

ジョブ間並列同期転送機能(PREST) の開発と評価

長須賀 弘文¹ 新井 利明¹ 今居 和男² 吉澤 康文¹
¹日立製作所システム開発研究所 ²日立製作所ソフトウェア開発本部

オンライン業務の長時間化に伴い、バッチ処理業務に割当てられる時間が制限されつつある。そこで、バッチ処理業務の高速化を目的に、ジョブ間並列同期転送機能(PREST)を提案し、開発した。PRESTは、ジョブ実行の並列度を向上させ、ジョブ間でのデータ受渡しのための入出力操作を削除する機能である。PRESTは、並列に実行しているジョブ間でのデータの受渡しを可能とし、データ受渡しのための入出力操作をデータ転送でシミュレーションすることで、バッチ処理の処理経過時間を短縮する。性能評価をした結果、データ入出力操作にかかるCPUオーバーヘッドを約1/6に削減し、処理時間を約1/2に短縮した。

THE DESIGN AND EVALUATION OF PREST

Hirofumi Nagasuka¹ Toshiaki Arai¹ Kazuo Imai² Yasufumi Yoshizawa¹

¹System Development Laboratory, Hitachi, Ltd.
(1099 Ohzenji Asao-ku Kawasaki-shi, 215 Japan)

²Software Development Center, Hitachi, Ltd.
(5030 Totsuka-cho Totsuka-ku Yokohama-shi, 244 Japan)

We proposed Parallel REference and Synchronous Transfer feature (PREST) which shortens batch execution time. PREST simulates I/O operation which is requested by a job (output job) to transfer temporary data to another job (input job). PREST allows the input job to execute in parallel with the output job. We have developed PREST, and got the following result: (1) CPU overhead for I/O operation is reduced to 1/6. (2) The execution time of a typical batch job, which sorts out and merges data, was shortened to 1/2.

1. はじめに

金融機関等では、オンライン業務のサービス時間延長に伴い、集計処理等のバッチ処理業務に割当てることができる時間帯が制限されつつある。一方で、データ量の増大、業務の複雑化にともないバッチ処理量が増加している。そのため、翌日のオンライン業務の開始までに、バッチ処理業務が終了しない危険性が大きくなりつつある。これを回避するために、バッチ処理業務の高速化を図ることが急務である。

一般的に、集計処理等のバッチジョブは、複数の処理で構成されている。各処理では、入力したデータを加工して後続の処理に渡しなが、全体の処理を進めている。大量のデータを取り扱うバッチジョブでは、ある処理とある処理でのデータの受け渡し（中間データ引継ぎ処理）のための入出力操作に時間がかかる。この傾向は、取扱うデータが大量になるにつれ、顕著である。また、中間データを引継ぐために、ジョブ実行を逐次的にさせる必要がある。

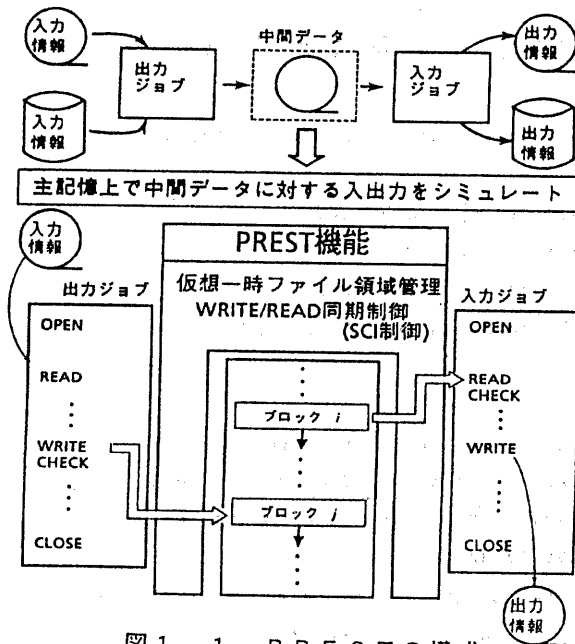


図 1. 1 PREST の構成

本稿で報告するジョブ間並列同期転送機能 (PREST: Parallel REference and Synchronous Transfer feature) は、ジョブ実行の並列度を向上させると共に、ジョブ間でのデータ引継ぎに伴う入出力操作を削除することで、バッチジョブの処理時間を短縮する機能である。(図 1. 1、図 1. 2 参照)

