

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3681415号
(P3681415)

(45) 発行日 平成17年8月10日(2005.8.10)

(24) 登録日 平成17年5月27日(2005.5.27)

(51) Int. Cl.⁷

F I

G06F 9/46

G06F 9/46 340G

G06F 11/30

G06F 11/30 305G

請求項の数 8 (全 26 頁)

| | | | |
|--------------|------------------------------|-----------|--|
| (21) 出願番号 | 特願平6-61771 | (73) 特許権者 | 000005223 富士通株式会社 |
| (22) 出願日 | 平成6年3月30日(1994.3.30) | | 神奈川県川崎市中原区上小田中4丁目1番1号 |
| (65) 公開番号 | 特開平6-337798 | (74) 代理人 | 100090516 弁理士 松倉 秀実 |
| (43) 公開日 | 平成6年12月6日(1994.12.6) | (74) 代理人 | 100113608 弁理士 平川 明 |
| 審査請求日 | 平成12年6月23日(2000.6.23) | (74) 代理人 | 100105407 弁理士 高田 大輔 |
| 審査番号 | 不服2004-13022(P2004-13022/J1) | (74) 代理人 | 100089244 弁理士 遠山 勉 |
| 審査請求日 | 平成16年6月24日(2004.6.24) | (72) 発明者 | 藤田 和彦 神奈川県川崎市中原区上小田中1015番地 富士通株式会社内 |
| (31) 優先権主張番号 | 特願平5-96928 | | |
| (32) 優先日 | 平成5年3月30日(1993.3.30) | | |
| (33) 優先権主張国 | 日本国(JP) | | |

最終頁に続く

(54) 【発明の名称】 デッドロック検出装置

(57) 【特許請求の範囲】

【請求項1】

複数のタスクが共通の資源を利用するマルチタスクシステムにおいて前記複数のタスクが互いに占有している資源を待ち合せて停止してしまうデッドロックを検出するためのデッドロック検出装置であって、複数のタスクを並列実行するために、前記タスクの実行を管理するタスク管理部と、各タスクがどの資源をロックしているかを管理するロック管理部と、一のタスクが他のタスクがロックしている資源を獲得要求した場合には、前記一のタスクが前記他のタスクを待っているとしてこの各タスクの「待ち関係」を登録する待ち管理テーブルと、前記ロック管理部と非同期で動作するとともに、前記待ち管理テーブルに登録された「待ち関係」からデッドロックを検出するデッドロック検出部とを備え、

10

前記マルチタスクシステムは、複数のシステムを有する分散システム上に実現され、各システムがそれぞれ前記タスク管理部、前記ロック管理部、前記待ち管理テーブル、前記デッドロック検出部を備え、

一方のシステムのタスクが他方のシステムのタスクに対して「待ち」の状態にあるときには、その「待ち関係」を前記一方のシステムの待ち管理テーブルに登録し、デッドロックを検出しないときは前記一方のシステムから前記他方のシステムに通信して前記他方のシステムの待ち管理テーブルに登録し、前記他方のシステムにおいてその待ち管理テーブルを見てデッドロックの有無を判定することを特徴とするデッドロック検出装置。

【請求項2】

前記ロック管理部には、各タスクとそれによりロックされている資源との関係が登録され

20

るロック管理テーブルを有していることを特徴とする請求項1記載のデッドロック検出装置。

【請求項3】

タスクにおいて「待ち関係」が発生したとき、前記デッドロック検出部は、前記待ち管理テーブルに「待ち関係」を登録するとともに、前記待ち管理テーブルを検索してデッドロックの有無を検出することを特徴とする請求項1記載のデッドロック検出装置。

【請求項4】

前記デッドロック検出部は、要求キュー受付部を有し、この要求キュー受付部で、前記デッドロック検出命令である「待ち関係の登録要求」を受け付けることを特徴とする請求項1記載のデッドロック検出装置。

10

【請求項5】

前記待ち関係にある2以上のタスクのそれぞれが別のシステム内にあるときには、各システムに設けた前記待ち管理テーブルに各タスクにおける「待ち関係」を登録するとともに、自己のシステムのタスクが他のシステムのタスクに対して「待ち」の状態にある場合には、各システムのデッドロック検出部は、その待ち先のシステムの待ち管理テーブルに通信でアクセスして自己のシステムの前記「待ち」の状態を示す「待ち関係」を送信するとともに、前記待ち先のシステムのデッドロック検出部に対して、自己のシステムの前記「待ち関係」と待ち先のシステムの待ち管理テーブルの登録内容とを突き合わせてデッドロックの有無を判定させることを特徴とする請求項1記載のデッドロック検出装置。

【請求項6】

20

前記マルチタスクシステムは、複数のシステムを有する分散システム上に実現され、あるタスクについて「待ち関係」が一定時間継続している場合に、そのタスクに対して再度資源獲得要求を出す待ち時間監視部を備えたことを特徴とする請求項1記載のデッドロック検出装置。

【請求項7】

複数のタスクが共通の資源を利用するマルチタスクシステムにおいて前記複数のタスクが互いに占有している資源を待ち合ってしまうデッドロックを検出するためのデッドロック検出装置であって、複数のタスクを並列実行するために、前記タスクの実行を管理するタスク管理部と、タスクとそれによりロックされている資源との関係を登録する第1のテーブルと、各タスクがどの資源をロックしているかを管理し、他のタスクによりロックされているものとして前記第1のテーブルに登録されている資源を一のタスクが獲得要求した場合には、前記一のタスクが前記他のタスクを待っている「待ち関係」の発生を検出するロック管理部と、前記待ち関係を登録する第2のテーブルと、前記ロック管理部とは非同期に動作するとともに、この第2のテーブルに登録された「待ち関係」から、デッドロックを検出するデッドロック検出部とを備え、

30

前記マルチタスクシステムは、複数のシステムを有する分散システム上に実現され、各システムがそれぞれ前記タスク管理部、前記ロック管理部、前記第2のテーブル、前記デッドロック検出部を備え、

一方のシステムのタスクが他方のシステムのタスクに対して「待ち」の状態にあるときには、その「待ち関係」を前記一方のシステムの第2のテーブルに登録し、デッドロックを検出しないときは前記一方のシステムから前記他方のシステムに通信して前記他方のシステムの第2のテーブルに登録し、前記他方のシステムにおいてその第2のテーブルを見てデッドロックの有無を判定することを特徴とするデッドロック検出装置。

40

【請求項8】

複数のタスクが共通の資源を利用するマルチタスクシステムにおいて前記複数のタスクが互いに占有している資源を待ち合ってしまうデッドロックを検出するためのデッドロック検出方法であって、

前記マルチタスクシステムは、複数のシステムを有する分散システム上に実現され、各システムがそれぞれタスク管理部、ロック管理部、待ち管理テーブル、デッドロック検出部を備え、

50

前記ロック管理部が、各タスクがどの資源をロックしているかを認識し、一のタスクが獲得要求している資源が既に他のタスクによってロックされているか否かを検知し、資源が既に他のタスクによってロックされていることを検知した場合には、前記一のタスクが前記他のタスクを待っていると認識し、この「待ち関係」を前記待ち管理テーブルに登録し、

前記デッドロック検出部が、登録された「待ち関係」が各タスクが互いに待ち合っていることを示す場から、デッドロックとして検出し、

一方のシステムのタスクが他方のシステムのタスクに対して「待ち」の状態にあるときには、その「待ち関係」を前記一方のシステムの待ち管理テーブルに登録し、デッドロックを検出しないときは前記一方のシステムから前記他方のシステムに通信して前記他方のシステムの待ち管理テーブルに登録し、前記他方のシステムにおいてその待ち管理テーブルを見てデッドロックの有無を判定することを特徴とするデッドロック検出方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】

本発明はマルチタスクシステムにおけるデッドロックの検出装置に関する。

【0002】

【従来の技術】

近年、コンピュータを用いた情報処理システムにおいて、複数のタスクあるいはトランザクションを同時に実行するマルチタスクシステムが発達してきた。タスクとは、CPU内部における仕事の単位である。トランザクションとは、ひとつの完結したデータ操作を行うオペレーションの集まりである。マルチタスクとは、複数のプログラム（タスク、トランザクション）が、単一のコンピュータシステム、又は相互に情報交換可能に接続された複数のコンピュータシステム上で、同時に並行して実行される状態である。

【0003】

このマルチタスクシステムでは、2以上のタスクが資源を共用する場合がある。その場合において、2以上のタスクの夫々が、そのタスクの実行に必要であり且つ他方のタスクの実行にも必要な複数の資源を、一部づつを占有（ロック）し合うケースが生じ得る。そのケースでは、お互いに他方のタスクが占有（ロック）している資源を待ち合うので、双方のタスクが停止し、それ以上プロセスを実行できない状態となってしまう。このような状態は、デッドロックの状態と呼ばれている。

【0004】

図3に、デッドロックの状態の例を示す。図3の例は、2つのコンピュータシステム*i*、*j*から構成される分散システムにおける例を示している。一方のコンピュータシステム*i*ではタスク*x*が実行され、他方のコンピュータシステム*j*ではタスク*y*が実行されている。また、各コンピュータシステム*i*、*j*でアクセスできる資源として、*A*、*B*という2つの資源があるとす。なお、資源とは、タスクに割り当てられるプログラム、ファイル、データ等のソフトウェアを指す。ここでは、各コンピュータシステム*i*、*j*外に存在するデータベースの中身（ページ、レコード等）として説明する。

【0005】

図3において、タスク*x*は資源*A*をロックしており、タスク*y*は*B*をロックしている。同時に、タスク*y*は資源*A*をも必要としているので、資源*A*をロックすることを待っている。同様に、タスク*x*は資源*B*をも必要としているので、資源*B*がロックできるようになることを待っている。この場合、タスク*x*が資源*A*のロックを解除しない限り、タスク*y*は資源*A*をロックできない。一方、タスク*y*が資源*B*のロックを解除しない限り、タスク*x*は資源*B*をロックできない。この結果、タスク*x*、*y*は互いがロックしている*A*、*B*を待ち合って停止する。両タスク*x*、*y*が停止すると、各々が既にロックしている資源*A*、*B*の解除もできなくなってしまうので、この状態は永遠に続くことになる。よって、各タスクはそれ以上のプロセスを実行できない。

【0006】

このようなデッドロックは、コンピュータシステムがマルチプロセッサ方式のシステムであるかシングルプロセッサ方式のシステムであるか、あるいは、コンピュータシステムがスタンドアロンで運用されるのか分散処理システムを構成するのかに拘らず、システムがマルチタスクシステムであれば生じ得る問題である。

