

並行制御アルゴリズムのシミュレーション

中津 梢男

愛知教育大学総合科学課程情報科学コース

あらまし 構文情報だけを利用した並行制御アルゴリズムのうちで、2相ロックアルゴリズムと直列化可能性グラフを用いたBBアルゴリズムの性能比較をシミュレーションによって行った。

まず、アルゴリズムの不動点集合（並行制御アルゴリズムによって直列化可能と判断される履歴集合）の大きさを比較し、次に、Agrawalらのモデルに基づいて並行処理のシミュレーションを行った。不動点集合には顕著な差がみられ、処理結果では、競合の少ない場合には両者に差がなく、競合が増えるにつれBBアルゴリズムが高いスループットを示すことが確かめられた。

Simulation Results of Database Concurrency Control Algorithms

Narao NAKATSU

Department of Information and Computer Sciences, Aichi University of Education

Abstract Two kinds of concurrency control algorithms using only syntactical information are compared. One is the two phase lock algorithm and the other is a BB algorithm, which is one of serializability test algorithms. The BB algorithm outperforms a 2PL algorithm in highly conflicted cases.

1. はじめに

並行制御のアルゴリズムはデータベースシステムの効率に大きな影響を及ぼすため、これまでも多くの研究が継続して行われており、そうしたアルゴリズムの比較研究も数多く報告されている。並行制御アルゴリズムは、直列化可能な履歴を生成するものであり、2相ロックに代表されるロックを用いるアルゴリズム、時刻印を用いるアルゴリズム、楽観的アルゴリズムなどが代表的なものとして知られている^[2,4]。またこうしたアルゴリズムにはさまざまな変形が存在し、細部にわたる記述が不足しているため、それらを共通の土俵で比較することは困難である。

さらに並行制御のアルゴリズムが利用できると仮定している情報にもさまざまな仮定があって、一般にトランザクションに関する情報が多ければ多いほど能率の良いアルゴリズムが得られるものと考えられる。たとえば情報量の少ない順に、

(1) トランザクションによって発行される命令の種類がその都度わかる場合、(2) トランザクションがアクセスする項目とアクセスする形態が前もってわかっている場合、(3) トランザクションがアクセスする項目とアクセスする形態を及びその順序が前もってわかっている場合 と分類することができる。さらに、トランザクションをREAD/WRITEの列ではなく、意味をもつ小命令の列として考えた^[3]、トランザクションをより小さいトランザクションに分解する^[5]などのように、トランザクションの持つ意味を利用して並行性を高める研究も広く行われている。また、直列化可能性という条件はかなり厳しい条件であるため、アルゴリズムの適用分野を限定して、直列化可能性の定義を拡張する試みがなされている^[6,8]。

並行制御アルゴリズムの性能は、各アルゴリズムの不動点集合（そのアルゴリズムによって実行順を変更せずに通過させられる履歴の集合）の大

きさを比較したり、シミュレーションによって比較される場合が多い。これまでの比較結果では、時刻印方式よりは2相ロック方式の方が有効であると報告されている。

本稿では十分なディスクキャッシュがある場合において、2相ロックアルゴリズム^[2,4]とSGTアルゴリズムの1種であるBBアルゴリズム^[9]をシミュレーションによって比較する。

2. 並行制御アルゴリズム

本稿で用いるアルゴリズムを説明する。ここではトランザクションに関して最小の情報だけを利用するアルゴリズムを考える。2相ロックアルゴリズムではトランザクションは共有ロック、排他ロック、アンロックの列と考えられている。その詳細は省略するが、情報量の仮定から、2相ロックアルゴリズムの中でもaggressive 2相ロックを用いる。このアルゴリズムはREADやBLIND WRITEを行う直前にロックを行い、コミットと同時にアンロックを行うアルゴリズムである。また2相ロックにおいて、デッドロックが生じた場合のアポート処理の方法によって効率が変化する。2相ロックにおけるデッドロックはwait-forグラフに閉路ができることで検出できる。できた閉路にを構成する頂点に対応したトランザクションの内より1つ選んで後退復帰させるのが通常のやり方である。本稿では、後退復帰させるトランザクションの選び方として、1) そのトランザクションの終了を待っているトランザクション数が最大のもの、2) これまで発行したI/O命令が最小のもの、の2種類を考えた。後退復帰するトランザクションの選び方として1)、2)の方法を採用した2相ロックアルゴリズムをそれぞれ2PLv1, 2PLv2と呼ぶ。

