

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5452765号  
(P5452765)

(45) 発行日 平成26年3月26日 (2014. 3. 26)

(24) 登録日 平成26年1月10日 (2014. 1. 10)

(51) Int. Cl.		F I			
<b>G06F 12/00</b>	<b>(2006.01)</b>	G06F 12/00	531R		
<b>G06F 3/06</b>	<b>(2006.01)</b>	G06F 3/06	306H		
		G06F 3/06	306K		
		G06F 3/06	304E		
		G06F 3/06	301X		

請求項の数 15 (全 53 頁)

(21) 出願番号 特願2013-501477 (P2013-501477)  
 (86) (22) 出願日 平成22年12月14日 (2010. 12. 14)  
 (65) 公表番号 特表2013-532314 (P2013-532314A)  
 (43) 公表日 平成25年8月15日 (2013. 8. 15)  
 (86) 国際出願番号 PCT/JP2010/007254  
 (87) 国際公開番号 W02012/081050  
 (87) 国際公開日 平成24年6月21日 (2012. 6. 21)  
 審査請求日 平成25年1月11日 (2013. 1. 11)

(73) 特許権者 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 110000176  
 一色国際特許業務法人  
 (72) 発明者 雑賀 信之  
 神奈川県小田原市中里322番2号 株式  
 会社日立製作所 R A I D システム事業部内

審査官 池田 聡史

最終頁に続く

(54) 【発明の名称】 情報処理システムにおける障害復旧方法、及び情報処理システム

(57) 【特許請求の範囲】

【請求項1】

第1ファイルシステムを有し、データI/O要求を受け付ける第1サーバ装置と、  
 第2ファイルシステムを有し、前記第1サーバ装置に通信可能に接続された第2サーバ  
 装置と、を備え、

前記第1サーバ装置は、前記データI/O要求の対象であるファイルのデータを第1ス  
 トレージ装置に記憶し、

前記第2サーバ装置は、前記データI/O要求の対象であるファイルのデータを第2ス  
 トレージ装置に記憶し、

前記第1サーバ装置が、前記第1ストレージ装置に記憶しているファイルのデータを前  
 記第2サーバ装置に送信し、

前記第2サーバ装置が、前記第1サーバ装置から送られてくる前記データを前記第2ス  
 トレージ装置に記憶している情報処理システムにおける障害の復旧方法であって、

前記第2サーバ装置は、障害からの復旧に際して前記第1サーバ装置が前記データI/  
 O要求の受け付けを開始するのに先立ち、前記第2ストレージ装置に記憶しているディレ  
 クトリイメージのうち最上位のディレクトリから所定の低位階層までのディレクトリイメ  
 ージを、前記第1サーバ装置に送信し、

前記第1サーバ装置は、前記第2サーバ装置から送られてくる前記ディレクトリイメ  
 ージを前記第1ストレージ装置に復元した後、前記データI/O要求の受け付けを再開し、

前記第1サーバ装置は、前記データI/O要求の受け付け再開後において、受け付けた

10

20

前記データ I / O 要求を処理するために必要となるディレクトリイメージが前記第 1 ストレージ装置に復元されていない場合には、前記ディレクトリイメージを前記第 2 サーバ装置に要求し、

前記第 2 サーバ装置は、前記第 1 サーバ装置から送られてくる前記要求に応じて前記ディレクトリイメージを前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信し、

前記第 1 サーバ装置は、前記第 2 ストレージ装置から送られてくる前記ディレクトリイメージに基づき前記データ I / O 要求を処理するとともに、前記ディレクトリイメージを前記第 1 ストレージ装置に復元する

情報処理システムにおける障害復旧方法。

10

#### 【請求項 2】

請求項 1 に記載の情報処理システムにおける障害復旧方法であって、

前記第 2 サーバ装置は、前記第 1 サーバ装置から前記要求が送られてきた場合、前記要求の対象である前記ディレクトリイメージを前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信するとともに、所定の選出方法に従って選出した、前記ディレクトリイメージとは異なる追加のディレクトリイメージを、前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信し、

前記第 1 サーバ装置は、前記第 2 サーバ装置から送られてくる前記ディレクトリイメージに基づき前記データ I / O 要求を処理するとともに、前記ディレクトリイメージ及び前記第 2 サーバ装置から送られてくる前記追加のディレクトリイメージを、前記第 1 ストレージ装置に復元する

20

情報処理システムにおける障害復旧方法。

#### 【請求項 3】

請求項 2 に記載の情報処理システムにおける障害復旧方法であって、

前記第 2 サーバ装置は、前記第 1 サーバ装置が受け付けた前記データ I / O 要求が所定の条件を満たす場合に、所定の選出方法に従って前記追加のディレクトリイメージを選出する

情報処理システムにおける障害復旧方法。

#### 【請求項 4】

請求項 3 に記載の情報処理システムにおける障害復旧方法であって、

前記所定の条件は、

前記データ I / O 要求の対象になっているファイルのデータサイズが、現在から遡って所定時間内に発生したデータ I / O 要求が対象としていたファイルのデータサイズの平均値よりも小さいこと、

前記データ I / O 要求の対象になっているファイルのデータサイズが予め設定された閾値よりも小さいこと

のうちの少なくともいずれかである

情報処理システムにおける障害復旧方法。

30

#### 【請求項 5】

請求項 3 に記載の情報処理システムにおける障害復旧方法であって、

前記所定の選出方法は、

前記第 1 ストレージ装置に既に復元されているディレクトリの配下に存在するファイルのメタデータ及び / 又は実体を前記追加のディレクトリイメージとする方法、

前記第 1 ストレージ装置に既に復元されているディレクトリの配下に存在するディレクトリのメタデータを前記追加のディレクトリイメージとする方法、

前記障害が発生する前に前記第 1 ストレージ装置に実体が格納されていたファイルの実体を前記追加のディレクトリイメージとする方法、

重要度が高く設定されているファイルのメタデータ及び / 又は実体を前記追加のディレクトリイメージとする方法、

前記障害の発生時点から遡って所定時間内におけるアクセス頻度が高いファイルを前記

40

50

追加のディレクトリイメージとする方法、  
のうちの少なくともいずれかである  
情報処理システムにおける障害復旧方法。

【請求項 6】

請求項 2 に記載の情報処理システムにおける障害復旧方法であって、

前記第 1 サーバ装置は、前記第 1 ファイルシステムに割り当てられている前記第 1 ストレージ装置の記憶領域の残容量が予め設定された閾値未満になっている場合に、前記第 1 ストレージ装置に記憶しているファイルのデータのうち、所定の選出基準に従って選出したファイルのデータのメタデータについては前記第 1 ストレージ装置に残すとともに当該データの实体については前記第 1 ストレージ装置から削除して前記残容量を確保するスタ

10

ブ化を行い、  
前記第 2 サーバ装置は、前記第 1 サーバ装置に送信した、前記ディレクトリイメージ又は前記追加のディレクトリイメージが再びスタブ化される現象である再スタブ化の発生頻度が予め設定された閾値以上になっているか否か、もしくは、再スタブ化の発生時間間隔が予め設定された閾値未満になっているか否かを監視し、

前記第 2 サーバ装置は、再スタブ化の発生頻度が予め設定された閾値以上になっているか、もしくは、再スタブ化の発生時間間隔が予め設定された閾値未満になっている場合に、前記第 1 サーバ装置への前記ディレクトリイメージ又は前記追加のディレクトリイメージの送信を抑制する

情報処理システムにおける障害復旧方法。

20

【請求項 7】

請求項 6 に記載の情報処理システムにおける障害復旧方法であって、

前記抑制は、

データ I/O 要求がファイルのメタデータのみを対象としている場合に当該ファイルの实体を復元しないようにする方法、

前記所定の選出方法の一つ以上を用いた前記追加のディレクトリイメージの選出を行っている場合に、更に別の選出方法を重複して適用するようにする方法、

前記データ I/O 要求の対象になっているファイルよりも高い重要度が設定されているファイルのメタデータ及び/又は实体を、前記追加のディレクトリイメージとして選出するようにしている場合に、選出基準としている前記重要度を更に高く設定するようにする

30

方法、  
前記障害の発生時点から遡って所定時間内におけるアクセス頻度が前記データ I/O 要求の対象になっているファイルよりも高いファイルを、前記追加のディレクトリイメージとして選出するようにしている場合に、選出基準としている前記アクセス頻度を更に高く設定するようにする方法、

のうちの少なくともいずれかの方法で行われる

情報処理システムにおける障害復旧方法。

【請求項 8】

請求項 6 に記載の情報処理システムにおける障害復旧方法であって、

前記第 1 サーバ装置は、前記第 1 ストレージ装置に記憶しているファイルが現在スタブ化されているか否かを示す情報を前記第 2 サーバ装置に随時送信し、

40

前記第 2 サーバ装置は、前記ディレクトリイメージ又は前記追加のディレクトリイメージの前記第 1 サーバ装置への送信履歴を管理し、

前記第 2 サーバ装置は、前記情報及び前記送信履歴に基づき、前記再スタブ化の発生頻度、もしくは、前記再スタブ化の発生時間間隔を取得する

情報処理システムにおける障害復旧方法。

【請求項 9】

請求項 1 に記載の情報処理システムにおける障害復旧方法であって、

前記第 1 サーバ装置は、前記第 1 ストレージ装置に記憶されているファイルのデータと、前記第 2 ストレージ装置に記憶されているファイルのデータと、を一致させる旨を要求

50

する同期要求が発生すると、前記第 1 ストレージ装置に記憶されているファイルのデータ又は前同期時からの更新差分を前記第 2 サーバ装置に送信し、

前記第 2 サーバ装置は、前記ファイルのデータ又は前記更新差分に基づき、前記第 1 ストレージ装置に記憶されているファイルのデータの複製として前記第 2 ストレージ装置に記憶しているファイルのデータを更新する

情報処理システムにおける障害復旧方法。

【請求項 10】

請求項 8 に記載の情報処理システムにおける障害復旧方法であって、

前記第 1 サーバ装置は、前記第 1 ストレージ装置に記憶されているファイルのデータと、前記第 2 ストレージ装置に記憶されているファイルのデータとを一致させる旨を要求する同期要求が発生すると、前記第 1 ストレージ装置に記憶されているファイルのデータ又は前同期時からの更新差分を前記第 2 サーバ装置に送信し、

前記第 2 サーバ装置は、前記ファイルのデータ又は前記更新差分に基づき、前記第 1 ストレージ装置に記憶されているファイルのデータの複製として前記第 2 ストレージ装置に記憶しているファイルのデータを更新し、

ファイルが現在スタブ化されているか否かを示す前記情報は、前記第 1 ストレージ装置又は前記第 2 ストレージ装置に記憶されているファイルのデータのうちメタデータに含まれている

情報処理システムにおける障害復旧方法。

【請求項 11】

請求項 1 に記載の情報処理システムにおける障害復旧方法であって、

前記ディレクトリイメージには、ディレクトリの階層構造、ディレクトリのメタデータ、ファイルのメタデータ、ファイルの実体のうちの少なくともいずれかが含まれている

情報処理システムにおける障害復旧方法。

【請求項 12】

請求項 1 に記載の情報処理システムにおける障害復旧方法であって、

前記第 2 サーバ装置は、前記第 1 サーバ装置から前記要求が送られてきた場合、前記要求の対象である前記ディレクトリイメージを前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信するとともに、所定の選出方法に従って選出した、前記ディレクトリイメージとは異なる追加のディレクトリイメージを、前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信し、

前記第 1 サーバ装置は、前記第 2 サーバ装置から送られてくる前記ディレクトリイメージに基づき前記データ I/O 要求を処理するとともに、前記ディレクトリイメージ及び前記第 2 サーバ装置から送られてくる前記追加のディレクトリイメージを、前記第 1 ストレージ装置に復元し、

前記第 2 サーバ装置は、前記第 1 サーバ装置が受け付けた前記データ I/O 要求が所定の条件を満たす場合に、所定の選出方法に従って前記追加のディレクトリイメージを選出し、

前記所定の条件は、

前記データ I/O 要求の対象になっているファイルのデータサイズが、現在から遡って所定時間内に発生したデータ I/O 要求が対象としていたファイルのデータサイズの平均値よりも小さいこと、

