

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2014-119964

(P2014-119964A)

(43) 公開日 平成26年6月30日(2014.6.30)

(51) Int.Cl.	F I	テーマコード (参考)
G 0 6 F 11/30 (2006.01)	G 0 6 F 11/30 3 0 5 G	5 B 0 4 2
G 0 6 F 9/48 (2006.01)	G 0 6 F 9/46 4 5 2 J	
G 0 6 F 9/52 (2006.01)	G 0 6 F 9/46 4 7 2 A	

審査請求 未請求 請求項の数 5 O L (全 15 頁)

(21) 出願番号	特願2012-274665 (P2012-274665)	(71) 出願人	000005108
(22) 出願日	平成24年12月17日 (2012.12.17)		株式会社日立製作所
			東京都千代田区丸の内一丁目6番6号
		(74) 代理人	110000350
			ポレール特許業務法人
		(72) 発明者	今野 功
			神奈川県横浜市戸塚区戸塚町216番地
			株式会社日立製作所通信ネットワーク事業
			部内
		(72) 発明者	松木 譲介
			神奈川県横浜市戸塚区戸塚町216番地
			株式会社日立製作所通信ネットワーク事業
			部内
		Fターム(参考)	5B042 GA23 JJ06 JJ14 KK05 MC07

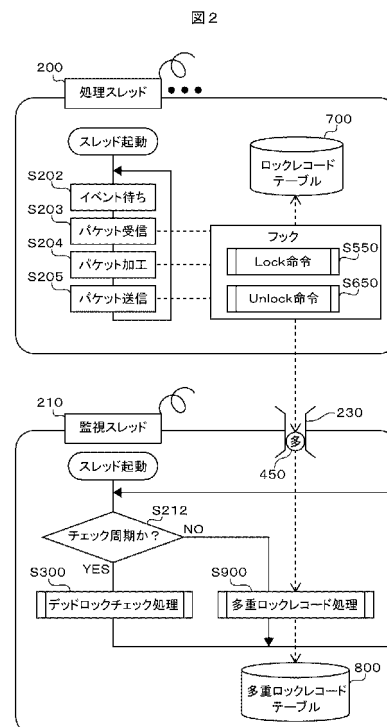
(54) 【発明の名称】 計算機システムおよびプログラム

(57) 【要約】

【課題】マルチスレッドで動作するソフトウェアにおいて、デッドロック解消のためにカーネルの改造せずにデッドロックを自動検出し、フェールオーバーなしでデッドロック状態を解消する方法を提供する。

【解決手段】スレッドの排他制御情報を収集、蓄積し、その情報からデッドロックを検出し、排他制御を追加することによりデッドロックを回避し、デッドロックが発生した場合はデッドロック発生前の状態に戻すことでデッドロックを解消する。

【選択図】図2



【特許請求の範囲】**【請求項 1】**

複数の処理スレッドと、監視スレッドとを並列実行する計算機システムであって、
プロセッサと、メモリとを備え、

前記処理スレッドは、排他制御情報を収集および蓄積し、多重排他制御情報を前記監視スレッドに送信し、

前記監視スレッドは、前記多重排他制御情報を蓄積し、前記多重排他制御情報間にデッドロックを引き起こす組み合わせを検出したとき、当該組み合わせの多重排他制御情報の復帰場所に排他制御を追加することを特徴とする計算機システム。

【請求項 2】

請求項 1 に記載の計算機システムであって、

前記処理スレッドは、前記排他制御情報を、ロック取得命令とロック解放命令へのフックを利用して収集することを特徴とする計算機システム。

【請求項 3】

請求項 1 に記載の計算機システムであって、

前記監視スレッドは、デッドロック要因となるロック取得命令を含む関数の開始時と、ロック解放命令を含む関数の終了時とにフックして、復帰場所を指定することを特徴とする計算機システム。

【請求項 4】

請求項 1 に記載の計算機システムであって、

デッドロック発生前にメモリ内容のスナップショットを取得し、

デッドロック発生後に、前記スナップショットからメモリ内容を復元し、

プログラムカウンタ位置を復帰場所に変更し、

デッドロックの要因となった排他制御を解放することによって、デッドロック状態から復元することを特徴とする計算機システム。

【請求項 5】

コンピュータを、

排他制御情報を収集および蓄積し、多重排他制御情報を前記監視スレッドに送信する処理スレッド、

前記多重排他制御情報を蓄積し、前記多重排他制御情報間にデッドロックを引き起こす組み合わせを検出したとき、当該組み合わせの多重排他制御情報の復帰場所に排他制御を追加する監視スレッド、

として機能させるプログラム。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、計算機システムおよびプログラムに係り、特に、マルチスレッドで動作する計算機システムにおいて、デッドロックを検出し、デッドロック発生前の回避とデッドロック発生後の解消を行う計算機システムおよびプログラムに関する。

【背景技術】**【0002】**

近年、様々なシステムにおいてマルチスレッド環境で稼働するソフトウェアが普及している。マルチスレッド環境では、OS (Operating System) がソフトウェア内の実行単位をスレッドで管理する。OS は、スレッドに対し CPU (Central Processing Unit) が時間割り当てのスケジューリングを行う。スレッドは、OS から CPU 時間を割り当てられた時間だけ動作することができる。

【0003】

したがって、CPU を複数備えるマルチスレッド環境では、スレッドを並列に実行し高速化できる。しかし、処理を並列に実行することにより、デッドロックと呼ばれる問題が発生する。

10

20

30

40

50

【 0 0 0 4 】

ここで、デッドロックとは複数のスレッドまたはプロセスなどの処理単位がリソースを共有する場合において、リソース解放漏れやリソースを占有 (L o c k) 中の処理単位が互いにリソースの解放 (U n l o c k) を待ち続け、処理を再開できずに停止している状態である。ここで、プロセスとはプログラムの実行単位であり、プログラム内で利用している変数や状態を保持し、1つ以上のスレッドから構成される。

