

データベース・マシンの同時処理制御

A Concurrency Control in a Database Machine

水 摩 正 行 牧 野 武 則 梅 村 護 日 吉 茂 樹
 Masayuki MIZUMA Takemoni MAKINO Mamoru UMEMURA Shigeeki HIYOSHI
 日 本 電 気 (株) 中 央 研 究 所
 Central Research Lab. Nippon Electric Co., Ltd.

1. はじめに

データの共同利用を主目的とするデータベース管理システム(以下DBMSと略す)では、そのシステムが採用している同時処理制御方式が、システム性能、システムの使い易さに大きな影響を与える。

DBMSの制御方式に対して、種々の方式が提案され、実際に実現されてきた。同時処理制御の特徴は、制御の対象(ロック対象)と制御モードに現われる。システム設計において、それらをどのように決定するかは、ロック対象への競合度、システムオーバヘッド、使い易さ、実現の容易性などを十分考慮したうえで行うべきである。

現存するDBMSの中で同時処理制御機能の高いシステムはシステムRである[1]。システムRでは、tuple、データアイテム値など、非常に細かいレベルでの排他制御を行っており、その制御モードもshared lock, exclusive lockなどを任意に選べる。しかしながら、実用システムとして、このようなきめ細かい制御機能を組み込むのは、それによってもたらされる効果と制御の複雑さ、オーバヘッド等を考慮すると必ずしも最適でない。

一方、DBMS1100[2]、IDSL3[3]などでは、ロックの最小単位をページでおさえている。この方式は、システムRに比較して同時動作性、制御の柔軟性に劣るが制御が単純になる。しかしながら、ページに対するロックモソフトウェアにより実現しているのでオーバヘッドが

向強である。

本資料では、後置プロセッサタイプのデータベースマシン(Generalized Database Subsystem; 以下GDS^[4]と略す)の同時処理制御の設計思想、実現方式について報告する。

GDSの同時処理制御はDBMS1100のアドレスに近いが、次のような特徴をもっている。

- (1) ユーザインタフェースレベルでは、二段階(エリア、レコードオカレンス)の排他制御を採用している。しかし、レコードレベルの排他制御はシステム内部でページに置換している。
- (2) ページの排他制御は、ロックオーバヘッドを減少させるため、ハードウェアによって実現されたアドレス変換機構を利用して行われる。
- (3) テッドロックフリー方式を採用している。テッドロック検出機構はハードウェア化される。

以下、データベースのロック対象の大きさとシステム性能との関係を解析する。そして、その解析結果を参考にして設計されたGDSの同時処理制御方式について紹介する。

2. ロック単位の選択基準

DBMSの設計において、同時処理制御の最小単位をどの程度にするか、が重要な問題である。ロックの対象単位が大きければなる程トランザクション間の競争が発生し、多重処理の効果が減小する。一方、ロックの最小単位を細かくすると、競争率は減少するが、時間的・空間的オーバーヘッドが増大し、インプリメンテーションも困難になってくる。

ここで、ロック単位と競争性との関係を考察する。

はじめに、記号の定義を行う。

λ ; 平均実行トランザクション数
 μ_r ; トランザクション当りの要求レコード数
 d_r ; データベース(DB)の総レコード数
 g ; DBブロック数(DBの分割数)
 s ; トランザクション当りの要求ブロック数
 このとき、特定ブロックに*i*個のトランザクションが競争する確率 C_i は、次の二項分布で表現される。

$$C_i = \binom{\lambda}{i} \left(\frac{s}{g}\right)^i \left(1 - \frac{s}{g}\right)^{\lambda-i} \quad (1)$$

平均専有ブロック数 B_c は、

$$B_c = \sum_{i=1}^{\lambda} g \cdot C_i = g(1 - C_0) = g\left(1 - \left(1 - \frac{s}{g}\right)^{\lambda}\right) \quad (2)$$

平均実行可能トランザクション数 $\tilde{\lambda}$ は、

$$\tilde{\lambda} = (B_c / s) \lambda = B_c / s \quad (3)$$

とれる。

図1(実線)に g , s , $\tilde{\lambda}$ の関係を示した。ここでは、 $\lambda = 10$, $\mu_r = 50$ と仮定している。 μ_r と s の関係は、一般にトランザクションのデータ参照特性(ローカリティ特性)に影響される。ここでは、次の3ケースによって示した。

