

1. はじめに

分散型データベースシステム (DDBMS)は、通信コストの低減化、ハードウェアや通信回線等の故障に対する信頼性の向上など、優れた特徴を有している。[BeG 81][Koh 81][CeP 85]。しかし、解決すべき課題も多く、その1つに、トランザクションの並行処理制御の問題がある。

先に我々は、各サイトのスケジューラに先読みスケジューラ [IKK 85]を用い、時刻印を併用する並行処理制御アルゴリズムを提案した [Hal 85]。このアルゴリズムは、入力されたトランザクションは、決して拒否されず、再実行の必要もないという特徴を持ち、並行処理の効果を高め得ると思われる。

しかし、[Hal 85]ではサイト間でのデータの重複のないモデルを対象としていた。このことは、参照主体のトランザクションが様々な場所で発生するようなシステムにおいては(データを複数のサイトに置くことによって効率を上げることが出来るので)、実用上大きな制約となる。そこで本稿では、サイト間にデータの重複がある場合について、分散型データベースシステムのモデル、スケジュールおよびその直列可能性を考察し、アルゴリズムの一般化を行う。

2. 分散型データベースシステムモデル

ここでは、本アルゴリズムが対象とする分散型データベースシステムについて述べる。次節で述べる先読みスケジューラを用いた n 個のサイトのデータベースシステム(各データベースシステムをサイトと呼ぶ)が、通信制御機能を通じて結ばれている。各サイトは独立したデータベースシステムとして機能するが、必要に応じて他のサイトのデータにもアクセス出来るものとする。形式的に述べると、分散型データベースシステム (DDBMS)は データ項目の集合 $D = \{X, Y, \dots, Z\}$ と外部からシステムへ到着するトランザクションの集合 $T = \{T_0, T_1, T_2, \dots, T_n, T_r\}$ から成る。各データ項目 X について、そのコピー(実データ)が 適当な K 個のサイトに配置されている (K はデータ項目によって異なる)。簡単のためその様な実データの集合をデータ項目 X と同じ文字を用い、 $X = \{x_1, x_2, \dots, x_k\}$ と記す。今後、データ項目の集合を D 、そ

のコピー(実データ)の集合を \bar{D} と記す。各サイト a には、いくつかのデータ項目のコピーが配置されているが、そのデータ項目の集合を D_a とし、実データの集合を \bar{D}_a と記す。

トランザクション T_i は、いくつかの読取りステップ(read step) $R_i[D]$ ($B \subseteq D$) と書込みステップ(write step) $W_i[D]$ ($B \subseteq D$) とからなる系列である。読取りステップ $R_i[D]$ (書込みステップ $W_i[D]$) に対し、各 $X \in B$ の $R_i[X]$ ($W_i[X]$) を読取り操作(書込み操作)という。ある $x_k \in X$ に対して、 $\bar{R}_i[x_k]$ ($\bar{W}_i[x_k]$) は、実データ x_k の値を T_i に渡す (x_k の値を更新する) 操作で、実読取り操作(実書込み操作)と呼ばれる。トランザクション T_i の各ステップは、システムによって適当な実操作の集合に置き換えられた後(置き換えの詳細については後述)、実行に移される。 T_i は、その1つのステップの実行に際し、置き換えられた全ての実操作の実行が終了後、次のステップの実行要求を発することが出来る。 $T_0 = W_0[D]$ ($= \bar{W}_0[D]$) および、 $T_r = R_r[D]$ ($= \bar{R}_r[D]$) は 初期値を書き、最終値を読取る架空のトランザクションで、各々初期トランザクション最終トランザクションと呼ばれる。

3. スケジュール

本節では、(分散型を含む)データベースシステムにおけるスケジュールの直列可能性について述べる。

T 上のステップの集合を $STEP(T)$ とする。 T 上のスケジュール s とは、 $STEP(T)$ 内のすべてのステップからなる系列である。系列は時間の経過に応じて左から右へ並べて記す。ただし、 s の最初のステップは $W_0[D]$ 、最後のステップは $R_r[D]$ に決められている。 $STEP(T)$ 内のすべてのステップは、システムによって実操作の集合に置き換えられるが、得られた実操作の集合を $STEP(s)$ とする。 T 上の実スケジュール \bar{s} とは、 $STEP(T)$ 内の全ての実操作からなる系列である。

スケジュール s の読取り操作 $R_j[X]$ に対して、 $R_j[X]$ より左に位置する X への書込み操作の中で一番右のものが $W_i[X]$ である時、 T_j は X を T_i から読取ると言い、 $RF_i(R_j[X]) = i$ と記す。書込み操作 $W_i[X]$ が更新した X の値を読取る $R_j[X]$ が存在しない時、 $W_i[X]$ は無効書込みであると言う。また、 T 上の2つのスケジュール s

と s' において、全ての読取り操作 $R_j[X]$ に対して

$RF_s(R_j[X]) = i \Leftrightarrow RF_{s'}(R_j[X]) = i$
 が成立するならば、 s と s' は等価であると言う。

