

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3754393号

(P3754393)

(45) 発行日 平成18年3月8日(2006.3.8)

(24) 登録日 平成17年12月22日(2005.12.22)

(51) Int. Cl.

G06F 9/46 (2006.01)

F I

G06F 9/46 420Z

請求項の数 8 (全 13 頁)

(21) 出願番号	特願2002-153004 (P2002-153004)	(73) 特許権者	000003078
(22) 出願日	平成14年5月27日(2002.5.27)		株式会社東芝
(65) 公開番号	特開2003-345613 (P2003-345613A)		東京都港区芝浦一丁目1番1号
(43) 公開日	平成15年12月5日(2003.12.5)	(74) 代理人	100076233
審査請求日	平成15年9月19日(2003.9.19)		弁理士 伊藤 進
		(72) 発明者	佐藤 記代子
			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	前田 誠司
			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	崎山 伸夫
			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

最終頁に続く

(54) 【発明の名称】 分散ファイル装置及びそのプロセスマイグレーション方法並びにコンピュータ装置

(57) 【特許請求の範囲】

【請求項1】

ネットワーク上に接続された複数の計算機ノードに夫々設けられる2次記憶装置と、前記各計算機ノードに夫々設けられ、自計算機ノードの計算機が実行するプロセスに従って退避ファイルを作成する退避ファイル作成手段と、前記退避ファイル作成手段で該プロセス毎に作成された前記退避ファイルを前記ネットワーク上の任意の計算機ノードの2次記憶装置に記憶させると共に、前記任意の計算機ノードの2次記憶装置とは異なる他の計算機ノードの2次記憶装置に前記退避ファイルを複製した複製退避ファイルを記憶させる記憶制御手段と、実行中のプロセスを他の計算機ノードに移送する場合に、前記実行中のプロセスに従って作成された退避ファイルの複製退避ファイルが記憶された2次記憶装置が属する計算機ノードを、前記プロセスの移送先に決定する決定手段とを具備したことを特徴とする分散ファイル装置。

10

【請求項2】

前記ネットワーク上の任意の計算機ノードの2次記憶装置は、前記プロセスを実行中の計算機ノードの2次記憶装置であることを特徴とする請求項1に記載の分散ファイル装置。

【請求項3】

前記記憶制御手段との間でデータの授受を行って、前記退避ファイル及び複製退避ファイルの前記ネットワーク上の記憶位置を示す管理テーブルを生成・管理する管理手段を更に備え、

20

前記決定手段は、前記管理テーブルを用いて、プロセスの移送先を決定することを特徴とする請求項 1 に記載の分散ファイル装置。

【請求項 4】

前記記憶制御手段及び前記決定手段は、移送するプロセスを実行する計算機ノード以外のネットワーク上に設けられることを特徴とする請求項 1 に記載の分散ファイル装置。

【請求項 5】

ネットワーク上に接続された複数の計算機ノードに夫々設けられる 2 次記憶装置にアクセスする処理と、

自計算機ノードの計算機が実行するプロセスに従って退避ファイルを作成する処理と、プロセス毎に作成された前記退避ファイルを前記ネットワーク上の任意の計算機ノードの 2 次記憶装置に記憶させると共に、前記任意の計算機ノードとは異なる他の計算機ノードの 2 次記憶装置に前記退避ファイルを複製した複製退避ファイルを記憶させる処理と、実行中のプロセスを他の計算機ノードに移送する場合に、前記実行中のプロセスに従って作成された退避ファイルの複製退避ファイルが記憶された 2 次記憶装置が属する計算機ノードを、前記プロセスの移送先に決定する処理とを具備したことを特徴とする分散ファイル装置のプロセスマイグレーション方法。

10

【請求項 6】

ネットワーク上に接続された複数の計算機ノードのうちの所定の計算機ノードが実行するプロセスに従って退避ファイルを作成するステップと、

プロセス毎に作成された前記退避ファイルを前記ネットワーク上の任意の計算機ノードの 2 次記憶装置に記憶させると共に、前記任意の計算機ノードとは異なる他の計算機ノードの 2 次記憶装置に前記退避ファイルを複製した複製退避ファイルを記憶させるステップと、

20

実行中のプロセスを、前記実行中のプロセスに従って作成された複製退避ファイルが記憶された 2 次記憶装置が属する計算機ノードに対して移送するステップとを具備したことを特徴とする分散ファイル装置のプロセスマイグレーション方法。

【請求項 7】

2 次記憶装置を備えた複数のコンピュータ装置とネットワークで接続される、2 次記憶装置を備えたコンピュータ装置であって、

実行中のプロセスに従って退避ファイルを作成する退避ファイル作成手段と、前記退避ファイル作成手段で作成された前記退避ファイルを自装置の 2 次記憶装置に記憶させると共に、前記ネットワークと接続される複数のコンピュータ装置の何れかのコンピュータ装置の 2 次記憶装置に記憶させるために前記退避ファイルを複製した複製退避ファイルを送信する分散ファイル手段と、

