

(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2010-97533
(P2010-97533A)

(43) 公開日 平成22年4月30日(2010.4.30)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 9/46 (2006.01)	G06F 9/46 350	5B082
G06F 12/00 (2006.01)	G06F 12/00 501B	
	G06F 12/00 545A	

審査請求 未請求 請求項の数 11 O L (全 51 頁)

(21) 出願番号 特願2008-269539 (P2008-269539)
(22) 出願日 平成20年10月20日(2008.10.20)

(71) 出願人 000005108
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 100075513
弁理士 後藤 政喜
(74) 代理人 100114236
弁理士 藤井 正弘
(74) 代理人 100120260
弁理士 飯田 雅昭
(72) 発明者 鈴木 友彦
神奈川県小田原市中里322番2号 株式会社日立製作所SANソリューション事業部内
Fターム(参考) 5B082 CA11 CA19 HA08

(54) 【発明の名称】パーティションで区切られた計算機システムにおけるアプリケーション移動及び消費電力の最適化

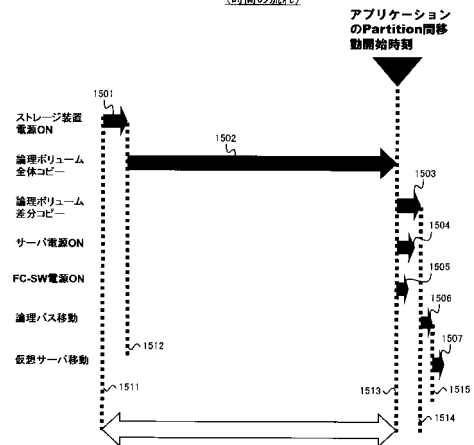
(57) 【要約】

【課題】論理分割された計算機システムにおいて、アプリケーション移動時の消費電力を削減する。

【解決手段】アプリケーションの移動元の論理ボリュームを含むストレージ装置は、前記移動元の論理ボリュームに格納されたデータを、アプリケーションの移動先の論理ボリュームにコピーし、前記コピーが開始された後に前記移動元の論理ボリュームに書き込まれたデータを前記移動元の論理ボリュームに格納せずに差分データとして保持し、管理計算機は、前記移動元の論理ボリュームに格納されたデータのコピーが終了すると、前記差分データのコピーを開始し、前記移動元の論理ボリュームに格納されたデータのコピーが終了してから前記差分データのコピーが終了するまでの間に、前記アプリケーションの移動先の計算機の電源を投入する。

【選択図】図15

Partitionの電源ONとアプリケーションの移動を組み合わせた動作の例 (時間の流れ)



【特許請求の範囲】

【請求項 1】

一つ以上の計算機と、ネットワークを介して前記一つ以上の計算機に接続される一つ以上のストレージ装置と、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される管理計算機と、を備える計算機システムであって、

前記各計算機は、ハードウェア資源として、前記ネットワークに接続される第 1 インタフェースと、前記第 1 インタフェースに接続される第 1 プロセッサと、前記第 1 プロセッサに接続される第 1 メモリと、前記管理計算機からの要求に従って前記ハードウェア資源への電源の投入及び遮断を制御する電源制御部と、を備え、

前記各計算機は、さらに、前記一つ以上の計算機のハードウェア資源に基づいて複数の仮想領域を提供する仮想化部を備え、

前記複数の仮想領域は、第 1 仮想領域及び第 2 仮想領域を含み、

前記第 1 仮想領域は、アプリケーションプログラムを実行する仮想計算機として動作し、

前記各ストレージ装置は、前記計算機によって書き込まれたデータを格納する記憶領域を提供する記憶媒体と、前記記憶媒体へのデータの入出力を制御するコントローラと、前記管理計算機からの要求に従って前記各ストレージ装置の電源の投入及び遮断を制御する第 2 電源制御部と、を備え、

前記一つ以上のストレージ装置の前記コントローラは、前記一つ以上のストレージ装置の前記記憶領域を複数の論理ボリュームとして前記計算機に提供し、

前記複数の論理ボリュームは、第 1 論理ボリューム及び第 2 論理ボリュームを含み、

前記第 1 論理ボリュームには、前記仮想計算機によって書き込まれたデータが格納され、

前記管理計算機は、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される第 2 インタフェースと、前記第 2 インタフェースに接続される第 2 プロセッサと、前記第 2 プロセッサに接続される第 2 メモリと、を備え、

前記管理計算機は、

前記第 1 論理ボリュームを含む前記ストレージ装置に、前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求を送信し、

前記第 1 論理ボリュームを含む前記ストレージ装置は、

前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求を受信した後に前記仮想計算機から前記第 1 論理ボリュームへのデータの書き込み要求を受信すると、書き込みを要求されたデータを前記第 1 論理ボリュームに書き込まずに、差分データとして保持し、

前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求に従って、前記第 1 論理ボリュームに格納されたデータを読み出して、前記第 2 論理ボリュームを含む前記ストレージ装置に送信し、

前記管理計算機は、

前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーの終了を検出すると、前記保持された差分データを前記第 2 論理ボリュームにコピーする要求を、前記第 1 論理ボリュームを含む前記ストレージ装置に送信し、

前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーの終了を検出してから、前記差分データの前記第 2 論理ボリュームへのコピーが終了するまでの間に、前記一つ以上の計算機のハードウェア資源のうち、前記第 2 仮想領域に割り当てられたハードウェア資源の電源を投入する要求を、前記第 2 仮想領域に割り当てられたハードウェア資源を含む前記計算機に送信し、

前記差分データの前記第 2 論理ボリュームへのコピーが終了した後、前記仮想計算機を前記第 2 仮想領域に移動する要求を送信し、

前記仮想計算機は、前記第 2 仮想領域に移動した後、前記第 2 論理ボリュームへのデータ入出力を実行することを特徴とする計算機システム。

10

20

30

40

50

【請求項 2】

前記第 1 論理ボリュームを含む前記ストレージ装置は、前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求に従う前記データの読み出し及び送信が終了すると、終了メッセージを前記管理計算機に送信し、

前記管理計算機は、前記終了メッセージを受信したときに、前記コピーの終了を検出することを特徴とする請求項 1 に記載の計算機システム。

【請求項 3】

前記管理計算機は、前記第 1 論理ボリュームの容量、及び、前記第 1 論理ボリュームから前記第 2 論理ボリュームへのデータの転送速度に基づいて、前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーに要する時間を算出し、

前記算出された時間が経過したときに、前記コピーの終了を検出することを特徴とする請求項 1 に記載の計算機システム。

【請求項 4】

前記ネットワークは、前記一つ以上の計算機と前記一つ以上のストレージ装置との間のデータ転送を中継する複数のスイッチ装置を備え、

前記各スイッチ装置は、複数のポートと、前記管理計算機に接続される第 3 インタフェースと、前記管理計算機からの要求に従って前記各スイッチ装置の電源の投入及び遮断を制御する第 3 電源制御部と、を備え、

前記各ポートは、前記計算機又は前記ストレージ装置に接続され、

前記複数のスイッチ装置は、第 1 スイッチ装置及び第 2 スイッチ装置を含み、

前記第 1 仮想領域で稼動する前記仮想計算機は、前記第 1 スイッチ装置を経由して前記第 1 論理ボリュームへのデータの入出力を実行し、

前記管理計算機は、

前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーの終了を検出してから、前記差分データの前記第 2 論理ボリュームへのコピーが終了するまでの間に、前記第 2 スイッチ装置に電源を投入する要求を送信し、

前記仮想計算機が前記第 2 スイッチ装置を経由して前記第 2 論理ボリュームへのデータの入出力を実行するように経路を切り替える要求を、前記差分データの前記第 2 論理ボリュームへのコピーが終了した後、前記仮想計算機を前記第 2 仮想領域に移動する要求を送信する前に送信し、

前記仮想計算機は、前記第 2 仮想領域に移動した後、前記第 2 スイッチ装置を経由して前記第 2 論理ボリュームへのデータ入出力を実行することを特徴とする請求項 1 に記載の計算機システム。

【請求項 5】

一つ以上の計算機と、ネットワークを介して前記一つ以上の計算機に接続される一つ以上のストレージ装置と、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される管理計算機と、を備える計算機システムであって、

前記各計算機は、ハードウェア資源として、前記ネットワークに接続される第 1 インタフェースと、前記第 1 インタフェースに接続される第 1 プロセッサと、前記第 1 プロセッサに接続される第 1 メモリと、前記管理計算機からの要求に従って前記ハードウェア資源への電源の投入及び遮断を制御する電源制御部と、を備え、

前記仮想計算機は、さらに、前記一つ以上の計算機のハードウェア資源に基づいて複数の仮想領域を提供する仮想化部を備え、

前記複数の仮想領域は、第 1 仮想領域及び第 2 仮想領域を含み、

前記第 1 仮想領域は、アプリケーションプログラムを実行する仮想計算機として動作し、

前記各ストレージ装置は、前記計算機によって書き込まれたデータを格納する記憶領域を提供する記憶媒体と、前記記憶媒体へのデータの入出力を制御するコントローラと、前記管理計算機からの要求に従って前記各ストレージ装置の電源の投入及び遮断を制御する第 2 電源制御部と、を備え、

10

20

30

40

50

前記一つ以上のストレージ装置の前記コントローラは、前記一つ以上のストレージ装置の前記記憶領域を複数の論理ボリュームとして前記計算機に提供し、

前記複数の論理ボリュームは、第1論理ボリューム及び第2論理ボリュームを含み、

前記第1論理ボリュームには、前記仮想計算機によって書き込まれたデータが格納され、

前記管理計算機は、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される第2インタフェースと、前記第2インタフェースに接続される第2プロセッサと、前記第2プロセッサに接続される第2メモリと、を備え、

前記管理計算機は、

前記第1論理ボリュームを含む前記ストレージ装置に、前記第1論理ボリュームに格納されたデータを前記第2論理ボリュームにコピーする要求を送信し、

前記第1論理ボリュームを含む前記ストレージ装置は、

前記第1論理ボリュームに格納されたデータを前記第2論理ボリュームにコピーする要求を受信した後に前記仮想計算機から前記第1論理ボリュームへのデータの書き込み要求を受信すると、書き込みを要求されたデータを前記第1論理ボリュームに書き込まずに、差分データとして保持し、

前記第1論理ボリュームに格納されたデータを前記第2論理ボリュームにコピーする要求に従って、前記第1論理ボリュームに格納されたデータを読み出して、前記第2論理ボリュームを含む前記ストレージ装置に送信し、

前記管理計算機は、

前記第1論理ボリュームに格納された全データのうち、所定の割合のデータの前記第2論理ボリュームへのコピーが終了すると、前記所定の割合のデータのコピーが終了した時点の前記差分データの量に基づいて、前記第1論理ボリュームに格納された全データのコピーが終了する時点の前記差分データの前記第2論理ボリュームへのコピーが終了する時刻を予測し、

前記予測されたコピーの終了時刻に前記計算機の起動処理が終了するように前記計算機の電源投入時刻を算出し、

前記算出された前記計算機の電源投入時刻が到来すると、前記一つ以上の計算機のハードウェア資源のうち、前記第2仮想領域に割り当てられたハードウェア資源の電源を投入する要求を、前記第2仮想領域に割り当てられたハードウェア資源を含む前記計算機に送信し、

前記第1論理ボリュームに格納されたデータの前記第2論理ボリュームへのコピーの終了を検出すると、前記保持された差分データを前記第2論理ボリュームにコピーする要求を、前記第1論理ボリュームを含む前記ストレージ装置に送信し、

前記差分データの前記第2論理ボリュームへのコピーが終了した後、前記仮想計算機を前記第2仮想領域に移動する要求を送信し、

前記仮想計算機は、前記第2仮想領域に移動した後、前記第2論理ボリュームへのデータ入出力を実行することを特徴とする計算機システム。

【請求項6】

前記ネットワークは、前記一つ以上の計算機と前記一つ以上のストレージ装置との間のデータ転送を中継する複数のスイッチ装置を備え、

前記各スイッチ装置は、複数のポートと、前記管理計算機に接続される第3インタフェースと、前記管理計算機からの要求に従って前記各スイッチ装置の電源の投入及び遮断を制御する第3電源制御部と、を備え、

前記各ポートは、前記計算機又は前記ストレージ装置に接続され、

前記複数のスイッチ装置は、第1スイッチ装置及び第2スイッチ装置を含み、

前記第1仮想領域で稼働する前記仮想計算機は、前記第1スイッチ装置を経由して前記第1論理ボリュームへのデータの入出力を実行し、

前記管理計算機は、

前記予測されたコピーの終了時刻に前記第2スイッチ装置の起動処理が終了するように

10

20

30

40

50

前記第 2 スイッチ装置の電源投入時刻を算出し、

前記算出された前記第 2 スイッチ装置の電源投入時刻が到来すると、前記第 2 スイッチ装置に電源を投入する要求を送信し、

前記仮想計算機が前記第 2 スイッチ装置を経由して前記第 2 論理ボリュームへのデータの入出力を実行するように経路を切り替える要求を、前記差分データの前記第 2 論理ボリュームへのコピーが終了した後、前記仮想計算機を前記第 2 仮想領域に移動する要求を送信する前に送信し、

前記仮想計算機は、前記第 2 仮想領域に移動した後、前記第 2 スイッチ装置を経由して前記第 2 論理ボリュームへのデータ入出力を実行することを特徴とする請求項 5 に記載の計算機システム。

10

【請求項 7】

前記管理計算機は、

前記計算機の起動処理に要する時間を示す情報、及び、前記第 2 スイッチ装置の起動処理に要する時間をあらかじめ保持し、

前記所定の割合のデータのコピーが終了した時点の前記差分データの量を、前記所定の割合の逆数を乗ずることによって、前記前記第 1 論理ボリュームに格納された全データのコピーが終了する時点の前記差分データの量を予測し、

前記予測された差分データの量と、前記第 1 論理ボリュームを含む前記ストレージ装置から前記第 2 論理ボリュームを含む前記ストレージ装置へのデータ転送速度と、に基づいて、前記第 1 論理ボリュームに格納された全データのコピーが終了する時点の前記差分データの

20

前記第 2 論理ボリュームへのコピーが終了する時刻を予測し、

前記予測されたコピーの終了時刻から前記計算機の起動処理に要する時間を減算することによって、前記計算機の電源投入時刻を算出し、

前記予測されたコピーの終了時刻から前記第 2 スイッチ装置の起動処理に要する時間を減算することによって、前記第 2 スイッチ装置の電源投入時刻を算出することを特徴とする請求項 6 に記載の計算機システム。

【請求項 8】

一つ以上の計算機と、ネットワークを介して前記一つ以上の計算機に接続される一つ以上のストレージ装置と、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される管理計算機と、を備える計算機システムを制御する方法であって、

30

前記各計算機は、ハードウェア資源として、前記ネットワークに接続される第 1 インタフェースと、前記第 1 インタフェースに接続される第 1 プロセッサと、前記第 1 プロセッサに接続される第 1 メモリと、前記管理計算機からの要求に従って前記ハードウェア資源への電源の投入及び遮断を制御する電源制御部と、を備え、

前記各計算機は、さらに、前記一つ以上の計算機のハードウェア資源に基づいて複数の仮想領域を提供する仮想化部を備え、

前記複数の仮想領域は、第 1 仮想領域及び第 2 仮想領域を含み、

前記第 1 仮想領域は、アプリケーションプログラムを実行する仮想計算機として動作し、

前記各ストレージ装置は、前記計算機によって書き込まれたデータを格納する記憶領域を提供する記憶媒体と、前記記憶媒体へのデータの入出力を制御するコントローラと、前記管理計算機からの要求に従って前記各ストレージ装置の電源の投入及び遮断を制御する第 2 電源制御部と、を備え、

40

前記一つ以上のストレージ装置の前記記憶領域は、前記一つ以上のストレージ装置の前記コントローラによって、複数の論理ボリュームとして前記計算機に提供され、

前記複数の論理ボリュームは、第 1 論理ボリューム及び第 2 論理ボリュームを含み、

前記第 1 論理ボリュームには、前記仮想計算機によって書き込まれたデータが格納され、

前記管理計算機は、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される第 2 インタフェースと、前記第 2 インタフェースに接続される第 2 プロセッサと、前

50

記第 2 プロセッサに接続される第 2 メモリと、を備え、

前記方法は、

前記管理計算機が、前記第 1 論理ボリュームを含む前記ストレージ装置に、前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求を送信する手順と、

前記第 1 論理ボリュームを含む前記ストレージ装置が、前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求を受信した後に前記仮想計算機から前記第 1 論理ボリュームへのデータの書き込み要求を受信すると、書き込みを要求されたデータを前記第 1 論理ボリュームに書き込まずに、差分データとして保持する手順と、

前記第 1 論理ボリュームを含む前記ストレージ装置が、前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求に従って、前記第 1 論理ボリュームに格納されたデータを読み出して、前記第 2 論理ボリュームを含む前記ストレージ装置に送信する手順と、

前記管理計算機が、前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーの終了を検出すると、前記保持された差分データを前記第 2 論理ボリュームにコピーする要求を、前記第 1 論理ボリュームを含む前記ストレージ装置に送信する手順と、

前記管理計算機が、前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーの終了を検出してから、前記差分データの前記第 2 論理ボリュームへのコピーが終了するまでの間に、前記一つ以上の計算機のハードウェア資源のうち、前記第 2 仮想領域に割り当てられたハードウェア資源の電源を投入する要求を、前記第 2 仮想領域に割り当てられたハードウェア資源を含む前記計算機に送信する手順と、

前記管理計算機が、前記差分データの前記第 2 論理ボリュームへのコピーが終了した後、前記仮想計算機を前記第 2 仮想領域に移動する要求を送信する手順と、

前記仮想計算機が、前記第 2 仮想領域に移動した後、前記第 2 論理ボリュームへのデータ入出力を実行する手順と、を含むことを特徴とする方法。

【請求項 9】

前記方法は、さらに、前記第 1 論理ボリュームを含む前記ストレージ装置が、前記第 1 論理ボリュームに格納されたデータを前記第 2 論理ボリュームにコピーする要求に従う前記データの読み出し及び送信が終了すると、終了メッセージを前記管理計算機に送信する手順を含み、

前記保持された差分データを前記第 2 論理ボリュームにコピーする要求を送信する手順は、前記管理計算機が、前記終了メッセージを受信したときに、前記コピーの終了を検出する手順を含むことを特徴とする請求項 8 に記載の方法。

【請求項 10】

前記保持された差分データを前記第 2 論理ボリュームにコピーする要求を送信する手順は、前記管理計算機が、前記第 1 論理ボリュームの容量、及び、前記第 1 論理ボリュームから前記第 2 論理ボリュームへのデータの転送速度に基づいて、前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーに要する時間を算出する手順と、前記算出された時間が経過したときに、前記コピーの終了を検出する手順と、を含むことを特徴とする請求項 8 に記載の方法。

【請求項 11】

前記ネットワークは、前記一つ以上の計算機と前記一つ以上のストレージ装置との間のデータ転送を中継する複数のスイッチ装置を備え、

前記各スイッチ装置は、複数のポートと、前記管理計算機に接続される第 3 インタフェースと、前記管理計算機からの要求に従って前記各スイッチ装置の電源の投入及び遮断を制御する第 3 電源制御部と、を備え、

前記各ポートは、前記計算機又は前記ストレージ装置に接続され、

前記複数のスイッチ装置は、第 1 スイッチ装置及び第 2 スイッチ装置を含み、

10

20

30

40

50

前記第 1 仮想領域で稼動する前記仮想計算機は、前記第 1 スイッチ装置を経由して前記第 1 論理ボリュームへのデータの入出力を実行し、

前記方法は、さらに、

前記管理計算機が、前記第 1 論理ボリュームに格納されたデータの前記第 2 論理ボリュームへのコピーの終了を検出してから、前記差分データの前記第 2 論理ボリュームへのコピーが終了するまでの間に、前記第 2 スイッチ装置に電源を投入する要求を送信する手順と、