PRESTは、ジョブ間で保存する必要のないデータを、逐次アクセス法を利用して引継ぐ処理を高速の対象とする。データを出力するジョブ（出力ジョブ）と、その出力されたデータを入力するジョブ（入力ジョブ）を並列に実行させることで、処理時間を短縮する。データの引継ぎのための入出力操作は、全てメモリ上でのデータ転送でシミュレーションする。これにより、入出力操作を削除する。

2. PRESTの特徴

PRESTと同様の機能として、UNIXのPIPE機能がある[2]。

PRESTは、汎用大型計算機用OSの上で実現した。汎用大型計算機用のOSには、ファイルへのアクセス方法が複数ある。その中で、バッチ

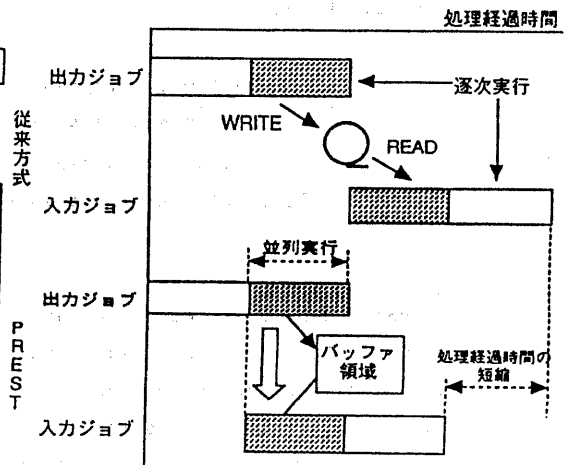


図 1. 2 PREST の効果

処理で頻繁に利用される逐次アクセス法に着目しこのアクセス法を用いたデータ引継ぎ処理の高速化を図ることとした。また、既存プログラムを改造することなしに、PRESTの適用を可能とすることで、汎用性をもたせた。さらに、PRESTの指定を、ジョブ起動用のジョブ制御文で指定するようにしたことで、柔軟性をもたせた。[1]

内部処理では、データ引継ぎのための出力ジョブと入力ジョブとの待ち合わせ処理（同期処理）回数の削減を図るSCI制御(Synchronous Control for Intermediate data)や、メモリの有効活用を図った動的バッファ容量決定方式に特徴がある。

3. PREST

PRESTの実現にあたっての課題と、その解決策について述べる。

3.1 データの転送方式

PRESTでは、入出力操作をデータ転送でシミュレーションする。そこで、この転送用のマクロ命令を新設し、ユーザプログラムにマクロ命令を発行させる方式を考られるが、この方式をとると、PRESTを利用するために既存のユーザプログラムを改造しなければならない。そこで、ユーザプログラムが発行した入出力命令（マクロ命令）をトラップして、データ転送に置き換える方式をとった。これにより、既存のユーザプログラムを書き換えることなしに、PRESTを利用することができる。

また、出力ジョブと入力ジョブは、異なるアドレス空間上で動作している。空間間で直接データ転送を行うと、空間切り替え処理などでCPUオーバーヘッドが大きくなる恐れがある。そこで、出力ジョブと入力ジョブが共に参照できる領域にバッファを設け、出力ジョブが出力命令を発行するとそのバッファ領域にデータを書き込み、入力ジョブが入力命令を発行するとそのバッファ領域からデータを読み込む方式をとった。

3.2 ジョブ間の同期処理

並列に実行しているジョブ間でデータの受け渡しをする場合、例えば、入力ジョブが入力命令を発行したときに、出力ジョブがデータを出力していなければ、出力ジョブがデータを出力するまで、入力ジョブを待たせなければならない。このような待ち合わせ処理を同期処理と呼ぶ。

PRESTでは、データ引継ぎのための出力ジョブと入力ジョブの同期処理を行い、データ引継ぎを保証する。

3.3 バッファ領域の管理

出力ジョブから入力ジョブへのデータ引継ぎを、バッファを介したデータ転送で実現する。そこで、割当てたバッファの容量が問題となる。最初に、取り扱うデータ量分のバッファを確保すると、大量データを取り扱う場合にメモリアーヘッドが膨大になり、取り扱えるデータ量も制限される。

