

オンライン高速履歴情報取得制御方式

松 沢 尚 光

(株)日立製作所ソフトウェア工場

銀行・証券業界を中心とした第3次オンラインシステムの構築の動きに対して、システムの性能に重要な要因となる履歴情報の取得方法を構成・制御方式から改善検討した結果について述べる。履歴情報の目的別分割取得の概念を導入し、目的にそった記憶装置に、適切な時点で分割取得する方法を検討した。具体的には外部記憶装置に半導体記憶装置とカートリッジ型磁気テープ装置を適用し、記憶装置の特徴を生かしたことによって、高性能な履歴情報取得機能を実現した。この方法により履歴情報取得性能を向上でき、設備規模の面からも都市銀行の大規模システムに対応可能なものとなった。

On High Speed Method of Logging for Online Systems

Nobumitsu Matsuzawa

Software Works, Hitachi, Ltd.

5030 Totsuka-machi, Totsuka-ku, Yokohama, 244, Japan

This paper deals with a new method of constructing and controlling log information which improves the performance of online systems, and is especially suitable for large-scale banking and stock marketing systems. In this method, log information is divided into several categories and collected on separate storage devices respectively at appropriate time. It improves the performance of logging, and permits feasible size of hardware equipment supporting large-scale banking systems. IC-disk memory units and magnetic tape units of cartridge type are adopted to realize them.

1 まえがき

近年、銀行・証券業界は第3次オンラインシステムに向けシステムの再構築を図っている。特に都市銀行を中心とした大規模システムでは、システム規模の拡大に対応した信頼性の向上、性能の確保がDB/DC製品に要求されている。今回、システム全体の性能に対して重要な要因となる履歴情報の取得方式について大規模システムを対象とした新しい方式を考案し履歴情報取得能力の向上を実現した。この方式の構想から実現方式への展開、開発結果及び評価について述べる。

オンラインシステムの履歴情報としては、多様な種別があるが、ここでは特にトランザクションの実行内容を表わすデータベースの更新情報、メッセージの入出力情報などの時系列な順序を保証すべき情報を対象として考える。以下の文中では「ジャーナル」として説明する。

2 従来のジャーナル取得方式の問題点

高トラフィックなオンラインシステムではジャーナルの取得時間が性能上のボトルネックとなり易い。この問題点に対する解決策としては、ジャーナルを複数台の外部記憶装置に情報発生順にラウンドロビン方式で取得し負荷を分散する方式がある。この方式を用いても超大規模なシステムでは以下の問題点が発生する事が分った。

- (1) ジャーナル取得能力を向上するには、取得媒体（磁気ディスク装置）の台数が膨大になる。
- (2) 磁気ディスク装置に取得したジャーナルを磁気テープに保存するためのバックアップ操作が頻繁に発生し、オペレータの負担が大きい。
- (3) (2) のオペレータ運用中に不測の障害によりジャーナル取得媒体が満杯となり、オンラインダウンを招く危険性が高い。

上記の問題点が発生する原因は、ジャーナルの

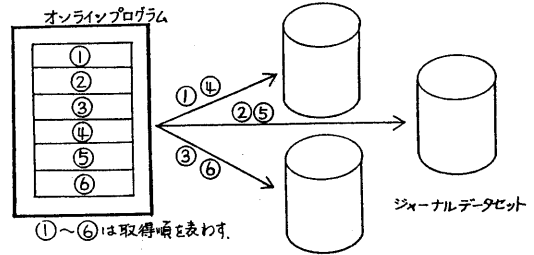


図-1 磁気ディスク装置の負荷分散

持つ複数の目的を一種類の記憶装置にトランザクション実行の特定時点で取得している事に係わっていると考えられる。

3 問題点の解決法

3.1 ジャーナルの目的

ジャーナルの使用される目的はオンライン、オフラインの両方で次のように分けられる。

(1) オンラインダウン時の処理途中トランザクションの回復

トランザクションの回復は、定期的を取得しているチェックポイントダンプを起点としているため、チェックポイントダンプ取得以降の必要な最低限の量が有れば十分である。しかしトランザクションの動作に同期して媒体に出力する必要がある。