T上の実スケジュール \bar{s} の実読取り操作 $\bar{R}_j[x_k]$ に対して、 $\bar{R}_j[x_k]$ より左に位置する x_k への実書き込み操作の中で一番右のものが $\bar{w}_i[x_k]$ である時、 T_j は x_k を (または、 X を) T_i から読取ると言い、 $RF_{\bar{s}}(\bar{R}_j[x_k]) = i$ と記す。実書き込み操作 $\bar{w}_i[x_k]$ に対し、無効実書き込み操作を上と同様に定義する。T上の2つの実スケジュール \bar{s} と \bar{s}' において、全ての読取り操作 $\bar{R}_j[x_k]$ に対して

$RF_{\bar{s}}(\bar{R}_j[x_k]) = i \Leftrightarrow RF_{\bar{s}'}(\bar{R}_j[x_k]) = i$
 が成立するならば \bar{s} と \bar{s}' は等価であると言う。

直列なスケジュールとは、各ステップをトランザクション毎にまとめて、順番に並べたステップの系列である。 s と s' が等価で s が直列なスケジュールであるとき、 s' を直列可能なスケジュールと呼ぶ。いうまでもなく、直列可能性は、データベースの一貫性を保証する条件であって、データベース・スケジューラの目的は、外部から入力されたステップ系列を直列可能なスケジュールに変換することにある。

T上の実スケジュール \bar{s} とスケジュール s が以下の条件を満たすとき、実スケジュール \bar{s} はスケジュール s を実現するという。

(1) \bar{s} は $\bar{w}_0[D]$ で始まり、 $\bar{r}_r[D]$ で終わる。

但し、 $\bar{w}_0[D]$ ($\bar{r}_r[D]$) はすべての $x_k \in D$ に対する $\bar{w}_0[x_k]$ ($\bar{r}_r[x_k]$) を任意の順序で並べたものである。

(2) s 中の、 $j \neq 0, f$ なるトランザクションの各操作 $R_j[X]$ ($W_j[X]$) に対して、 \bar{s} 中にいくつかの $x_k \in X$ に対する実操作 $\bar{R}_j[x_k]$ ($\bar{W}_j[x_k]$) が存在する。

(3) \bar{s} 内の任意の $\bar{R}_j[x_k]$ に対し、

$RF(\bar{R}_j[x_k]) = i \Leftrightarrow RF(R_j[X]) = i$
 である。

\bar{s} が s を実現するとき、 \bar{s} と s においてトランザクション間の入出力関係は等しい。即ち、各トランザクションが読取る値と書込む値は、 \bar{s} と s において等しい。

4. TIOグラフと実TIOグラフ

本節では、DBMSの各実スケジュールにおいて、トランザクション間のデータの入出力関係を示す実TIOグラフと呼ばれるグラフを定義し、それに基づいて実スケジュールの直列可能性の判定法を示す。

T上のスケジュール s に対して、次の様なTIOグラフ (Transaction Input/Output graph)、 $TIO(s)$ を考える。 $TIO(s)$ の節点はT内のトランザクションに対応し、 s において $RF_s(R_j[X]) = i$ の時、トランザクション T_i の節点から T_j の節点へ X でラベル付けされた有向枝 $(T_i, T_j):X$ を付す。 T_i が X に関する無効書き込み $w_i[X]$ を有する時は、追加枝 $(T_i, T_i'):X$ を追加節点 T_i' と共に導入する。 $TIO(s)$ に他の節点、枝は存在しない。 $TIO(s)$ の節点間に次の2つの性質を持つ様な全順序を定めることが出来るとき、 $TIO(s)$ はDITSを持つという [IKM 83]。

(1) $T_i \ll T_j$ ならば T_j から T_i への有向路は、存在しない。

(2) 排除規則: $g \neq i$ を満たし、同じラベルを持つ2つの有向枝 (T_g, T_h) と (T_i, T_j) に対し、 $T_g \ll T_j$ ならば $T_h \ll T_i$ が成立する。

スケジュール s と $TIO(s)$ に関する次の定理は重要である。

定理4.1 [IKM 83] s が直列可能であることの必要十分条件は、 $TIO(s)$ が T_0 を最初として T_f を最後とする、DITSを持つことである。□

T上の実スケジュール \bar{s} に対して、次の様な実TIOグラフ、 $\bar{TIO}(s)$ を考える。 $\bar{TIO}(s)$ はやはりT内のトランザクションに対応する節点を持つ。 \bar{s} において $RF_{\bar{s}}(\bar{R}_j[x_k]) = i$ の時、トランザクション T_i の節点から T_j の節点へ x_k でラベル付けされた有向枝 $(T_i, T_j):x_k$ を付す。 T_i が x_k に関する無効実書き込み $\bar{w}_i[x_k]$ を有する時は、追加枝 $(T_i, T_i'):x_k$ を追加節点 T_i' と共に導入する。 $\bar{TIO}(s)$ に他の節点、枝は存在しない。 $\bar{TIO}(s)$ の節点間に次の2つの性質を持つ様な全順序を定めることが出来るとき、 $\bar{TIO}(s)$ はDITSを持つという。

(1) $T_i \ll T_j$ ならば T_j から T_i への有向路は、存在しない。

(2) 排除規則: $x_n, x_n \in X$ および、 $g \neq i$ を満たす2つの有向枝 $(T_g, T_h):x_n$ と $(T_i, T_j):x_n$ を考える。このとき、 $T_g \ll T_j$ ならば T_h