【0007】

このようなデッドロックが生じたとき、これを修復する手段を講じなければならない。そのためには、前提としてデッドロックが生じたことを検出しなければならない。

【0008】

デッドロック検出には、実用性を向上させる理由から、以下のスペックを満たすことが要求される。

第1に、実際はデッドロックではないにも拘らずデッドロックと誤認してしまう現象、すなわち疑似デッドロック (phantom deadlock) の検出が防止されていなければならない (第1の要求)。

【0009】

第2に、全てのデッドロックが検出されなければならない。換言すれば、現実にはデッドロックが生じている場合には、デッドロックを検出できるときと検出できないときがあってはならず、全てデッドロックであると検出されなければならない (第2の要求)。

【0010】

第3に、デッドロック検出を行うことによるシステムへの影響を小さく抑えなければならない。即ち、デッドロックを検出するためにタスクを停止するようなことは、できるだけ避けなければならない (第3の要求)。

【0011】

なお、マルチタスクシステムを分散処理システム上で実現する場合には、上記各要求の他に次のスペックが要求される。即ち、デッドロックを検出するためにシステム間で通信を行う必要があるが、この通信のオーバーヘッドをできるだけ削減しなければならない (第4の要求)。

【0012】

従来のデッドロック検出装置では、以下のようなような条件を満足させることによって、上述した第1乃至第3の要求を満足してデッドロックを検出しようとしていた。その条件とは、

- (a) トランザクションの非同期 abort (異常終了) が発生しないこと、
- (b) マルチタスクシステムを分散処理システム上で実現する場合には、各システム間の通信メッセージの遅延・消失が発生しないこと、
- (c) デッドロック検出中のトランザクション待ち関係の変更がないこと。

【0013】

(d) マルチタスクシステムを分散処理システム上で実現する場合には、システムの非同期ダウンが発生していないこと、である。

【0014】

上述の第1乃至第3の要求と(a)乃至(d)の条件との関係の関係を説明する。

(a)の条件に関し、トランザクションの非同期 abort (異常終了) が発生すると、前記第1の要求における疑似デッドロック (phantom deadlock) の検出防止を図ることができない。例えば、タスク (トランザクション) x が資源 A をロックしており、タスク (トランザクション) y が資源 B をロックしている場合、タスク y が資源 A を待ち、タスク x が資源 B に待ち要求を出した時点で、タスク y が非同期に異常終了してしまったとする。この場合、タスク y の非同期終了によってタスク y による資源 B のロックが解除されるので、タスク x は資源 B をロックできる。従って、本来ならばデッドロックは発生しないはずである。しかし、タスク y の非同期異常終了による資源 B のロック解除は直ちに検出できないので、現実には、デッドロックが発生していないにもかかわらずデッドロックが発生したものとして扱われてしまう。

【0015】

10

20

30

40

50

(b)の条件に関し、システムの非同期ダウンが生じると、デッドロックを検出できるときと検出できないときが生じ、前記第2の要求における全デッドロックの検出を行うことができない。なぜならば、分散処理システムにおいてはシステム間の通信によって共通資源へのアクセスをするわけであるが、その通信メッセージの伝達が遅れる場合や通信異常により消失する場合があると、デッドロックが発生したこと自体不明となるからである。

【0016】

同様に、(d)の条件に関し、システム間の通信メッセージの遅延・消失が生じると、デッドロックを検出できるときと検出できないときが生じ、前記第2の要求における全デッドロックの検出を行うことができない。なぜならば、一方のシステム作動中に他方のシステムがダウンしてしまうと、システムにおけるタスクに関する管理情報が失われ、デッドロック検出の判定ができなくなるからである。

10

【0017】

さらに、(c)の条件に関し、デッドロック検出中にトランザクション待ち関係の変更があると、現実にはどのタスクがどの資源をロックしているかに関する情報が混乱してしまうので、第1又は第2の要求を満たすことができない。

【0018】

【発明が解決しようとする課題】

しかしながら、(c)の条件は、デッドロック検出中における新たなトランザクション(タスク)の発生や待ち関係の発生を全て禁止することを内容とするものである。すなわち、デッドロック検出のためには、システムにおける要求受付を一旦停止する必要があるとする条件である。従って、本来デッドロックに関係ない資源(その資源を仮に資源Zとする。)に要求を出しているタスクがあっても、その要求を停止しなければならないことになる。従って、この条件(c)を追求すると、かえってシステムの円滑な運用が図れなくなり、第3の要求を満足できない結果となる。

20

【0019】

なお、条件(a)に起因する疑似デッドロックの検出を防ぐことは、現実には不可能である。すなわち、各システムにどのような異常が生ずるかを予想しこれをすべて回避することは不可能であって、タスクが非同期に異常終了することは防止することはできないからである。

【0020】

そこで、本発明の第1の技術的課題は、以上の問題点を鑑み、デッドロック検出中のトランザクション待ち関係の変更があってもデッドロック検出を継続でき、それによりデッドロック検出を行うことによるシステムへの影響を小さくすることができるデッドロック検出装置を提供することである。

30

【0021】

なお、本発明の第2の技術的課題は、分散処理システムを対象としたデッドロック検出装置において、システム間の通信メッセージの遅延・消失が発生した場合、デッドロック検出中のトランザクション待ち関係の変更が有った場合、及びシステムの非同期ダウンが生じた場合の何れにおいても、全てのデッドロックを検出でき、疑似デッドロックを検出せず、通信のオーバーヘッドをできるだけ削減でき、デッドロック検出を行うことによるシステムへの影響を小さくすることができるデッドロック検出装置を提供することをである。

40

【0022】

【課題を解決するための手段】

本発明は、前記第1の課題を解決するために、図1の原理図のように、以下の手段を採用した。

【0023】

<本発明の要旨>

即ち、複数のタスク100が共通の資源101を利用するマルチタスクシステムにおいて前記複数のタスク100が互いに占有している資源100を待ち合って停止してしまうデ

50

ッドロックを検出するためのデッドロック検出装置であって、複数のタスク100を並列実行するために、前記タスク100の実行を管理するタスク管理部(TM)102と、各タスクがどの資源100をロックしているかを管理するロック管理部(LM)103と、一のタスクが他のタスクがロックしている資源を獲得要求した場合には、前記一のタスクが前記他のタスクを待っているとしてこの各タスクの「待ち関係」を登録する待ち管理テーブル(LT)105と、前記ロック管理部(LM)103と非同期で動作するとともに、前記待ち管理テーブル(LM)103に登録された「待ち関係」からデッドロックを検出するデッドロック検出部(DD)104とを備えたことを特徴とする。

【0024】

以下に、本発明の構成要素の概要と、そのポイントを簡単にまとめる。

10

【0025】

〔タスク〕

“タスク”とは、通常CPU内部における仕事の単位を意味する。本発明においては、“タスク”を“トランザクション”と言い替えることができる。この“トランザクション”とは、ひとつの完結したデータ操作を行うオペレーションの集まりを意味し、“タスク”に含まれる概念であり、プログラムによって実行されるものである。要するに、本発明は、複数のプログラムが同時に並行して実行されるとき、各プログラムが資源を共有してロック状態となるのを検出しようとするものである。よって、“タスク”との用語を用いても、“トランザクション”との用語を用いても、単にプログラムの実行単位との用語を用いても、本発明においては、用語の差異は特に問題とはならない。以下、“タスク”=“トランザクション”と理解しても本発明の実施において何等の問題もない。また、本発明において、各タスクが共有する資源とは、データの集合であるファイルやファイルの中の下層的に記録されたレコードなどである。本発明でロックとは、或るタスク又はトランザクションがファイル全体を占有すること、あるいは、ファイルの下の或るレコードを占有することをいう。

20

【0026】

〔デッドロック検出〕

デッドロック検出部(DD)104におけるデッドロック検出は、例えば、次の通りに行うことができる。即ち、タスク(トランザクション)と資源の占有関係、即ち「待ち関係」を前記待ち管理テーブル(LT)105により登録する。この待ち管理テーブル(LT)105をデッドロック検出部(DD)104が見て、デッドロックを検出する。この検出は、前記ロック管理部(LM)103によるロック管理とは別個に行われる。好ましくは、前記ロック管理部(LM)103が、「あるタスクがある資源について「待ち関係」となった」ことを検出したとき、デッドロック検出部(DD)104に「待ち関係」を待ち管理テーブル(LT)105に登録するよう要求する。そして、その登録内容を参照することでデッドロックの有無を判定する。

30

【0027】

デッドロック検出のためにトランザクションの待ち関係を待ち管理テーブル(LT)105に登録する方法としては、以下の方法が好適である。即ち、トランザクションの待ち関係をグラフによって表現する。このグラフを、ここではウェイトフォークラフ(WFG: Wait-for-graph)と呼ぶ。このグラフを前記待ち管理テーブル(LT)105に登録するのである。

40

【0028】

このグラフにおいて、システム*i*で発生したトランザクション*x*を $T(i, x)$ で定義し、システム*j*で発生したトランザクション*y*を $T(j, y)$ で定義する。また、 $T(i, x)$ が $T(j, y)$ について待つこと、即ち $T(j, y)$ がロックしている資源を解放するのを $T(i, x)$ が待つことを、

$T(i, x) \rightarrow T(j, y)$

と表すこととする。この場合、 $T(j, y)$ が終了しない限り、 $T(i, x)$ はそれ以上プロセスを進めることができない。

50

【0029】

$T(i, x) \quad T(j, y)$

と

$T(j, y) \quad T(i, x)$

とが同時に成立したときには、

$T(i, x) \quad T(j, y) \quad T(i, x)$

というループが形成される。この場合はデッドロック状態であるので、このループを検出することにより、デッドロックの検出をすることができる。

【0030】

ところで、以上は2つのシステム間でのデッドロック検出例の説明であるが、自システム内でのデッドロックは、

$T(i, x) \quad T(i, y) \quad T(i, x)$

で表現される。

【0031】

本発明の特徴点は、タスク管理部(TM)102やロック管理部(LM)103等、タスクの実行に必要なブロックから独立したデッドロック検出部(DD)104をシステム(i, j)内に設け、このロック管理部(LM)103とは非同期にデッドロック検出部(DD)104を作動させる点にある。

