

TR-I-0319

用例主導型機械翻訳における  
負荷分散による並列検索手法の検討

“A Study on the Parallel Implementation by Data Distribution

in Example-based Machine Translation”

Takashi Okada

1993.3

概要

用例主導型機械翻訳では、大量の用例データからの用例検索処理の高速化が課題である。そこで今回、負荷分散による用例の並列検索の検討とその実現を行った。本項では、並列化手法を、動的/静的負荷分散、及び、並列機構を考慮したデータ分配型/プロセス分配型負荷分散(筆者の提案する区分)に分けて検討した。結果として、疎結合型マシンでプロセッサ数倍、密結合型マシンで約4倍の用例検索処理の高速化を実現した。

ATR自動翻訳電話研究所

ATR Interpreting Telephony Research Laboratories

© ATR Interpreting Telephony Research Laboratories

## 1. はじめに

文献[1]の用例主導型機械翻訳（EBMT）では、入力文と類似な用例（原言語表現とその訳の対）を検索し、それを利用して翻訳を行う。この翻訳過程において、全用例について用例単位に入力文との属性距離計算を実行する類似検索処理が、翻訳時間上のボトルネックになっている。したがって、この類似検索処理の高速化が問題となり、その解決手段として、以下の二通りが考えられる。

### （アプローチ1）並列計算機上での負荷分散による並列化（全件検索）

### （アプローチ2）検索対象の動的絞り込み手法による検索（部分検索）

そこで今回、「アプローチ1」を選択し、MIMD型並列計算機上で、負荷分散による並列化手法の検討・実験をし、類似検索処理の高速化を実現した。以下で、その成果報告をする。

## 2. 用例主導型機械翻訳システムと類似検索処理の概要

### 2.1 用例主導型機械翻訳（EBMT）システム

今回の類似検索処理の高速化の実験で使用するEBMTシステムは、「名詞句”N1のN2”」の翻訳システムであり、全体の処理フローを図1に示す。このEBMTシステムは、三つの処理過程（解析・用例検索・変換）からなり、用例検索では千～十万件の用例からなる用例データベースと3千件程度の類語からなるシソーラス・データベースを使用する。用例は、「国際会議の登録」に関する対話の日英対訳データベースから抽出したもので、シソーラスは、大野、浜西の体系に準拠していて、その階層は三段階である。

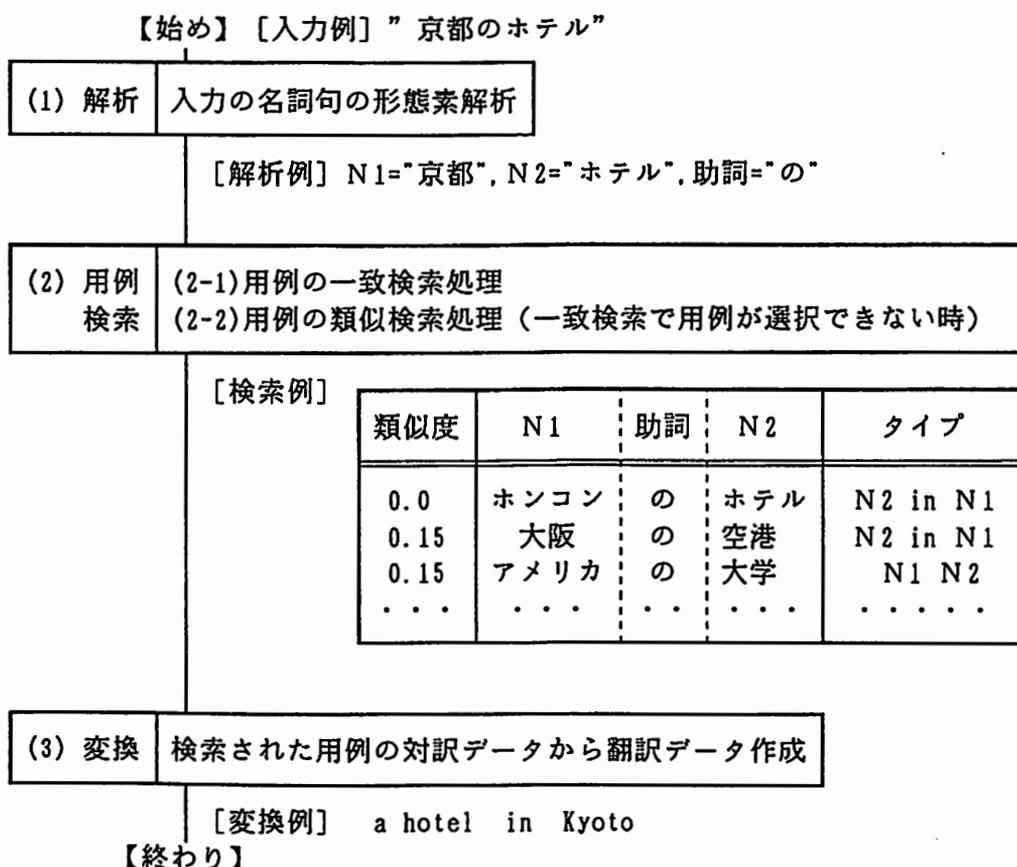


図1. EBMTシステムでの「名詞句”N1のN2”」の翻訳処理過程

## 2.2 類似検索処理の概略

類似検索処理は、用例検索部において一致検索処理で用例が選択できない場合に実行される処理で、用例と入力文とで属性距離計算を行い、それによって算出された類似度を比較して類似用例を選択する処理である。この処理の検索キーワードは類似度であり、それは入力の”N1のN2”と各用例間との属性距離計算で決まるため、全用例の類似度の算出を翻訳のたびに実行する必要がある。1回の属性距離計算にかかる処理時間は小さいが、用例数回(千~十万のオーダー)実行するため、類似検索処理全体に大きな処理時間を費やす。それゆえ、類似検索処理がEBMTシステムの処理時間上のボトル・ネックになっている。

次に、類似検索処理の処理内容(図2)、及び、属性距離計算の方法と処理実行時の特徴を、具体的に説明する。

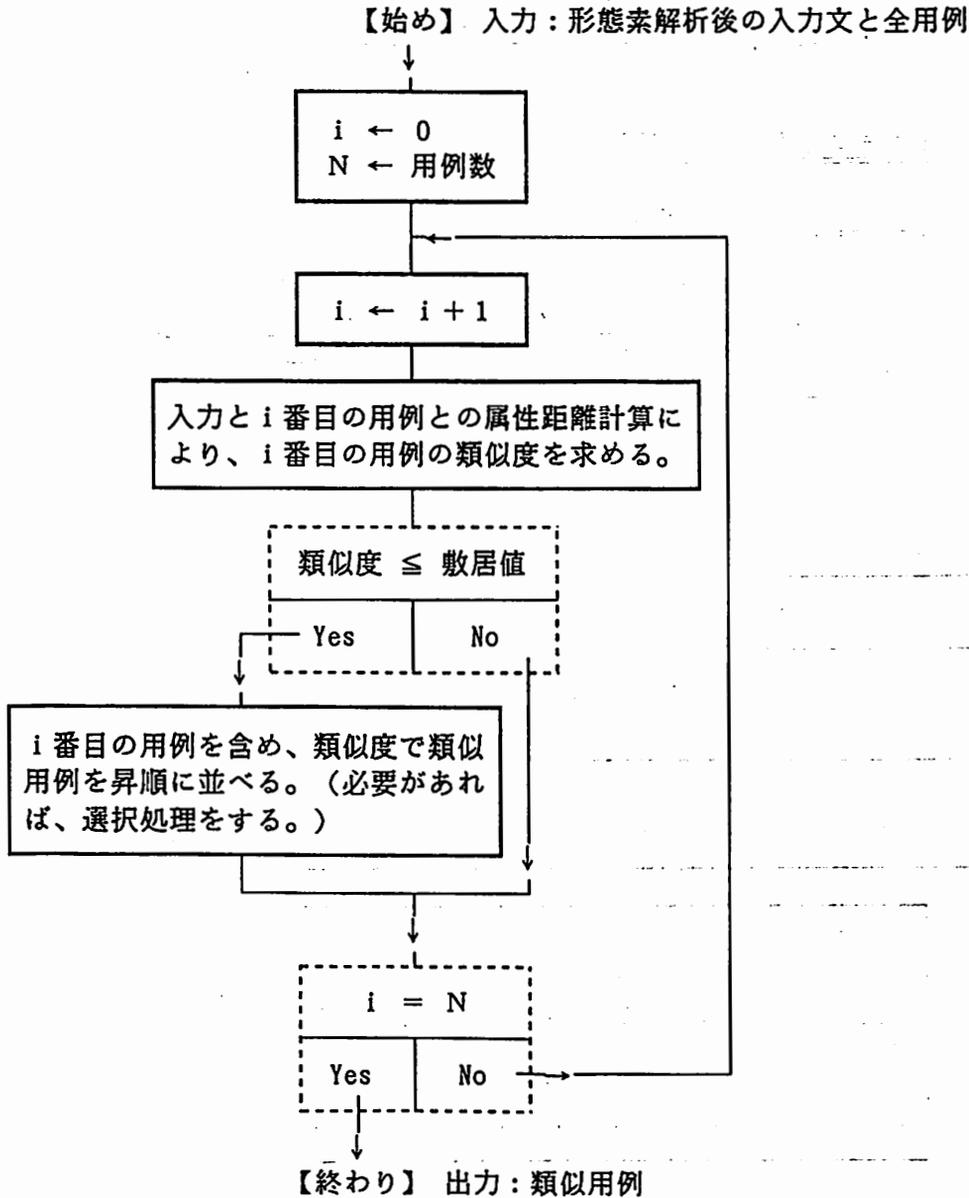


図2. 類似検索処理の処理フロー図

## 【属性距離計算】

属性として、「品詞」・「接辞」・「助詞」・「類語」を用い、「品詞」・「接辞」・「助詞」の距離は、等しければ距離を0 ( $r_i=0$ )、異なっていれば距離を1 ( $r_i=1$ )とする。「類語」は、角川類語辞典を使用し、そのソーラスは基本的に三段階のクラスから構成されている。(基本的に、類語コードは三桁の数値)によって、「類語」の距離は、 $(1 - \text{一致する階層数}/3)$ で求める。そして、各属性距離を加算する時、各属性距離には予め決められた重み ( $w_i$ ) が付加され、その合計を類似度とする。