〈 T_i が成立する。

スケジュール s の場合と同様、実スケジュール \bar{s} とその実TIOグラフ $\bar{TIO}(s)$ の間に次の関係を示すことができる。

定理4.2 実スケジュール \bar{s} が実現するスケジュール s が直列可能であることの必要十分条件は、 $\bar{TIO}(s)$ が T_0 が最初で T_r が最後となるような、DITSを持つことである。

証明(⇒) s が直列可能ならば $TIO(s)$ はDITSを持つ。 \bar{s} は s を実現しているから、 $TIO(s)$ における有向枝 $(T_i, T_j):X$ に対し、 $\bar{TIO}(s)$ には有向枝 $(T_i, T_j):x_k$ (ただし、 $x_k \in X$) が1本以上存在する。また逆に、 $\bar{TIO}(s)$ の各有向枝 $(T_i, T_j):x_k$ に対し $TIO(s)$ には有向枝 $(T_i, T_j):X$ が存在する。ところで、 \bar{s} では、 $W_i[X]$ に応じて複数の x_k に対して $\bar{W}_i[x_k]$ を実行し、その一部のみを読取ることがあるので、 $\bar{TIO}(s)$ には $TIO(s)$ にはない追加枝と追加節点が存在する事がある。しかし、これら以外の有向枝は存在しない。従って、 $TIO(s)$ のDITSはそのまま $\bar{TIO}(s)$ のDITSになっている。(余分な追加枝および追加節点は、その始点の節点にまとめておけばDITSの条件を乱すことはない。)

(⇐) $\bar{TIO}(s)$ の各節点 T_i に対し、データ項目 X 毎に次の手続きを実行する。

- ① T_i からの、 $x_k \in X$ を満たす全ての有向枝 $(T_i, T):x_k$ が追加節点 T_i' へ入っている場合は、それらを1本の有向枝 $(T_i, T_i'):X$ に置き換える。
- ② $T_j (\neq T_i)$ からある $x_k \in X$ をラベルを持つ有向枝 $(T_j, T_i):x_k$ が存在する場合は、 T_i に入る全ての有向枝 $(T, T_i):x_1, x_1 \in X$ を1本の有向枝 $(T_i, T_j):X$ に置き換える。(この時、実現の条件(3)において、 $T_h (\neq T_j)$ から有向枝 $(T_h, T_i):x_1$ が存在する事はないので、 T_j を一意的に定めることができる。)

以上の手続きによってできたグラフは、 \bar{s} が実現するスケジュール s のTIOグラフ $TIO(s)$ である。 $\bar{TIO}(s)$ がDITSを持てば $TIO(s)$ もDITSをもつことは、手続きから明らかである。□

我々の目的は、直列可能なスケジュールを実現する実スケジュールを出力するアルゴリズムを構成することにある。直列可能なスケジュールを実現するという以上の概念は、1-SR(One-copy serializability) [BeG 83] と基本的に同一のものである。

5. 先読みスケジューラ

集中型データベースシステムの先読みスケジューラについては[IKK 85]等に詳しいが、ここでは、その動作の概要を述べ、8節で提案するアルゴリズムを特徴づけているいくつかの性質を示す。先読みスケジューラ(Cautious Scheduler)CS(C)は、全ての直列可能なスケジュールの集合SRのある部分集合Cに対し定義される。

各トランザクションTはシステム到着時に、Tの読取り集合(read set:Tが読むデータ項目の集合)と書込み集合(write set:Tが更新するデータ項目の集合)を宣言し、登録を行う。登録が済むとTはステップの実行要求をスケジューラに発することが出来る。各トランザクションから発せられた実行要求に対して、スケジューラは最終的にクラスCに属する直列可能なスケジュールが得られるように、それらの順序を定めて出力する(実行を許可する)。各ステップの実行可能性の判定を、完成テストと呼ぶ。

完成テストは、定理4.1に基づいてなされるが、用いられるクラスCによってはNP-完全になる場合もある。実的には、完成テストが多項式時間で実行できるクラスCが重要であって、例えば、 $C=WW$ をあげることが出来る[IKK 85][CaB 80][KaP 85]。このとき、具体的な手順は以下ようになる。

ある時点で既に出力されたステップの系列(部分スケジュール)をP、CS(C)に登録されているがまだ実行されていないステップの集合をPEND、判定すべきステップをqとする。このとき、まずPqとPENDに関する活性TIOグラフ(active TIO graph)を構成する。活性TIOグラフとはTIOグラフと類似の方法で以下のように定義されるグラフである。既に登録済みのトランザクション T_i を節点とし、Pq部分においてTjがXを T_i から読取っているならば有向枝 $(T_i, T_j):X$ を付す。Pq部分の無効書込み $W_i[X]$ は追加節点 T_i' への有向枝 $(T_i, T_i'):X$ で示される。PEND内の $W_i[X]$ は、追加節点 T_i' への有向枝 $(T_i, T_i'):X$ で示される。PEND内の $R_i[X]$ は、節点 T_i への端を持たない有向枝で示される(完成テストの手順では無視される)。得られたグラフをATIO(Pq, PEND)と記す。