【0032】

従来のデッドロック検出方法においては、デッドロックを検出するために、各タスクの実行を一旦停止させていた。そして、その間に、ロック管理部(LM)103におけるロック情報に基づいてデッドロックの有無を判定していた。しかしながら、これでは、タスクの円滑な実行を確保できない。

【0033】

これに対し、本発明では、タスク管理部(TM)102やロック管理部(LM)103から分離した待ち管理テーブル(LT)を設けて、前記ロック管理部(LM)103からのロック情報を待ち管理テーブル(LT)105に登録しておく。そして、タスクの実行とは独立して作動するデッドロック検出部(DD)104を設けた。このデッドロック検出部(DD)104は、あるタスクが資源を待つ状態に入ったという情報をロック管理部(LM)103が受けた時、そのタスクの待ちの関係も踏まえて、前記待ち管理テーブル(LT)105の登録内容を見てデッドロックの有無を判定するように構成することができる。

【0034】

このように、本発明は、タスク管理部(TM)102やロック管理部(LM)103等、タスクの実行に必要なシステムから分離して、デッドロック検出部(DD)104を設け、独立して作動させるため、デッドロック検出の為にタスクの実行を停止する必要がない。

【0035】

ところで、デッドロック検出部(DD)104によるデッドロック検出は、前記ロック管理部(LM)103によってタスクの待ち関係が検出された時に行うのが好ましい。すなわち、ロック管理部(LM)103によりタスクの待ち関係が検出された場合には、デッドロック検出部(DD)104がその待ち関係を管理テーブル(LT)105に登録する。この登録は、デッドロック検出の開始用トリガーとなる。デッドロック検出部(DD)104は、この登録の通知を受けることを契機に、待ち管理テーブル(LT)105を参照してデッドロックの有無検出を行う。

【0036】

デッドロックが検出されたとき、いずれかのタスクを強制的に異常終了させなければデッドロックを修復できない。いずれのタスクを強制的に異常終了させるかはシステムにより異なる。例えば、タスクの開始時刻の遅い方が仕事量が少ないとみてそのタスクを終了させても良い。あるいは、仕事量を実際に計上して少ない仕事量のタスクを終了させるよう

10

20

30

40

50

にしても良い。

【0037】

<分散システムへの適用>

本発明によるデッドロック検出装置は、複数のシステムを有する分散システム上に実現することができる。この様な分散システムを採用する場合には、上記した第1の課題に加えて、第2の課題の達成を考慮しなければならない。この場合のデッドロック検出は以下の様になる。

【0038】

即ち、前記待ち管理テーブル(LT)105に前記「待ち関係」を登録する場合において、2以上のタスク(x, y)が同一システム内のものであれば、そのシステムに設けた待ち管理テーブル(LT)105に各タスクにおける「待ち関係」を登録すれば足りる。

10

【0039】

一方、あるシステムのタスクが他のシステムのタスクに対して「待ち」の状態にある場合には、その「待ち関係」を一方のシステムから他方のシステムの待ち管理テーブル(LT)105へ通知すれば良い。この通知を受けた待ち管理テーブル(LT)105は、この「待ち関係」を登録する。この登録と同時に、他方のシステムのデッドロック検出部(DD)104がその待ち管理テーブル(LT)105を見に行き、デッドロックの有無を判定することができる。逆に、「待ち関係」を他方のシステムから一方のシステムの待ち管理テーブル(LT)105へ通知すれば、一方のシステムで、その待ち管理テーブル(LT)105を見てデッドロックの有無を判定できる。以上は、自己のシステムの待ち管理

20

【0040】

これとは別に、自己のシステムの「待ち関係」のみを自己の管理テーブル(LT)105に登録するにしても良い。この場合には、デッドロック検出をする際に、待ち先のシステムの待ち管理テーブル(LT)105に通信でアクセスする。そして、自己のシステムのタスクの「待ち関係」と、その待ち先のタスクの「待ち関係」とを突き合わせ、上述したループが形成されていればデッドロックとして検出することができる。

【0041】

本発明を分散システムに適用した場合、自己システム内でのデッドロックに対しては、複数のシステム相互間でデッドロック検出のための情報の通信は行わない。即ち、複数のシステム間でデッドロックが発生した場合のみ情報の通信を行う。但し、デッドロックは二者間で生じることがほとんどであるので、1回の通信でデッドロックを検出できる。従って、上記第2の課題を達成することができる。

30

【0042】

<待ち時間管理テーブルの付加>

上記した本発明の必須の構成要件に、図2の原理図に示すように、待ち時間監視部(WT)106を設けてもよい。この待ち時間監視部(WT)106は、あるタスク(トランザクション)について「待ち関係」が一定時間継続している場合に、そのタスク(トランザクション)について再度資源獲得要求を出すブロックである。この待ち時間監視部(WT)106を設ける目的は、以下に説明する通りである。

40

【0043】

即ち、デッドロックが生じても、デッドロックを検出できないと、デッドロックを修復できない。デッドロックが発生しているにも拘らずデッドロックを検出することができない原因としては、通信の欠落により管理テーブル(LT)105に「待ち関係」の情報が登録されていないことが考えられる。そこで、あるタスクにつき「待ち関係」が一定時間継続している場合には、待ち時間監視部(WT)106が再度資源獲得要求を出すようにするのである。これにより、「待ち関係」の情報をその待ち先のシステムの待ち管理テーブル(LT)105に再度送信する契機を与えることができる。従って、確実にデッドロックを検出することができる。

50

【 0 0 4 4 】

【作用】

本発明によるデッドロック検出装置では、タスクの実行状況は、タスク管理部 (T M) 1 0 2 によって管理される。この際、各タスクが資源を占有する場合には、ロック管理部 (L M) 1 0 3 によって、どのタスクがどの資源を占有したかについての情報が管理される。そして、他のタスクが占有している資源を一のタスクが獲得要求すると、この一のタスクは他のタスクの終了を待たねばならない。この「待ち関係」は、待ち管理テーブル (L T) 1 0 5 において登録管理される。

【 0 0 4 5 】

この待ち管理テーブル (L T) 1 0 5 をデッドロック検出部 (D D) 1 0 4 が見て、デッドロックを検出する。この検出は、ロック管理部 (L M) 1 0 3 によるロック管理とは別個に行われる。従って、デッドロック検出部 (D D) 1 0 4 によるデッドロックが行われていても、ロック管理部 (L M) 1 0 3 はその動作を行うことができる。よって、新たなタスクの発生や待ち関係の発生を禁ずる必要がなくなる。そのため、デッドロック検出を行うことによるシステムへの影響を小さくすることができるのである。

10

【 0 0 4 6 】

【実施例】

以下、本発明の好適実施例を、図面を参照して説明する。ここでは、今まで使用した“タスク”という言葉を“トランザクション”で置き換えて説明する。また、この好適実施例は、本発明を分散処理システムにおいて実施する場合の具体例である。

20

【 0 0 4 7 】

<システムの概要>

図 4 には分散処理システムの構成が示されている。この分散処理システムにおいては、二つのコンピュータシステム (システム i 及びシステム j) が分散して設けられ、相互にネットワーク (N W) 3 0 によって接続されている。また、両コンピュータシステム (i , j) とネットワーク (N W) 3 0 によって接続され、且つ両コンピュータシステム (i , j) からアクセス可能なデータベース (D B) 2 0 が設けられている。このようなシステムは、例えば、預金システムに利用される。

【 0 0 4 8 】

図 4 から明かなように、各コンピュータシステム (システム i 及びシステム j) は、トランザクション管理部 (T M) 1 0 , 資源管理部 (R M) 1 1 , ロック管理部 (L M) 1 2 , デッドロック検出部 (D D) 1 5 , 待ち管理テーブル T 3 , 及びウォッチドックタイマ (W T) 1 3 を備えている。なお、システム j はシステム i と全く同じ構成を有している。そのため、図 4 においては、システム i についてのみその詳細な構成を示し、システム j についてはその詳細な構成の図示を省略した。

30

【 0 0 4 9 】

データベース (D B) 2 0 には、資源としてのファイル又はレコードが複数個格納されている。図 4 においては、これら資源として、資源 A 及び資源 B を例示した。

【 0 0 5 0 】

以下、各構成ブロックを詳細に説明する。

40

【 0 0 5 1 】

<トランザクション管理部 (T M) >

トランザクション管理部 (T M) 1 0 は、複数のトランザクションの実行を管理している。トランザクション管理部 (T M) 1 0 を、タスク管理部 (T M) 1 0 と言っても良い。

【 0 0 5 2 】

このトランザクション管理部 (T M) 1 0 は、応用プログラムからのトランザクション開始・正常終了 (commit) ・異常終了 (abort) の通信を受付け、システム内でのトランザクションを管理するブロックである。より詳しく言うと、例えばシステム i においてトランザクション x が開始されたときには、 T (i , x) という形式のデータを登録し、トランザクション x が終了・異常終了したときにはこの T (i , x) という形式のデータを削

50

除するのである。

【 0 0 5 3 】

トランザクション管理部 (T M) 1 0 は、トランザクションから資源の要求を受け付けると、その要求を資源管理部 (R M) 1 1 に渡し、その応答 (o k / n o) をもらう。また、トランザクション管理部 (T M) 1 0 は、デッドロック検出部 (D D) 1 5 から送信されたトランザクションのデッドロック通知を受け付けて、トランザクションを終了させる。また、トランザクション管理部 (T M) 1 0 は、デッドロック検出部 (D D) 1 5 から送信されたリトライ通知を受け付けて、資源管理部 (R M) 1 1 に資源獲得要求を再発行する。また、トランザクション管理部 (T M) 1 0 は、他のコンピュータシステム (システム j) のトランザクションの正常終了や異常終了の通信を受信し、デッドロック検出部 (D D) 1 5 にグラフの登録・削除を要求する。

10

【 0 0 5 4 】

図 4 において、“ s t a r t ” はトランザクションの実行開始を意味し、“ a b o r t ” はトランザクションの異常終了を意味し、“ c o m m i t ” はトランザクションの正常終了を意味する。

【 0 0 5 5 】

トランザクション管理部 (T M) 1 0 における、資源獲得要求・資源解放要求は、二相ロック (2 P L : T w o P h a s e L o c k) 方式で実行される。これは、疑似デッドロック (p h a n t o m d e a d l o c k) の検出を防止するのに効果的である。