前記データ I/O 要求の対象になっているファイルのデータサイズが予め設定された閾値よりも小さいこと

のうちの少なくともいずれかであり、

前記所定の選出方法は、

前記第 1 ストレージ装置に既に復元されているディレクトリの配下に存在するファイルのメタデータ及び/又は実体を前記追加のディレクトリイメージとする方法、

前記第 1 ストレージ装置に既に復元されているディレクトリの配下に存在するディレク

10

20

30

40

50

トリのメタデータを前記追加のディレクトリイメージとする方法、

前記障害が発生する前に前記第1ストレージ装置に実体が格納されていたファイルの実体を前記追加のディレクトリイメージとする方法、

重要度が高く設定されているファイルのメタデータ及び/又は実体を前記追加のディレクトリイメージとする方法、

前記障害の発生時点から遡って所定時間内におけるアクセス頻度が高いファイルを前記追加のディレクトリイメージとする方法、

のうちの少なくともいずれかであり、

前記第1サーバ装置は、前記第1ファイルシステムに割り当てられている前記第1ストレージ装置の記憶領域の残容量が予め設定された閾値未満になっている場合に、前記第1ストレージ装置に記憶しているファイルのデータのうち、所定の選出基準に従って選出したファイルのデータのメタデータについては前記第1ストレージ装置に残すとともに当該データの実体については前記第1ストレージ装置から削除して前記残容量を確保するスタブ化を行い、

10

前記第2サーバ装置は、前記第1サーバ装置に送信した、前記ディレクトリイメージ又は前記追加のディレクトリイメージが再びスタブ化される現象である再スタブ化の発生頻度が予め設定された閾値以上になっているか否か、もしくは、再スタブ化の発生時間間隔が予め設定された閾値未満になっているか否かを監視し、

前記第2サーバ装置は、再スタブ化の発生頻度が予め設定された閾値以上になっているか、もしくは、再スタブ化の発生時間間隔が予め設定された閾値未満になっている場合に、前記第1サーバ装置への前記ディレクトリイメージ又は前記追加のディレクトリイメージの送信を抑制し、

20

前記抑制は、

データI/O要求がファイルのメタデータのみを対象としている場合に当該ファイルの実体を復元しないようにする方法、

前記所定の選出方法の一つ以上を用いた前記追加のディレクトリイメージの選出を行っている場合に、更に別の選出方法を重複して適用するようにする方法、

前記データI/O要求の対象になっているファイルよりも高い重要度が設定されているファイルのメタデータ及び/又は実体を、前記追加のディレクトリイメージとして選出するようにしている場合に、選出基準としている前記重要度を更に高く設定するようにする方法、

30

前記障害の発生時点から遡って所定時間内におけるアクセス頻度が前記データI/O要求の対象になっているファイルよりも高いファイルを、前記追加のディレクトリイメージとして選出するようにしている場合に、選出基準としている前記アクセス頻度を更に高く設定するようにする方法、

のうちの少なくともいずれかの方法で行われ、

前記第1サーバ装置は、前記第1ストレージ装置に記憶しているファイルが現在スタブ化されているか否かを示す情報を前記第2サーバ装置に随時送信し、

前記第2サーバ装置は、前記ディレクトリイメージ又は前記追加のディレクトリイメージの前記第1サーバ装置への送信履歴を管理し、

40

前記第2サーバ装置は、前記情報及び前記送信履歴に基づき、前記再スタブ化の発生頻度、もしくは、前記再スタブ化の発生時間間隔を取得し、

前記第1サーバ装置は、前記第1ストレージ装置に記憶されているファイルのデータと、前記第2ストレージ装置に記憶されているファイルのデータと、を一致させる旨を要求する同期要求が発生すると、前記第1ストレージ装置に記憶されているファイルのデータ又は前同期時からの更新差分を前記第2サーバ装置に送信し、

前記第2サーバ装置は、前記ファイルのデータ又は前記更新差分に基づき、前記第1ストレージ装置に記憶されているファイルのデータの複製として前記第2ストレージ装置に記憶しているファイルのデータを更新し、

ファイルが現在スタブ化されているか否かを示す前記情報は、前記第1ストレージ装置

50

又は前記第 2 ストレージ装置に記憶されているファイルのデータのうちメタデータに含まれており、

前記ディレクトリイメージには、ディレクトリの階層構造、ディレクトリのメタデータ、ファイルのメタデータ、ファイルの実体のうちの少なくともいずれかが含まれている  
情報処理システムにおける障害復旧方法。

【請求項 1 3】

第 1 ファイルシステムを有し、データ I / O 要求を受け付ける第 1 サーバ装置と、  
第 2 ファイルシステムを有し、前記第 1 サーバ装置に通信可能に接続された第 2 サーバ装置と、を備え、

前記第 1 サーバ装置は、前記データ I / O 要求の対象であるファイルのデータを第 1 ストレージ装置に記憶し、

前記第 2 サーバ装置は、前記データ I / O 要求の対象であるファイルのデータを第 2 ストレージ装置に記憶し、

前記第 1 サーバ装置が、前記第 1 ストレージ装置に記憶しているファイルのデータを前記第 2 サーバ装置に送信し、

前記第 2 サーバ装置が、前記第 1 サーバ装置から送られてくる前記データを前記第 2 ストレージ装置に記憶している情報処理システムであって、

前記第 2 サーバ装置が、障害からの復旧に際して前記第 1 サーバ装置が前記データ I / O 要求の受け付けを開始するのに先立ち、前記第 2 ストレージ装置に記憶しているディレクトリイメージのうち最上位のディレクトリから所定の低位階層までのディレクトリイメージを、前記第 1 サーバ装置に送信し、

前記第 1 サーバ装置が、前記第 2 サーバ装置から送られてくる前記ディレクトリイメージを前記第 1 ストレージ装置に復元した後、前記データ I / O 要求の受け付けを再開し、

前記第 1 サーバ装置が、前記データ I / O 要求の受け付け再開後において、受け付けた前記データ I / O 要求を処理するために必要となるディレクトリイメージが前記第 1 ストレージ装置に復元されていない場合には、前記ディレクトリイメージを前記第 2 サーバ装置に要求し、

前記第 2 サーバ装置が、前記第 1 サーバ装置から送られてくる前記要求に応じて前記ディレクトリイメージを前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信し、

前記第 1 サーバ装置が、前記第 2 ストレージ装置から送られてくる前記ディレクトリイメージに基づき前記データ I / O 要求を処理するとともに、前記ディレクトリイメージを前記第 1 ストレージ装置に復元する

情報処理システム。

【請求項 1 4】

請求項 1 3 に記載の情報処理システムであって、

前記第 2 サーバ装置は、前記第 1 サーバ装置から前記要求が送られてきた場合、前記要求の対象である前記ディレクトリイメージを前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信するとともに、所定の選出方法に従って選出した、前記ディレクトリイメージとは異なる追加のディレクトリイメージを、前記第 2 ストレージ装置から読み出して前記第 1 サーバ装置に送信し、

前記第 1 サーバ装置は、前記第 2 サーバ装置から送られてくる前記ディレクトリイメージに基づき前記データ I / O 要求を処理するとともに、前記ディレクトリイメージ及び前記第 2 サーバ装置から送られてくる前記追加のディレクトリイメージを、前記第 1 ストレージ装置に復元する

情報処理システム。

【請求項 1 5】

請求項 1 4 に記載の情報処理システムであって、

前記第 1 サーバ装置は、前記第 1 ファイルシステムに割り当てられている前記第 1 スト

10

20

30

40

50

レージ装置の記憶領域の残容量が予め設定された閾値未満になっている場合に、前記第1ストレージ装置に記憶しているファイルのデータのうち、所定の選出基準に従って選出したファイルのデータのメタデータについては前記第1ストレージ装置に残すとともに当該データの実体については前記第1ストレージ装置から削除して前記残容量を確保するスタブ化を行い、

前記第2サーバ装置は、前記第1サーバ装置に送信した、前記ディレクトリイメージ又は前記追加のディレクトリイメージが再びスタブ化される現象である再スタブ化の発生頻度が予め設定された閾値以上になっているか否か、もしくは、再スタブ化の発生時間間隔が予め設定された閾値未満になっているか否かを監視し、

前記第2サーバ装置は、再スタブ化の発生頻度が予め設定された閾値以上になっているか、もしくは、再スタブ化の発生時間間隔が予め設定された閾値未満になっている場合に、前記第1サーバ装置への前記ディレクトリイメージ又は前記追加のディレクトリイメージの送信を抑制する

情報処理システム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、情報処理システムにおける障害復旧方法、及び情報処理システムに関する

【背景技術】

【0002】

特許文献1には、階層記憶装置の復旧に要する時間を低減し、階層記憶装置の復旧を高速に行うべく、オペレーティングシステム上で稼動する階層記憶装置において、ファイルの属性情報を含むiノードを有し且つそのiノード番号で当該ファイルを一意に識別するファイルシステムが構築された第1の記憶装置と、ファイルシステムのバックアップデータを含むデータを格納する第2の記憶装置とを有する階層記憶装置の復旧方法において、第2の記憶装置上のバックアップデータから第1の記憶装置上にファイルシステムが復旧される時に、バックアップデータに含まれるiノード番号を用いて、復旧対象ファイルのiノード番号を指定し、指定されたiノード番号をファイルシステムの復旧対象ファイルに割り当てることが記載されている。

【0003】

特許文献2には、HSMにおける名前空間のバックアップの世代管理を効率的に行うべく、一次ストレージと二次ストレージを有するHSMの制御を行うHSM制御方法において、HSMのバックアップ毎に、該バックアップの世代番号を含む情報である世代情報を作成し、HSMにおけるファイル毎の名前空間に関する情報である名前空間情報と、世代情報作成ステップにより作成された世代番号を用いて該名前空間に関する情報が有効である世代番号の範囲を示す有効世代番号範囲とを名前空間情報履歴として管理することが記載されている。

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2005-316708号公報

【特許文献2】特開2008-040699号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

ところで、企業における支店や事業所等に設けられている情報処理装置のデータのバックアップをデータセンタ等に設けられているバックアップ装置で管理している情報処理システムにおいて、情報処理装置側に障害が発生した場合には、その復旧に際してバックアップ装置側のバックアップデータの全てを情報処理装置に復元した後に情報処理装置のサ

10

20

30

40

50

ービスを再開するため、例えば、バックアップデータのサイズが大きい場合はサービスが再開されるまでに長時間を要し、ユーザの業務等への影響が生じる場合もあった。

【0006】

本発明はこのような背景に鑑みてなされたもので、障害復旧に際してサービスを迅速に再開することが可能な、情報処理システムにおける障害復旧方法、及び情報処理システムを提供することを主たる目的とする。

【課題を解決するための手段】

【0007】

上記目的を達成するための本発明の一つは、第1ファイルシステムを有し、データI/O要求を受け付ける第1サーバ装置と、第2ファイルシステムを有し、前記第1サーバ装置に通信可能に接続された第2サーバ装置と、を備え、前記第1サーバ装置は、前記データI/O要求の対象であるファイルのデータを第1ストレージ装置に記憶し、前記第2サーバ装置は、前記データI/O要求の対象であるファイルのデータを第2ストレージ装置に記憶し、前記第1サーバ装置は、前記第1ストレージ装置に記憶しているファイルのデータを前記第2サーバ装置に送信し、前記第2サーバ装置は、前記第1サーバ装置から送られてくる前記データを、前記第1ファイルシステムにおけるディレクトリイメージを保持しつつ前記第2ストレージ装置に記憶している情報処理システムにおける障害の復旧方法であって、前記第2サーバ装置が、障害からの復旧に際し、前記第1サーバ装置が前記データI/O要求の受け付けを開始するのに先立ち、前記第2ストレージ装置に記憶しているディレクトリイメージのうち最上位のディレクトリから所定の下位階層までのディレクトリイメージを前記第1サーバ装置に送信し、前記第1サーバ装置が、前記第2サーバ装置から送られてくる前記ディレクトリイメージを前記第1ストレージ装置に復元した後、前記データI/O要求の受け付けを再開し、前記第1サーバ装置は、前記データI/O要求の受け付け再開後、受け付けた前記データI/O要求を処理するために必要となるディレクトリイメージが前記第1ストレージ装置に復元されていない場合に前記ディレクトリイメージを前記第2サーバ装置に要求し、前記第2サーバ装置は、前記第1サーバ装置から前記要求が送られてくると前記ディレクトリイメージを前記第2ストレージ装置から読み出して送信し、前記第1サーバ装置は、前記第2ストレージ装置から送られてくる前記ディレクトリイメージに基づき前記データI/O要求を処理するとともに、前記ディレクトリイメージを前記第1ストレージ装置に復元することとする。

