステムをPEに分割している。そのため、物理的な系で見たときある点の右隣の点が欲しいといっても、その点は同じPEの上にいるとは限らない。そこで通信の必要が起こる。言い換えると、系全体を表すグローバルな変数(X, Y, Z, T)を使ったとき、例えばy方向の隣の点は(X, Y+1, Z, T)とすぐ分かるのだが、PEの上のローカルな変数(x, y, z, t)で見たとき物理的にその点からみてある関係にある点の変数を持って来いと言われたとき、通信が起こる可能性がある。我々のタイプの計算では、通信は本質的にこの場合しか起こらない。変数をPEにどのように分割したか、ユーザーは知っているのだから、上のような場合だけを1つのルーティンにまとめておけばあとは並列計算のことは忘れてしまってもよい。問題は、1つのPE上でローカルなインデックスiを使って記述される

$$Y(i) = X(M(i))$$

というオペレーションが現れ、M(i)という点と同じPE上にない場合である。このとき、通信をしてよそからデータを持ってきてくれるルーティンを作り、XGATHERと名づけた(図6)。上の逆の機能を持ったルーティンをXSCATTERと呼んでいる。

$$Y(M(i)) = X(i)$$

(\*) 図5の配置は最善の解ではない。通信路の距離の総和を最短にするものは、下図のように配置した場合に実現される。ここで1は格子のz=1-4が置かれるPE、2はz=5-8が置かれるPE等々である。富士通計算科学研究センター柴田一哉氏、金澤宏幸氏のご指摘による。

```
1 2 3 4  1 2 3 4  1 2 3 4  1 2 3 4  1 2 3 4  1 2 3 4  ...
5 6 7 8  5 6 7 8  5 6 7 8  5 6 7 8  5 6 7 8  5 6 7 8  ...
```

```
1 2 3 4  1 2 3 4  1 2 3 4  1 2 3 4  1 2 3 4  1 2 3 4  ...
5 6 7 8  5 6 7 8  5 6 7 8  5 6 7 8  5 6 7 8  5 6 7 8  ...
```

....

前節で述べたように、繰り込み変換の計算では、格子のPEへの配置は複雑なものになり、このXGATHER/XSCATTERはたいへん強力なツールとなった。しかし以下のような問題点を持っている。