【 0 0 5 6 】

疑似デッドロックは、一般にグラフの登録と削除が競合した場合に発生する。例えば、 $T(i, x)$ $T(j, y)$ のグラフが既に登録されている場合において、 $T(i, x)$ $T(j, y)$ の削除要求と $T(j, y)$ $T(i, x)$ の登録要求とが、デッドロック検出部 (D D) 1 5 に対して同時に発生したとする。この際、削除要求が先に受理された場合にはデッドロックが発生しない。これに対して、登録要求が先に受理された場合には、疑似デッドロックとなる。

20

【 0 0 5 7 】

$T(i, x)$ $T(j, y)$ の削除要求が発生するのは、このグラフで表される待ち関係がなくなった場合である (即ち、 $T(j, y)$ が資源のロックを解除した場合である。) 。ロックの解除は、トランザクションが自ら資源のロックを解除する場合が非同期に a b o r t (異常終了) する場合に行われる。

30

【 0 0 5 8 】

二相ロック方式は、ある処理がデータのロック (占有) を始めたらロックし続け、ロックを解除し始めたら解除し続けるという 2 つの相 (フェーズ) からなるロック方式である。この方式によれば、複数のタスクやトランザクションがそれぞれ逐次実行されたのと同じ結果となる。

40

【 0 0 5 9 】

トランザクションが非同期に a b o r t (異常終了) しないとすれば、この方式により、一旦ロックが解除されると、新たなロック獲得の要求が発生しないことを保証する。つまり、一旦発生したグラフ

$T(i, x)$ $T(j, y)$

は、 $T(j, y)$ が終了するまで削除されることはない。従って、上記の

$T(i, x)$ $T(j, y)$

の削除要求が発生した後においては、 $T(j, y)$ が終了しているので、 $T(j, y)$ が新たな資源獲得要求をすることはない。よって、

$T(j, y)$ $T(i, x)$

50

の登録要求が発生しないことを保証できる。

【0060】

以上の理由により、この方式を採用することによって、疑似デッドロックの発生原因をトランザクションの非同期 abort (異常終了) のみに限定することができる。

【0061】

<資源管理部 (RM)>

資源管理部 (RM) 11 は、このトランザクション管理部 (TM) 10 に双方向で接続されている。資源管理部 (RM) 11 は、資源管理テーブル T1 を有している。資源管理部 (RM) 11 は、トランザクション管理部 (TM) 10 からの資源獲得要求及び資源解放要求の内容に基づいて、トランザクションとそのトランザクションが要求している資源との対応関係を、資源管理テーブル T1 上にマッピングして管理している。

10

【0062】

また、資源管理部 (RM) 11 は、トランザクション管理部 (TM) 10 からの資源獲得要求に応じてロック要求をロック管理部 (LM) 12 に対して行い、トランザクション管理部 (TM) 10 からの資源解放要求に応じて、ロック解放要求をロック管理部 (LM) 12 に対して行う。

【0063】

<ロック管理部 (LM)>

ロック管理部 (LM) 12 は、資源管理部 (RM) 11 に双方向で接続されている。ロック管理部 (LM) 12 は、ロック管理テーブル T2 を有している。即ち、ロック管理部 (LM) 12 は、このロック管理テーブル T2 によりロック状態の管理を行う制御部である。

20

【0064】

トランザクション x, y と資源 A, B がある場合において、トランザクション x が資源 A をロック (占有) し、トランザクション y が B をロックしたときには、ロック管理部 (LM) 12 は、この関係をロック管理テーブル T2 に登録する。即ち、図 4 に示すように、x が A をロックした状態を例えば (x : A) と定義し、y が B をロックした状態を例えば (y : B) と定義し、この情報をロック管理テーブル T2 に登録する。

【0065】

なお、このロック管理テーブル T2 は、そのコンピュータシステム (i 又は j) におけるトランザクションについてのロック情報を管理するばかりでなく、他のコンピュータシステム (j 又は i) におけるトランザクションについてのロック情報をも管理する。この他のシステムにおけるトランザクションについてのロック情報は、コンピュータシステム間で通信を行うことにより獲得することができる。但し、各コンピュータシステム (i, j) によって共用される共用メモリ上に単一のロック管理テーブル T2 を作成し、全コンピュータシステム (i, j) における全トランザクションに関するロック情報を一括管理させれば、各コンピュータシステム間における通信の必要はなくなる。

30

【0066】

いま、上述した状態において、更に資源管理部 (RM) 11 から、トランザクション y による資源 A のロック要求がなされ、トランザクション x による資源 B のロック要求がなされるとする。そうすると、ロック管理部 (LM) 12 はロック管理テーブル T2 の情報を参照し、このようなロックができないことを認識する。この場合には、トランザクション y はトランザクション x による資源 A のロック解放を待ち、トランザクション x はトランザクション y による資源 B のロック解放を待たねばならない。この待ち状態は、それぞれ、(x : B)、(y : A) と定義される。このような定義が発生したとき、ロック管理部 (LM) 12 は「待ち」が発生したと判断するのである。

40

【0067】

本実施例では、このような「待ち関係」を、ロック管理部 (LM) 12 とは切り離して、待ち管理テーブル T3 に登録して管理する。即ち、「待ち関係」が生じたときには、ロック管理部 (LM) 12 は、上記定義に基づいて、ウェイトフォークラフ登録をデッドロッ

50

ク検出部 (DD) 15 に要求する。この要求の際には、上記定義におけるトランザクション (x, y) がどのコンピュータシステムにおけるトランザクションであるのか、及び、上記定義における資源 (A, B) が現在どのコンピュータシステムのどのトランザクションによってロックされているのかの情報も、デッドロック検出部 (DD) 15 に通知する。

【0068】

ロック管理部 (LM) 12 は、資源管理部 (RM) 11 からのロック要求が「待ち」にならない場合には、資源管理部 (RM) 11 に対してすぐに応答 (ok) を返す。ロック管理部 (LM) 12 が判断して「待ち」が発生した場合のみ、待ち関係を示すウェイトフォークラフを、待ち管理テーブル T3 に登録する登録要求キューを発行する。従って、「待ちが発生しない資源獲得要求」に関しては、デッドロック検出中か否かに拘らず、その資源獲得要求を行ったトランザクションの実行処理は停止されない。

10

【0069】

<デッドロック検出部 (DD) >

デッドロック検出部 (DD) 15 は、待ち管理テーブル T3 の登録内容からデッドロックの有無を判定する部分である。

【0070】

デッドロック検出部 (DD) 15 は、要求キュー受付部 (QR) 14 を有する。この要求キュー受付部 (QR) 14 は自システムのロック管理部 (LM) 12 から待ち関係 (ウェイトフォークラフ) の登録・削除要求キューを受け付ける。また、他のコンピュータシステムからの待ち関係 (ウェイトフォークラフ) の登録・削除要求キューを受け付ける。さらに、他システムのトランザクションが異常終了 (abort) 又は正常終了 (commit) した場合には、トランザクション管理部 (TM) 10 からの待ち関係 (ウェイトフォークラフ) の削除要求キューを受け付ける。

20

【0071】

デッドロック検出部 (DD) 15 は、これら要求キューに従い、先ずウェイトフォークラフの登録又は削除を、待ち管理テーブル T3 に対して行う。このウェイトフォークラフ (Wait-for-graph) の形式は以下の通りである。即ち、例えば、システム i で発生したトランザクション $x = T(i, x)$ 、システム j で発生したトランザクション $y = T(j, y)$ としたとき、 $T(i, x)$ が $T(j, y)$ について待つことを $T(i, x) \rightarrow T(j, y)$ と表す。

30

【0072】

ウェイトフォークラフの登録を行う際には、デッドロック検出部 (DD) 15 は、通知された情報に基づいて、予めウェイトフォークラフを作成する。ウェイトフォークラフの登録がなされると、デッドロック検出部 (DD) 15 はデッドロック検出を開始する。デッドロックが検出されたときは、デッドロック検出部 (DD) 15 は、トランザクション管理部 (TM) 10 にデッドロック通知を行う。ウェイトフォークラフの削除要求を受け付けた場合は、当該ウェイトフォークラフを削除し、動作できるトランザクションに対しリトライ通知を行う。

40

【0073】

<待ち管理テーブル T3 >

待ち管理テーブル T3 には、ウェイトフォークラフが登録される。上述した通り、システム i で発生したトランザクション x (資源 A を占有中) = $T(i, x)$ 、システム j で発生したトランザクション y (資源 B を占有中) = $T(j, y)$ としたとき、 $T(i, x)$ が $T(j, y)$ について待つことを $T(i, x) \rightarrow T(j, y)$

と表す。この場合、トランザクション y が占有している資源 B を更新して $T(j, y)$ が終了しない限り、トランザクション x は獲得しようとしている資源 B を使用できない。この状態を「待ち関係」といい、「 $T(i, x)$ が $T(j, y)$ について待つ」という。

50

【 0 0 7 4 】

デッドロック検出部 (D D) 1 5 は、待ち管理テーブル T 3 に、この

$T (i , x) \quad T (j , y)$

のグラフを、「待ち関係」として登録する。

【 0 0 7 5 】

一方、この

$T (i , x) \quad T (j , y)$

のグラフの成立と同時に

$T (j , y) \quad T (i , x)$

が成立していることがある。この場合、トランザクション x が占有している資源 A を更新して T (i , x) が終了しない限り、トランザクション y は獲得しようとしている資源 A を使用できない。この 2 つの待ち関係を突き合わせると、

$T (i , x) \quad T (j , y) \quad T (j , y) \quad T (i , x)$

というループが形成される。よって、このループが検出されればデッドロックが発生しているといえることができるのである。

【 0 0 7 6 】

この待ち関係は、システム i のトランザクション x とシステム j のトランザクション y との間で生じている。そして、

$T (i , x) \quad T (j , y)$

のグラフはシステム i の待ち管理テーブル T 3 に登録され、

$T (j , y) \quad T (i , x)$

のグラフはシステム j の待ち管理テーブル (T 3) に登録される。このため、両者を突き合わせるためには、いずれかを他方に送信しなければならない。ここでは、「待ち関係」が登録される時、その待ち先に「待ち関係」を送信する。

【 0 0 7 7 】

すなわち、

$T (i , x) \quad T (j , y)$

がシステム i の待ち管理テーブル (T 3) に登録されたとき、システム i のデッドロック検出部 (D D) 1 5 は、システム j の待ち管理テーブル T 3 に同一の内容のグラフを送信して登録する。