【0008】

その他本願が開示する課題やその解決方法については、発明の実施形態の欄及び図面により明らかにされる。

【発明の効果】

【0009】

本発明によれば、障害復旧に際してサービスを迅速に再開することができる。

【図面の簡単な説明】

【0010】

【図1】情報処理システム1の概略的な構成を示す図である。

【図2】クライアント装置2のハードウェアの一例である。

【図3】第1サーバ装置3a又は第2サーバ装置3bとして利用可能な情報処理装置のハードウェアの一例である。

【図4】第1ストレージ装置10a又は第2ストレージ装置10bのハードウェアの一例である。

【図5】チャンネル基板11のハードウェアの一例である。

【図6】プロセッサ基板12のハードウェアの一例である。

【図7】ドライブ基板13のハードウェアの一例である。

【図8】ストレージ装置10が備える基本的な機能を示す図である。

【図9】書き込み処理S900を説明するフローチャートである。

【図10】読み出し処理S1000を説明するフローチャートである。

10

20

30

40

50

- 【図 1 1】クライアント装置 2 が備える主な機能を示す図である。
- 【図 1 2】第 1 サーバ装置 3 a が備える主な機能、及び第 1 サーバ装置 3 a において管理される主な情報（データ）を示す図である。
- 【図 1 3】レプリケーション情報管理テーブル 3 3 1 の一例である。
- 【図 1 4】ファイルアクセスログ 3 3 5 の一例である。
- 【図 1 5】第 2 サーバ装置 3 b が備える主な機能、及び第 2 サーバ装置 3 b において管理される主な情報（データ）を示す図である。
- 【図 1 6】リストアログ 3 6 5 の一例である。
- 【図 1 7】抑制フラグ管理テーブル 3 6 6 の一例である。
- 【図 1 8】リコールログ 3 6 7 の一例である。 10
- 【図 1 9】i n o d e を説明する図である。
- 【図 2 0】i n o d e の概念を説明する図である。
- 【図 2 1】i n o d e の概念を説明する図である。
- 【図 2 2】一般的な i n o d e 管理テーブル 1 9 1 2 の一例である。
- 【図 2 3】本実施形態の i n o d e 管理テーブル 1 9 1 2 の一例である。
- 【図 2 4】レプリケーション開始処理 S 2 4 0 0 を説明する図である。
- 【図 2 5】スタブ化候補選出処理 S 2 5 0 0 を説明する図である。
- 【図 2 6】スタブ化処理 S 2 6 0 0 を説明する図である。
- 【図 2 7】レプリケーションファイル更新処理 S 2 7 0 0 を説明する図である。
- 【図 2 8】レプリケーションファイル参照処理 S 2 8 0 0 を説明する図である。 20
- 【図 2 9】同期処理 S 2 9 0 0 を説明する図である。
- 【図 3 0】メタデータアクセス処理 S 3 0 0 0 を説明する図である。
- 【図 3 1】スタブ化ファイル実体参照処理 S 3 1 0 0 を説明する図である。
- 【図 3 2】スタブ化ファイル実体更新処理 S 3 2 0 0 を説明する図である。
- 【図 3 3】仮想マシン復旧処理 S 3 3 0 0 を説明する図である。
- 【図 3 4】ディレクトリイメージ事前回復処理 S 3 4 0 0 を説明する図である。
- 【図 3 5】オンデマンド復元処理 S 3 5 0 0 を説明する図である。
- 【図 3 6】第 1 ストレージ装置 1 0 a にディレクトリイメージが復元されていく様子を説明する図である。
- 【図 3 7】オンデマンド復元処理（復元対象追加有）S 3 7 0 0 を説明する図である。 30
- 【図 3 8】再スタブ化回避処理 S 3 8 0 0 を説明する図である。
- 【図 3 9】レプリケーション開始処理 S 2 4 0 0 の詳細を説明するフローチャートである。
- 【図 4 0】スタブ化候補選出処理 S 2 5 0 0 の詳細を説明するフローチャートである。
- 【図 4 1】スタブ化処理 S 2 6 0 0 の詳細を説明するフローチャートである。
- 【図 4 2】レプリケーションファイル更新処理 S 2 7 0 0 の詳細を説明するフローチャートである。
- 【図 4 3】レプリケーションファイル参照処理 S 2 8 0 0 の詳細を説明するフローチャートである。
- 【図 4 4】同期処理 S 2 9 0 0 の詳細を説明するフローチャートである。 40
- 【図 4 5】メタデータアクセス処理 S 3 0 0 0 の詳細を説明するフローチャートである。
- 【図 4 6】スタブ化ファイル実体参照処理 S 3 1 0 0 の詳細を説明するフローチャートである。
- 【図 4 7】スタブ化ファイル実体更新処理 S 3 2 0 0 の詳細を説明するフローチャートである。
- 【図 4 8】仮想マシン復旧処理 S 3 3 0 0 及びディレクトリイメージ事前回復処理 S 3 4 0 0 の詳細を説明するフローチャートである。
- 【図 4 9】オンデマンド復元処理 S 3 5 0 0 の詳細を説明するフローチャートである。
- 【図 5 0】オンデマンド復元処理（復元対象追加有）S 3 7 0 0 の詳細を説明するフローチャートである。 50

【図51】オンデマンド復元処理（復元対象追加有）S3700の詳細を説明するフローチャート（図50の続き）である。

【図52】再スタブ化回避処理S3800の詳細を説明するフローチャートである。

【発明を実施するための形態】

【0011】

以下、発明を実施するための形態について図面とともに説明する。

【0012】

図1に実施形態として説明する情報処理システム1の概略的な構成を示している。同図に示すように、本実施形態として例示する情報処理システム1は、商社や電機メーカ等の企業における支点や事業所などのように、ユーザが実際に業務を行う拠点（以下、エッジ50（Edge）と称する。）に設けられるハードウェアと、データセンタのように、情報処理システム（アプリケーションサーバ/ストレージシステム等）の管理やクラウドサービスの提供などを行う拠点（以下、コア51（Core）と称する。）に設けられるハードウェアと、を備える。

10

【0013】

同図に示すように、エッジ50には、第1サーバ装置3a、第1ストレージ装置10a、及びクライアント装置2が設けられている。またコア51には、第2サーバ装置3b及び第2ストレージ装置10bが設けられている。

【0014】

エッジに設けられている第1サーバ装置3aは、例えば、エッジに設けられているクライアント装置2に対してファイルを単位としたデータの管理機能を提供するファイルシステムを備えたファイルストレージ装置である。またコアに設けられている第2サーバ装置3bは、例えば、エッジに設けられている第1ストレージ装置10aに対してデータのアーカイブ（書庫）先として機能するアーカイブ装置である。

20

【0015】

同図に示すように、クライアント装置2と第1サーバ装置3aとは、通信ネットワーク5を介して通信可能に接続している。また第1サーバ装置3aと第1ストレージ装置10aとは、第1ストレージネットワーク6aを介して通信可能に接続している。また第2サーバ装置3bと第2ストレージ装置10bとは、第2ストレージネットワーク6bを介して通信可能に接続している。また第1サーバ装置3aと第2サーバ装置3bとは、通信ネットワーク7を介して通信可能に接続している。

30

【0016】

通信ネットワーク5及び通信ネットワーク7は、例えば、LAN（Local Area Network）、WAN（Wide Area Network）、インターネット、公衆通信網、専用線などである。第1ストレージネットワーク6a及び第2ストレージネットワーク6bは、例えば、LAN、WAN、SAN（Storage Area Network）、インターネット、公衆通信網、専用線などである。

【0017】

通信ネットワーク5、通信ネットワーク7、第1ストレージネットワーク6a、及び第2ストレージネットワーク6bを介して行われる通信は、例えば、TCP/IP、iSCSI（internet Small Computer System Interface）、ファイバーチャネルプロトコル（Fibre Channel Protocol）、FICON（Fibre Connection）（登録商標）、ESCON（Enterprise System Connection）（登録商標）、ACONARC（Advanced Connection Architecture）（登録商標）、FIBARC（Fibre Connection Architecture）（登録商標）などのプロトコルに従って行われる。

40

【0018】

クライアント装置2は、第1サーバ装置3aを介して第1ストレージ装置10aが提供する記憶領域を利用する情報処理装置（コンピュータ）であり、例えば、パーソナルコンピュータ、オフィスコンピュータなどである。クライアント装置2では、ファイルシステム、カーネルやドライバなどのソフトウェアモジュールによって実現されるオペレーティ

50

ングシステム (Operating System)、及びアプリケーションなどが機能している。

【0019】

図2にクライアント装置2のハードウェアを示している。同図に示すように、クライアント装置2は、CPU21、揮発性または不揮発性のメモリ22 (RAMまたはROM)、記憶装置23 (例えばハードディスクドライブ、半導体記憶装置 (SSD (Solid State Drive)))、キーボードやマウス等の入力装置24、液晶モニターやプリンタ等の出力装置25、及びNIC (Network Interface Card) (以下、LANアダプタ261と称する。)等の通信インタフェース (通信I/F26と称する。)を備えている。

【0020】

第1サーバ装置3aは、第1ストレージ装置10aが提供する記憶領域を利用してクライアント装置2に情報処理サービスを提供する情報処理装置である。第1サーバ装置3aは、パーソナルコンピュータ、メインフレーム (Mainframe)、オフィスコンピュータなどを用いて構成されている。第1サーバ装置3aは、第1ストレージ装置10aが提供する記憶領域へのアクセスに際し、データI/O要求 (データ書き込み要求、データ読み出し要求等)を含んだデータフレーム (以下、フレームと略記する。)を、第1ストレージネットワーク6aを介して第1ストレージ装置10aに送信する。尚、上記フレームは、例えば、ファイバチャネルのフレーム (FCフレーム (FC: Fibre Channel)) である。

10

【0021】

第2サーバ装置3bは、第2ストレージ装置10bが提供する記憶領域を利用して情報処理を行う情報処理装置である。第2サーバ装置3bは、パーソナルコンピュータ、メインフレーム、オフィスコンピュータなどを用いて構成されている。第2サーバ装置3bは、第2ストレージ装置10bが提供する記憶領域へのアクセスに際し、データI/O要求を含んだフレームを、第2ストレージネットワーク6bを介して第2ストレージ装置10bに送信する。

20

【0022】

図3に第1サーバ装置3aのハードウェアを示している。同図に示すように、第1サーバ装置3aは、CPU31、揮発性または不揮発性のメモリ32 (RAMまたはROM)、記憶装置33 (例えばハードディスクドライブ、半導体記憶装置 (SSD))、キーボードやマウス等の入力装置34、液晶モニターやプリンタ等の出力装置35、NIC (以下、LANアダプタ361と称する。)やHBA (以下、FCアダプタ362と称する。)等の通信インタフェース (通信I/F36と表記する。)及びタイマ回路やRTC等を用いて構成される計時装置37を備えている。尚、コア側に存在する第2サーバ装置3bも第1サーバ装置3aと同一又は類似のハードウェア構成を有する。

30

【0023】

図4に第1ストレージ装置10aのハードウェアを示している。第1ストレージ装置10aは、例えばディスクアレイ装置である。尚、コア側に存在する第2ストレージ装置10bも第1ストレージ装置10aと同一又は類似のハードウェア構成を有する。ストレージ装置10は、サーバ装置3 (第1サーバ装置3a又は第2サーバ装置3b。以下同様) から送られてくるデータI/O要求を受け付け、受け付けたデータI/O要求に応じて記録媒体にアクセスしてサーバ装置3にデータやレスポンスを送信する。

40

【0024】

同図に示すように、ストレージ装置10は、一つ以上のチャンネル基板11、一つ以上のプロセッサ基板12 (Micro Processor)、一つ以上のドライブ基板13、キャッシュメモリ14 (Cache Memory)、共有メモリ15 (Shared Memory)、内部スイッチ16、記憶装置17、及び保守装置18 (SVP: Service Processor)を備える。チャンネル基板11、プロセッサ基板12、ドライブ基板13、キャッシュメモリ14、及び共有メモリ15は、内部スイッチ16を介して互いに通信可能に接続されている。

【0025】

チャンネル基板11は、サーバ装置3から送られてくるフレームを受信し、受信したフレ

50

ームに含まれているデータ I/O 要求についての処理の応答（例えば読み出したデータ、読み出し完了報告、書き込み完了報告）を含んだフレームをサーバ装置 3 に送信する。

