

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3604176号

(P3604176)

(45) 発行日 平成16年12月22日(2004.12.22)

(24) 登録日 平成16年10月8日(2004.10.8)

(51) Int. Cl.<sup>7</sup>

F I

G 0 6 F 9/46

G 0 6 F 9/46 3 6 0 C

G 0 6 F 12/08

G 0 6 F 12/08 5 2 5 A

請求項の数 5 (全 18 頁)

(21) 出願番号	特願平6-219655	(73) 特許権者	000003078
(22) 出願日	平成6年9月14日(1994.9.14)		株式会社東芝
(65) 公開番号	特開平8-83253		東京都港区芝浦一丁目1番1号
(43) 公開日	平成8年3月26日(1996.3.26)	(74) 代理人	100083161
審査請求日	平成12年9月13日(2000.9.13)		弁理士 外川 英明
		(72) 発明者	岡本 利夫
			神奈川県川崎市幸区小向東芝町1番地 株
			式会社東芝 研究開発センター内
		(72) 発明者	木村 哲郎
			神奈川県川崎市幸区小向東芝町1番地 株
			式会社東芝 研究開発センター内
		(72) 発明者	白川 健治
			神奈川県川崎市幸区小向東芝町1番地 株
			式会社東芝 研究開発センター内

最終頁に続く

(54) 【発明の名称】 仮想空間管理方法及び分散処理システム

(57) 【特許請求の範囲】

【請求項1】

複数の計算機と、前記複数の計算機で共有される単一の仮想空間を管理する管理計算機と、前記複数の計算機と前記管理計算機とのそれぞれの間で相互に通信を行うための通信手段とから成るシステムにおける仮想空間管理方法であって、

前記管理計算機は、前記仮想空間から記憶領域を各計算機に割り当てる場合に、この割り当てる記憶領域が既に割り当て済の記憶領域と重複していないことを保証する重複回避機構を備えており、

前記複数の計算機のそれぞれは、記憶領域の確保が必要になったとき、必要量より大きい量の記憶領域を、前記管理計算機の前記重複回避機構を介して確保し、

前記複数の計算機のそれぞれは、確保した前記記憶領域を、他の計算機と通信することなく自計算機内で管理するようにし、

前記複数の計算機のそれぞれは、自計算機以外の計算機が管理する記憶領域を使用したい場合は、この前記記憶領域を管理する計算機と前記通信手段により通信を行い、使用の許可を得るようにしたことを特徴とする仮想空間管理方法。

【請求項2】

複数の計算機と、前記複数の計算機で共有される単一の仮想空間を管理する管理計算機と、前記複数の計算機と前記管理計算機とのそれぞれの間で相互に通信を行うための通信手段とを備える分散処理システムであって、

前記管理計算機は、

10

20

前記仮想空間から記憶領域を各計算機に割り当てる場合に、この割り当てる記憶領域が既に割り当て済の記憶領域と重複していないことを保証する重複回避機構を備え、  
前記複数の計算機のそれぞれは、  
前記通信手段を介して前記管理計算機へ記憶領域の確保の要求を行い、必要量より大きい量の記憶領域を前記重複回避機構を介して確保する確保手段と、  
前記確保手段で確保した前記記憶領域を、他の計算機と通信することなく自計算機内で管理する管理手段と、  
自計算機以外の計算機が管理する記憶領域を使用したい場合に、前記使用したい記憶領域を管理する計算機と前記通信手段により通信を行い、使用の許可を得るよう取得手段とを備えたことを特徴とする分散処理システム。

10

【請求項 3】

前記複数の計算機のそれぞれは、割り当てられた前記記憶領域をさらに複数のサブ領域に分け、このサブ領域ごとに保護属性を設定して管理することを特徴とする請求項 1 記載の仮想空間管理方法。

【請求項 4】

各計算機はそれぞれディスク装置を備えており、自計算機に割り当てられた前記記憶領域に配置したプログラムあるいはデータを、自計算機の備えるディスク装置に対して退避あるいは復旧することを特徴とする請求項 1 記載の仮想空間管理方法。

【請求項 5】

前記重複回避機構は、既に割り当て済の記憶領域と割り当てた計算機を示す情報とを対応付けた管理テーブルを備えるようにしたことを特徴とする請求項 1 記載の仮想空間管理方法。

20

【発明の詳細な説明】

【0001】

【産業上の利用分野】

本発明は、仮想記憶空間管理機構を有する計算機システムが、相互にネットワーク等を介して接続され、分散システムを構築している場合の仮想空間を管理する方式に関する。

【0002】

【従来の技術】

近年、計算機性能の向上とネットワーク技術の進歩により、複数の計算機上でプログラムを動作させ、通信しながら協調して処理を進めていく分散処理が一般化してきている。

30

【0003】

従来からの計算機を用いての処理形態では、1 台の計算機で、単独にすべての処理を行っていた。したがって計算機資源を管理し、アプリケーションプログラムにサービスを提供する基本プログラムであるオペレーティングシステム（OS）は、単独システム用のものが古くから開発され利用されてきた。たとえば AT & T 社が開発した UNIX や IBM 社の MVS などの OS がその代表例である。

【0004】

ここで、新しい分散処理の形態に拡張するために、従来からの OS を拡張して対応させることは、比較的容易である。つまり、従来からの機能を維持しつつ、新しい分散機能を追加提供していく。しかし、拡張には自ずと限度があり、将来の新しい機能の要望に対して柔軟に対応できない可能性がある。また、OS の上ですべてのアプリケーションプログラムが動作するので、OS 自身の実行効率、信頼性などの点が、拡張を繰り返し、OS の内部が複雑になるにつれて問題になってくる。

40

【0005】

従来の単独システム用の OS では、アプリケーションプログラムの実行する場であるアドレス空間は、計算機ごとに管理されてきた。したがって、従来の OS を分散システムに拡張するには、専用のシステムコールを設けるのが常であった。たとえば、Unix 系の OS である 4.3 BSD では、ソケットインタフェースという新しいシステムコールを設けて、従来のアドレス空間上でアプリケーションプログラムを動かす。従来の（分散を使わ

50

ない)アプリケーションを変更無しに動作できるが、アプリケーションプログラマは、分散対応のアプリケーションを設計する際に、そのシステム構成を理解しなくてはならない欠点があった。また、OSの種類によって、分散化への拡張方法が異なるので、アプリケーションプログラムを他のOS上で動作するように変更しなければならないという欠点もあった。

【0006】

以上まとめると、従来のOSを分散環境に対応するように拡張すると、分散対応のアプリケーションを実行する仮想空間の対応がばらばらになり、プログラムの再利用、共有データの利用などが非常に煩雑になる欠点があった。

【0007】

【発明が解決しようとする課題】

本発明においては、上記の欠点を解決するために、単一仮想空間、つまり、すべてのアプリケーションを実行する仮想空間を同一にしたOSに関し、単一仮想空間の分散環境への拡張方法を単一仮想空間をネットワークワイドに広げることで、ネットワーク透過な仮想空間をOSが提供し、アプリケーションレベルに分散化を意識させないようにする。特に、本発明は、上記のようなネットワーク透過な単一仮想空間を実現する際に問題となる仮想空間の管理法について解決するものである。