【 0 0 7 8 】

これにより、待ち先のシステム (即ち、システム j) において、待ち管理テーブル T 3 を参照すれば、デッドロックを検出できる。なお、トランザクション x の「待ち関係」が解消したとき、トランザクション x に関する「待ち関係」の情報を待ち先のシステム (即ち、システム j) から回収 (削除) しないと、いつまでもデッドロックを検出してしまう。そこで、「待ち関係」が解消した場合には、待ち先のシステム (即ち、システム j) の待ち管理テーブル T 3 から「待ち関係」を示すグラフを削除・あるいは回収しなければならない。デッドロック検出部 (D D) 1 5 は、このようなグラフ削除・回収機能をも有する。

【 0 0 7 9 】

ウェイトフォークグラフのループが検出されない場合、ウェイトフォークグラフの先端が他のシステムのトランザクションであれば、当該他のシステムにデッドロックの可能性のあることになる。そこで、当該他のシステムにグラフ登録を要求する。他のシステムからのグラフ登録の要求を受け付けた場合は、デッドロック検出部 (D D) 1 5 は必要なグラフを登録し、ループ検出を行う。

【 0 0 8 0 】

この待ち関係は、自システム内で生じるときがある。例えば、システム i で発生したトランザクション $x_1 = T (i , x_1)$ が自システム (即ち、システム i) で発生したトランザクション $y_1 = T (i , y_1)$ にいて待つとき、自システム (即ち、システム i) の待ち管理テーブル T 3 に、

10

20

30

40

50

$T(i, x1)$ $T(i, y1)$

が登録される。ここで、自システム（即ち、システム i ）の待ち管理テーブル $T3$ に、

$T(i, y1)$ $T(i, x1)$

が登録されているなら、

$T(i, x1)$ $T(i, y1)$ $T(i, y1)$ $T(i, x1)$

というループが形成されるので、デッドロックが検出できる。

【0081】

ところで、分散処理システムにあっては、あるシステムで発生したトランザクションに
関係するウェイトフォークラフ等の登録内容を、そのシステムのローカルウェイトフォーク
ラフという。また、分散処理システム全体での待ち関係を表現したグラフ、即ち、その分
散システムにおける全ローカルウェイトフォークラフの集合を、グローバルウェイトフォ
ークラフという。図5は、ローカルウェイトフォークラフとグローバルウェイトフォーク
ラフの関係の例を示したものである。

10

【0082】

各コンピュータシステムでは、ローカルウェイトフォークラフのみを待ち管理テーブル $T3$
で管理する。ここでは、前記したように、デッドロック検出部（ DD ）15は、自コン
ピュータシステム内のトランザクション間の待ち関係を表すウェイトフォークラフ等の登
録内容を他のコンピュータシステムに送信しない。一方、デッドロック検出部（ DD ）1
5は、他のコンピュータシステム内のトランザクションとの間の待ち関係を表すウェイト
フォークラフを、関係を持った他システムにのみ送信する。統計的に見てデッドロックの
90%以上が2つのタスク間で発生することを考慮すると、他システムのトランザクショ
ンに関連するデッドロックであっても、ほとんど1回の通信で検出することができる。し
かも、自コンピュータシステム内でのデッドロックであれば、通信なしで検出できる。従
って、デッドロック検出のための通信のオーバーヘッドを削減できる。

20

【0083】

<待ち時間監視部（ WT ）>

次に、待ち時間監視部（ WT ）13は、待ち管理テーブル $T3$ を監視するタイマーである
。このタイマーは、待ち管理テーブル $T3$ に登録されている「待ち関係」を監視する。そ
して、その「待ち関係」が登録されてから一定時間経過した時点で、なおその「待ち関係
」が継続しているならば、その「待ち関係」にある待ち元のトランザクションが資源獲得
要求を再発行するように、リトライ通知を発行する。このリトライ通知は、要求キュー受
付部（ QR ）14に投入される。このリトライ通知が要求キュー受付部（ QR ）14に投
入されると、デッドロック検出部（ DD ）15は、トランザクション管理部（ TM ）10
にリトライ通知を送る。トランザクション管理部（ TM ）10は、このリトライ信号を受
けて、待ち関係にあるトランザクションに対し、再度資源獲得要求を出す。

30

【0084】

システム間の通信メッセージの遅延・消失やコンピュータシステムの非同期ダウンが発生
すると、実際はデッドロック状態であるのにこれを検出できない場合がある。待ち時間監
視部（ WT ）13は、このような不都合を防止する。すなわち、分散処理システムにおい
て各コンピュータシステム間で通信するとき、通信の欠落によりデッドロック状態を表示
するグラフが欠落することがある。すると、デッドロック検出ができなくなる。これを防
止するために、待ち関係にあるトランザクションを監視する機構として、前記待ち時間監
視部（ WT : Watchdog Timer）13を設けたのである。

40

【0085】

このタイマーは、一定時間以上ウェイトフォークラフの待ち先の関係にあるトランザク
ションに対し、前述した様に、デッドロック検出部（ DD ）15を介して、再度資源獲得要
求することを促す。すると、トランザクション監視部（ TM : Transaction
Manager）10から、再度資源獲得要求が出される。このとき、既に「待ち」が解
消されているなら、この資源獲得要求は満たされる。これに対して、待ちが解消されてい
ないなら、他のコンピュータシステムに対して再度ウェイトフォークラフが送信される。

50

これによりウェイトフォークラフ欠落が補われ、デッドロックが検出できる。

【0086】

<各部の動作例>

以下、前記各部の動作をフローチャート図に従って説明する。

〔トランザクション管理部(TM)の動作〕

図6に示したフローチャートのように、トランザクション管理部(TM)10は、start(開始要求)、abort(異常終了)、commit(正常終了)、資源獲得要求、デッドロック通知、リトライ通知などの各種要求を待つ(ステップS101)。なお、ここで言うabort(異常終了)、commit(正常終了)には、他のコンピュータシステムから通知されたものも含む。何れかの要求を受け付ける(ステップS102)と

10

【0087】

ステップS102で受け付けた要求がstart(開始要求)の場合、トランザクション管理部(TM)10自身にそのトランザクション(ここでは、仮にT(i,x)とする。)を登録し(ステップS103)、その後の要求を待つ。

【0088】

ステップS102で受け付けた要求がabort(異常終了)又はcommit(正常終了)である場合、まず、その終了するトランザクション(ここでは、仮にT(i,x)とする。)を削除する(ステップS104)。次に、資源管理部(RM)11に対し、資源解放要求を発行する(ステップS105)。その資源開放要求に対する応答を資源管理部(RM)11から受けると(ステップS106)、デッドロック検出部(DD)15にグラフ削除要求を出す(ステップS107)。その後、その要求が自コンピュータシステムからの要求か否かを判定する(ステップS108)。他コンピュータシステムからの要求であればそのままとする。これに対して、自コンピュータシステムからの要求であれば、他のコンピュータシステムに、commit又はabortを通知する(ステップS109)。通知を受けた他のコンピュータシステムでは、ステップS104乃至107の処理を行う。

20

【0089】

ステップS102で受け付けた要求が資源獲得要求である場合、資源管理部(RM)11に資源獲得要求を出す(ステップS110)。その要求に対する応答を資源管理部(RM)11から受けたら(ステップS111)、トランザクションに応答を返す(ステップS112)。

30

【0090】

ステップS102で受け付けた要求がデッドロック通知である場合、まず、デッドロックとなっているトランザクションの中からabortさせるべきトランザクションを選択する(ステップS120)。即ち、デッドロック通知には、デッドロックの関係にある全トランザクション(ここでは、仮にT(i,x)、T(j,y)とする。)の特定が含まれている。トランザクション管理部(TM)10は、このデッドロック通知に含まれているトランザクション名からabortさせるべきトランザクションを選択するのである。従って、トランザクション管理部(TM)10は、他のコンピュータシステムのトランザクションをも、abort対象として特定することができる。次いで、トランザクション管理部(TM)10は、選択されたトランザクションにabortすべき旨の通知をする(ステップS121)。選択されたトランザクションが他のコンピュータシステムのものである場合には、当該他のコンピュータシステムのトランザクション管理部(TM)10を介して、選択されたトランザクションにabortすべき旨を通知する。

40

【0091】

ステップS102で受け付けた要求がリトライ通知である場合、まず、資源管理部(RM)11に資源獲得要求を出す(ステップS130)。その要求に対する応答を資源管理部(RM)11からもらったら(ステップS131)、トランザクションに応答を返す(ステップS131)。

50

【 0 0 9 2 】

〔 資源管理部 (R M) の動作 〕

図 7 に示したフローチャートのように、資源管理部 (R M) 1 1 は、まず、資源獲得要求及び資源解放要求を待つ (ステップ S 2 0 1) 。何れかの要求があり、それが受理されると (ステップ S 2 0 2) 、テーブル上に示された資源をロックしようとし、その関係を資源獲得テーブル T 1 に登録する (ステップ S 2 0 3) 。即ち、どのトランザクションがどの資源をロックしようとするのかを登録する。

【 0 0 9 3 】

その後、要求が資源獲得要求か資源解放要求かを判定する (ステップ S 2 0 4) 。要求が資源獲得要求の場合、ロック管理部 (L M) 1 2 にロック獲得要求を出す (ステップ S 2 0 5) 。これに対して、要求が資源解放要求の場合、ロック管理部 (L M) 1 2 にロック解放要求を出す (ステップ S 2 0 6) 。

10

【 0 0 9 4 】

そして、ロック獲得要求又はロック解放要求に対する応答 (o k / n o) をロック管理部 (L M) 1 2 から受けた後 (ステップ S 2 0 7) 、トランザクション管理部 (T M) 1 0 に応答 (o k / n o) を返す (ステップ S 2 0 8) 。

【 0 0 9 5 】

〔 ロック管理部 (L M) の動作 〕

図 8 に示したフローチャートのように、資源管理部 (R M) 1 1 におけるステップ S 2 0 5 又はステップ S 2 0 6 の要求があると (ステップ S 3 0 1) 、ロック管理部 (L M) 1 2 は、要求を受け付ける (ステップ S 3 0 2) 。その後、要求がロック獲得要求かロック解放要求かを判定する (ステップ S 3 0 3) 。要求がロック獲得要求である場合には、ロック管理部 (L M) 1 2 は、ロック獲得が可能か否かを判定する (ステップ S 3 0 4) 。