(2) データベース障害回復・統計情報解析

データベース障害時のオフライン回復、統計情報解析の目的ではオンライン稼働中の全ての期間のジャーナルが必要である。ジャーナルの取得順で抜けが発生しなければ、媒体への出力タイミングはトランザクションの動作に同期している必要はない。

(3) ユーザ業務処理

オンライン中の業務処理結果をジャーナルとして出力しオフライン業務処理の入力に使用する。この目的では(2)と同様、ジャーナルの取得順で抜けが発生しなければ媒体への

出力タイミングはトランザクションの動作に同期している必要はない。

(4) ジャーナルの保存

トランザクションの実行状況を表わすジャーナルは一定期間保存しておき後で使用することがある。この目的では(2)、(3)の目的と重なり、オンライン稼働中の全ての期間のジャーナルが必要である。

上記のようにジャーナルの目的別にその特性を考慮すると、記憶装置の必要な特性が決まってくる。

(1) オンラインダウン時の処理途中トランザクション回復

高速小容量……不揮発性半導体記憶装置が適する

(2) データベース障害回復、統計情報解析

(3) ユーザ業務処理

(4) ジャーナルの保存

中速大容量……磁気テープ装置／磁気ディスク装置が適する

3.2 ジャーナルの構成

3.1の内容を実際のジャーナル構成へ反映することを考える。ジャーナルの使用目的別に適した装置へ、適した取得タイミングで分割取得する

(1) システム回復ジャーナル (SYJ)

オンライン中のトランザクションの回復に使用するために取得する。半導体記憶装置にラップアラウンド方式で回復に必要な範囲のジャーナルを保存しながら小容量のデータセットに取得する。小容量・高価な半導体記憶装置の適用が可能となる。又、トランザクション非同期出力方式の記録ジャーナル (LGJ) / アプリケーションジャーナル (APJ) のジャーナル情報が失われることがないように LGJ / APJ の未出力分を回復できるように管理する。

(2) 記録ジャーナル (LGJ)

LGJはデータベースの障害回復プログラムあるいは統計情報解析プログラムの入力用として取得する。トランザクションの実行状態とは同期せずに複数のトランザクション分のジャーナルをブロッキングして出力する。ブロッキングの効果によりSYJより低速の装置でも適用可能である。取得媒体としては磁気テープ装置が適し、オンライン稼働中の全期間のジャーナルを取得し磁気テープで保存する。

(3) アプリケーションジャーナル (APJ)

APJは業務処理固有の任意の情報を取得しバッチ業務処理の入力に使用する。磁気テープ装置へLGJと同じ方式で出力するか、磁気ディスク装置へ順編成出力する。

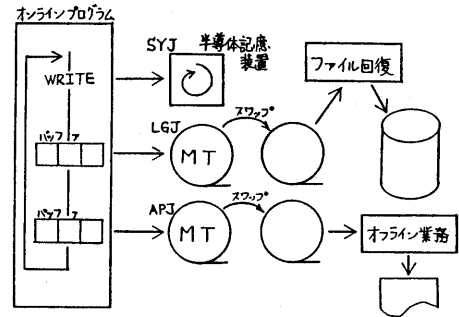


図-2 ジャーナルの構成

3.3 性能の検討

性能の検討にあたって次のトランザクションモデルと記憶装置仕様を基にした。

目標トランザクション件数 40万件/時
ジャーナル出力量

SYJ: 10KB/トランザクション
バッファ長 4KB×3

LGJ: 5KB/トランザクション
バッファ長 32KB

APJ: 5KB/トランザクション
バッファ長 32KB

ハードウェア

磁気ディスク装置

転送速度 3 MB/秒
 回転待ち時間(最大) 16.6ミリ秒

半導体記憶装置

転送速度 6 MB/秒
 アクセス時間先頭CI 0.76ミリ秒
 2番目以降CI 0.52ミリ秒

磁気テープ装置

カートリッジ型
 転送速度 2.5 MB/秒
 内部処理時間 1.0ミリ秒
 オープンリール型
 転送速度 1.39 MB/秒
 内部処理時間 1.94ミリ秒

磁気ディスク装置と半導体記憶装置のI/O時間を比較すると(SYJに使用した場合)