【0008】

即ち、本発明は、実際には単体で処理を行う計算機が複数で共通して1つの仮想空間を利用しても、相互の計算機がどのようにこの空間を利用しているか等の管理を簡単にでき、OSやアプリケーションを効率よく動作させることのできる仮想空間管理装置及び方法の提供を目的とする。

【0009】

【課題を解決するための手段】

本発明は、複数の計算機と、前記複数の計算機で共有される単一の仮想空間を管理する管理計算機と、前記複数の計算機と前記管理計算機とのそれぞれの間で相互に通信を行うための通信手段とから成るシステムにおける仮想空間管理方法であって、前記管理計算機は、前記仮想空間から記憶領域を各計算機に割り当てる場合に、この割り当てる記憶領域が既に割り当て済の記憶領域と重複していないことを保証する重複回避機構を備えており、前記複数の計算機のそれぞれは、記憶領域の確保が必要になったとき、必要量より大きい量の記憶領域を、前記管理計算機の前記重複回避機構を介して確保し、前記複数の計算機のそれぞれは、確保した前記記憶領域を、他の計算機と通信することなく自計算機内で管理するようにし、前記複数の計算機のそれぞれは、自計算機以外の計算機が管理する記憶領域を使用したい場合は、この前記記憶領域を管理する計算機と前記通信手段により通信を行い、使用の許可を得るようにした。

【0010】

また、本発明は、複数の計算機と、前記複数の計算機で共有される単一の仮想空間を管理する管理計算機と、前記複数の計算機と前記管理計算機とのそれぞれの間で相互に通信を行うための通信手段とを備える分散処理システムであって、前記管理計算機は、前記仮想空間から記憶領域を各計算機に割り当てる場合に、この割り当てる記憶領域が既に割り当て済の記憶領域と重複していないことを保証する重複回避機構を備え、前記複数の計算機のそれぞれは、前記通信手段を介して前記管理計算機へ記憶領域の確保の要求を行い、必要量より大きい量の記憶領域を前記重複回避機構を介して確保する確保手段と、前記確保手段で確保した前記記憶領域を、他の計算機と通信することなく自計算機内で管理する管理手段と、自計算機以外の計算機が管理する記憶領域を使用したい場合に、前記使用したい記憶領域を管理する計算機と前記通信手段により通信を行い、使用の許可を得るよう取得手段とを備えた。

【0014】

【作用】

本発明では、各計算機が処理のために必要とする仮想空間の一部を獲得する(あるいは解

10

20

30

40

50

放する)頻度は高いが、処理のために利用する部分は単一仮想空間の全てではなくごく一部である(仮想空間のローカリティ)という事実を利用して、仮想空間を管理する。つまり、仮想空間を比較的大きな単位で分割し、その分割された領域ごとに仮想空間の一部を各計算機に割り当てて管理する。

**【0015】**

本発明によれば、全てのアプリケーションプログラムを一つの仮想空間に配置する単一仮想空間のOSを分散環境に拡張するにあたり、分散環境を構成する全ての計算機で同一の単一仮想空間を共有する分散単一仮想空間方式を取り、しかも上述した領域の割り当てによりこの分散単一仮想空間の管理が簡単化されているため、単一計算機上で動作していたプログラムを変更なしに動作させることが可能であり、アプリケーションプログラムを作成するプログラマに分散環境の複雑な仮想空間管理を意識させることなくプログラムを作成可能になる。

10

**【0016】**

また、上述した各計算機への領域の割り当てを集中管理する仮想空間管理装置を分散システム内に設け、割り当ての管理と、各計算機ごとに独立に行う計算機資源の管理とを分離することにより、全体の動作を効率化できる。この仮想空間管理装置(実施例ではチャプターサーバと呼ぶものに相当)に対し、各計算機は新たに仮想空間上の領域を必要とする際にこれを要求し、割り当てられた領域を確保し、その領域にプログラムやデータを配置して実行することになる。

**【0017】**

ここで、複数の計算機により共有された前記仮想空間の一部を各計算機に割り当てる単位である領域と、各計算機が個々に、割り当てられた領域内でアクセス保護を実現する単位であるサブ領域とを、独立にすることで、さらに全体の動作を効率化できる。

20

**【0018】**

また、分散単一仮想空間の管理を効率よく実現するために、プログラムの実行上の性質を利用して、仮想空間の領域は割り当てられた計算機が一元的に管理するようにすると、多数の計算機が一つの仮想空間上のプログラムを実行しても管理が容易で、効率が低下しない。

**【0019】**

ここで、仮想空間の領域に配置されたプログラムやデータのバッキングストレージを、その領域が割り当てられた計算機のディスク装置で行うようにすれば、領域内の管理が割り当てられた計算機内に閉じて行えるため、さらに効率化できる。

30

**【0020】****【実施例】**

以下、図面を参照しながら本発明の一実施例について説明する。

**(実施例1)**

図1に本実施例に係る分散システムの全体構成を示す。高速LAN(Local Area Network)に結合された複数のホストから分散システムが構成されている。ただし、ホストとは(1つ以上の)CPUとメモリから構成される計算機で、かつLAN上(OSIのネットワーク層)で識別のためのアドレス等のような識別子が与えられる単位である。さらに、多くの場合、ホストは、プログラム、データ、ファイルなどを記憶しておくディスク装置を持っている。ホストでの実行時は、ディスクの内容は、後述するNSVSに裏打ちされた物理メモリにロードされ実行される。

40

**【0021】**

ホストには、分散システムを構成する為にオペレーティングシステム(OS)がのっている。このOSは、プログラムやデータなどを配置する仮想記憶空間の管理方法として単一仮想記憶空間管理方式を採用する。つまり、OSの機能によって、ホスト上で動作するプログラムや、それを利用したりプログラムを開発したりするユーザからは、1つの仮想記憶空間しか存在していないように見せる。その仮想記憶空間を分散システム全体で共有している。

50

## 【 0 0 2 2 】

次に、分散システムで管理する仮想記憶空間の構成を図 2 に示す。本実施例では、分散システム全体で 1 つの巨大な仮想記憶空間を管理し、その上に本システムで扱うすべてのプログラム、データ、ファイルなどを配置する。ここで、この分散システム全体に渡って共有される単一の仮想記憶空間のことをネットワーク単一仮想記憶空間 ( N S V S ) と呼ぶ。本実施例では空間を 6 4 ビットで表現するものとする。したがって空間の大きさは ( 2 の 6 4 乗 ) b y t e となる。

## 【 0 0 2 3 】

この空間は OS が管理して作り出している仮想的なものである。実際には、各ホストの上での仮想記憶空間機構で実現する。つまり、各ホストは、それぞれに仮想記憶空間を別々に保持しており、その中のプログラム配置はすべてのホストから共通である。ただし、すべてのホストが、すべてのプログラムをもっているわけではなく、各ホストの保持する仮想記憶空間内にはそのホストにおける実行に必要なプログラムが配置されている。つまり、同じプログラムは同じ仮想記憶空間上の位置 ( アドレス ) に存在することを各ホストの OS が管理している。