前記管理計算機が、前記仮想計算機が前記第 2 スイッチ装置を経由して前記第 2 論理ボリュームへのデータの入出力を実行するように経路を切り替える要求を、前記差分データの前記第 2 論理ボリュームへのコピーが終了した後、前記仮想計算機を前記第 2 仮想領域に移動する要求を送信する前に送信する手順と、

前記仮想計算機が、前記第 2 仮想領域に移動した後、前記第 2 論理ボリュームへのデータ入出力を実行する手順は、前記第 2 スイッチ装置を経由して実行されることを特徴とする請求項 8 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本願明細書に開示される技術は、計算機システムにおけるリソースの管理に関し、特に、パーティションで区切られた計算機システムにおけるアプリケーションの移動及び電源の制御に関する。

【背景技術】

【0002】

情報処理システムに対しても地球環境への配慮が求められる中、仮想化によって計算機及びストレージ装置を集約し、それによってリソースの有効利用及び消費電力の削減を実現する技術が提案されている。

【0003】

特許文献 1 には、仮想化されたサーバ、スイッチ及びストレージ装置を物理サーバ、物理スイッチ及び物理ストレージ装置に最適に配置することによって負荷を分散させる技術が開示されている。

【0004】

特許文献 2 には、複数の物理サーバを備える計算機システムにおいて、仮想サーバを少数の物理サーバに集約し、残りの物理サーバの電源を遮断する技術が開示されている。これによって計算機システムの消費電力が削減される。

【0005】

特許文献 3 には、仮想計算機と、その仮想計算機が使用する仮想ストレージ装置と、の対応を管理し、使用されていない仮想計算機及びそれに対応する仮想ストレージ装置の電源を遮断する技術が開示されている。これによって計算機システムの消費電力が削減される。

【特許文献 1】特開 2007 - 47986 号公報

【特許文献 2】特開 2007 - 310791 号公報

【特許文献 3】特開 2008 - 102667 号公報

【発明の開示】

【発明が解決しようとする課題】

【0006】

所定の利用方針に基づいてパーティションが区切られた計算機システムにおいては、アプリケーションを所望の時刻に所望のパーティションで実行することが求められる。そのために、仮想サーバ及び論理ボリュームを一つのパーティションから別のパーティションに移動させる場合がある。この移動の際に消費される電力を抑制するためには、これらの移動のタイミングを適切に制御する必要があるが、従来はそのための技術がなかった。

【課題を解決するための手段】

10

20

30

40

50

【 0 0 0 7 】

本願で開示する代表的な発明は、一つ以上の計算機と、ネットワークを介して前記一つ以上の計算機に接続される一つ以上のストレージ装置と、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される管理計算機と、を備える計算機システムであって、前記各計算機は、ハードウェア資源として、前記ネットワークに接続される第1インタフェースと、前記第1インタフェースに接続される第1プロセッサと、前記第1プロセッサに接続される第1メモリと、前記管理計算機からの要求に従って前記ハードウェア資源への電源の投入及び遮断を制御する電源制御部と、を備え、前記各計算機は、さらに、前記一つ以上の計算機のハードウェア資源に基づいて複数の仮想領域を提供する仮想化部を備え、前記複数の仮想領域は、第1仮想領域及び第2仮想領域を含み、前記第1仮想領域は、アプリケーションプログラムを実行する仮想計算機として動作し、前記各ストレージ装置は、前記計算機によって書き込まれたデータを格納する記憶領域を提供する記憶媒体と、前記記憶媒体へのデータの入出力を制御するコントローラと、前記管理計算機からの要求に従って前記各ストレージ装置の電源の投入及び遮断を制御する第2電源制御部と、を備え、前記一つ以上のストレージ装置の前記コントローラは、前記一つ以上のストレージ装置の前記記憶領域を複数の論理ボリュームとして前記計算機に提供し、前記複数の論理ボリュームは、第1論理ボリューム及び第2論理ボリュームを含み、前記第1論理ボリュームには、前記仮想計算機によって書き込まれたデータが格納され、前記管理計算機は、前記一つ以上の計算機及び前記一つ以上のストレージ装置に接続される第2インタフェースと、前記第2インタフェースに接続される第2プロセッサと、前記第2プロセッサに接続される第2メモリと、を備え、前記管理計算機は、前記第1論理ボリュームを含む前記ストレージ装置に、前記第1論理ボリュームに格納されたデータを前記第2論理ボリュームにコピーする要求を送信し、前記第1論理ボリュームを含む前記ストレージ装置は、前記第1論理ボリュームに格納されたデータを前記第2論理ボリュームにコピーする要求を受信した後に前記仮想計算機から前記第1論理ボリュームへのデータの書き込み要求を受信すると、書き込みを要求されたデータを前記第1論理ボリュームに書き込まずに、差分データとして保持し、前記第1論理ボリュームに格納されたデータを前記第2論理ボリュームにコピーする要求に従って、前記第1論理ボリュームに格納されたデータを読み出して、前記第2論理ボリュームを含む前記ストレージ装置に送信し、前記管理計算機は、前記第1論理ボリュームに格納されたデータの前記第2論理ボリュームへのコピーの終了を検出すると、前記保持された差分データを前記第2論理ボリュームにコピーする要求を、前記第1論理ボリュームを含む前記ストレージ装置に送信し、前記第1論理ボリュームに格納されたデータの前記第2論理ボリュームへのコピーの終了を検出してから、前記差分データの前記第2論理ボリュームへのコピーが終了するまでの間に、前記一つ以上の計算機のハードウェア資源のうち、前記第2仮想領域に割り当てられたハードウェア資源の電源を投入する要求を、前記第2仮想領域に割り当てられたハードウェア資源を含む前記計算機に送信し、前記差分データの前記第2論理ボリュームへのコピーが終了した後、前記仮想計算機を前記第2仮想領域に移動する要求を送信し、前記仮想計算機は、前記第2仮想領域に移動した後、前記第2論理ボリュームへのデータ入出力を実行することを特徴とする。

10

20

30

40

【 発明の効果 】

【 0 0 0 8 】

本発明の一実施形態によれば、アプリケーションの移動の際のリソースの電源投入タイミングを制御することによって、アプリケーションの移動の際に消費される電力を抑制することができる。

【 発明を実施するための最良の形態 】

【 0 0 0 9 】

本発明の実施形態を、図面を用いて詳細に説明する。

【 0 0 1 0 】

図1は、本発明の実施形態におけるパーティションの説明図である。

50

【 0 0 1 1 】

パーティションとは、ユーザの業務形態及び情報システムの利用方針に基づいて、情報システムを区切ることによって定義された領域である。各パーティションは、システムの設計者からは物理パーティションとして、ユーザからは論理パーティションとして認識される。

【 0 0 1 2 】

物理パーティションとは、物理的なリソースを区切ることによって定義された物理的な領域である。設計者は、各パーティションに割り当てる物理的なリソースの量を決めることができる。物理的なリソースの量とは、例えば、CPUのコア数、メモリの記憶領域の容量、ネットワークを構成するスイッチ装置のゾーンの大きさ、及び、ストレージ装置の記憶領域の容量等である。

10

【 0 0 1 3 】

論理パーティションとは、実行される業務（アプリケーション）が配置される領域である。ユーザは、アプリケーションが実行されるパーティション、及び、一つのパーティションで実行される複数のアプリケーションの組み合わせを決めることができる。

【 0 0 1 4 】

図2は、本発明の実施形態の計算機システムの構成を示すブロック図である。

【 0 0 1 5 】

本実施形態の計算機システムは、一つ以上のストレージ装置100、一つ以上のサーバ120、一つ以上のファイバーチャネルスイッチ（FC-SW）140、管理サーバ130及び管理ネットワーク150を備える。

20

【 0 0 1 6 】

ストレージ装置100は、サーバ120によって書き込まれたデータを格納する。具体的には、ストレージ装置100は、ディスクドライブ110A～110C及びコントローラ101を備える。

【 0 0 1 7 】

ディスクドライブ110A～110Cは、サーバ120によって書き込まれたデータを格納する記憶媒体を含む記憶デバイスである。以下、ディスクドライブ110A～110Cに共通する説明をする場合、これらを総称してディスクドライブ110とも記載する。

【 0 0 1 8 】

図2には三つのディスクドライブ110を示すが、ストレージ装置100はいくつのディスクドライブ110を備えてもよい。複数のディスクドライブ110がRAID（Redundant Arrays of Inexpensive Disks）を構成してもよい。

30

【 0 0 1 9 】

本実施形態のディスクドライブ110は、記憶媒体として磁気ディスクを備えるハードディスクドライブであるが、これはいかなる種類の装置によって置き換えられてもよい。例えば、ディスクドライブ110は、フラッシュメモリのような不揮発性の半導体記憶装置によって置き換えられてもよい。

【 0 0 2 0 】

コントローラ101は、ディスクドライブ110へのデータの書き込み及びディスクドライブ110からのデータの読み出しを制御する。コントローラ101は、相互に接続されたチャネルアダプタ（CHA）102、ディスクアダプタ（DA）103、インタフェース（I/F）104及び電源制御部105を備える。

40

【 0 0 2 1 】

CHA102は、FC-SW140のいずれかのポート（例えばポート141B）に接続され、サーバ120からのデータ入出力（I/O）要求、すなわち、データ書き込み要求及びデータ読み出し要求を処理する。

【 0 0 2 2 】

DA103は、ディスクドライブ110に接続され、ディスクドライブ110に対する

50

データの書き込み及びディスクドライブ 110 からのデータの読み出しを制御する。

【0023】

CHA102 及び DA103 は、それぞれ、要求された処理を実行するための CPU (図示省略) 及びローカルメモリ (図示省略) を備えてもよい。

【0024】

I/F104 は、管理ネットワーク 150 に接続され、管理ネットワーク 150 を介して管理サーバ 130 と通信する。

【0025】

電源制御部 105 は、ストレージ装置 100 の電源の投入及び遮断 (すなわちストレージ装置 100 への電力の供給の開始及び停止) を制御する。より詳細には、電源制御部 105 は、I/F104 を介して管理サーバ 130 から受信した制御情報に従って、ストレージ装置 100 内の電源制御部 105 以外の部分への電力の供給を制御する。電源制御部 105 は、ストレージ装置 100 のリソースの部分ごと (例えば、ディスクドライブ 110 ごと) の電力の供給を制御してもよい。電源制御部 105 は、例えば、いわゆる Baseboard Management Controller (BMC) によって実現されてもよい。

10

【0026】

コントローラ 101 は、さらに、共有メモリ (図示省略) 及びキャッシュメモリ (図示省略) を備えてもよい。共有メモリには、種々の制御情報が格納される。キャッシュメモリには、ディスクドライブ 110 に書き込まれるデータ及びディスクドライブ 110 から読み出されたデータが一時的に格納される。

20

【0027】

コントローラ 101 は、ディスクドライブ 110 の記憶媒体によって実現される物理的な記憶領域を、複数の論理ボリューム 111 として管理する。図 2 に示す論理ボリューム 111A 及び 111B は、それぞれ、複数の論理ボリューム 111 の一つである。コントローラ 101 は、任意の数の論理ボリューム 111 を管理することができる。

【0028】

本実施形態の計算機システムは、複数のストレージ装置 100 を備えてもよい。

【0029】

サーバ 120 は、相互に接続された CPU 121、メモリ 122、電源制御部 123、ホストバスアダプタ (HBA) 124、I/F 125 及び仮想化部 129 を備える計算機である。

30

【0030】

CPU 121 は、メモリ 122 に格納されたプログラムを実行するプロセッサである。図 2 には一つの CPU 121 のみを示すが、サーバ 120 は複数の CPU 121 を備えてもよい。

【0031】

メモリ 122 には、CPU 121 によって実行されるプログラム及び CPU 121 によって参照されるデータが格納される。本実施形態のメモリ 122 には、少なくとも、オペレーティングシステム (OS) 128、バス管理プログラム 127 及びアプリケーションプログラム 126 が格納される。

40

【0032】

OS 128 は、サーバ 120 を管理する基本ソフトウェアであり、例えば Windows (登録商標) 又は Unix (登録商標) のようなものであってもよい。後述のように、サーバ 120 上で複数の OS 128 が実行される場合がある。その場合、同一の種類の複数の OS 128 が実行されてもよいし、異なる種類の複数の OS 128 が実行されてもよい。

【0033】

バス管理プログラム 127 は、サーバ 120 から論理ボリューム 111 へのアクセス経路 (バス) を制御する。

50

【 0 0 3 4 】

アプリケーションプログラム 1 2 6 は、種々の業務（アプリケーション）を実現するプログラムである。サーバ 1 2 0 のユーザは、所望のアプリケーションを実現するアプリケーションプログラム 1 2 6 をサーバ 1 2 0 にインストールし、実行することができる。アプリケーションプログラム 1 2 6 は、必要に応じて、論理ボリューム 1 1 1 へのデータ I / O 要求を発行する。

【 0 0 3 5 】

上記の各ソフトウェア（プログラム）は、CPU 1 2 1 によって実行される。したがって、以下の説明において上記の各ソフトウェア（プログラム）が実行する処理は、実際には CPU 1 2 1 によって実行される。

【 0 0 3 6 】

仮想化部 1 2 9 は、計算機システム内の一つ以上のサーバ 1 2 0 のリソース（すなわち、CPU 1 2 1 及びメモリ 1 2 2 等）を用いて、複数の仮想的な領域（すなわち仮想的な計算機）を提供する。仮想化部 1 2 9 は、一つのサーバ 1 2 0 に複数の OS 1 2 8 を実行させることができる。あるいは、仮想化部 1 2 9 は、複数のサーバ 1 2 0 を一つの仮想的な計算機として使用する、いわゆるクラスタリングを実現することもできる。

【 0 0 3 7 】

仮想化部 1 2 9 は、サーバ 1 2 0 内に実装されたハードウェアであってもよいし、メモリ 1 2 2 に格納されたプログラム（例えば、いわゆる仮想マシンモニタ又はハイパーバイザ等の仮想化ソフトウェア）であってもよい。仮想化部 1 2 9 がメモリ 1 2 2 に格納されたプログラムである場合、CPU 1 2 1 が仮想化ソフトウェアを実行することによって仮想化部 1 2 9 の機能が実現される。

【 0 0 3 8 】

HBA 1 2 4 は、FC - SW 1 4 0 のいずれかのポート（例えばポート 1 4 1 A）に接続され、ストレージ装置 1 0 0 との通信を実行するインタフェースである。

【 0 0 3 9 】

I / F 1 2 5 は、管理ネットワーク 1 5 0 に接続され、管理ネットワーク 1 5 0 を介して管理サーバ 1 3 0 と通信する。

【 0 0 4 0 】

電源制御部 1 2 3 は、サーバ 1 2 0 の電源を制御する。具体的には、電源制御部 1 2 3 は、I / F 1 2 5 を介して管理サーバ 1 3 0 から受信した制御情報に従って、サーバ 1 2 0 の電源の投入及び遮断を制御する。電源制御部 1 2 3 は、電源制御部 1 0 5 と同様のものであってもよい。電源制御部 1 2 3 は、サーバ 1 2 0 に含まれるリソースの部分ごと（例えばサーバ 1 2 0 が複数の CPU 1 2 1 を備える場合、CPU 1 2 1 ごと）に電源の投入及び遮断を制御してもよい。

【 0 0 4 1 】

CPU 1 2 1 は、アプリケーションプログラム 1 2 6 を実行し、必要に応じて、論理ボリューム 1 1 1 に対するデータ書き込み要求及びデータ読み出し要求を、HBA 1 2 4 を介してストレージ装置 1 0 0 に送信する。これらの要求の送信先（言い換えると、これらの要求によるデータ I / O に使用される論理パス）は、バス管理プログラム 1 2 7 によって制御される。

【 0 0 4 2 】

本実施形態の計算機システムは、複数のサーバ 1 2 0 を備えてもよい。

【 0 0 4 3 】

FC - SW 1 4 0 は、サーバ 1 2 0 とストレージ装置 1 0 0 との間のデータ I / O を中継するネットワークを構成する。FC - SW 1 4 0 は、サーバ 1 2 0 とストレージ装置 1 0 0 との間のデータ I / O 経路を切り替えることができる。本実施形態において、サーバ 1 2 0 とストレージ装置 1 0 0 との間のデータ I / O は、ファイバーチャネル（FC）プロトコルに基づいてやり取りされる。

【 0 0 4 4 】

10

20

30

40

50

FC-SW140は、複数のポート141（図2の例では、ポート141A～ポート141D）、I/F142及び電源制御部143を備える。

【0045】

各ポートは、サーバ120のHBA124又はストレージ装置100のCHA102に接続される。

【0046】

FC-SW140は、各ポート141間の接続を設定することによって、サーバ120とストレージ装置100との間のデータ通信経路を設定することができる。さらに、FC-SW140は、互いに独立したいわゆるゾーンを設定することができる。

【0047】

I/F142は、管理ネットワーク150に接続され、管理ネットワーク150を介して管理サーバ130と通信する。

【0048】

電源制御部143は、FC-SW140の電源を制御する。具体的には、電源制御部143は、I/F142を介して管理サーバ130から受信した制御情報に従って、FC-SW140の電源の投入及び遮断を制御する。電源制御部143は、電源制御部105と同様のものであってもよい。電源制御部143は、FC-SW140に含まれるリソースの部分ごと（例えばポート141ごと）に電源の投入及び遮断を制御してもよい。

【0049】

本実施形態の計算機システムは、複数のFC-SW140を備えてもよい。

【0050】

管理サーバ130は、相互に接続されたCPU131、メモリ132、データベース133、及びI/F134を備える計算機である。

【0051】

CPU131は、メモリ132に格納されたプログラムを実行するプロセッサである。

【0052】

メモリ132には、CPU131によって実行されるプログラム及びCPU131によって参照されるデータが格納される。本実施形態のメモリ132には、少なくとも、管理プログラム135が格納される。

【0053】

データベース133には、計算機システムを管理するための情報が格納される。データベース133は、例えば、管理サーバ130に接続された（又は内蔵された）ディスクドライブに格納されてもよい。

【0054】

本実施形態のデータベース133には、管理テーブル136が格納される。管理テーブル136の全部又は一部が、必要に応じてメモリ132にコピーされ、CPU131によって参照されてもよい。管理テーブル136の内容については後述する（図16参照）。

【0055】

I/F134は、管理ネットワーク150に接続され、管理ネットワーク150を介してストレージ装置100、サーバ120及びFC-SW140と通信する。例えば、ストレージ装置100、サーバ120及びFC-SW140の各々の電源を制御するための信号は、I/F134から管理ネットワーク150を介して送信される。

【0056】

管理ネットワーク150は、いかなる種類のものであってもよい。典型的には、管理ネットワーク150は、いわゆるLAN（Local Area Network）のようなIP（Internet Protocol）ネットワークである。その場合、I/F104、I/F125及びI/F134は、いわゆるネットワークインタフェースカードであってよい。

【0057】

図3は、本発明の実施形態における物理パーティションの第1の例の説明図である。

10

20

30

40

50

【 0 0 5 8 】

図 3 は、三つの物理パーティション、すなわち、パーティション 1 __ 3 0 0 A、パーティション 2 __ 3 0 0 B 及びパーティション 3 __ 3 0 0 C が定義される例を示す。この例において、各パーティションは、一つ以上のサーバ 1 2 0、一つ以上の FC - SW 1 4 0 及び一つ以上のストレージ装置 1 0 0 を含む。

【 0 0 5 9 】

具体的には、パーティション 1 __ 3 0 0 A は、三つのサーバ 1 2 0、三つの FC - SW 1 4 0 及び三つのストレージ装置 1 0 0 を含む。これらの三つのサーバ 1 2 0 は、仮想化部 1 2 9 によって実現された一つのクラスタである。

【 0 0 6 0 】

同様に、パーティション 2 __ 3 0 0 B は、二つのサーバ 1 2 0、二つの FC - SW 1 4 0 及び二つのストレージ装置 1 0 0 を含む。これらの二つのサーバ 1 2 0 は、仮想化部 1 2 9 によって実現された一つのクラスタである。