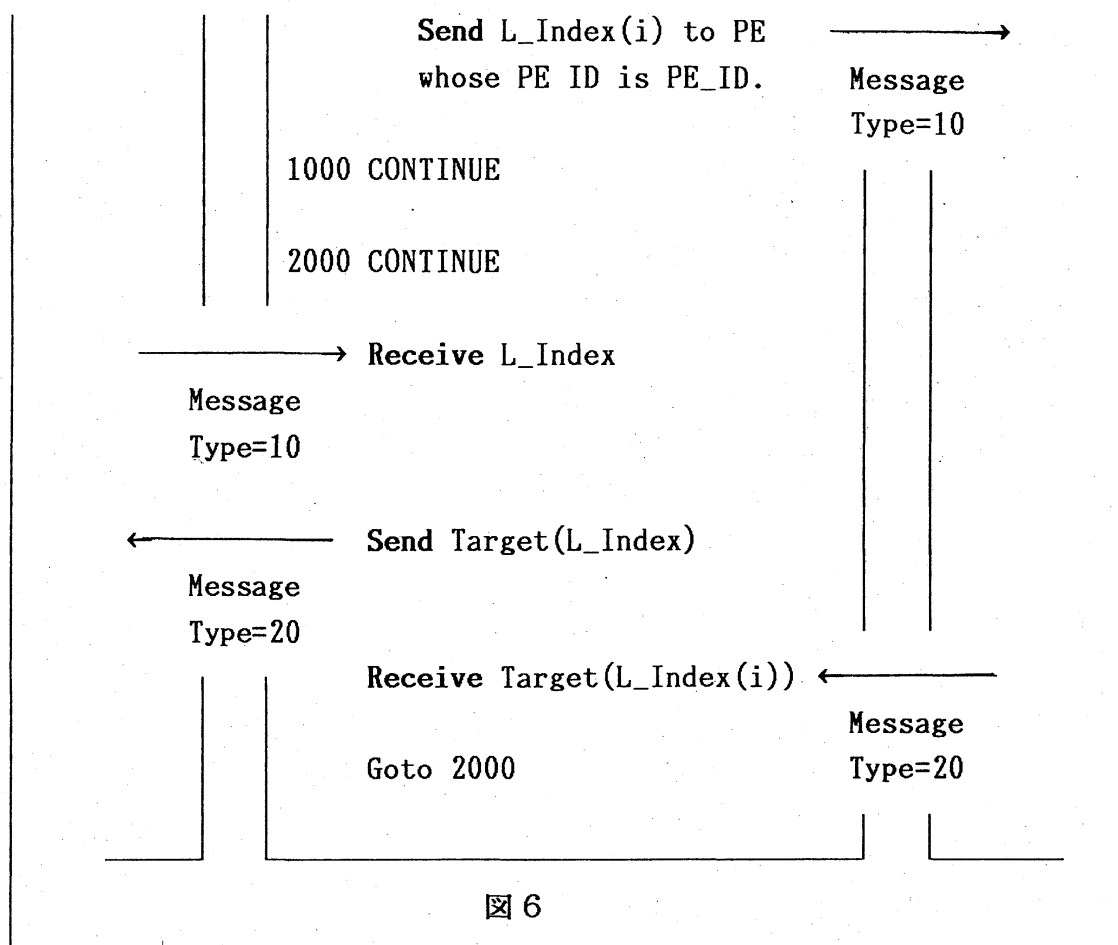
- i) グローバルなインデックスは、表にすると巨大になり、各PEに持たせることができなかった。そこで、配置のルールからそのたびに計算で求めている。このためXGATHER/XSCATTERはたいへん時間を喰うルーティンとなってしまった。
- ii) 図6に示したように、これらのルーティンは
  - ① 必要なデータを持っているPEに、欲しいデータのローカル・インデックスを送る。
  - ② 他のPEからの要求を受け取る。
  - ③ 他のPEへ要求されたローカル・インデックスを持ったデータを送る。
  - ④ 他のPEから送られてきたデータを受け取る。
 他のPEからのメッセージを受け取るためには、Receiveを発行しなければならないので、②の段階でどれだけの回数の要求がくるか知っていなければならない。裏タスク等の技巧を使わずにプログラムを書いたため、他のPEに要求するデータの数は、自分が要求されるデータの数と等しいということ仮定してしまった。このため、完全に一般性をもったgather/scatterルーティンにはなっていない。

```
XGATHER(Target,Source,G_Index,N)
```

```
Target = Source(G_Index(i)), i=1,2, ..., N
G_Index : Global Index
```

```
call MAP(G_Index,L_Index,PE_ID,N)
G_Index(i) → L_Index(i) at PE whose
PE ID is PE_ID.
i=1,2, ..., N
L_Index : Local Index
```

```
DO 1000 i=1, N
```



## 5. 結論

幾何学的な分割による並列化が可能であるような大規模問題では、超並列計算機は大変コスト・パフォーマンスのよいマシンである。実用にあたってユーザーの観点から気が付いた点を、順不同であるが以下にあげてみる。

- 1) 並列プログラミングの経験が浅いためか、人間が本質的に並列に考えることが困難なのか、バグが発生したとき、発見がむずかしい。新しいタイプのプログラミング・スタイルあるいはデバッガーが必要なのもかもしれない。(デバッガーで各PE上の状態を追うことができても、多数のPEの動きを把握するのは容易ではない。
- 2) 並列計算機ごとに並列化のストラテジー、プログラムを変えなければいけないのはつらい。たとえば、VPP型の並列プログラムとAP1000型の並列プログラムではまったく頭を切り替えなければならない。また、ベクトル演算装置がPEについているものとスカラーではストラテジーがま

まったく違う。

- 3) AP1000のグローバル・サム(全PE上の和)を計算するライブラリ-はたいへん有用であった。このように並列計算において非常に有用になるものを抽出して行ってライブラリ-化していくことは、並列計算機の普及には大きな意味をもつであろう。

#### 謝辞

本報告は富士通(株)・田子精男、奥田基、藤崎正英、川添明美、ハイデルベルグ大学・I.O.Stamatescu、スイス連邦工科大学・Ph.deForcrand、ベルリン計算センター・H.C.Hege、広島大学・宮村修、日置慎治、高石哲弥、福井大学・橋本貴明の諸氏との共同研究に基づくものです。

#### 参考文献

- K.Akemi 他、 "QCD on Parallel Computer AP1000",  
Nucl.Phys. B(Proc.Suppl.)26, 1992, 644-646.
- K.Akemi 他、 "Quantum Chromodynamics Simulations on a Nondedicated  
Highly Parallel Computer",  
プレプリント YAMAGATA-HEP-93-11  
(Comp.Phys.Comm. に投稿中)
- K.Akemi 他、 "Scaling Study of Pure Gauge Lattice QCD by Monte Carlo  
Renormalization Group Method",  
Phys.Rev.Letters 71, 1993, 3063-3066.