そこで、比較的小容量のバッファ領域を確保し、ラップアラウンドに使用することとした。これにより、取り扱えるデータ量にも制限がなくなる。

3.4 PREST 指定方法

PRESTを利用するジョブに関する情報等を予めシステムに登録する方法も考えられるが、PRESTの利用を容易とするため、ジョブ制御文で指定することとした。（図3.1参照）

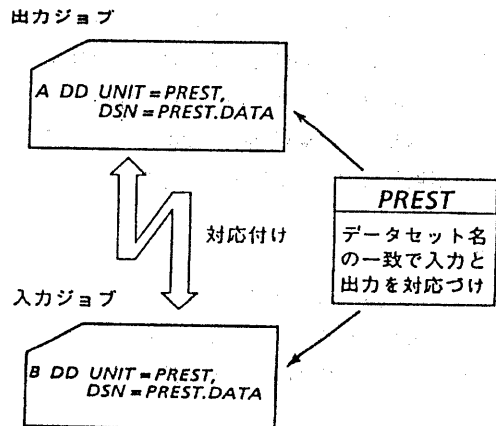


図3.1 PRESTの指定方法

4. PRESTの制御方式

本章では、同期処理方式とバッファ容量の管理方式について述べる。

同期処理に関しては、出力ジョブと入力ジョブが入出力命令を発行する度に同期をとると、そのオーバーヘッドは膨大になるという問題点が残っている[2]。また、バッファ容量に関しては、どの程度のバッファ領域を用意すればよいのかが明確でないという問題点が残っている。

そこで、同期処理の頻度が小さく、かつ小容量のバッファを用意することで、効果的にデータ引継ぎ処理を実現できる同期処理方式とバッファ容量の管理方式を求めるために、次節に示すモデルを用いて解析した。

4. 1 解析モデル

【モデル】

- (1) バッファ領域は、先頭ブロックから順次使用される。
- (2) 出力ジョブがバッファ領域に書き込んだデータを、入力ジョブがすべて読み込んだ後、入力ジョブは、出力ジョブがmブロック分のデータを書き込むまで待つ。
- (3) 出力ジョブがバッファ領域のすべてのブロックにデータを書き込み、入力ジョブが、バッファ領域内のすべてのデータを読み込んでいない場合、出力ジョブは、入力ジョブがmブロック分のデータを読み込むまで待つ。
- (4) 出力ジョブの書き込み要求と、入力ジョブの読み込み要求には、相互関係はない。
- (5) 目標とする同期処理の頻度は、 α 回のデータ引継ぎに対して1回未満である。

【記号】

- λ : 書き込み頻度 (回/sec)
 μ : 読み込み頻度 (回/sec)
 n : バッファ領域のブロック数 (ブロック)

【検証1】書き込み頻度が読み込み頻度よりも多い場合 ($\lambda > \mu$)

出力ジョブが待ち状態になり、入力ジョブがmブロックのデータを読み込んだ後、出力ジョブの待ち状態を解除したとする。待ち状態を解除した後、再び出力ジョブが待ち状態になるまでの時間 t は、式4. 1で求められる。

$$\begin{aligned} \lambda t &= \mu t + m \\ \therefore t &= m / (\lambda - \mu) \quad \dots \text{式4. 1} \end{aligned}$$

この時間 t が長ければ、同期処理を行う頻度が小さくなる。バッファのブロック数は n であるので、 t の最大値は、 $m = n$ の時である。

これは、一度待ち状態になった出力ジョブは、入力ジョブがバッファ内のデータをすべて読み込むまで待ち状態になっていることが、同期処理の頻度を最小にすることを意味する。

さらに、時間 t の間に、 α 回以上のデータ書き込みを行うことができれば、 α 回のデータ書き込みに対して、1回未満のジョブジョブ間の同期処理で済ませることができる。この目標を満たす最小のバッファ容量 (ブロック数) n^* は、式4. 2で求められる。

$$\begin{aligned} \lambda t &= \lambda n / (\lambda - \mu) > \alpha \\ \therefore n^* &\geq \{ \alpha (\lambda - \mu) / \lambda \} + 1 \quad \dots \text{式4. 2} \end{aligned}$$