【 0 0 6 1 】

パーティション 3 __ 3 0 0 C は、一つのサーバ 1 2 0、一つの FC - SW 1 4 0 及び一つのストレージ装置 1 0 0 を含む。

【 0 0 6 2 】

なお、図 3 において、サーバ 1 2 0、FC - SW 1 4 0 及びストレージ装置 1 0 0 の間の接続の図示は省略されている。

【 0 0 6 3 】

各パーティションに含まれる一つ以上のサーバ 1 2 0 では、他のパーティションとは独立に OS 1 2 8 が実行され、その OS 1 2 8 の上で、他のパーティションとは独立にアプリケーションプログラム 1 2 6 が実行される。各パーティションで実行されるアプリケーションプログラム 1 2 6 は、そのパーティションに含まれる一つ以上の FC - SW 1 4 0 を介して、そのパーティションに含まれる一つ以上のストレージ装置 1 0 0 内の論理ボリュームへのデータ I / O を実行する。

【 0 0 6 4 】

図 4 は、本発明の実施形態における物理パーティションの第 2 の例の説明図である。

【 0 0 6 5 】

図 4 は、一つのサーバ 1 2 0、一つの FC - SW 1 4 0 及び一つのストレージ装置 1 0 0 からなる計算機システムに、三つの物理パーティション、すなわち、パーティション 1 __ 3 0 0 D、パーティション 2 __ 3 0 0 E 及びパーティション 3 __ 3 0 0 F が定義される例を示す。

【 0 0 6 6 】

具体的には、図 4 に示すサーバ 1 2 0 は、6 個の CPU 1 2 1 を備え、それらのうち 3 個がパーティション 1 __ 3 0 0 D に、2 個がパーティション 2 __ 3 0 0 E に、1 個がパーティション 3 __ 3 0 0 F に割り当てられる。

【 0 0 6 7 】

図 4 に示す FC - SW 1 4 0 は、1 2 個のポート 1 4 1 を備え、それらのうち 4 個がパーティション 1 __ 3 0 0 D に割り当てられ、残りの 8 個のうち 4 個がパーティション 2 __ 3 0 0 E に割り当てられ、残りの 4 個がパーティション 3 __ 3 0 0 F に割り当てられる。各パーティションに割り当てられた 4 個のポート 1 4 1 は、ゾーン 4 0 1 を形成する。

【 0 0 6 8 】

図 4 に示すストレージ装置 1 0 0 が提供する記憶領域のうち、一部（例えば、一部のボリュームプール 4 0 2）がパーティション 1 __ 3 0 0 D に割り当てられ、別のボリュームプール 4 0 2 がパーティション 2 __ 3 0 0 E に割り当てられ、残りがパーティション 3 __ 3 0 0 F に割り当てられる。なお、ボリュームプール 4 0 2 とは、一つ以上の論理ボリューム 1 1 1 からなる記憶領域の管理単位である。

【 0 0 6 9 】

この例において、各パーティションに割り当てられた 1 個以上の CPU 1 2 1 において

10

20

30

40

50

、OS 128が他のパーティションとは独立に実行され、それらのOS 128の上でアプリケーションプログラム126が他のパーティションとは独立に実行される。各パーティションで実行されるアプリケーションプログラムは、そのパーティションに割り当てられた記憶領域へのデータI/Oを実行する。

【0070】

上記のように、パーティションは、複数の装置のリソースを結合することによって定義されてもよいし(図3参照)、一つの装置のリソースを分割することによって定義されてもよい(図4参照)。いずれの場合も、システム的设计者は、各パーティションを、物理的なリソースを区切る物理パーティションとして認識する。

【0071】

図5は、本発明の実施形態における論理パーティションの例の説明図である。

【0072】

計算機システムのユーザは、各パーティションを、実行する業務が配置される領域として認識する。以下の説明において、「業務」を「アプリケーション」と記載する。例えば、図3に示すパーティション1__300A、パーティション2__300B及びパーティション3__300Cが、それぞれ、論理パーティション、すなわち、パーティション1__300G、パーティション2__300H及びパーティション3__300Iとしてユーザに認識されてもよい。あるいは、図4に示すパーティション1__300D、パーティション2__300E及びパーティション3__300Fが、それぞれ、論理パーティション、すなわち、パーティション1__300G、パーティション2__300H及びパーティション3__300Iとしてユーザに認識されてもよい。

【0073】

図5の例では、パーティション1__300Gで三つの業務(アプリケーション)500が実行され、パーティション2__300Hで二つのアプリケーション500が実行され、パーティション3__300Iで一つのアプリケーション500が実行される。

【0074】

図6は、本発明の実施形態におけるアプリケーション500の構成の説明図である。

【0075】

図5に示す各アプリケーション500は、図6に示すように、仮想サーバ(VM)601、論理パス602及び論理ボリューム(LU)603によって構成される。

【0076】

仮想サーバ601は、仮想化部129によって実現される仮想的な計算機である。

【0077】

論理パス602は、仮想サーバ601が論理ボリューム603にアクセスするために使用される経路である。論理パス602は、HBA124からFC-SW140を經由してストレージ装置100に至る物理的な経路によって実現される。

【0078】

論理ボリューム603は、ストレージ装置100が仮想サーバ601に提供する論理的な記憶領域である。仮想サーバ601は、一つの論理ボリューム603を一つの記憶デバイスとして認識する。例えば、一つの論理ボリューム111が一つの論理ボリューム603として提供されてもよいし、複数の論理ボリューム111が一つの論理ボリューム603として提供されてもよい。あるいは、論理ボリューム603の記憶領域へのデータ書き込み要求を受信したときに、その記憶領域に、論理ボリューム111の記憶領域が割り当てられてもよい。

【0079】

仮想サーバ601上でOS128が実行され、そのOS128上でアプリケーションプログラム126が実行される。そのアプリケーションプログラム126は、論理パス602を介して、論理ボリューム603へのデータI/Oを実行する。これによって、アプリケーション500が実現される。

【0080】

10

20

30

40

50

各アプリケーション500は、一つのパーティションから別のパーティションへ移動することができる。このようなアプリケーション500の移動については後述する。

【0081】

図7は、本発明の実施形態の全体システム構成の説明図である。

【0082】

図7に示すパーティション1__300A、パーティション2__300B及びパーティション3__300Cは、図3に示したものと同様の物理パーティションとしてシステム管理者に認識される。

【0083】

一方、ユーザは、パーティション1__300A、パーティション2__300B及びパーティション3__300Cを論理パーティションとして認識する。パーティション1__300Aには、二つのアプリケーション(A P P)500(すなわち、アプリケーション500A及び500B)が配置される。パーティション2__300Bには、二つのアプリケーション500(すなわち、アプリケーション500C及び500D)が配置される。パーティション3__300Cには一つのアプリケーション500(すなわち、アプリケーション500E)が配置される。

10

【0084】

アプリケーション500Aを構成する仮想サーバ601、論理バス602及び論理ボリューム603を、それぞれ、仮想サーバ601A、論理バス602A及び論理ボリューム603Aと記載する。同様に、アプリケーション500Bは、仮想サーバ601B、論理バス602B及び論理ボリューム603Bによって構成される。アプリケーション500Cは、仮想サーバ601C、論理バス602C及び論理ボリューム603Cによって構成される。アプリケーション500Dは、仮想サーバ601D、論理バス602D及び論理ボリューム603Dによって構成される。アプリケーション500Eは、仮想サーバ601E、論理バス602E及び論理ボリューム603Eによって構成される。

20

【0085】

以下の説明において、アプリケーション500A~500Eに共通する説明をする場合には、これらを総称してアプリケーション500と記載する。仮想サーバ601A~601Eに共通する説明をする場合、これらを総称して仮想サーバ601と記載する。論理バス602A~602Eに共通する説明をする場合、これらを総称して論理バス602と記載する。論理ボリューム603A~603Eに共通する説明をする場合、これらを総称して論理ボリューム603と記載する。

30

【0086】

管理サーバ130の管理プログラム135は、各パーティションを管理する。具体的には、管理サーバ130は、管理ネットワーク150を介して各パーティションに含まれるサーバ120、FC-SW140及びストレージ装置100に接続され、それらの装置を管理する。

【0087】

さらに、管理プログラム135は、各パーティションへのハードウェアリソースの割り当て、及び、各パーティションへの各アプリケーション500の配置を管理する。例えば、管理プログラム135は、アプリケーション500のパーティション間の移動(具体的には、論理ボリューム603の移動、論理バス602の切り替え、及び、仮想サーバ601の移動)を制御することができる。

40

【0088】

ここで、本実施形態において実行されるアプリケーション500のパーティション間の移動について説明する。

【0089】

管理プログラム135は、アプリケーション500をあるパーティションから別のパーティションに移動することができる。このような移動は、種々の目的、例えば、消費電力の削減又は負荷分散のために実行される。

50

【0090】

具体的には、例えば、図7に示すアプリケーション500A～500Eが夜間にはほとんど使用されない場合、アプリケーション500A～500Eが夜間には一つのパーティションに配置され、昼間には複数のパーティションに分散されるように、アプリケーション500の移動が制御されてもよい。例えば、昼間（例えば8：00から翌日の0：00まで）は図7に示すように各アプリケーション500が配置され、夜間（例えば0：00から8：00まで）はアプリケーション500A～500Eがすべてパーティション3__300Cに配置されてもよい。

【0091】

この場合、0：00にアプリケーション500A及び500Bをパーティション1__300Aからパーティション3__300Cに、同じく0：00にアプリケーション500C及び500Dをパーティション2__300Bからパーティション3__300Cに移動する処理を管理プログラム135が実行する。その後、パーティション1__300A及びパーティション2__300Bの電源を遮断する（すなわち、それらのパーティションに割り当てられた物理的なリソースへの電力の供給を遮断する）ことによって、計算機システムの消費電力を削減することができる。

10

【0092】

その後、8：00にアプリケーション500A及び500Bをパーティション3__300Cからパーティション1__300Aに、同じく8：00にアプリケーション500C及び500Dをパーティション3__300Cからパーティション2__300Bに移動する処理を管理プログラム135が実行する。0：00の移動の後に移動先のパーティションの電源が遮断されている場合、8：00の移動の前にそれらのパーティションの電源を投入する必要がある。

20

【0093】

図8は、本発明の実施形態において実行されるアプリケーション500のパーティション間の移動処理を示すフローチャートである。

【0094】

管理プログラム135は、時刻等をトリガーにして、アプリケーション500の移動処理を開始する（ステップ801）。例えば、上記の図7の例の場合、0：00及び8：00に移動処理が開始されてもよい。あるいは、移動処理は、0：00及び8：00から終了処理に要する時間を減算した時刻開始されてもよい。その場合、移動処理は、0：00及び8：00に終了する。

30

【0095】

次に、管理プログラム135は、これから移動するアプリケーション500が現在どのパーティションに配置されているかを確認する（ステップ802）。この確認のために、後に説明する管理テーブル136が参照される。

【0096】

次に、管理プログラム135は、移動先のパーティションを選択する（ステップ803）。上記の図7の例のようにあらかじめ移動先のパーティションが決定されている場合、そのパーティションが選択されてもよい。移動先として選択可能なパーティションが複数存在する場合には、それらのうち一つが選択されてもよい。

40

【0097】

次に、管理プログラム135は、移動先として選択されたパーティションの状態を確認する（ステップ804）。具体的には、管理プログラム135は、移動先として選択されたパーティションに既に配置されているアプリケーション500の数、それらのアプリケーション500によって既に使用されている論理ボリューム603の容量、及び、そのパーティションを使用可能な時間帯、等の情報を参照する。ステップ804で実行される詳細な処理については、後に図27を参照して説明する。

【0098】

次に、管理プログラム135は、ステップ804で参照した情報に基づいて、選択され

50

たパーティションへのアプリケーション500の移動が可能であるか否かを判定する(ステップ805)。ステップ804の参照及びステップ805の判定のために、管理テーブル136が参照される。

【0099】

ステップ805において、選択されたパーティションへのアプリケーション500の移動が可能でないと判定された場合、管理プログラム135は、ステップ803に戻って、別のパーティションを移動先として選択する。なお、選択可能な全てのパーティションについて、移動が可能でないと判定された場合、管理プログラム135は、アプリケーション500の移動を実行せずに図8の処理を終了してもよい。

【0100】

ステップ805において、選択されたパーティションへのアプリケーション500の移動が可能であると判定された場合、管理プログラム135は、アプリケーション500の移動を実行する(ステップ806)。この移動については、後に図9及び図10を参照して詳細に説明する。

【0101】

次に、管理プログラム135は、ステップ806の移動が反映されるように、アプリケーション500の配置に関する情報(具体的には、管理テーブル136に含まれる情報)を更新する(ステップ807)。この更新については、後に管理テーブル136を参照して詳細に説明する。

【0102】

以上で、アプリケーション500の移動処理が終了する(ステップ808)。

【0103】

図9は、本発明の実施形態において実行されるアプリケーション500のパーティション間の移動処理の詳細な手順を示す説明図である。

【0104】

図9は、例として、図7の計算機システムにおいてアプリケーション500Aをパーティション1__300Aからパーティション2__300Bに移動する手順を示す。ただし、説明を簡単にするため、説明に必要なない部分の図示は省略されている。

【0105】

図9(A)は、初期状態、すなわち、アプリケーション500Aが移動する前の状態を示す。この状態において、仮想サーバ601Aは、パーティション1__300A内のサーバ120で稼働している。論理ボリューム603Aは、パーティション1__300A内のストレージ装置100によって管理されている。仮想サーバ601Aが論理ボリューム603Aにアクセスするために使用される論理パス602Aは、パーティション1__300A内のFC-SW140を経由する。

【0106】

この状態において、アプリケーション500Aの移動が開始されると、最初に、論理ボリューム603Aの移動が実行される。具体的には、論理ボリューム603Aの複製がパーティション2__300B内のストレージ装置100に作成される。

【0107】

図9(B)は、論理ボリューム603Aの移動が終了した状態を示す。図9(B)に示す論理ボリューム603Fは、論理ボリューム603Aの複製である。この複製を作成するために、論理ボリューム603Aに格納された全てのデータを読み出して、それを論理ボリューム603Fにコピーする処理が実行される。図9(B)の状態において、コピーは終了しているが、仮想サーバ601Aはまだ論理ボリューム603Aにアクセスし、論理ボリューム603Fはまだ使用されない。

【0108】

次に、論理パス602Aの移動が実行される。具体的には、論理パス602Aがパーティション2__300B内のFC-SW140を経由して論理ボリューム603Fへのアクセスを可能とするように、論理パス602Aの設定が切り替えられる。この切り替えは、

10

20

30

40

50

パーティション 1__300A 内のサーバ 120 のパス管理プログラム 127 が、論理ボリューム 603F を新たに論理ボリューム 603A として認識し、さらに、論理ボリューム 603A へのデータ I/O 要求の送信先を、パーティション 1__300A 内の FC-SW 140 のポート 141 から、パーティション 2__300B 内の FC-SW 140 のポート 141 に変更することによって実行されてもよい。

【0109】

図 9 (C) は、論理パス 602A の移動が終了した状態を示す。この状態において、論理パス 602A は、パーティション 1__300A 内の仮想サーバ 601A がパーティション 2__300B 内の FC-SW 140 を経由して論理ボリューム 603F にアクセスできるように設定されている。

10

【0110】

次に、仮想サーバ 601A の移動が実行される。この移動は、例えば、パーティション 1__300A 内のサーバ 120 のメモリ 122 のイメージを、パーティション 2__300B 内のサーバ 120 のメモリ 122 にコピーすることによって実行される。このような移動は、仮想化部 129 の機能によって実現されてもよい。

【0111】

図 9 (D) は、仮想サーバ 601A の移動が終了した状態を示す。この状態において、パーティション 2__300B 内のサーバ 120 に移動した仮想サーバ 601A は、パーティション 2__300B 内の FC-SW 140 を経由する論理パス 602A を用いて、論理ボリューム 603F (すなわち新たな論理ボリューム 603A) へのデータ I/O を実行することができる。これによって、アプリケーション 500A がパーティション 2__300B に移動する。

20

【0112】

後述するように、論理ボリューム 603A の移動は、アプリケーション 500A の稼働を停止することなく実行することができる。一方、論理パス 602A 及び仮想サーバ 601A の移動を実行するためには、アプリケーション 500A の稼働を停止する必要がある。しかし、その停止時間は十分に短いため、ユーザの利便性を損ねることなくアプリケーション 500A を移動することができる。

【0113】

図 10 は、本発明の実施形態において実行されるアプリケーション 500 のパーティション間の移動処理の詳細な手順を示すフローチャートである。

30

【0114】

図 10 に示す移動処理は、図 8 のステップ 806 において実行を開始される (ステップ 1001)。すなわち、この処理の実行が開始される時点で、移動先のパーティションが既に選択されている。選択されたパーティションを識別する情報が引数として与えられる。

【0115】

次に、管理プログラム 135 は、選択されたパーティションへのアプリケーション 500 の移動を実行する (ステップ 1002)。具体的には、管理プログラム 135 は、図 9 に示したように、最初に論理ボリューム 603 を移動し、次に論理パス 602 を移動し、最後に仮想サーバ 601 を移動する。ステップ 1002 において実行される詳細な処理については、後に図 17 等を参照して説明する。

40

【0116】

次に、管理プログラム 135 は、移動の結果を反映するように管理テーブル 136 を更新する (ステップ 1003)。

【0117】

以上で、図 10 に示す移動の処理が終了する (ステップ 1004)。その後、処理は図 8 のフローチャートに戻り、ステップ 807 以降の処理が実行される。

【0118】

次に、論理ボリューム 603 の移動の詳細な手順を説明する。

50

【 0 1 1 9 】

図 1 1 は、本発明の実施形態において実行される論理ボリューム 6 0 3 の移動の説明図である。

【 0 1 2 0 】

上記のように、論理ボリューム 6 0 3 の移動は、アプリケーション 5 0 0 を停止することなく実行される。このため、論理ボリューム 6 0 3 の移動のためのコピーが開始された後、それが終了するまでの間に、移動中の論理ボリューム 6 0 3 へのデータ I / O が実行される場合がある。本実施形態では、コピーが開始された後のデータ I / O は、論理ボリューム 6 0 3 に反映されず、差分ストックとして保持される。そして、論理ボリューム 6 0 3 の全データのコピーが終了した後、差分ストックとして保持されているデータ I / O がコピー先の論理ボリューム 6 0 3 に反映される。

10

【 0 1 2 1 】

例えば、図 9 に示すように、論理ボリューム 6 0 3 A のデータが論理ボリューム 6 0 3 F にコピーされる場合について説明する。ストレージ装置 1 0 0 のコントローラ 1 0 1 は、コピーが開始された後に論理ボリューム 6 0 3 A へのデータ書き込み要求を受信すると、そのデータを論理ボリューム 6 0 3 A 及び論理ボリューム 6 0 3 F のいずれにも書き込まずに差分ストックとして保持する。この差分ストックは、例えばストレージ装置 1 0 0 内のいずれかの論理ボリューム 1 1 1 に格納されてもよい。

【 0 1 2 2 】

論理ボリューム 6 0 3 A に格納されている全てのデータの論理ボリューム 6 0 3 F へのコピー（以下、全体コピーと記載する）が終了すると、次に、差分ストックとして格納されたデータを、論理ボリューム 6 0 3 F に書き込む。差分ストックとして格納されたデータのコピー先への書き込みを、以下、差分コピーと記載する。

20

【 0 1 2 3 】

このような全体コピー及び差分コピーは、コントローラ 1 0 1 によって制御される。例えば、コントローラ 1 0 1 内の制御プロセッサ（図示省略）が、制御メモリ（図示省略）に格納された I / O 制御プログラム（図示省略）を実行することによって全体コピー及び差分コピーが実現されてもよい。