磁気ディスク装置 21.6ミリ秒

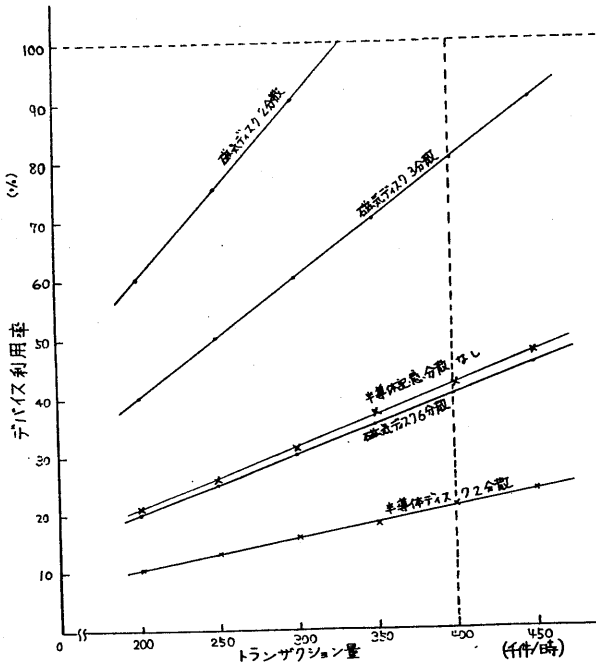


図-3 SYJのデバイス利用率

半導体記憶装置 3.8ミリ秒

半導体記憶装置では負荷分散を行わなくとも、磁気ディスク装置の6分散と同等の性能を確保できる。磁気ディスク装置を用いた場合、2分散では目標の40万件/時のトラフィック量でデバイス利用率が100%を越えてしまう。3分散ではデバイス利用率が80%程度で危険である。実用上5分散以上は必要と考えると半導体記憶装置の適用効果は大きく、目標のトラフィック量にも対応できる。

磁気テープ装置のI/O時間をLGJ/APJについて計算すると

カートリッジ型 13.8ミリ秒
 オープンリール型 25.0ミリ秒

LGJ/APJを同一チャネルに接続して使用すると目標の40万件/時のトラフィック量で、オープンリール型ではチャネル利用率が90%を超えてしまい危険である。チャネルをLGJ/APJで使い分けると、カートリッジ型磁気テープ装置の使用が必須となってくる。

このように3.2で検討したジャーナル構成で目標とするトラフィック量に対応可能なジャーナル取得処理を実現できる見通しがあった。

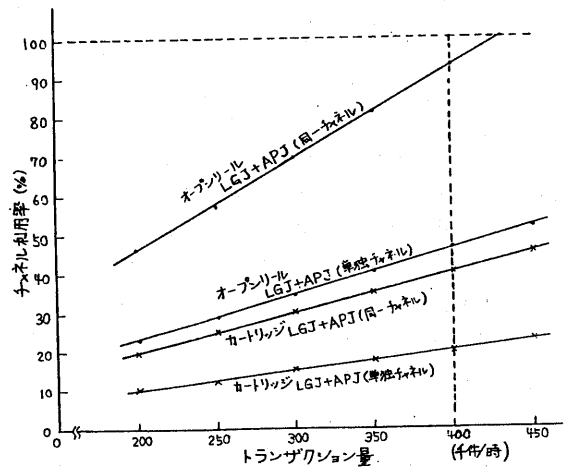


図-4 LGJ/APJのチャネル利用率

3.4 設備規模の検討

磁気ディスク装置に分散取得する従来のジャーナル取得方式と、目的別にジャーナル取得装置を使い分ける新しい方式とで記憶装置の容量を比較し設備規模の大きさについて検討した。

