

## Cenju における並列プリミティブの実現

5 L - 5

松下 智, 中田 登志之, † 浅野 由裕, 小池 誠彦

日本電気 (株) C&C システム研究所, † 日本電気技術情報システム開発 (株)

### 1 はじめに

我々は並列回路シミュレーション向けの64台構成のマルチプロセッサシステム Cenju を開発した [1]. Cenju の通信網はバスと多段ネットワークからなる階層構成であり, 疑似分散共有メモリのユーザビューを提供している. Cenju ではクラスタ間アクセスのうち write のみハードウェア化している. ユーザはソフトウェア的に提供されるクラスタ間 read, および排他制御 / 非同期事象管理のための遠隔手続き呼び出し (RPC: Remote Procedure Call) を用いてアプリケーションを記述することができる [2].

さらに我々は, C 言語, Fortran77 からライブラリルーチンとして使用可能な並列プリミティブを整備した.

本稿では Cenju の並列プリミティブについて, 設計指針, 内部インプリメント, 記述性, 評価結果について報告する.

### 2 並列プリミティブの設計指針

#### Cenju 並列プリミティブの利点

- 高水準化
  - 記述の容易化.
  - バグ混入の防止. クリティカルセクションの見落としを減らす, 予期せぬ非決定性によるバグの混入を防止する.
  - アーキテクチャへの依存性を減らし, 後継機に対するソースの互換性を上げる.
- 高速化
  - 並列プリミティブはハードウェア構造を意識し, 最適なインプリメントを行なう. これにより, ユーザは容易に Cenju の性能を充分引きださう.

並列プリミティブの選択に当たっては, 既存のプリミティブ [3] を参考に, (1) 記述性が良いものであること, さらに, (2) Cenju でハードウェア化されているクラスタ間 write のみを用いて記述できるものを選択した. このため, ユーザはハードウェア構成を熟知しなくても最大限の性能を引きだすことが可能になる.

### 3 並列プリミティブの種類

Cenju の並列プリミティブのうち, 代表的なものを以下に示す.

Implementation of Parallel Primitives on Cenju  
Satoshi MATSUSHITA, Toshiyuki NAKATA, † Yoshihiro ASANO and Nobuhiko KOIKE  
NEC Corporation and †NEC Scientific Information System Development Ltd.

CJForki 並列実行の開始, Mach[4], Dynix[3] の fork と同様に, コンテキストを Cenju 内のプロセッサにブロードキャストし並列実行を開始する. ただし, データ, コードは共有できない.
CJSysBarrier 全てのプロセッサ間で同期をとる.
CJBarrier 設定したプロセッサグループで同期をとる.
CJRpc 他 PE の手続きを呼び出す, 呼びだされたルーチンの終了まで呼び出し側はブロックする.
CJNbRpc 他 PE の手続きをポインタで指定し呼び出す. 呼び出し側はブロックしない.
CJSend, CJRecv 領域 (配列) の転送. 転送が終了するまで双方はブロックする (同期型転送).
CJASend, CJAREcv ブロックしない (非同期型) 転送.

### 4 並列プリミティブの内部インプリメント

第2章で述べた性能向上の理由から, 本並列プリミティブはインプリメントにおいて, PE 間のデータ転送をできる限り write で記述している. 以下, barrier を例にとり概要を示す.

barrier barrier を行なう PE は, 同期テーブルをもつマスター PE とその他のスレーブ PE に分類される. barrier を行なう前に, 並列プリミティブ barron(master) によって, マスタ PE を指定する. マスタ PE ごとに 同期グループを分割することが可能である.

図1に示すように, マスタ PE には mask テーブルと sync テーブルが作成される. マスタプロセッサが同期の成立をスレーブに伝達するため, スレーブには ack フラグが置かれる. mask テーブルは, barron ディレクティブによりスレーブからセットされ同期の対象になる. 全てのスレーブがセットを終了したことを確認するために, 動作中の全てのプロセッサを同期対象にした sys-barrier を用意しこれを解決する. 以下, C ライクにアルゴリズムを記述する. 図中, SET, RESET は定数, 記号 <- は, クラスタ間 write アクセスを示す.

#### マスタ側

```
foreach(pe)
    if ((mask[pe] == SET) and (sync[pe] != SET))
        continue;
foreach (pe)
    ack <- SET;
```

#### スレーブ側

```
ack = RESET
sync[pe] <- SET
while (ack == RESET)
    ;
```

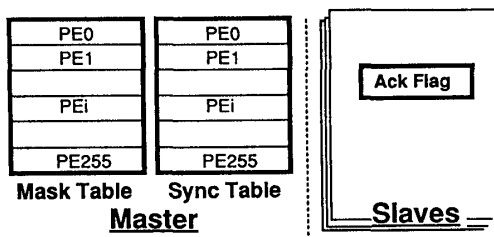


図 1: Barrier のデータ構造

## 5 記述力

本プリミティブを用いた場合、(1)barrier と非ブロック型 send, receive を用いた記述例を下に示す。このほか、(2)barrier と RPC による動的データの扱い、(3)send, receive によるパイプライン処理、(4) RPC による分散型のスケジューラ (5) 非ブロック型 RPC による集中型のスケジューラなど、さまざまなスタイルが容易に記述可能である。