10

## 【 0 0 2 4 】

なお、上記のように OS は、どのユーザに対しても同じ一つの空間しか提供しないが、実際には、ホストのハードウェア機能として複数の仮想記憶空間を管理できるのが通常である。したがって、図 3 に示すように、複数の仮想記憶空間を利用しながら、各仮想記憶空間の同じ位置 ( アドレス ) に同じプログラムを割り当てる管理を OS が実現すること

20

## 【 0 0 2 5 】

このように分散システム内の各ホストが協調して単一の仮想記憶空間を構築する場合、各ホストが個別に記憶領域 ( 仮想記憶空間の一部 ) を確保すると、領域が重複してしまう可能性があり、重複した部分の整合性がとれなくなり単一仮想記憶空間の管理を困難にする。よって、あるホストが記憶領域を確保する場合に、確保された領域が他のホストの確保している領域と重複していないことを保証する機構 ( 重複回避機構 ) が必要となる。

## 【 0 0 2 6 】

本実施例では、各ホストに割り当てる記憶領域が重複しないように記憶領域の分配状況を管理する機構を担ったホスト ( 以下チャプターサーバと呼ぶ ) を分散システム内に 1 台設ける。すなわち、分散システムは図 1 に示すように、多数のホストとチャプターサーバが LAN で接続された構成をとる。

30

## 【 0 0 2 7 】

ホストが必要に応じて ( 随時 ) 、必要量分だけの記憶領域の確保を行うと、領域の重複回避のための機構のオーバーヘッドが増大し、性能に重大な影響を及ぼす。また、予め記憶領域を各ホストに分配する方式を採用すると、ホストの追加や削除といった分散システムの運用上頻繁に起こる事態に対し柔軟に対応できない。そこで、各ホストは、チャプターサーバが提供する重複回避機構を介して記憶領域の確保を行う時には、必要量より十分大きい量を一度に確保し ( 以下この領域をメモリチャプターと呼ぶ ) 、その領域から随時必要量だけ切り出し使用する。予め確保した領域を使い果たした場合、再度、重複回避機構を介してメモリチャプターの確保を行う。この方法により、重複回避のオーバーヘッドの性能への影響を低減できる。

40

## 【 0 0 2 8 】

以上で本実施例の構成の概要を説明したので、以下に、チャプターサーバと各ホストの内部構造と動作を分けて詳しく説明する。

本実施例では仮想記憶空間の管理を効率良く行うために、メモリチャプターの大きさを固定とし ( 例えば ( 2 の 4 8 乗 ) b y t e ) 、仮想記憶空間アドレスの上位 1 6 b i t をメモリチャプター番号とする。

[ チャプターサーバ側の内部構造 ]

50

上記のチャプターサーバの構成を図4に示す。チャプターサーバは、分散システム内の他のホストと通信するための通信部41を備える。また、チャプターサーバ内でも他のホストと同様にネットワーク単一仮想記憶空間をアプリケーションに対して提供するために、後述する仮想記憶空間管理部42を備えても良い。

#### 【0029】

さらに、分散システム内の各ホストに対して、記憶領域の割当を行うNSVS管理部43を備える。NSVS管理部43内には、使用中のメモリチャプターとそれが割り当てられたホスト(これをチャプターオーナーと呼ぶ)との対応関係を記録するNSVS管理テーブル44を持つ。NSVS管理テーブル44の構造は、図5に示すように、メモリチャプター番号、ホスト番号、状態の組から成り、いずれかのホストに割り当てられたチャプターの状態は「使用中」となる。

10

ただしチャプターサーバは、各メモリチャプター内のメモリの利用状況に関する情報は保持せず、メモリチャプター内部についてはそれを確保したチャプターオーナーが個別に管理する。

#### [ホスト側内部構造]

分散システム内の各ホストの構成を図6に示す。各ホストは、他のホストやチャプターサーバと通信するための通信部61を備える。また、ユーザプログラム(アプリケーション)からのサービス要求(システムコール)を受け付けるサービス要求受付部64を備える。さらに、ネットワーク単一仮想記憶空間をアプリケーションに対して提供するために、仮想記憶空間管理部62(42と同じもの)を備える。

20

#### 【0030】

仮想記憶空間管理部62内には、そのホストが確保した複数のメモリチャプターを管理するためにメモリチャプター管理リスト63を持つ。メモリチャプター管理リスト63は図7に示すような構造を有し、リスト構造の要素はメモリチャプター管理テーブルと呼ばれるテーブルである。メモリチャプター管理テーブルは、そのホストが確保した複数のメモリチャプターのうちの1つについて、そのメモリチャプター中のメモリの割り当て状況を保持する。メモリチャプター管理テーブルの構造を図8に示す。

#### 【0031】

以下、図7と図8について説明する。各メモリチャプター管理テーブルにはその属性として、そのテーブルがどのメモリチャプターを管理するものを示すためのメモリチャプター番号と、他のメモリチャプター管理テーブルとをリストの形で結合するためのポインタが登録されている。テーブルの内部には、メモリチャプター内の利用状況を記録するために、領域(開始番地と大きさ)とその利用状況(使用中または未使用)が登録されている。

30

#### [動作]

以下、分散システム内のあるホストが記憶領域を確保する場合を例にとり、システム全体の動作を図9を用いながら順を追って説明する。

#### 【0032】

あるホスト上でユーザプログラムが動作し、システムコールによって、メモリ獲得要求が発行されると、まず、そのホストの仮想記憶空間管理部62が、既に獲得しているメモリチャプター内の空き領域を、メモリチャプター管理テーブルを用いて検索する。メモリ獲得要求を満たす大きさの空き領域がある場合は、そこからメモリ領域を確保し割り当てを行い、メモリチャプター管理テーブルを更新する。要求を満たす空き領域が、ホストの保持するメモリチャプター内に存在しない場合は、通信部61を介してチャプターサーバに対し新規メモリチャプターの獲得要求を発行する。

40

#### 【0033】

通信部41を介して新規メモリチャプターの獲得要求を受けとった(S1及びS2Yes)チャプターサーバ内のNSVS管理部43は、空きメモリチャプターをNSVS管理テーブル44を用いて検索する(S3)。即ち、どのホストにも割り当てられていない(状態が「未使用」である)メモリチャプターを検索する。そして、これを要求を発行した

50

ホストに割り当てる。即ち、検索により得られた空きメモリチャプターの番号を通信部 41 を介して要求元のホストへ返す (S4)。さらに、NSVS 管理テーブル 44 内の割り当てたメモリチャプターに該当する欄に、割り当て先のホスト番号を記入し状態を使用中に変更する (S5)。

【0034】

チャプターサーバから新規メモリチャプターを受けとったホストの仮想記憶空間管理部 62 は、受けとったメモリチャプターの管理のために、メモリチャプター管理テーブルを新たに作成し、メモリチャプター管理リスト 62 に加える。そして、メモリ獲得要求を発行したユーザプログラムに対し、新規メモリチャプターの中から記憶領域を確保し割り当てを行い、その割り当て状況を表す情報を該当するメモリチャプター管理テーブルに記録する。最後に要求元へ処理が戻る。