BBアルゴリズムではトランザクションは、READonly, READ followed by WRITE およびWRITEの列と考える。READを2種類に区別しているが、いずれの種類のREADであるかはアルゴリズムにわかっているものと仮定している。2相ロック方式においても、共有ロックと排他ロックの区別がアルゴリズムにわかっているものと仮定しており、両

者が利用する情報は同じである。

SGTアルゴリズムは競合する命令のすべての組み合わせに対して、その出現順をトランザクション間の順序として利用するものである。トランザクション間の順序を、トランザクションを頂点とするグラフで表現したとき、そのグラフが閉路を持たない限り、その履歴は直列化可能であることが知られている。SGTアルゴリズムでは、1つの命令を実行許可した場合に、そのグラフを更新し、グラフに閉路ができなければその命令を許可し、閉路ができればどれか1つのトランザクションをアポートする。

BBアルゴリズムは、各実体ごとに、トランザクション間の順序を示すグラフを管理し、こうしたグラフによって規定される順序に閉路ができるかどうかを調べる。実体ごとに順序関係を表現するためより細かな制御ができ、blind WRITE間の出現順序を利用するためオンラインでグラフの更新ができる。また、実体毎のグラフは実際には特殊な形のグラフになるため、比較的単純な処理で実現できる。以下にBBアルゴリズムを示す。

データベースに初期値を書き込む仮想的トランザクション T_0 の存在を仮定する。

[定義1] 与えられた履歴 h において、特に実体 x にアクセスする命令だけを考える。read-fromグラフ $G(h, x) = (N, E)$ を次のように定義する。 N は頂点集合で、1つの頂点は x にアクセスするトランザクションに対応している。 E は有向枝で、 h が $w_i[x] \cdot R_j[x]$ を部分列として含む(ただし、その間に $w[x]$ はない)とき、 T_j read-from T_i といい、

(T_i, T_j)が死の要素になる。以後、read-fromグラフの頂点とトランザクションを区別しない。

一般に $G(h, x)$ はいくつかの連結成分から構成され、各連結成分の根は x にblind WRITEを行ったトランザクションになる。このとき次の補題が成立する。

[補題1] h が直列化可能ならば、任意の x に関するread-fromグラフ $G(h, x)$ の連結成分 C_i は木であり、 C_i の同一レベルにある頂点のうち x にWRITEする頂点は高々1つしかない。また、頂点 T_j の子を $T_{j1}, T_{j2}, \dots, T_{jk}$ とし、その内、 T_{j1} だけが x にWRI

TEしているとする。このときhと等価な直列履歴において、 T_{j2}, \dots, T_{jk} はどれも T_{j1} に先行する。

〔補題2〕履歴hが直列化可能とする。任意のG(h, x)とその2つの連結成分 C_1, C_2 を考える。hと等価な直列履歴h'において、 $T_m \rightarrow T_n$ ($T_m \in C_1, T_n \in C_2$)であれば、任意の $T_u \in C_1$ と任意の $T_v \in C_2$ に対して $T_u \rightarrow T_v$ が成り立つ。

実体xに関する直列化可能性判定を行うためのBB判定グラフを定義する。

〔定義2〕履歴hと実体xに対して、read-fromグラフG(h, x)の連結成分間の順序を、各成分の根に対応したトランザクションが発行したblind WRITEのhにおける出現順とする。このG(h, x)において、補題1、2によってえられる枝集合を付加して得られるグラフをxのBB判定グラフといいBB(h, x)と書き、hに出現するすべての実体xに関してxのBB判定グラフの和集合をhのBB判定グラフという。(グラフG(N, E)とG'(N', E')の和集合は(N∪N', E∪E')というグラフである)。

〔例1〕 $h=R_1[x]R_2[y]R_3[x]W_2[y]W_1[x]R_5[y]R_2[x]R_4[y]R_5[x]R_6[y]W_4[x]R_6[x]W_6[y]$ とする。read-fromグラフG(h, x)、xのBB判定グラフBB(h, x)、およびhのBB判定グラフを図1～3に示す。

このとき次の補題が成り立つ。

〔補題3〕履歴hのBB判定グラフが閉路を持た

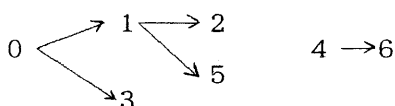


図1 例1の履歴に対するG(h, x)。

Fig.1 The read-from graph on x for the history of Example1.

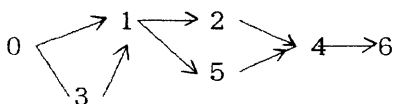


図2 例1の履歴に対するBB(h, x)。

Fig.2 The BB decision graph on x for the history of Example1.