【 0 0 0 5 】

デッドロックが発生したプロセスは、異常終了しない。このため、ソフトウェアの使用者、開発者が異常の発生直後に気付かない。また、ソフトウェアの開発者が異常に気付いたとしても、スレッドの挙動確認や解析に時間がかかる。そこで、デッドロックを防ぎ、解決する技術が重要となってくる。

10

【 0 0 0 6 】

特許文献 1 には、デッドロックが二重以上のロック取得するスレッド同士で発生する点に着目し、大域的なロックを追加している。つまり、二重以上のロックを取得するスレッドに対し、大域的なロックを取得する制約を設ける。これによって、二重以上のロックを取得するスレッドを制限し、デッドロックを防止する。さらに、大域的なロックを一律に取得するのではなく、条件を付けることで性能低下を抑えている。

【 0 0 0 7 】

また、特許文献 2 では、ロック取得および解除に対して A P I 関数フックを行い、スレッド毎に共有リソースに対して「ロック中」および「ロック取得中」の情報を取得し、その組み合わせによってデッドロックを検出する計算機システムおよびプログラムが提案されている。ここで、A P I 関数とは、O S やミドルウェアが提供するアプリケーションおよびソフトウェア開発向けのインタフェースであり、共通ライブラリの形で提供される。また、A P I 関数フックとは、A P I 関数の処理を横取りし、利用者が独自に定義した処理を行うことである。また、「ロック中」とは、スレッドがリソースを占有中により他のスレッドがリソースにアクセスできない状態である。一方、「ロック取得中」とは他のスレッドが占有中のリソースが解放されるのを待ってからロックしようとしている状態である。

20

【 0 0 0 8 】

特許文献 2 は、これらの情報を組み合わせ、具体的には、リソース A を「ロック中」でリソース B を「ロック取得中」のスレッド 1 と、リソース A を「ロック取得中」でリソース B を「ロック中」のスレッド 2 が同時に存在すると、互いにリソースが解放されるのを待ち続けるデッドロックとして検出できる。デッドロック検出後は、発生経路を記憶して冗長系にフェールオーバーすることで、フェールオーバー先でデッドロックの再発を防止する。ここで、フェールオーバーとは、正常時に冗長構成の計算機でデータの同期を取っておき、障害が発生した計算機から残りの正常な計算機にデータを引継ぎ、処理を継続する動作である。

30

【 先行技術文献 】

【 特許文献 】

【 0 0 0 9 】

40

【 特許文献 1 】 特開平 0 7 - 1 9 1 9 4 4 号公報

【 特許文献 2 】 特開 2 0 0 9 - 2 7 1 8 5 8 号公報

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 1 0 】

特許文献 1 に記載の技術は、ロックの取得を連続で行う必要があり、多重ロックを取得する時点で予め必要なロックが全て既知であることが前提条件となる。実際のプログラムでは、経路によって必要なロックが異なるため、予め必要なロックを全て取得することは困難であり、必要なロックを全て取得できなかった場合、デッドロックが発生する。

【 0 0 1 1 】

50

特許文献 2 に記載の技術は、デッドロックの再発を防げる。しかし、別経路でデッドロックが発生する度にフェールオーバを行う必要がある。フェールオーバ実行中は、本来実行しなければならない処理を実行できず、本来提供したいサービスに影響が出てしまう。また、切替え可能な計算機が無ければフェールオーバできないため、可能な限りフェールオーバ実施を避けるべきである。

特許文献 1 と特許文献 2 を組み合わせたとしても、デッドロック発生後の復帰手段がないため、容易には解決できない。

【課題を解決するための手段】

【0012】

本発明の代表的な一形態によると、複数のスレッドを実行する計算機システムであって、計算機システムは、少なくとも一つのプロセッサと、メモリとを備え、計算機システムで動作するソフトウェア内で第一のスレッドと第一のスレッドを監視するための監視スレッドを実行し、第一のスレッドの排他制御の情報を保持するための監視情報領域と、第一のスレッド用にデッドロック回避用の領域および情報と、メモリ内容を復元するための退避領域と、をメモリに保持し、第一のスレッドと監視スレッドは監視情報領域に排他制御の情報を格納し、第一のスレッドはメモリ内容を退避領域に格納し、監視スレッドが監視情報領域からデッドロックを検出し、監視スレッドはデッドロックが発生する可能性がある第一のスレッドに対し、デッドロック回避用領域に排他制御を追加し、第一のスレッドでデッドロックが発生した場合、第一のスレッドの処理を退避位置に戻し、かつ、退避領域からメモリ内容を復元することでデッドロックを解消する。

10

20

【0013】

上述した課題は、複数の処理スレッドと、監視スレッドとを並列実行する計算機システムであって、プロセッサと、メモリとを備え、前記処理スレッドは、排他制御情報を収集および蓄積し、多重排他制御情報を前記監視スレッドに送信し、前記監視スレッドは、前記多重排他制御情報を蓄積し、前記多重排他制御情報間にデッドロックを引き起こす組み合わせを検出したとき、当該組み合わせの多重排他制御情報の復帰場所に排他制御を追加する計算機システムにより、達成できる。

【0014】

また、コンピュータを、排他制御情報を収集および蓄積し、多重排他制御情報を前記監視スレッドに送信する処理スレッド、前記多重排他制御情報を蓄積し、前記多重排他制御情報間にデッドロックを引き起こす組み合わせを検出したとき、当該組み合わせの多重排他制御情報の復帰場所に排他制御を追加する監視スレッド、として機能させるプログラムにより、達成できる。