10

【0035】

チャプターサーバ側の処理は、メモリチャプター要求が同時に複数のホストから発行されることを考慮すると、上記検索・割り当ては排他的に行わなければならない。そこで、チャプターサーバは複数のホストからの割り当て要求を並行に処理することはせず、逐次的に処理を行う。

(実施例 2)

ここでは、実施例 1 で説明したメモリチャプターに加え、新たにメモリセクションを加え、2 種類の単位を用いてネットワーク単一仮想記憶空間 (以下、NSVS と呼ぶ) の管理を行う実施例について述べる。

20

【0036】

実施例 1 で述べたように、空間全体が複数のメモリチャプターと呼ばれる領域に分割されている。このメモリチャプターは、仮想記憶空間を複数のホストで分担して管理する時の単位として用いられる。

【0037】

各メモリチャプターはさらに複数のメモリセクションと呼ばれる領域に分割されている。このメモリセクションは、仮想記憶空間内に配置されるファイルやプログラムテキスト、データ等 (以下プログラム等という) の記憶空間内の内容の違いによる割り当ての単位として用いられる。つまり、記憶領域の確保は、メモリセクション単位に行うため、ユーザプログラムは、新規のメモリセクションを確保するためのシステムコール (OS に対するサービス要求) を発行することで獲得できる。

30

【0038】

アクセス保護もメモリセクションを単位として行われる。アクセス保護とは、あるプログラム等へのアクセス (例えば読み出し、書き込み、実行) を、制限する機構のことをいう。つまり、このアクセスをどのスレッド (またはプロセス) が行おうとしているのかという情報や、このアクセスがどのプログラムから行われようとしているのかという情報に基づいて許可 / 禁止することを言う。スレッドとは、プログラムを実行する主体のことをいい、この実施例では、アクセス保護が許される範囲で、複数のメモリセクションをまたがってプログラムを実行し、また、複数のメモリセクション内のデータを読み書きすることができる。

40

【0039】

このようなアクセス制御の具体的機構には、あるスレッド (またはプロセス) によるあるプログラム (またはメモリセクション) からのアクセスに対する許可 / 禁止を表す情報を、各メモリセクションに対応して登録し、プログラム実行時にこの情報を元にアクセス制御することにより、プログラム等に対するアクセス制御を実現する。このアクセス制御の機構は、特願平 5 - 3937 号に詳しい。

【0040】

なお、仮想記憶空間の主記憶 (物理メモリ) への写像は、一般的な仮想記憶空間の実現法であるページング方式を採用しているものとするが、他の方式であっても一般性を失わない。ページング方式の場合は、上記の各メモリセクションは複数のページから成る。つま

50

り、メモリセクションの大きさはページサイズの整数倍であるとする。

【0041】

本実施例では仮想記憶空間の管理を効率良く行うために、メモリチャプター、メモリセクション、ページの各大きさを固定とし（それぞれ（2の48乗）byte、（2の32乗）byte、（2の12乗）byte）、仮想記憶空間アドレスの上位16bitをメモリチャプターID、上位32bitをメモリセクションID、上位52bitをページIDとする。この様子を図10に示す。

【0042】

また、各ホスト上で管理されるスレッドIDは、NSVSを形成するホストの範囲中でユニーク（お互いのスレッドIDが重複しない）であるとする。

10

まず、最初に本実施例の説明で必要となるデータ構造の管理に関して述べ、さらに、メモリ領域の管理手順について述べる。

[データ構造]

NSVSが有効に機能するためには、メモリチャプターやメモリセクションがどのように割り当てられているかの管理が重要である。ここでは、メモリチャプターやメモリセクションの割り当て状況等の管理情報の配置について説明する。

【0043】

NSVS内の各メモリチャプターがどのホストに割り当てられているかという情報の管理は、実施例1と同様、チャプターサーバと呼ばれるホストで集中管理している。この管理は図5に示すような表を用いて行っている。この表には、現在のチャプターIDとチャプターオーナーのホスト名および関連情報の対応が登録されている。

20

【0044】

各ホストには、自ホストで現在利用しているメモリチャプターの情報をメモリチャプター管理リストで保持している。メモリチャプター管理リスト63は、図11に示す構造を採用し、図8のメモリ領域（開始番地と大きさ）の代わりにメモリセクションIDを用い、利用状況（使用中または未使用）のフィールドを拡張して、このメモリセクションを所有するスレッドID（またはプロセスID）を用いるようになっている。この値が0の時は、空きセクションを表す。また、そのメモリチャプターを管理しているオーナーホスト名を表すフィールドを付加した。

【0045】

30

あるメモリチャプター内の全てのメモリセクションに関する管理情報は、チャプターオーナーがこのメモリチャプター管理テーブルで保持する。

さらに、現在利用中のメモリセクションに関する各種の情報は、メモリセクション管理テーブルに記憶されている。メモリセクションの管理情報としては、メモリセクションID、そのメモリチャプターを管理するオーナースレッドのID、オーナーホスト名、アクセス保護情報、バッキングストア（そのメモリセクションの内容を持っているディスク装置等）との対応関係（対応するディスク内のアドレス等）、メモリセクションが他のホストと共有された場合（後述する）の内容のコピーの分配状況、およびコピー間の一貫性維持のための制御情報などがある。このメモリセクション管理テーブルの例（図11のオーナーホストH-Bのもの）を図12に示す。

40

[管理手順]

次に、NSVSに基づく分散システムにおけるメモリ領域の管理手順について説明する。

【0046】

管理手順は、大きく2つに分けられる。1つは、新規にメモリ領域を確保する場合であり、他方は、既に存在しているメモリ領域を利用する場合で他のホストと共有する場合である。

[新規メモリ領域確保]

新規にあるホスト上でメモリ領域を確保する場合、ユーザプログラムはメモリセクション確保のシステムコールを使ってOSに要求する。OSでは、メモリチャプター管理リストをたどりながら、メモリチャプター管理テーブルのチャプターオーナーが自ホスト名であるテ

50



ブルを探し、その中から空きメモリセクションを探す。あったら、要求したスレッドIDをオーナスレッドとして空きメモリセクションのオーナスレッドのところに登録する。要求元には、確保できたメモリセクションIDを返す。さらに、メモリセクション管理テーブルへも確保できたメモリセクションIDに関して新たなエントリを作成し、データを登録する。オーナスレッドは、新しいメモリセクションを作成すると同時に、以後、バッキングストアの管理、アクセス保護情報の管理、他ホストとのメモリ共有に関する管理などを行う。

#### 【0047】

もしそのホストに割り当てられたメモリチャプター内の全てのメモリセクションが使用中である場合には、そのホストは新たなメモリチャプターを確保する。メモリチャプターの確保は、実施例1と同様の方法である。したがって、各ホストに確保されたメモリチャプターの領域は、他のホストが確保した領域と重複しないことが保証されているので、このメモリチャプター内のすべてのメモリセクションは未使用であることが保証される。そして確保されたメモリチャプターに関して新規にメモリチャプター管理テーブルを作成し、データを初期化し、登録する。そこから上記の手順でメモリセクションの割当を行う。