【 0 0 2 6 】

プロセッサ基板 1 2 は、チャンネル基板 1 1 が受信したフレームに含まれている上記データ I/O 要求に応じて、チャンネル基板 1 1、ドライブ基板 1 3、及びキャッシュメモリ 1 4 の間で行われるデータ転送（DMA（Direct Memory Access）等を用いた高速大容量のデータ転送）に関する処理を行う。プロセッサ基板 1 2 は、キャッシュメモリ 1 4 を介して行われる、チャンネル基板 1 1 とドライブ基板 1 3 との間のデータ（記憶装置 1 7 から読み出したデータ、記憶装置 1 7 に書き込むデータ）の転送（引き渡し）や、キャッシュメモリ 1 4 に格納されるデータのステージング（記憶装置 1 7 からのデータの読み出し）及びデステージング（記憶装置 1 7 への書き出し）等を行う。

10

【 0 0 2 7 】

キャッシュメモリ 1 4 は、高速アクセスが可能な RAM（Random Access Memory）を用いて構成されている。キャッシュメモリ 1 4 には、記憶装置 1 7 に書き込まれるデータ（以下、書き込みデータと称する。）や記憶装置 1 7 から読み出されたデータ（以下、読み出しデータと記載する。）等が格納される。共有メモリ 1 5 には、ストレージ装置 1 0 の制御に用いられる様々な情報が格納される。

【 0 0 2 8 】

ドライブ基板 1 3 は、記憶装置 1 7 からのデータの読み出しや記憶装置 1 7 へのデータの書き込みの際に記憶装置 1 7 と通信を行う。内部スイッチ 1 6 は、例えば高速クロスバースイッチ（Cross Bar Switch）を用いて構成される。尚、内部スイッチ 1 6 を介して行われる通信は、例えば、ファイバーチャネル、iSCSI、TCP/IP 等のプロトコルに従って行われる。

20

【 0 0 2 9 】

記憶装置 1 7 は、複数の記憶ドライブ 1 7 1 を備えて構成されている。記憶ドライブ 1 7 1 は、例えば、SAS（Serial Attached SCSI）、SATA（Serial ATA）、FC（Fibre Channel）、PATA（Parallel ATA）、SCSI 等のタイプのハードディスクドライブ、半導体記憶装置（SSD）等である。

【 0 0 3 0 】

記憶装置 1 7 は、記憶ドライブ 1 7 1 を、例えば、RAID（Redundant Arrays of Inexpensive（or Independent）Disks）等の方式で制御することによって提供される論理的な記憶領域を単位としてサーバ装置 3 に記憶装置 1 7 の記憶領域を提供する。この論理的な記憶領域は、例えば RAID グループ（パリティグループ（Parity Group））を用いて構成される論理装置（LDEV 1 7 2（LDEV: Logical Device））である。

30

【 0 0 3 1 】

またストレージ装置 1 0 は、サーバ装置 3 に対して LDEV 1 7 2 を用いて構成される論理的な記憶領域（以下、LU（Logical Unit、Logical Volume、論理ボリューム）と称する。）を提供する。ストレージ装置 1 0 は、LU と LDEV 1 7 2 との対応（関係）を管理しており、ストレージ装置 1 0 は、この対応に基づき LU に対応する LDEV 1 7 2 の特定、もしくは、LDEV 1 7 2 に対応する LU の特定を行う。

40

【 0 0 3 2 】

図 5 にチャンネル基板 1 1 のハードウェア構成を示している。同図に示すように、チャンネル基板 1 1 は、サーバ装置 3 と通信するためのポート（通信ポート）を有する外部通信インタフェース（以下、外部通信 I/F 1 1 1 と表記する。）、プロセッサ 1 1 2（フレーム処理チップ及びフレーム転送チップを含む）、メモリ 1 1 3、プロセッサ基板 1 2 と通信するためのポート（通信ポート）を有する内部通信インタフェース（以下、内部通信 I/F 1 1 4 と表記する。）を備えている。

【 0 0 3 3 】

外部通信 I/F 1 1 1 は、NIC（Network Interface Card）や HBA（Host Bus Adaptor）等を用いて構成されている。プロセッサ 1 1 2 は、CPU（Central Processing U

50

nit)、MPU (Micro Processing Unit) 等を用いて構成されている。メモリ 113は、RAM (Random Access Memory)、ROM (Read Only Memory) である。メモリ 113にはマイクロプログラムが格納されている。プロセッサ 112が、メモリ 113から上記マイクロプログラムを読み出して実行することにより、チャンネル基板 11が提供する各種の機能が実現される。内部通信 I/F 114は、内部スイッチ 16を介してプロセッサ基板 12、ドライブ基板 13、キャッシュメモリ 14、共有メモリ 15と通信する。

#### 【0034】

図6にプロセッサ基板 12のハードウェア構成を示している。プロセッサ基板 12は、内部通信インタフェース(以下、内部通信 I/F 121と表記する。)、プロセッサ 122、及び共有メモリ 15に比べてプロセッサ 122からのアクセス性能が高い(高速アクセスが可能な)メモリ 123(ローカルメモリ)を備えている。メモリ 123にはマイクロプログラムが格納されている。プロセッサ 122が、メモリ 123から上記マイクロプログラムを読み出して実行することにより、プロセッサ基板 12が提供する各種の機能が実現される。

#### 【0035】

内部通信 I/F 121は、内部スイッチ 16を介してチャンネル基板 11、ドライブ基板 13、キャッシュメモリ 14、及び共有メモリ 15と通信を行う。プロセッサ 122は、CPU、MPU、DMA (Direct Memory Access) 等を用いて構成されている。メモリ 123は、RAM又はROMである。プロセッサ 122は、メモリ 123及び共有メモリ 15のいずれにもアクセスすることができる。

#### 【0036】

図7にドライブ基板 13のハードウェア構成を示している。ドライブ基板 13は、内部通信インタフェース(以下、内部通信 I/F 131と表記する。)、プロセッサ 132、メモリ 133、及びドライブインタフェース(以下、ドライブ I/F 134と表記する。)を備えている。メモリ 133にはマイクロプログラムが格納されている。プロセッサ 132が、メモリ 133から上記マイクロプログラムを読み出して実行することにより、ドライブ基板 13が提供する各種の機能が実現される。内部通信 I/F 131は、内部スイッチ 16を介して、チャンネル基板 11、プロセッサ基板 12、キャッシュメモリ 14、及び共有メモリ 15と通信する。プロセッサ 132は、CPU、MPU等を用いて構成されている。メモリ 133は、例えば、RAM、ROMである。ドライブ I/F 134は、記憶装置 17との間で通信を行う。

#### 【0037】

図4に示した保守装置 18は、ストレージ装置 10の各構成要素の制御や状態監視を行う。保守装置 18は、パーソナルコンピュータやオフィスコンピュータ等である。保守装置 18は、内部スイッチ 16又はLAN等の通信手段を介してチャンネル基板 11、プロセッサ基板 12、ドライブ基板 13、キャッシュメモリ 14、共有メモリ 15、及び内部スイッチ 16等のストレージ装置 10の構成要素と随時通信を行い、各構成要素から稼働情報等を取得して管理装置 19に提供する。また保守装置 18は管理装置 19から送られてくる制御情報や稼働情報に基づき、各構成要素の設定や制御、保守(ソフトウェアの導入や更新を含む)を行う。

#### 【0038】

管理装置 19は、保守装置 18とLAN等を介して通信可能に接続されるコンピュータである。管理装置 19はストレージ装置 10の制御や監視のためのGUI (Graphical User Interface) やCLI (Command Line Interface) 等を用いたユーザインタフェースを備える。

#### 【0039】

図8にストレージ装置 10が備える基本的な機能を示している。同図に示すように、ストレージ装置 10は、I/O処理部 811を備えている。I/O処理部 811は、記憶装置 17への書き込みに関する処理を行うデータ書き込み処理部 8111、記憶装置 17からのデータの読み出しに関する処理を行うデータ読み出し処理部 8112を備えている。

10

20

30

40

50

## 【 0 0 4 0 】

尚、I/O処理部811の機能は、ストレージ装置10のチャンネル基板11やプロセッサ基板12、ドライブ基板13が備えるハードウェアにより、又はプロセッサ112, 122, 132が、メモリ113, 123, 133に格納されているマイクロプログラムを読み出して実行することにより実現される。

## 【 0 0 4 1 】

図9は、ストレージ装置10（第1ストレージ装置10a又は第2ストレージ装置10b。以下同様。）がサーバ装置3（第1サーバ装置3a又は第2サーバ装置3b）からデータ書き込み要求を含んだフレームを受信した場合に、I/O処理部81のデータ書き込み処理部8111によって行われる基本的な処理（以下、書き込み処理S900と称する。）を説明するフローチャートである。以下、同図とともに書き込み処理S900について説明する。尚、以下の説明において、符号の前に付している「S」の文字は処理ステップを意味している。

10

## 【 0 0 4 2 】

同図に示すように、まずサーバ装置3から送信されてくるデータ書き込み要求のフレームが、ストレージ装置10のチャンネル基板11によって受信される（S911, S912）。

## 【 0 0 4 3 】

チャンネル基板11は、サーバ装置3からデータ書き込み要求を含んだフレームを受信すると、その旨をプロセッサ基板12に通知する（S913）。

20

## 【 0 0 4 4 】

プロセッサ基板12は、チャンネル基板11から上記通知を受信すると（S921）、当該フレームのデータ書き込み要求に基づくドライブ書き込み要求を生成し、書き込みデータをキャッシュメモリ14に格納し、チャンネル基板11に上記通知の受信通知を応答する（S922）。またプロセッサ基板12は、生成したドライブ書き込み要求をドライブ基板13に送信する（S923）。

## 【 0 0 4 5 】

一方、チャンネル基板11は、プロセッサ基板12からの上記応答を受信すると、サーバ装置3に完了報告を送信し（S914）、サーバ装置3は、チャンネル基板11から完了報告を受信する（S915）。

30

## 【 0 0 4 6 】

ドライブ基板13は、プロセッサ基板12からドライブ書き込み要求を受信すると、受信したドライブ書き込み要求を書き込み処理待ちキューに登録する（S924）。

## 【 0 0 4 7 】

ドライブ基板13は、書き込み処理待ちキューからドライブ書き込み要求を随時読み出し（S925）、読み出したドライブ書き込み要求に指定されている書き込みデータをキャッシュメモリ14から読み出して、読み出した書き込みデータを記憶装置（記憶ドライブ171）に書き込む（S926）。そしてドライブ基板13は、ドライブ書き込み要求について書き込みデータの書き込みが完了した旨の報告（完了報告）をプロセッサ基板12に通知する（S927）。

40

## 【 0 0 4 8 】

プロセッサ基板12は、ドライブ基板13から送られてきた完了報告を受信する（S928）。

## 【 0 0 4 9 】

図10は、ストレージ装置10が、サーバ装置3からデータ読み出し要求を含んだフレームを受信した場合に、ストレージ装置10のI/O処理部811の読み出し処理部8112によって行われるI/O処理（以下、読み出し処理S1000と称する。）を説明するフローチャートである。以下、同図とともに読み出し処理S1000について説明する。

## 【 0 0 5 0 】

50

同図に示すように、まずサーバ装置 3 から送信されてくるフレームが、ストレージ装置 10 のチャンネル基板 11 によって受信される (S1011, S1012)。

【0051】

チャンネル基板 11 は、サーバ装置 3 からデータ読み出し要求を含んだフレームを受信すると、その旨をプロセッサ基板 12 及びドライブ基板 13 に通知する (S1013)。

【0052】

ドライブ基板 13 は、チャンネル基板 11 から上記通知を受信すると (S1014)、当該フレームに含まれているデータ読み出し要求に指定されているデータ (例えば LBA (Logical Block Address) によって指定される) を、記憶装置 (記憶ドライブ 171) から読み出す (S1015)。尚、キャッシュメモリ 14 に読み出しデータが存在する場合 (キャッシュヒットした場合) には、記憶装置 17 からの読み出し処理 (S1015) は省略される。

10

【0053】

プロセッサ基板 12 は、ドライブ基板 13 によって読み出されたデータをキャッシュメモリ 14 に書き込む (S1016)。そしてプロセッサ基板 12 は、キャッシュメモリ 14 に書き込んだデータをチャンネル基板 11 に随時転送する (S1017)。

【0054】

チャンネル基板 11 は、プロセッサ基板 12 から随時送られてくる読み出しデータを受信すると、それらをサーバ装置 3 に順次送信する (S1018)。読み出しデータの送信が完了すると、チャンネル基板 11 は、サーバ装置 3 に完了報告を送信する (S1019)。サーバ装置 3 は、読み出しデータ及び完了報告を受信する (S1020, S1021)。