つぎに、ATIO(Pq, PEND)に以下の有向枝を加えATIO_{new}(Pq, PEND)を構成し、最後にその排

除閉包（加えるべき排除枝がなくなるまで、排除規則を繰り返し適用してできるグラフ）を $AT10^*_{ww}(Pq, PEND)$ と記す。

- (i) Pq において $W_i[X]$ が $W_j[X]$ に先行するとき枝 (T_i, T_j) を付す（ ww -制約枝）。
- (ii) $W_k[X]$ を Pq における X への最後の書込み操作とする。このとき、すべての $W_j[X] \in PEND$ に対し枝 (T_k, T_j) を付す（ ww -制約枝）。
- (iii) $W_k[X]$ を(ii)と同様に定める。すべての $R_j[X] \in PEND$ に対し (T_k, T_j) を付す（ wr -制約枝）。

定理5.1 [IKK 85] 与えられた Pq と $PEND$ に対し、 $CS(WW)$ の完成テストは、 $AT10^*_{ww}(Pq, PEND)$ が有向閉路を持たないとき、かつその時に限り成功する。□

有向グラフが有向閉路を持つかどうかは多項式時間で判定できるので、 $C=WW$ における完成テストは多項式時間で実行できる。完成テストが成功すれば、 $CS(C)$ は q の実行を指示し、失敗すれば q の実行を遅延し、遅延待ち行列の最後に入れられる。なお、完成テストが成功した場合は、その都度遅延待ち行列内のすべての操作に完成テストを試みる。その結果、遅延待ち行列に置かれたステップも最終的に必ず実行される。

完成テストが多項式時間で実行可能なクラスとしては、他にも WRW [IKK 85]、 MWW 、 $MWRW$ [IKK 86]などが知られており、これらでは、適当に定義された活性 $T10$ グラフが有向閉路を持たないことが完成テスト成功の必要十分条件となっている。以後、これらのクラスにもとづく先読みスケジューラを仮定して議論を進める。

単一のサイトに用いられる集中型先読みスケジューラは、次の性質を持つ。

- (1) システム中でデッドロックを生じない。
- (2) ステップの再実行の必要がない。
- (3) 要求されたステップは、必ず有限時間内で実行される。

6. 分散型データベースシステムモデルの仮定

2節でその概略を述べた $DBMS$ をここで具体化する。

- (1) サイト数 $= n$ 。各サイトには1から n までのサイト番号が割り当てられている。
- (2) 各サイト a は先読みスケジューラ $CS_a(C)$ によって制御され、同一の機能をもつ。

- (3) 各サイトはクロックをもち、それに基づいて図6.1のような時刻印を発行することが出来る。但し、複数のトランザクションに同じ時刻印は与えないものとする。

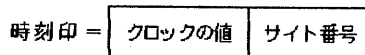


図6.1 時刻印

- (4) 各サイトは、他のサイトがどのような実データを持っているかを知っている。
- (5) サイト a で受け付けた各実操作は、それが実行されれば有限時間で終了する。実行が終了すると、それが他のサイト b から送られてきたものであるならサイト b にその結果を伝える。
- (6) サイト $a \neq b$ に対し、 $D_a \cap D_b \neq \emptyset$ を許す。
- (7) 各サイト a の構成は、図6.2の通り。

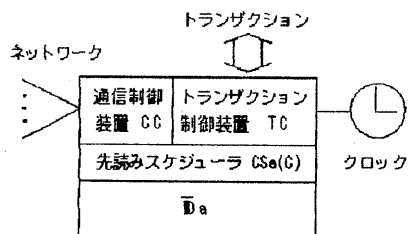


図6.2 各サイトaの構成

7. トランザクションの受けとめとその処理

前節のモデルにおいて、ユーザは次のようにトランザクション T_i を実行する。

- ① まず最寄りのサイトへトランザクション T_i をその読取り集合と書込み集合と共に登録する。
- ② トランザクション T_i の各ステップの実行要求を出す。この時、登録が完了して初めて最初のステップの実行要求を、また前のステップの実行が終了してはじめて、次のステップの実行要求を出すことが出来る。これに対して、トランザクション T_i を受け付けたサイト a は、その読取り集合と書込み集合を実データの集合へ次のように置き換える。
 - (1) $R_i[X]$ は $x_k \in X$ なる適当な実データ x_k 1個を選び、割り付ける（例えば、 $x_k \in X$ なる

実データをもつサイトの中で、通信費用最小のものを選ぶ)。

- (2) 書き込みデータXについてはすべての $x_k \in X$ を割り付ける。

このようにして得られた実データの集合を $\bar{D}(Ti)$ とする。あるbに対し $\bar{D}(Ti) \subseteq \bar{D}_b$ の時($a=b$ の場合もある)、Tiを局所的なトランザクション、そうでなければ全域的なトランザクションと呼ぶ。なお、このとき実データの集合 $\bar{D}(Ti)$ をうまく定めれば、局所的トランザクションと成り得るものに付いては、必ずその様な選択を行うものとする(局所的トランザクションは並行処理上有利であるため)。その後、ユーザからの①、②の要求に対してシステムは次節のように動作する。

8. DDBMSの並行処理制御アルゴリズム

A. 各サイトaにおけるトランザクション制御装置TCの動作

A-1 新しいトランザクションTiの受け付け:

TCは7節のような実データへの割り付けを行った後、以下の動作を行う。

- (1) 局所的なトランザクションの場合:
 $\bar{D}(Ti) \subseteq \bar{D}_b$ を満たすサイトbの $CS_b(C)$ にTiを登録する。

(2) 全域的なトランザクションの場合:

- (2-1) Tiにサイトaの時刻印 t' を与え、関係するサイトへTiと共に t' を送り、登録可能性を問い合わせる。

- (2-2) 関係する全サイトからの返答のなかで最大の時刻印を見つける。

それを t_{max} とすると、 t_{max} をTiの登録時刻印と定める。この時点で、サイトaのクロックの値が t_{max} による値より小さければ、そこまでクロックを進める。

A-2 トランザクションからのステップの実行要求の受け付け:

先に定めた実データへの割付に従って、実操作への変換を行った後以下の動作を行う。

- (1) 局所的なトランザクションの場合:
そのステップの実操作のすべてを関係するサイトbへ送る。
- (2) 全域的なトランザクションの場合:

そのステップの実操作それぞれを関係するサイトへ送る。この時、それがそのサイトへ送る初めての実操作であるならば、 t_{max} を共に送る。

A-3 他のサイトbからのトランザクションTiの登録可能性の問い合わせの受け付け:

自己のサイトaのクロックによる時刻印tとTiに付与された時刻印 t' を比べる。

- (1) $t' > t$ のとき: 自己のサイトのクロックを t' による値まで進め、実行可能であると返答し、 $CS_a(C)$ に (Ti, t') で仮登録。

- (2) $t' < t$ のとき: 自己の $CS_a(C)$ に t' で登録可能かどうか問い合わせる。可能ならその旨返答し、 (Ti, t') で $CS_a(C)$ に仮登録する。不可能ならt以上の時刻印であれば可能であると返答し、 (Ti, t) で $CS_a(C)$ に仮登録する。

A-4 他のサイトからの仮登録済みトランザクションTiの実操作の実行要求の受け付け:

実操作と共に時刻印 t_{max} が送られてくれば (Ti, t_{max}) で本登録し、自己のサイトの $CS_a(C)$ に渡す。実操作のみ送られてくれば本登録は済んでいるので、そのまま $CS_a(C)$ に渡す。

B. 先読みスケジューラ $CS_a(C)$ の動作

基本的には集中型データベースの場合と同一。但し $CS_a(C)$ が処理しているトランザクションのうち、全域的なトランザクション (Ti, t') と (Tj, t') に対しては、 $t' < t'$ ならば、TiをTjの前に直列化するという制約を付した上で(具体的には $CS_a(C)$ が保持する活性TIOグラフにおいて節点TiからTjへ制約枝を付す)、次のように扱う。

B-1 トランザクション (Ti, t') の(他サイトあるいは自己のサイトからの)登録要求の処理:

Tiの実操作のうち自己のサイトaに関するものを、 $CS_a(C)$ のPENDに入れる。このとき、活性TIOグラフに有向閉路ができなければ (Ti, t') のまま仮登録しておく。そうでなければ、時刻印 t' での登録を拒否し、サイトaの現在のクロックによる時刻印tを用いて、 (Ti, t) を仮登録する。

B-2 実操作qの実行要求の受け付け:

先読みスケジューラ $CS_a(C)$ の通常の手順に従い、4節で述べた完成テストをqに関

して行い活性 $\Pi_0(s)$ グラフに有向閉路ができなければ、 q を実行する。そうでなければ、即時の実行を拒否し q を遅延操作の待ち行列の最後に入れる。

9. アルゴリズムの正当性

8節の手順を6、7節で定義した分散型データベースのモデルに適用したとき、システムはデッドロックに陥ることなく直列可能なスケジュールを実現することを示す。

9.1 直列可能なスケジュール s の実現について

上述の DDBMS で実現される 実スケジュール \bar{s} とは、システム内で実行されたすべての実操作をそれが実行された絶対時刻に従って並べた系列であると 考えられる。このとき、複数の（異なるサイトで実行された）実操作が、同じ時刻の所に 並ぶならば、それらの順序は 任意に定めて一列に並べる。これは、分散型データベースシステムを、全ユーザから 1つの データベースシステムとして ながめることに 相当する。この実スケジュール \bar{s} 中の全ての 実読取り操作 $\bar{R}i[x_k]$ に対して2節の定義に従って、 $RF_i(\bar{R}i[x_k])$ の値を 定めることができる。同時に 実行された複数の 実操作は、すべて異なるサイトで 実行されたものであることから、このように 定められた RF の値は、各サイト a における $CS_a(C)$ が定めた値に一致する。 RF が定まると、 \bar{s} に対する $\Pi_0(s)$ を描くことができる。この $\Pi_0(s)$ が DITS を持つことを示せば、上述のアルゴリズムが出力する実スケジュールは、直列可能なスケジュールを実現する事がいえる。

定理9.1 8節のアルゴリズムによって 制御されている分散型データベースシステムが実スケジュール \bar{s} を出力したとすれば、その実 Π_0 グラフ $\Pi_0(s)$ は、 T_0 を最初とし T_r を最後とする DITS をもつ。