なければ、hは直列化可能である。BB判定グラフによる直列化可能性の判定は多項式時間で実行でき、BB判定グラフによって直列化可能と判定される履歴はクラスCSR^[2,4]を含む。

BB判定グラフによる並行制御アルゴリズムは次のようである。

- 1) 実体xへのアクセス要求があればBB(h, x)に必要があれば枝を追加し、それに応じてBB判定グラフにも枝を追加する。
- 2) BB判定グラフに閉路ができなければその要求を許可し、閉路ができれば、条件によって、その要求を待たせるか、あるいは、その要求を許可し、かつ閉路に含まれるトランザクションの1つを後退復帰させる。

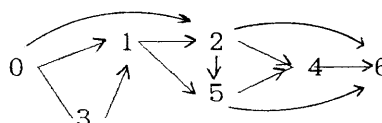


図3 例1の履歴に対するBB判定グラフ。

Fig.3 The BB decision graph for the history of Example1.

3. 不動点集合の比較

任意に発生させた実行履歴のうち、2相ロックアルゴリズム、BBアルゴリズムによってそれぞれ直列化可能と判定された履歴の割合を測定した。直列化可能と判定された履歴はこれらのアルゴリズムによってそのまま実行されるが、そうでないと判定された履歴は(実際には直列化可能であっても)、実行順序を変更して実行される。並行制御アルゴリズムによって、実行順序を変更せずに実行される履歴集合をそのアルゴリズムの不動点集合と呼ぶ。

シミュレーションパラメータは次の通りである。

トランザクションの平均長: AV_LENGTH

トランザクション長の標準偏差: AV_LENGTH*0.2

トランザクション数: TR_NO

読み専用トランザクションの比率: RATE_OF_R
読み書きを行うトランザクションにおける

READアクションの比率: RATE_OF_R_IN_RW
 blind WRITEの比率: RATE_OF_B_IN_RW
 データ数は1000、その内ホットスポットデータを200、ホットスポットデータのアクセス割合を80%とし、1000の履歴を発生させ、不動点集合に入った履歴の数を数えた。

AV_LENGTHを変化させた場合の実験結果を表1に示す。この実験ではTR_NO=10, RATE_OF_R=50%, RATE_OF_R_IN_RW=30%, RATE_OF_B_IN_RW=20%とした。

次にTR_NOを変化させた場合、各クラスのサイズがどの程度になったかの実験結果を表2に示す。表2の実験ではAV_LENGTH=10, RATE_OF_R=50%, RATE_OF_R_IN_RW=30%, RATE_OF_B_IN_RW=20%とした。

不動点集合の大きさで比較する限りでは、BBアルゴリズムの方が2PLより格段に優れている。しかし競合が激しくなるにつれ、不動点集合は急速に小さくなるのがわかる。

表1 トランザクションの平均長を変化させた場合の不動点集合の大きさ
 (トランザクション数=10)

Table1 Size of fixed points sets.
 (The number of transactions is 10).

AV_LENGTH	5	10	15	20	25
2PL	819	373	53	6	2
BB	957	729	330	79	13

表2 トランザクション数を変化させた場合の不動点集合の大きさ
 (トランザクション平均長=10)

Table2 Size of fixed points sets.
 (The average length of transactions is 10).

TR_NO	10	20	30	40
2PL	373	18	0	0
BB	729	223	38	1

4. シミュレーションモデル

Agrawal^[1]らはシミュレーションモデルによって全く逆の結果が生じ得ることが指摘し、並行制御アルゴリズムのためのシミュレーションモデルを提案し、その後このモデルが事実上の標準モデルと考えられてきている。ここでも原則的にこのモデルに従ってシミュレーションを行った。

トランザクション: WRITEを伴わないREAD、WRITEを伴うREAD、WRITEの列であり、それぞれの命令はある平均時間間隔をおいて発行されるものとした。新しいトランザクションは、同じ端末で実行された以前のトランザクションが終了した後、ある時間間隔をおいて発生するものとしている。

このシミュレーションで用いている論理モデルを図4に示す。物理モデルとしてCPUは1台と仮定した。従って、トランザクションの処理は逐次的であり、1つの処理を終えた後次の処理が行われる。I/Oサーバは無限と仮定した。時刻 t_1 で $A_1[x]$ 、時刻 t_2 で $A_2[y]$ が実行されるとき、トランザクション T_1 の次の命令は $t_1 + I/O$ 時間以後に実行可能となり、 T_2 の次の命令は $t_2 + I/O$ 時間以後に実行可能となる。つまり2つの入出力は重畳して実行される。ただし x と y が同じ場合は、 $t_1 \leq t_2 \leq t_1 + I/O$ であれば、 T_2 の実行は $t_1 + I/O$ までblockされる。ここで I/O とはデータの入出力にかかる時間である。