#### [既存メモリ領域利用]

次に、NSVSに基づく分散システムにおける情報共有について考える。NSVSに基づく分散システムでは、全てのホストから同じアドレス上に同じデータが存在するように見せているため、情報の共有は単に同じアドレスをアクセスすることをきっかけに行われる。

#### 【0048】

本実施例では、メモリセクション毎にアクセス制御が行われているため、情報の共有に際しては、まずメモリセクションの共有を行わなければならない。メモリセクションの共有は、あるホストが上述の手続きにより確保したメモリセクションの内容を、他のホストが参照さらには更新できるようにすることにより実現される。

#### 【0049】

ここで、ホストH-AがメモリセクションMS-Aを確保し単独で利用している状態で、ホストH-Bとそのメモリセクションを共有するまでの具体的手順について説明する。メモリセクションMS-Aのチャプターは、MC-Aであり、そのオーナーホスト名は、ホストH-Aである。

#### 【0050】

まず、ホストH-B上で動いているプログラムがメモリセクションMS-A内のあるアドレスをアクセスする。アクセス手段は、ホストH-B上で動作しているスレッドがそのアドレスへコールやジャンプした場合や、そのアドレスのデータを参照した場合などである。しかしメモリセクションMS-AはまだホストH-Bにアタッチされていない、つまり、ホストH-Bで、そのメモリセクションはまだ利用実績がないので、ホストH-Bの仮想記憶管理部62には、メモリセクションMS-Aが登録されていない。したがって、アクセス時にページフォルトが発生し、OSのページフォルト処理手続きが呼び出される。

#### 【0051】

ページフォルトが生じると、OSはまずフォルトを起こしたページがバッキングストアに退避されているページかどうかを調べる。

これは、メモリセクション管理テーブルを調べることで判明する。つまり、自ホストがチャプターオーナである利用中のメモリセクションの場合には、メモリセクション管理テーブルに登録されており、そのメモリセクションの内容が保管されているバッキングストアに登録されているからである。登録されている際は、バッキングストアに退避されているページを、通常のページング方式に従い主記憶(物理メモリ)にロードし、フォルトを起こしたスレッド(またはプロセス)を再開させるための手続きを行う。

#### 【0052】

また、他のホストがチャプターオーナであるメモリセクションの場合には、登録されている分配状況、一貫性制御情報などにしたがって、オーナホストからデータを転送し、主記

10

20

30

40

50

憶（物理メモリ）にロードし、フォルトを起こしたスレッド（またはプロセス）を再開させるための手続きを行う。

【0053】

しかし、メモリセクション管理テーブルに登録されていない場合、他のホストが保持するページに対するアクセスと解釈され、チャプターサーバに対してページID（または、メモリセクションID）を送り、そのページが含まれるメモリチャプターのオーナーの検索を依頼する。

【0054】

具体的には、H-BのOSはページフォルトを起こしたアドレスから、メモリチャプターIDを求める。これは、アドレス上位16ビットがメモリチャプターIDなので求められる。

10

【0055】

次に、自ホスト内のメモリチャプター管理テーブルをサーチして、このチャプターIDが登録されているか調べる。あれば、そのテーブルのオーナーホスト名を調べる。なければ、自分の知らないメモリセクションであると判断し、他のホストが保持するメモリセクションへのアクセスであると見なして、そのメモリセクションが含まれるメモリチャプターのオーナーがどのホストであるかをチャプターサーバに問い合わせる。チャプターサーバはチャプターIDと、それを保持し管理するホストの識別子の対応を保持するNSVS管理テーブル（図5）を持っており、その表からチャプターオーナーを調べてホストH-Bに伝達する。

【0056】

20

検索の結果、そのメモリチャプターが未使用でオーナーがいらない場合、フォルトを起こしたメモリアクセスは不当として処理される。

チャプターオーナーIDが得られた場合には、そのホストに対してアクセス要求を送るとともに、アクセス保護のチェックを受けるために要求したスレッドID（またはプロセスID）などの情報を送る。チャプターオーナーでアクセスが許可された場合には、オーナーホストから該当ページ（またはそれを含むメモリセクション全体）のコピーが送られてくるので、それを受けとり、フォルトを起こした（要求した）スレッド（またはプロセス）を再開させるための手続きを行う。さらに、そのメモリチャプターやメモリセクションを管理するためのメモリチャプター管理テーブルを必要なら要求し、送られてきたテーブルをメモリチャプター管理リストに加える。既に、存在している場合は、更新される。

30

【0057】

また、そのメモリセクションに関する情報をメモリセクション管理テーブルに追加する。具体的には、ホストH-Bは得られたチャプターオーナー（この場合ホストH-A）に対して、アクセスしたスレッドID（またはプロセスID）と共に、メモリセクションMS-Aの共有要求を送り、メモリセクションの内容（コピー）と必要に応じてそのメモリチャプター管理テーブルを受けとる。

【0058】

一方、要求を受けた側のチャプターオーナーH-Aは、要求元のスレッドIDをメモリセクション管理テーブルのアクセス保護情報に従ってアクセス許可を判断し、許可されると、該当メモリセクションのコピーをホストH-Bに送り、自分の保持するメモリセクションのコピーがホストH-Bにできたことを自ホスト内のメモリセクション管理テーブルの分配状況に登録するとともに、コピー間の一貫性維持のための一貫性制御情報に登録して、ロックの管理や、更新されたデータの配布、ホストH-Bに対するメモリセクションMS-Aのバッキングストアのサービス等を提供する。データの更新が生じた時やロック/アンロックなどのイベントが生じた時に、コピーを持つホストH-Bに伝達する。要求に応じて、メモリチャプター管理テーブルの内容を送る。

40

【0059】

オーナーホストにて使用中のメモリセクションが不要になり、メモリセクション解放要求が発行されると、OSはメモリチャプター管理リスト（図11）の中を検索し、解放され

50

たメモリセクションの状態を未使用に変更し、ページテーブルからそのメモリセクションに属するページを全て削除する。さらに、メモリセクション管理テーブルからも該当メモリセクションのエントリの開放作業を行う。他ホストと共有関係にある場合は、それを加味した開放作業を行う。

#### [チャプターサーバ]

上述のようにチャプターサーバは、新たなメモリチャプターを割り当てる処理と、アドレスまたはメモリセクションIDまたはメモリチャプターIDからチャプターオーナーがどのホストであるかを検索・応答する処理の2種類のサービスを提供する。

#### 【0060】

このサービスを提供するために、チャプターサーバはNSVS管理テーブル(図5)を内部的に保持している。このテーブルはメモリチャプターIDをキーとするテーブルであり、チャプターオーナーの検索が高速に行えるように構成されている。さらに、このテーブルをリスト構造とすれば、未使用領域の検索を高速に行える。

10

#### 【0061】

つぎにチャプターサーバの処理手順について図13を用いて説明する。チャプターサーバは分散システム内のホストが発行するサービス要求を受信し(S11)、その要求が新規メモリチャプター要求である場合には(S12 Yes)、NSVS管理テーブル44から未使用メモリチャプターを検索し(S15)、そのIDを要求を発行したホストに返す(S16)。尚、空きメモリチャプターがない場合は、その旨を要求元のホストに伝える。この時、NSVS管理テーブル44内の未使用メモリチャプターがリストを構成していると、この検索が高速に行える。