30

実行中のプロセスを他のコンピュータ装置へ移送する場合に、前記実行中のプロセスに従って作成された退避ファイルの複製退避ファイルを前記分散ファイル手段で送信した前記他のコンピュータ装置を、前記プロセスの移送先に決定する決定手段とを具備したことを特徴とするコンピュータ装置。

【請求項 8】

前記分散ファイル手段との間でデータの授受を行って、前記退避ファイルに対応する複製退避ファイルの送信先を示す管理テーブルを生成・管理する管理手段を更に備え、前記決定手段は、前記管理テーブルを用いて、プロセスの移送先を決定することを特徴とする請求項 7 に記載のコンピュータ装置。

40

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、複数の計算機ノードで構成される計算機クラスタシステムに採用される分散ファイル装置及びそのプロセスマイグレーション方法並びにコンピュータ装置に関する。

【0002】

【従来の技術】

50

従来、計算機においては、プロセスを実行するために、プロセスの実行プログラムに基づくデータを自計算機の主記憶装置上に展開する必要がある。しかし、計算機が実装している主記憶装置の容量は有限であることから、同時に複数のプロセスを実行しようとする場合等においては、主記憶装置の容量不足が発生する。そこで、一般的には、主記憶装置の容量以上の記憶空間を使用可能にして、実メモリの制約を越えたプロセスを実行可能にする仮想記憶が採用される。

【 0 0 0 3 】

仮想記憶技術では、プログラムが特定するデータの位置を仮想アドレスによって表し、例えばオペレーションシステム（OS）によって、この仮想アドレスを主記憶装置上の位置を表す実アドレスに変換する。従って、異なる複数のプログラムで同一の仮想アドレスが使用されている場合であっても、各プログラムの同一の仮想アドレスを相互に異なる実アドレスに変換することで、これらの複数のプログラムを同時に実行することが可能となる。また、仮想アドレスが有限の実アドレスに制限されることもない。

10

【 0 0 0 4 】

仮想記憶では、プロセスの実行プログラムのうちプロセス実行中に必要なデータのみを主記憶装置に読み込み、主記憶の容量が不足してくると、不要になったデータは主記憶装置から除去する。この結果、主記憶装置上に無駄なスペースが無くなり、主記憶装置を最大限有効に使用することができる。

【 0 0 0 5 】

また、仮想記憶は、主記憶装置上で不要になったデータのうち、プロセス実行中に更新されたデータについては、2次記憶装置の退避スペースに書き出して退避する処理（ページアウトともいう）を行う。

20

【 0 0 0 6 】

主記憶装置上の未変更のデータは退避されない。除去して主記憶装置上にないデータが再度必要になると、退避したものについては退避スペースからデータを主記憶装置に読み込み（ページインともいう）、単に除去しただけのデータは実行ファイルからデータを主記憶装置に読みこむ。なお、仮想記憶によって使用される退避スペースは、通常、自計算機専用のスペースである2次記憶装置上に確保するようになっている。

【 0 0 0 7 】

ところで、従来、複数台のコンピュータを連携して使用して、1台のコンピュータと同様の使用を可能にする計算機クラスタシステムが採用されることがある。計算機クラスタシステムは、比較的安価なコンピュータを用いた場合でも、高度な業務処理が可能であり、しかも、システムの拡張が極めて容易である。

30

【 0 0 0 8 】

このようなクラスタシステムにおいて、実行中のプロセスを他計算機に移動させ、移動した計算機上でプロセスの実行を継続させるプロセスマイグレーションが採用されることがある。

【 0 0 0 9 】

このプロセスマイグレーションによれば、クラスタシステム内の任意の計算機が故障した場合において、故障した計算機で実行していたプロセスをクラスタシステム内の他の正常な計算機に移すことで、プロセス実行をそのまま継続することができる。これにより、クラスタシステム全体では安定した動作を続けることが可能となる。また、実行中のプロセスをクラスタシステム内で負荷が小さい他の計算機に移動して、プロセス実行を継続することで、クラスタシステム内の負荷分散を可能にすることもできる。

40

【 0 0 1 0 】

このようなプロセスマイグレーションを実現するには、

(1) プロセスの移動元となる計算機において移動させるプロセスのプロセス状態を取得し、このプロセス状態をプロセスの移動先となる計算機に送る。

【 0 0 1 1 】

(2) プロセス移動先の計算機は、(1)でプロセス移動元から送られたプロセス状態を

50

自計算機上に復元する。

【0012】

(3) プロセス移動先の計算機でプロセスの実行を再開する。

【0013】

という過程を経る。

【0014】

この場合において、プロセス状態は、プロセスを実行するために使用していた主記憶装置の全内容(退避スペースの内容を含む)、CPUのレジスタ値を伝達することによって、移動させることができる。

【0015】

【発明が解決しようとする課題】

プロセス状態の移動に際して伝送する情報のうち、主記憶装置の全内容(退避スペースの内容を含む)及びCPUのレジスタ値については、夫々プロセス移動先の計算機内の主記憶装置及びレジスタに格納する。プロセスマイグレーションに要する時間は、プロセス状態の伝送に必要な情報の伝送に要する時間の制約を受ける。