20

【0055】

図 11 にクライアント装置 2 が備える主な機能を示している。同図に示すように、クライアント装置 2 は、アプリケーション 211、ファイルシステム 212、及びカーネル/ドライバ 213 の各機能を備える。尚、これらの機能は、クライアント装置 2 の CPU 21 が、メモリ 22 や記憶装置 23 に格納されているプログラムを読み出して実行することにより実現される。

【0056】

ファイルシステム 212 は、クライアント装置 2 に対してファイル単位又はディレクトリ単位での論理ボリューム (LU) への I/O 機能を実現する。ファイルシステム 213 は、例えば、FAT (File Allocation Table)、NTFS、HFS (Hierarchical File System)、ext2 (second extended file system)、ext3 (third extended file system)、ext4 (fourth extended file system)、UDF (Universal Disk Format)、HPFS (High Performance File system)、JFS (Journaled File System)、UFS (Unix File System)、VTOC (Volume Table Of Contents)、XFS 等である。

30

【0057】

カーネル/ドライバ 213 は、オペレーティングシステムのソフトウェアを構成しているカーネルモジュールやドライバモジュールを実行することにより実現される。カーネルモジュールには、クライアント装置 2 において実行されるソフトウェアについて、プロセスの管理、プロセスのスケジューリング、記憶領域の管理、ハードウェアからの割り込み要求のハンドリング等、オペレーティングシステムが備える基本的な機能を実現するためのプログラムが含まれている。ドライバモジュールには、クライアント装置 2 を構成しているハードウェアやクライアント装置 2 に接続して用いられる周辺機器とカーネルモジュールとが通信するためのプログラム等が含まれている。

40

【0058】

図 12 に第 1 サーバ装置 3a が備える主な機能、及び第 1 サーバ装置 3a において管理される主な情報 (データ) を示している。同図に示すように、第 1 サーバ装置 3a では、仮想環境を提供する仮想化制御部 305、及びこの仮想化制御部 305 の制御の下で動作する 1 つ以上の仮想マシン 310 が実現されている。

【0059】

50

各仮想マシン310では、ファイル共有処理部311、ファイルシステム312、データ操作要求受付部313、データ複製/移動処理部314、ファイルアクセスログ取得部317、及びカーネル/ドライバ318の各機能が実現されている。

【0060】

尚、仮想環境は、第1サーバ装置3aのハードウェアと仮想化制御部305との間にオペレーティングシステムを介在させるいわゆるホストOS型、もしくは第1サーバ装置3aのハードウェアと仮想化制御部305との間にオペレーティングシステムを介在させないハイパーバイザー型のいずれの方式で実現されていてもよい。またデータ操作要求受付部313、データ複製/移動処理部314、ファイルアクセスログ取得部317の各機能は、ファイルシステム312の機能として実現してもよいし、ファイルシステム312とは独立した機能として実現してもよい。

10

【0061】

同図に示すように、各仮想マシン310は、レプリケーション情報管理テーブル331、ファイルアクセスログ335などの情報(データ)を管理している。これらの情報は、第1ストレージ10aから第1サーバ装置3aに随時読み出されて第1サーバ装置3aのメモリ32や記憶装置33に格納される。

【0062】

同図に示す機能のうち、ファイル共有処理部311は、クライアント装置2にファイルの共有環境を提供する。ファイル共有処理部311は、例えば、NFS(Network File System)、CIFS(Common Internet File System)、AFS(Andrew File System)等のプロトコルに従った機能を提供する。

20

【0063】

ファイルシステム312は、クライアント装置2に対して、第1ストレージ装置10aによって提供される論理ボリューム(LU)に管理されるファイル(又はディレクトリ)に対するI/O機能を提供する。ファイルシステム312は、例えばFAT(File Allocation Table)、NTFS、HFS(Hierarchical File System)、ext2(second extended file system)、ext3(third extended file system)、ext4(fourth extended file system)、UDF(Universal Disk Format)、HPFS(High Performance File system)、JFS(Journaled File System)、UFS(Unix File System)、VTOC(Volume Table Of Contents)、XFS等である。

30

【0064】

データ操作要求受付部313は、クライアント装置2から送信されてくるデータの操作に関する要求(以下、データ操作要求と称する。)を受け付ける。データ操作要求には、後述する、レプリケーション開始要求、レプリケーションファイルに対する更新要求、レプリケーションファイルに対する参照要求、同期要求、メタデータに対するアクセス要求、ファイルの実体に対する参照要求、リコール要求、スタブ化ファイルの実体に対する更新要求などがある。

【0065】

尚、スタブ化とは、ファイル(又はディレクトリ)のデータのメタデータについては第1ストレージ装置10aにて保持するが、ファイル(又はディレクトリ)のデータの実体については第1ストレージ装置10a側では管理せず、第2ストレージ装置10bにおいてのみ保持するようにすることをいう。スタブ化されたファイル(又はディレクトリ)に対してそのファイル(又はディレクトリ)の実体が必要になるようなデータI/O要求を第1サーバ装置3aが受け付けた場合には、そのファイル(又はディレクトリ)の実体が第2ストレージ装置10bから第1ストレージ装置10aに送信される(書き戻される(以下、リコールと称する。))。

40

【0066】

データ複製/移動処理部314は、後述する、レプリケーション開始処理S2400、スタブ化候補選出処理S2500、同期処理S2900、スタブ化ファイル実体参照処理S3100、スタブ化ファイル実体更新処理S3200、仮想マシン復旧処理S3300

50

、ディレクトリイメージ事前回復処理 S 3 4 0 0、オンデマンド復元処理 S 3 5 0 0、オンデマンド復元処理（復元対象追加有） S 3 7 0 0、再スタブ化回避処理 S 3 8 0 0 などにおいて、第 1 サーバ装置 3 a と第 2 サーバ装置 3 b との間、もしくは、第 1 ストレージ装置 1 0 a と第 2 ストレージ装置 1 0 b との間での制御情報（フラグやテーブルを含む）の授受やデータ（ファイルのメタデータや実体を含む）の転送、レプリケーション情報管理テーブル 3 3 1、メタデータ 3 3 2 などの各種テーブルの管理を行う。

【 0 0 6 7 】

図 1 2 に示すカーネル/ドライバ 3 1 8 は、オペレーティングシステムのソフトウェアを構成しているカーネルモジュールやドライバモジュールを実行することにより実現される。カーネルモジュールには、第 1 サーバ装置 3 a において実行されるソフトウェアについて、プロセスの管理、プロセスのスケジューリング、記憶領域の管理、ハードウェアからの割り込み要求のハンドリング等、オペレーティングシステムが備える基本的な機能を実現するためのプログラムが含まれる。ドライバモジュールには、第 1 サーバ装置 3 a を構成しているハードウェアや第 1 サーバ装置 3 a に接続して用いられる周辺機器とカーネルモジュールとが通信するためのプログラムが含まれる。

10

【 0 0 6 8 】

図 1 2 に示すファイルアクセスログ取得部 3 1 7 は、ストレージ装置 1 0 の論理ボリューム（LU）に格納されているファイルへのアクセス（ファイルの更新（Write, Update）、ファイルの読み出し（Read）、ファイルのオープン（Open）、ファイルのクローズ（Close）等）が行われると、そのアクセスの内容（履歴）を示す情報（以下、アクセスログと称する。）を、計時装置 3 7 から取得される日時情報に基づくタイムスタンプを付与してファイルアクセスログ 3 3 5 として記憶する。

20

【 0 0 6 9 】

図 1 3 にレプリケーション情報管理テーブル 3 3 1 の一例を示している。同図に示すように、レプリケーション情報管理テーブル 3 3 1 には、レプリケーション先となるホスト名 3 3 1 1（例えば IP アドレス等のネットワークアドレス。）、スタブ化するか否かの判断に用いる閾値 3 3 1 2（後述するスタブ化閾値）が設定される。

【 0 0 7 0 】

図 1 4 にファイルアクセスログ 3 3 5 の一例を示している。同図に示すように、ファイルアクセスログ 3 3 5 には、アクセス日時 3 3 5 1、ファイル名 3 3 5 2、及びユーザ ID 3 3 5 3 の各項目からなる一つ以上のレコードで構成されるアクセスログが記録される。

30

【 0 0 7 1 】

このうちアクセス日時 3 3 5 1 には、そのファイル（又はディレクトリ）に対するアクセスがされた日時が設定される。ファイル名 3 3 5 2 には、アクセス対象となったファイル（又はディレクトリ）のファイル名（又はディレクトリ名）が設定される。ユーザ ID 3 3 5 3 には、そのファイル（又はディレクトリ）にアクセスしたユーザのユーザ ID が設定される。

【 0 0 7 2 】

図 1 5 に第 2 サーバ装置 3 b が備える主な機能、及び第 2 サーバ装置 3 b において管理される主な情報（データ）を示している。同図に示すように、第 2 サーバ装置 3 b は、ファイル共有処理部 3 5 1、ファイルシステム 3 5 2、データ複製/移動処理部 3 5 4、及びカーネル/ドライバ 3 5 8 の各機能を備えている。尚、データ複製/移動処理部 3 5 4 の機能は、ファイルシステム 3 5 2 の機能として実現してもよいし、ファイルシステム 3 5 2 とは独立した機能として実現してもよい。

40

【 0 0 7 3 】

また同図に示すように、第 2 サーバ装置 3 b は、リストアログ 3 6 5、抑制フラグ管理テーブル 3 6 6、リコールログ 3 6 7、及びファイルアクセスログ 3 6 8 を管理している。

【 0 0 7 4 】

50

ファイル共有処理部 351 は、第 1 サーバ装置 3a に対してファイルの共有環境を提供する。ファイル共有処理部 351 は、例えば NFS、CIFS、AFS 等のプロトコルを用いて実現されている。

【0075】

ファイルシステム 352 は、第 2 ストレージ装置 10b によって提供される論理ボリューム (LU) を用い、第 1 サーバ装置 3a に対してファイル単位又はディレクトリ単位での論理ボリューム (LU) への I/O 機能を提供する。ファイルシステム 352 は、例えば FAT、NTFS、HFS、ext2、ext3、ext4、UDF、HPFS、JFS、UFS、VTOC、XFS などである。

【0076】

データ複製 / 移動処理部 354 は、第 1 ストレージ装置 10a と第 2 ストレージ装置 10b との間でのデータの移動や複製に関する処理を行う。

【0077】

カーネル / ドライバ 358 は、オペレーティングシステムのソフトウェアを構成しているカーネルモジュールやドライバモジュールを実行することにより実現される。カーネルモジュールには、第 2 サーバ装置 3b において実行されるソフトウェアについて、プロセスの管理、プロセスのスケジューリング、記憶領域の管理、ハードウェアからの割り込み要求のハンドリング等、オペレーティングシステムが備える基本的な機能を実現するためのプログラムが含まれる。ドライバモジュールには、第 2 サーバ装置 3b を構成しているハードウェアや第 2 サーバ装置 3b に接続して用いられる周辺機器とカーネルモジュールとが通信するためのプログラムが含まれる。

【0078】

図 16 にリストアログ 365 の一例を示している。リストアログ 365 は、後述するディレクトリイメージの復元 (リストア) が行われた場合に、第 1 サーバ装置 3a 又は第 2 サーバ装置 3b によって復元に関する処理の内容が記録される。同図に示すように、リストアログ 365 は、日時 3651、イベント 3652、及びリストア対象ファイル 3653 の各項目を有する一つ以上のレコードで構成される。

【0079】

このうち日時 3651 には、リストアに関するイベントが実行された日時が設定される。イベント 3652 には、実行されたイベントの内容を示す情報 (リストア開始、リストア実行等) が設定される。リストア対象ファイル 3653 には、リストアの対象となったファイル (又はディレクトリ) を特定する情報 (パス名、ファイル名 (又はディレクトリ名) 等) が設定される。

【0080】

図 17 に抑制フラグ管理テーブル 366 の一例を示している。抑制フラグ管理テーブル 366 の内容は第 2 サーバ装置 3b によって管理される。同図に示すように、抑制フラグ管理テーブル 366 には、後述する再スタブ化回避処理 S3800 において用いられる抑制フラグ 3661、及び抑制フラグ 3661 の最終更新日時 3662 が管理されている。

【0081】

図 18 にリコールログ 367 の一例を示している。リコールログ 367 の内容は第 2 サーバ装置 3b によって生成される。リコールログ 367 には、第 2 サーバ装置 3b が第 1 サーバ装置 3a から受け付けたリコール要求の履歴が管理される。同図に示すように、リコールログ 367 は、日時 3671 及びリコール対象ファイル 3672 の各項目を有する一つ以上のレコードで構成されている。日時 3671 には、リコール要求を受け付けた日時が設定される。リコール対象ファイル 3672 には、受け付けたリコール要求に指定されているリコール対象となるファイル (又はディレクトリ) を特定する情報 (パス名、ファイル名等) が設定される。