次のモデルシステムによる試算値を示す。

トランザクション件数 40万件/時

ジャーナル出力量

全ジャーナル：10KB/トランザクション

LGJ：5KB/トランザクション

APJ：5KB/トランザクション

ホストオンライン数 3オンライン

記憶装置容量

磁気ディスク装置

1200MB/ボリューム

カートリッジ型磁気テープ装置

約200MB/ボリューム

半導体記憶装置

64MB/ボリューム

〈従来方式〉 磁気ディスク装置分散使用

- a 磁気ディスク分散数……………5
- b 正副二重取得……………2
- c データセット数/ボリューム…………2
- d 世代数(現用・交代用・コピー待ち)……………6
- e 予備ボリューム数……………1
- f オンライン数……………3

$$\{a \times (d/c) \times b + e\} \times f = 93$$

(ボリューム)

ジャーナルデータセット用装置のみでこのような大規模な構成をとることは現実的には難しい。

〈新方式〉 SYJ

- a 半導体記憶装置分散数……………4
- b 正副二重取得……………2
- c データセット数/ボリューム…………2
- d 世代数(現用)……………1

- e 予備ボリューム数……………2
- f オンライン数……………3

$$\{a \times (d/c) \times b\} \times f + e = 14$$

(ボリューム)

LGJ

カートリッジ型磁気テープ装置

- a 正副二重取得……………2
- b 世代数(現用・交代用)……………4
- c オンライン数……………3

$$a \times b \times c = 24 \text{ (デッキ)}$$

APJ

カートリッジ型磁気テープ装置

- a 正副二重取得……………2
- b 世代数(現用・交代用)……………4
- c オンライン数……………3

$$a \times b \times c = 24 \text{ (デッキ)}$$

このように従来方式と新方式の装置構成の比較をすると大幅に低減を図ることが可能であり、設置スペースの削減に効果が期待できる。

3.5 ジャーナル制御方式

3.2のジャーナル構成にてジャーナルを取得する時の制御方式について述べる。

(1) SYJのラップアラウンド管理

SYJはラップアラウンド方式で古いデータに新しいデータを重ね書きしていく。このためにはオンラインダウン時のトランザクション回復に必要なジャーナルを消去することがないよう保証しなければならない。トランザクションの回復にはメモリの情報をチェックポイントダンプを基準点として回復しチェックポイントダンプ取得以降のジャーナルを用いて回復を行う

このためチェックポイントダンプ取得時点前のジャーナルは消去可能としてラップアラウンド管理する。SYJの容量はチェックポイントダンプ取得間のジャーナル量に依存する。

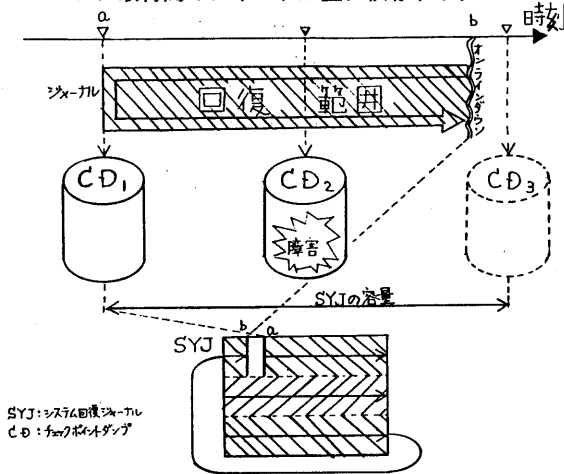


図-5 SYJのラップアラウンド管理

(2) カートリッジ型磁気テープ装置の適用

カートリッジ型磁気テープ装置はオープンリール型に比較し装置・媒体の小容積化、操作性の向上を図ることができるのと、磁気テープ制御装置(MTC)に大容量のバッファを内蔵しCPUの入出力要求とは非同期に媒体との入出力を行うことによって高速な入出力を実現できる特徴がある。