【0016】

そこで、Fred Douglass及びJohn Ousterhoutは、文献1(「Transparent Process Migration: Design Alternatives and the Sprite Implementation」)において、プロセスマイグレーション時に、主記憶装置の全内容をプロセス移動先の計算機内の主記憶装置に全て伝送する代わりに、プロセス実行中に更新されたページのみを、プロセス毎の退避スペースとして用意した退避ファイルにページアウトする。そして、プロセス状態の伝送に必要な情報の伝送量を低減して、プロセスマイグレーションに要する時間を短縮した技術を提案している。即ち、この提案においては、退避ファイルを転送元及び転送先の計算機が属するネットワーク上のファイルサーバ内の2次記憶装置に記憶させる。

【0017】

この場合には、退避ファイルは、転送元及び転送先の計算機によってアクセス可能である。退避ファイルはプロセス毎に作成されるので、プロセスの移動にともなって、移動するプロセスに対応した退避ファイルの使用権を転送元から転送先に移動させればよく、退避ファイルの転送は不要である。

【0018】

ところで、計算機システムにおいて、ファイルアクセスは、頻繁に発生する動作であり、システム全体の性能のボトルネックになりやすい項目である。一般的に、計算機ノード間の通信帯域は単一の計算機内のデバイス間の通信帯域に比べて狭い。従って、計算機ノード間の通信を伴う処理は、計算機ノード内で閉じた処理に比べて極めて低速である。このため、ファイルアクセスのたびに発生する計算機ノード間のデータ通信量の多さはシステム全体の性能を低下させる要因となる。

【0019】

ところが、上述した文献1では、退避ファイルがネットワーク上の他のノードに存在することから、ページイン毎にファイルサーバ内の退避ファイルをプロセス実行中の計算機に転送する必要があり、結果としてページインに長時間を要してしまい、プロセスが低速になってしまうという欠点がある。また、プロセスマイグレーション直後には、退避ファイル内のデータの多くを移動先の計算機内の主記憶装置に転送することが多い。このためプロセス移動先において、実際にプロセスが稼働するまでに比較的長時間を要してしまうという問題もあった。

【0020】

本発明はかかる問題点に鑑みてなされたものであって、ページインに要する時間を短縮すると共に、プロセスマイグレーションに要する時間を短縮することができる分散ファイル装置及びそのプロセスマイグレーション方法並びにコンピュータ装置を提供することを目的とする。

【0021】

10

20

30

40

50

【課題を解決するための手段】

本発明の請求項 1 に係る分散ファイル装置は、ネットワーク上に接続された複数の計算機ノードに夫々設けられる 2 次記憶装置と、前記各計算機ノードに夫々設けられ、自計算機ノードの計算機が実行するプロセスに従って退避ファイルを作成する退避ファイル作成手段と、前記退避ファイル作成手段で該プロセス毎に作成された前記退避ファイルを前記ネットワーク上の任意の計算機ノードの 2 次記憶装置に記憶させると共に、前記任意の計算機ノードの 2 次記憶装置とは異なる他の計算機ノードの 2 次記憶装置に前記退避ファイルを複製した複製退避ファイルを記憶させる記憶制御手段と、実行中のプロセスを他の計算機ノードに移送する場合に、前記実行中のプロセスに従って作成された退避ファイルの複製退避ファイルが記憶された 2 次記憶装置が属する計算機ノードを、前記プロセスの移送先に決定する決定手段とを具備したものであり、

本発明の請求項 7 に係るコンピュータ装置は、2 次記憶装置を備えた複数のコンピュータ装置とネットワークで接続される、2 次記憶装置を備えたコンピュータ装置であって、実行中のプロセスに従って退避ファイルを作成する退避ファイル作成手段と、前記退避ファイル作成手段で作成された前記退避ファイルを自装置の 2 次記憶装置に記憶させると共に、前記ネットワークと接続される複数のコンピュータ装置の何れかのコンピュータ装置の 2 次記憶装置に記憶させるために前記退避ファイルを複製した複製退避ファイルを送信する分散ファイル手段と、実行中のプロセスを他のコンピュータ装置へ移送する場合に、前記実行中のプロセスに従って作成された退避ファイルの複製退避ファイルを前記分散ファイル手段で送信した前記他のコンピュータ装置を、前記プロセスの移送先に決定する決定手段とを具備したものである。

【0022】

本発明の請求項 1 において、ネットワーク上に接続された複数の計算機ノードには夫々 2 次記憶装置が設けられる。退避ファイル作成手段は、自計算機ノードの計算機が実行するプロセスに従って退避ファイルを作成する。この退避ファイルは、記憶制御手段によって、ネットワーク上の任意の計算機ノードの 2 次記憶装置に記憶される。更に、記憶制御手段は、退避ファイルの複製ファイルを他の計算機ノードの 2 次記憶装置に記憶させる。決定手段は、実行中のプロセスを他の計算機ノードに移送する場合には、移送先として退避ファイルの複製ファイルを記憶した 2 次記憶装置が属する計算機ノードを決定する。これにより、プロセスマイグレーション時に、退避ファイルの転送は不要である。また、移送先の計算機ノードにおいては、退避ファイルは自計算機ノードの 2 次記憶装置から読出せばよい。