証明 上述の議論から、 $\Pi_0(s)$ は、各サイト a の実 Π_0 グラフ $\Pi_{0a}(s)$ のすべてを、各全域的トランザクション T_i の（いくつかのサイトに存在する）節点を一つにまとめつつ重ね合わせることで実現されることがわかる。ただし、 \bar{s}_a はサイト a の実スケジュールである。各サイト a の $\Pi_{0a}(s)$ は $CS_a(C)$ の動作から T_0 を最初に、 T_r を最後に置く DITS をもつことが保証されている。

これらの DITS と矛盾せず、かつ全域的トランザクションの順序は時刻印にしたがって並ぶように、 $\Pi_0(s)$ の全節点上に全順序 \ll を定めよう。これが可能なことは、

- (1) 全域的な トランザクション間には 時刻印による順序が全サイト共通に定められていて、各サイト a では、 a に関与する全域的トランザクション間に、時刻印による制約枝を付してスケジューリングしていることから、その順序に矛盾しないことが保証されていること、

および

- (2) 各局所的トランザクションは一つのサイトにのみ登場すること、
- を考えると明らかである。この \ll が T_0 を最初に、 T_r を最後におくことも明らかである。そこで \ll が $\Pi_0(s)$ の DITS であることを示そう。

定理4.2の前に示した DITS の性質(1)は明らかに成立する。なぜなら、 $\Pi_0(s)$ の枝は時刻印による制約枝の一部を除いて全てどれかの $\Pi_{0a}(s)$ に存在していたものであり、 \ll はその方向に矛盾しないように定められている。また、時刻印による制約枝は、上述のように必ず左から右へ向かうことが保証されているから、これも問題ない。

つぎに、DITS の性質(2)を示すために、 $\Pi_0(s)$ 内の有向枝 $(T_g, T_h):x_a$ と $(T_i, T_j):x_b$ を考える。ただし、 $x_a, x_b \in X$ かつ $x_a, (x_b)$ は サイト a (サイト b) の実データとする。 $a=b$ であれば、これら2本の有向枝は $\Pi_{0a}(s)$ に存在し、 $\Pi_{0a}(s)$ の DITS に矛盾しないように \ll が定められていることから、性質(2)が示せる。 $a \neq b$ であれば、書込み集合の実データへの割り付け方法(7節参照)からわかるように、データ項目 X への書込みを行っている T_g と T_i は全域的トランザクションであり、どちらもサイト a とサイト b の両方に関係する。従って、 $\Pi_{0a}(s)$ と $\Pi_{0b}(s)$ のどちらにおいても T_g と T_i の間には時刻印による制約枝がすでに付されている。一般性を失う事なく制約枝は (T_g, T_i) であるとしよう。すなわち、 $\Pi_0(s)$ において T_g から T_j へ有向路が存在し、 $T_g \ll T_j$ である。さて、 T_g と T_i の両節点は $\Pi_{0a}(s)$ に存在するから、DITS の性質(2)より、 $\Pi_{0a}(s)$ では T_g, T_h, T_i の順に左から右へ並んでいる。同様に $\Pi_{0b}(s)$ では T_g, T_i, T_j の順に左から右へ並んでいる。従って、 \ll はこの順序に矛盾しない

ので、 $(T_g \ll) T_h \ll T_i (\ll T_j)$ となり、 $T_i \in S$ において 性質(2)が成立する。□

9.2 デッドロックについて

各サイト個々の動作については、先読みスケジューラの性質として、デッドロックが生じないことが示されているので省略する。ここではシステム全体としてのデッドロックのみを考える。

いま、トランザクション T_i と T_j が $T_i = O_{i1} \dots O_{ir} \dots O_{iq} \dots$ 、 $T_j = O_{j1} \dots O_{jp} \dots O_{js} \dots$ であるとする。 O_{jk} 等は実操作を表す。このとき、サイト a において O_{jp} が O_{iq} の実行を待っていたとする。すなわち、サイト a の先読みスケジューラが O_{jp} の実行の前に O_{iq} の実行を要求した。(換言すると、 O_{jp} の前に O_{iq} を置かないと、先読みスケジューラの完成テストを満足しない)しかし、 T_i は O_{iq} の前に位置する O_{ir} (つまり $r < q$) の実行をサイト b で要求中であり、それが済まなければ O_{iq} の実行要求を出せない。そして、サイト b ではこの O_{ir} が何等かの実操作の実行を待ち、その実操作はまた別の実操作の実行を待つというふうにしてその先頭が O_{js} であったとする。ここで、 $p < s$ を仮定すると、 T_j の O_{js} は O_{jp} が実行されないと実行できず、結局、図9.1に示すようにこれらの実操作のそれぞれは自分自身の実行を待っていることになり、デッドロックを生じている。確かに、一般の先読みスケジューラでは、このような状況が考えられ、デッドロックが生じることがある。

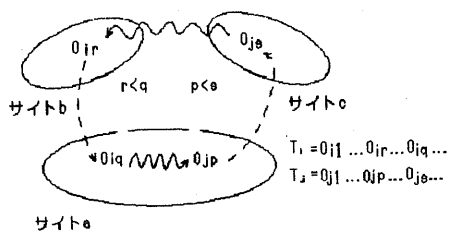


図9.1 システム全体にかかわるデッドロック

ところで、一般に先読みスケジューラでは、ユーザが宣言した書込み集合や読取り集合の一部を途中で取り消すと、スケジューラ自体がデッドロックに陥るといった取り消し異常を引き起こすことがある。そこで、実用的な先読みスケ