また、重み ( $w_i$ ) は、属性値の持つ翻訳時の変換の相違におけるばらつきを考慮したもので、この値は、用例データのロード時に全て計算され、類似検索処理の実行前に参照される。

$$\text{類似度} = \sum_i w_i \times r_i$$

注)  $w_i$ :  $i$  番目の属性の重み  $r_i$ :  $i$  番目の属性の距離 (「類語」のみ 0, 1/3, 2/3, 1、その他 0, 1)

$$w_i = \frac{1}{\sqrt{\sum_j (k_{ij}/n)^2}}$$

注)  $n$ : 該当する用例数  $k_{ij}$ : 変換後のものが共通な用例数

## 【処理上の特徴】

- (1) 1用例当たりの類似検索処理時間が小さい。  
(Intel社「ipsc/2」のLisp環境上で、平均  $4.1 \times 10^{-2}$  sec )  
(Sequent社「Symmetry S2000」のLisp環境上で、平均  $1.7 \times 10^{-3}$  sec )
- (2) 用例のデータ量があらかじめ決まっていて、そのデータ量が十分大きい。
- (3) 処理において、他の用例の類似検索処理と同期関係がなく、用例が個々独立に処理される。

類似検索処理の処理上の特徴から、類似検索処理の高速化を目的として、負荷分散による並列処理を検討するのが適当であることが分かる。マルチプロセッサ・システムは、疎結合型と密結合型の2種類に分類できる。疎結合型では、Intel社「ipsc/2」を、密結合型では、Sequent社「Symmetry S2000」を使用して、負荷分散による並列化を考えていく。

## 3. MIMD型並列計算機上での負荷分散のためのプログラム並列化手法

まず最初に、マルチプロセッサ・システム上での負荷分散によるプログラム並列化のための方法論が、理論的にカテゴリー化されていないため、疎結合型と密結合型の並列化機構を考慮して、負荷分散によるプログラム並列化を分類し、典型的ないくつかの方法論にまとめてみる。

### 3.1 MIMD型並列計算機の特徴

疎結合型では、Intel社「ipsc/2」を、密結合型では、Sequent社「Symmetry S2000」を参考にして、マルチプロセッサ・システムの主な特徴を説明する。それらのマルチプロセッサ・システムの概略図を、図3に示す。

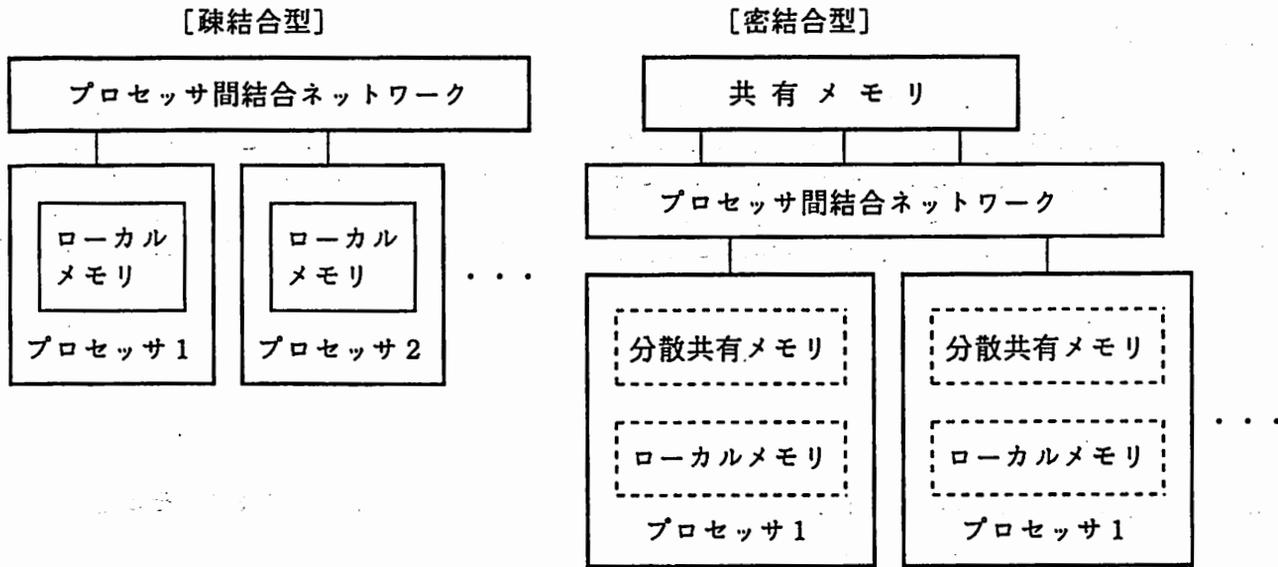


図3. マルチプロセッサ・システムの概略図

**疎結合型** : 分散メモリ・メッセージ交換方式

**【特徴】**

- ・ プロセッサ単位に、ローカルメモリを持つプロセスが静的に割り当てられている。
- ・ プロセス間の同期処理（グローバル・デッドロックの危険が伴う）・共有データの取り扱いは、メッセージ交換方式で行うため、比較的大きなオーバーヘッドがともない、且つ、他プロセスの能動的な制御を実行するのが難しい。

**密結合型** : 共有メモリ方式

**【特徴】**

- ・ 共有変数機能（共有メモリ上の変数）とそれに対する排他制御機能（lock・unlock）がある。
- ・ 高々数十程度のプロセッサによるマルチプロセッサ機構とプロセス制御機能（プロセス作成・起動・中止・続行・終了）がある。
- ・ アクセス時間は、共有変数よりもローカル変数の方が小さい。

3. 2. 負荷分散のためのプログラム並列化手法の分類

一般に、プログラムの並列化には、次の二種類がある。

- ・ 負荷分散（データ分散）：データを分割して複数のプロセッサに割り振り、並列実行する方法
- ・ 機能分散（コード分散）：コードを分割して複数のプロセッサに割り振り、並列実行する方法

以下では、「負荷分散（データ分散）」のみを具体的に検討する。負荷分散のためのプログラム並列化をMIMD型並列計算機上で実現するため、並列化機構上でのデータとプロセスとの関係から、以下の二つに分類した。（筆者が考案した分類）

- ・ データ分配型 : プロセッサ上のプロセスヘデータを分配して並列実行する負荷分散
- ・ プロセス分配型 : 個々のデータ（または、データの塊）を分割所有するプロセスを作成し、実行可能なプロセッサがそれを手分けして並列実行する負荷分散