30

【発明の効果】

【0015】

本発明によって、デッドロックを自動検出し、フェールオーバを実施せずにデッドロック状態から正常な状態に復帰できる。

【図面の簡単な説明】

【0016】

【図 1】通信装置のハードウェアブロック図である。

40

【図 2】デッドロック監視を行う装置のスレッド構成と処理内容を説明する図である。

【図 3】ロックレコード、多重ロックレコードを説明する図である。

【図 4 A】Lock 命令のフローチャートである。

【図 4 B】フック有り Lock 命令のフローチャートである。

【図 5 A】Unlock 命令のフローチャートである。

【図 5 B】フック有り Unlock 命令のフローチャートである。

【図 6】ロックレコードテーブルを説明する図である。

【図 7】多重ロックレコードテーブルを説明する図である。

【図 8】多重ロックレコードの処理フローチャートである。

【図 9】デッドロック検出、回避、解消の流れを示すフローチャートである。

50

【図 10】デッドロック検出のフローチャートである。

【図 11】デッドロック回避のフローチャートである。

【図 12】デッドロック解消のフローチャートである。

【図 13】メモリ内容をデッドロック発生前の状態に復元する処理を説明する図である。

【発明を実施するための形態】

【0017】

以下、本発明の実施の形態について、実施例を用い図面を参照しながら詳細に説明する。なお、実質同一部位には同じ参照番号を振り、説明は繰り返さない。

【実施例 1】

【0018】

実施例 1 は、通信装置のように特定の処理を繰り返す計算機を説明する。

図 1 を参照して、最小構成のハードウェアの通信装置を説明する。図 1 において、通信装置 100 は、プロセッサ 101 と、メモリ 102 と、N I F (Network InterFace) 103 とを具備する。

【0019】

プロセッサ 101 は、C P U 等の演算装置であり、O S やアプリケーションプログラム等のソフトウェアを実行する。メモリ 102 は、主記憶装置であり、プロセッサ 101 がソフトウェア実行時に、プログラム実行バイナリおよびプログラムが使用するデータを格納する。N I F 103 は、通信装置 100 とは別の装置とパケットを送受信するためのインタフェースである。各ハードウェアブロックは、バス 110 - 1、110 - 2 によって相互に接続するため、互いに命令メッセージ及びデータを送信することが可能である。

【0020】

なお、本実施例を適用するためには、プロセッサ 101 がマルチタスクに対応している必要がある。ここで、マルチタスクとは演算装置が複数の処理を切換えながら複数の処理を実行する方式であり、汎用計算機のプロセッサでは、そのほとんどがマルチタスクに対応している。

【0021】

また、通信装置 100 は、物理的に一つの計算機によって実装されてもよいし、少なくとも一つの計算機が提供する仮想的な計算機によって実装されてもよい。

【0022】

図 2 を参照して、ソフトウェアのスレッド構成と、ロックを記録するためのフック処理を説明する。図 2 のサブルーチンに関しては、別途、後述の図を用いて詳述する。図 2 において、スレッドは、通信装置の呼制御などを行う複数の処理スレッド 200 と、処理スレッド 200 が正常に動作しているか否かを判定する監視スレッド 210 から成る。

【0023】

処理スレッド 200 は、起動後、パケット受信などのイベントを待つ (S 202)。処理スレッド 200 は、パケットを受信し (S 203)、送信するパケットを加工し (S 204)、送信する (S 205) 動作を繰り返す。スレッド起動からパケット送信までは、ソフトウェアバイナリの変更は行わない。

【0024】

処理スレッド 200 は、フック処理 220 とロックレコードテーブル 700 を用意しておく。フック処理 220 は、パケット処理のステップ 203 とパケット加工のステップ 204 とパケット送信のステップ 205 で呼び出されるすべての L o c k 命令 500 および U n l o c k 命令 600 を横取りするフック関数の集合である。フック処理 220 では、L o c k 命令 500 が呼び出された際に、別途定義するサブルーチン L o c k 命令 550 を呼び出す。また、U n l o c k 命令 600 が呼び出された際に、別途定義するサブルーチン U n l o c k 命令 650 を呼び出す。なお、サブルーチン L o c k 命令 550 は、図 4 で後述し、サブルーチン U n l o c k 命令 650 は、図 5 で後述する。フック処理 220 は、一時的にロックレコードをロックレコードテーブル 700 に保存する。このうち、多重ロックレコードと判断した場合は、多重ロックレコード 450 を、監視スレッド 21

10

20

30

40

50

0 のキュー 2 3 0 に追加する。

【 0 0 2 5 】

監視スレッド 2 1 0 は、プログラムカウンタが、サブルーチン多重ロックレコード処理 9 0 0 に到達した際に、キュー 2 3 0 に多重ロックレコード 4 5 0 が入っているか調べ、多重ロックレコード 4 5 0 がキュー 2 3 0 に入っていた場合は、多重ロックレコードテーブル 8 0 0 に、多重ロックレコード 4 5 0 を追加する。なお、サブルーチン多重ロックレコード処理 9 0 0 は、図 8 を用いて詳述する。

【 0 0 2 6 】

上述した手順により、スレッド起動とイベント待ちのステップ 2 0 2 とパケット受信のステップ 2 0 3 とパケット加工のステップ 2 0 4 とパケット送信のステップ 2 0 5 を、変更することなく、ロック情報をロックレコードテーブル 7 0 0 に記録することができる。

10

【 0 0 2 7 】

また、監視スレッド 2 1 0 は、起動後、チェック周期であるか否かを判定し、チェック周期であればデッドロックチェック処理 (S 3 0 0) を行い、チェック周期でなければ多重ロックレコード処理 (S 9 0 0) を行う動作を繰返す。ここで、チェック周期は開発者や使用者が指定する時間や回数以外にも、排他制御の統計から自動的に算出しても良い。デッドロックチェック処理は、図 9 を用いて、後述する。

【 0 0 2 8 】

本実施例では、監視スレッド 2 1 0 がデッドロックチェック処理を繰り返し実行する。これにより、デッドロックを検出し、デッドロック発生前に検出した場合は、デッドロックを回避し、デッドロック発生と同時にデッドロックを検出した場合はデッドロック状態からの解消処理を行い、デッドロックを自動的に解決する。なお、本実施例では監視スレッドを監視専用としているが、監視スレッドに別の役割を持たせることや、処理スレッド側に監視動作を追加することも可能である。このことにより、スレッド数を増やす必要が無く、メモリ量や C P U リソースを節約する利点がある。