【0023】

本発明の請求項 7 に係るコンピュータ装置は、2 次記憶装置を備えた複数のコンピュータ装置とネットワークで接続される、2 次記憶装置を備えたコンピュータ装置であって、実行中のプロセスに従って退避ファイルを作成する退避ファイル作成手段と、前記退避ファイル作成手段で作成された前記退避ファイルを自装置の 2 次記憶装置に記憶させると共に、前記ネットワークと接続される複数のコンピュータ装置の何れかのコンピュータ装置の 2 次記憶装置に記憶させるために前記退避ファイルを複製した複製退避ファイルを送信する分散ファイル手段と、実行中のプロセスを他のコンピュータ装置へ移送する場合に、前記実行中のプロセスに従って作成された退避ファイルの複製退避ファイルを前記分散ファイル手段で送信した前記他のコンピュータ装置を、前記プロセスの移送先に決定する決定手段とを具備したものである。

【0024】

本発明の請求項 7 において、2 次記憶装置を備えた複数のコンピュータ装置とはネットワークを介して接続される。退避ファイル作成手段は、実行中のプロセスに従って退避ファイルを作成する。この退避ファイルは、記憶制御手段によって、ネットワーク上の複数のコンピュータ装置のいずれかのコンピュータ装置の 2 次記憶装置に記憶される。更に、記憶制御手段は、退避ファイルの複製退避ファイルを他のコンピュータ装置の 2 次記憶装置に記憶させる。決定手段は、実行中のプロセスを他のコンピュータ装置に移送する場合に

は、移送先として退避ファイルの複製ファイルを記憶した2次記憶装置を備えたコンピュータ装置を決定する。これにより、プロセスマイグレーション時に、退避ファイルの転送は不要である。また、移送先のコンピュータ装置においては、退避ファイルは自コンピュータ装置の2次記憶装置から読出せばよい。

【0025】

なお、装置に係る本発明は、プロセスマイグレーションの方法に係る発明としても成立する。

【0026】

【発明の実施の形態】

以下、図面を参照して本発明の実施の形態について詳細に説明する。図1は本発明の一実施の形態に係る分散ファイル装置を示すブロック図である。本実施の形態は本発明を計算機クラスタシステムに適用した例である。

【0027】

計算機クラスタシステムでは、ファイルを計算機ノード間で分散して保持することがある。このようなシステムにおいて、プロセスがどの計算機ノードで動作していても、全てのファイルへのアクセスを同様に可能とするために、分散ファイルシステムが用いられる。

【0028】

分散ファイルシステムを使用すると、全ての計算機ノードの全プロセスが、クラスタシステム内の計算機ノードに分散して格納されているファイルを一意に指定することができる。分散ファイルシステムの代表例としては、AFS (Andrew File System) がある。分散ファイルシステムは、ファイルの実体であるマスターファイルをシステム内のいずれかの計算機ノード上の記憶装置に格納し、ファイルがどの計算機ノードに格納されているかという情報をシステム内のデータベースに登録する。ファイルを使用する場合には、システム内のデータベースからマスターファイルが実際に格納されている計算機ノードを検索し、この検索結果を利用することで、いずれの計算機ノードにおいてもマスターファイルの読み出しを可能にしている。

【0029】

この場合において、システム内のファイルを保護して、システムの信頼性を向上させるために、ファイルの多重化が行われる。即ち、分散ファイルシステムを用いた計算機クラスタシステムにおいては、ファイルの実体であるマスターファイルと同一の内容を持った複製ファイルを作成し、マスターファイルが格納されている計算機ノードとは別の計算機ノードに複製ファイルを格納する多重化を採用する。このような高信頼型の分散ファイルシステムでは、ファイルに対する更新は、マスターファイルに行うと同時に逐一複製ファイルに対しても行い、ファイルの多重度を維持する。この方法によれば、マスターファイル及び複製ファイルのいずれか一方が壊れた場合でも、ファイルの内容を他方から復元することができる。

【0030】

本実施の形態においては、プロセスマイグレーションに際して、分散ファイルシステムによって作成される退避ファイルの複製ファイルを格納する計算機ノードに、プロセスマイグレーション先を設定することで、プロセスマイグレーションに要する時間を短縮すると共に、ページインに要する時間を短縮するようになっている。

【0031】

図1において、LAN (ローカルエリアネットワーク) 等の所定のネットワーク13上には複数の計算機A, B, ... が接続されている。なお、図1では2台の計算機A, Bのみを示している。各計算機A, B, ... は略同一構成であり、各計算機A, B, ... 及び後述する各2次記憶装置1a, 1b, ... によって計算機クラスタシステムの各計算機ノードが構成されている。