ジューラとして、取り消し異常を生じないものが提案されており、先に言及した $CS(WW)$ 、 \overline{CS} (WRW)、 $CS_m(MWW)$ 、 $CS_m(MWRW)$ 等は、そのような性質を持っている。これらの、先読みスケジューラでは、補題9.1に示すように (T_i, t') 、 (T_j, t'') 、 $t' < t''$ を考えたとき、 T_i の実操作 O_i は T_j の実操作 O_j の実行を決して待たないことがいえる。従って、各サイトに補題9.1の条件を満たす先読みスケジューラを用いれば、時刻印の小さなトランザクションの実操作がより大きな時刻印を持つトランザクションの実操作の実行を待つことはなく、デッドロックを生じないこと(定理9.2)が証明できる。

補題9.1 P を現在までにスケジューラ $CS(C)$ によって、出力された部分スケジュール、また、現在の $ATIO$ グラフにおける $DITS$ の順に従って、 $PEND$ 中のステップをトランザクション毎に直列に並べた系列を s' とする。このとき、 $Ps' \in C$ が常に成立するような先読みスケジューラでは、 (T_i, t') 、 (T_j, t'') 、 $t' < t''$ なる2つの全域的トランザクションを考えたととき、 T_i のステップ O_i は T_j のステップ O_j の実行を待つことはない。

証明 T_i から T_j へ時刻印による制約枝が付されているので、 $ATIO$ グラフのどのような $DITS$ においても、 $T_i \ll T_j$ が成立する。従って、 s' では O_i は O_j の前に位置する。 $Ps' \in C$ の性質によって、先読みスケジューラは s' の順序で実操作を処理でき、これは O_i が O_j の実行を待たないことを意味する。□

この補題9.1の性質 $Ps' \in C$ は、すべての先読みスケジューラで成立するわけではないが、上記の $CS(WW)$ 、 $\overline{CS}(WRW)$ 、 $CS_m(MWW)$ 、 $CS_m(MWRW)$ 等では成立することが知られている [IKK 85] [IKK 86]。

定理9.2 各サイトに補題9.1の条件を満たす先読みスケジューラを用いるとき、システムにデッドロックは生じない。

証明 上述の様なデッドロックが生じたとすれば、閉路中に存在する全域的トランザクション (T_i, t') と (T_j, t'') で、 $t' < t''$ にもかかわらず、 T_i が T_j を待つ (T_i の実操作が T_j の実操作の実行を待つ) という事態が必ず起こっている。これは、補題9.1に示したスケジューラの動作に矛盾する。□

9.3 時刻印の付与について

同一時刻印が、複数のトランザクションへ割当てられることはない。しかし、同一トランザクションへ複数の時刻印が割当てられることは起こり得る。この状態は、仮登録の手続き中に生じる。しかし、実行時（本登録）の時刻印はそのトランザクションを最初に受け付けたサイトにおいて決定され、各サイトでの実操作実行時にはその時刻印が知らされるので、実操作実行時の時刻印は各サイトで同一である。

各サイト a の $CS_a(C)$ は、トランザクション T_i の登録可能性を問われたとき、サイト a のその時点でのクロックによる時刻印を t とすれば、必ず (T_i, t) を登録する事ができる。これは T_i を受け付けたサイトでも同様である。なぜなら、サイト a におけるクロックの修正法 (A-1(2-2), A-3(1)) から、時刻印 t を発行したとき、 $CS_a(C)$ 内の全てのトランザクションは t より小さな時刻印を持つので、 $CS_a(C)$ 内では、 T_i を最後の位置に直列化するような DITS が必ず存在するからである。既に登録されているトランザクションの実操作を全て実行してから、 T_i の実操作を実行すれば必ず直列可能なスケジュールを実現する実スケジュールを出力できるからである（ただし、実際にその様にスケジュールするとは限らない）。

次に、 t' なる時刻印を持つ全域的なトランザクション T_i を既に仮登録している $CS_a(C)$ に対し、時刻印を $t (> t')$ に修正し、本登録を要求することがある (A-1(2-2), A-2(2))。このとき $CS_a(C)$ は (T_i, t) を受理する事を示す。

定理 9.3 $CS_a(C)$ として $CS(WW)$ 、 $CS(WRW)$ 、 $CS(MWW)$ 、 $CS(MWRW)$ を考えるとき、 t' なる時刻印を持つ全域的なトランザクション T_i を既に受理している $CS_a(C)$ に対し、時刻印を $t (> t')$ に修正し、本登録を要求するとき、 $CS_a(C)$ は (T_i, t) を受理する。

証明 最初に、これらのスケジューラにおける完成テストの方法を示す。

5 節に $CS_a(WW)$ の場合を述べたが、他の場合も含めて、 $ATIO$ グラフの排除閉包 $ATIO^*(Pq, PEND)$ が閉路を持たないことで完成テストの結果を判定できる。

このようなテストにおいて、 (T_i, t') の (T_i, t) への変更を考える。ただし、 $t > t'$ である。時刻印の変更によって、全域的トランザクシ