【 0 1 2 4 】

本実施形態では、全体コピー処理の終了を契機として、サーバ 1 2 0 及び F C - S W 1 4 0 の電源が投入される。これは、計算機システムによる消費電力を削減するためである。電源投入のタイミングの制御については、後に詳細に説明する（図 1 4、図 1 5 等参照）。

30

【 0 1 2 5 】

なお、従来のストレージ装置では、全体コピーと差分コピーとの組が一つのコピー処理として扱われるため、ストレージ装置の外部の装置（例えば管理サーバ）は、全体コピーのみが終了したことを知ることができなかった。このため、全体コピー処理の終了（又は差分コピー処理の開始）とその他の処理とを関連付けて制御することができなかった。これに対して、本実施形態では、上記のような制御を実現するため、全体コピー処理が終了したことを示すメッセージをストレージ装置 1 0 0 が管理サーバ 1 3 0 に送信してもよいし（後述する図 2 0 参照）、管理サーバ 1 3 0 が全体コピー処理の終了時刻を算出してもよい。

40

【 0 1 2 6 】

図 1 2 は、本発明の実施形態において実行される論理ボリューム 6 0 3 の移動の別の例の説明図である。

【 0 1 2 7 】

ストレージ装置 1 0 0 がボリュームコピー機能（いわゆるリモートコピー機能又はローカルコピー機能）を持つ場合、そのボリュームコピー機能を用いて論理ボリューム 6 0 3 を移動することができる。

【 0 1 2 8 】

50

例えば、論理ボリューム 603A を正ボリューム、論理ボリューム 603F を副ボリュームとするボリュームペアが構成される場合、論理ボリューム 603A に格納されたデータの更新が、論理ボリューム 603F にも反映される。

【0129】

論理ボリューム 603A の更新が直ちに論理ボリューム 603F に反映されるように設定されている場合、論理ボリューム 603F には、ほぼ常に論理ボリューム 603A と同一のデータが格納されている。このため、図 9 のようにアプリケーション 500A が移動する場合の論理ボリューム 603A の移動に要する時間はほぼ 0 となる。

【0130】

一方、論理ボリューム 603A の更新が直ちに論理ボリューム 603F に反映されないように設定されている場合もある。例えば、定期的に更新が反映されるように設定されていてもよいし、論理ボリューム 603A と論理ボリューム 603F との間のデータ転送経路のトラフィック量が所定の閾値より小さいときに更新が反映されるように設定されていてもよい。このような場合、図 9 のようなアプリケーション 500A の移動における論理ボリューム 603A の移動は、前回の反映が実行された後で更新された論理ボリューム 603A のデータ（すなわち差分データ）のみがコピーされる。

【0131】

アプリケーション 500A がパーティション 1__300A に配置されているときには論理ボリューム 603A が、アプリケーション 500A がパーティション 2__300B に配置されているときには論理ボリューム 603F が使用される。

【0132】

上記のように、ボリュームコピー機能を使用することによって、図 11 に示すような全体コピー及び差分コピーを実行する場合と比較して短時間で論理ボリューム 603 の移動を実行することができる。

【0133】

図 13 は、本発明の実施形態における物理パーティションの第 3 の例の説明図である。

【0134】

図 13 の例では、図 3 と同様に、一つ以上のサーバ 120 が各パーティションに割り当てられるが、図 4 と同様に、一つの FC - SW 140 及び一つのストレージ装置 100 が複数のパーティションに共有される。

【0135】

このような場合であっても、FC - SW 140 及びストレージ装置 100 が部分ごとに電源を制御することができれば、上記と同様の電源制御を実行することができる。例えば、ストレージ装置 100 が複数のディスクドライブ 110 を備え、かつ、ディスクドライブ 110 ごと（又は、複数のディスクドライブからなる RAID グループ（図示省略）ごと）の電源の投入、切断を制御できる場合、移動元の論理ボリューム 603 を格納するディスクドライブ 110 及び移動先の論理ボリューム 603 が属するディスクドライブ 110 の電源を適切なタイミングで制御することによって、後述するように消費電力を削減することができる。

【0136】

ただし、計算機システム全体の消費電力に対して、上記のようにストレージ装置 100 等の部分ごとの制御によって削減される電力の量が十分に小さい場合、そのような部分ごとの制御による消費電力削減効果は小さい。このような場合、複数のパーティションによって共有されるストレージ装置 100 等の電源は常時投入し、サーバ 120 のみの電源を制御してもよい。

【0137】

図 14 は、本発明の実施形態において実行されるアプリケーション 500 の移動及び電源制御の説明図である。

【0138】

図 14 は、図 9 を参照して説明したアプリケーション 500 の移動の手順に加えて、そ

10

20

30

40

50

の移動の際に実行される電源制御の手順を示す。ただし、図9はアプリケーション500がパーティション1__300Aからパーティション2__300Bに移動する例を示したが、図14は、アプリケーション500がパーティション2__300Bからパーティション1__300Aに移動する例を示す。例えば、図9の処理によってパーティション2__300Bに移動したアプリケーション500Aが、パーティション1__300Aに戻るときに図14の処理が実行される。なお、図9の処理が実行される場合にも図14と同様の電源制御が実行されてよい。

【0139】

図14に示すサーバ120等のハードウェア及び仮想サーバ601A等のアプリケーションの構成は図9と同様であるため説明を省略する。ただし、図14のパーティション1__300Aのサーバ120、FC-SW140及びストレージ装置100に表示された「サーバP1」、「FC-SWP1」及び「ストレージP1」、並びに、パーティション2__300Bのサーバ120、FC-SW140及びストレージ装置100に表示された「サーバP2」、「FC-SWP2」及び「ストレージP2」は、それぞれ、各ハードウェアの識別子である。これらの識別子は、後述するテーブルに登録される(図16等参照)。

10

【0140】

図14の処理が実行される直前の時点において、パーティション1__300Aに割り当てられたサーバ120、FC-SW140及びストレージ装置100の電源は切断されている。

20

【0141】

図14の例では、最初に、移動先であるパーティション1__300Aのストレージ装置100の電源が投入される。

【0142】

次に、パーティション2__300Bからパーティション1__300Aへの論理ボリューム603の移動のための全体コピーが開始される。図14の例では、論理ボリューム603Fに格納されている全データが論理ボリューム603Aにコピーされる。このコピーも、図11に示したように実行される。すなわち、この全体コピーが開始された後の論理ボリューム603Fの更新は禁止され、その更新の内容は差分ストックとして保持される。

30

【0143】

次に、全体コピーが終了する。

【0144】

全体コピーの終了を契機として、パーティション1__300Aのサーバ120及びFC-SW140の電源が投入される。本実施形態では仮想サーバ601Aを移動する前に論理パス602Aを移動する必要があるため、サーバ120の起動処理がFC-SW140の起動処理より先に終了する必要はない。このため、サーバ120の電源を投入する前にFC-SW140の電源を投入してもよい。しかし、一般に、サーバ120の起動に要する時間は、FC-SW140の起動に要する時間より長い。このため、図14に示すように、先にサーバ120の電源を投入してもよい。

40

【0145】

次に、論理ボリューム603Fから論理ボリューム603Aへの差分コピーが実行される。具体的には、差分ストックとして保持されているデータが論理ボリューム603Aに書き込まれる。

【0146】

上記のような論理ボリューム603の移動は、例えば、従来のオンラインマイグレーションの機能によって実現される。

【0147】

差分コピーが終了すると、論理パス602Aがパーティション1__300Aに移動し、続いて、仮想サーバ601Aがパーティション1__300Aに移動する。

【0148】

50

図15は、本発明の実施形態において実行されるアプリケーション500の移動及び電源制御のタイミングの説明図である。

【0149】

具体的には、図15は、図14に示した各処理の実行のタイミング及びその実行に要する時間を示す。図15に示された片方向の矢印の基点(図15の例では矢印の左端)、先端(図15の例では矢印の右端)及び長さは、それぞれ、各処理の開始時刻、終了時刻及び処理時間に対応する。

【0150】

例えば、矢印1501は、ストレージ装置100の起動処理を示す。矢印1501の基点は、ストレージ装置100の電源が投入された時刻(時刻1511)を示し、矢印1501の先端は、ストレージ装置100の起動処理が終了した時刻(時刻1512)を示す。同様に、矢印1502は論理ボリューム603Fの全体コピー処理を、矢印1503は論理ボリューム603Fの差分コピー処理を、矢印1504はサーバ120の起動処理を、矢印1505は、FC-SW140の起動処理を、矢印1506は論理パス602Aの移動処理を、矢印1507は仮想サーバ601Aの移動処理を示す。

10

【0151】

図14にも示したように、最初にストレージ装置100の電源が投入され(時刻1511)、ストレージ装置100が起動すると(時刻1512)、論理ボリューム603Fの全体コピーが実行される。全体コピーに要する時間は、論理ボリューム603Fに格納されているデータ量、及び、コピーに使用されるデータ転送経路の転送性能に依存するが、一般には、各ハードウェアの起動に要する時間及び差分コピーに要する時間のいずれと比較しても十分に長い場合が多い。

20

【0152】

全体コピーの終了(時刻1513)を契機として、サーバ120及びFC-SW140の電源が投入され、さらに、論理ボリューム603Fの差分コピーが開始される。

【0153】

サーバ120及びFC-SW140の起動が終了し、さらに差分コピーも終了すると(時刻1514)、論理パス602Aの移動が実行される。論理パス602Aの移動が終了すると(時刻1515)、仮想サーバ601Aの移動が実行される。

30

【0154】

従来は、全体コピー処理及び差分コピー処理が一つのコピー処理として実行されたため、管理サーバ130が全体コピー処理の終了時刻(すなわち時刻1513)を知ることができなかった。このため、全体コピーの終了を契機とする制御を実行することができなかった。その場合、制御の基準として、例えば時刻1511又は時刻1514が使用される。

【0155】

時刻1511にサーバ120及びFC-SW140の電源が投入された場合、図15の矢印1504及び矢印1505は、それらの基点が時刻1511となるように移動する。しかし、論理パス602A及び仮想サーバ601Aの移動が可能になるのは時刻1514以降であるので、サーバ120及びFC-SW140は、それぞれの起動が終了した後、時刻1514までの時間、何も実行せずにただ電力を消費しながら待機する。すなわち、その時間、サーバ120及びFC-SW140は電力を浪費する。

40

【0156】

一方、時刻1514(すなわち差分コピーが終了した時刻)にサーバ120及びFC-SW140の電源が投入された場合、図15の矢印1506及び1507が後の時刻(すなわち右方向)に移動する。この場合、サーバ120及びFC-SW140の起動処理が終了するまで論理パス602A及び仮想サーバ601Aの移動を実行することができない。その結果、ストレージ装置100の電源が投入されてからアプリケーション500Aの移動が終了するまでの時間が図15に示す場合より長くなる。すなわち、差分コピー処理が終了してからサーバ120及びFC-SW140の起動処理が終了するまでの時間、ス

50

ストレージ装置 100 が電力を浪費する。

【0157】

計算機システム全体の消費電力量を最小にするには、差分コピー処理、サーバ 120 の起動処理及び FC-SW140 の起動処理が同時刻に終了するように（すなわち、矢印 1503、1504 及び 1505 の先端が同一の時刻に対応するように）、サーバ 120 及び FC-SW140 の電源投入を制御することが望ましい。しかし、既に述べたように、差分コピー処理に要する時間を正確に予測することは困難である。このため、本実施形態では、全体コピー処理の終了を契機としてサーバ 120 及び FC-SW140 の電源が投入される。

【0158】

差分ストックとして格納されているデータ量が多ければ、図 15 に示すように、サーバ 120 及び FC-SW140 の起動が終了してもまだ差分コピーが終了しない場合がある。その場合、サーバ 120 及び FC-SW140 の起動が終了した後、差分コピーが終了するまでの時間、サーバ 120 及び FC-SW140 が電力を浪費する。一方、差分ストックとして格納されているデータ量が少なければ、サーバ 120 及び FC-SW140 の起動が終了する前に差分コピーが終了する場合がある。その場合、差分コピーが終了した後、サーバ 120 及び FC-SW140 の起動が終了するまでの間、ストレージ装置 100 が電力を浪費する。

【0159】

このように、全体コピー処理の終了時刻は、サーバ 120 及び FC-SW140 の電源投入のタイミングとして厳密に最適であるとは限らない。しかし、一般には、全体コピーに要する時間と差分コピーに要する時間とをあわせた論理ボリューム 603 の移動時間は、サーバ 120 及び FC-SW140 の起動処理に要する時間より大幅に長い。このため、サーバ 120 及び FC-SW140 の電源の最適な投入時刻は、論理ボリューム 603 の移動処理の終了間際であることが多い。これに対して、一般に、全体コピーに要する時間は、差分コピー処理に要する時間より大幅に長いため、全体コピー処理の終了時刻は、論理ボリューム 603 の移動処理の終了間際であることが多い。このため、全体コピー処理の終了時刻（時刻 1513）を、サーバ 120 及び FC-SW140 の電源を投入する近似的に最適な時刻として用いることができる。

【0160】

ストレージ装置 100 と同時にサーバ 120 及び FC-SW140 の電源を投入した場合、又は、差分コピー処理が終了した後にサーバ 120 及び FC-SW140 の電源を投入した場合のいずれの場合と比較しても、全体コピー処理の終了を契機として（具体的には、全体コピー処理が終了した後、差分コピー処理が終了する前に）サーバ 120 及び FC-SW140 の電源を投入することによって、計算機システム全体の消費電力を削減することができる。

【0161】

以下、上記のようなアプリケーション 500 の移動の処理について詳細に説明する。

【0162】

図 16 は、本発明の実施形態の管理テーブル 136 の説明図である。

【0163】

本実施形態の管理テーブル 136 は、パーティション管理テーブル 136A 及びアプリケーション管理テーブル 136B を含む。説明を簡単にするため、図 16 には、図 14 の例においてアプリケーション 500A の移動が実行される前の時点のパーティション管理テーブル 136A 及びアプリケーション管理テーブル 136B を示す。

【0164】

図 16(A) は、パーティション管理テーブル 136A の説明図である。パーティション管理テーブル 136A は、管理サーバ 130 が管理する計算機システム上に定義された各パーティションを管理するための情報を含む。

【0165】

10

20

30

40

50

具体的には、パーティション管理テーブル 1 3 6 A は、パーティション番号 1 6 0 1、ハードウェア種別 1 6 0 2、ハードウェア名 1 6 0 3、リソース量 1 6 0 4、残りリソース量 1 6 0 5、電源 1 6 0 6、配置アプリケーション (A P P) 1 6 0 7 及びアプリケーションリソース量 1 6 0 8 を含む。

【 0 1 6 6 】

パーティション番号 1 6 0 1 は、計算機システム上に定義された各パーティションを識別する情報である。図 1 6 (A) の例では、パーティション番号 1 6 0 1 として、「 1 」及び「 2 」が保持される。この例において、「 1 」は、パーティション 1 _ 3 0 0 A の識別子であり、「 2 」は、パーティション 2 _ 3 0 0 B の識別子である。

【 0 1 6 7 】

ハードウェア種別 1 6 0 2 は、各パーティションに割り当てられたハードウェアの種類を識別する情報である。具体的には、ハードウェア種別 1 6 0 2 は、各パーティションに割り当てられた各ハードウェアがサーバ 1 2 0、F C - S W 1 4 0 又はストレージ装置 1 0 0 のいずれであるかを識別する情報である。

【 0 1 6 8 】

図 1 6 (A) の例において、「サーバ 1」、「F C - S W 1」及び「ストレージ 1」は、それぞれ、パーティション 1 _ 3 0 0 A のサーバ 1 2 0、F C - S W 1 4 0 及びストレージ装置 1 0 0 を示す。一方、「サーバ 2」、「F C - S W 2」及び「ストレージ 2」は、それぞれ、パーティション 2 _ 3 0 0 B のサーバ 1 2 0、F C - S W 1 4 0 及びストレージ装置 1 0 0 を示す。

【 0 1 6 9 】

ハードウェア名 1 6 0 3 は、各パーティションに割り当てられたハードウェアを識別する情報である。図 1 6 (A) の例では、パーティション 1 _ 3 0 0 A に対応するハードウェア名 1 6 0 3 として、「サーバ P 1」(エントリ 1 6 1 1)、「F C - S W P 1」(エントリ 1 6 1 2) 及び「ストレージ P 1」(エントリ 1 6 1 3) が保持される。この例において、「サーバ P 1」、「F C - S W P 1」及び「ストレージ P 1」は、それぞれ、パーティション 1 _ 3 0 0 A に割り当てられたサーバ 1 2 0、F C - S W 1 4 0 及びストレージ装置 1 0 0 の識別子である。

【 0 1 7 0 】

さらに、図 1 6 (A) の例では、パーティション 2 _ 3 0 0 B に対応するハードウェア名 1 6 0 3 として、「サーバ P 2」(エントリ 1 6 1 4)、「F C - S W P 2」(エントリ 1 6 1 5) 及び「ストレージ P 2」(エントリ 1 6 1 6) が保持される。この例において、「サーバ P 2」、「F C - S W P 2」及び「ストレージ P 2」は、それぞれ、パーティション 2 _ 3 0 0 B に割り当てられたサーバ 1 2 0、F C - S W 1 4 0 及びストレージ装置 1 0 0 の識別子である。

【 0 1 7 1 】

図 1 6 (A) には、図 1 4 に示したものと整合するように、一つのパーティションに一つのサーバ 1 2 0、一つの F C - S W 1 4 0 及び一つのストレージ装置 1 0 0 が割り当てられる例を示す。しかし、実際には、一つのパーティションに複数のサーバ 1 2 0、複数の F C - S W 1 4 0 及び複数のストレージ装置 1 0 0 が割り当てられてもよい。その場合、それらの複数のハードウェアの識別子がハードウェア名 1 6 0 3 として保持される。

【 0 1 7 2 】

例えば、パーティション 1 _ 3 0 0 A に二つのサーバ 1 2 0 が割り当てられ、それぞれの識別子が「サーバ P 1」及び「サーバ P 1 0」(図示省略) である場合、ハードウェア種別 1 6 0 2 の値「サーバ 1」に二つのエントリが対応付けられ、それぞれのエントリのハードウェア名 1 6 0 3 として「サーバ P 1」及び「サーバ P 1 0」が保持される。

【 0 1 7 3 】

リソース量 1 6 0 4 は、各ハードウェアに含まれる全リソースの量である。リソースは、どのように計量されてもよい。例えば、サーバ 1 2 0 のリソース量は、そのサーバ 1 2 0 が備える C P U 1 2 1 の数であってもよいし、C P U 1 2 1 の使用率であってもよい。

10

20

30

40

50

FC - SW 140のリソース量は、そのFC - SW 140が備えるポート141の数であってもよいし、接続することができる論理パス602の数であってもよい。ストレージ装置100のリソース量は、論理ボリューム111として提供することができる記憶容量であってもよい。

【0174】

図16(A)の例では、各ハードウェアのリソース量1604として「10」が保持される。図16(A)の例では単位が省略されているが、実際には、各リソース量の単位(例えば、ストレージ装置100のリソース量の場合、「テラバイト」等)が明示されてもよい。

【0175】

残りリソース量1605は、各ハードウェアに含まれる全リソースのうち、まだいずれのアプリケーション500にも割り当てられていないものの量である。言い換えると、残りリソース量1605は、各ハードウェアに含まれる全リソースの量から、既にいずれかのアプリケーション500に割り当てられたものの量を減算した残りの量である。

【0176】

電源1606は、各ハードウェアの電源の状態を示す情報である。図16(A)において、電源1606の値「ON」は電源が投入されている状態、「OFF」は電源が遮断されている状態を示す。

【0177】

配置アプリケーション1607は、各パーティションに配置されたアプリケーション500を識別する情報である。

【0178】

アプリケーションリソース量1608は、各ハードウェアに含まれる全リソースのうち、アプリケーション500に割り当てられたリソースの量である。

【0179】

図16(A)の例は、図14においてアプリケーション500Aがパーティション2__300Bからパーティション1__300Aに移動する前の時点に対応する。すなわち、この時点においてパーティション1__300A内の各ハードウェアのリソースはまだいずれのアプリケーション500にも割り当てられていない。