本来、ジャーナルの取得処理は磁気ディスク/磁気テープ等の外部記憶装置へジャーナルを出力要求すると媒体上へ反映されたことを保証して成立している。しかしカートリッジ型磁気テープ装置はCPUのプログラム処理とは非同期に媒体へ出力される。このためLGJ/APJに適用するにあたって、出力タイミングがトランザクションの動作に同期しないLGJ/APJの出力方式とは合致するが、次の問題点がある。

MTC内のバッファに残留していた分のデータが磁気テープ媒体へ出力されずに消失する可能性がある。

磁気テープ媒体/磁気テープ装置(MTU)障害の場合、MTC内に残留しているデータをMTCからCPUに読み戻せばデータの消失は防止できる。しかし、MTC/チャンネル障害の場合、MTC内の残留データを読み戻せないで、上記のデータ消失の問題が発生する。

この問題点の対策としてはLGJ/APJのバッファを3つの状態に分けて管理することにより実現できる。

- ① バッファ上にジャーナルデータをバッファリング中であり、MTCに対して出力要求していない。
 - ② MTCへ出力要求したが、MTCのバッファ上に残留して、まだ媒体上へ出力されていない。
 - ③ MTCへ出力要求しMTCのバッファ上から媒体上へ出力済になっている。
- ③のバッファは①に変更して新しいデータのバッファリングを続けられるが②のバッファは残留させて置くように管理する。ジャーナル出力

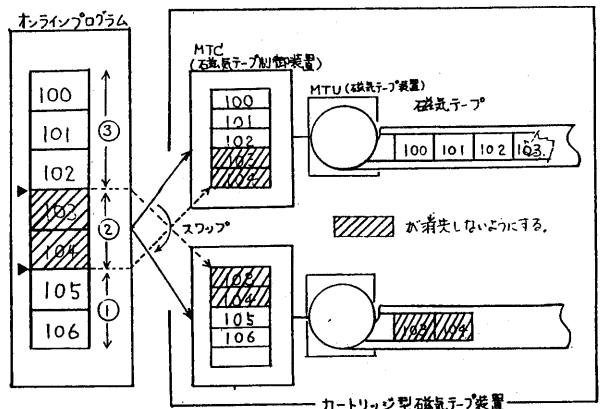


図-6 カートリッジ型磁気テープの管理

中にMTCより障害の報告を受けると、②のバッファをスワップ先に再出力する。これによってジャーナルデータの欠落を防ぐ。この方式の特長は、磁気テープのスワップ時間を短くすることができる点と、MTC/チャンネルの障害に対しても対応が可能な点である。

4 あとがき

高トラフィックなオンラインシステムに対しジャーナルの使用目的別に、適した記憶装置へ適した取得タイミングで分割取得する考え方は、性能の確保、設備規模の低減を図る有効な方法といえる。(株)日立製作所においては、TMS-4V/SP/J1 (Transaction Management System-4V/System Product/Extended Journal Function 1) の製品名で製品化を行い、昭和63年秋に本格稼働した都市銀行のシステム他で、性能面・設備面での実績を得ている。

社会のニーズはさらに大規模なオンラインシステムの実現へと拡大しているが、それに対応するための履歴情報の取得方式を本稿の内容を基に発展させていきたいと思う。

参 考 文 献

- 1) 金居他：履歴情報取得方式，公開特許公報，特願 昭64-21649号 (1989)。
- 2) TMS-4V/SP 解説，(株)日立製作所，(1988)。
- 3) TMS-4V/SP 操作，(株)日立製作所，(1988)。