ョン T_g との順序が $T_i \ll T_g$ から、 $T_g \ll T_i$ になったとする。時刻印変更前の状態を、図 9.2(a) に示す。ただし、 T_i と T_g はともに実データ x_m への書き込み操作を持つ（その様な x_m がなければ、証明は簡単）。 T_i の時刻印が変更されたことは、その実操作がどれも未実行であることを意味する。また、 $ATIO^*(Pq, PEND)$ に閉路が存在しないことから、 T_g は $\bar{W}g[x_m]$ を未実行 ($C=WW$ の場合)、あるいは実行済みであるが、それを読んでいるものはない ($C=WRW$ の場合)。すなわち、 T_g から出ている x_m の枝は追加節点に終わっている。

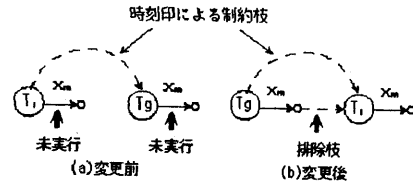


図 9.2 (T_i, t') から (T_i, t) への変更

時刻印を t へ変更すると、 T_i から出る枝は、追加節点へのものと時刻印によるもののみであるから、変更以前に比べると減少する。 T_i へ入る枝は (T_g, T_i) (時刻印による制約枝) が増える。その結果、図 9.2(b) に示すように、排除枝で T_i に入るものが出来ることがある。しかし、同図に示したように、この排除枝は T_g の追加節点からのものであるため、 T_g から T_i への経路が出来たにすぎず、制約枝 (T_g, T_i) と同じ効果を与えるだけである。 t への時刻印変更によって $ATIO^*(Pq, PEND)$ に閉路が出来たとする。上の議論から、この閉路は T_i に入る時刻印による制約枝と T_i から出る時刻印による制約枝を含んでいる。そこで、この閉路中、 T_i の前に T_g 、後ろに T_h があるとすると、つまり閉路は時刻印による制約枝 (T_g, T_i) と (T_i, T_h) を持つ。しかし、これは (T_g, T_h) が、 T_i の時刻印の変更以前から存在していた（つまり閉路が存在していた）ことになり矛盾が生じる。よって、 $CS_a(C)$ は (T_i, t) を必ず受理する。□

10. おわりに

本報告ではサイト間にデータの重複がある場合について、分散型データベースシステムモデルの構築、そこでのスケジュールの直列可能性

の意味づけ、およびトランザクションの並行処理制御アルゴリズムの提案を行った。

本稿の冒頭で述べたように、データのコピーを複数のサイトにおく目的として、参照主体のトランザクションが多数生ずるような環境の下での効率向上、サイト故障からのデータの保護等を上げることが出来る。これらに対して、ここで提案したアルゴリズムは

- (1) 全域的なトランザクションが生起しない場合、各サイトは独立に集中型の場合とほぼ同様の効率で動作する、
- (2) デッドロックを生じないため全域的なトランザクションの介入による効率の低下が比較的少ない、

といった望ましい特徴を有する。

今後、スケジュールの効率評価や障害時における対処、回復アルゴリズムを含めて、実用システムへの適用性の検討が必要と思われる。

謝辞

本研究を行うに当たり、熱心に御討論頂いた山下雅史助教授（広島大）、滝沢誠助教授（東京電機大）、橋口攻三郎助教授（豊橋技科大）、今井正治講師（豊橋技科大）、平田富夫講師（名大）に深謝します。なお、本研究は一部文部省科学研究費によるものである。

参考文献

- [BeG 81] P.A.Bernstein, N.Goodman: Concurrency control in distributed database systems, ACM Computing Surveys, vol.13 No.2, June 1981, pp.185-224.
- [BeG 82] P.A.Bernstein, N.Goodman: Multiversion concurrency control: Theory and Algorithms, ACM TODS, vol.8 No.4, December 1983, pp.465-483.
- [CaB 80] M.A.Casanova, P.A.Bernstein: General purpose schedulers for database systems, Acta Informatica 14, 195-220, 1980.
- [CeP 85] S.Ceri, G.Pelagatti: Distributed Databases Principles & Systems, McGraw-Hill International Student Editions, 1985.
- [Hal 85] 原嶋秀次, 茨木俊秀: Cautious Scheduler (先読みスケジューラ)を用いた分散型データベースシステムの並行処理制御, 情報処理学会データベース研究会 1985年7月.
- [IKK 85] T.Ibaraki, T.Kameda, N.Katoh: Cautious transaction schedulers for database concurrency control, Technical Report LCCR TR85-6, Simon Fraser University, 1985.
- [IKK 86] T.Ibaraki, T.Kameda, N.Katoh: Multiversion cautious schedulers for database concurrency control, Technical Report LCCR TR86-2, Simon Fraser University, 1986.
- [IKM 83] T.Ibaraki, T.Kameda, T.Minoura: Disjoint-interval topological sort: A useful concept in serializability theory, Proc. 9th International Conf. on VLDB, Florence Italy, 89-91, Oct/Nov 1983.
- [KaP 85] P.C.Kanellakis, C.H.Papadimitriou: The complexity of distributed concurrency control, SIAM J.Comput, Vol.14, No.1, February 1985.
- [Koh 81] W.H.Kohler: A survey of techniques for synchronization and recovery in decentralized computer systems, ACM Computing Surveys, vol.13 No.2, June 1981, pp.149-183.