20

【 0 0 2 9 】

図 3 を参照して、デッドロックチェック動作の検出処理に利用するロックレコード 4 0 0 および多重ロックレコード 4 5 0 を説明する。

【 0 0 3 0 】

図 3 (a) において、ロックレコード 4 0 0 は、ネスト数 4 0 1、ロックハンドル 4 0 2、復帰場所 4 0 3 から成る。ロックレコード 4 0 0 は、処理スレッド 2 0 0 が個々に管理する排他制御の記録である。ここで、ネスト数 4 0 1 は、排他制御の深さである。ロックハンドル 4 0 2 は、排他制御に必要な識別子である。復帰場所 4 0 3 は、デッドロック回避用領域である。

30

【 0 0 3 1 】

図 3 (b) において、多重ロックレコード 4 5 0 は、項番 4 5 1、ロック順序 4 5 2、復帰場所 4 0 3、追加ロックハンドル 4 5 3 から成る。多重ロックレコード 4 5 0 は、監視スレッド 2 1 0 が管理する排他制御組合せパターンの記録である。ここで、項番 4 5 1 は、多重ロックレコードを管理するための識別子である。ロック順序 4 5 2 は、ロックハンドル 4 0 2 を取得した順序である。追加ロックハンドル 4 5 3 は、デッドロック回避用のロックハンドルである。

40

【 0 0 3 2 】

ロックレコード 4 0 0 は、図 6 に示すロックレコードテーブル 7 0 0 で管理する。多重ロックレコード 4 5 0 は、図 7 に示す多重ロックレコードテーブル 8 0 0 で管理する。ロックレコードテーブル 7 0 0 は、処理スレッド 2 0 0 が個々に管理するテーブルである。ロックレコードテーブル 7 0 0 は、処理スレッド 2 0 0 がその時点で取得している排他制御を示す。ロックレコードテーブル 7 0 0 は、ロック取得命令 (以下、L o c k 命令、本実施例では、p t h r e a d _ m u t e x _ l o c k) とロック解放命令 (以下、U n l o c k 命令、本実施例では、p t h r e a d _ m u t e x _ u n l o c k) にフックを行って取得する。ここで、フックとは処理を横取りし、使用者が定義した処理をプロセッサ

50

101に行わせる処理である。

【0033】

具体的には、Linux（登録商標）系OSのLD__PRELOAD環境変数と共有ライブラリを用いたAPI関数フックで実現できる。LD__PRELOAD環境変数を指定しない状態では、pthread__mutex__lock関数およびpthread__mutex__unlock関数は、libpthread.soという共有ライブラリの関数が実行される。しかし、LD__PRELOAD環境変数で、共有ライブラリファイル名を指定すると、関数名から実行する関数へのアドレス解決順序を入れ替え、LD__PRELOAD環境変数で指定した、共有ライブラリを優先的に検索する。指定する共有ライブラリに、pthread__mutex__lock関数とpthread__mutex__unlock関数を定義することにより、libpthread.soのpthread__mutex__lock関数およびpthread__mutex__unlock関数を実行することができる。共有ライブラリのpthread__mutex__lock関数およびpthread__mutex__unlock関数にて、libpthread.soのpthread__mutex__lock関数およびpthread__mutex__unlock関数を呼び出すことにより、Lock命令へのフックおよびUnlock命令へのフックを実現することができる。

10

【0034】

図4を参照して、Lock命令のフローチャートとフック有りLock命令のフローチャートを説明する。また、図5を参照して、Unlock命令のフローチャートとフック有りUnlock命令のフローチャートを説明する。

20

【0035】

図4Aにおいて、Lock命令を開始すると、処理スレッド200は、ロックハンドルが取得可能か判定する(S501)。取得可能であれば、処理スレッド200は、ロックハンドルを取得し(S504)、処理を終了する。ステップ501でロックハンドルが取得不可であれば、ロックハンドルが解放されるまで待ち(S503)、ロックハンドルを取得する(S504)。

【0036】

図4Bにおいて、フック有りLock命令を開始すると、処理スレッド200は、ロックレコードの記録追加処理を実施する(S551)。ロックレコードの記録追加処理は、取得しようとしているロックハンドルをロックレコードに登録し、ロックレコードテーブル700に追加する処理である。追加後、処理スレッド200は、ロックレコードテーブル700のネスト数401が2以上か判定する(S552)。判定結果が2以上の場合、処理スレッド200は、ロックレコードテーブル700に登録中の全ロックレコード400から多重ロックレコード450を作成し、監視スレッド210のキュー230にエンキューする(S553)。エンキュー完了後およびネスト数401が2未満の場合、処理スレッド200は、Lock命令を実行し(S500)、フック有りLock関数を完了する。

30

【0037】

図5Aにおいて、Unlock命令を開始すると、処理スレッド200は、ロックハンドルを解放し(S601)、Unlock命令を終了する。

40

図5Bにおいて、フック有りUnlock命令650を開始すると、処理スレッド200は、Unlock命令を実行する(S600)。処理スレッド200は、処理完了後にロックレコードの記録削除処理を行い(S651)、フック有りUnlock命令を完了する。

【0038】

フック有りLock命令のロックレコードの記録追加処理のステップ551と、フック有りUnlock命令のロックレコードの記録削除処理のステップ651によって、処理スレッド200は、取得しているロックハンドルを最新の状態として反映することができる。

50

【 0 0 3 9 】

図 6 を参照して、ロックハンドルテーブル 7 0 0 の状態遷移を説明する。図 6 において、フック 2 2 0 によって、ロックハンドル A を取得した状態 7 0 1 から、フック有り L o c k (B) でロックハンドル B を取得すると状態 7 0 2 になり、さらにフック有り U n l o c k (B) でロックハンドル B を解放すると状態 7 0 3 になる。

【 0 0 4 0 】

状態 7 0 2 ではネスト数 4 0 1 が 2 以上になるため、フック有り L o c k 命令において、処理スレッド 2 0 0 は、多重ロックレコード 4 5 0 を作成する。ロックレコードテーブル 7 0 0 が状態 7 0 2 であれば、作成する多重ロックレコード 4 5 0 のロック順序 4 5 2 は A B、復帰場所 4 0 3 は F u n c A B となり、作成した多重ロックレコード 4 5 0 を監視スレッド 2 1 0 のキュー 2 3 0 にエンキューする。なお、多重ロックレコード 4 5 0 の項番 4 5 1 と追加ロックハンドル 4 5 3 は、監視スレッド 2 1 0 が後で登録する。なお、デッドロックは、二重以上のロックハンドル 4 0 2 を取得した処理スレッド 2 0 0 同士で発生するため、ネスト数 4 0 1 が 2 以上の場合のみ、多重ロックレコード 4 5 0 を作成すれば良い。