【0082】

第 2 サーバ装置 3b が管理しているファイルアクセスログ 368 の内容は、第 1 サーバ装置 3a におけるファイルアクセスログ 335 の内容に基本的に一致している。第 1 サー

10

20

30

40

50

バ装置 3 a からファイルアクセスログ 3 3 5 の内容を第 2 サーバ装置 3 b に随時通知することにより、両者の同一性が確保されている。

【 0 0 8 3 】

次に第 1 サーバ装置 3 a が備えるファイルシステム 3 1 2 ( 第 2 サーバ 3 b が備えるファイルシステム 3 5 2 も同様である。 ) について詳細に説明する。

【 0 0 8 4 】

図 1 9 は、ファイルシステム 3 1 2 が論理ボリューム ( L U ) 上に管理するデータの構造 ( 以下、ファイルシステム構造 1 9 0 0 と称する。 ) の一例である。同図に示すように、ファイルシステム構造 1 9 0 0 には、スーパーブロック 1 9 1 1、i n o d e 管理テーブル 1 9 1 2、ファイルの実体 ( データ ) が格納されるデータブロック 1 9 1 3 の各記憶領域が含まれている。

10

【 0 0 8 5 】

このうちスーパーブロック 1 9 1 1 には、ファイルシステム 3 1 2 に関する情報 ( ファイルシステムが取り扱う記憶領域の容量、使用量、空き容量等 ) が格納される。スーパーブロック 1 9 1 1 は、原則としてディスク区画 ( 論理ボリューム ( L U ) 上に設定されるパーティション ) ごとに設けられる。スーパーブロック 1 9 1 1 に格納される上記情報の具体例として、区画内のデータブロック数、ブロックサイズ、空きブロック数、空き i n o d e 数、当該区画のマウント数、最新の整合性チェック時からの経過時間等がある。

【 0 0 8 6 】

i n o d e 管理テーブル 1 9 1 2 には、論理ボリューム ( L U ) に格納されるファイル ( 又はディレクトリ ) の管理情報 ( 以下、i n o d e と称する。 ) が格納される。ファイルシステム 3 1 2 は、1 つのファイル ( 又はディレクトリ ) に対して 1 つの i n o d e を対応させて管理している。i n o d e のうちディレクトリに関する情報のみを含むものはディレクトリエントリと称される。ファイルに対するアクセスが行われる場合にはディレクトリエントリを参照してアクセス対象のファイルのデータブロックにアクセスする。例えば、「 /home/user-01/a.txt 」というファイルにアクセスする場合には、図 2 0 に示すように、i n o d e 番号 2 1 0 1 5 1 0 0 とディレクトリエントリを順に辿ってアクセス対象ファイルのデータブロックにアクセスする。

20

【 0 0 8 7 】

図 2 1 に一般的なファイルシステム ( 例えば、U N I X ( 登録商標 ) 系のオペレーティングシステムが備えるファイルシステム ) における i n o d e の概念を示している。また図 2 2 に i n o d e 管理テーブル 1 9 1 2 の一例を示している。

30

【 0 0 8 8 】

これらの図に示すように、i n o d e には、個々の i n o d e を区別する識別子である i n o d e 番号 2 2 1 1、当該ファイル ( 又はディレクトリ ) の所有者 2 2 1 2、当該ファイル ( 又はディレクトリ ) について設定されているアクセス権 2 2 1 3、当該ファイル ( 又はディレクトリ ) のファイルサイズ 2 2 1 4、当該ファイル ( 又はディレクトリ ) の最終更新日時 2 2 1 5、当該 i n o d e がディレクトリエントリである場合に設定される当該ディレクトリの親ディレクトリ 2 2 1 6、当該 i n o d e がディレクトリエントリである場合に設定される当該ディレクトリの子ディレクトリ 2 2 1 7、当該ファイルのデータの实体が格納されるデータブロックを特定する情報 ( 以下、ブロックアドレス 2 2 1 8 と称する。 ) などの情報が含まれる。

40

【 0 0 8 9 】

図 2 3 に示すように、本実施形態のファイルシステム 3 1 2 は、図 2 2 に示した通常の一般的なファイルシステムにおける i n o d e 管理テーブル 1 9 1 2 の内容に加えて、更にスタブ化フラグ 2 3 1 1、メタデータ同期要フラグ 2 3 1 2、実体同期要フラグ 2 3 1 3、レプリケーションフラグ 2 3 1 4、リンク先 2 3 1 5、及び重要度 2 3 1 6 を管理している。

【 0 0 9 0 】

尚、レプリケーションによる管理方式やスタブ化による管理方式により、第 1 ストレー

50

ジ装置 10 a に記憶されているファイルのメタデータ ( 図 23 に示した各種のフラグを含む。 ) の複製が第 2 ストレージ装置 10 b にも記憶されている場合 ( レプリケーションされている場合 ) には、後述する同期処理 S 2900 によって一方の装置のメタデータが更新されるとその旨が他方の装置にも通知される。これにより第 1 ストレージ装置 10 a のメタデータと第 2 ストレージ装置 10 b のメタデータの内容の同一性がほぼリアルタイムに確保される。

【 0091 】

同図において、スタブ化フラグ 2311 には、当該 *inode* に対応するファイル ( 又はディレクトリ ) がスタブ化されているか否かを示す情報が設定される。ここでスタブ化とは、第 1 ストレージ装置 10 a から第 2 ストレージ装置 10 b にファイル ( 又はディレ

10

【 0092 】

尚、スタブ ( stub ) とは、その場合に第 1 ストレージ装置 10 a に残留するメタデータのことをいう。当該 *inode* に対応するファイル ( 又はディレクトリ ) がスタブ化されている場合はスタブ化フラグ 2311 に ON が設定され、スタブ化されていない場合はスタブ化フラグ 2311 に OFF が設定される。

【 0093 】

メタデータ同期要フラグ 2312 には、複製元の第 1 ストレージ装置 10 a のファイル ( 又はディレクトリ ) のメタデータと複製先の第 2 ストレージ装置 10 b のファイル ( 又はディレクトリ ) のメタデータとの間で同期をとる必要 ( 内容を一致させる必要 ) が有るか否かを示す情報が設定される。メタデータの同期が必要な場合はメタデータ同期要フラグ 2312 に ON が設定され、同期が不要な場合はメタデータ同期要フラグ 2312 に OFF が設定される。

20

【 0094 】

実体同期要フラグ 2313 には、複製元の第 1 ストレージ装置 10 a のファイルのデータの実体と複製先の第 2 ストレージ装置 10 b のファイルのデータの実体との間で同期をとる必要 ( 内容を一致させる必要 ) が有るか否かを示す情報が設定される。ファイルのデータの実体について同期が必要な場合は実体同期要フラグ 2313 に ON が設定され、同

30

【 0095 】

メタデータ同期要フラグ 2312 及び実体同期要フラグ 2313 は、後述する同期処理 S 2900 において随時参照される。メタデータ同期要フラグ 2312、もしくは実体同期要フラグ 2313 が ON になっている場合は、第 1 ストレージ装置 10 a のメタデータ又は実体と、その複製である第 2 ストレージ装置 10 b のメタデータ又は実体とが自動的に同期される。

【 0096 】

レプリケーションフラグ 2314 には、その *inode* に対応するファイル ( 又はディレクトリ ) が、現在後述するレプリケーション管理方式による管理の対象になっているか否かを示す情報が設定される。当該 *inode* に対応するファイルが、現在、レプリケーション管理方式による管理の対象になっている場合はレプリケーションフラグ 2314 に ON が設定され、レプリケーションによる管理の対象になっていない場合はレプリケーションフラグ 2314 に OFF が設定される。

40

【 0097 】

リンク先 2315 には、その *inode* に対応するファイルが、後述するレプリケーション管理方式で管理されている場合は、そのファイルの複製先を示す情報 ( 例えば、格納先を特定するパス名、RAID グループの識別子、ブロックアドレス、URL ( Uniform Resource Locator )、LU 等 ) が設定される。

【 0098 】

50

重要度 2 3 1 6 には、ファイルの重要度が設定される。重要度 2 3 1 6 の内容は、例えばユーザがクライアント装置 2 において設定する。また重要度 2 3 1 6 は負荷分散等を目的として設定される場合もある。

【 0 0 9 9 】

= 概略的な動作 =

次に、以上の構成からなる情報処理システム 1 の動作について説明する。

【 0 1 0 0 】

図 2 4 は、第 1 サーバ装置 3 a が、1 ストレージ装置 1 0 a に格納されているファイルを対象としたレプリケーションを開始させる旨の要求（以下、レプリケーション開始要求と称する。）を受け付けた場合に、情報処理システム 1 において行われる処理（以下、レプリケーション開始処理 S 2 4 0 0 と称する）を説明する図である。

10

【 0 1 0 1 】

第 1 サーバ装置 3 a は、クライアント装置 2 からレプリケーション開始要求を受け付けると、当該要求に対象として指定されているファイルについてレプリケーションによる管理方式による管理を開始する。尚、第 1 サーバ装置 3 a は、通信ネットワーク 5 を介してクライアント装置 2 からレプリケーション開始要求を受け付ける以外に、例えば、当該第 1 サーバ装置 3 a において内部的に発生するレプリケーション開始要求も受け付ける。

【 0 1 0 2 】

ここでレプリケーションによる管理方式とは、ファイルのデータ（メタデータ及び実体）を、第 1 ストレージ装置 1 0 a 及び第 2 ストレージ装置 1 0 b の双方において管理する方式である。

20

【 0 1 0 3 】

レプリケーションによる管理方式では、第 1 ストレージ装置 1 0 a に格納されているファイルの実体又はメタデータが更新されると、このファイルの複製（又はアーカイブファイル）として管理されている、第 2 ストレージ装置 1 0 b 側のファイルのメタデータ又は実体が、同期又は非同期に更新される。レプリケーションによる管理方式が行われることにより、第 1 ストレージ装置 1 0 a に格納されているファイルのデータ（メタデータ又は実体）と、その複製として第 2 ストレージ装置 1 0 b に格納されているファイルのデータ（メタデータ又は実体）の一致性が、同期又は非同期に確保（保証）される。

【 0 1 0 4 】

30

尚、第 2 ストレージ装置 1 0 b 側のファイル（アーカイブファイル）におけるメタデータはファイルの実体として管理してもよい。そのようにすることで、第 1 サーバ装置 3 a のファイルシステム 3 1 2 と第 2 サーバ装置 3 b のファイルシステム 3 5 2 の仕様が相違する場合でもレプリケーションによる管理方式を実現することができる。

【 0 1 0 5 】

同図に示すように、レプリケーション開始要求を受け付けると（S 2 4 1 1）、第 1 サーバ装置 3 a は、受け付けたレプリケーション開始要求に指定されているファイルのデータ（メタデータ及び実体）を第 1 ストレージ装置 1 0 a から読み出し、読み出したファイルのデータを第 2 サーバ装置 3 b に送信する（S 2 4 1 2）。

【 0 1 0 6 】

40

第 2 サーバ 3 b は、第 1 サーバ装置 3 a から送られてくる、上記ファイルのデータを受信すると、受信したデータを第 2 ストレージ装置 1 0 b に格納する（S 2 4 1 3）。

【 0 1 0 7 】

尚、上記転送に際し、第 1 サーバ装置 3 a のデータ複製 / 移動処理部 3 1 4 は、転送元ファイルのレプリケーションフラグ 2 3 1 4 を ON に設定する（S 2 4 1 4）。

【 0 1 0 8 】

図 2 5 は、第 1 ストレージ装置 1 0 a に格納されている、レプリケーション管理方式により管理されているファイル（レプリケーションフラグ 2 3 1 4 が ON に設定されているファイル。以下、レプリケーションファイルと称する。）を前述したスタブ化の候補として設定する際に、情報処理システム 1 において行われる処理（以下、スタブ化候補選出処