【0180】

このため、パーティション1__300A内の各ハードウェアに対応する残りリソース量1605の値は、リソース量1604の値と同じである。パーティション1__300A内の各ハードウェアの電源はまだ投入されていないため、対応する電源1606として「OFF」が保持される。さらに、パーティション1__300A内の各ハードウェアに対応する配置アプリケーション1607及びアプリケーションリソース量1608は空白である。

【0181】

一方、図14のアプリケーション500Aが移動する前の時点において、パーティション2__300Bにはアプリケーション500Aが配置されている。このため、パーティション2__300Bに対応する配置アプリケーション1607として、アプリケーション500Aの識別子(図16(A)の例では「APP1」)が保持される。この状態において、パーティション2__300B内の各ハードウェアの電源が投入されているため、各ハードウェアに対応する電源1606として「ON」が保持される。

【0182】

さらに、パーティション2__300Bにおいて、サーバ120の全リソース量「10」のうち「5」、FC - SW 140の全リソース量「10」のうち「1」、ストレージ装置100の全リソース量「10」のうち「6」がアプリケーション500に割り当てられている場合、パーティション2__300Bのサーバ120、FC - SW 140及びストレージ装置100に対応するアプリケーションリソース量1608として、それぞれ「5」、「1」及び「6」が保持される。そしてそれらに対応する残りリソース量1605として、それぞれ「5」、「9」及び「4」が保持される。

10

20

30

40

50

【 0 1 8 3 】

図 1 6 (B) は、アプリケーション管理テーブル 1 3 6 B の説明図である。アプリケーション管理テーブル 1 3 6 B は、管理サーバ 1 3 0 が管理する計算機システム上で稼動する各アプリケーション 5 0 0 を管理するための情報を含む。

【 0 1 8 4 】

具体的には、アプリケーション管理テーブル 1 3 6 B は、アプリケーション名 1 6 2 1、配置パーティション番号 1 6 2 2、配置サーバ名 1 6 2 3、配置 FC - SW 名 1 6 2 4、配置ストレージ名 1 6 2 5、サーバリソース量 1 6 2 6、FC - SW リソース量 1 6 2 7 及びストレージリソース量 1 6 2 8 を含む。

【 0 1 8 5 】

アプリケーション名 1 6 2 1 は、計算機システム内で稼動するアプリケーション 5 0 0 を識別する情報である。

【 0 1 8 6 】

配置パーティション番号 1 6 2 2 は、各アプリケーション 5 0 0 が配置されているパーティションを識別する情報である。

【 0 1 8 7 】

配置サーバ名 1 6 2 3 は、各アプリケーション 5 0 0 の仮想サーバ 6 0 1 が配置されているサーバ 1 2 0 を識別する情報である。

【 0 1 8 8 】

配置 FC - SW 名 1 6 2 4 は、各アプリケーション 5 0 0 の論理パス 6 0 2 が配置されている FC - SW 1 4 0 を識別する情報である。

【 0 1 8 9 】

配置ストレージ名 1 6 2 5 は、各アプリケーション 5 0 0 の論理ボリューム 6 0 3 が配置されているストレージ装置 1 0 0 を識別する情報である。

【 0 1 9 0 】

サーバリソース量 1 6 2 6 は、各アプリケーション 5 0 0 に割り当てられているサーバ 1 2 0 のリソース量を示す。

【 0 1 9 1 】

FC - SW リソース量 1 6 2 7 は、各アプリケーション 5 0 0 に割り当てられている FC - SW 1 4 0 のリソース量を示す。

【 0 1 9 2 】

ストレージリソース量 1 6 2 8 は、各アプリケーション 5 0 0 に割り当てられているストレージ装置 1 0 0 のリソース量を示す。

【 0 1 9 3 】

図 1 6 (B) の例は、図 1 4 においてアプリケーション 5 0 0 A がパーティション 2 __ 3 0 0 B からパーティション 1 __ 3 0 0 A に移動する前の時点に対応する。すなわち、この時点においてアプリケーション 5 0 0 A は、パーティション 2 __ 3 0 0 B に配置されている。このため、図 1 6 (B) の例において、アプリケーション名 1 6 2 1 の値「APP 1」に対応する配置パーティション番号 1 6 2 2、配置サーバ名 1 6 2 3、配置 FC - SW 名 1 6 2 4 及び配置ストレージ名 1 6 2 5 として、それぞれ、「2」、「サーバ P 2」、「FC - SW P 2」及び「ストレージ P 2」が保持される。これらの値は、図 1 6 (A) の例と整合する。さらに、図 1 6 (A) に示すように、サーバリソース量 1 6 2 6、FC - SW リソース量 1 6 2 7 及びストレージリソース量 1 6 2 8 として、それぞれ、「5」、「1」及び「6」が保持される。

【 0 1 9 4 】

次に、アプリケーション 5 0 0 の移動及びそれに伴う電源制御の処理について、フローチャートを参照して説明する。以下の説明において、処理の具体例として、図 1 4 及び図 1 5 に示したアプリケーション 5 0 0 A の移動処理及びそのための電源制御処理を参照する。

【 0 1 9 5 】

10

20

30

40

50

図 17 は、本発明の実施形態において実行されるアプリケーション 500 の移動及び電源制御の処理の全体を示すフローチャートである。

【0196】

この処理は、図 10 のステップ 1002 において実行される。

【0197】

管理サーバ 130 の管理プログラム 135 は、アプリケーション 500 の移動処理を開始すると（ステップ 1701）、最初に、ストレージ電源 ON 処理を実行する（ステップ 1702）。これは、移動先のストレージ装置 100（図 14 の例では、パーティション 1_300A のストレージ装置 100）の電源を投入する処理である。

【0198】

ステップ 1702 において実行されるストレージ電源 ON 処理については、後に図 18 を参照して詳細に説明する。

【0199】

ストレージ電源 ON 処理が終了する（すなわち、移動先のストレージ装置 100 の起動処理が終了する）と、管理プログラム 135 は、論理ボリューム 603 の全体コピー処理を実行する（ステップ 1703）。ステップ 1703 において実行される全体コピー処理については、後に図 19 及び図 20 を参照して詳細に説明する。

【0200】

ステップ 1703 の全体コピー処理が終了すると、管理プログラム 135 は、サーバ電源 ON 処理を実行する（ステップ 1704）。これは、移動先のサーバ 120 の電源を投入する処理である。ステップ 1704 において実行されるサーバ電源 ON 処理については、後に図 21 を参照して詳細に説明する。

【0201】

ステップ 1703 の全体コピー処理が終了すると、管理プログラム 135 は、さらに、FC-SW 電源 ON 処理を実行する（ステップ 1705）。これは、移動先の FC-SW の電源を投入する処理である。ステップ 1705 において実行される FC-SW 電源 ON 処理については、後に図 22 を参照して詳細に説明する。

【0202】

ステップ 1703 の全体コピー処理が終了すると、管理プログラム 135 は、さらに、論理ボリューム 603 の差分コピー処理を実行する（ステップ 1706）。ステップ 1706 において実行される差分コピー処理については、後に図 23 及び図 24 を参照して詳細に説明する。

【0203】

管理プログラム 135 は、全体コピー処理の終了を、ストレージ装置 100 から受信した終了メッセージに基づいて判定してもよいし、管理プログラム 135 が算出した全体コピー処理時間に基づいて判定してもよい。

【0204】

前者の場合、管理サーバ 130 がストレージ装置 100 から終了メッセージを受信した時刻が全体コピー処理の終了時刻として扱われ、その時刻を契機にステップ 1704 ~ ステップ 1706 が実行される。終了メッセージについては後述する（図 20 参照）。

【0205】

後者の場合、全体コピー処理が開始されてから、全体コピー処理時間が経過した時刻が全体コピー処理の終了時刻として扱われ、その時刻を契機にステップ 1704 ~ ステップ 1706 が実行される。全体コピー処理時間は、コピー元の論理ボリューム 603 F の容量を、論理ボリューム 603 F からコピー先の論理ボリューム 603 A へのデータ転送速度によって除算することによって算出される。データ転送速度は、ハードウェアのスペックに基づいて算出されてもよいし、実測された値であってもよい。

【0206】

なお、図 17 は、便宜上、ステップ 1704 からステップ 1706 が順次実行されるように記載されているが、この順序は一例にすぎない。これらの三つのステップは、ステッ

10

20

30

40

50

ブ 1703 の全体コピー処理の終了を契機として実行される必要があるが、どのような順序で実行されてもよい。可能であれば、これらの三つのステップが同時に実行されてもよい。ただし、図 15 を参照して説明したように、本実施形態による消費電力削減の効果を得るためには、ステップ 1704 及びステップ 1705 は、遅くとも差分コピー処理が終了する前に開始する必要がある。

【0207】

サーバ電源 ON 処理、FC - SW 電源 ON 処理及び差分コピー処理（ステップ 1704 ~ ステップ 1706）が全て終了すると、管理プログラム 135 は、論理パス 602 の移動処理を実行する（ステップ 1707）。この処理については、後に図 25 を参照して詳細に説明する。

10

【0208】

ステップ 1707 の論理パス 602 の移動処理が終了すると、管理プログラム 135 は、仮想サーバ 601 の移動処理を実行する（ステップ 1708）。この処理については、後に図 26 を参照して詳細に説明する。

【0209】

ステップ 1708 の仮想サーバ 601 の移動処理が終了すると、管理プログラム 135 は、アプリケーション 500 の移動処理を終了する（ステップ 1709）。なお、仮想サーバ 601 の移動処理が終了した後（例えばステップ 1709 において）、移動元のパーティションのハードウェア（図 14 の例では、パーティション 2 __ 300 B のサーバ 120、FC - SW 140 及びストレージ装置 100）の電源が遮断されてもよい。

20

【0210】

図 18 は、本発明の実施形態において実行されるストレージ電源 ON 処理を示すフローチャートである。

【0211】

この処理は、図 17 のステップ 1702 において実行される。

【0212】

管理プログラム 135 は、ストレージ電源 ON 処理が開始されると（ステップ 1801）、対象ハードウェアを確認する（ステップ 1802）。具体的には、管理プログラム 135 は、パーティション管理テーブル 136 A を参照して、移動先のパーティションに含まれるストレージ装置 100（すなわち移動先のストレージ装置）の状態を確認する。

30

【0213】

例えば、図 14 のようにパーティション 1 __ 300 A が移動先として指定されている場合、管理プログラム 135 は、パーティション管理テーブル 136 A を参照して、パーティション番号 1601 の値「1」及びハードウェア種別 1602 の値「ストレージ 1」に対応する全てのエントリを特定する。そして、管理プログラム 135 は、特定された各エントリの電源 1606 の値を確認する。

【0214】

次に、管理プログラム 135 は、移動先のストレージ装置 100 の電源が遮断された状態（すなわち「OFF」状態）であるか否かを判定する（ステップ 1803）。具体的には、管理プログラム 135 は、ステップ 1802 で取得した値が「ON」又は「OFF」のいずれであるかを判定する。

40

【0215】

移動先のストレージ装置 100 の電源が既に投入されている場合、さらに電源を投入する処理を実行する必要はないため、管理プログラム 135 は、ストレージ電源 ON 処理を終了する（ステップ 1806）。

【0216】

移動先のストレージ装置 100 の電源が遮断されている場合、管理プログラム 135 は、管理ネットワーク 150 を介して、移動先のストレージ装置 100 に電源 ON 命令を送信する（ステップ 1804）。ストレージ装置 100 の電源制御部 105 は、I / F 104 を介して電源 ON 命令を受信すると、ストレージ装置 100 の電源を投入する。なお、

50

このような処理を実行するために、ストレージ装置 100のうち、少なくともI/F 104及び電源制御部 105の電源は、ストレージ電源ON処理が実行される時点で既に投入されている必要がある。

【0217】

ストレージ装置 100は、電源ON命令に従って電源が投入され、起動処理が終了すると、起動処理の終了の報告を管理サーバ 130に送信してもよい。

【0218】

次に、管理プログラム 135は、ストレージ電源ON処理の結果が反映されるように、管理テーブル 136の電源に関する情報を更新する(ステップ 1805)。

【0219】

ステップ 1802で複数のエントリが特定された場合(すなわち、移動先のパーティションが複数のストレージ装置 100を含む場合)、各エントリについてステップ 1803~ステップ 1805が実行される。

【0220】

以上でストレージ電源ON処理が終了する(ステップ 1806)。

【0221】

図 16(A)の例では、ステップ 1802において、ハードウェア名「ストレージ P1」を含むエントリ 1613のみが特定される。エントリ 1613の電源 1606の値が「OFF」であるため、ステップ 1804において、管理プログラム 135は、移動先のストレージ装置 100(すなわち「ストレージ P1」によって識別されるストレージ装置 100)に電源ON命令を送信する。これによって移動先のパーティション 1__300Aに含まれるストレージ装置 100の電源が投入され、ステップ 1805において、エントリ 1613の電源 1606の値が「ON」に更新される。

【0222】

図 19は、本発明の実施形態の管理サーバ 130が実行する論理ボリューム 603の全体コピー処理のフローチャートである。

【0223】

論理ボリューム 603の全体コピー処理は、図 15を参照して説明したように、論理ボリューム 603の移動処理の一部である。この処理は、図 17のステップ 1703において実行される。

【0224】

管理プログラム 135は、論理ボリューム 603の移動処理を開始すると(ステップ 1901)、論理ボリューム 603の全体コピー処理を開始し(ステップ 1902)、全体コピーを実行して(ステップ 1903)、全体コピー処理を終了する(ステップ 1904)。

【0225】

具体的には、管理プログラム 135は、ステップ 1903において、全体コピー命令を移動元のストレージ装置 100に送信する。この命令を受信したストレージ装置 100が実行する処理については、後に図 20を参照して詳細に説明する。

【0226】

図 20は、本発明の実施形態のストレージ装置 100が実行する論理ボリューム 603の全体コピー処理のフローチャートである。

【0227】

この処理は、図 19のステップ 1903において送信された命令を受信したストレージ装置 100(図 14の例では、パーティション 2__300Bのストレージ装置 100)のコントローラ 101によって実行される。

【0228】

コントローラ 101は、論理ボリューム 603の全体コピー処理を開始すると(ステップ 2001)、差分I/Oを格納する記憶領域を確保する(ステップ 2002)。この記憶領域は、例えば、ストレージ装置 100が管理するいずれかの論理ボリューム 111の

10

20

30

40

50

空き記憶領域に確保されてもよい。

【0229】

次に、コントローラ101は、これから移動する論理ボリューム603へのデータI/Oの実行先を、ステップ2002において確保した記憶領域に切り替える(ステップ2003)。その後、コントローラ101は、これから移動する論理ボリューム603(図14の例では、論理ボリューム603F)へのデータI/Oを受信すると、それによって更新されるデータを論理ボリューム603に反映せず、ステップ2002において確保した記憶領域に差分ストックとして格納する(図11参照)。

【0230】

次に、コントローラ101は、論理ボリューム603の全体コピーを実行する(ステップ2004)。図14の例では、コントローラ101は、論理ボリューム603Fの全データを読み出し、読み出したデータを論理ボリューム603Aに書き込む要求をパーティション1_300Aのストレージ装置100に送信する。なお、コピー元及びコピー先の論理ボリューム603は、図19のステップ1903において送信された命令に含まれる引数によって指定される。

10

【0231】

読み出した全データの送信が終了すると、コントローラ101は、論理ボリューム603の全体コピー処理を終了する(ステップ2005)。このとき、コントローラ101は、全体コピー処理が終了したことを示す終了メッセージを管理サーバ130に送信する。

【0232】

管理プログラム135は、ステップ2005で送信された終了メッセージを受信したことを契機として、ステップ1704~ステップ1706を実行する。

20

【0233】

なお、図17を参照して説明したように、管理プログラム135は、ハードウェアスペック等に基づいて全体コピー処理時間を算出し、その算出された全体コピー処理時間に基づいてステップ1704~ステップ1706の実行を制御してもよい。その場合、コントローラ101は、ステップ2005において終了メッセージを送信しなくてもよい。

【0234】

図21は、本発明の実施形態において実行されるサーバ電源ON処理を示すフローチャートである。

30

【0235】

この処理は、図17のステップ1704において実行される。すなわち、この処理は、図20のステップ2005において送信された終了メッセージを管理サーバ130が受信したことを契機として実行される。

【0236】

管理プログラム135は、ストレージ電源ON処理が開始されると(ステップ2101)、対象ハードウェアを確認する(ステップ2102)。具体的には、管理プログラム135は、パーティション管理テーブル136Aを参照して、移動先のパーティションに含まれるサーバ120(すなわち移動先のサーバ)の状態を確認する。

【0237】

例えば、図14のようにパーティション1_300Aが移動先として指定されている場合、管理プログラム135は、パーティション管理テーブル136Aを参照して、パーティション番号1601の値「1」及びハードウェア種別1602の値「サーバ1」に対応する全てのエントリを特定する。そして、管理プログラム135は、特定された各エントリの電源1606の値を確認する。

40

【0238】

次に、管理プログラム135は、移動先のサーバ120の電源が遮断された状態(すなわち「OFF」状態)であるか否かを判定する(ステップ2103)。具体的には、管理プログラム135は、ステップ2102で取得した値が「ON」又は「OFF」のいずれであるかを判定する。

50

【0239】

移動先のサーバ120の電源が既に投入されている場合、さらに電源を投入する処理を実行する必要はないため、管理プログラム135は、サーバ電源ON処理を終了する(ステップ2106)。

【0240】

移動先のサーバ120の電源が遮断されている場合、管理プログラム135は、管理ネットワーク150を介して、移動先のサーバ120に電源ON命令を送信する(ステップ2104)。サーバ120の電源制御部123は、I/F125を介して電源ON命令を受信すると、サーバ120の電源を投入する。なお、このような処理を実行するために、サーバ120のうち、少なくともI/F125及び電源制御部123の電源は、サーバ電源ON処理が実行される時点で既に投入されている必要がある。

10

【0241】

サーバ120は、電源ON命令に従って電源が投入され、起動処理が終了すると、起動処理の終了の報告を管理サーバ130に送信してもよい。

【0242】

次に、管理プログラム135は、サーバ電源ON処理の結果が反映されるように、管理テーブル136の電源に関する情報を更新する(ステップ2105)。

【0243】

ステップ2102で複数のエントリが特定された場合(すなわち、移動先のパーティションが複数のサーバ120を含む場合)、各エントリについてステップ2103~ステップ2105が実行される。

20

【0244】

以上でサーバ電源ON処理が終了する(ステップ2106)。

【0245】

図16(A)の例では、ステップ2102において、ハードウェア名「サーバP1」を含むエントリ1611のみが特定される。エントリ1611の電源1606の値が「OFF」であるため、ステップ2104において、管理プログラム135は、移動先のサーバ120(すなわち「サーバP1」によって識別されるサーバ120)に電源ON命令を送信する。これによって移動先のパーティション1__300Aに含まれるサーバ120の電源が投入され、ステップ2105において、エントリ1611の電源1606の値が「ON」に更新される。

30

【0246】

なお、図14の例のように、各パーティションに一つ以上のサーバ120が割り当てられている場合、ステップ2104において、移動先のパーティションのサーバ120の電源を投入する命令(すなわち、そのサーバ120に含まれるリソース全体の電源を投入する命令)が送信される。しかし、例えば図4に示すように、一つのサーバ120のリソースの各部分(例えば各CPU121)が各パーティションに割り当てられ、かつ、その部分ごとに電源を投入することができる場合、ステップ2104では、移動先のパーティションに含まれるサーバ120のリソースの電源を投入する命令が、そのリソースを含むサーバ120に送信される。この命令を受信したサーバ120の電源制御部123は、命令によって指定されたリソースの部分(例えば指定されたCPU)の電源を投入する。