20

#### 【0062】

サービス要求がチャプターオーナーの問い合わせの場合には(S13 Yes)、要求とともに送られてきたチャプターIDをキーとしてNSVS管理テーブル44を引くことにより、簡単かつ高速にチャプターオーナーのIDを得ることができ、それを要求を発行したホストに返す(S14)。以上の処理を繰り返すことによりチャプターサーバの機能が実現される。また、サービス要求が発行されない間は、分散システムの他のホストと同様、普通の仕事もすることができる。

#### (実施例3)

本実施例に係るホストの構成を図14に示す。本実施例で実施例1の構成と異なる点は、ディスク管理部66及びその内部のディスク領域管理テーブル67が付加される点である。

30

#### 【0063】

この実施例では、自己がそのチャプターオーナーであるチャプター(即ち、ネットワーク単一仮想記憶空間中の自ホストが管理している領域)に割り当てられているデータを実際に格納する領域として、自ホストに接続されているディスクの領域を利用する。

#### 【0064】

以下、図15のフローチャートに従って説明する。あるホスト上で、動作中のユーザ、アプリケーションプログラムから仮想空間の一部領域の獲得要求(システムコール)が発行されると、その要求は、サービス要求受付部64に入る(S1)。つぎに、要求内容がメモリ領域の獲得要求かどうか調べ(S2)、そうであると、仮想記憶空間管理部62に処理が進む。そうでないと、その他の処理実行部65に処理が進む(S3)。

40

#### 【0065】

仮想記憶空間管理部62では、まず、すでに獲得してあるメモリチャプター内の空き領域から、メモリチャプター管理リスト63を用いて要求サイズを満たす大きさのものを検索する(S4)。この検索の結果を判定し(S5)、メモリ獲得要求を満たす大きさの空き領域がある場合は、そこからメモリ領域を確保し割り当てを行い、メモリチャプター管理テーブルを更新する(S6)。次に、確保した領域に対するディスク領域の割当を要求するために、ディスク管理部66へ、ディスクの割当を要求する(S7)。ディスク管理部66からの返答により、ディスクの割当が確保できたか判定し(S8)、確保できないと

50

エラー処理をして終える（S 9）。確保できると、そのディスクの領域（ディスクアドレス）をメモリチャプター管理テーブルにあわせて登録する（S 10）。

【0066】

この場合のメモリチャプター管理リスト63の構造は図16のようになり、ディスク内の記憶領域の位置（ディスクアドレス）など必要な情報を、仮想記憶空間上の領域と対応させて記憶するようになっている。そして最後に、メモリ獲得サービスを要求を発行元（スレッドまたはプロセス、タスク）に対して獲得できた仮想記憶空間上の領域の開始アドレス他必要な情報をサービス要求受付部64を介して伝え、処理を終える（S 11）。

【0067】

要求を満たす大きさの空き領域が、ホストの保持するメモリチャプター内に存在しない場合は、チャプターサーバ通信部61を介してチャプターサーバに対し新規メモリチャプターの獲得要求を発行する（S 12）。チャプターサーバをもつホストでは、通信部41を介して新規メモリチャプターの獲得要求を受けとったチャプターサーバ内のNSVS管理部43が、空きメモリチャプターをNSVS管理テーブル44を用いて検索し、要求を発行したホストに割り当てる。さらに、NSVS管理テーブル44内の割り当てたメモリチャプターに該当する欄に、割り当てたホストIDを記入し状態を使用中に変更する。以上のチャプターサーバ内の処理は、実施例1と同様である。

10

【0068】

要求元のホスト側では、チャプターサーバ通信部61を介してチャプターサーバから返答がくると、この返答から、新規メモリチャプターが獲得できたかチェックする（S 13）。獲得できなかったときは、エラー処理をして終える（S 15）。新規メモリチャプターを受けとったホストは、受けとったメモリチャプターの管理のために、メモリチャプター管理テーブルを新たに作成し、メモリチャプター管理リスト63に加える（S 14）。そして、メモリ獲得要求を発行したスレッド（またはプロセス、タスク）と、新規メモリチャプターに関して上記S 6以降の処理を行う。

20

【0069】

以下に、ディスク管理部66の動作を図17を用いながら説明する。ディスク管理部66には、ディスクの利用状況を管理するディスク領域管理テーブル67がある。その構成は図18に示すように、ディスクのセクター単位に管理されたディスクアドレスとその利用状況を記憶するようになっている。テーブルの大きさは、全物理ディスク容量をカバーできる範囲である。

30

【0070】

まずディスク管理部66内にある要求受付部が要求を受け付ける（S 21）。要求内容は、ディスク領域の要求、解放などである。要求の場合は、パラメータとして容量を添える。また、解放の場合は、解放するディスクアドレスを添える。

【0071】

以下の説明では、領域要求の場合（S 22 Yes）で説明する。要求された容量を、まず、テーブル67の空き容量に格納されている値と比較する（S 24）。足りる場合は、要求量を引いて空き容量のデータを更新する（S 26）。不足する場合は、要求元（仮想記憶空間管理部62）にその旨を伝えて処理を終了する（S 25）。

40

【0072】

足りた場合、次に空きセクターを探して確保する。空きポイントによって空きセクターがリスト構造でつながっているため、必要な量のセクターをリストをたどって取り出す。つまり必要な数のセクター分、ポイントをたどりながらポイント先のアドレスを記憶し、このアドレスのエントリの状態フィールドを使用中にかえていく（S 27, S 28, S 29 No, S 31）。

【0073】

必要な数のセクタを確保したら（S 29 Yes）、記憶しておいたアドレスを要求元に返して処理を終了する（S 30）。尚、ディスクの解放の処理は、上記と逆の手順をふみ、返されたセクターの分だけ空き容量を増やし、空きポイントに返されたセクターのアド

50

レスをつなぎ、そのアドレスの状態を未使用にする（S 2 3 の例）。

【0074】

【発明の効果】

以上、説明したように、本発明によれば、複数の計算機をネットワークを介して結合した分散環境において、メモリーチャプターという仮想空間を比較的大きく分割する管理単位を導入し、この単位を用いて分散環境内の各計算機に記憶領域を割り当てる。

【0075】

一度、計算機にメモリチャプターを割り当てられると、その内の管理は、自計算機内だけで独占して行えるので、分散環境下でも、プログラム実行時のプログラムコードやデータなどを格納するところである仮想空間の領域の管理が他の計算機と交信して決めることなく独立して行えるので、通信量を減らせ、オーバヘッド少ない仮想空間管理が可能になる。

10

【0076】

よって、分散環境下でのオペレーティングシステム（OS）が行っている仮想空間の管理機能において、どの計算機からでも同じ仮想空間を共有する仮想空間の管理法の効率化が実現できる。