【 0 0 4 1 】

処理スレッド 2 0 0 が多重ロックレコード 4 5 0 をキュー 2 3 0 にエンキュー後、監視スレッド 2 1 0 は、多重ロックレコード処理 9 0 0 で多重ロックレコードテーブル 8 0 0 へ登録する。多重ロックレコードテーブル 8 0 0 は、図 7 に示す多重ロックレコード 4 5 0 を蓄積するテーブルであり、多重ロックレコード 4 5 0 を持たない状態 8 0 1 から始まる。なお、図 7 は、フローチャートの中で説明する。

【 0 0 4 2 】

図 8 を参照して、多重ロックレコード処理によって、未登録の多重ロックレコードを追加する処理を説明する。図 8 において、多重ロックレコード処理を開始すると、監視スレッド 2 1 0 は、キュー 2 3 0 に多重ロックレコードがあるか判定する (S 9 0 1)。多重ロックレコードがある場合、監視スレッド 2 1 0 は、多重ロックレコードテーブル 8 0 0 に未登録か判定する (S 9 0 2)。未登録であれば、監視スレッド 2 1 0 は、多重ロックレコードテーブルに登録して (S 9 0 3)、多重ロックレコード処理を終了する。なお、多重ロックレコードが未登録か否かの判定は、ロック順序および復帰場所が一致しているか否かで判定する。多重ロックレコードがキュー 2 3 0 に無い場合、または多重ロックレコードが登録済みの場合、監視スレッド 2 1 0 は、多重ロックレコード処理を終了する。この操作によって、図 7 において、多重ロックレコード 4 5 0 を保持していない状態 8 0 1 から、未登録の多重ロックレコード (A B) を受取り、受取った順に項番 4 5 1 を登録することで、状態 8 0 2 になる。

【 0 0 4 3 】

図 2 に示したように、監視スレッド 2 1 0 は、多重ロックレコードテーブル 8 0 0 に対し、チェック周期以外では多重ロックレコード処理 9 0 0 を行い、チェック周期であれば図 9 に示すデッドロックチェック処理 3 0 0 を行う。

【 0 0 4 4 】

図 9 において、デッドロックチェック処理を開始すると、監視スレッド 2 1 0 は、デッドロック検出処理を実行する (S 3 1 0)。監視スレッド 2 1 0 は、デッドロックを検出したか判定する (S 3 2 0)。Y E S のとき、監視スレッド 2 1 0 は、デッドロック回避処理を実行する (S 3 3 0)。監視スレッド 2 1 0 は、検出したデッドロックが発生しているか判定する (S 3 4 0)。Y E S のとき、監視スレッド 2 1 0 は、デッドロック解消処理を実行して (S 3 5 0)、終了する。ステップ 3 2 0 またはステップ 3 4 0 で N O のとき、監視スレッド 2 1 0 は、デッドロック検出処理を終了する。

【 0 0 4 5 】

図 1 0 を参照して、デッドロック検出処理を説明する。図 1 0 において、デッドロック検出処理を開始すると、監視スレッド 2 1 0 は、比較未実施の多重ロックレコードがあるか判定する (S 3 1 1)。ここで、比較未実施の多重ロックレコードとは、多重ロックレ

10

20

30

40

50

コードテーブル 800 に含まれる多重ロックレコードの組合せのうち、デッドロック検出の為の比較を実施していない多重ロックレコードの集合である。デッドロック検出処理の開始時点では全多重ロックレコードが比較未実施となる。未比較の多重ロックレコードがある場合、監視スレッド 210 は、その中から候補 X と Y を選択する (S312)。監視スレッド 210 は、候補 X と Y にデッドロックの可能性があるか判定する (S313)。デッドロックの可能性がある場合、監視スレッド 210 は、デッドロックパターンとして登録し (S314)、ステップ 311 へ戻る。全ての比較が完了後 (S311: NO)、監視スレッド 210 は、終了する。

【0046】

ステップ 313 でデッドロックの可能性あるか否かを判定する単純な方法としては、候補 X と Y のロック順序 452 から一致するロックハンドルを抽出し、その順序が逆転しているか否かを判定する方法が挙げられる。

【0047】

図 7 に戻って、多重ロックレコードテーブルが状態 803 であれば、ロックハンドル A のみが一致し、順序逆転はないため、デッドロックは発生しないと判定する。また、多重ロックレコードテーブルが状態 804 であれば、3 種類の多重ロックレコードの全組合せ 3 通りで比較し、項番 1 と項番 3 の多重ロックレコードでロックハンドル A と B が一致し、かつ順序逆転が発生しているため、デッドロックとして検出できる。

【0048】

なお、多重ロックレコードテーブル 800 は、多重ロックレコードを蓄積していくため、過去の多重ロックレコードから発生確率が低いデッドロックでも検出することが可能になる。しかし、多重ロックレコードを蓄積し続けると記憶領域を圧迫するため、必要に応じて古い情報を削除しても良い。また、未比較の多重ロックレコード 450 は、多重ロックレコードテーブル 800 に新規登録する際に発生するため、多重ロックレコード処理 900 のステップ 903 後にデッドロック検出処理を実行しても良い。

【0049】

図 11 を参照して、デッドロック回避処理を説明する。図 11 において、監視スレッド 210 は、デッドロックを検出した多重ロックレコードの組合せ (最低でも 2 つ) を受取り、その多重ロックレコードの中で追加ロックハンドル 453 が登録されているかを確認する (S331)。追加ロックハンドル 453 が全ての多重ロックレコードに登録されていない場合 (NO)、監視スレッド 210 は、追加ロックハンドル 453 を新規確保 (Global AB とする) し、多重ロックレコードの追加ロックハンドル 453 に Global AB を登録して (S332)、デッドロック回避処理を終了する。