またマルチプログラミングのレベルも無限と仮定

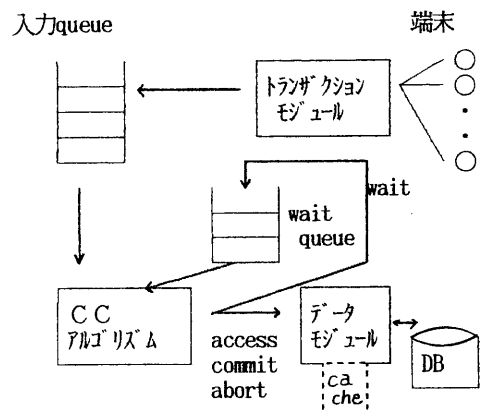


図4 論理モデル

Fig.4 The logical simulation model

した。従って、すべてのトランザクションは同じ優先度で並行して実行される。

5. 実験結果

前節のモデルに基づいて、シミュレーションを行った。端末数とトランザクションの平均長以外のパラメータの値は次の通りである。

○端末当りのトランザクションの平均発生時間間隔 10秒

○1つのトランザクション内での、平均I/O発生時間間隔 3ミリ秒

トランザクションに関する上記以外のパラメータは3.と同じである。トランザクションはその値の書き込みをキャッシュに対して行い、コミットした時点でデータベースに一括して書き込むことにした。また、BBアルゴリズムでは、どのトランザクションもそれに先行するトランザクションがコミットするまでコミットできないものとした。BBアルゴリズムにおけるアボートは2PLv2と同じく、ループに含まれるトランザクションのなかで、もっとも処理の進んでいないトランザクションをアボートさせることとした。

○キャッシュサイズ 300

○ディスクI/O時間 100ミリ秒

○キャッシュI/O時間 1ミリ秒

○トランザクションモジュール処理時間 1ミリ秒

○2PLアルゴリズムの処理時間 2ミリ秒

○BBアルゴリズムの処理時間 3ミリ秒

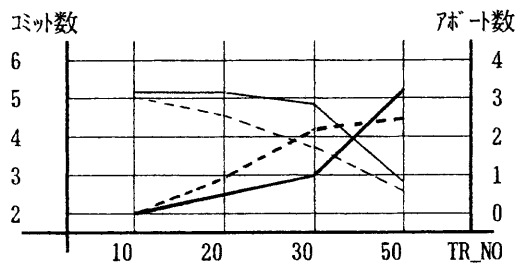


図5 1端末あたりのコミットしたトランザクションの数
1つのトランザクション当たりのアボート回数
(平均トランザクション長=20)

Fig.5 The number of committed transactions and aborted transactions (average length=20)

この仮定のもとで、あらかじめ各端末ごとに発生させておいたトランザクション集合に対して、並行制御アルゴリズムを適用し、30秒間のシミュレーションを行って、その間にコミットされたトランザクション数とアボートの回数を測定した。BBアルゴリズムではアボートの連鎖が起りえるが、その回数も算入している。以下で破線は2PLを実線はBB、細線はコミット数、太線はアボート数を示す。

図5は平均トランザクション長が20で端末数を変化させた場合、1端末あたりのコミットしたトランザクションの数と1つのトランザクション当たりのアボート回数を示している。図6は端末数を30とし、平均トランザクション長を変化させた場合、1端末あたりのコミットしたトランザクションの数と1つのトランザクション当たりのアボート回数を示している。図7は端末数を変化させた場合の、システム内のコミットトランザクション数を示している。

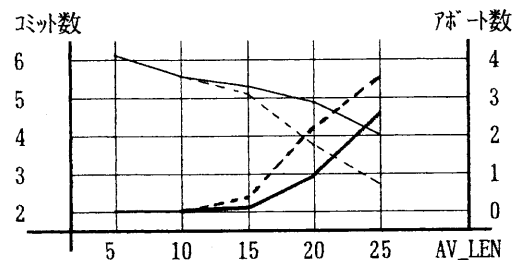


図6 1端末あたりのコミットしたトランザクションの数
1つのトランザクション当たりのアボート回数
(端末数=30)

Fig.6 The number of committed transactions and aborted transactions (terminals=30)