並列プリミティブの記述例 (静的なデータ交換をするもの)

```

Forki(); /* 並列実行の開始 */
if (groupA) {
    CJBarron(peA); /* PE A をマスタにした,
                  同期グループ A */
} else {
    CJBarron(peB); /* 同期グループ B */
}

CJSysBarrier(); /* 同期グループの設定終了待ち */
/* 配列 vtx の計算 */
:
/* vtx を近隣のプロセッサに配布 */
/* CJASendDb1(送り先, 配列の先頭, 要素の個数,
              要素間の飛び幅) */
CJASendDb1(ilpid, &vtx[0][nrl], ld+1, nd+1);
CJASendDb1(iupid, &vtx[0][nru], ld+1, nd+1);
/* 同期 */
CJBarrier(groupA ? peA: peB);
/* CJAREcvDb1(送り元, 受け配列の先頭, 要素の個数,
              要素間の飛び幅) */
CJAREcvDb1(ilpid, &vtx[0][nrl-1], ld+1, nd+1);
CJAREcvDb1(iupid, &vtx[0][nru+1], ld+1, nd+1);

```

## 6 性能評価

メモリアクセスおよび RPC の性能を表 1 に示す。クラスター間 read は、ソフトウェア的に記述しているためクラスター間 write に比較して遅い。

第 4 章で示した、barrier, send の性能を表 2 に示す。64 台構成時に、最後の barrier 発行後、もっとも早い PE が barrier から抜けるまでの時間を First-done、最後の PE が抜けるまでの時間を Last-done で示す。First-done の値がテーブルサーチのオーバーヘッドを First-done と Last-done の差がマスタからスレーブへの ack 配布のオーバーヘッドにあたる。これまで、Cenju に実装されてきたアプリケーションは粗粒度のものが多かったため、同期のオーバーヘッドは問題にならなかったが、今後、データ構造の工夫により高速化しより細粒度の処理にも対応していく予定である。

表 1: 通信性能

Access Operation	Access Time
Local Access	250 nsec
Intra-Cluster Global Access	750 nsec
Inter-Cluster (Access)	5.2 $\mu$ sec
Global Write Access (Throughput)	800 nsec
(... + Pipe Clog)	2.4 $\mu$ sec
Inter-Cluster Read Access	240 $\mu$ sec
Inter-Cluster RPC (Latency)	640 $\mu$ sec to 1.2msec
(Throughput)	200 $\mu$ sec

表 2: barrier, send-receive の性能

Access Operation	Waiting Time
CJSysBarrier (First-done)	770 $\mu$ sec
(Last-done)	1500 $\mu$ sec
CJSend, CJRecv (per Long Word)	13 $\mu$ sec

send については、ブロック型の send, receive を用いて先に receive が発行された条件での転送時間をバスサイクル 1 回当たり換算したものを示す。表 1 と比較すれば、write アクセスに近い性能が引きだしていることがわかる。

## 7 終わりに

Cenju は write アクセスのみを強化した、回路シミュレータに特化したアーキテクチャを採用しているが、今回、汎用的な並列プリミティブを効率良くこの上にマップできることを示した。

さらに、我々は現在、Cenju 上で回路シミュレータ [5]、故障シミュレータ、LSI ルータ [6] など CAD 応用、モンテカルロシミュレータ、偏微分方程式など科学技術計算、ニューラルネットワークシミュレータなどの並列化を通して本プリミティブの評価整備を行なっている。

今後はアプリケーションの評価をもとに並列プリミティブの整備を進め、記述性の良い並列言語 / 並列 OS に発展させていくとともに、Cenju のアーキテクチャをより汎用な並列処理に適したものに発展させていきたいと考えている。

## 参考文献

- [1] 松下, 中田, 梶原, 浅野, 小池: 並列シミュレーションマシン, 情報 37 回全国大会予稿集, (1988)
- [2] 松下, 中田, 梶原, 浅野, 小池: 並列 CAD マシン Cenju のソフトウェア体系, 情報研報, Vol.89, No.99, pp. 93-100 (Nov. 1989)
- [3] Robert G.Babb II: *Programming Parallel Processors*, Addison-Wesley, (1987)
- [4] Michael B.Jones, Richard F.Rashid: *Mach and Machmaker; Object-Oriented Distributed Systems*, CMU Technical Report, CMU-CS-87-150, (1987)
- [5] Nakata, Tanabe, Onozuka, Kurobe, Koike: *A Multiprocessor System for Modular Circuit Simulation*, Proc. ICCAD '87, pp.364-367, (1987)
- [6] 山内, 中田, 石塚, 西口, 小池: 並列シミュレーションマシン Cenju 上の LSI ルータの評価, 情報 40 回全国大会予稿集, (1989)