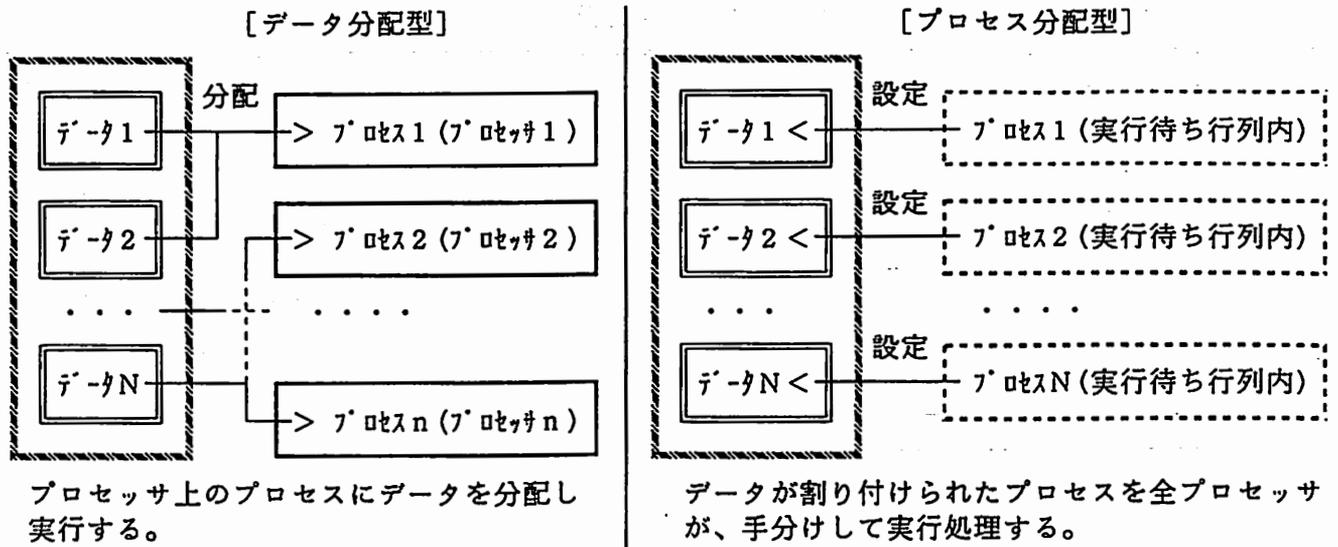


図4. MIMD型並列計算機の並列化機構を考慮した分類（データ分配型とプロセス分配型）

また、通常、負荷分散は、データ処理を分散させるタイミング（実行前・実行中）から、静的負荷分散と動的負荷分散とに分類される。動的負荷分散のデータ分配型は、1個のデータ分配管理用プロセス（masterプロセス）と1個以上のデータ受信用プロセス（workerプロセス）とで構成され、並列実行時に、それらのプロセス間でデータの授受が行われる。このデータ送受信のタイミングの違いから、同期型と非同期型の二種類に分類できる。

図5に、負荷分散のためのプログラム並列化手法における分類をまとめて示す。以下では、それらに対応したプログラム並列化方式を説明する。

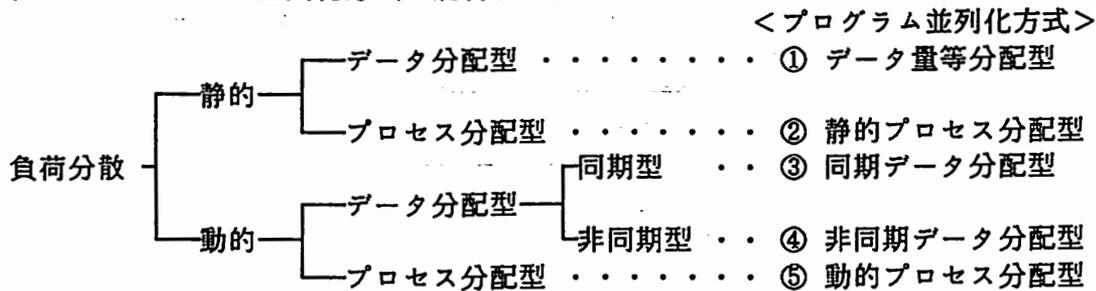
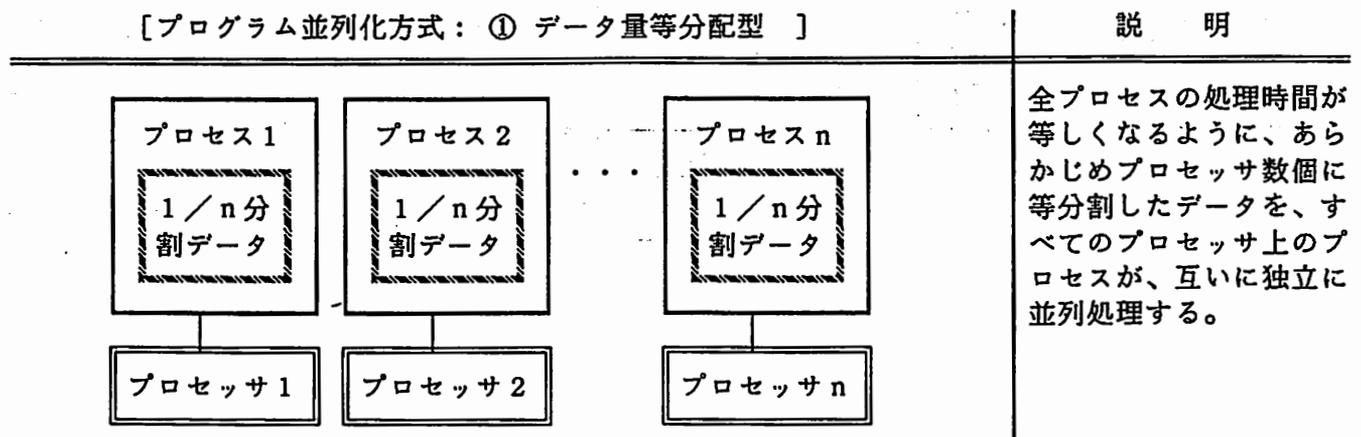
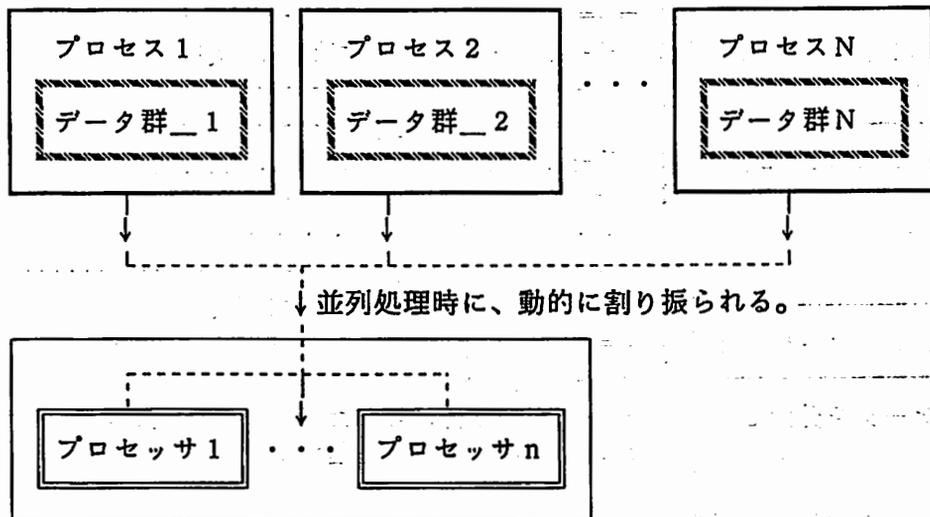


図5. 負荷分散のためのプログラム並列化手法における分類



[プログラム並列化方式：② 静的プロセス分配型]

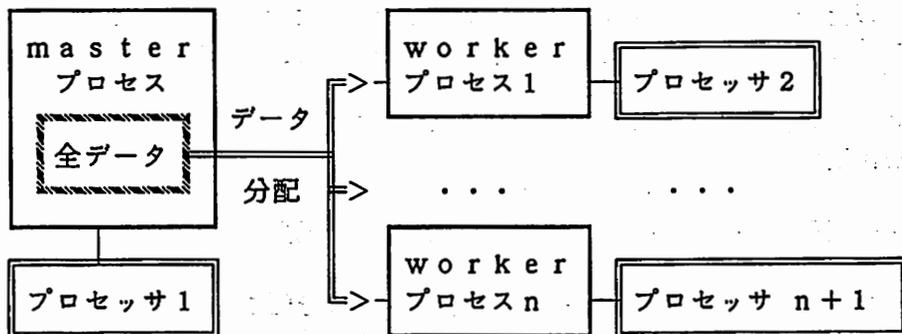
説明



すべてのプロセッサの処理時間が均等になるように、あらかじめ適当な数にプロセスを分割しておき、並列処理時に、すべてのプロセッサがそれらのプロセスを手分けして処理していく。

[プログラム並列化方式：③ 同期データ分配型]