【0050】

ステップ 331 で追加ロックハンドル 453 を登録済みの場合 (YES)、監視スレッド 210 は、登録している追加ロックハンドル 453 が 1 種類かを確認する (S333)。登録済みの追加ロックハンドル 453 が 1 種類 (Global XX とする) の場合 (YES)、監視スレッド 210 は、未登録の追加ロックハンドル 453 に Global XX を登録して (S334)、デッドロック回避処理を終了する。

【0051】

ステップ 333 で登録されている追加ロックハンドル 453 が複数 (Global XX、Global YY とする) の場合 (NO)、監視スレッド 210 は、ロックハンドルを新規確保 (Global XY とする) し、未登録の追加ロックハンドル 453 と、受取った多重ロックレコードの Global XX と Global YY を Global XY で登録または上書き登録する (S335)。さらに、監視スレッド 210 は、多重ロックテーブル 800 の中から、追加ロックハンドル 453 に Global XX または Global YY を登録しているレコードを Global XY で上書き登録する (S336)。監視スレッド 210 は、上書き登録によって不要となった Global XX と Global YY のメモリ領域を削除する (S337)。追加ロックハンドル 453 の登録が完了後、監視スレッド 210 は、デッドロック回避処理を終了する。

10

20

30

40

50

なお、追加ロックハンドル 4 5 3 が複数となるのは、ネスト数が 3 以上の場合である。

【 0 0 5 2 】

デッドロック回避処理により、具体的には、図 7 において、多重ロックレコードテーブル 8 0 0 が状態 8 0 4 であれば、項番 1 と項番 3 の多重ロックレコードを検出し、両方とも追加ロックハンドル 4 5 3 を登録していないため、追加ロックハンドル (G l o b a l A B) を新規確保し、登録することで状態 8 0 5 となる。

【 0 0 5 3 】

デッドロック回避処理で登録した追加ロックハンドル 4 5 3 は、多重ロックレコードに登録している復帰場所にてスコープドロックし、デッドロック発生の可能性がある箇所を追加した排他制御で覆い込む。状態 8 0 4 では (A) (B) と (B) (A) の順番でロックハンドルを取得していたが、状態 8 0 5 で追加ロックハンドルにより (G l o b a l A B) (A) (B) と (G l o b a l A B) (B) (A) という順番になるため、G l o b a l A B によってデッドロックを防止できる。なお、スコープドロックとは、L o c k 命令を開始後、関数などの有効範囲 (スコープ) 終了時に自動で U n l o c k 命令を実行するための仕組みであり、リソース解放漏れを防ぐことが可能になる。

【 0 0 5 4 】

また、追加ロックハンドル 4 5 3 を上書き登録することにより、過去に検出したデッドロックパターン (具体的には、A B C、D C B) と新たに検出したデッドロックパターン (E C B) で同一の追加ロックハンドル 4 5 3 を指定するため、デッドロックを全て防止できる。また、追加ロックハンドル 4 5 3 を一種類だけとし、全てのデッド

【 0 0 5 5 】

ロックレコードテーブル 7 0 0 や多重ロックレコードテーブル 8 0 0 の復帰場所 4 0 3 は、デッドロックを防止できる箇所 (プログラム上のアドレス) を指定する必要がある。その指定する方法は、チェックポイント方式と関数フック方式がある。チェックポイント方式は、復帰場所をソフトウェアの開発者がチェックポイントとして明示的に指定する方式である。一方、関数フック方式は、C 言語や C + + 言語で記述したソフトウェアのソースコードビルド時に「 - f i n s t r u m e n t - f u n c t i o n s 」オプションを指定することにより、関数の開始時と終了時にフックできるため、デッドロック要因となる L o c k 命令を含む関数の開始箇所を、自動的に指定する方式である。

【 0 0 5 6 】

デッドロック検出と同時にデッドロックが発生している場合は回避できないため、デッドロック解消処理を実施する。なお、デッドロックが発生しているか判定するには、チェック周期の際に処理スレッド 2 0 0 に応答があるか確認し、応答がない、かつ処理スレッド 2 0 0 のロックレコードテーブル 7 0 0 がデッドロックの可能性を検出した多重ロックレコードを含む場合に、デッドロック状態と判定する。

【 0 0 5 7 】

図 1 2 を参照して、デッドロック解消処理を説明する。デッドロック解消処理を開始すると、監視スレッド 2 1 0 は、プログラムカウンタに復帰場所を設定する (S 3 5 1) 。監視スレッド 2 1 0 は、メモリ内容をデッドロック発生前の状態に復元する (S 3 5 2) 。監視スレッド 2 1 0 は、デッドロック要因のロックハンドルを解放して (S 3 5 3) 、デッドロック解消処理を終了する。

【 0 0 5 8 】

プログラムカウンタは、プロセッサ 1 0 1 がプログラム上で次に実行する命令のアドレスを指す。各スレッドは、各々プログラムカウンタを保持しているため、プログラムカウンタを書き換えることにより、任意のスレッドが次に実行する命令のアドレスを任意のアドレスに変更することができる。したがって、プログラムカウンタを多重ロックレコードテーブル 8 0 0 に登録した復帰場所 4 0 3 を指定する。処理スレッド 2 0 0 は、追加ロックハンドル (G l o b a l A B) 取得から処理を再開し、デッドロック再発を防止する。

【 0 0 5 9 】

次に、デッドロック発生時点で既にメモリ内容が書き換わり、メモリに更新途中のデータが保持されている可能性がある。このため、メモリ内容をデッドロック発生前の状態に戻すことで、データの整合性を維持する。

【0060】

図13を参照して、メモリ復元を説明する。メモリ復元では、まず、メモリ102内でデッドロック発生前メモリ内容1300のスナップショットを取り、メモリ内容をコピーする(状態1301)。ここで、スナップショットの際に名前付きメモリマップを利用することで、デッドロック発生後に必要なメモリを検索することができる。その後、デッドロック発生を検出した時点で、メモリ内容が状態1302のように値が書き換わっていたとしても、スナップショット1301をメモリ状態1302に上書きすることで、メモリ内容を復元できる。なお、スナップショットは、フック有りLock命令内で実行すれば、デッドロック発生前のメモリ内容を保持することができる。