$s = 5\lambda$; 完全分散

$s = g\left(1 - \left(\frac{g-1}{g}\right)^{\lambda_r}\right)$; ランダム

$s = 10$; ローカリティがある程度存在

また、Closed Queueing Network Model [5] を用いて、ロックによりシステム性能がどの程度低下するかを調べた。その結果も図1(破線)に示している。なお、システムのモデルおよび値の設け方は付録に示す。

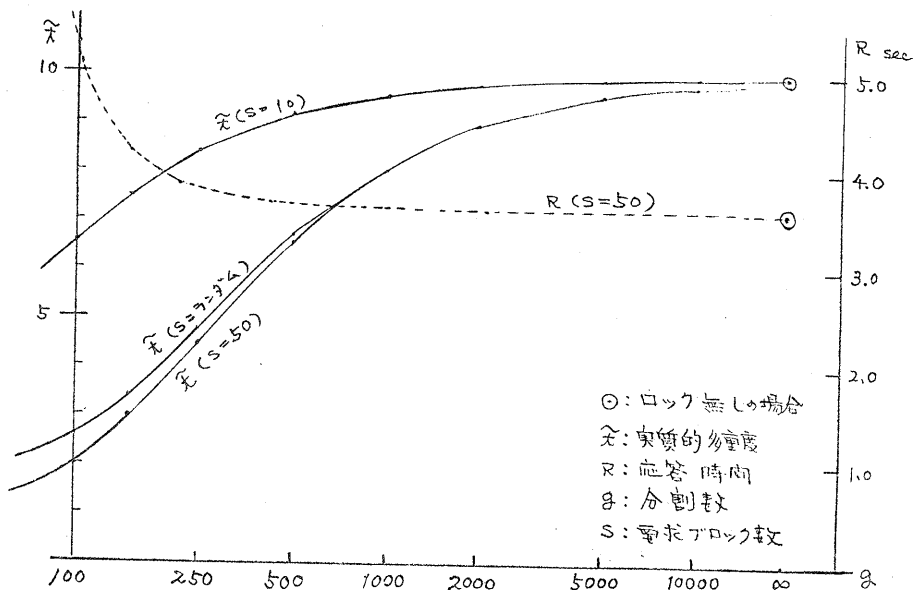


図1. ロック単位と競争性の変化

図から、DBの分割数が 500 以上のオメガにはれば、ロック競合による実質的多重度の低下はそれ程認められないことが判る。

一般に、DBのロックの対象候補として、DB全体、エリア、レコードクラス、レコードオカレニス、データアイテム値、ページなどが考えられるが、前記の解析結果より、大規模DBを取り扱うDBMSでは、ロックの最小単位はページで十分であると云えよう。

上記の背景をもとに、GDSでは、ロックの最小単位をページと定めシステム設計を行った。以下、GDSの同時処理制御について紹介する。

3. データベースの同時処理制御

GDSの同時処理制御は、ユーザが隔に要求を出すものと、システムが自動的に行うものがある。

前者は、ユーザのプログラムロックが他のユーザの更新によって乱されるのを防ぐために行われる。後者は、DB上での矛盾を防ぐために行われる。矛盾が生じる主な原因は、複数ユーザが同時に出す更新要求の干渉、および更新処理の途中でトランザクションの撤回である。

3.1 ユーザインタフェース レベルでの保護

ユーザインタフェース レベルでデータの一致性 (Consistency) を保障するため、GDSでは、エリアレベルの保護とレコードレベルでの保護機能を提供している。

(1) エリアの保護

エリアはCOBASYL集[6]のエリアの概念に相当する。一般にエリアに含まれるデータ量は非常に多い。このレベルでは、トランザクション向の競合を緩和させるために柔軟な排他制御モードが必要になる。

GDSでは、COBASYL仕様に従った排他制御モードを全てサポートしている。

エリアに対する排他制御の宣言は、OPEN/CLUSE命令によって行われる。ユーザは

この時、エリアの活用モード (RETRIEVE/UPDATE) と排他制御モード (EXCLUSIVE/PROTECTED/SHARED) の宣言を行う。表1に同時処理の規則を示す。