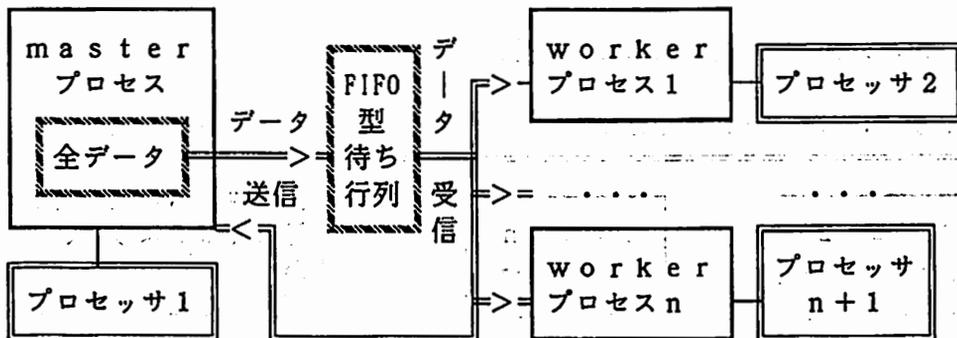
説明



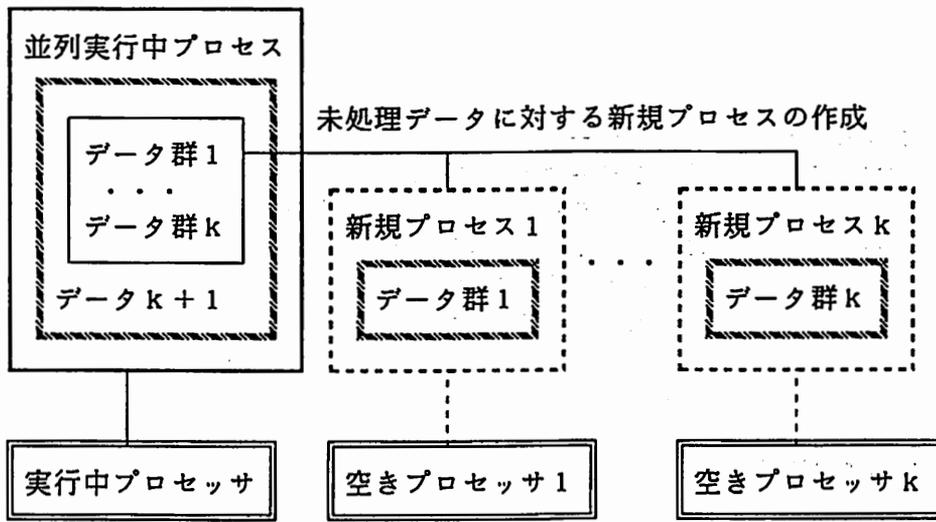
並列処理中に、master プロセスが、全 worker プロセスと直接同期をとってデータを分配し、worker プロセスにデータ処理させる。

[プログラム並列化方式：④ 非同期データ分配型]

説明



並列処理中に、master プロセスが、データ、または、その塊を媒体（FIFO型待ち行列）に送信し、worker プロセスが、それを非同期に受信し、データ処理を実行する。また、master プロセスもすべてのデータを送信した後、worker プロセスと同じ処理を実行する。



並列処理中に、プロセッサ上のプロセスが、データ処理の合間に、空きプロセッサを検索する処理を実行する。もし、空きプロセッサを見つけたら、自プロセス担当の未処理データの一部を処理するプロセスを作成する。その新しいプロセスを、空きプロセッサが処理する。

上述したプログラム並列化方式とそのデータ割り振り方法の一覧を、表1に示す。

表1. プログラム並列化方式とそのデータ割り振り方法

	プログラム並列化方式	データ割り振り方法
静的 負荷 分散	データ量等分配型 (静的データ分配型) ①	プロセッサ数個の等量データを作成し、並列処理の開始直後に、それらを各プロセッサ上のプロセスに割り振る。
	静的プロセス分配型 ②	並列処理開始以前に、データ量とデータ処理内容によるデータ区分に従って、データを分割し、おのおのにプロセスを割りつけておく。
動的 負荷 分散	同期データ分配型 ③	並列実行中に、masterプロセスが、すべてのworkerプロセスへのデータ分配を、直接管理する。
	非同期データ分配型 ④	並列実行中に、masterプロセスがデータ、または、その塊を受信先プロセスの指定をしないでmailboxに送信し、workerプロセスが、それを非同期に受信する。また、masterプロセスもデータ送信後、workerプロセスと同一の処理を行う。
	動的プロセス分配型 ⑤	並列実行中に、空きプロセッサを検知したプロセスが、自プロセス担当の未処理データの一部を処理するプロセスを作成する。

次に、プログラム並列化方式の処理時間の概算式を述べる。尚、「⑤動的プロセス分配型」は、並列実行時の空きプロセッサの発生状況やデータ分割・プロセス作成にかかるオーバーヘッドに大きく依存するため、ここでは省略する。

[以下で使用する記号の説明]

- D : 処理データ量
- n : 使用可能な最大プロセッサ数 (動的なデータ分配型の場合、workerプロセス数 = n - 1)
- Tseq : コード内の並列化できる処理部分の1データ当りの平均実行時間
- Tsynch : コード内の並列化不可能な処理部分の1データ当りの平均実行時間
- Tto : データのプロセスへの割り付けに於ける1データ当りの平均オーバーヘッド時間
- Tpro : プロセス制御 (作成・起動・終了) に於ける1プロセス当りの平均オーバーヘッド時間
- T'pro : プロセス制御 (起動・終了のみ) に於ける1プロセス当りの平均オーバーヘッド時間

[処理時間の概算式]

- ・ シーケンシャル版 :  $T_{seq} \times D + T_{synch} \times D$
- ① データ量等分配型 :  $\max_{i=1, n} (T_{seq_i}) \times D / n + T_{pro} + T_{synch} \times D$
- ③ 同期データ分配型 :  $T_{seq} \times D / (n - 1) + T_{pro} + T_{synch} \times D$
- ④ 非同期データ分配型 :  $T_{seq} \times D / n + T_{to} \times D / n + T_{pro} + T_{synch} \times D$
- ② 静的プロセス分配型 : 分割数が小さい (高々、プロセッサ数の数倍程度 :  $h \times n$ ) 場合、  

$$\max_{i=1, n} (\sum_{j=1, h} T_{seq_{i,j}}) \times D / (h \times n) + T'_{pro} \times h + T_{synch} \times D$$

分割数 (K) が大きい場合、

$$T_{seq} \times D / n + T'_{pro} \times K / n + T_{synch} \times D$$

注) 「①データ量等分配型」・「②静的プロセス分配型」の場合、分割されたデータの処理時間がすべて同じである保証がないため、1データ当りの平均処理時間を個別に取り扱う。

注) 「③同期データ分配型」の場合、「Tto」は、masterプロセスの各workerプロセス用の個別変数への値更新操作のみであるので無視できるものとした。

注) 「④非同期データ分配型」の場合、具体的に言えば、[資料-2]の式g1から、 $T_{seq} = t_{ave} \cdot \frac{T_{send} + T_{receive}}{m}$ 、 $T_{to} = \frac{T_{send} + T_{receive}}{m} + \Delta t$  という置き換えで導くことができる。

## 4. 「用例検索処理」の高速化実験とその結果

### 4.1 疎結合型マシン (Intel社「ipsc/2」)

#### [A] 特徴と性能

##### <特徴>

- ・ ローカル変数機能 (各要素プロセッサ内の個別メモリ上の変数)
- ・ ハイパーキューブ方式で要素プロセッサを結合したマルチプロセッサ機構
- ・ CPUは、i80386 (Vector Processor) で、そのキャッシュ・メモリのサイズは、8Kバイトである。また、各プロセッサ上で、NX/2というOSが走っていて、それらが要素プロセッサ間のメッセージ応答やファイル・アクセス、及び、並列処理用Lisp環境を個々に管理している。

注) 要素プロセッサ：ローカルメモリの付随したプロセッサで、一個のプロセスが対応する。

##### <性能>

- ・ プロセスへの関数実行指示の平均時間 : 0.18 sec
- ・ 非同期のメッセージ交換の平均時間 : 送信側: 0.06 sec 受信側: 0.08 sec
- ・ 同期のメッセージ交換の平均時間 : 送信側: 0.08 sec 受信側: 0.08 sec
- ・ 10000回 do-loop 処理の平均時間: 5.72 sec
- ・ 処理結果の回収 (join) 処理の方法は、[資料-1] より、階段方式による非同期型のメッセージ交換が効果的である。

#### [B] プログラム並列化方式の実現性

- ・ 実現に支障のないもの

##### ① データ量等分配型

- ・ 実現不可能なもの

- ② 静的プロセス分配型 : プロセスがプロセッサ固定で、プロセスの作成が不可能
- ⑤ 動的プロセス分配型 : プロセスがプロセッサ固定で、プロセスの作成が不可能

- ・ 実現できるが、並列効果を望めない。

- ③ 同期データ分配型 : プロセス間の応答はメッセージ交換のみで、これで同期処理やデータ分配を行うと、大きなオーバーヘッドがかかる。
- ④ 非同期データ分配型 : プロセス間の応答はメッセージ交換のみで、この方式の媒体 (FIFO型待ち行列) を使ったデータ送受信処理を実現するのは難しい。

#### [C] 実験結果

Intel社「ipsc/2」の「ipsc/2 lisp Version 1.3」(Common LISP + IPSC/2 Concurrent Constructs) という並列処理用Lisp環境上で、「①データ量等分配型」のプログラム並列化方式を実現し、用例検索処理の負荷分散による高速化実験を行った。尚、用例データベース内には2550件の用例しかなかったため、その以上の用例数を必要とする場合、それら2550件を複写して使用した。