【0032】

各計算機A, B, ... には、夫々、プロセスマイグレーション実現部6a, 6b, ... (以下、代表してプロセスマイグレーション実現部6という)、主記憶装置7a, 7b, ... (以

10

20

30

40

50

下、代表して主記憶装置 7 という)、仮想記憶管理部 8 a , 8 b , ... (以下、代表して仮想記憶管理部 8 という)及び分散ファイルシステム 2 a , 2 b , ... (以下、代表して分散ファイルシステム 2 という)が含まれると共に、各計算機ノードは、ネットワーク 1 3 よりも高速な通信が可能なローカルの 2 次記憶装置 1 a , 1 b , ... (以下、代表して 2 次記憶装置 1 という)が接続されている。なお、2 次記憶装置 1 としては、各計算機内の内部バスによって接続されたものであってもよく、他の通信ケーブル等によって接続されたものでもよい。

【 0 0 3 3 】

2 次記憶装置 1 a , 1 b , ... は、計算機クラスタシステム内の各計算機ノードによってアクセス可能であり、2 次記憶装置 1 a , 1 b , ... に格納されるファイルは、分散ファイルシステム 2 によって一元管理されるようになっている。

10

【 0 0 3 4 】

図 2 は分散ファイルシステム 2 のファイル管理に用いるファイル管理テーブルを示す説明図である。

【 0 0 3 5 】

図 2 に示すように、各ファイルは、ファイル ID によって管理され、各ファイル ID 毎に、1 つのマスターファイルと複数の複製ファイルとが設定される。マスターファイル及び複数の複製ファイルは、夫々ネットワーク内の各計算機ノードに分散して記憶されるようになっており、各ファイル毎に、保存先の計算機ノードが決定されるようになっている。

【 0 0 3 6 】

即ち、各ファイルは、ファイル ID によって特定され、各ファイル ID 毎にマスターファイルが格納される計算機ノードの ID (マスターノード ID) と 1 つ以上の複製ファイルが夫々格納される計算機ノードの ID (レプリカノード ID) が対応付けられる。ファイル管理テーブルは、マスターファイルと 1 つ以上の複製ファイルとのネットワーク上の位置を記述している。

20

【 0 0 3 7 】

プロセスマイグレーション実現部 6 は、実行中のプロセスを他計算機に移し実行を継続させるための処理を行う部分である。即ち、自計算機で実行中のプロセスを他計算機に移送する場合は、プロセスを一旦停止してその状態を保存し、移送先の計算機のプロセスマイグレーション実現部 6 へ送出する。また、他計算機で実行されていたプロセスを自計算機に移送して実行を継続させる場合には、移送元の計算機のプロセスマイグレーション実現部 6 から受け取ったプロセス状態を自計算機で復元する処理を行う。

30

【 0 0 3 8 】

主記憶装置 7 は、プロセスの実行に必要なデータを展開するメモリ領域である。仮想記憶管理部 8 は、仮想記憶管理のための処理を行う部分である。即ち、仮想記憶管理部 8 は、仮想アドレスから実アドレスへの変換や、主記憶装置 7 上の領域のうち、プロセス実行中に書きかえられた領域のみを退避ファイルとして、自ノードの 2 次記憶装置 1 にページアウトし、退避ファイル上の必要な領域のみをその領域が必要とされた場合に自計算機の主記憶装置 7 にページインさせるための処理を行う。なお、主記憶装置 7 の退避スペースとして利用する退避ファイルは、プロセス毎に固有のファイルである。

40

【 0 0 3 9 】

なお、退避ファイルは、他ノードの 2 次記憶装置 1 に設けてもよいが、高速なページインを可能にするためには、自ノードの 2 次記憶装置 1 に退避ファイルを記憶させた方がよい。

【 0 0 4 0 】

各仮想記憶管理部 8 a , 8 b , ... は、夫々、メモリ管理テーブルを用いることによって仮想記憶を実現する。図 3 は仮想記憶管理部 8 が記憶保持しているメモリ管理テーブルを示す説明図である。

【 0 0 4 1 】

図 3 に示すように、メモリ管理テーブルは、仮想アドレス、実アドレス及び退避ファイル

50

のオフセットの関係を記述したものであり、プロセス毎に設けられる。仮想アドレスは、プログラムが特定するデータの位置を示すアドレスであり、この仮想アドレスは、実際の主記憶装置 7 上の位置を表す実アドレスに変換される。メモリ管理テーブルは、この場合の仮想アドレスと実アドレスとの間の対応を示している。

【 0 0 4 2 】

また、ページアウトが発生した場合には、仮想アドレスによって与えられるデータが退避ファイル上のいずれの位置のデータであるかの対応を取ることができる。

【 0 0 4 3 】

本実施の形態においては、上述したように、プロセス毎に退避ファイルを設定することができる。プロセス毎に退避ファイルを設定した仮想記憶を実現するために、仮想記憶管理部 8 a , 8 b , ... は夫々仮想記憶管理テーブル 9 a , 9 b , ... (以下、代表して仮想記憶管理テーブル 9 という) を有している。