50

理 S 2 5 0 0 と称する。) を説明する図である。以下、同図とともにスタブ化候補選出処理 S 2 5 0 0 について説明する。

【 0 1 0 9 】

第 1 サーバ装置 3 a は、ファイル格納領域の残容量を随時 (リアルタイム、定期的、予め設定されたタイミング等) 監視している。

【 0 1 1 0 】

第 1 サーバ装置 3 a は、ファイルシステム 3 1 2 に対してファイルの格納領域として割り当てられている第 1 ストレージ装置 1 0 a の記憶領域 (以下、ファイル格納領域と称する。) の残容量が、予め設定された閾値 (以下、スタブ化閾値と称する。) 未満になると、所定の選出基準に従い、第 1 ストレージ装置 1 0 a に格納されているレプリケーション  
10  
ファイルの中からスタブ化の候補を選出する (S 2 5 1 1)。尚、上記所定の選出基準としては、例えば、最終更新日時が古い順、アクセス頻度の低い順などがある。

【 0 1 1 1 】

次に第 1 サーバ装置 3 a は、スタブ化の候補を選出すると、選出したレプリケーションファイルのスタブ化フラグ 2 3 1 1 を ON、レプリケーションフラグ 2 3 1 4 を OFF、メタデータ同期要フラグ 2 3 1 2 を ON に夫々設定する (S 2 5 1 2)。尚、第 1 サーバ装置 3 a は、ファイル格納領域の残容量を、例えば、ファイルシステム 3 1 2 が管理している情報から取得する。

【 0 1 1 2 】

図 2 6 は、スタブ化候補選出処理 S 2 5 0 0 によってスタブ化候補として選出されたファイルを実際にスタブ化する際、情報処理システム 1 において行われる処理 (以下、スタブ化処理 S 2 6 0 0 と称する。) を説明する図である。スタブ化処理 S 2 6 0 0 は、例えば、予め設定されたタイミング (例えばスタブ化候補選出処理 S 2 5 0 0 が行われたのに  
20  
続いて) で行われる。以下、同図とともにスタブ化処理 S 2 6 0 0 について説明する。

【 0 1 1 3 】

同図に示すように、第 1 サーバ装置 3 a は、第 1 ストレージ装置 1 0 a のファイル格納領域に格納されているファイルの中から、スタブ化候補として選出されているファイル (スタブ化フラグ 2 3 1 1 が ON に設定されているファイル) を一つ以上抽出する (S 2 6  
30  
1 1)。

【 0 1 1 4 】

そして第 1 サーバ装置 3 a は、抽出したファイルの実体を第 1 ストレージ装置 1 0 a から削除するとともに、抽出したファイルのメタデータからそのファイルの第 1 ストレージ装置 1 0 a の格納先を示す情報に無効な値を設定し (例えば、メタデータの当該ファイルの格納先を設定する欄 (例えばブロックアドレス 2 2 1 8 を設定する欄) に NULL 値や  
40  
ゼロを設定する。)、スタブ化候補として選出されているファイルを実際にスタブ化する。またこのとき第 1 サーバ装置 3 a は、メタデータ同期要フラグ 2 3 1 2 を ON に設定する (S 2 6 1 2)。

【 0 1 1 5 】

図 2 7 は、第 1 サーバ装置 3 a が、クライアント装置 2 から第 1 ストレージ装置 1 0 a のファイル格納領域に格納されているレプリケーションファイルに対する更新要求を受け  
40  
付けた場合に、情報処理システム 1 において行われる処理 (以下、レプリケーションファイル更新処理 S 2 7 0 0 と称する。) を説明する図である。以下、同図とともにレプリケーションファイル更新処理 S 2 7 0 0 について説明する。

【 0 1 1 6 】

第 1 サーバ装置 3 a は、レプリケーションファイルに対する更新要求を受け付けると (S 2 7 1 1)、その第 1 ストレージ装置 1 0 a に格納されているそのレプリケーションファイルのデータ (メタデータ、実体) を、受け付けた更新要求に従って更新する (S 2 7  
50  
1 2)。

【 0 1 1 7 】

そして第 1 サーバ装置 3 a は、メタデータを更新した場合はそのレプリケーションファ

イルのメタデータ同期要フラグ2312をONに設定し、レプリケーションファイルの実体を更新した場合はそのレプリケーションファイルの実体同期要フラグ2313をONに設定する(S2713)。

【0118】

図28は、第1サーバ装置3aのファイルシステム312が、クライアント装置2から、第1ストレージ装置10aのファイル格納領域に格納されているレプリケーションファイルに対する参照要求を受け付けた場合に情報処理システム1において行われる処理(以下、レプリケーションファイル参照処理S2800と称する。)を説明する図である。以下、同図とともにレプリケーションファイル参照処理S2800について説明する。

【0119】

第1サーバ装置3aのファイルシステム312は、レプリケーションファイルに対する更新要求を受け付けると(S2811)、そのレプリケーションファイルのデータ(メタデータ又は実体)を第1ストレージ装置10aから読み出し(S2812)、読み出したデータに基づきクライアント装置2に対して応答する情報を生成し、生成した応答情報をクライアント装置2に送信する(S2813)。

【0120】

図29は、第1サーバ装置3aが、クライアント装置2から、第1ストレージ装置10aに格納されているレプリケーションファイルとその第2ストレージ装置10b側のファイルの内容とを一致させる要求(以下、同期要求と称する。)を受け付けた際に情報処理システム1において行われる処理(以下、同期処理S2900と称する。)を説明する図である。以下、同図とともに同期処理S2900について説明する。

【0121】

尚、同期処理S2900は、クライアント装置2からの同期要求を受け付けた場合以外の事象を契機として開始するようにしてもよい。例えば、予め設定されたタイミング(リアルタイム、定期的等)が到来したことを契機として、第1サーバ装置3aが自発的に同期処理S2900を開始するようにしてもよい。

【0122】

第1サーバ装置3aは、クライアント装置2からレプリケーションファイルの同期要求を受け付けると(S2911)、第1ストレージ装置10aのファイル格納領域に格納されているレプリケーションファイルのうち、メタデータ同期要フラグ2312又は実体同期要フラグ2313の少なくともいずれかがONに設定されているファイルを取得する(S2912)。

【0123】

そして第1サーバ装置3aは、取得したファイルのメタデータ又は実体を第2サーバ装置3bに送信するとともに、当該レプリケーションファイルのメタデータ同期要フラグ2312又は実体同期要フラグ2313をOFFに設定する(S2913)。

【0124】

第2サーバ装置3bは、メタデータ又は実体を受信すると(S2913)、受信したメタデータ又は実体に対応する、第2ストレージ装置10bに格納されているファイルのメタデータ又は実体を、受信したメタデータ又は実体に基づき更新する(S2914)。  
尚、第1サーバ装置3aから第2サーバ装置3bにメタデータ又は実体の全体を送信するのではなく、前回同期時からの更新差分のみを送信するようにしてもよい。

【0125】

以上に説明した同期処理S2900が行われることにより、第1ストレージ装置10aに格納されているファイルのデータ(メタデータ及び実体)と、当該ファイルに対応する第2ストレージ装置10bに格納されているファイルのデータ(メタデータ及び実体)とが同期される。

【0126】

図30は、第1サーバ装置3aのファイルシステム312が、クライアント装置2等から、スタブ化されているファイル(スタブ化フラグ2311がONに設定されているファ

10

20

30

40

50

イル)のメタデータに対するアクセス要求(参照要求又は更新要求)を受け付けた場合に情報処理システム1において行われる処理(以下、メタデータアクセス処理S3000と称する。)を説明する図である。以下、同図とともにメタデータアクセス処理S3000について説明する。

【0127】

同図に示すように、第1サーバ装置3aは、スタブ化されているファイルのメタデータに対するアクセス要求を受け付けると(S3011)、アクセス要求の対象になっている第1ストレージ装置10aのメタデータを取得し、アクセス要求の内容に従い参照(読み出したメタデータに基づく応答情報のクライアント装置2への送信)、又はメタデータの更新を行う(S3012)。またメタデータの内容を更新した場合には、そのファイルのメタデータ同期要フラグ2312にONを設定する(S3013)。

10

【0128】

このように、第1サーバ装置3aは、スタブ化されているファイルに対するアクセス要求が発生し、そのアクセス要求がそのファイルのメタデータのみを対象としている場合には、第1ストレージ装置10aに格納されているメタデータを用いてアクセス要求を処理する。このため、アクセス要求がそのファイルのメタデータのみを対象としている場合にはクライアント装置2に迅速に応答を返すことができる。

【0129】

図31は、第1サーバ装置3aが、クライアント装置2から、スタブ化されているファイル(スタブ化フラグ2311がONに設定されているファイル。以下、スタブ化ファイルと称する。)の実体に対する参照要求を受け付けた場合に情報処理システム1において行われる処理(以下、スタブ化ファイル実体参照処理S3100と称する。)を説明する図である。以下、同図とともにスタブ化ファイル実体参照処理S3100について説明する。

20

【0130】

第1サーバ装置3aは、クライアント装置2からスタブ化ファイルの実体に対する参照要求を受け付けると(S3111)、取得したメタデータを参照して当該スタブ化ファイルの実体が第1ストレージ装置10aに格納されているか否かを判断する(S3112)。ここでこの判断は、例えば、取得したメタデータにスタブ化ファイルの実体の格納先を示す情報(例えばブロックアドレス2218)に有効な値が設定されているか否かに基づいて行う。

30

【0131】

上記判断の結果、スタブ化ファイルの実体が第1ストレージ装置10aに格納されている場合には、第1サーバ装置3aは、第1ストレージ装置10aから当該スタブ化ファイルの実体を読み出し、読み出した実体に基づきクライアント装置2に応答する情報を生成し、生成した応答情報をクライアント装置2に送信する(S3113)。

【0132】

一方、上記判断の結果、スタブ化ファイルの実体が第1ストレージ装置10aに格納されていない場合には、第1サーバ装置3aは、第2サーバ装置3bに対してスタブ化ファイルの実体の提供を要求する(以下、リコール要求と称する。)(S3114)。尚、実体の取得要求は、必ずしも一度の取得要求によって実体の全体を取得する要求でなくてもよく、例えば実体の一部のみを複数回要求するようにしてもよい。

40

【0133】

第1サーバ装置3aは、上記取得要求に応じて第2サーバ装置3bから送られてくるスタブ化ファイルの実体を受信すると(S3115)、受信した実体に基づき応答情報を生成し、生成した応答情報をクライアント装置2に送信する(S3116)。

【0134】

また第1サーバ装置3aは、上記第2サーバ装置3bから受信した実体を第1ストレージ装置10aに格納し、当該スタブ化ファイルのメタデータの当該ファイルの実体の格納先を示す情報(例えばブロックアドレス2218)に、当該ファイルの第1ストレージ装

50

置 1 0 a における格納先を示す内容を設定する。また第 1 サーバ装置 3 a は、当該ファイルのスタブ化フラグ 2 3 1 1 を OFF に、レプリケーションフラグ 2 3 1 4 を ON に、メタデータ同期要フラグ 2 3 1 2 を ON に、夫々設定する（当該ファイルをスタブ化ファイルからレプリケーションファイルに変更する。）（S 3 1 1 7）。

【 0 1 3 5 】

尚、メタデータ同期要フラグ 2 3 1 2 を ON に設定するのは、第 1 ストレージ装置 1 0 a と第 2 ストレージ装置 1 0 b との間で、当該スタブ化ファイルのスタブ化フラグ 2 3 1 1 及びレプリケーションフラグ 2 3 1 4 の内容を事後的に自動的に同期させるためである。

【 0 1 3 6 】

図 3 2 は、第 1 サーバ装置 3 a が、クライアント装置 2 からスタブ化ファイルの実体に対する更新要求を受け付けた場合に、情報処理システム 1 において行われる処理（以下、スタブ化ファイル実体更新処理 S 3 2 0 0 と称する。）を説明する図である。以下、同図とともにスタブ化ファイル実体更新処理 S 3 2 0 0 について説明する。

【 0 1 3 7 】

第 1 サーバ装置 3 a は、スタブ化ファイルの実体に対する更新要求を受け付けると（S 3 2 1 1）、更新要求の対象になっているスタブ化ファイルのメタデータを取得し、取得したメタデータに基づき当該スタブ化ファイルの実体が第 1 ストレージ装置 1 0 a に格納されているか否かを判断する（S 3 2 1 2）。尚、判断方法はスタブ化ファイル実体参照処理 S 3 1 0 0 の場合と同様である。

【 0 1 3 8 】

上記判断の結果、スタブ化ファイルの実体が第 1 ストレージ装置 1 0 a に格納されていた場合、第 1 サーバ装置 3 a は、第 1 ストレージ装置 1 0 a に格納されている当該スタブ化ファイルの実体を更新要求の内容に従って更新するとともに、当該スタブ化ファイルの実体同期要フラグ 2 3 1 3 を ON に設定する（S 3 2 1 3）。