【検証2】読み込み頻度が書き込み頻度よりも高い場合 ($\lambda < \mu$)

入力ジョブが待ち状態になり、出力ジョブがmブロックのデータを書き込んだ後、入力ジョブの待ち状態を解除したとする。待ち状態を解除した後、再び入力ジョブが待ち状態になるまでの時間 t は、式4. 3で求められる。

$$\begin{aligned} \mu t &= \lambda t + m \\ \therefore t &= m / (\mu - \lambda) \quad \dots \text{式4. 3} \end{aligned}$$

この時間 t が長ければ、同期処理を行う頻度が小さくなる。バッファのブロック数は n であるので、 t の最大値は、 $m = n$ の時である。

これは、一度待ち状態になった入力ジョブは、出力ジョブがバッファ内のすべてのブロックにデータを書き込むまで待ち状態になっていることが、同期処理の頻度を最小にすることを意味する。

さらに、この時間 t の間に入力ジョブが、 α 回

以上のデータ読み込みを行うことができれば、 α 回のデータ読み込みに対して、1 回未満のジョブ間の同期処理で済ませることができる。この目標値を満たす最小のバッファ容量（ブロック数） n^* は、式 4. 4 で求められる。

$$\mu t = \mu n / (\mu - \lambda) > \alpha$$

$$\therefore n^* \geq (\alpha (\mu - \lambda) / \mu) + 1 \quad \dots \text{式 4. 4}$$

式 4. 2 と式 4. 4 から、書き込み要求頻度と読み込み要求頻度に応じて、式 4. 5 に示す容量（ブロック数）のバッファ領域を確保すればよいことになる。

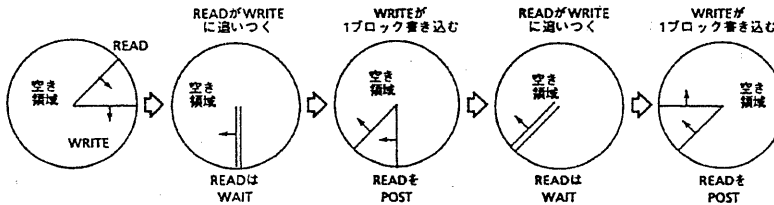
$$n^* \geq \frac{\alpha |\lambda - \mu|}{\max(\lambda, \mu)} + 1 \quad \dots \text{式 4. 5}$$

以上のモデルを用いた解析をもとに、次節以降で示す同期処理の方式と、バッファ容量の管理方式を採った。PREST で用いた同期処理の方式を SCI 制御と呼び、バッファ容量の管理方式を動的バッファ容量決定方式と呼ぶ。

4. 1 SCI 制御

SCI 制御では、与えられたバッファ容量のもとで、同期処理の頻度が最小になるように制御する。すなわち、ジョブの待ち状態の解除を、以下に示すタイミングで行う。（図 4. 1 参照）

[問題点] 最悪の場合、データの引継ぎの度に同期処理を実行



SCI 制御による解決

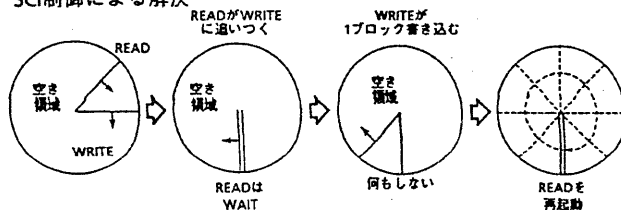


図 4. 1 SCI 制御

- (1) 出力ジョブがデータ出力要求を出した時に、バッファ領域に空きがなかった場合、入力ジョブが、バッファ内のすべてのデータを入力するまで出力ジョブを待ち状態にする。
- (2) 出力ジョブがまだ出力していないデータに対して入力ジョブが入力要求を発行した場合、出力ジョブがバッファ内のすべての領域にデータを出力するまで、入力ジョブを待ち状態にする。

4. 3 動的バッファ容量決定方式

出力ジョブと入力ジョブの入出力要求頻度が大きく異なる場合、バッファが小容量であれば、SCI 制御を施しても、ジョブ間で同期をとらなければならない機会が多くなる。そこで、動的バッファ容量決定方式では、出力ジョブと入力ジョブの入出力要求頻度の割合に応じて、バッファ容量を決定することで、メモリの有効活用を図ると共に、同期処理にかかるオーバーヘッドをある一定値以下に抑える。