20

【 0 0 9 6 】

資源のロックが可能であれば、そのロック状態をロック管理テーブル T 2 に登録する (ステップ S 3 0 5) 。資源のロックが不可能であれば、「待ち関係」であるので、要求側のトランザクションと待ち先のトランザクションとの関係をウェイトフォークラフとして待ち管理テーブル T 3 に登録する旨を、デッドロック検出部 (D D) 1 5 に対して要求する (ステップ S 3 0 6) 。その後、資源管理部 (R M) 1 1 にロックできなかった旨 (n o) を返答する (ステップ S 3 0 9) 。

30

【 0 0 9 7 】

ステップ S 3 0 3 において要求が資源解放要求であると判定された場合、ロック管理テーブル T 2 からロックの登録を削除する (ステップ S 3 0 7) 。ロックの登録 (ステップ S 3 0 5) とその削除 (ステップ S 3 0 7) の後は、その完了 (o k) を示す応答を、資源管理部 (R M) 1 1 に返す (ステップ S 3 0 8) 。

【 0 0 9 8 】

〔 デッドロック検出部 (D D) の動作 〕

図 9 に示したフローチャート図のように、デッドロック検出部 (D D) 1 5 には、グラフ登録要求、グラフ削除要求、及びリトライ通知が、要求キュー受付部 (Q R) 1 4 に受け付けられる。従って、その要求があると (ステップ S 4 0 1) 、要求キュー受付部 (Q R) 1 4 から要求を取り出し (ステップ S 4 0 2) 、要求の種別を判定する (ステップ S 4 0 3) 。

40

【 0 0 9 9 】

要求が、グラフ登録である場合には、先ず、待ち管理テーブル T 3 にウェイトフォークラフを登録する (ステップ S 4 0 4) 。但し、待ち管理テーブル T 3 を検索した結果同一グラフが既に登録されていれば、そのグラフは登録しない。次いで、登録したグラフにつき、グラフの先端までたどる (ステップ S 4 0 5) 。たどった結果によって、ループが形成されているか判断する (ステップ S 4 0 6) 。ループが形成されていれば、トランザクション管理部 (T M) 1 0 にデッドロックを通知する (ステップ S 4 0 7) 。ループが形成されていなければ、グラフの先端が自コンピュータシステムか否かを判定する (ステップ

50

S 4 0 8)。自コンピュータシステムであればそのままステップ S 4 0 1 に戻る。これに対して、他コンピュータシステムであれば、その他コンピュータシステムのデッドロック検出部 (D D) 1 5 に当該ウェイトフォークラフを送信して、その他システムの待ち管理テーブル T 3 に当該グラフを登録させる (ステップ S 4 0 9)。

【 0 1 0 0 】

次に、ステップ S 4 0 3 において、要求がグラフの削除であるときは、待ち管理テーブル T 3 を検索して、該当するグラフを探す (ステップ S 4 1 0)。該当グラフを探しあてたら、該当グラフを削除する (ステップ S 4 1 1)。その後、待ち関係が解除されたトランザクションを動作させるため、トランザクション管理部 (T M) 1 0 にリトライ通知をする (ステップ S 4 1 2)。

10

【 0 1 0 1 】

ステップ S 4 0 3 で要求がリトライ通知であるとき、まず、リトライするトランザクションのウェイトフォークラフを削除する (ステップ S 4 2 0)。次いで、トランザクション管理部 (T M) 1 0 にリトライ通知をする (ステップ S 4 2 1)。

【 0 1 0 2 】

〔 待ち時間監視部 (W T) の動作 〕

図 1 0 に示したフローチャートのように、待ち時間監視部 (W T) 1 3 は、待ち管理テーブル T 3 に登録されている各トランザクション (T (i , x) 等) を順次検索する (ステップ S 5 0 1)。次いで、そのトランザクション x が「待ち関係」にあるか否かを判断する (ステップ S 5 0 2)。検索されたトランザクションが「待ち関係」でなければ、ステップ S 5 0 1 に戻り、次のトランザクションを検索する。

20

【 0 1 0 3 】

これに対して、検索されたトランザクションが「待ち関係」であれば、タイムカウントを開始し、タイムアウトとなったら (ステップ S 5 0 3)、デッドロック検出部 (D D) 1 5 にリトライ通知を行う (ステップ S 5 0 4)。ステップ S 5 0 3 でタイムアウトになる前に「待ち関係」が解消されたなら、ステップ S 5 0 1 に戻る (ステップ S 5 0 2)。

【 0 1 0 4 】

< 具体的なデッドロック検出の例 >

次に、以上の構成におけるデッドロック検出例を、3通りの場合に沿って説明する。

〔 例 1 自コンピュータシステム内におけるデッドロック検出 〕

30

例 1 は、自コンピュータシステム内におけるデッドロック検出の例で、具体的には以下の動作を行う。ここでは、他のコンピュータシステムとの間に通信が発生しないことが解る。

【 0 1 0 5 】

(1) 先ず、トランザクション T (1 , 1) がトランザクション管理部 (T M) 1 0 に対してトランザクションの開始を通知したとする。

(2) すると、トランザクション管理部 (T M) 1 0 はデッドロック検出部 (D D) 1 5 にトランザクション T (1 , 1) の登録を要求する。

(3) 一方、トランザクション T (1 , 2) がトランザクション管理部 (T M) 1 0 に対してトランザクションの開始を通知したとする。

40

【 0 1 0 6 】

(4) すると、トランザクション管理部 (T M) 1 0 はデッドロック検出部 (D D) 1 5 にトランザクション T (1 , 2) の登録を要求する。

(5) いま、トランザクション T (1 , 1) がトランザクション管理部 (T M) 1 0 に資源 A を要求したとする。

(6) すると、トランザクション管理部 (T M) 1 0 が資源管理部 (R M) 1 1 に資源 A の獲得を要求する。

【 0 1 0 7 】

(7) すると、資源管理部 (R M) 1 1 がロック管理部 (L M) 1 2 に資源 A のロック獲得を要求する。

50

(8) 資源 A が未ロックであれば、ロック管理部 (LM) 12 が資源管理部 (RM) 11 に OK を応答する。

(9) すると、資源管理部 (RM) 11 がトランザクション管理部 (TM) 10 に OK を応答する。

【0108】

(10) 一方、トランザクション T (1, 2) がトランザクション管理部 (TM) 10 に資源 B を要求したとする。

(11) すると、トランザクション管理部 (TM) 10 が資源管理部 (RM) 11 に資源 B を獲得を要求する。

(12) すると、資源管理部 (RM) 11 がロック管理部 (LM) 12 に資源 B のロック獲得を要求する。 10

【0109】

(13) 資源 B が未ロックであれば、ロック管理部 (LM) 12 が資源管理部 (RM) 11 に OK を応答する。

(14) すると、資源管理部 (RM) 11 がトランザクション管理部 (TM) 10 に OK を応答する。

(15) この状態において、トランザクション T (1, 1) がトランザクション管理部 (TM) 10 に資源 B を要求したとする。

【0110】

(16) すると、トランザクション管理部 (TM) 10 が資源管理部 (RM) 11 に資源 B の獲得を要求する。 20

(17) すると、資源管理部 (RM) 11 がロック管理部 (LM) 12 に資源 B のロック獲得を要求する。

【0111】

(18) ところが、資源 B はトランザクション T (1, 2) によって既にロック済みであるので、T (1, 1) がトランザクション T (1, 2) に対して待つことになる。そこで、ロック管理部 (LM) 12 がデッドロック検出部 (DD) 15 にグラフ

T (1, 1) T (1, 2)

の登録を要求する。要求を受けたデッドロック検出部 (DD) 15 は、このグラフを待ち管理テーブル T3 に登録する。 30

【0112】

(19) 一方、トランザクション T (1, 2) がトランザクション管理部 (TM) 10 に資源 A を要求したとする。

(20) すると、トランザクション管理部 (TM) 10 が資源管理部 (RM) 11 に資源 A の獲得を要求する。

(21) すると、資源管理部 (RM) 11 がロック管理部 (LM) 12 に資源 A のロック獲得を要求する。

【0113】

(22) ところが、資源 A はトランザクション T (1, 1) によって既にロック済みであるので、T (1, 2) がトランザクション T (1, 1) に対して待つことになる。そこで、ロック管理部 (LM) 12 がデッドロック検出部 (DD) 15 にグラフ

T (1, 2) T (1, 1)

の登録を要求する。要求を受けたデッドロック検出部 (DD) 15 は、このグラフを待ち管理テーブル T3 に登録する。 40

(23) デッドロック検出部 (DD) 15 がループを検出し、デッドロック発生をトランザクション管理部 (TM) 10 に通知する。

【0114】

[例2 2つのコンピュータシステム間における2つのトランザクションのデッドロック検出]

例2は、2つのコンピュータシステム(システム1, システム2)間でデッドロックが発 50

生する場合を示している。ここでは、デッドロック検出のための通信が1回で済むことがわかる。

(1) 先ず、システム1におけるトランザクションT(1,1)が、システム1のトランザクション管理部(TM)10に対しトランザクションの開始を通知したとする。

(2) すると、システム1のトランザクション管理部(TM)10は、デッドロック検出部(DD)15にトランザクションT(1,1)の登録を要求する。

(3) いま、トランザクションT(1,1)が、システム1のトランザクション管理部(TM)10に、資源Aを要求したとする。

【0115】

(4) すると、システム1のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Aの獲得を要求する。 10

(5) すると、システム1の資源管理部(RM)11が、ロック管理部(LM)12に資源Aのロック獲得を要求する。

(6) 資源Aが未ロックであれば、システム1のロック管理部(LM)12が、資源管理部(RM)11にOKを応答する。

【0116】

(7) すると、システム1の資源管理部(RM)11が、トランザクション管理部(TM)10にOKを応答する。

(1)' 一方、システム2におけるトランザクショントランザクションT(2,1)が、システム2のトランザクション管理部(TM)10に対しトランザクションの開始を通知したとする。 20

(2)' すると、システム2のトランザクション管理部(TM)10は、デッドロック検出部(DD)15にトランザクションT(2,1)の登録を要求する。

【0117】

(3)' いま、トランザクションT(2,1)が、システム2のトランザクション管理部(TM)10に、資源Bを要求したとする。

(4)' すると、システム2のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Bの獲得を要求する。

(5)' すると、システム2の資源管理部(RM)11が、ロック管理部(LM)12に資源Bのロック獲得を要求する。 30

(6)' 資源Bが未ロックであれば、システム2のロック管理部(LM)12が、資源管理部(RM)11にOKを応答する。

【0118】

(7)' すると、システム2の資源管理部(RM)11が、トランザクション管理部(TM)10にOKを応答する。

(8) 以上の状況下において、トランザクションT(1,1)がシステム1のトランザクション管理部(TM)10に資源Bを要求したとする。

(9) すると、システム1のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Bの獲得を要求する。

【0119】

(10) すると、システム1の資源管理部(RM)11が、ロック管理部(LM)12に資源Bのロック獲得を要求する。 40

(11) ところが、資源Bはシステム2のトランザクションT(2,1)によって既にロック済みであるので、トランザクションT(1,1)がトランザクションT(2,1)に対して待つことになる。そこで、システム1のロック管理部(LM)12が、デッドロック検出部(DD)15にグラフ