10

【 0 0 4 4 】

図 4 は仮想記憶管理テーブル 9 の内容を示す説明図である。

【 0 0 4 5 】

仮想記憶管理テーブル 9 は、各プロセス毎に退避ファイルとメモリ管理テーブルとの対応を示すものである。

【 0 0 4 6 】

各プロセスは、プロセス ID によって特定され、退避ファイルは退避ファイルの ID (退避ファイル ID) によって特定される。仮想記憶管理テーブル 9 によって、プロセスと、そのプロセスに用いる退避ファイルとそのプロセスに利用する仮想記憶のためのメモリ管理テーブルとの対応が記述される。

20

【 0 0 4 7 】

仮想記憶管理部 8 は、プロセスマイグレーション実現部 6 がプロセスを他計算機に移送する際に、移送するプロセスの仮想記憶管理テーブル 9 をプロセス移送先の仮想記憶管理部 8 に送る。そして、プロセス移送先の仮想記憶管理部 8 は、受け取ったプロセスの仮想記憶管理テーブル 9 を用いてプロセスの仮想記憶管理を引き継いで行うようになっている。

【 0 0 4 8 】

本実施の形態においては、仮想記憶管理部 8 は、分散ファイルシステム 2 との間でデータの授受を行って、退避ファイルについての複製ファイル(複製退避ファイル)の位置の情報を得て、仮想記憶管理テーブル 9 に書き込むようになっている。

30

【 0 0 4 9 】

そして、本実施の形態においては、プロセスマイグレーション実現部 6 は、仮想記憶管理部 8 に記憶されている仮想記憶管理テーブル 9 と分散ファイルシステム 2 が用いるファイル管理テーブルとを用いて、複製退避ファイルが格納されている 2 次記憶装置 1 を有する計算機ノードの情報を得て、この計算機ノードをプロセスマイグレーションの移送先に設定するようになっている。

【 0 0 5 0 】

次に、このように構成された実施の形態の動作について図 5 及び図 6 のフローチャートを参照して説明する。図 5 は計算機 A で実行中のプロセスを計算機 B に移送する場合の、計算機 A におけるプロセスマイグレーション実現部 6 a 及び仮想記憶管理部 8 a の処理手順を示すフローチャートであり、図 6 は計算機 A で実行中のプロセスを計算機 B に移送する場合の、計算機 B におけるプロセスマイグレーション実現部 6 b 及び仮想記憶管理部 8 b の処理手順を示すフローチャートである。

40

【 0 0 5 1 】

いま、図 1 の計算機 A においてプロセス 1 0 a , 1 1 a を実行中であり、また、計算機 B においてプロセス 1 2 b を実行中であるものとする。

【 0 0 5 2 】

即ち、計算機 A の仮想記憶管理部 8 a は、プロセス 1 0 a の実行に伴って、主記憶装置 7 a にプロセス 1 0 a の実行に必要なデータを展開する。また、仮想記憶管理部 8 a は、プ

50

プロセス10aが書き換えた主記憶装置7aの領域をプロセス10a用の退避ファイル4として、分散ファイルシステム2aを介して、2次記憶装置1aに転送して格納させる(ステップS1)。

【0053】

仮想記憶管理部8は、各プロセス毎にプロセスIDを割当て、各プロセスID毎に退避ファイルを作成して、プロセスと退避ファイルとの対応を仮想記憶管理テーブル9に記述する。例えば、図1の例では、仮想記憶管理部8aは、プロセス10aについてプロセスIDを割当て、このプロセスIDについて退避ファイルを作成する。

【0054】

一方、分散ファイルシステム2は、各ファイル毎にファイルIDを割当て、各ファイルIDで示されるファイルのマスターファイルを記憶させた計算機ノードのIDとその複製ファイルを記憶させた計算機ノードのIDとの関連を、ファイル管理テーブルに記述している。本実施の形態においては、退避ファイルについても、分散ファイルシステム2によって管理される。

10

【0055】

仮想記憶管理部8aは、分散ファイルシステム2aとの間で通信を行って、作成した退避ファイルのファイルIDを取得し、プロセスID、退避ファイルID及びメモリ管理テーブルからなる仮想記憶管理テーブル9aを作成する。(ステップS2)

なお、分散ファイルシステム2aは、退避ファイルを計算機クラスタシステム内のいずれのノードの計算機に接続された2次記憶装置に記憶させることも可能であるが、退避ファイルについては、プロセスを実行中の自ノードの2次記憶装置に記憶させた方が、処理を高速化させることが可能である。

20

【0056】

また、分散ファイルシステム2aは、他の計算機ノードの分散ファイルシステム2b,...と通信を行って、2次記憶装置1aに記憶させた退避ファイルをマスターファイルとし、このマスターファイルの複製である複製退避ファイル5を他の計算機ノードに接続された2次記憶装置に記憶させるようになっている。