すなわち、出力ジョブと入力ジョブの入出力要求頻度の割合から、同期処理と同期処理の間に引継がれるデータ数がある設計値以上となるために必要なバッファ容量を求め、増量する。

したがって、データの引継ぎ処理を行いながら、PREST は、下記の処理を行い、メモリの有効

活用を図り、同期処理の頻度を小さくする。

- (1) 初期割当て量として、一定容量のバッファ領域を割当てる。
- (2) 同期処理と同期処理との間に引継がれたデータ数をカウントする。
- (3) (2) でカウントしたデータ数がある設計値未満で

あれば、設計値を満たす最小のバッファ容量に増量する。

本開発では、最初に10ブロックのバッファを割当て、同期処理の頻度が、50回のデータ転送に対して1回未満となるような、最小のバッファ容量に増量するようにした。これにより、1回の入出力操作に占める同期処理のCPUオーバーヘッドが十分小さくなると考えたからである。

表4. 1は、50回のデータ転送に対して、同期処理の頻度を1回未満にするために必要なバッファ容量を示したものである。

表4. 1 入出力頻度とTRAD領域のブロック数の関係

λ	μ	10	20	30	40	50	60
10	1	1	26	35	39	41	43
20	26	1	18	26	31	35	
30	35	18	1	14	21	26	
40	39	26	14	1	11	18	
50	41	31	21	11	1	10	
60	43	35	26	18	10	1	

ここで、 λ : 出力要求頻度(1/sec)

μ : 入力要求頻度(1/sec)

5. 評価

5. 1 単体入出力にかかるCPUオーバーヘッド

1ブロックのデータの入出力操作をシミュレーションするために要するCPUオーバーヘッド(I/O処理ステップ数)と、入出力装置へのI/O処理ステップ数を測定した。その結果、I/O処理ステップ数を約1/6に削減した。(図5. 1参照)

図5. 1において、同期処理用のCPUオーバーヘッドは、SCI制御による同期処理を実行した時のオーバーヘッドである。バッファ管理用のCPUオーバーヘッドは、データ引継ぎの実行中にバッファ領域を拡張するか否かを判断する時に、出力ジョブと入力ジョブの実行をシリアルライズさせるために必要なCPUオーバーヘッドである。

図5. 1で示した同期処理用CPUオーバーヘッ

ドとバッファ管理用CPUオーバーヘッドは、1回の入出力操作当りに換算したものである。

これら、同期処理用のCPUオーバーヘッドとバッファ管理用のCPUオーバーヘッドは、PRESTのI/O処理ステップの約35%を占めている。

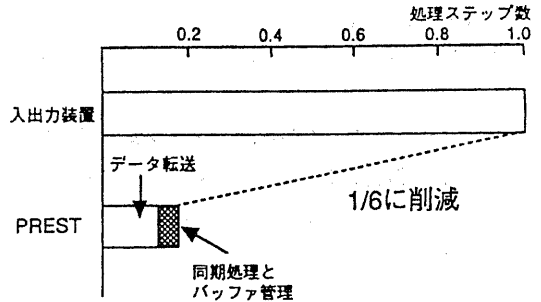


図5. 1 I/O処理ステップ数の比較

5. 2 SCI制御と動的バッファ容量決定方式の評価

マルチプロセッサ環境下とユニプロセッサ環境下において、PRESTを用いたデータ引継ぎ処理を行うパッチジョブを実行させ、同期処理の頻度と、データ引継ぎ中のバッファ容量の変化の様子を観測した。最終的なバッファ容量と、同期処理を実行した頻度を表5. 2に示す。

表5. 2 バッファ容量の変化と同期処理頻度

	ユニプロセッサ	マルチプロセッサ
バッファ容量	10⇒50	10
同期処理頻度	1/48.9	1/174.8

同期処理の頻度は、設計値に近い値かもしくはは設計値を満たす値となっている。ユニプロセッサでは、データ引継ぎ中にバッファを50ブロックに増量したのに対し、マルチプロセッサは、初期割当てである10ブロックのまま、データ引継ぎを完了した。この違いは、以下のことが要因となっていると思われる。