——— : データ量等分配型

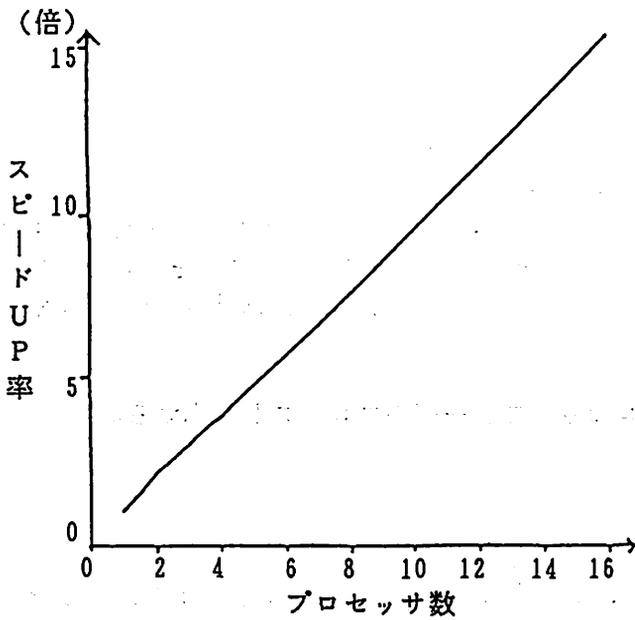


図5. 用例データ数=5000件での処理 (疎結合型 Intel社「ipsc/2」)

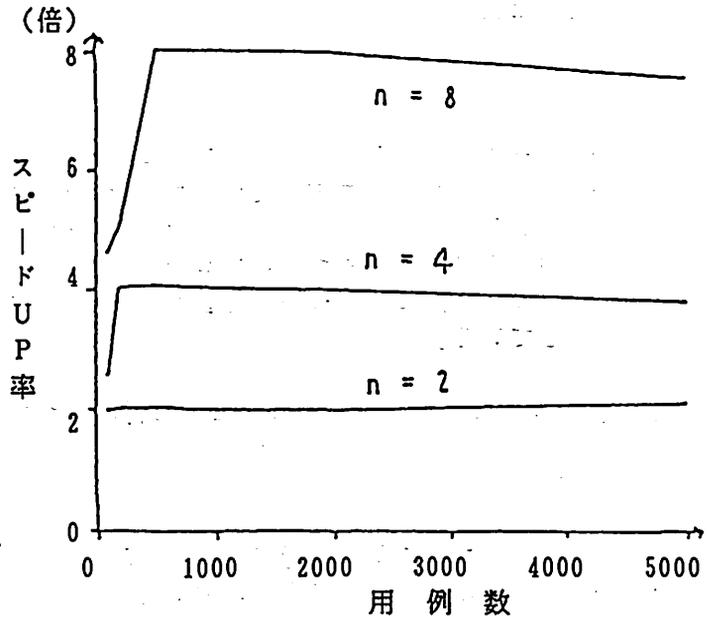


図6. プロセッサ数=2, 4, 8での処理 (疎結合型 Intel社「ipsc/2」)

表2. 「用例データ数=5000件での処理」の処理時間 (CPU sec) 一覧

プログラム並列化方式	1プロセッサ時	8プロセッサ時
① データ量等分配型	206.12 sec	27.21 sec

実験の結果、図5に示すように、用例数=5000件の時、プロセッサ数倍に近いスピードUPを実現できた。また、図6では、用例数が1000件以上で、プロセッサ数倍に近いスピードUPを実現できることが分かる。これらのことから、次の二点が言える。

- (1) 「①データ量等分配型」で、且つ、用例数が約1000件以上の用例検索処理において、並列処理のためのオーバーヘッド時間が、相対的に十分小さい。
- (2) 「①データ量等分配型」で、且つ、用例数が約1000件以上の用例検索処理において、すべての並列処理プロセスの処理時間が、近似している。

[まとめ]

1000件以上の用例に対する用例検索処理において、「①データ量等分配型」で、常に、プロセッサ数倍に近い高速化を実現できた。

## 4. 2 密結合型マシン (Sequent社「Symmetry S2000」)

### [A] 特徴と性能

#### <特徴>

- ・ 共有変数機能 (共有メモリ上の変数)
- ・ 排他制御機能 (lock・unlock)
- ・ プロセス管理機能
- ・ 高々数十個程度のマルチプロセッサ機構 (CPUは、i80486) で、DYNIXというOSがこのマルチプロセッサ機構上で走っていて、この並列処理用Lisp環境を管理している。

#### <性能>

- ・ プロセス起動/終了の平均時間 :  $1.5 \times 10^{-2}$  sec  
(1プロセッサ使用時の10プロセスでの平均時間)
- ・ lock/unlockの平均時間 :  $0.65 \times 10^{-2}$  sec ← sleepy-lock  
 $0.17 \times 10^{-2}$  sec ← spin-lock
- ・ 10000回do-loop処理の平均時間 :  $1.6 \times 10^{-2}$  sec
- ・ プロセス数、または、プロセッサ数を増加させると、1プロセスあたりのプロセス制御のオーバーヘッド時間が増加する。
- ・ アクセス・スピードは、共有変数よりもローカル変数の方が速い。プロセッサ数が増加するほど、その傾向が増加する。
- ・ 1用例あたりの用例検索処理時間が、媒体 (FIFO型待ち行列) の排他制御 (lock/unlock) にかかる時間より十分小さいので、「④非同期データ分配型」では、[資料-2] で求めたデータ送信単位を使用する。

### [B] プログラム並列化方式の実現性

- ・ 実現に支障のないもの

- ① データ量等分配型 , ② 静的プロセス分配型 , ③ 同期データ分配型 ,  
④ 非同期データ分配型 , ⑤ 動的プロセス分配型

### [C] 実験結果

Sequent社「Symmetry S2000」の「Allegro Clip」(Version 3.0.3) という並列処理用Lisp環境上で、上述したすべてのプログラム並列化方式を実現し、用例検索処理の負荷分散による高速化実験を行った。また、この実験では、用例はすべて共有メモリ上に置き、プロセスに分配、または、分割されるデータとしては、アドレスやデータ領域指定のインデックスを使用した。

表3. 「用例データ数=10000件での処理」の処理時間 (CPU sec) 一覧

プログラム並列化方式	1プロセッサ時	8プロセッサ時
① データ量等分配型	15.93 sec	4.42 sec
② 静的プロセス分配型	15.33 sec	16.62 sec
③ 同期データ分配型	-	4.48 sec
④ 非同期データ分配型	-	5.63 sec
⑤ 動的プロセス分配型	15.28 sec	4.85 sec

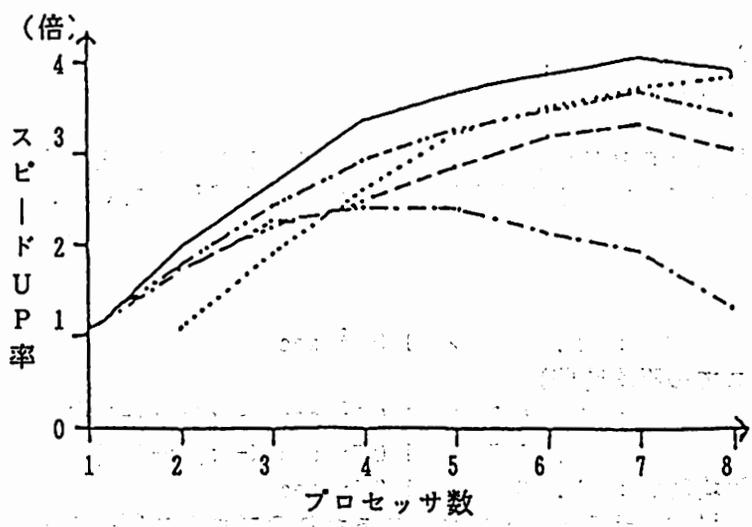
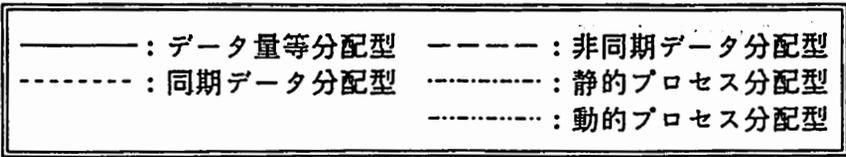


図7. 用例データ数=10000件での処理  
(密結合型 Sequent社「Symmetry S2000」)

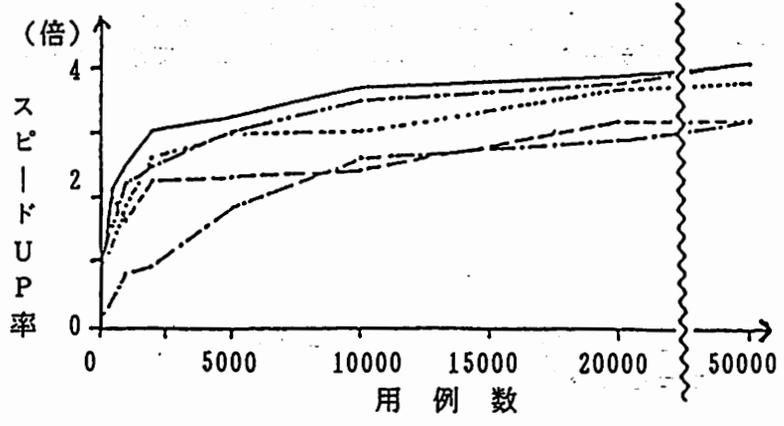


図8. プロセッサ数=8での処理  
(密結合型 Sequent社「Symmetry S2000」)