40

【0247】

図22は、本発明の実施形態において実行されるFC-SW電源ON処理を示すフローチャートである。

【0248】

この処理は、図17のステップ1705において実行される。すなわち、この処理は、図20のステップ2005において送信された終了メッセージを管理サーバ130が受信したことを契機として実行される。

【0249】

管理プログラム135は、FC-SW電源ON処理が開始されると(ステップ2201

50

)、対象ハードウェアを確認する(ステップ2202)。具体的には、管理プログラム135は、パーティション管理テーブル136Aを参照して、移動先のパーティションに含まれるFC-SW140(すなわち移動先のFC-SW)の状態を確認する。

【0250】

例えば、図14のようにパーティション1__300Aが移動先として指定されている場合、管理プログラム135は、パーティション管理テーブル136Aを参照して、パーティション番号1601の値「1」及びハードウェア種別1602の値「FC-SW1」に対応する全てのエントリを特定する。そして、管理プログラム135は、特定された各エントリの電源1606の値を確認する。

【0251】

次に、管理プログラム135は、移動先のFC-SW140の電源が遮断された状態(すなわち「OFF」状態)であるか否かを判定する(ステップ2203)。具体的には、管理プログラム135は、ステップ2202で取得した値が「ON」又は「OFF」のいずれであるかを判定する。

【0252】

移動先のFC-SW140の電源が既に投入されている場合、さらに電源を投入する処理を実行する必要はないため、管理プログラム135は、FC-SW電源ON処理を終了する(ステップ2206)。

【0253】

移動先のFC-SW140の電源が遮断されている場合、管理プログラム135は、管理ネットワーク150を介して、移動先のFC-SW140に電源ON命令を送信する(ステップ2204)。FC-SW140の電源制御部143は、I/F142を介して電源ON命令を受信すると、FC-SW140の電源を投入する。なお、このような処理を実行するために、FC-SW140のうち、少なくともI/F142及び電源制御部143の電源は、FC-SW電源ON処理が実行される時点で既に投入されている必要がある。

【0254】

FC-SW140は、電源ON命令に従って電源が投入され、起動処理が終了すると、起動処理の終了の報告を管理サーバ130に送信してもよい。

【0255】

次に、管理プログラム135は、FC-SW電源ON処理の結果が反映されるように、管理テーブル136の電源に関する情報を更新する(ステップ2205)。

【0256】

ステップ2202で複数のエントリが特定された場合(すなわち、移動先のパーティションが複数のFC-SW140を含む場合)、各エントリについてステップ2203~ステップ2205が実行される。

【0257】

以上でFC-SW電源ON処理が終了する(ステップ2206)。

【0258】

図16(A)の例では、ステップ2202において、ハードウェア名「FC-SWP1」を含むエントリ1612のみが特定される。エントリ1612の電源1606の値が「OFF」であるため、ステップ2204において、管理プログラム135は、移動先のFC-SW140(すなわち「FC-SWP1」によって識別されるFC-SW140)に電源ON命令を送信する。これによって移動先のパーティション1__300Aに含まれるFC-SW140の電源が投入され、ステップ2205において、エントリ1612の電源1606の値が「ON」に更新される。

【0259】

図23は、本発明の実施形態の管理サーバ130が実行する論理ボリューム603の差分コピー処理のフローチャートである。

【0260】

10

20

30

40

50

論理ボリューム 603 の差分コピー処理は、図 15 を参照して説明したように、論理ボリューム 603 の移動処理の一部である。この処理は、図 17 のステップ 1706 において実行される。

【0261】

管理プログラム 135 は、論理ボリューム 603 の差分コピー処理を開始すると（ステップ 2301）、差分コピーを実行する（ステップ 2302）。具体的には、管理プログラム 135 は、差分コピー命令を移動元のストレージ装置 100 に送信する。この命令を受信したストレージ装置 100 が実行する処理については、後に図 24 を参照して詳細に説明する。

【0262】

次に、管理プログラム 135 は、実行されたコピーの結果が反映されるように、管理テーブル 136 を更新する（ステップ 2303）。例えば、図 14 に示すように論理ボリューム 603 F から論理ボリューム 603 A へのコピーが実行された場合、パーティション管理テーブル 136 A のうち、ストレージ装置 100 に関するエントリ 1613 及び 1616 が更新される。

【0263】

図 16 (A) の例では、論理ボリューム 603 A を含むストレージ装置 100 がエントリ 1613 に対応し、論理ボリューム 603 F を含むストレージ装置 100 がエントリ 1616 に対応する。この場合、コピー前の論理ボリューム 603 F が使用しているリソース量が「6」である。論理ボリューム 603 F から論理ボリューム 603 A へのコピーによって、アプリケーション 500 A の一部としての論理ボリューム 603 F はリソースを使用しなくなり、一方、論理ボリューム 603 A は、アプリケーション 500 A の一部として新たにリソースを使用するようになる。

【0264】

このため、ステップ 2303 において、エントリ 1613 の残りリソース量 1605、配置アプリケーション 1607 及びアプリケーションリソース量 1608 が、それぞれ、「4」、「APP1」及び「6」に更新される。一方、エントリ 1616 の残りリソース量 1605 は「10」に更新され、エントリ 1616 の配置アプリケーション 1607 及びアプリケーションリソース量 1608 は空白に更新される。

【0265】

さらに、ステップ 2303 において、アプリケーション管理テーブル 136 B も更新される。具体的には、アプリケーション名「APP1」に対応する配置ストレージ名 1625 が「ストレージ P1」に更新される。

【0266】

管理プログラム 135 は、ステップ 2303 を実行すると、論理ボリューム 603 の差分コピー処理を終了する（ステップ 2304）。これによって、論理ボリューム 603 の移動処理が終了する（ステップ 2305）。

【0267】

図 24 は、本発明の実施形態のストレージ装置 100 が実行する論理ボリューム 603 の差分コピー処理のフローチャートである。

【0268】

この処理は、図 23 のステップ 2302 において送信された命令を受信したストレージ装置 100（図 14 の例では、パーティション 2__300 B のストレージ装置 100）のコントローラ 101 によって実行される。

【0269】

コントローラ 101 は、論理ボリューム 603 の差分コピー処理を開始すると（ステップ 2401）、コピー元の論理ボリューム 603（図 14 の例では論理ボリューム 603 F）へのデータ I/O を停止する（ステップ 2402）。ただし、既に図 20 のステップ 2003 においてデータ I/O の実行先が切り替えられているため、ステップ 2402 では、差分データ I/O の記憶領域へのデータ I/O が停止される。

10

20

30

40

50

【0270】

次に、コントローラ101は、差分コピーを実行する(ステップ2403)。具体的には、コントローラ101は、差分データI/Oの記憶領域からデータを読み出し、読み出したデータをコピー先の論理ボリューム603Aに書き込む要求をパーティション1__300Aのストレージ装置100に送信する。

【0271】

全ての差分データの送信が終了すると、コントローラ101は、確保した差分データI/Oの記憶領域を削除する(ステップ2404)。

【0272】

以上で、コントローラ101は、論理ボリューム603の差分コピー処理を終了する(ステップ2405)。

【0273】

なお、ステップ2402において停止されたデータI/Oは、後述する論理パス602の移動処理によって再開される(図25参照)。

【0274】

図25は、本発明の実施形態において実行される論理パス602の移動処理のフローチャートである。

【0275】

この処理は、図17のステップ1707において実行される。

【0276】

管理プログラム135は、論理パス602の移動処理を開始すると(ステップ2501)、論理パス602を移動する(ステップ2502)。例えば、図14に示すように論理パス602Aを移動する場合、管理プログラム135は、論理パス602Aを移動させる命令を、パーティション2__300Bのサーバ120に送信してもよい。この命令を受信したサーバ120は、論理パス602Aがパーティション1__300AのFC-SW140を経由して論理ボリューム603Aに接続されるようにバス管理プログラム127の設定を変更してもよい。

【0277】

ステップ2502において論理パス602Aが移動すると、パーティション1__300Aのストレージ装置100は、論理ボリューム603AへのデータI/Oを再開する。

【0278】

次に、管理プログラム135は、ステップ2502における論理パス602の移動が反映されるように管理テーブル136を更新する(ステップ2503)。例えば、図14に示すように論理パス602Aが移動した場合、パーティション管理テーブル136Aのうち、FC-SW140に関するエントリ1612及び1615が更新される。

【0279】

図16(A)の例では、パーティション1__300AのFC-SW140がエントリ1612に対応し、パーティション2__300BのFC-SW140がエントリ1615に対応する。この場合、移動前の論理パス602Aが使用しているパーティション2__300BのFC-SW140のリソース量が「1」である。ステップ2502の移動の結果、論理パス602Aは、パーティション2__300BのFC-SW140のリソースを使用しなくなり、代わりにパーティション1__300AのFC-SW140のリソースを使用するようになる。

【0280】

このため、ステップ2503において、エントリ1612の残りリソース量1605、配置アプリケーション1607及びアプリケーションリソース量1608が、それぞれ、「9」、「APP1」及び「1」に更新される。一方、エントリ1615の残りリソース量1605は「10」に更新され、エントリ1615の配置アプリケーション1607及びアプリケーションリソース量1608は空白に更新される。

【0281】

10

20

30

40

50

さらに、ステップ 2503 において、アプリケーション管理テーブル 136B も更新される。具体的には、アプリケーション名「APP1」に対応する配置 FC-SW 名 1624 が「FC-SWP1」に更新される。

【0282】

以上で論理パス 602 の移動処理が終了する（ステップ 2504）。

【0283】

図 26 は、本発明の実施形態において実行される仮想サーバ 601 の移動処理のフローチャートである。

【0284】

この処理は、図 17 のステップ 1708 において実行される。

10

【0285】

管理プログラム 135 は、仮想サーバ 601 の移動処理を開始すると（ステップ 2601）、仮想サーバ 601 を移動する（ステップ 2602）。この移動は、例えば、図 9 において説明したように、メモリ 122 のイメージをコピーすることによって実行される。

【0286】

次に、管理プログラム 135 は、ステップ 2602 における論理パス 602 の移動が反映されるように管理テーブル 136 を更新する（ステップ 2603）。例えば、図 14 に示すように仮想サーバ 601A が移動した場合、パーティション管理テーブル 136A のうち、サーバ 120 に関するエントリ 1611 及び 1614 が更新される。

【0287】

20

図 16 (A) の例では、パーティション 1__300A のサーバ 120 がエントリ 1611 に対応し、パーティション 2__300B のサーバ 120 がエントリ 1614 に対応する。この場合、移動前の仮想サーバ 601A が使用しているパーティション 2__300B のサーバ 120 のリソース量が「1」である。ステップ 2602 の移動の結果、仮想サーバ 601A は、パーティション 2__300B のサーバ 120 のリソースを使用しなくなり、代わりにパーティション 1__300A のサーバ 120 のリソースを使用するようになる。

【0288】

このため、ステップ 2603 において、エントリ 1611 の残りリソース量 1605、配置アプリケーション 1607 及びアプリケーションリソース量 1608 が、それぞれ、「5」、「APP1」及び「5」に更新される。一方、エントリ 1614 の残りリソース量 1605 は「10」に更新され、エントリ 1614 の配置アプリケーション 1607 及びアプリケーションリソース量 1608 は空白に更新される。

30

【0289】

さらに、ステップ 2603 において、アプリケーション管理テーブル 136B も更新される。具体的には、アプリケーション名「APP1」に対応する配置サーバ名 1623 が「サーバ P1」に更新される。

【0290】

以上で仮想サーバ 601 の移動処理が終了する（ステップ 2604）。

【0291】

図 27 は、本発明の実施形態において実行される移動先パーティションの状態確認処理のフローチャートである。

40

【0292】

この処理は、図 8 のステップ 804 において実行される。

【0293】

管理プログラム 135 は、移動先パーティションの状態確認処理を開始すると（ステップ 2701）、これから移動しようとするアプリケーション 500 が現在配置されているパーティションを確認する（ステップ 2702）。この確認のために、アプリケーション管理テーブル 136B の配置パーティション番号 1622 から配置ストレージ名 1625 までが参照される。

【0294】

50

次に、管理プログラム 135 は、移動先のリソース、すなわち、移動先のパーティションを選択する（ステップ 2703）。

【0295】

次に、管理プログラム 135 は、パーティション管理テーブル 136A を参照して、移動先として選択されたパーティションに含まれるストレージ装置 100（以下、選択されたストレージ装置）の状態を確認する（ステップ 2704）。

【0296】

次に、管理プログラム 135 は、ステップ 2704 で確認した結果に基づいて、選択されたストレージ装置 100 へのアプリケーション 500 の移動（すなわち、論理ボリューム 603 のデータのコピー）が可能であるか否かを判定する（ステップ 2705）。

10

【0297】

選択されたストレージ装置 100 へのアプリケーション 500 の移動が可能であると判定された場合、管理プログラム 135 は、パーティション管理テーブル 136A を参照して、移動先として選択されたパーティションに含まれる FC-SW 140（以下、選択された FC-SW）の状態を確認する（ステップ 2706）。

【0298】

次に、管理プログラム 135 は、ステップ 2706 で確認した結果に基づいて、選択された FC-SW 140 へのアプリケーション 500 の移動（すなわち論理パス 602 の移動）が可能であるか否かを判定する（ステップ 2707）。

【0299】

20

選択された FC-SW 140 へのアプリケーション 500 の移動が可能であると判定された場合、管理プログラム 135 は、パーティション管理テーブル 136A を参照して、移動先として選択されたパーティションに含まれるサーバ 120（以下、選択されたサーバ）の状態を確認する（ステップ 2708）。

【0300】

次に、管理プログラム 135 は、ステップ 2708 で確認した結果に基づいて、選択されたサーバ 120 へのアプリケーション 500 の移動（すなわち仮想サーバ 601 の移動）が可能であるか否かを判定する（ステップ 2709）。

【0301】

ステップ 2709 において、選択されたサーバ 120 へのアプリケーション 500 の移動が可能であると判定された場合、結局、移動先として選択されたパーティションへのアプリケーション 500 の移動が可能である。このため、管理プログラム 135 は、アプリケーションが移動可能であることを示す応答を返して処理を終了する（ステップ 2710）。この場合、図 8 のステップ 805 において、選択したパーティションへのアプリケーション 500 の移動が可能であると判定される。

30

【0302】

一方、ステップ 2705、ステップ 2707 又はステップ 2709 の少なくとも一つにおいて、アプリケーション 500 の移動が可能でないと判定された場合、結局、移動先として選択されたパーティションへアプリケーション 500 を移動することはできない。このため、管理プログラム 135 は、アプリケーションを移動できないことを示す応答を返して処理を終了する（ステップ 2711）。この場合、図 8 のステップ 805 において、選択したパーティションへのアプリケーション 500 の移動が可能でないと判定される。

40

【0303】

ここで、図 14 に示すアプリケーション 500A の移動を例として、図 27 の処理の具体例を説明する。

【0304】

図 14 に示すアプリケーション 500A が移動する前の時点の管理テーブル 136 は、図 16 に示すとおりである。この場合、ステップ 2702 において、配置パーティション番号 1622 の値「2」、配置サーバ名 1623 の値「サーバ P2」、配置 FC-SW 名 1624 の値「FC-SWP2」及び配置ストレージ名 1625 の値「ストレージ P2」

50

から、アプリケーション 500 A が、パーティション 2 __ 300 B に配置されていることが確認される。

【0305】

そして、ステップ 2703 において、パーティション番号 1601 が参照され、現在アプリケーション 500 A が配置されているパーティション 2 __ 300 B 以外のパーティション、すなわち、パーティション 1 __ 300 A が移動先として選択される。

【0306】

ステップ 2704 において、パーティション 1 __ 300 A に含まれるストレージ装置 100 に対応するエントリ (図 16 (A) の例ではエントリ 1613) の残リソース量 1605 が参照される。

10

【0307】

ステップ 2705 において、ステップ 2704 で参照された残リソース量 1605 の値とストレージリソース量 1628 の値とが比較される。残リソース量 1605 の値がストレージリソース量 1628 の値より小さい場合には、アプリケーション 500 A を移動すると、移動先のパーティション 1 __ 300 A のリソース (例えばストレージ装置 100 の容量) が不足するため、アプリケーション 500 A を移動することはできない。図 16 の例では、残リソース量 1605 の値「10」がストレージリソース量 1628 の値「6」より大きいため、アプリケーション 500 A を移動可能であると判定される。

【0308】

ステップ 2706 において、パーティション 1 __ 300 A に含まれる FC - SW 140 に対応するエントリ (図 16 (A) の例ではエントリ 1612) の残リソース量 1605 が参照される。

20

【0309】

ステップ 2707 において、ステップ 2706 で参照された残リソース量 1605 の値と FC - SW リソース量 1627 の値とが比較される。残リソース量 1605 の値が FC - SW リソース量 1627 の値より小さい場合には、アプリケーション 500 A を移動すると、移動先のパーティション 1 __ 300 A のリソース (例えば FC - SW 140 のポート 141) が不足するため、アプリケーション 500 A を移動することはできない。図 16 の例では、残リソース量 1605 の値「10」が FC - SW リソース量 1627 の値「1」より大きいため、アプリケーション 500 A を移動可能であると判定される。

30

【0310】

ステップ 2708 において、パーティション 1 __ 300 A に含まれるサーバ 120 に対応するエントリ (図 16 (A) の例ではエントリ 1611) の残リソース量 1605 が参照される。

【0311】

ステップ 2709 において、ステップ 2708 で参照された残リソース量 1605 の値とサーバリソース量 1626 の値とが比較される。残リソース量 1605 の値がサーバリソース量 1626 の値より小さい場合には、アプリケーション 500 A を移動すると、移動先のパーティション 1 __ 300 A のリソース (例えばサーバ 120 の CPU 121) が不足するため、アプリケーション 500 A を移動することはできない。図 16 の例では、残リソース量 1605 の値「10」がサーバリソース量 1626 の値「5」より大きいため、アプリケーション 500 A を移動可能であると判定される。

40

【0312】

図 28 は、本発明の実施形態においてアプリケーション 500 の移動が実行された後の管理テーブル 136 の説明図である。

【0313】

具体的には、図 28 は、図 16 に示すパーティション管理テーブル 136 A 及びアプリケーション管理テーブル 136 B が、図 17 から図 26 に示す処理によって更新された結果を示す。すなわち、図 28 は、図 14 に示すアプリケーション 500 A の移動が実行された後のパーティション管理テーブル 136 A 及びアプリケーション管理テーブル 136

50

Bを示す。

【0314】

「APP1」によって識別されるアプリケーション500Aは、パーティション1__300Aに移動した。このため、アプリケーション管理テーブル136Bにおいて、「APP1」に対応する配置パーティション番号1622、配置サーバ名1623、配置FC-SW名1624及び配置ストレージ名1625は、それぞれ、「1」、「サーバP1」、「FC-SWP1」及び「ストレージP1」に更新される。

【0315】

さらに、パーティション管理テーブル136Aにおいて、パーティション1__300Aに対応するエントリ1611~1613の電源1606は全て「ON」、それらの配置アプリケーション1607は全て「APP1」に更新される。さらに、「サーバP1」に対応する残リソース量1605及びアプリケーションリソース量1608がそれぞれ「5」及び「5」に、「FC-SWP1」に対応する残リソース量1605及びアプリケーションリソース量1608がそれぞれ「9」及び「1」に、「ストレージP1」に対応する残リソース量1605及びアプリケーションリソース量1608がそれぞれ「4」及び「6」に更新される。

10

【0316】

一方、パーティション2__300Bに対応するエントリ1614~1616の残リソース量1605はリソース量1604と同じ「10」に、電源1606は全て「OFF」に、配置アプリケーション1607及びアプリケーションリソース量1608は全て空白に更新される。