(1) ユニプロセッサの場合

CPUが一つのため、ある瞬間に実行可能なジョブは、出力ジョブか入力ジョブのどちらか一方である。したがって、出力ジョブが、全てのバッファ

領域にデータを書き込んで待ち状態になり、次に入力ジョブがバッファ領域の全てのデータを読み込んで、出力ジョブの待ち状態を解除するといった処理を繰り返していた。このため、同期処理を行う頻度が多くなり、途中でバッファ容量を増量する処理を行った。

(2) マルチプロセッサの場合

CPUが複数あるため、出力ジョブと入力ジョブが、並列に実行することが可能である。したがって、同期処理を行う頻度が小さく、初期割当てのバッファ容量のまま、データの引継ぎ処理を完了した。

5.3 バッチジョブのモデルを用いた評価

バッチ処理モデルを用いて、従来の処理方式の場合とPRESTを用いた場合との処理経過時間を比較した。

(1) 評価モデル

評価に用いたジョブは、店舗での商品別売上げの集計処理を想定したものである。(図5.2参照)

取り扱うデータは、2種類である。

(a) 入力データ：商品毎の売上げ数

(b) マスタデータ：商品の在庫数

処理内容は、以下の通りである。

(a) ソート処理：入力データをソートする。

(b) マージ処理：入力データとマスタデータから、現在の在庫数を求め、新しいマスタデータとして出力する。

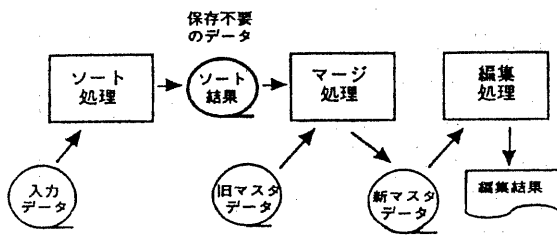


図5.2 評価モデルの構成

(c) 編集処理：新しいマスタデータから、品切れ状態にある商品を探し、編集結果を出力する。

このバッチ処理において、ソート処理とマージ処理との間でのデータ引継ぎ処理にPRESTを適用し、処理経過時間を比較した。

(2) 処理経過時間の比較

取り扱うデータ量と、データ引継ぎ用に使用する記憶媒体を変えて、各々の処理経過時間を比較した結果を図5.3に示す。

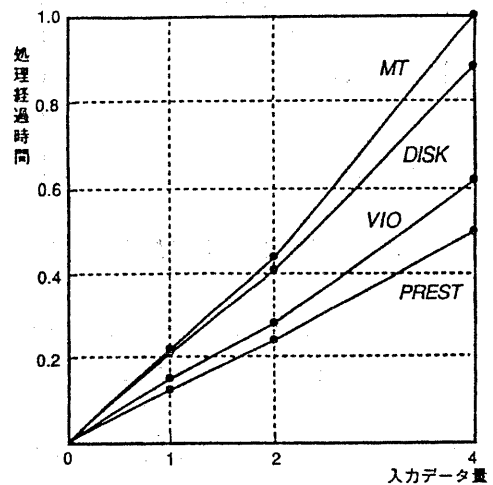


図5.3 処理経過時間の比較

図5.3において、データ量4の場合の処理経過時間の内訳を図5.4に示す。

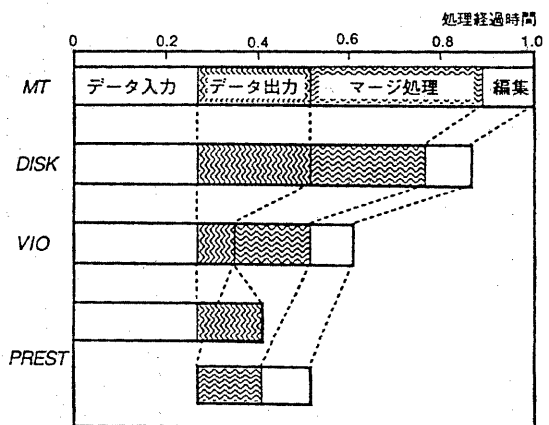


図5.4 処理経過時間の内訳

(3) データ引継ぎ処理部分の処理経過時間の比較

PRESTを適用したデータ引継ぎ処理部分に着目して、処理経過時間を比較する。(図5.5参照)