【図面の簡単な説明】

【図1】本発明の一実施例のシステム構成を示す図。

【図2】本実施例システムで管理する仮想記憶空間の構成例を示す図。

【図3】図2とは別の構成例を示す図。

【図4】チャプターサーバの内部構成を示す図。

20

【図5】NSVS管理テーブルの構造例を示す図。

【図6】各ホストの内部構成を示す図。

【図7】メモリチャプター管理リストの構造例を示す図。

【図8】メモリチャプター管理テーブルの構造例を示す図。

【図9】実施例1におけるチャプターサーバの動作の流れを示すフローチャート。

【図10】仮想記憶空間アドレスの使用法を表す図。

【図11】メモリチャプター管理リストの別の構造例を示す図。

【図12】メモリセクション管理テーブルの構造例を示す図。

【図13】実施例2におけるチャプターサーバの動作の流れを示すフローチャート。

【図14】実施例3における各ホストの内部構成を示す図。

30

【図15】実施例3のシステムの動作の流れを示すフローチャート。

【図16】メモリチャプター管理リストのさらに別の構造例を示す図。

【図17】ディスク管理部の動作の流れを示すフローチャート。

【図18】ディスク領域管理テーブルの構造例を示す図。

【符号の説明】

4 1、6 1 ... 通信部

4 2、6 2 ... 仮想記憶空間管理部

4 3 ... NSVS 管理部

4 4 ... NSVS 管理テーブル

6 3 ... メモリチャプター管理リスト

40

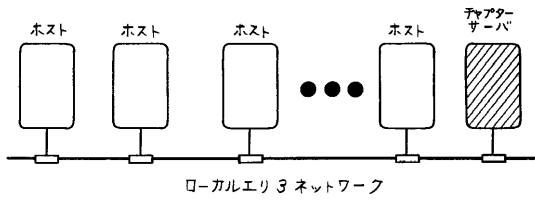
6 4 ... サービス要求受付部

6 5 ... その他の処理実行部

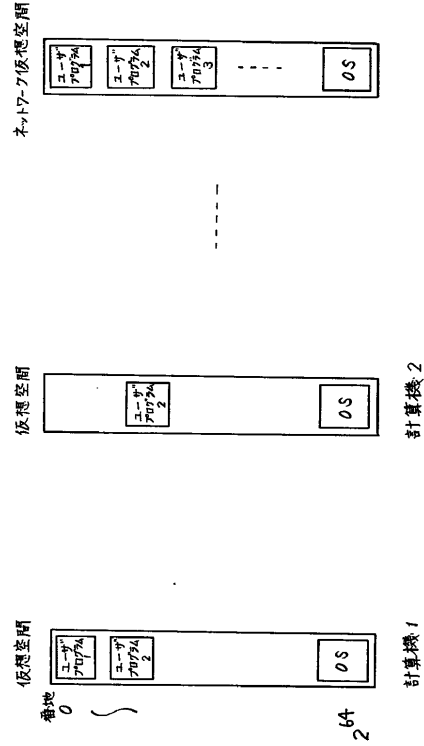
6 6 ... ディスク管理部

6 7 ... ディスク領域管理テーブル

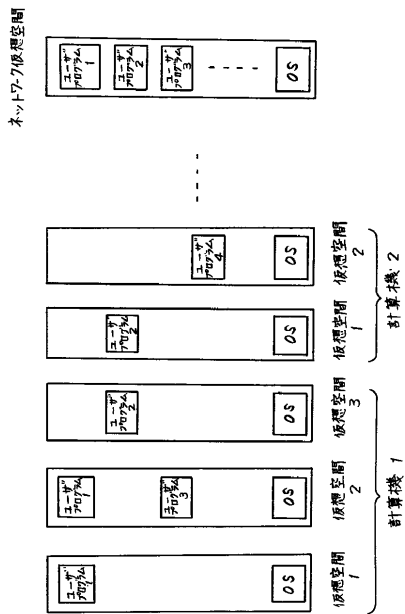
【図 1】



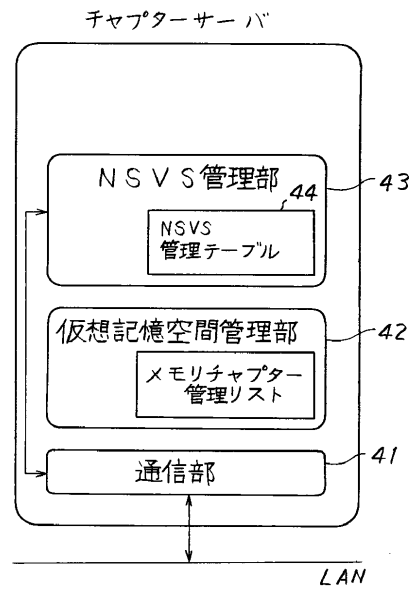
【図 2】



【図 3】



【図 4】

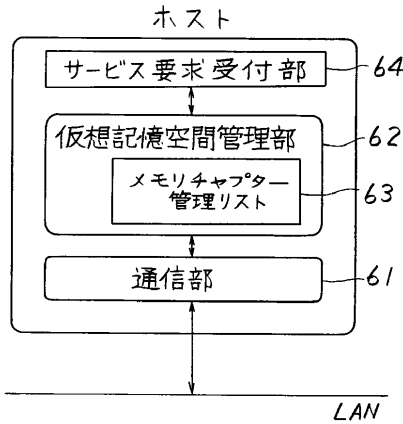


【 図 5 】

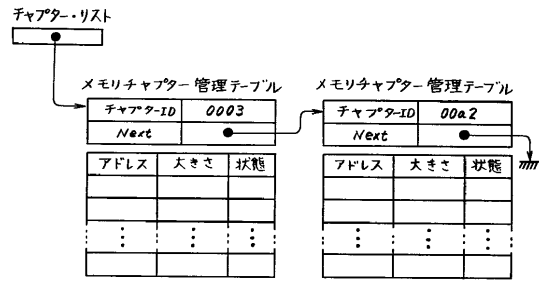
(各IDは16進数表記)

チャプターID	ホストID	状態
0000	—	未使用
0001	1c02	使用中
0002	—	未使用
0003	930f	使用中
⋮	⋮	⋮
ffff	—	未使用

【 図 6 】



【 図 7 】



【 図 8 】

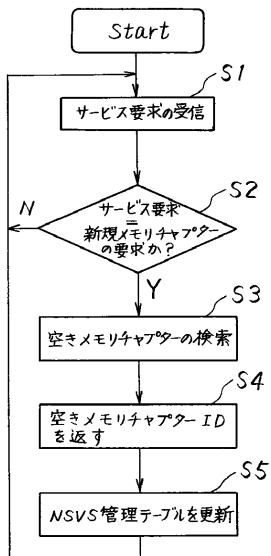
(各IDは16進数表記)

チャプターID	次のチャプターへのポインタ
0003	0007

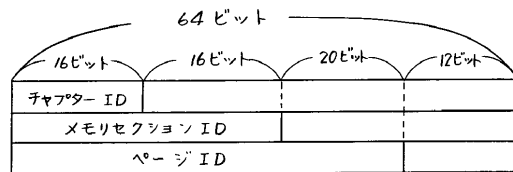
  

開始アドレス	大きさ	状態
0003 0000 0000 0000	1 0000 0000	未使用
0003 0001 0000 0000	1 0000 0000	使用中
0003 0002 0000 0000	3 5000 0000	使用中
0003 0005 5000 0000	5000 0000	未使用
⋮	⋮	⋮
0003 ffff 0000 0000	1 0000 0000	未使用