表1. エリアの排他制御規則

	S.R	S.U	P.R	P.U	EX
S.R	O	O	O	O	X
S.U	O	O	X	X	X
P.R	O	X	O	X	X
P.U	O	X	X	X	X
EX	X	X	X	X	X

O: 実行可 X: 不可

S (Shared): エリアの共有使用

P (Protected): 検索処理ユーザに対しては許されるが更新処理ユーザは共用不可。

EX (Exclusive): 完全排他使用

R (Retrieve): 検索処理

U (Update): 更新処理

なお、GDSでは、エリアレベルでの排他制御によってシステムがデッドロック状態に陥るのを防ぐため、トランザクションはエリアの使用を一括要求するようにしている。

(2) レコードオカレニスの保護

2章の解析結果から判るように、レコードオカレニス単位ではトランザクション向の競合は余り生じない。そこで、GDSでは、このレベルでの排他制御は共用モード (Shared) と専有モード (Lock) に限定している。

ロック要求は次のように行われる。GDSでは、種々のアクセスパス (PRIMARYパス, IMAGEパス, LINKパス) に従って、DB処理を行うことができるが、ユーザはその処理に先立ってアクセスパスの生成を要求する[4]。この時、必要があればロック要求を指定する。ロック要求が指定された場合、アクセスパスのカレントレコードは自動的にロックされる。カレントレコードが切り替わると旧カレントレコードはUNLOCKされ、新カレントレコードはLOCKされる。カレントレコードが切り替わった後に旧レコードをロックしたい場合は

LOCK 命令を使って陽に指定できる。

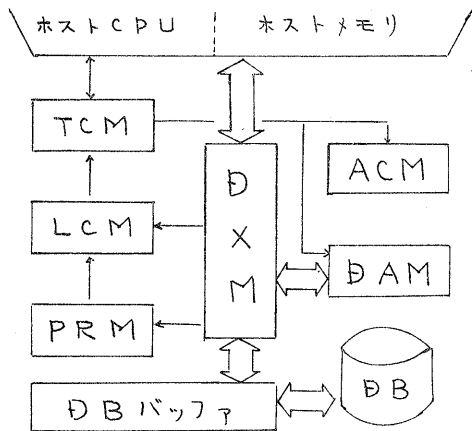
3.2 システムレベルでの保護

GDS では、1レコードアクセスを基本としているため、DBの更新は複数の命令を使って手続的に行われる。更新処理が完全に終了するまでの間に、他のトランザクションの更新が発生した場合DBの論理的な矛盾が発生する。また、更新の途中でトランザクションが撤回される可能性もある。このような状況からDBを守るため、DBレコードの更新されたページは更新トランザクションが終了するまで、GDSが自動的にロックをかける。

また、DBの更新に伴って、インテックスやディスクリフク等のシステム制御情報の修正が必要になってくる。これらの制御情報に対しては、GDSが必要に応じて該当ページを一時的にロックする。

4. システム構成

GDSの多重処理制御は、1つのトランザクションにつき、1つのGDIコマンド (Generalized



ACM: Area Control Module ⇒ データフロー
 DAM: Data Access Module → 制御フロー
 DXM: Data Transfer Module
 LCM: Lock Control Module
 TCM: Transaction Control Module
 PRM: Page Replacement Module

図2. システム構成

Database Interface)のみ実行する方式を採用している。

以下、多重処理制御機能に注目し、GDSのシステム構成を説明する。

GDSは処理の形態に応じて、いくつかの機能モジュール (FM) に分割されている。図2はその構成を示したものである。

(1) トランザクション コントロール モジュール (TCM)

TCMはホストプロセッサとの通信、およびトランザクション処理のスケジューリングを行う。はじめに、START TRANS 命令を受けると、トランザクションIDを作成する。トランザクションの動きは、トランザクションコントロールブロック (TCB) によって管理される。TCMは汎用計算機のものに比べ単純化されており、少ないオーバーヘッドで制御できる。

多重処理の制御は、3章で示したDBの同時処理制御イベント、I/O処理イベント、ホストからのCALL GDI 命令を契機に行われる。図3は、GDS内部でのトランザクションの制御状態の動きを示している。

(2) エリア コントロール モジュール

GDSに認識されたトランザクションは、最初にエリアの使用宣言を行う。ACMはトランザクションが要求した使用モードと排他制御モードを現在のエリアの使用状態と付き合わせて、トランザクションのキュー管理を行う。