T(1,1) T(2,1)

の登録を要求する。

【0120】

(12) この要求を受けて、システム1のデッドロック検出部(DD)15は、グラフT(50

1, 1) T(2, 1)を待ち管理テーブルT3に登録し、待ち管理テーブルT3に登録されたグラフにループが形成されているかどうか(デッドロックが発生しているかどうか)を判断し、デッドロックを検出しないときにシステム2にグラフT(1, 1) T(2, 1)を送信する。(13)システム2のデッドロック検出部(DD)15が、このグラフT(1, 1) (2, 1)を受信し、これをシステム2の待ち管理テーブルT3に登録する。

【0121】

(14) この後で、トランザクションT(2, 1)がシステム2のトランザクション管理部(TM)10に資源Aを要求したとする。

(15) すると、システム2のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Aの獲得を要求する。 10

(16) すると、システム2の資源管理部(RM)11が、ロック管理部(LM)12に資源Aのロック獲得を要求する。

【0122】

(17) ところが、資源Aはシステム1のトランザクションT(1, 1)によって既にロック済みであるので、トランザクションT(2, 1)がトランザクションT(1, 1)に対して待つことになる。そこで、システム2のロック管理部(LM)12が、デッドロック検出部(DD)15にグラフT(2, 1) T(1, 1)の登録を要求する。

(18) システム2のデッドロック検出部(DD)15は、グラフT(2, 1) T(1, 1)を待ち管理テーブルT3に登録し、待ち管理テーブルT3に登録されたグラフにループが形成されているか判断する。その結果、システム2のデッドロック検出部(DD)15は、ループを検出し、デッドロック発生をシステム2のトランザクション管理部(TM)10に通知する。 20

【0123】

[例3 2つのコンピュータシステム間における2つのトランザクションのデッドロック検出中に、メッセージの消失発生]

例3では、2つのコンピュータシステム(システム1, システム2)間での通信エラーにより、メッセージ消失が発生した場合の例である。

(1) 先ず、システム1におけるトランザクショントランザクションT(1, 1)が、システム1のトランザクション管理部(TM)10に対しトランザクションの開始を通知したとする。 30

【0124】

(2) すると、システム1のトランザクション管理部(TM)10は、デッドロック検出部(DD)15にトランザクションT(1, 1)の登録を要求する。

(3) いま、トランザクションT(1, 1)が、システム1のトランザクション管理部(TM)10に、資源Aを要求したとする。

(4) すると、システム1のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Aの獲得を要求する。

(5) すると、システム1の資源管理部(RM)11が、ロック管理部(LM)12に資源Aのロック獲得を要求する。 40

【0125】

(6) 資源Aが未ロックであれば、システム1のロック管理部(LM)12が、資源管理部(RM)11にOKを応答する。

(7) すると、システム1の資源管理部(RM)11が、トランザクション管理部(TM)10にOKを応答する。

(1)'一方、システム2におけるトランザクションT(2, 1)が、システム2のトランザクション管理部(TM)10に対しトランザクションの開始を通知したとする。

【0126】

(2)'すると、システム2のトランザクション管理部(TM)10は、デッドロック検出部(DD)15にトランザクションT(2, 1)の登録を要求する。 50

(3) 'いま、トランザクションT(2,1)が、システム2のトランザクション管理部(TM)10に、資源Bを要求したとする。

(4) 'すると、システム2のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Bの獲得を要求する。

【0127】

(5) 'すると、システム2の資源管理部(RM)11が、ロック管理部(LM)12に資源Bのロック獲得を要求する。

(6) '資源Bが未ロックであれば、システム2のロック管理部(LM)12が、資源管理部(RM)11にOKを応答する。

【0128】

(7) 'すると、システム2の資源管理部(RM)11が、トランザクション管理部(TM)10にOKを応答する。

(8) 以上の状況下において、トランザクションT(1,1)がシステム1のトランザクション管理部(TM)10に資源Bを要求したとする。

(9) 'すると、システム1のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Bの獲得を要求する。

【0129】

(10) 'すると、システム1の資源管理部(RM)11が、ロック管理部(LM)12に資源Bのロック獲得を要求する。

(11) 'ところが、資源Bはシステム2のトランザクションT(2,1)によって既にロック済みであるので、トランザクションT(1,1)がトランザクションT(2,1)に対して待つことになる。そこで、システム1のロック管理部(LM)12が、デッドロック検出部(DD)15にグラフ

T(1,1) T(2,1)

の登録を要求する。

【0130】

(12) 'この要求を受けて、システム1のデッドロック検出部(DD)15は、このグラフをシステム1の待ち管理テーブルT3に登録する。これと同時に、システム1のデッドロック検出部(DD)15は、システム2にグラフ

T(1,1) T(2,1)

を送信する。

(13) 'ただし、その送信内容は、通信エラーにより消失して、システム2に届かなかった。

【0131】

(14) 'この後で、トランザクションT(2,1)がシステム2のトランザクション管理部(TM)10に資源Aを要求したとする。

(15) 'すると、システム2のトランザクション管理部(TM)10が、資源管理部(RM)11に資源Aの獲得を要求する。

(16) 'すると、システム2の資源管理部(RM)11が、ロック管理部(LM)12に資源Aのロック獲得を要求する。

【0132】

(17) 'ところが、資源Aはシステム1のトランザクションT(1,1)によって既にロック済みであるので、トランザクションT(2,1)がトランザクションT(1,1)に対して待つことになる。そこで、システム2のロック管理部(LM)12が、デッドロック検出部(DD)15にグラフ

T(2,1) T(1,1)

の登録を要求する。この時点で、実際にはデッドロック状態が生じている。

【0133】

しかしながら、メッセージ消失によりデッドロック状態は検出できないので、デッドロック状態が持続することになる。

10

20

30

40

50

(18) 一定時間後、システム 1 の待ち時間監視部 (WT) 13 が起動され、システム 1 のデッドロック検出部 (DD) 15 に対してリトライ通知を行う。

(19) すると、システム 1 のデッドロック検出部 (DD) 15 は、トランザクション管理部 (TM) 10 にトランザクション T (1, 1) のリトライを通知する。

【0134】

(20) リトライ通知に従って、システム 1 のトランザクション管理部 (TM) 10 が、資源管理部 (RM) 11 に資源 B の獲得を要求する。

(21) すると、システム 1 の資源管理部 (RM) 11 が、ロック管理部 (LM) 12 に資源 B のロック獲得を再度要求する。

(22) ところが、資源 B はシステム 2 のトランザクション T (2, 1) によって既にロック済みであるので、トランザクション T (1, 1) がトランザクション T (2, 1) に対して待つことになる。そこで、システム 1 のロック管理部 (LM) 12 が、デッドロック検出部 (DD) 15 にグラフ

T (1, 1) T (2, 1)

の登録を再度要求する。

【0135】

(23) この要求を受けて、システム 1 のデッドロック検出部 (DD) 15 は、このグラフをシステム 1 のウェイトフォークラフテーブル T3 に登録する。これと同時に、システム 2 にグラフ

T (1, 1) T (2, 1)

を再度送信する。

(24) システム 2 のデッドロック検出部 (DD) 15 が、このグラフ

T (1, 1) T (2, 1)

を受信し、これをシステム 2 の待ち管理テーブル T3 に登録する。これにより、グラフの欠落が補われる。

(25) システム 2 のデッドロック検出部 (DD) 15 がループを検出し、デッドロック発生をシステム 2 の TM に通知する。

【0136】

【発明の効果】

本発明では、以上説明したように、タスク (トランザクション) による資源のロック状態を管理するロック管理部 (LM) 103 と、デッドロック検出部 (DD) 104 とを分離し、双方が非同期に動作するようにした。そのため、タスク (トランザクション) が新たに発生して資源を要求しても、待ちが発生せずロックが獲得できる場合は、デッドロック検出部 (DD) 104 を介せず動作できる。従って、システムの円滑な運用が図れ、処理の高速化を図れる。また、ロックが獲得できない場合でも、ロック状態 (グラフ) の登録やデッドロックの検出はロックの要求とは非同期に動作するので影響は小さい。

【0137】

特に、デッドロックを少なくするように設計されたシステムで、デッドロック検出の与える影響は極めて小さくなる。

本発明が分散システムに適用された場合、システムは他のシステムと待ち関係に陥った場合にのみデッドロックのための通信を行う。したがって、自システム内のみのデッドロックの検出では通信は発生しない。他システムとの関連があった場合も、デッドロックの 90% 以上が 2 者間で発生することから、ほとんどの場合、1 回の通信でデッドロックは検出される。このため、通信のオーバーヘッドを削減でき、効率のよいシステム運用を図ることができる。

【0138】

また、本発明で、待ち時間監視部 (WT) 13 を設けた場合、分散システムでのメッセージ通信中にメッセージが遅延した場合や消失した場合でも、待ち時間監視部 (WT) 13 が再びデッドロック検出の契機を与えるため、すべてのデッドロックを検出することができる。

10

20

30

40

50

【図面の簡単な説明】

- 【図 1】 本発明の原理図 1
- 【図 2】 本発明の原理図 2
- 【図 3】 デッドロックを示す説明図
- 【図 4】 実施例を示すブロック図
- 【図 5】 ローカルWFGとグローバルWFGの関係を示す図
- 【図 6】 トランザクション管理部の動作を示すフローチャート
- 【図 7】 資源管理部の動作を示すフローチャート
- 【図 8】 ロック管理部の動作を示すフローチャート
- 【図 9】 デッドロック検出部の動作を示すフローチャート
- 【図 10】 待ち時間監視部の動作を示すフローチャート