10

【0061】

図12に戻って、ステップ353において、処理スレッド200の取得済みロックハンドルの内、デッドロック要因となっているロックハンドル(A)および(B)と、その間に取得したロックハンドルを解放することで、Unlock待ちで停止している処理スレッド200が復帰場所FuncABから処理を再開することが可能になる。

【0062】

以上により、デッドロック検出、回避、解消が可能となった。

本実施例によれば、多重ロックを取得する時点で予め必要なロックが全て既知である前提条件が不要となり、どのようなソフトウェアにも汎用的に適用可能となる。また、デッドロックが発生しても、フェールオーバを実施することなく処理を再開できるため、ソフトウェアの信頼性および稼働率を向上することができる。

20

【符号の説明】

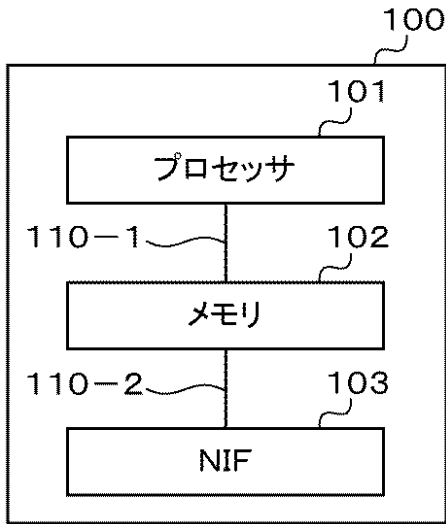
【0063】

100...通信装置、101...プロセッサ、102...メモリ、103...NIF、110...バス、200...処理スレッド、210...監視スレッド、220...フック、230...キュー、400...ロックレコード、450...多重ロックレコード、500...Lock命令、550...フック有りLock命令、600...Unlock命令、650...フック有りUnlock命令、700...ロックレコードテーブル、800...多重ロックレコードテーブル、1301...スナップショット

30

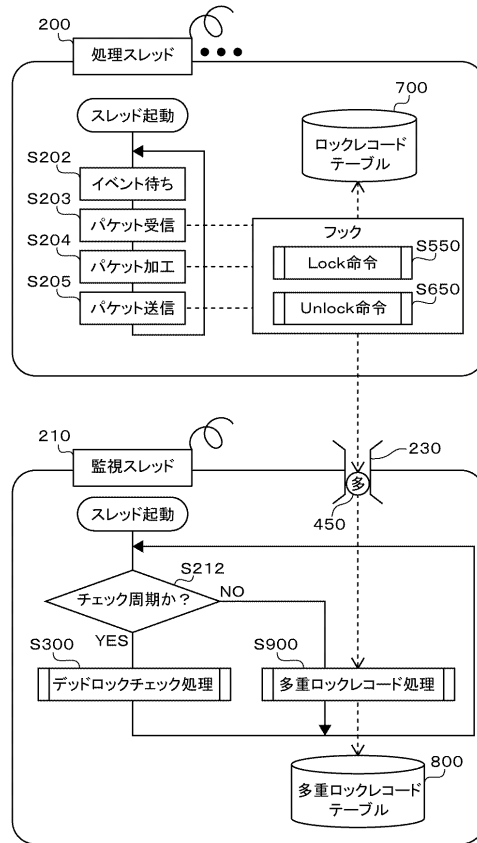
【図 1】

図 1



【図 2】

図 2



【図 3】

図 3

(a)

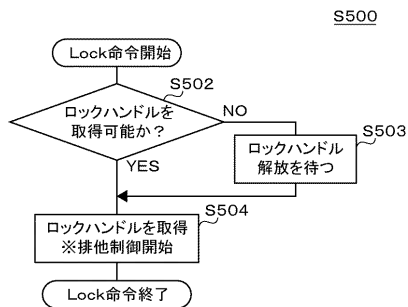
401	402	403
ネスト数	ロックハンドル	復帰場所
—	—	—

(b)