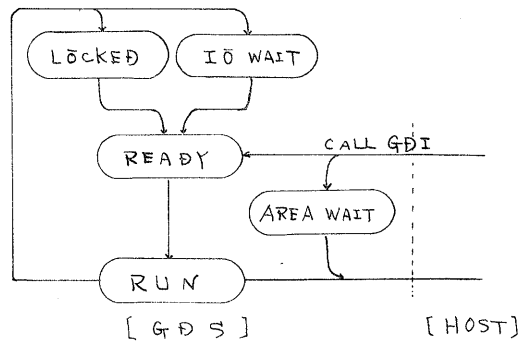


図3 制御状態の変動

(3) データ転送モジュール(DXM)

DXMは、GDSの中心的な存在で、DBバッファと、他のFMのワークメモリ、ホストプロセッサの主メモリとの間のデータ転送を行う。データ転送はレコードイメージで直接行われる。データ転送の指定パラメータは、論理レコードアドレス(エリア番号、ページ番号、ライン番号)、レコード内オフセットと長さ、で表現される。このデータ転送中にロックのチェックを行い、他のトランザクションによってロックされている場合にはLCMに競合指令を発生する。

(4) ロックコントロールモジュール(LCM)

LCMは、DXMが確認したデータ競合に従って、トランザクションのキュー管理を行う。また、キュー管理と同時にデッドロック状態の検出を行う。

(5) ページリプレイメントコントロールモジュール(PRM)

DBバッファ内のページ管理を行う。即ち、DXMがデータ参照時にミッシングフォルトを確認した場合、毎ページを二次記憶より転送するため、ページの置き換え制御を行う。置き換え制御では、ロックされたページは出来る限り避けるようにしている。ロックされたページが避けられない場合には、ロック情報をセーブするためLCMに制御を渡す。

(6) データアクセスモジュール(DAM)

GDSがサポートしている各種のアクセスパスに従って、レコードの検索、更新処理を実行する。

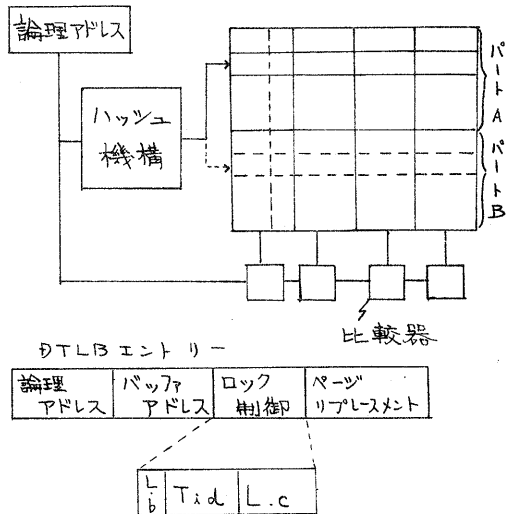
5. アドレス変換機構とロック

レコード、ページレベルの排他制御は従来ソフトウェアによって実現されていた。このため、ロック制御のためのオーバーヘッドが多であった。また、そのオーバーヘッドを無視することができないため、参照データに關係する制御情

報を更新する場合(たとえば、カレンシーインジケータの修正時ほど)のみ、アクセス権のチェックをする、というように方式で予め設計時に一定の規則で定めその規則に従ってチェックを行っていた。しかしながら、このような方式では、正規の手続きを経ないデータ参照に対してアクセス違反を摘発することができない。

GDSでは、このレベルの排他制御を行うために、ハードウェアによって実現されるアドレス変換機構を有効に利用している。DBデータを参照する場合には、常にこのアドレス変換機構を介するので、ロック機能を組み込めば、常時アクセス権のチェックを行うことができ、より強力はロック制御機能を実現できる。また、ロックオーバーヘッドも無視できる。なお、レコードに対する排他制御は、競合が少ないため、すべてページで代行している。

アドレス変換機構はDXMの1コンポーネントであり、論理アドレスをI/Oバッファ上の物理アドレスに変換する。図4にアドレス変換機構の構成を示す。



T.id; トランザクションID
L.b; ロックビット
L.c; ロックカウント

図4. アドレス変換機構(DTLB)

高性能はアドレス変換機能を実現する場合には、連想メモリ素子の使用が望まれるが、低価格で実現するためハッシュ機構と比較器を組み合わせた疑似連想記憶方式で実現される(DTLBと呼ぶ)。

DTLBの各エントリーには、アドレス変換情報の他に、ロック制御情報やバッファ管理情報が含まれている。

他のFMからデータベースデータの参照要求が発生した場合、このアドレス変換機構を通して、目的データベースの物理アドレスを決定すると同時に、データに対するアクセス権のチェックを行う。