注) 実験での「⑤動的プロセス分配型」のプロセス数は、プロセッサ数×8個とした。

実験の結果、図7・図8に示すように、並列処理時のオーバーヘッドが最も小さい「①データ量等分配型」が、常に、最も高いスピードUP率を示した。また、どのプログラム並列化方式も、並列処理時のオーバーヘッドのため、スピードUP率が頭打ちの状態になった。

[まとめ]

「データ量等分配型」の並列化方式が、常に、最も高いスピードUP率を示し、そのスピードUP率は、高々、4倍程度である。スピードUP率が4倍程度に留まるのは、類似検索処理上での共有メモリ・アクセスに時間がかかっているからである。これは、各プロセスが頻繁に共有メモリをアクセスするために、プロセス間のバス競合によるアクセス渋滞が発生するからである。

## 5. 類似検索処理の負荷分散による並列化の並列効果についての考察

文献 [3] [4] [5] より、負荷分散における並列効果（ここでは、スピードUP）の要因として、次の二つの主要な要因を考える。尚、以下の議論では、並列実行前に負荷分散にかかるオーバーヘッド時間を除外して考える。

[要因-a] 使用可能なプロセッサのアイドル時間（休止・待ち時間）を最小にする。

[要因-b] 負荷分散にかかるオーバーヘッド時間を最小にする。

並列効果（ここでは、スピードUP）の第一義的目的は、[要因-a] である。また、[要因-b] について、並列実行時のオーバーヘッドは、静的負荷分散より動的負荷分散の方が大きいことが容易に分かる。これより、静的負荷分散のプログラム並列化方式で [要因-a] が満たされる場合、静的負荷分散による方式が最も効果的であると言える。もし、[要因-a] が満たされない場合、動的負荷分散のプログラム並列化方式の範囲内で [要因-b] を満たす方式を選択する必要がある。

[要因-a] を決定づける主要なものを、いくつか挙げる。

- (1) データ処理時間の分散
- (2) データ領域内でのデータ処理時間のかたより
- (3) データ数
- (4) プロセッサ数

[要因-b] を決定づける主要なものを、いくつか挙げる。

- (1) 並列処理機構の性能
- (2) データ数
- (3) プロセッサ数

次に、類似検索処理の場合、「①データ量等分配型」のプログラム並列化方式が最適であるのは、用例のデータ量を等分割しただけで、[要因-a] を満足するからである。[要因-a] を満足する根拠として、1 データあたりの処理時間が小さいだけでなく、以下のものが挙げられる。

[根拠]

データ領域内で、1 データ当りの処理時間が、独立で、かつ、予想される処理時間の期待値が等しいため、中心極限定理により、データ数が多いほど、等量のデータ数を処理するプロセスの実行時間のばらつきが、相対的に小さくなる。

注) 文献 [7] を参考に、データ処理時間のばらつきが正規分布をしている場合のデータ数 (n) と全体の処理時間 (W) ののばらつきとの関係を理論的に求める。

中心極限定理より、母集団の分布の形に関係なく、平均  $\mu$ 、分散  $\sigma^2$  が与えられると、標本  $X_1, X_2, \dots, X_n$  に基づく標本平均  $\bar{X}$  の標本分布は、正規分布  $N(\mu, \sigma^2/n)$  となる。  
即ち、 $\bar{X}$  は、母平均  $\mu$  を中心として分布し、しかも—これが重要なのであるが— $\mu$  に集中する傾向がデータ数 (n) と共に顕著になっている。

それゆえ、全体のデータ処理時間 (W) は、 $n \cdot \bar{X}$  であるから、その標本分布は、

$$\text{正規分布 } N(W_0, \sigma_w^2) \quad \therefore W_0 = n \cdot \mu \quad \sigma_w = \sqrt{n} \cdot \sigma$$

となる。これより、データ数 (n) が増加するに従って、全体の処理時間 (W) の分散 ( $\sigma_w$ ) は、相対的に小さくなる。

## 6. まとめ

類似検索処理の負荷分散による並列化の最適な高速化手法は、「データ量等分配型」のプログラム並列化方式である。その内容を、以下で述べる。

### 【手続き】

並列計算機上で、すべてのプロセッサ対応プロセスに対し、各プロセスのローカル・メモリ上にデータ量を等分割した用例をあらかじめロードしておいて並列実行する静的負荷分散による並列化方式

### 【並列効果】

疎結合型マシンでプロセッサ数倍、密結合型マシンで約4倍の高速化

## 7. 今後

今回、検索対象の用例数が処理時間を遅らせる原因である類似検索処理において、負荷分散による並列処理が効果的であることが示された。今後は、本研究で取り扱わなかった検索対象の動的絞り込み手法による検索（部分検索）を検討する必要がある。

## 参考文献

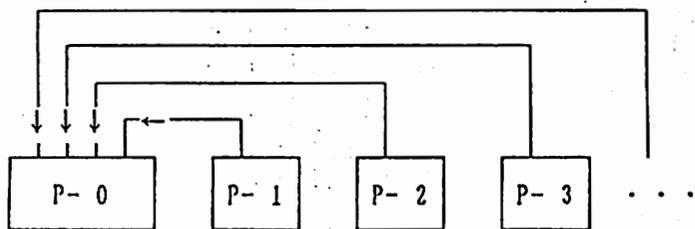
- [1] 隅田英一郎・飯田仁 「用例主導型機械翻訳」, 自然言語処理 82-5, 1991. 3. 15
- [2] Carriero and Gelern 「How to write Parallel Programs」
- [3] 新田淳 他, 「大規模DB/D.Cの高性能・高信頼化方式の考察」, アドバンスデータベースシステムのシンポジウム, 1885, 12月
- [4] 笠原博徳 著「並列処理技術」, コロナ社, P41-P45
- [5] 益田隆司 他, 「データベース及び知識ベースを対象としたストリーム指向型並列処理系の共有メモリ・マシン上への実現」, アドバンスデータベースシステムのシンポジウム, 1885, 12月
- [6] 牧野都治 著「待ち行列の応用」, 森北出版, 数学ライブラリー12
- [7] RAMON E. HENKEL 著「統計的検査」, 朝倉書店, 人間科学の統計学6

資料 1

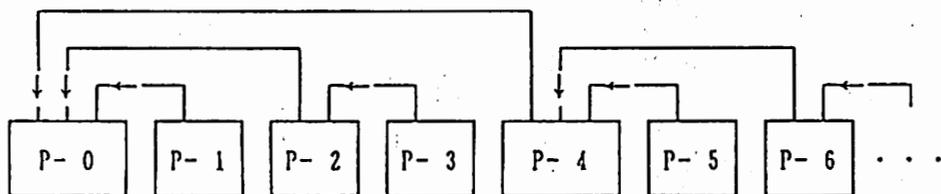
Intel社「ipsc/2」の「ipsc/2 lisp Version 1.3」という並列処理用Lisp環境上における  
 メッセージ交換による処理結果の回収 (join) 処理の性能テスト  
 (1プロセスが、他のすべてのプロセスから1Kバイトの文字列データを収集するメッセージ交換処理)

処理結果の回収のためのメッセージ交換方法として、2種類を設定した。

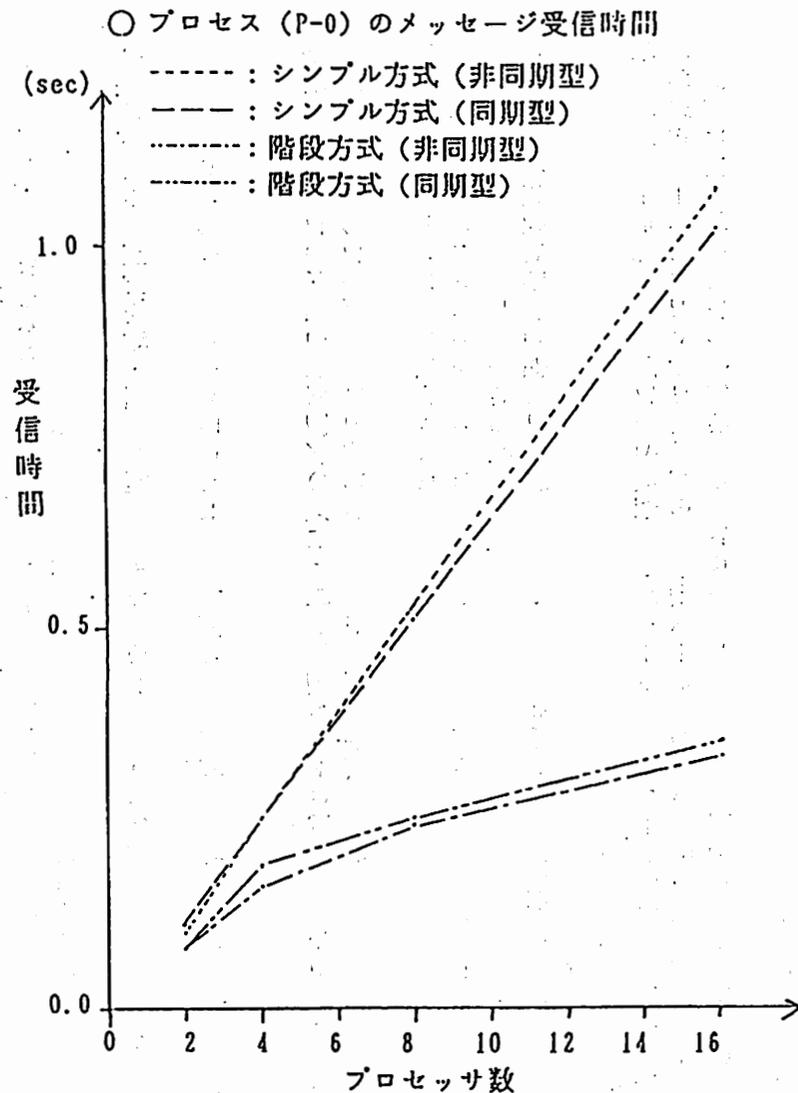
- 1) シンプル方式：単純に、1個のプロセス(P-0)に、他プロセスからのメッセージ送信を集中させる。



- 2) 階段方式：iを0～ $\lceil \log_2 n \rceil$ まで、(奇数 $\times 2^i$ )番目から(奇数 $\times 2^i - 1$ )番目へのメッセージ送信を繰り返す。(奇数：1, 3, 5, ...)