451	452	403	453
項番	ロック順序	復帰場所	追加ロックハンドル
—	—	—	—

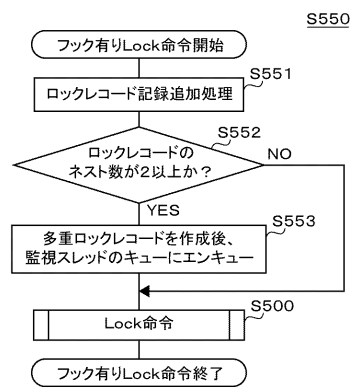
【図 4 A】

図 4 A



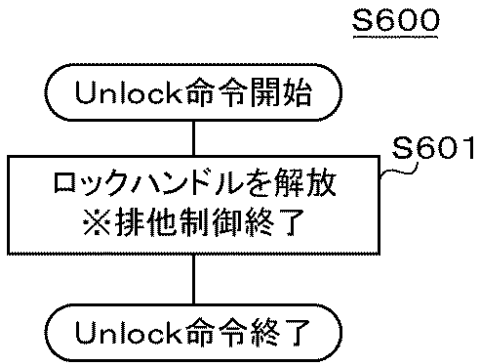
【図 4 B】

図 4 B



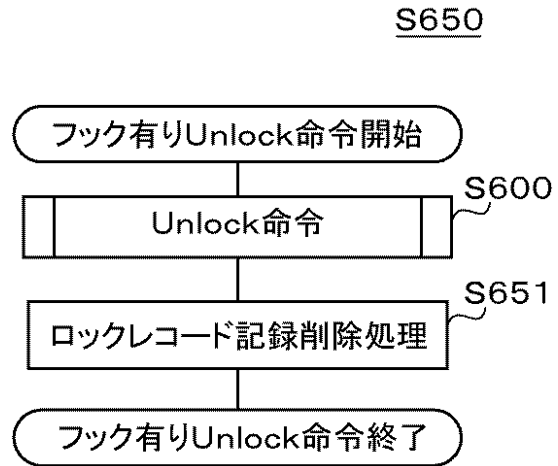
【図 5 A】

図 5 A



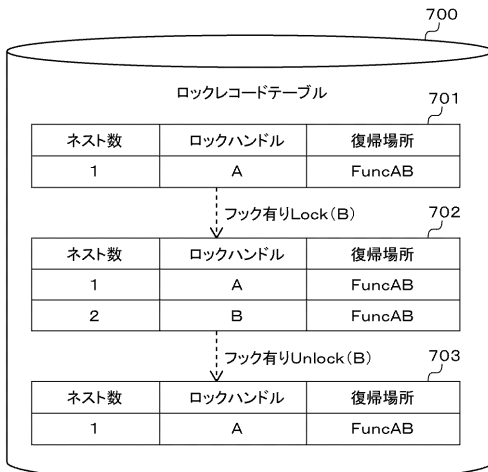
【図 5 B】

図 5 B



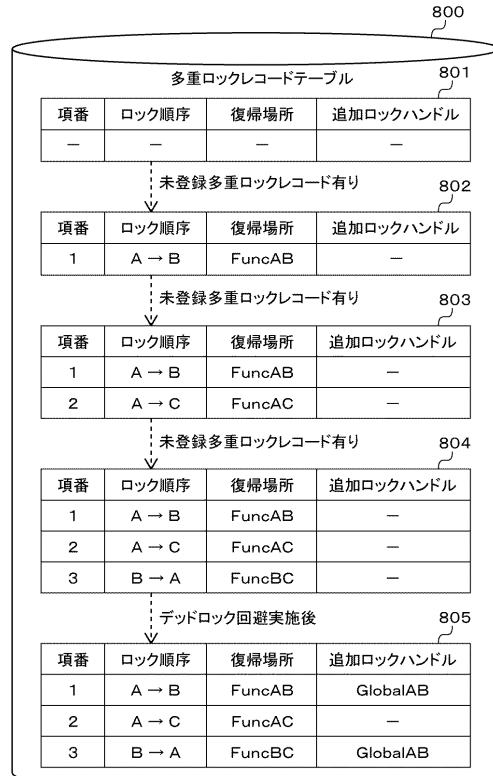
【図 6】

図 6



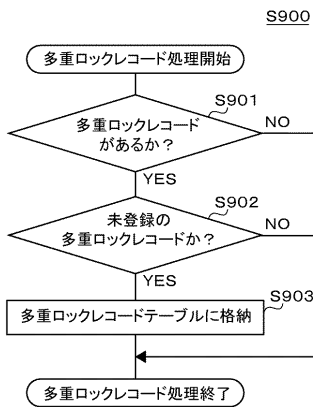
【図 7】

図 7



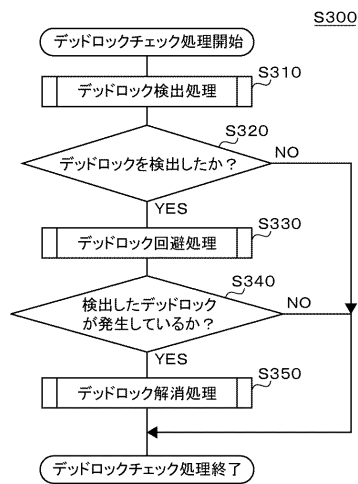
【 図 8 】

図 8



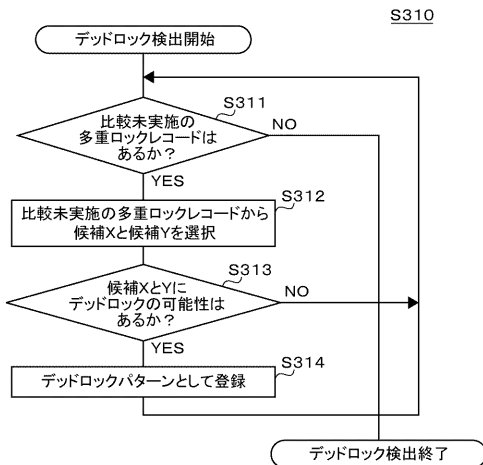
【 図 9 】

図 9



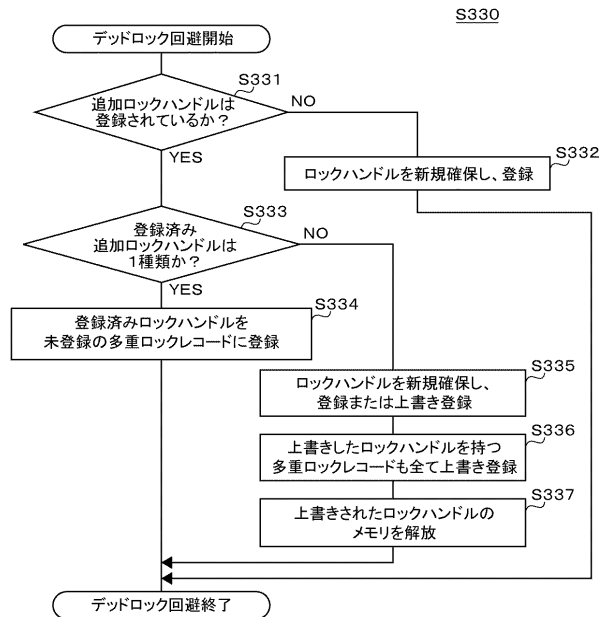
【 図 10 】

図 10



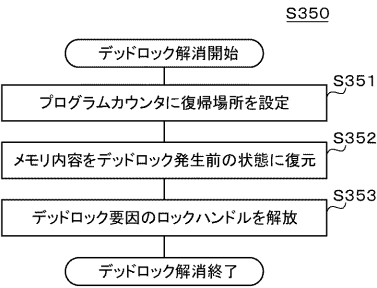
【 図 11 】

図 11



【 図 1 2 】

図 1 2



【 図 1 3 】

図 1 3