20

【0317】

図29は、本発明の実施形態において実行されるアプリケーション500の移動及び電源制御の処理の変形例の説明図である。

【0318】

既に説明した図17の処理では、サーバ電源ON処理(ステップ1704)及びFC-SW電源ON処理(ステップ1705)が、いずれも論理ボリューム603の全体コピー処理の終了を契機として開始される。これに対して、図29では、サーバ電源ON処理及びFC-SW電源ON処理の開始時刻が算出される。

【0319】

図15を参照して説明したように、差分コピー処理、サーバ120の起動処理及びFC-SW140の起動処理が同時刻に終了するように(すなわち、矢印1503、1504及び1505の先端が同一の時刻に対応するように)、サーバ電源ON処理及びFC-SW電源ON処理を開始することが望ましい。しかし、差分コピー処理に要する時間を正確に予測することは困難であるため、図17の処理では、全体コピー処理の終了時刻を、サーバ電源ON処理及びFC-SW電源ON処理を開始する近似的に最適な時刻として用いている。

30

【0320】

これに対して、図29の処理では、差分コピー処理に要する時間を近似的に算出することによって、サーバ電源ON処理及びFC-SW電源ON処理を開始する時刻が設定される。

40

【0321】

以下、図29の処理を説明する。図29の処理は、図17の処理の代わりに、すなわち、図10のステップ1002において実行される。

【0322】

管理サーバ130の管理プログラム135は、アプリケーション500の移動処理を開始すると(ステップ2901)、最初に、ストレージ電源ON処理を実行する(ステップ2902)。これらの処理は、それぞれ図17のステップ1701及びステップ1702と同様であるため、説明を省略する。

【0323】

50

ストレージ電源ON処理が終了すると、管理プログラム135は、論理ボリューム603の全体コピー処理を実行する(ステップ2903)。この処理は、図17のステップ1703と同様である。

【0324】

論理ボリューム603の全体コピー処理が終了すると、管理プログラム135は、論理ボリューム603の差分コピー処理を実行する(ステップ2904)。この処理は、図17のステップ1706と同様である。

【0325】

さらに、管理プログラム135は、ステップ2903及び2904と並行して、ステップ2905から2909を実行する。

10

【0326】

具体的には、管理プログラム135は、論理ボリューム603の全体コピー処理が開始された後、論理ボリューム603の差分コピー処理の終了予定時間を算出する(ステップ2905)。この処理は、全体コピー処理によってコピーされるべきデータのうち所定の割合(例えば80%、90%又は100%等)のデータのコピーが終了したことを契機に実行される。ここでは、所定の割合が90%である場合を例としてステップ2905を説明する。なお、1回の全体コピー処理において複数回(例えば、80%が終了した時点及び90%が終了した時点)ステップ2905が実行されてもよい。

【0327】

管理プログラム135は、全体コピー処理によってコピーされるべきデータのうち90%のコピーが終了した時点で差分ストックとして格納されているデータ(以下、差分データ)の量から、全体コピー処理が全て(すなわち100%)終了した時点の差分データ量を予測する。例えば、差分I/Oが一定の間隔で発生すると仮定した場合、全体コピー処理の90%が終了した時点の差分データ量を0.9で除算する(すなわち0.9の逆数を乗ずる)ことによって、全体コピー処理が全て終了した時点の差分データ量を算出することができる。一般には、差分I/Oの間隔は一定ではないが、上記のように算出された値を近似的な予測値として使用してもよい。

20

【0328】

そして、管理プログラム135は、差分データ量の予測値と、コピーに使用されるデータ転送経路の転送速度から、差分データのコピーに要する時間を算出する。この算出の際に、データ転送経路の転送速度として、既に行われた全体コピー処理の際に実測された値が用いられてもよい。

30

【0329】

このようにして算出された時間を全体コピー処理の終了時刻に加算することによって、差分コピー処理の終了予定時刻が算出される。

【0330】

なお、既に説明したように、論理ボリューム603の全体コピー処理の終了時刻は、論理ボリューム603のデータ量及びデータ転送経路の転送速度に基づいて算出することができる。

【0331】

40

次に、管理プログラム135は、算出された差分コピー処理の終了予定時刻に基づいて、サーバ電源ON処理を開始する時刻及びFC-SW電源ON処理を開始する時刻を設定する(それぞれステップ2906及びステップ2907)。

【0332】

具体的には、管理プログラム135は、差分コピー処理の終了予定時刻からサーバ120の起動処理に要する時間を減算することによってサーバ電源ON処理を開始する時刻を算出する。そして、その時刻が到来したら、サーバ電源ON処理を開始する(ステップ2908)。この処理は、図17のステップ1704と同様である。

【0333】

さらに、差分コピー処理の終了予定時刻からFC-SW140の起動処理に要する時間

50

を減算することによってFC - SW電源ON処理を開始する時刻を算出する。そして、その時刻が到来したら、FC - SW電源ON処理を開始する(ステップ2909)。この処理は、図17のステップ1705と同様である。

【0334】

なお、管理プログラム135は、サーバ120の起動処理に要する時間及びFC - SW140の起動処理に要する時間を示す情報をあらかじめ保持している。これらの時間は、各装置のハードウェアスペックから算出されてもよいし、過去に実測された値であってもよい。

【0335】

ステップ2904、ステップ2908及びステップ2909のいずれも終了すると、管理プログラム135は、論理パス602の移動処理(ステップ2910)及び仮想サーバ601の移動処理(ステップ2911)を実行し、アプリケーション500の移動処理を終了する(ステップ2912)。これらは、それぞれ、図17のステップ1707、ステップ1708及びステップ1709と同様である。

10

【0336】

差分I/Oの偏在、及び、転送速度の変動の程度が小さいほど、上記のような差分コピー処理の終了時刻の予測は正確になる。その結果、計算機システムの消費電力を効果的に削減することができる。

【0337】

なお、上記のように、図29の処理は、図17の処理の代わりに実行されるものである。図17の処理及び図29の処理は、いずれも、計算機システムの消費電力を削減するために、サーバ120及びFC - SW140の電源を近似的に最適な時刻に投入するためのものである。いずれの処理においても、サーバ120及びFC - SW140の電源を投入する近似的に最適な時刻が、全体コピー処理及び差分コピー処理に要する時間に基づいて算出される。このように、図17の処理及び図29の処理の技術的な意義は共通する。

20

【0338】

図30は、本発明の実施形態の計算機システムの規模と効果との関係を示す説明図である。

【0339】

図30(A)は、比較的小規模な計算機システム、例えば図3又は図4に示したように、一つの物理パーティションに含まれるサーバ120、FC - SW140及びストレージ装置100がそれぞれ1台~数台程度である計算機システムにおけるアプリケーションの移動及び電源制御のタイミングを示す。図30(A)の内容は、図15に示したものと同様であるため、説明を省略する。

30

【0340】

図30(B)は、比較的大規模なシステム、例えば、一つの物理パーティションに含まれるサーバ120、FC - SW140及びストレージ装置100がそれぞれ数十台~数百台程度である計算機システムにおけるアプリケーションの移動及び電源制御のタイミングを示す。このような物理パーティションは、例えばデータセンターとして使用される。図30(B)の内容も、基本的には図15に示したものと同様である。しかし、一般に、計算機システムの規模が大きくなるほど、各装置の起動及びアプリケーションの移動に要する時間は長くなる。

40

【0341】

例えば、論理ボリューム603の移動に要する時間は、前述のように転送されるデータ量等に依存するが、一般には小規模なシステムでは数時間程度である場合が多い。これに対して、データセンターのような大規模なシステムでは、論理ボリューム603の移動に数週間以上を要する場合もある。

【0342】

本実施形態では、図15を参照して説明したように、サーバ120及びFC - SW140の電源投入時刻を論理ボリューム603の移動の開始時刻より遅らせることによって、

50

論理ボリューム603の移動の終了時刻と、サーバ120及びFC-SW140の起動処理の終了時刻とが概ね一致するように制御される。すなわち、論理ボリューム603の移動が実行されている時間の大部分(例えば、全体コピー処理が実行されている時間)は、サーバ120及びFC-SW140の電源が投入されていない。このため、論理ボリューム603の移動に長時間を要する大規模な計算機システムに適用するほど、本実施形態の効果は大きくなる。

【0343】

図31は、本発明の実施形態における物理パーティションの第4の例の説明図である。

【0344】

図31に示す物理パーティションは、例えばデータセンタとして使用される大規模なものである。図31の例では、パーティション1__300A及びパーティション2__300Bそれぞれが20台のサーバ120、9台のFC-SW140及び3台のストレージ装置100を備える。なお、これらの台数は一例であり、大規模なデータセンタではさらに台数が多い場合もある。

10

【0345】

このような計算機システムにおいて、データセンタを移動(マイグレーション)するために、パーティション内に配置されたアプリケーションを移動する必要が生じる場合がある。例えば、移動すべきアプリケーションに含まれる論理ボリューム群(すなわち論理ボリューム603の集合)3101Aを移動先のパーティション内の論理ボリューム群3101Fに移動するために、図30において説明したように数週間以上を要する場合もある。このような計算機システムに本実施形態を適用することによって、サーバ120及びFC-SW140の消費電力を削減することができる。

20

【0346】

仮に、本実施形態を適用することによって、サーバ120及びFC-SW140の電源投入時刻を従来より1日遅らせることができた場合、1日分のサーバ120及びFC-SW140の消費電力が削減される。サーバ120及びFC-SW140の台数が多いほど、消費電力削減の効果が大きくなる。

【0347】

図32は、本発明の実施形態における物理パーティションの第5の例の説明図である。

【0348】

図32は、パーティションによって階層化されたデータセンタの例を示す。この例において、パーティション1__300Aは高い階層、パーティション2__300Bは低い階層に相当する。すなわち、パーティション1__300Aは、パーティション2__300Bより高い処理性能を有する。例えば、パーティション1__300Aには、パーティション2__300Bに割り当てられているものより高性能のハードウェアが割り当てられてもよい。

30

【0349】

例えば、あるプロジェクトで使用されるアプリケーション群を、ある期間のみ、高い階層の(すなわち高性能の)パーティションで実行する必要がある場合、その期間の開始時点にアプリケーション群の移動が終了するように、あらかじめそのアプリケーション群の論理ボリューム603のコピーを開始することができる。

40

【0350】

図32の例では、低い階層(Tier2)に相当するパーティション2__300Bに、プロジェクトAで使用されるアプリケーション群3201A及びプロジェクト(又は業務)Bで使用されるアプリケーション群3201Bが配置されている。このような場合、本実施形態によって、アプリケーション群ごとに、必要に応じて高い階層(Tier1)のパーティション1__300Aへの移動を制御することができる。これによって、計算機システムの消費電力を削減しつつ、所望の時刻におけるアプリケーション群の移動を実現することができる。

【図面の簡単な説明】

50

- 【 0 3 5 1 】
- 【 図 1 】 本発明の実施形態におけるパーティションの説明図である。
- 【 図 2 】 本発明の実施形態の計算機システムの構成を示すブロック図である。
- 【 図 3 】 本発明の実施形態における物理パーティションの第 1 の例の説明図である。
- 【 図 4 】 本発明の実施形態における物理パーティションの第 2 の例の説明図である。
- 【 図 5 】 本発明の実施形態における論理パーティションの例の説明図である。
- 【 図 6 】 本発明の実施形態におけるアプリケーションの構成の説明図である。
- 【 図 7 】 本発明の実施形態の全体システム構成の説明図である。
- 【 図 8 】 本発明の実施形態において実行されるアプリケーションのパーティション間の移動処理を示すフローチャートである。 10
- 【 図 9 】 本発明の実施形態において実行されるアプリケーションのパーティション間の移動処理の詳細な手順を示す説明図である。
- 【 図 10 】 本発明の実施形態において実行されるアプリケーションのパーティション間の移動処理の詳細な手順を示すフローチャートである。
- 【 図 11 】 本発明の実施形態において実行される論理ボリュームの移動の説明図である。
- 【 図 12 】 本発明の実施形態において実行される論理ボリュームの移動の別の例の説明図である。
- 【 図 13 】 本発明の実施形態における物理パーティションの第 3 の例の説明図である。
- 【 図 14 】 本発明の実施形態において実行されるアプリケーションの移動及び電源制御の説明図である。 20
- 【 図 15 】 本発明の実施形態において実行されるアプリケーションの移動及び電源制御のタイミングの説明図である。
- 【 図 16 】 本発明の実施形態の管理テーブルの説明図である。
- 【 図 17 】 本発明の実施形態において実行されるアプリケーションの移動及び電源制御の処理の全体を示すフローチャートである。
- 【 図 18 】 本発明の実施形態において実行されるストレージ電源 ON 処理を示すフローチャートである。
- 【 図 19 】 本発明の実施形態の管理サーバが実行する論理ボリュームの全体コピー処理のフローチャートである。
- 【 図 20 】 本発明の実施形態のストレージ装置が実行する論理ボリュームの全体コピー処理のフローチャートである。 30
- 【 図 21 】 本発明の実施形態において実行されるサーバ電源 ON 処理を示すフローチャートである。
- 【 図 22 】 本発明の実施形態において実行される FC - SW 電源 ON 処理を示すフローチャートである。
- 【 図 23 】 本発明の実施形態の管理サーバが実行する論理ボリュームの差分コピー処理のフローチャートである。
- 【 図 24 】 本発明の実施形態のストレージ装置が実行する論理ボリュームの差分コピー処理のフローチャートである。
- 【 図 25 】 本発明の実施形態において実行される論理バスの移動処理のフローチャートである。 40
- 【 図 26 】 本発明の実施形態において実行される仮想サーバの移動処理のフローチャートである。
- 【 図 27 】 本発明の実施形態において実行される移動先パーティションの状態確認処理のフローチャートである。
- 【 図 28 】 本発明の実施形態においてアプリケーションの移動が実行された後の管理テーブルの説明図である。
- 【 図 29 】 本発明の実施形態において実行されるアプリケーションの移動及び電源制御の処理の変形例の説明図である。
- 【 図 30 】 本発明の実施形態の計算機システムの規模と効果との関係を示す説明図である 50