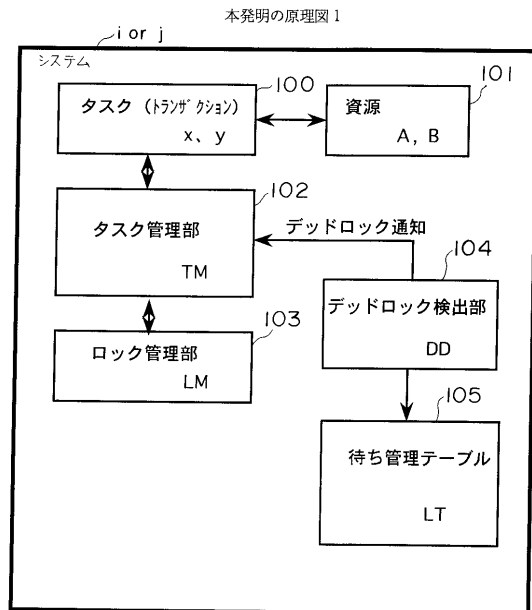
10

【符号の説明】

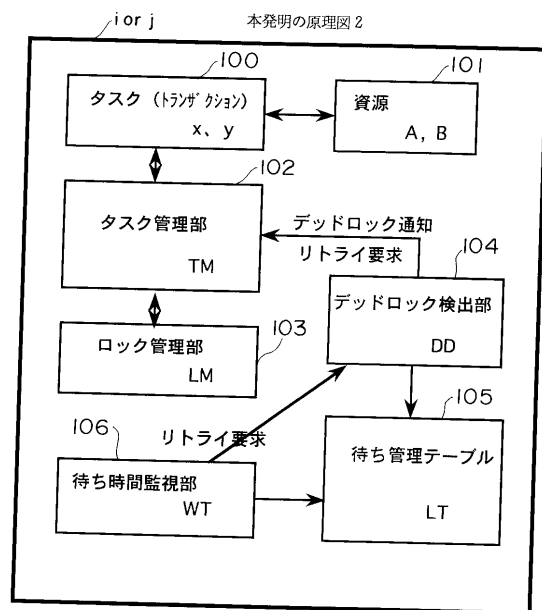
- 1 0 トランザクション管理部
- 1 1 資源管理部
- 1 2 ロック管理部
- 1 3 待ち時間監視部
- 1 4 要求キュー受付部
- 1 5 デッドロック検出部
- 2 0 データベース
- T 3 待ち管理テーブル

20

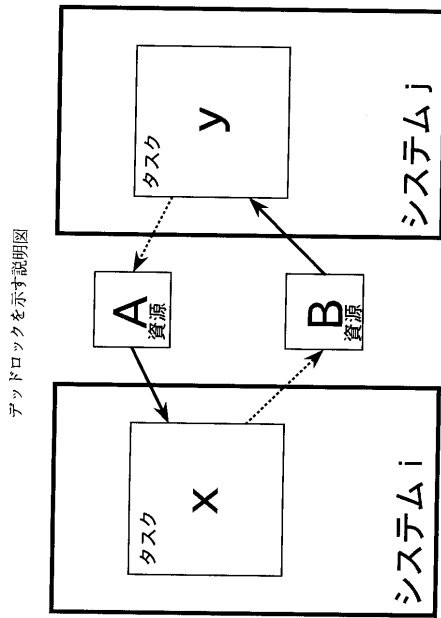
【図 1】



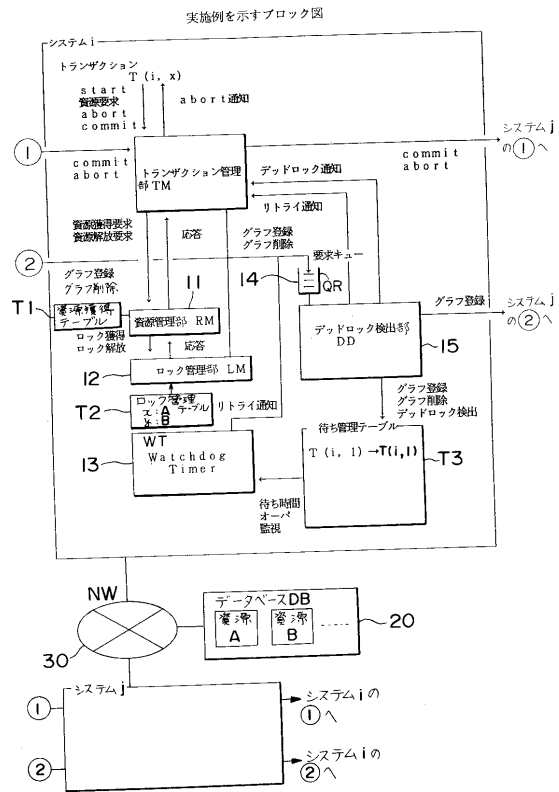
【図 2】



【図3】

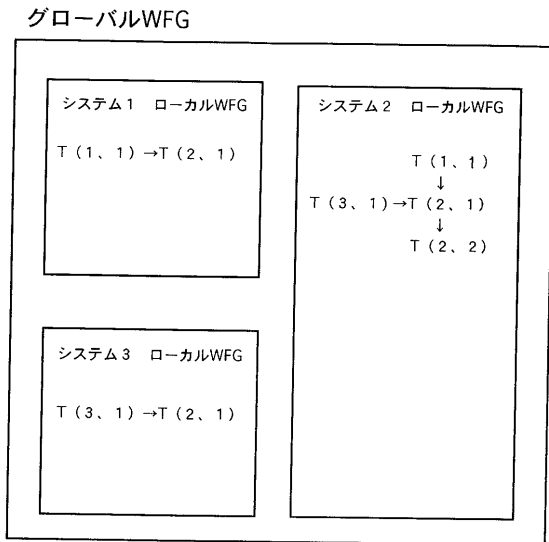


【図4】



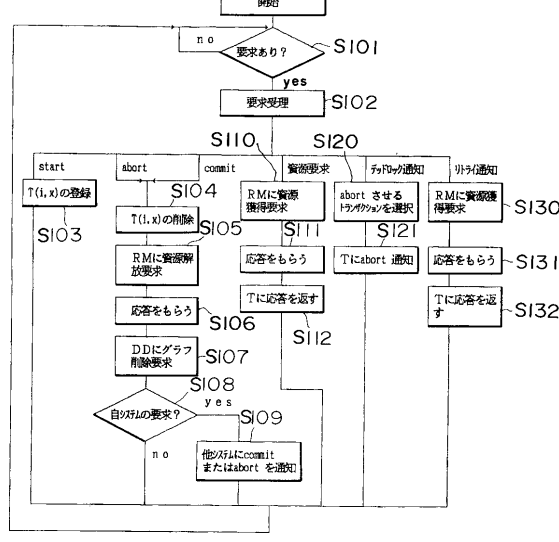
【図5】

ローカルWFGとグローバルWFGの関係を示す図

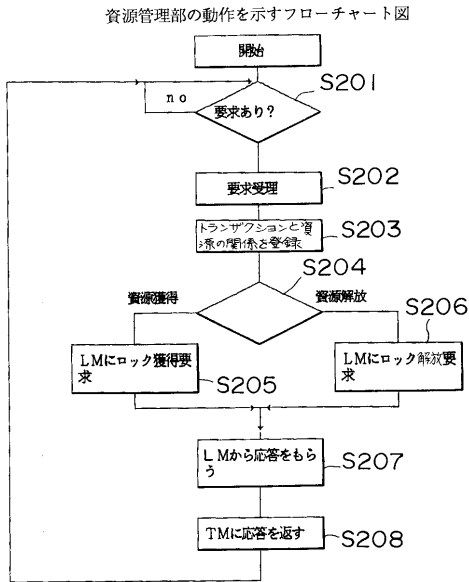


【図6】

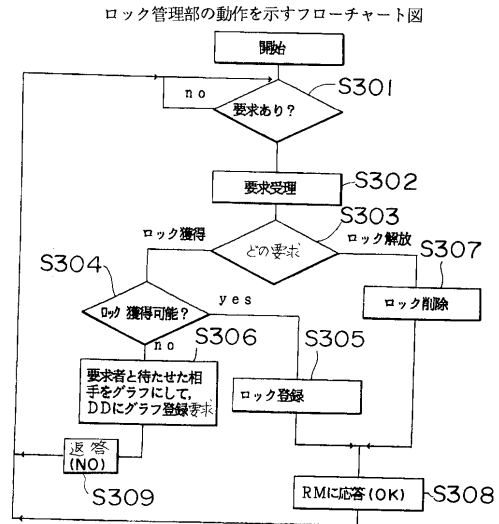
トランザクション管理部の動作を示すフローチャート図



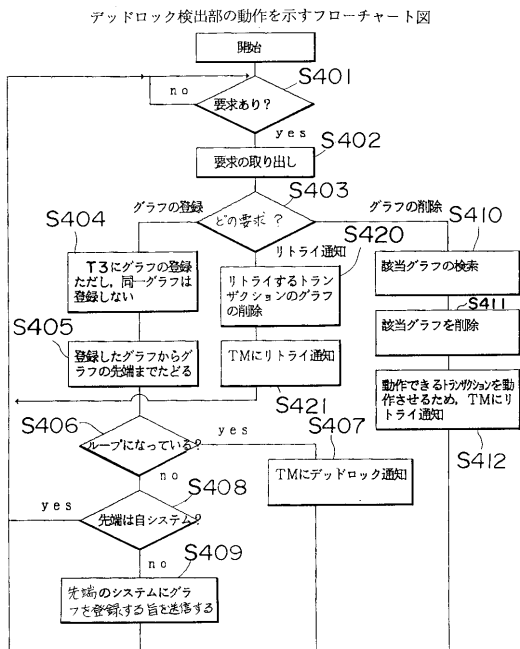
【 図 7 】



【 図 8 】

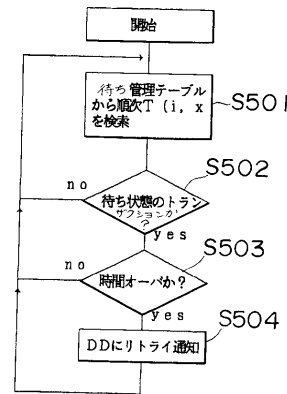


【 図 9 】



【 図 10 】

待ち時間監視部の動作を示すフローチャート図



フロントページの続き

合議体
審判長 吉岡 浩
審判官 松浦 功
審判官 山中 実

(56)参考文献 特開平2 - 2 4 7 7 7 0 (J P , A)

(58)調査した分野(Int.Cl.⁷, D B名)

G06F9/46

G06F15-16-15/177