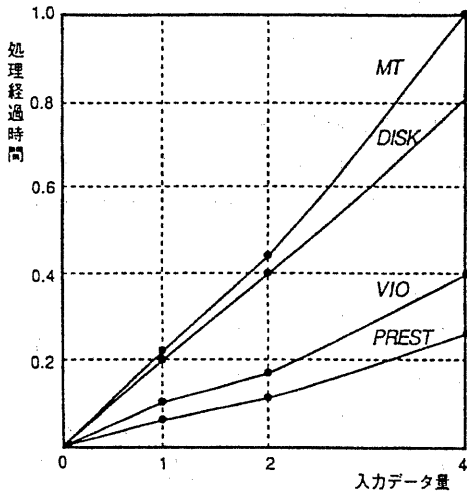


図5.5 データ引継ぎ処理の比較

PRESTとVIOを比較すると、約35%の処理経過短縮程度となっている。マージ処理では、入力データの入力、マスターデータの入力、および新マスターデータの出力と、3種類のデータの入出力を一緒に行っており、マスターデータの実入出力操作のために、処理時間の短縮がこの程度となっている。

6. バッチ処理業務への適用

開発したPRESTを、流通系のバッチ処理業務に適用し、その効果を実測した。

PRESTを適用したバッチ処理の規模は以下の通りである。

- (1) ジョブ本数：17
- (2) ジョブステップ数：487
- (3) PRESTの適用箇所：79

実測の結果、従来は約6時間20分かかっていたバッチ処理を、約4時間15分に短縮した。図6.1に、タイムチャートを示す。

7. おわりに

バッチ処理経過時間を大幅に短縮する機能であるPRESTを開発し、評価をした。その結果、ジョブ間でのデータ引継ぎに要するCPUオーバーヘッドを1/6に削減した。評価モデルとしたバッチ処理にPRESTを適用した結果、処理経過時間を約1/2(対MT比)に短縮した。

さらに、流通系のバッチ処理業務にPRESTを適用したところ、処理経過時間を従来の約30%短縮した。

【参考文献】

- [1] 長須賀他、ジョブ間並列同期転送機能の開発と評価 -VOS3/AS:PREST-、第41回情報処理学会講演論文集(4)、1990年9月
- [2] Maurice J. Bach, THE DESIGN OF THE UNIX OPERATING SYSTEM, 1986
- [3] P. J. Courtois, Concurrent Control with "Readers" and "Writers", Communication of the ACM, 1971

注)UNIXオペレーティングシステムは、UNIXシステムラボラトリ社が開発し、ライセンスしています。

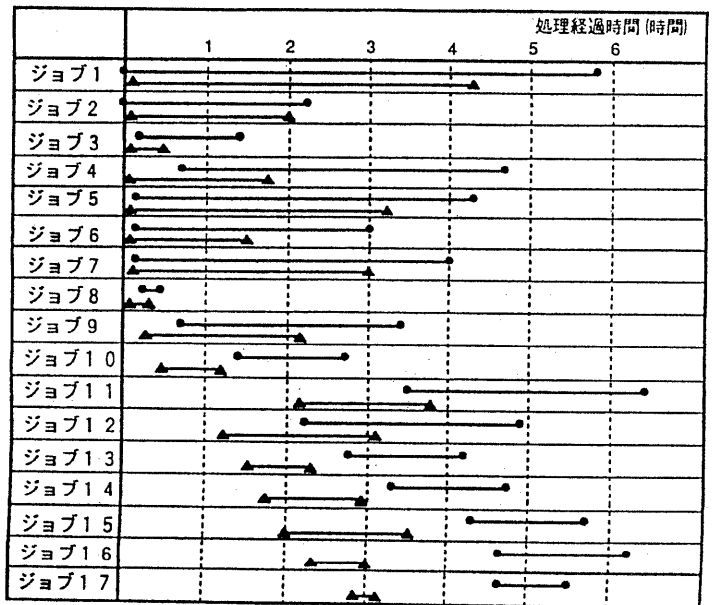


図6.1 バッチ処理業務のタイムチャート