注) P-n : n番目のプロセス (n : 0, 1, 2, ...)



「④非同期データ分配型」では、すべてのプロセスが競合して、媒体（FIFO型待ち行列）にアクセス（データ送受信）する。媒体とのデータ送受信を効率的に行うため、データを複数個まとめて送受信する方法を取り、そのデータ数を『データ送信単位』として取り扱う。データ送信単位の大小は、以下のような問題を引き起こす。

- 『データ送信単位』が小さい。 <-> 媒体に、空き状態が発生しやすくなり、データ送受信の回数が増加する。
- 『データ送信単位』が大きい。 <-> workerプロセスのデータ処理終了のタイミングにズレが生じ、プロセッサのアイドル時間が増える。

これより、適切なデータ送信単位（ $m$ ）の値が必要となる。以下で、媒体（FIFO型待ち行列）を mailbox と呼び、適切なデータ送信単位（ $m$ ）の決定条件を提示して、データ送信単位（ $m$ ）の算出方法を、理論的に求める。

まず、負荷分散における効果的な並列効果のための二つの要因から、5つの制約条件を提示する。

[要因1] 使用可能なプロセッサのアイドル時間（休止・待ち時間）を最小にする。

- ・ 制約条件 a > mailbox の空状態の回避
- ・ 制約条件 b > 各プロセスのデータ処理終了タイミングのズレが小さいことの保証  
(1回のデータ受信によるデータ処理時間が、データ処理全体からみて十分小さいことの保証)
- ・ 制約条件 c > 全 worker プロセスが最低1個のデータを処理することの保証

[要因2] 負荷分散にかかるオーバーヘッド時間を最小にする。

- ・ 制約条件 d > mailbox アクセスに於けるプロセスの平衡状態（使用率 < 1）の保証
- ・ 制約条件 e > データ送信/受信の回数をなるべく少なくする。

これらの制約条件を数式で表現するために、mailbox を介した master/worker プロセスによるデータ送受信をモデル化し、次の二つのアプローチを行う。

- (A) プロセスのアクセス頻度から、mailbox の空状態の回避条件を求める。
- (B) プロセス間のアクセス競合を考慮した mailbox へのデータ送信時間（受信時間）を求める。

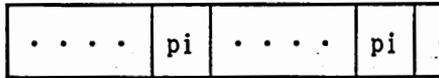
注) 以下で使用する記号の説明

- D : 全データ数
- n : プロセッサ数 (= worker プロセス数 + 1)
- m : 送信データ単位 (1回の送信で送られるデータ数)
- $\Delta t$  : 送信データ作成時間 (全てのアクセス範囲情報作成の1データ当りの消費時間、又は、データ・リスト作成の1データ当りの消費時間)
- Tsend : master プロセスから mailbox へのデータ送信時間 (Tsend  $\approx$  Treceive)
- Treceive : worker プロセスの mailbox からのデータ受信時間
- T0 : 単一プロセスでのデータ送信/データ受信の平均時間 (プロセス競合のない場合)
- tave : 一個のデータの平均処理時間

【 (A) プロセスのアクセス頻度から、mailboxの空状態の回避条件を求める。】

mailboxの空状態の回避を考える場合、mailboxへのアクセス頻度を問題にする。以下で、アクセス待ち行列を使って、この問題を考えていく。

[アクセス待ち行列]



プロセス  $i$  ( $p_i$ ) のアクセス頻度を、

順番: 1 2 3 ... k (k+1) ... m ...

$$a_i = 1 / \text{アクセス間隔 } i \quad \dots \text{ 式 } h1$$

とすると、

masterプロセスは、mailboxへのデータ格納処理を行い、workerプロセスは、mailboxからのデータ搬出処理を行うため、mailboxの空状態の回避のためには、次の式が成り立つ必要がある。

全workerプロセスのアクセス頻度 < masterプロセスのアクセス頻度

$$(n-1) \times a_{\text{worker}} < a_{\text{master}}$$

式 h1 より、

$$(n-1) \times \text{masterのアクセス間隔} < \text{workerのアクセス間隔}$$

$$(n-1) \times (\Delta t \times m + T_{\text{send}}) < m \times t_{\text{ave}} + T_{\text{receive}}$$

$$m \times (t_{\text{ave}} - (n-1) \times \Delta t) > (n-1) \times T_{\text{send}} - T_{\text{receive}} \quad \dots \text{ 式 } h2$$

$T_{\text{send}} \approx T_{\text{receive}}$ ,  $n \geq 2$  により、

$$(n-1) \times T_{\text{send}} - T_{\text{receive}} \approx (n-2) \times T_{\text{send}} \geq 0 \quad \dots \text{ 式 } h3$$

【 (B) プロセス間のアクセス競合を考慮した mailbox へのデータ送信時間 (受信時間) を求める 】

データ送信時間 ( $T_{\text{send}}$ ) において、mailbox へのプロセス競合による遅延の問題を、高々プロセッサ数個の待ち行列問題としてとらえ、プロセッサ数個のコインを投げる問題と同等と考えて、この理論標本分布を、二項分布とみなす。

[今回必要とする二項分布の方程式]



$0 < p < 1$ ,  $q = 1 - p$ ,  $N$ : 試行回数とすると、

$$(p+q)^N = \sum_{i=0}^N {}_N C_i p^i q^{N-i} = 1$$

注)  $p_i$ :  $i$  番目のプロセス □: 待ち状態のプロセス

$$\text{平均値: } \sum_{i=0}^N i \cdot {}_N C_i p^i q^{N-i} = N \cdot (p+q)^{N-1} \cdot p = N \cdot p \quad \text{分散: } N \cdot p \cdot q$$

注) データ受信時間: ここでは、データ送信時間と同じものとして取り扱う。

「④非同期データ分配型」では、masterプロセスの処理内容の相違（データ分配・データ処理）から、処理状態を二つに分け、データ送信時間（Tsend）を考える。

[状態-1] masterプロセスのmailboxへのデータ分配処理時で、masterプロセスは常時、mailboxをアクセスしているものとして取り扱い、workerプロセスと別途に取り扱う。この時のworkerプロセスのmailboxへのアクセス待ち状態は、二項分布をなすものと見なす。

[状態-2] masterプロセスのmailboxへのデータ分配処理終了後で、全プロセスがmailboxを均等にアクセスしているものと見なし、全プロセスのmailboxへのアクセス待ち状態は、二項分布をなすものと見なす。

(B-1) 状態-1 (masterプロセスのmailboxへのデータ分配処理時) の場合

masterプロセスのデータ送信処理は、常時、mailboxをアクセスしているとみなし、mailboxへのアクセス確率(p)を1とする。workerプロセスのmailboxへのアクセス確率(p)をデータ受信時間/(1回のデータ送信でのデータ処理時間) = Treceive/(m · tave) とする。

$$\begin{aligned} \text{アクセス待ち行列の列平均数 (k)} &= \sum_{i=0}^{n-1} (i+1) \cdot {}_{n-1}C_i p^i q^{n-1-i} \\ &= \sum_{i=0}^{n-1} i \cdot {}_{n-1}C_i p^i q^{n-1-i} + \sum_{i=0}^{n-1} {}_{n-1}C_i p^i q^{n-1-i} \\ &= (n-1) \cdot p + 1 \\ &= (n-1) \frac{\text{Treceive}}{m \cdot t_{ave}} + 1 \end{aligned}$$

よって、Treceive ≒ (k+1) T0 より

$$\text{Treceive} \approx \frac{2 \cdot T_0}{1 - \frac{(n-1) \cdot T_0}{m \cdot t_{ave}}} \quad \dots \text{式 h 4}$$