ロック制御情報はロックビット、トランザクションID、ロックカウントから成っている。

ロックビットは対応するページのロック状態を示す。もしロックビットがONであれば、他のトランザクションが対象ページをロックしていることを示す。またトランザクション名はトランザクションIDフィールドに示されている。この場合データ参照はトランザクションIDフィールドに一致するトランザクションのみ許可される。ロックビットがOFFの場合には、任意のトランザクションからのデータ参照が許可される。

ロックカウントはロック中のトランザクションが対象ページに要求を出したロックの回数とアンロックの回数の差を示す。対象ページをロックしている間はその値は正である。ロックカウントが0になった場合ロックビットがOFFになりページが解放される。このロックカウントは次の2つの目的で設けられた。

- (1) 従来では、レコードに対するロックをページで代行している。ページには一般に複数のレコードが存在し、ユーザはどのいくつかのレコードに対しロック要求を出す可能性がある。この時、ユーザが対象ページに含まれているロック中のすべてのレコードをアンロックするまでそのページを解放してはならない。
- (2) 一般にシステムを構成する場合 E.W. Dijkstra 等が提唱しているように階層設計が行われる。これは他の階層に対する余分な知識を不要にするためである。このように

設計においては、各階層で同一オブジェクトに対し異なった目的でロックをかける場合が発生する。このように、目的に応じて自由にロック、アンロックをかけ得るようになれば設計が楽になる。

上記の状況において、システムがページに対するロックの解放のタイミングを待つためにロックカウン트가設けられた。

6. デッドロック検出機構

デッドロックを防止しつつの方法として、資源を一度要求させる方式がある。エリア単位の排他制御を行うため、従来では、この方式を採用している。しかしながら、レコードオカレンスレベルの排他制御でこの方式を採用した場合には、ユーザの負担が大きくはり内題である。

このため、従来では、レコードレベルでの排他制御はデッドロックフリー方式を採用している。

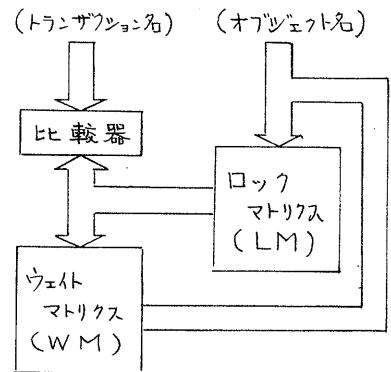


図5 デッドロック検出機構

システムのデッドロック状態を検出するためハードウェアにより実現されるデッドロック検出機構を備える。図5はその構成を示したものである。

デッドロック検出機構は、ロックマトリクス(LM)とウェイトマトリクス(WM)から成っている。LMの各要素は、ロック中のトランザクションとオブジェクトとの関係を示している。WMの各要素はオブジェクトを待っ

ているトランザクションの状態を示している。

5章で記述したアドレス変換機構により、あるトランザクションが目的ページのロック待ちイベントを確認すると本装置に通知される。

デッドロック検出機構は、ロック待ちトランザクション名とロックオブジェクト名を入力パラメータとして受け取る。最初に、LMより対応するオブジェクトをロックしているトランザクション名を検出する。そして、そのトランザクション名をWMに入力することによって待ち状態にあるオブジェクト名が抽出される。そのオブジェクト名を再度LMに入力させ、その出力が要求トランザクション名(入力パラメータのトランザクション名)と一致した場合にデッドロック状態が発生したと云える。一致しなければ、同じ動作を続ける。そして、最後まで一致するトランザクションが検出できなければ、デッドロック状態が発生していないこととなる。

7. おわりに

DBMSの設計において、DBの保護の最小単位をどの程度に設定すればよいか、は興味深い問題であろう。

本資料では、この問題を解析的に求め、ページを最小単位とすれば十分であることを示した。

また、データベースマシンでの同時処理制御方式について論じた。

本データベースマシンはメモリ共用型の後置プロセッサである。データベースの保護はページを基本単位としている。更に、核機能をハードウェア/ファームウェアにより実現することによって、オーバヘッドの減少をはかっている。

本方式の有効性を確かめるために、現在実験システムを構築中である。

最後に、本研究の機会を与えて下さった日本電気中央研究所コンピュータシステム研究部 藤野部長、苅津課長、さらばに、日頃直接御指導頂いている箱崎主任、久保主任に感謝致します。