【図31】本発明の実施形態における物理パーティションの第4の例の説明図である。

【図32】本発明の実施形態における物理パーティションの第5の例の説明図である。

【符号の説明】

【0352】

100 ストレージ装置

101 コントローラ

110A ~ 110C ディスクドライブ

111A、111B、603、603A ~ 603F 論理ボリューム

120 サーバ

130 管理サーバ

135 管理プログラム

136 管理テーブル

136A パーティション管理テーブル

136B アプリケーション管理テーブル

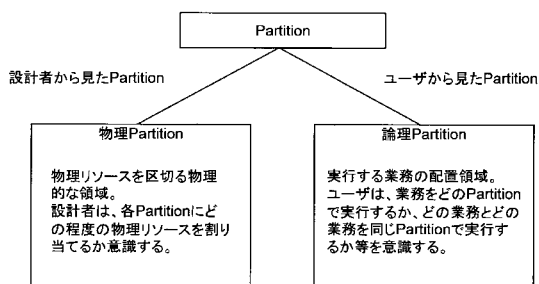
300A ~ 300I パーティション

500、500A ~ 500E アプリケーション

601、601A ~ 601E 仮想サーバ

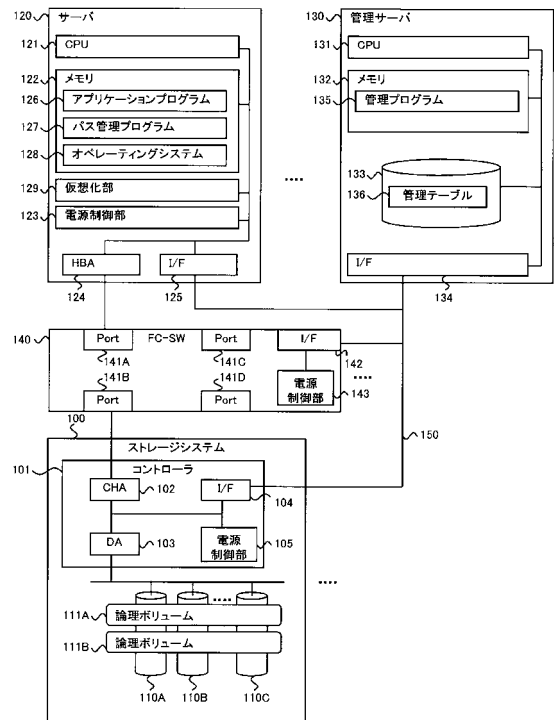
602、602A ~ 602E FC-SW

【図1】

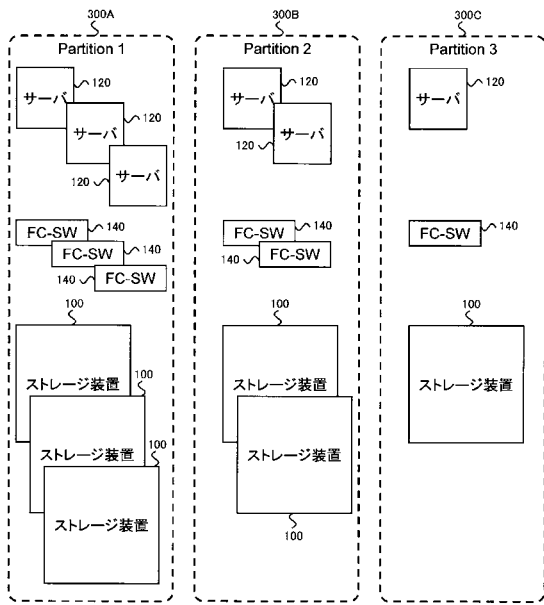


ITシステムのPartitionの定義

【図2】

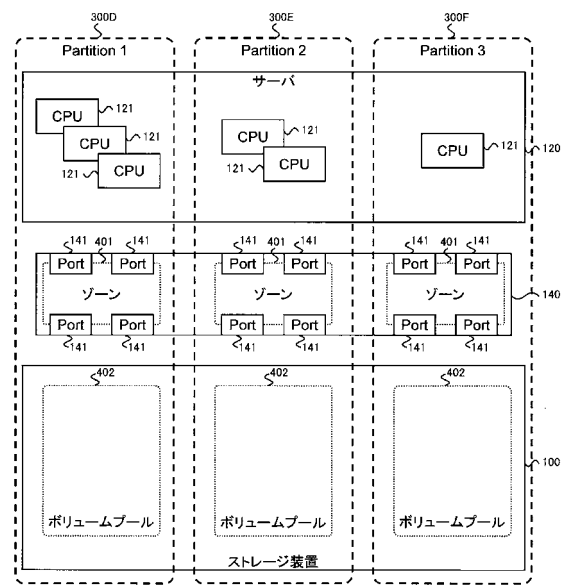


【 図 3 】



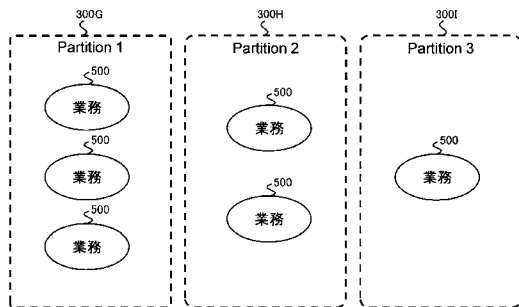
物理パーティションの第1の例

【 図 4 】



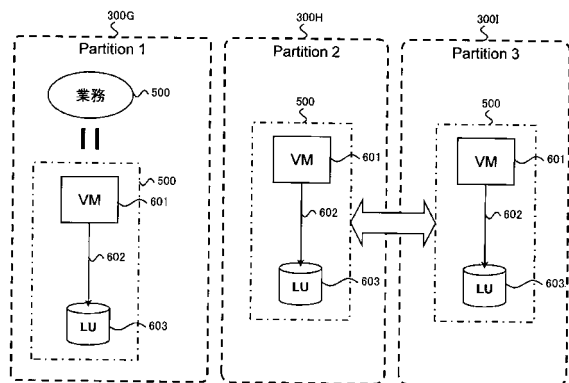
物理パーティションの第2の例

【 図 5 】

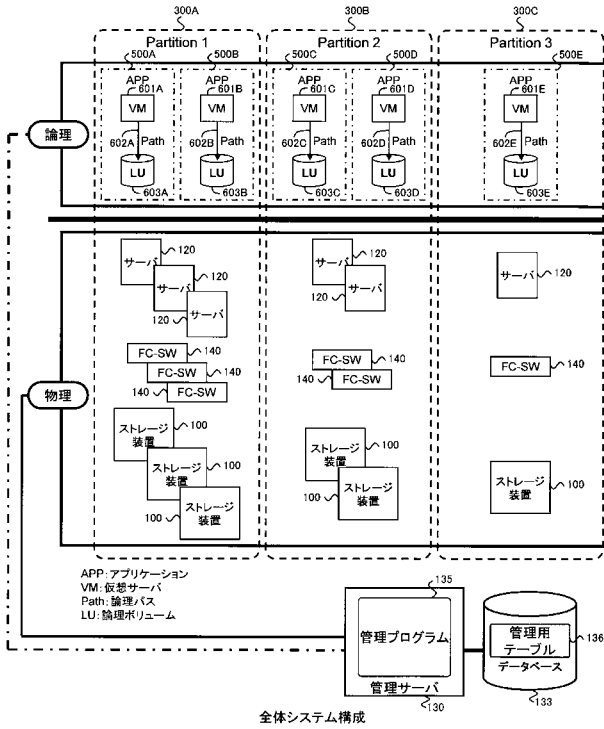


論理パーティション

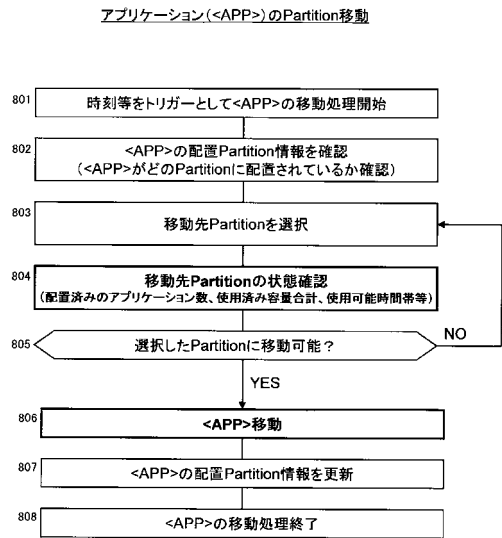
【 図 6 】



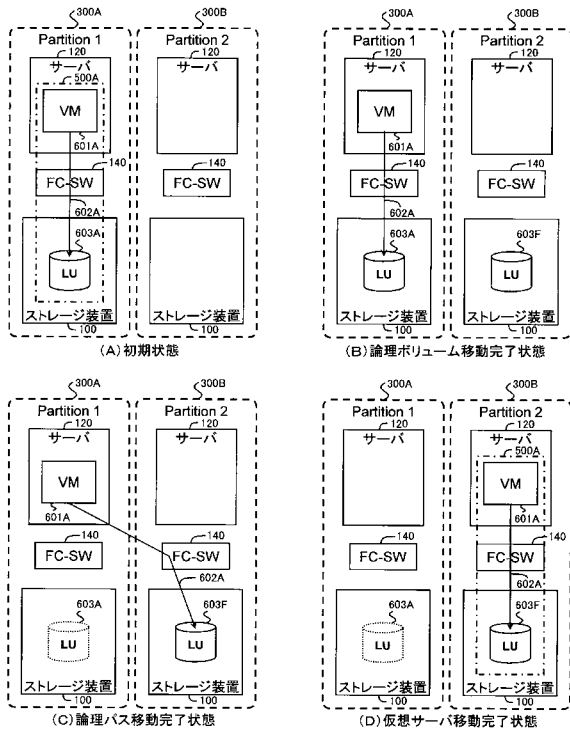
【 図 7 】



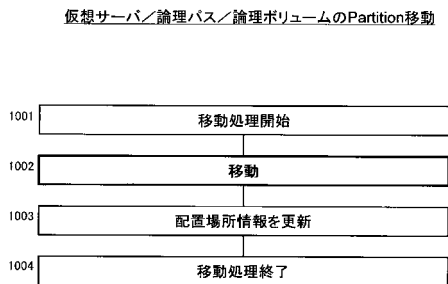
【 図 8 】



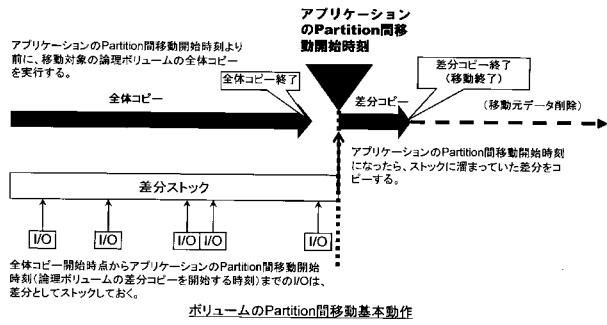
【 図 9 】



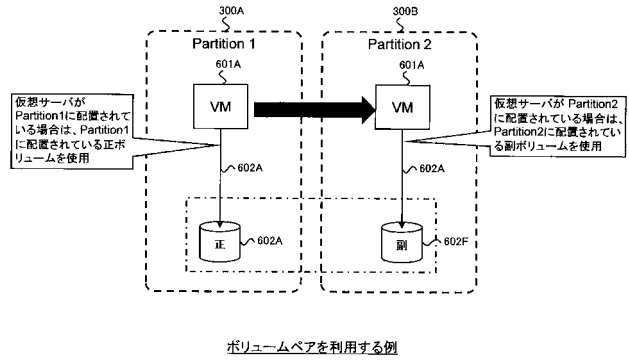
【 図 10 】



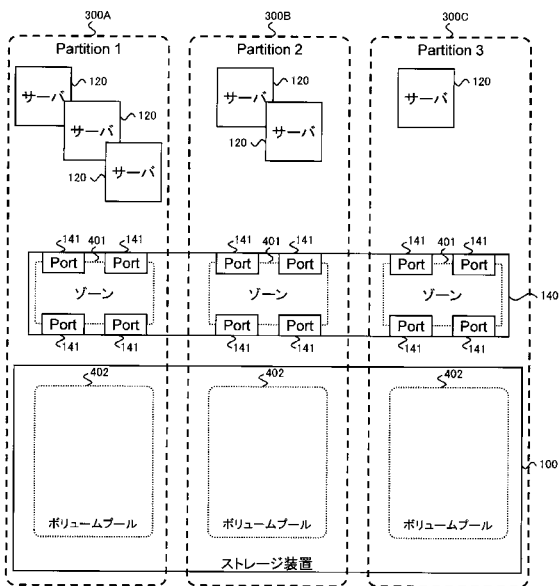
【図 1 1】



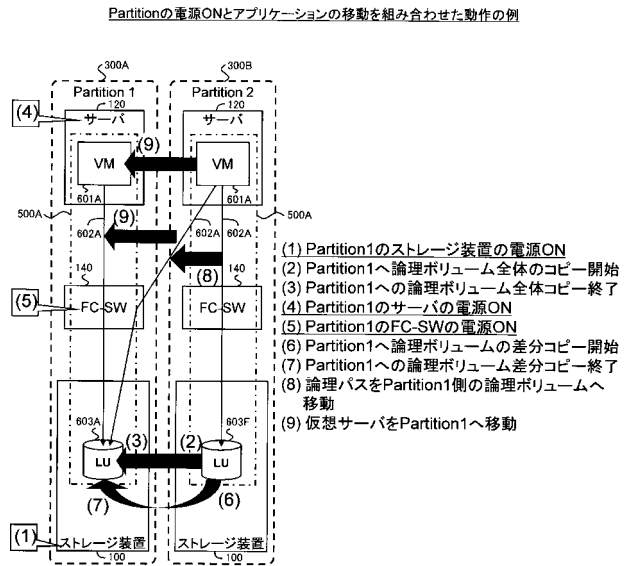
【図 1 2】



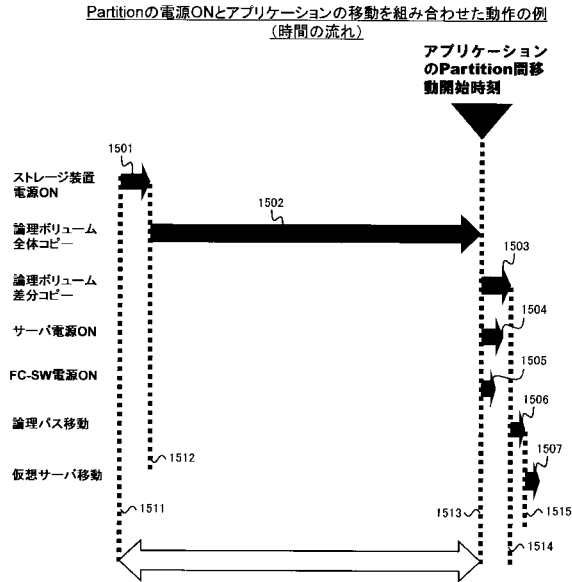
【図 1 3】



【図 1 4】



【 図 1 5 】



【 図 1 6 】

パーティション管理テーブル 136A

パーティション番号	ハードウェア種別	ハードウェア名	リソース量	残りリソース量	電源	配置APP	APPリソース量
1	サーバ1	サーバIP1	10	10	OFF		
	FC-SW1	FC-SWP1	10	10	OFF		
	ストレージ1	ストレージP1	10	10	OFF		
2	サーバ2	サーバIP2	10	5	ON	APP1	5
	FC-SW2	FC-SWP2	10	9	ON	APP1	1
	ストレージ2	ストレージP2	10	4	ON	APP1	6

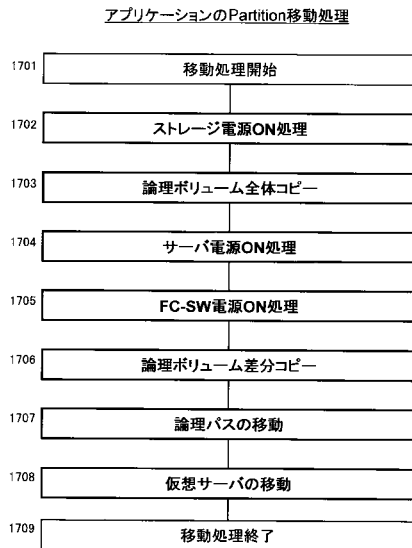
(A)

アプリケーション管理テーブル 136B

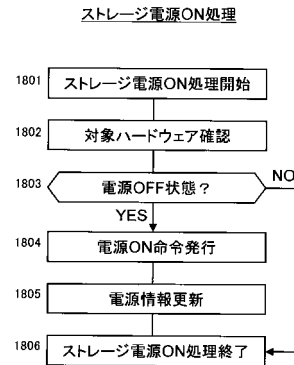
アプリケーション名	配置Partition#	配置サーバ名	配置FC-SW名	配置ストレージ名	サーバリソース量	FC-SWリソース量	ストレージリソース量
APP1	2	サーバIP2	FC-SWP2	ストレージP2	5	1	6

(B)

【 図 1 7 】

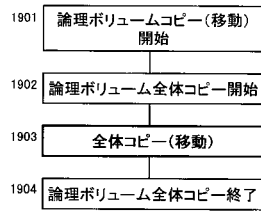


【 図 1 8 】



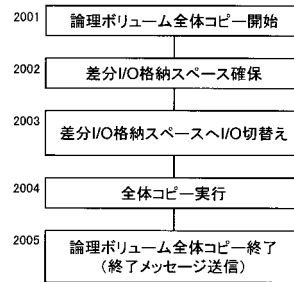
【 図 1 9 】

論理ボリュームのPartition移動(全体コピー)



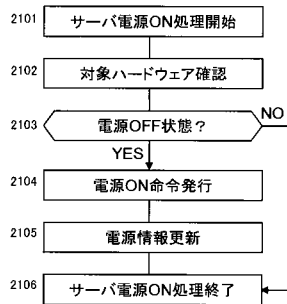
【 図 2 0 】

論理ボリュームの全体コピー(ストレージ側処理)



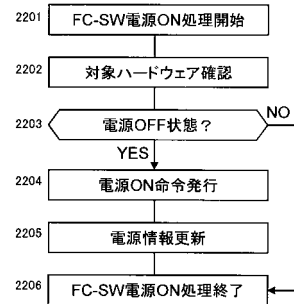
【 図 2 1 】

サーバ電源ON処理

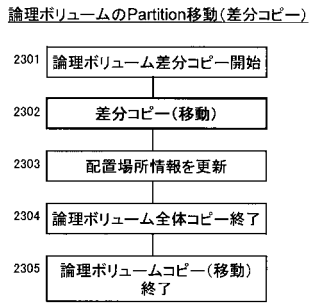


【 図 2 2 】

FC-SW電源ON処理

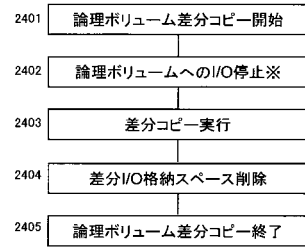


【 図 2 3 】



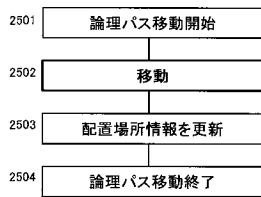
【 図 2 4 】

論理ボリュームの差分コピー(ストレージ側処理)



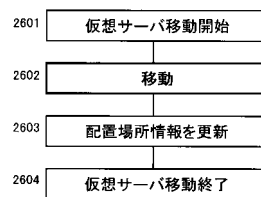
【 図 2 5 】

論理バスの移動

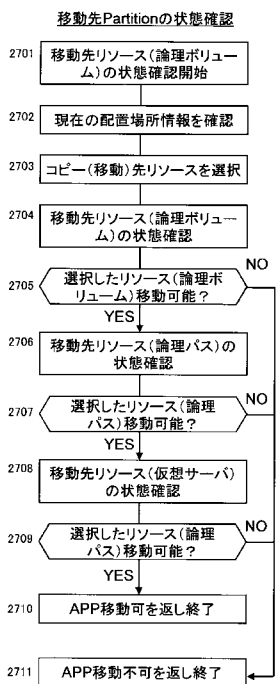


【 図 2 6 】

仮想サーバの移動



【 図 2 7 】



【 図 2 8 】

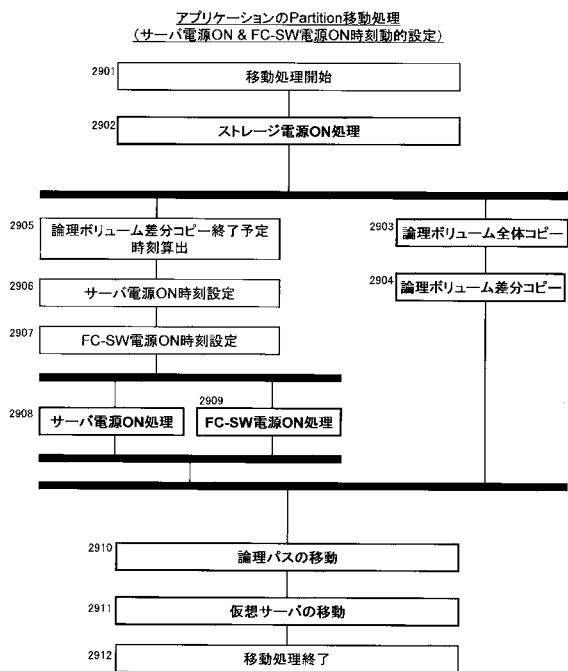
パーティション管理テーブル 136A

パーティション番号	ハードウェア種別	ハードウェア名	リソース量	残りリソース量	電源	配置APP	APPリソース量
1	サーバ1	サーバP1	10	5	ON	APP1	5
	FC-SW1	FC-SWP1	10	9	ON	APP1	1
	ストレージ1	ストレージP1	10	4	ON	APP1	6
2	サーバ2	サーバP2	10	10	OFF		
	FC-SW2	FC-SWP2	10	10	OFF		
	ストレージ2	ストレージP2	10	10	OFF		

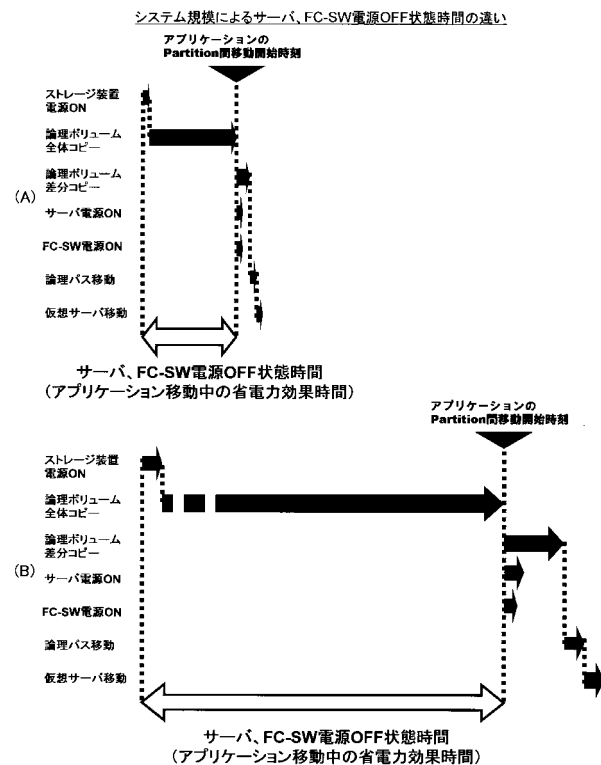
アプリケーション管理テーブル 136B

アプリケーション名	配置Partition#	配置サーバ名	配置FC-SW名	配置ストレージ名	サーバリソース量	FC-SWリソース量	ストレージリソース量
APP1	1	サーバP1	FC-SWP1	ストレージP1	5	1	6

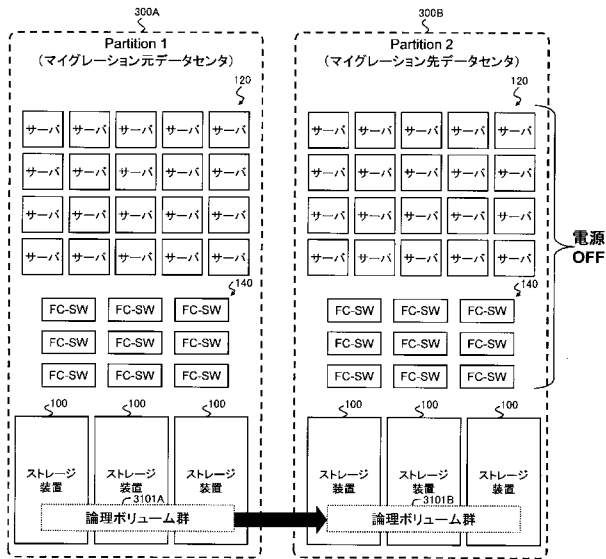
【 図 2 9 】



【 図 3 0 】

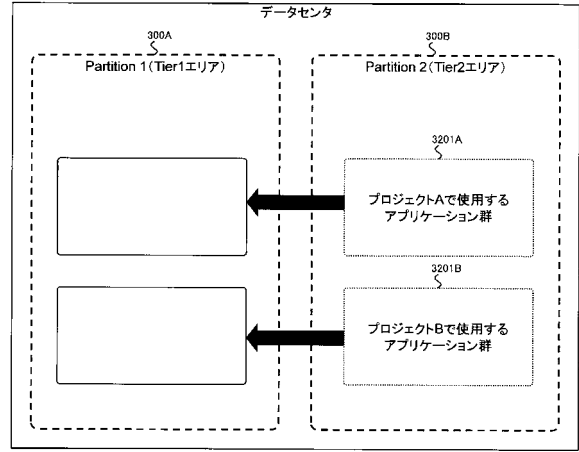


【図 3 1】



物理パーティションの第4の例
(データセンタを1つのPartitionと定義する例)

【図 3 2】



物理パーティションの第5の例
(データセンタをTierのエリアで区切る例)