【0057】

これらの退避ファイルのマスターファイル及び複製ファイルについても、分散ファイルシステム2によって管理される。例えば、図1の例では、分散ファイルシステム2aによって、プロセス10aの実行に伴う退避ファイルにファイルIDが付され、このファイルIDに関連付けて、マスターファイル4が格納されている計算機Aが属する計算機ノードのIDと、退避ファイルの複製ファイル5が格納されている計算機Bが属する計算機ノードのIDとが記述される。

30

【0058】

ここで、計算機Aで実行中のプロセス10aを他の計算機に移送するものとする。この場合には、プロセスマイグレーション実現部6aは、仮想記憶管理テーブル9aの内容を読み出して、プロセス移送先の計算機ノードを決定する(ステップS3)。即ち、プロセスマイグレーション実現部6aは、仮想記憶管理部8aに問い合わせを行って、仮想記憶管理テーブル9aの記述から、移送しようとするプロセス10aについての退避ファイルのファイルIDを取得する。そして、プロセスマイグレーション実現部6aは、取得したファイルIDを元に、分散ファイルシステム2aに問い合わせを行って、ファイル管理テーブルの記述から、プロセス10aについての退避ファイルの複製ファイル5が記憶されている計算機ノードの情報を得る。

40

【0059】

本実施の形態においては、プロセスマイグレーション実現部6aは、退避ファイルの複製ファイル5が作成されている2次記憶装置1が接続された計算機ノードをプロセス移送先に決定する。いま、プロセス10aの退避ファイルの複製ファイル5が計算機Bに接続された2次記憶装置1bに記憶されているものとする。この場合には、プロセスマイグレーション実現部6aによって、プロセス10aの移送先として計算機Bが選択される。

50

【 0 0 6 0 】

次に、仮想記憶管理部 8 a は、プロセスの移送先として選択された計算機 B の仮想記憶管理部 8 b にプロセス 1 0 a の仮想記憶管理テーブル 9 a を送る (ステップ S 4)。仮想記憶管理部 8 b は、仮想記憶管理テーブル 9 a の内容を仮想記憶管理テーブル 9 b に書き込む。

【 0 0 6 1 】

次に、プロセスマイグレーション実現部 6 a は、プロセス 1 0 a のプロセス状態を取得して計算機 B のプロセスマイグレーション実現部 6 b に出力する (ステップ S 5)。

【 0 0 6 2 】

一方、計算機 B においては、図 6 のステップ S 11 において、プロセスマイグレーション実現部 6 b が、計算機 A のプロセスマイグレーション実現部 6 a から受け取ったプロセス 1 0 a のプロセス状態を復元する。そして、計算機 B は、プロセス 1 0 a の実行が再開されると、仮想記憶管理部 8 b によって、計算機 A から受け取ったプロセス 1 0 a の仮想記憶管理テーブル 9 a (仮想記憶管理テーブル 9 b) を参照してプロセス 1 0 a の実行プログラムまたは計算機 A の仮想記憶管理部が退避したプロセス 1 0 a 用の退避ファイルから必要な領域のみを必要な時に主記憶装置 7 b に読み込む (ステップ S 12)。

10

【 0 0 6 3 】

本実施の形態においては、仮想記憶管理部 8 b は、分散ファイルシステム 2 b を介して、自ノードに接続された 2 次記憶装置 1 b に記憶されている退避ファイルの複製ファイル 5 を、退避ファイルとして読み込む。

20

【 0 0 6 4 】

即ち、本実施の形態においては、ページアウト及びページインは、プロセスマイグレーションの前後において、常に、自ノードに接続された 2 次記憶装置 1 に対して行われる。従って、高速なページアウト及びページインが可能である。しかも、プロセスマイグレーションの移送先として、退避ファイルの複製ファイルを保持する 2 次記憶装置が接続された計算機ノードを選択しており、退避ファイルについては、移送の必要がなく、しかも、プロセスマイグレーション後におけるページインを自ノードの 2 次記憶装置から行うことができ、高速なページインが可能である。

【 0 0 6 5 】

これにより、本実施の形態においては、プロセスマイグレーションに要する時間を短縮することができ、しかも、移送先の計算機ノードにおいて、ページインに要する時間を短縮することができる。

30

【 0 0 6 6 】

各計算機ノードのプロセスマイグレーション実現部 6 及び分散ファイルシステム 2 は、相互に協働してプロセスマイグレーション及び分散ファイルシステムを提供するもので、ネットワーク上のいずれの計算機ノードによって制御可能であり、また、ネットワーク上のいずれかの計算機ノードのみに設けて集中制御するように構成してもよい。