《参考文献》

- 1) M.M. Astrahan et al : System R: A relational approach to database management, ACM TODS, Vol.1, No.2 (1976).
- 2) UNIVAC 1100 Series, Data Management System (DMS1100) Data Manipulation Programmer Reference, UP-7908.
- 3) ACOS-6 データ管理 統合データベース (IDS) 説明書.
- 4) K. Hakegaki et al : A Conceptual Design of A Generalized Database Subsystem, 3rd International Conf. on VLDB, Oct. '77.
- 5) W.J. Gordon and G.F. Newell : Closed Queuing Systems with Exponential Servers, Operations Research, Vol.15, No.1-3, 1967.
- 6) CD DASYL Data Description Language Journal of Development, Jun. '73, NBS Handbook (1974).

[付録] システム評価モデル

Closed Queuing Network モデルでは、次のことと前提にしている。

- 各サーバのサービス時間: 指数分布に従う。
- サーバ間の遷移: 遷移確率に従う random walk
- システム内のトランザクションの数は一定である

上記の前提をもとに、システム構成、各サーバの平均サービス時間、サーバ間の遷移確率、システム内のトランザクションの数を入力パラメータとして、各サーバの平均利用率、平均待ち行列の長さ、システムルータンアタイム(応答時間)を計算できる。

Closed Queuing Network モデルを使って、DBの分割数の変化によってシステム性能がどのように影響されるかを近似的に示めるため、図A-1に示すようなシステム構成を仮定した。

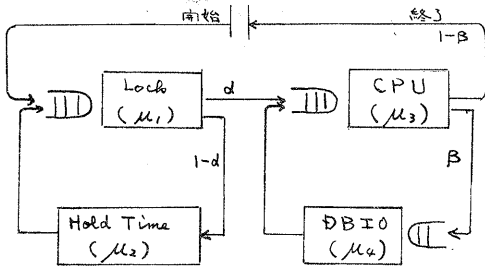


図 A-1. システム構成モデル

図において、

μ_1 : ロックオーバーヘッド時間。本来 CPU 処理であるが別サーバとした。

μ_2 : ロック競合によって待ちさせられる時間

μ_3 : CPU による 1 回当りのサービス時間

μ_4 : DB 入出力時間。DB I/O は複数サーバへの指定可能

d, β : サーバ間の遷移確率

を表わす。

入力パラメータは、次のようにして決めた。
 μ_3, β はトランザクションの DB 処理特性と反映させるため、ある保険会社のオンラインデータベースシステムのデータを参考に設定した。本システムでは、トランザクション当りの CPU 処理量の平均は 320 msec であり、平均 I/O 回数は 25.5 回であった (I/O 時間: 40 msec)。この例では、CPU 処理 / 2 msec に 1 回の割合で I/O が発生することになる。

これから

$$\mu_3 = 12 \text{ msec}$$

$$\beta = 25.5 / (25.5 + 1) = 0.96$$

が得られる。

システム内の全トランザクションの数は 10 に固定した。

d は 2 章の (3) 式とシステム内全トランザクションの関係から求めた。 ($d = \bar{x} / 10$)

μ_1 はロックオーバーヘッドを表わすので、固定的に 2 msec とした。

Hold Time (μ_2) は、トランザクションの競合による待ち時間を近似的に表わすため、マルチサーバ (待ちが全いのように $20 (> 10)$) と仮定した。

μ_4 の値を粗雑くトランザクションの処理の終了間隔と対応させるため、CPU と DB I/O のみからなる別の Queuing Network Model を作り多重度 \bar{x} のときの応答時間を計算した。そして、その応答時間を \bar{x} で割った値を μ_4 として与えた。

表 A-1. 入力パラメータの設定値

入力パラメータ	値
多重度	10
μ_2	(多重度 \bar{x} のときの CPU, DB I/O 系のみでの応答時間) / \bar{x}
μ_1	2 msec (サーバ 20)
μ_3	12 msec
μ_4	40 msec (サーバ 5)
d	$\bar{x} / 10$
β	0.96
S	50

なお、本モデルでは、各トランザクションは処理に関係するすべてのリソース (7 ロック) をトランザクション処理の開始時に確保し、終了時に解放することと仮定している。

結果は、図 1 の破線で示している。