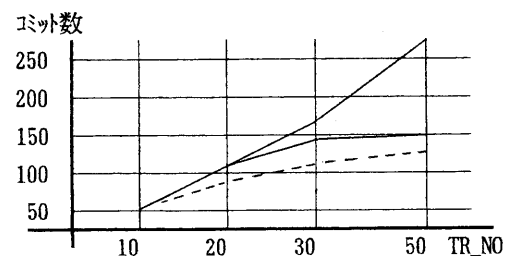


図7 システム内のコミットトランザクション数

Fig.7 Total number of committed transactions

図5、6から明らかなように、競合の少ない場合には2PLアルゴリズムもBBアルゴリズムもその効率に大差がない。むしろアルゴリズムの処理時間の仮定から2PLアルゴリズムの方がほんの少し良いといえる。しかし、競合が激しくなるにつれ、2PLアルゴリズムの効率が急速に低下するのに対して、BBアルゴリズムは比較的良い効率を示し、アルゴリズムの差が顕著に現れてくる。2PLv1と2PLv2とでは、個別のケースでは1長1短ではあるが、概して2PLv2の方が良い効率を示している。

BBと2PLを比較すれば、2PLの方がロックを用いてトランザクションの実行をブロックするケースが多いので、2PLの方がアボートが少ないと考えていたが、結果はBBアルゴリズムの方が少なかった。

平均トランザクション長を固定して、端末数を変化させてテストしたトランザクションの総数の変化を見てみると、初めは端末数の増加に比例して増加してゆくが、次第に増加率が鈍り、ついには飽和状態になるように見受けられる(図7)。

6. まとめ

2. で示したBBアルゴリズムはアクセスする各実体毎にBB判定グラフを作成するため、2相ロックアルゴリズムに比べて実体毎にNバイトほど余分の記憶量が必要となる(Nはトランザクション数)。従って、実体として余り小さな単位を仮定するのは現実的ではない。

並行制御ではどの単位で並行制御を行うかという議論があるが(粒度という)、粒度が小さい場合にBBアルゴリズムを適用するには多少工夫をする必要がある。1つはホットスポットデータについてはBBを、他の実体については2PLアルゴリズムを適用するというハイブリッドアルゴリズムである。この方法の効率はシミュレーションによれば両者の中間的な性能を示し、取り立てて見るべき結果にはならなかったが、実現面からみれば実際的かもしれない。もう1つの方法はアクセス頻度の小さい実体集合を分割し、その実体集合に対して1つのBB判定グラフを用いることである。

この場合、1つの集合に属する実体へのアクセスは実際には競合していなくても競合しているように扱われるため、1つのトランザクションでともにアクセスされる確率の高い実体同士を同じ集合に含めるように分割することが必要である。

分散データベースに対してBBアルゴリズムを適用することも可能である。特にBB判定グラフは1次元で表現できるため、全域的な競合情報を1箇所管理しても通信量は大きくならない。分散データベースでの効率測定は残されている。

[参考文献]

- [1] Agrawal, R., Carey, M. J. and Livny, M.: Concurrency Control Performance Modeling: Alternatives and Implications, ACM TODS, vol.12, No.4, 609-654, 1987.
- [2] Bernstein, P. A., Hadzilacos, V. and Goodman, N.: Concurrency Control and Recovery in Database Systems, Addison Wesley, 1987.
- [3] Fragg, A. A. and Ozsu, M. T.: Using Semantic Knowledge of Transactions to Increase Concurrency, ACM TODS, vol.14, No.4, 503-525, 1989.
- [4] Papadimitriou, "The Theory of Database Concurrency Control", Computer Sci. Press, 1986.
- [5] Shasha, D., Llibat, F., Simon, E. and Valduriez, P.: Transaction Chopping: Algorithms and Performance Studies, ACM TODS, vol.20, No.3, 325-363, 1995
- [6] Wu, K-L., Yu, P. S. and Pu, C.: Divergence Control for Epsilon-Serializability, IEEE 9th Inter. Conf. on Data Engineering, 506-515, 1992.
- [7] 加藤, 「分散データベース・システムにおける同時実行制御方式の性能評価」、情報処理学会論文誌、33巻11号、1361-1371, 1992年
- [8] 徐、古川、史, 「一貫性情報を用いたデータベースの並行処理制御」、情報処理学会論文誌、35巻12号、2752-2761, 1994年12月
- [9] 中津, 「直列化可能性の新しい判定法」信学論 vol. J76-D-I, no.6, 1993.