【 0 0 6 7 】**【 発明の効果 】**

以上説明したように本発明によれば、ページインに要する時間を短縮すると共に、プロセスマイグレーションに要する時間を短縮することができるという効果を有する。

40

【 図面の簡単な説明 】

【 図 1 】 本発明の一実施の形態に係る分散ファイル装置を示すブロック図。

【 図 2 】 分散ファイルシステム 2 のファイル管理に用いるファイル管理テーブルを示す説明図。

【 図 3 】 仮想記憶管理部 8 が記憶保持しているメモリ管理テーブルを示す説明図。

【 図 4 】 仮想記憶管理テーブル 9 の内容を示す説明図。

【 図 5 】 実施の形態の動作を説明するためのフローチャート。

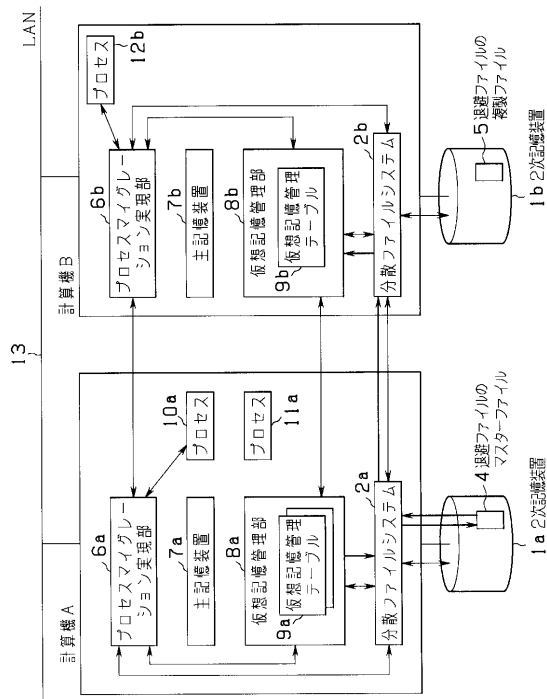
【 図 6 】 実施の形態の動作を説明するためのフローチャート。

【 符号の説明 】

50

1 a , 1 b ... 2 次記憶装置、 2 a , 2 b ... 分散ファイルシステム、 4 ... 退避ファイルのマスターファイル、 5 ... 退避ファイルの複製ファイル、 6 a , 6 b ... プロセスマイグレーション実現部、 7 a , 8 b ... 主記憶装置、 8 a , 8 b ... 仮想記憶管理部、 9 a , 9 b ... 仮想記憶管理テーブル、 1 0 a , 1 1 a , 1 2 b ... プロセス

【 図 1 】



【 図 2 】

ファイルID	マスター ノードID	レプリカ ノードID	レプリカ ノードID	...
...				

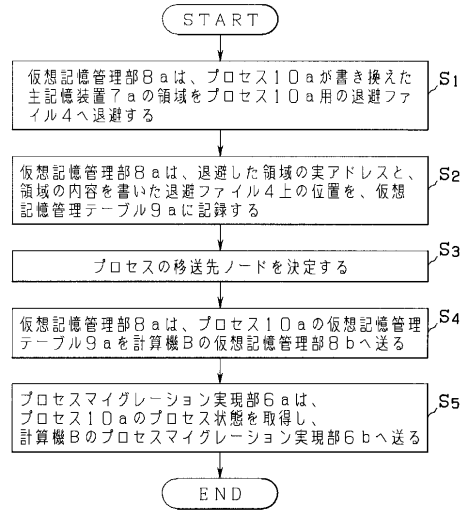
【 図 3 】

仮想 アドレス	実アドレス	退避ファイル のオフセット
...		

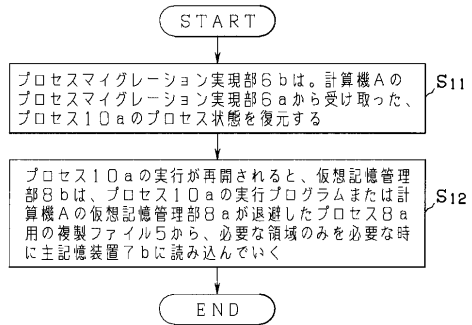
【 図 4 】

プロセスID	退避 ファイルID	メモリ 管理テーブル
⋮		

【 図 5 】



【 図 6 】



フロントページの続き

(72)発明者 矢野 浩邦

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 林 拓也

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

審査官 殿川 雅也

(56)参考文献 DOUGLIS, F., et al., Transparent Process Migration: Design Alternatives and the Sprite Implementation, SOFTWARE-PRACTICE AND EXPERIENCE, John Wiley & Sons, Ltd., 1991年8月, Vol. 21, No. 8, pp. 757 - 785

MALKAWI, M., Process Migration in Virtual Memory Multicomputer Systems, System Sciences, 1993, Proceedings of the 26th Hawaii International Conference on, IEEE, 1993年1月8日, Vol. 2, pp. 90 - 98

DE PAOLI, D., et al., A Copy on Reference Process Migration in RHODOS, Algorithms and Architectures for Parallel Proceedings, 1996. ICAPP'96. 2nd Int. Conference on, IEEE, 1996年6月13日, pp. 100 - 107

CHANCHIO, K., et al., SNOW: Software Systems for Process Migration in High-Performance, Heterogeneous Distributed, Parallel Processing Workshops, 2002. Proceedings. Int. Conference on, IEEE, 2002年8月21日, pp. 589 - 596, タイトルの続き: Environments

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54

G06F 15/16 - 15/177

G06F 11/20

G06F 12/00