(B-2) 状態-2 (masterプロセスのmailboxへのデータ分配後の全プロセスによるデータ処理時) の場合

全プロセスのmailboxへのアクセスは、平衡状態を保ち、アクセス確率(p)をデータ受信時間/(1回のデータ送信によるデータ処理時間) = Treceive/(m · tave) とする。

$$\begin{aligned} \text{アクセス待ち行列の列平均数 (k)} &= \sum_{i=0}^n i \cdot {}_n C_i p^i q^{n-i} \\ &= n \cdot p \\ &= \frac{n \cdot \text{Treceive}}{m \cdot t_{ave}} \end{aligned}$$

よって、 $T_{receive} \approx (k+1) T_0$  より

$$T_{receive} \approx \frac{T_0}{1 - \frac{n \cdot T_0}{m \cdot t_{ave}}} \quad \dots \text{式 h 5}$$

以下で、上述の二つのアプローチの結果（式 h 1～式 h 5）から、制約条件 a～e を式に展開する。

- ・ 制約条件 a > mailbox の空状態の回避

式 h 2・式 h 3 により、プロセッサ数 (n) と送信データ単位 (m) が次のような制約を受ける。

◆ プロセッサ数 (n) の上限値に対する制約

状態-1 に於て、一個のデータについて、平均処理時間 ( $t_{ave}$ ) が、worker プロセス数分の送信データを作成する時間 ( $(n-1) \times \Delta t$ ) よりも大きくなければ成らない。

$$t_{ave} > (n-1) \times \Delta t \quad \dots \text{式 f 1}$$

◆ 送信データ単位 (m) の下限値に対する制約

状態-1 に於て、master プロセスの worker プロセス数個のデータ送信にかかる全体の処理時間よりも、worker プロセスの mailbox アクセス間隔 (1 送信データ単位のデータ受信インターバル) の方が大きくなければならない。

$$m > \frac{(n-2) \times T_{send}}{t_{ave} - (n-1) \times \Delta t}$$

式 h 4,  $T_{send} \approx T_{receive}$  より、

$$m > \frac{(n-2) \times 2 \times T_0}{t_{ave} - (n-1) \times \Delta t} + (n-1) \times \frac{T_0}{t_{ave}} \quad \dots \text{式 f 2-1}$$

- ・ 制約条件 b > 各プロセスのデータ処理終了タイミングのズレが小さいことの保証  
(1 回のデータ受信によるデータ処理時間が、データ処理全体からみて十分小さいことの保証)

この制約条件を、「最大のズレの時間が、並列処理時間の高々 1% 程度であればよい」という条件に置き換えて考えると、送信データ単位 (m) が次のような制約を受ける。

◆ 送信データ単位 (m) に対する制約

プロセス内のデータ処理終了時間のズレの最大 (Max) は、1 回のデータ受信のよる処理にかかる処理時間 ( $Lag = m \times t_{ave} + T_{receive}$ ) である。データ処理全体の時間 (T) は、式 g 1 を使う。すると、そのズレ比率の逆数 (L) は、

$$L = \frac{T}{Lag} \approx \frac{D}{n} \left( t_{ave} + \frac{T_{send} + T_{receive}}{m} + \Delta t \right) \times \frac{1}{m \times t_{ave} + T_{receive}}$$

この場合、 $T_{send}$ ,  $T_{receive}$  は、 $m \times t_{ave}$  より十分小さい場合の問題であり、又、 $\Delta t$  は無視できるので、

$$\approx \frac{D}{n \times m}$$

となり、この比率 (L) の逆数がズレの比率になるから、以下のような条件式になる。

$$\frac{D}{n \times m} > 100 \quad \dots \text{式 f 3}$$

- ・ 制約条件 c > 全 worker プロセスが最低 1 個のデータを処理することの保証

送信回数が、worker プロセス数以上あればよいので、以下のようなになる。

$$(n-1) \times m \leq D \quad \dots \text{式 f 4}$$

- ・ 制約条件 d > mailbox アクセスに於けるプロセスの平衡状態 (使用率 < 1) の保証

アクセス待ち行列系の平衡状態を維持する (使用率 < 1) ための条件として、 $n \neq 2$  の場合、式 f 2-1 を十分条件として利用できるが、特に、状態 - 2 に於て、 $n = 2$  の場合 (1 個の master プロセスと 1 個の worker プロセスの場合)、式 h 5 により、送信データ単位 (m) が以下のような制約の追加を受ける。

◆ 送信データ単位 (m) の下限値に対する制約 (追加)

$$m > \frac{n \times T_0}{t_{ave}} \quad \dots \text{式 f 2-2}$$

- ・ 制約条件 e > データ送信/受信の回数をなるべく少なくする。

データ受信回数を減らすには、送信データ単位 (m) を最大にすればよい。

$$m = \text{MAX} (m \text{ 値の許される範囲}) \quad \dots \text{式 f 5}$$

また、これら 6 つの制約条件式の内、式 f 1・式 f 2-1・式 f 2-2・式 f 4 を必要条件として扱い、式 f 3・式 f 5 は、処理のスピード UP 率 (速度向上率) の安定化のために、おおまかな目安として付加した条件であり、他に単位データ当りの処理時間のばらつきなどを加味して、最適な値を決定付ける必要があるものであり、ここでは、この二式を十分条件として扱う。

したがって、データ送信単位 (m) を算出方法は、以下のようになる。

【 データ送信単位 (m) の算出方法 】

[入力] : プロセッサ数 (n), 平均データ処理時間 ( $t_{ave}$ ), データ量 (D),  
送信データ作成時間 ( $\Delta t$ ), 1 プロセスでのデータ送信/受信時間の平均 (T0)

[出力] : データ送信単位 (m), 変更されたプロセッサ数 (n), エラー (実行不可能) の有無

[処理]

- step1 => プロセッサ数 n が、式 f 1 を満たすか?
  - • • 満たさない場合、プロセッサ数 n の変更 or エラー終了
- step2 => 式 f 2 - 1 ・式 f 2 - 2 を満たす最小の送信データ単位 (m) を設定
- step3 => プロセッサ数 (n) ・送信データ単位 (m) ・データ量 (D) が、式 f 4 を満たすか?
  - • • 満たさない場合、プロセッサ数 n の変更 or エラー終了
- step4 => プロセッサ数 (n) ・送信データ単位 (m) ・データ量 (D) が、式 f 3 を満たすか?
  - • • 満たさない場合、警告
  - • • 満たす場合、式 f 5 により、送信データ単位 (m) を、 $D / (100 \times n)$  を超えない最大の整数値に変更する。

データ送信単位 (m) の算出において、最も並列効果の高い最適な値を求める為には、このデータ送信単位 (m) の算出方法内で、処理データに対して、「平均データ処理時間」・「データ量」以外に、「単位データの処理時間の分散 (ばらつき)」・「データ領域に於ける単位データの処理時間の分布状態」などを考慮する必要がある。

最後に、「④非同期データ分配型」の並列処理部分のデータ処理時間の期待値を求める。期待値を求めるにあたって、計算を簡単にするために、次の二点を前提とする。

- 前提 1 : データ受信時間 ( $T_{receive}$ ) ・データ送信時間 ( $T_{send}$ ) は、master プロセスの処理状態に関わらず一定であるとみなす。
- 前提 2 : 全プロセスの実行時間を全て同じとみなす。 (= T1) また、全データ数 (D) は、送信データ単位 (m) で割り切れるものとみなす。

これより、「④非同期データ分配型」の並列処理実行部分の実行時間 (T1) の期待値は、以下の通りになる。

$$\begin{aligned}
 \Delta t \times m + T_{send} &+ (m \times t_{ave} + T_{receive}) \times k_1 = T_1 \\
 (\Delta t \times m + T_{send}) \times 2 &+ (m \times t_{ave} + T_{receive}) \times k_2 = T_1 \\
 &\vdots \\
 (\Delta t \times m + T_{send}) \times (n-1) &+ (m \times t_{ave} + T_{receive}) \times k_{n-1} = T_1 \\
 (\Delta t \times m + T_{send}) \times \frac{D}{m} &+ (m \times t_{ave} + T_{receive}) \times k_n = T_1 \\
 k_1 + k_2 + k_3 + \dots + k_{n-1} + k_n &= \frac{D}{m}
 \end{aligned}$$

注)  $k_i$  ( $i=1 \sim n-1$ ): 各 worker プロセスのデータ受信回数  
 注)  $k_n$ : master プロセス内の worker プロセス処理でのデータ受信回数

$$T_1 = \left( \frac{D}{2} + \frac{D}{n \times m} \right) (\Delta t \times m + T_{send}) + \frac{D}{m} \left( t_{ave} + \frac{T_{receive}}{n} \right)$$

実行開始時の遅れによる時間 & master プロセスのデータ分配処理時間

(n × データ数が十分大きいので (D) > 0) の場合、

$$T_1 \approx \frac{D}{m} \left( t_{ave} + \frac{T_{send} + T_{receive}}{n} + \Delta t \right) \dots \text{式 g 1}$$

式 h 5,  $T_{send} \approx T_{receive}$  より、

$$T_1 \approx \frac{D}{m} \left( t_{ave} + \frac{2 \times T_0}{n} + \Delta t \right) \dots \text{式 g 2}$$