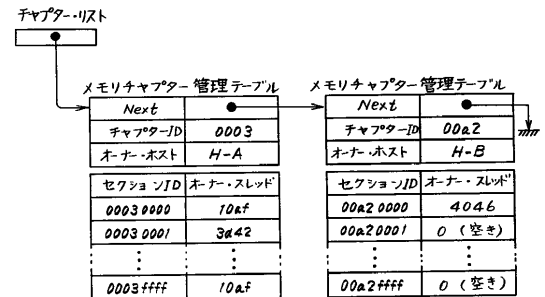
【 図 9 】



【 図 10 】



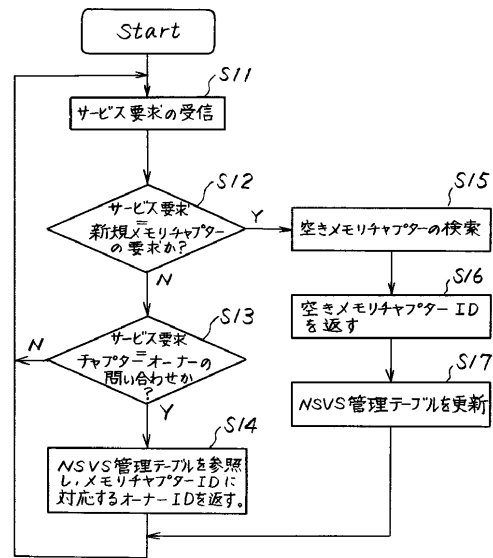
【 図 11 】



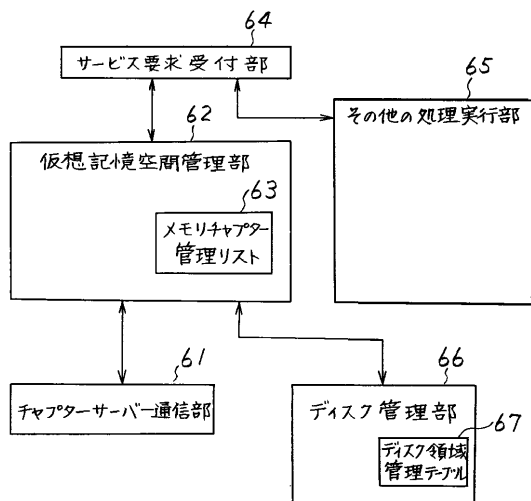
【 図 1 2 】

メモリチップID	0003 0000 ..... 00a2 0000	ホースレット	10af 4046	ホース	H-A H-B	アドレス保護	002 1W- 002 1-X	バックアップストア	ディスクTA 0100 ---	分割状況	H-B H-A, H-D	一貫性制御情報	Copy, パラメータ, バージョン, ...	その他	
----------	---------------------------------	--------	--------------	-----	------------	--------	--------------------	-----------	--------------------	------	-----------------	---------	----------------------------	-----	--

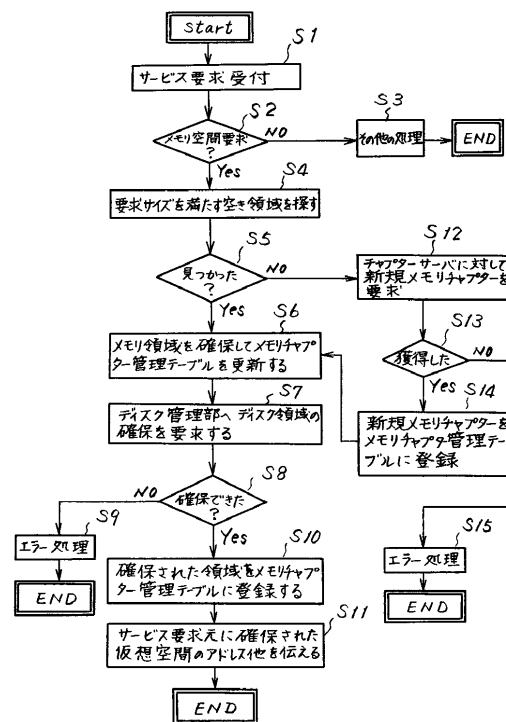
【 図 1 3 】



【 図 1 4 】

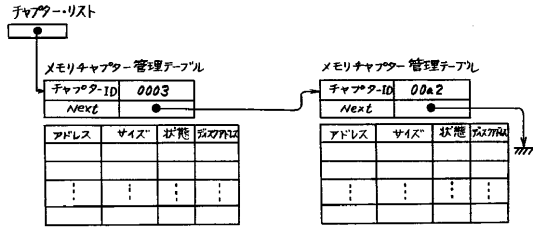


【 図 1 5 】

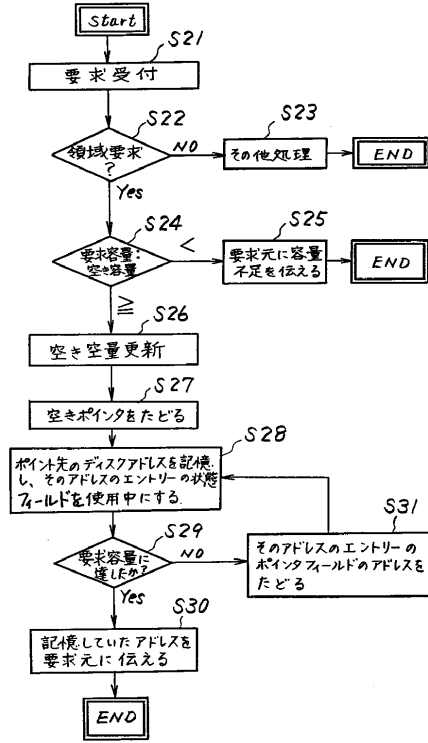




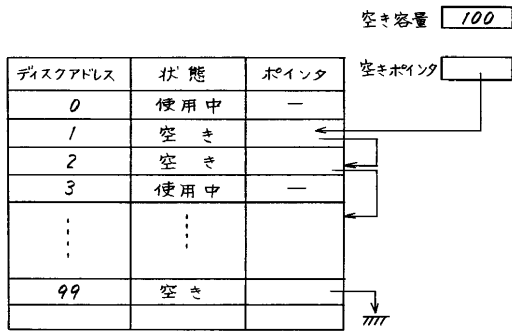
【 図 16 】



【 図 17 】



【 図 18 】



---

フロントページの続き

審査官 鳥居 稔

- (56)参考文献 特開平04 - 362751 (JP, A)  
特開平05 - 257803 (JP, A)  
特開平05 - 257811 (JP, A)  
特開平06 - 052054 (JP, A)

(58)調査した分野(Int.Cl.<sup>7</sup>, DB名)

G06F 9/46  
G06F 12/08-10  
G06F 15/16