【 0 1 3 9 】

一方、スタブ化ファイルの実体が第 1 ストレージ装置 1 0 a に格納されていない場合、第 1 サーバ装置 3 a は、第 2 サーバ装置 3 b に当該スタブ化ファイルの実体の取得要求（リコール要求）を送信する（S 3 2 1 4）。

【 0 1 4 0 】

第 1 サーバ装置 3 a は、上記要求に応じて第 2 サーバ装置 3 b から送られてきたファイルの実体を受信すると（S 3 2 1 5）、受信した実体の内容を更新要求の内容に従って更新し、更新後の実体を当該スタブ化ファイルの実体として第 1 ストレージ装置 1 0 a に格納する。また第 1 サーバ装置 3 a は、当該スタブ化ファイルのスタブ化フラグ 2 3 1 1 を OFF、レプリケーションフラグ 2 3 1 4 を OFF、メタデータ同期要フラグ 2 3 1 2 を ON に夫々設定する（S 3 2 1 6）。

【 0 1 4 1 】

< 障害回復時の処理 >

次に第 1 サーバ装置 3 a において何らかの障害が発生して情報処理システム 1 の機能が停止し、その後、第 1 サーバ装置 3 a が修復されて情報処理システム 1 の機能を再開させる場合に当該情報処理システム 1 において行われる処理について説明する。

【 0 1 4 2 】

図 3 3 は、修復した第 1 サーバ装置 3 a に仮想マシン 3 1 0 を復旧させる際に情報処理システム 1 において行われる処理（以下、仮想マシン復旧処理 S 3 3 0 0 と称する。）を説明する図である。以下、同図とともに仮想マシン復旧処理 S 3 3 0 0 について説明する。

【 0 1 4 3 】

尚、仮想マシン復旧処理 S 3 3 0 0 を実行する前提として、仮想マシン 3 1 0 を復旧させるための仮想マシンイメージ（仮想化制御部 3 0 5 に仮想マシン 3 1 0 を実現させるために必要な構成情報であり、例えば、CPU やメモリ等のハードウェア構成、記憶領域の

10

20

30

40

50

サイズ、ネットワーク仕様等の情報が含まれる。)が、第2ストレージ装置10bに予め格納されているものとする。

【0144】

同図に示すように、まず第1サーバ装置3aにおいて、ブートロード等を用いて記録媒体3310等に記録されているインストールプログラムを実行し、第1サーバ装置3aに仮想化制御部305をインストールし(S3311)、仮想化制御部305の機能を開始させる(S3312)。

【0145】

次に、機能を開始した仮想化制御部305が、第2サーバ装置3bに仮想マシンイメージの提供を要求する(S3313)。

10

【0146】

第2サーバ装置3bは、第1サーバ装置3aから上記要求を受信すると、上記要求に指定されている仮想マシンイメージを第2ストレージ装置10bから取得し(S3314)、取得した仮想マシンイメージを第1サーバ装置3aに送信する(S3315)。

【0147】

尚、第2サーバ装置3bは、仮想マシンイメージを、例えば、第1サーバ装置3aの識別子(以下、サーバIDと称する。)と第1サーバ装置3aにおいて実現される仮想マシン310の識別子(以下、仮想マシンIDと称する。)とに対応付けて管理しており、上記取得要求を受信すると、その取得要求に指定されているサーバID及び仮想マシンIDで特定される仮想マシンイメージを特定し、特定した仮想マシンイメージを第1ストレージ装置10aに送信する。

20

【0148】

第1サーバ装置3aは、第2サーバ装置3bから仮想マシンイメージを受信すると(S3316)、受信した仮想マシンイメージを第1ストレージ装置10aに格納し、受信した仮想マシンイメージに基づく仮想マシン310の動作を開始させる(S3317)。

【0149】

尚、以上に説明した仮想マシン復旧処理S3300は、原則として仮想マシンイメージに基づく仮想マシン310の再起動が必要になるような大きな障害が発生した場合に行われ、例えば、障害が仮想マシン310の再起動が必要になるような障害でない場合は必ずしも仮想マシン310を再起動させる必要はない。

30

【0150】

図34は、図33に示した仮想マシン復旧処理S3300により第1サーバ装置3aにおいて仮想マシン310が動作を開始した後、クライアント装置2からのデータI/O要求を受け付ける前に、情報処理システム1において行われる、ディレクトリイメージを復旧させる処理(以下、ディレクトリイメージ事前回復処理S3400と称する。)を説明する図である。以下、同図とともにディレクトリイメージ事前回復処理S3400について説明する。

【0151】

まず第1サーバ装置3aが、仮想マシン復旧処理S3300によって再起動された仮想マシン310のファイルシステム312が、障害が発生する前に第1ストレージ装置10aに構成していたディレクトリの構成(つまり第2ストレージ装置10bに記憶されているディレクトリの構成であり、ディレクトリの階層構造を示すデータ、ディレクトリのデータ(メタデータ)、及びファイルのデータ(メタデータ及び実体)を含む。以下、ディレクトリイメージと称する。)における、最上位ディレクトリ(以下、ルートディレクトリと称する。)に存在するディレクトリのメタデータ、及びルートディレクトリに存在するファイルのメタデータの取得要求を、第2サーバ装置3bに送信する(S3411)。

40

【0152】

尚、本実施形態において、ルートディレクトリに存在するディレクトリのメタデータ、及びルートディレクトリに存在するファイルのメタデータという場合には、ルートディレクトリに存在する(から見える)ディレクトリ及びファイルは含まれるが、ルートディレ

50

クトリに存在するディレクトリの更に配下に存在するディレクトリやそのディレクトリに存在するファイルは含まれない。

【0153】

第2サーバ装置3bは、上記取得要求を受信すると、要求されているルートディレクトリに存在するディレクトリのメタデータ、及びルートディレクトリに存在するファイルのメタデータを、第2ストレージ装置10bから取得し(S3412)、取得したメタデータを第1ストレージ装置10aに送信する(S3413)。

【0154】

尚、第2サーバ装置3bは、前述したレプリケーションによる管理方式の運用において、メタデータをサーバIDと仮想マシンIDとに対応づけて管理し、第2サーバ装置3bは、上記取得要求を受信すると、その取得要求に指定されているサーバID及び仮想マシンIDで特定されるメタデータを特定し、特定したメタデータを第2ストレージ装置10bから取得する。

10

【0155】

第1サーバ装置3aは、第2サーバ装置3bからメタデータを受信すると(S3413)、受信したメタデータに基づくディレクトリイメージを、第1ストレージ装置10aに復元する(S3414)。またこのとき、第1サーバ装置3aは、メタデータ同期要フラグ2312をON、実体同期要フラグ2313をONに夫々設定する。尚、復元されたファイルはいずれもメタデータのみに基づくものであるため、これらのファイルはいずれもスタブ化された状態になっており、スタブ化フラグ2311がONに設定されている。

20

【0156】

そして第1サーバ装置3aは、第1ストレージ装置10aにディレクトリイメージが復元されると、クライアント装置2へのサービスを開始する。

【0157】

図35は、図34に示したディレクトリイメージ事前回復処理S3400の後、クライアント装置2からのデータI/O要求の受け付けを開始した第1サーバ装置3aが、障害発生前に当該第1サーバ装置3aが管理していたディレクトリイメージを復元していく処理(以下、オンデマンド復元処理S3500と称する。)を説明する図である。以下、同図とともにオンデマンド復元処理S3500について説明する。

【0158】

30

第1サーバ装置3aは、サービスを開始した後、クライアント装置2からあるファイルについてデータI/O要求を受け付けると(S3511)、受け付けたデータI/O要求の対象になっているファイル(以下、アクセス対象ファイルと称する。)のメタデータが、第1ストレージ装置10aに存在するか否か(サービス開始後、既に第1ストレージ装置10aにメタデータが復元されているか否か)を調べる(S3512)。

【0159】

メタデータが第1ストレージ装置10aに復元されている場合、第1サーバ装置3aは、受け付けたデータI/O要求の対象(メタデータ又は実体)、データI/O要求の種類(参照要求か更新要求か)、レプリケーションによる管理方式で管理されているか否か(レプリケーションフラグ2314がONか否か)、スタブ化されているか否か(スタブ化フラグがONか否か)に応じて、受け付けたデータI/O要求に対応する処理(前述したレプリケーションファイル更新処理S2700、レプリケーションファイル参照処理S2800、メタデータアクセス処理S3000、スタブ化ファイル実体参照処理S3100、スタブ化ファイル実体更新処理S3200)を行い、クライアント装置2に応答を返す(S3518)。

40

【0160】

一方、アクセス対象ファイルのメタデータが復元されていない場合には、第1サーバ装置3aは、ルートディレクトリを起点としてアクセス対象ファイルが存在するディレクトリレベル(ディレクトリ階層)に至るまでのディレクトリイメージを復元するためのデータを第2サーバ装置3b(第2ストレージ装置10b)から取得し(S3513~S35

50

15)、取得したデータを用いてルートディレクトリから上記ディレクトリレベルに至るまでのディレクトリイメージを第1ストレージ装置10aに復元する(S3516)。

【0161】

また第1サーバ装置3aは、アクセス対象ファイルのスタブ化フラグ2311をONに、レプリケーションフラグ2314をOFFに、メタデータ同期要フラグ2312をONに、夫々設定する(S3517)。

【0162】

次に第1サーバ装置3aは、受け付けたデータI/O要求の対象や種類、管理方式、スタブ化の有無等に応じて、受け付けたデータI/O要求に対応する処理を行い、クライアント装置2に応答を返す(S3518)。

10

【0163】

図36に、データI/O要求が繰り返し発生することにより、以上に説明したオンデマンド復元処理S3500によって段階的に第1ストレージ装置10aにディレクトリイメージが復元されていく様子を示している。

【0164】

同図中、強調された文字列(アンダーラインが付与されている文字列)で示すディレクトリは、そのディレクトリのメタデータは復元されているがその配下のディレクトリのメタデータは未だ復元されていない。また強調されていない文字で示すディレクトリは、そのディレクトリの配下のディレクトリのメタデータも既に復元されている。また強調された文字で示すファイルは、そのファイルのメタデータは復元されているが、その実体は未だ復元されていない。また強調されていない文字で示すファイルは、そのファイルの実体が既に復元されている。

20

【0165】

図36の図(0)は、障害が発生する直前に第1サーバ装置3a(第1ストレージ装置10a)において管理されていたディレクトリイメージ(最終的に復元されるディレクトリイメージの全体)である。

【0166】

図36の図(A)は、ディレクトリイメージ事前回復処理S3400による回復直後(まだ第1サーバ装置3aがデータI/O要求を受け付けていない状態)におけるディレクトリイメージである。この段階ではルートディレクトリ「/」の直下に存在するディレクトリ「/dir1」及び「/dir2」のメタデータは復元されているが、その更に配下のディレクトリのメタデータについては未だ復元されていない。またルートディレクトリ「/」の直下に存在するファイル「a.txt」のメタデータは復元されているが、実体については未だ復元されていない。

30

【0167】

図36の図(B)は、図(A)の状態クライアント装置2からディレクトリ「/dir1」の配下に存在するファイル「c.txt」に対するデータI/O要求を受け付けた後の状態である。クライアント装置2からファイル「c.txt」に対するデータI/O要求を受け付けたため、ディレクトリ「/dir1」のメタデータと「/c.txt」のメタデータが復元されている。

40

【0168】

図36の図(C)は、図(B)の状態更にクライアント装置2からディレクトリ「/dir2」の配下に存在するファイル「b.txt」に対するデータI/O要求を受け付けた後の状態である。同図に示すように、クライアント装置2からファイル「b.txt」に対するデータI/O要求を受け付けたため、「/b.txt」のメタデータが復元されている。尚、「/dir2」の配下に存在する「/b.txt」のメタデータが復元されたため、「/dir2」は強調されていない文字で表記している。

【0169】

図36の図(D)は、図(C)の状態クライアント装置2からファイル「b.txt」に対するデータI/O要求(更新要求)を受け付けた後の状態である。クライアント装

50

置 2 からファイル「b . t x t」に対するデータ I / O 要求（更新要求）を受け付けたため、ファイル「b . t x t」の実体が復元されている。

【 0 1 7 0 】

以上に説明したように、本実施形態の情報処理システム 1 においては、第 1 サーバ装置 3 a に障害が発生した後、データ I / O 要求の受け付けを開始した時点では、ディレクトリイメージ事前回復処理 S 3 4 0 0 によってルートディレクトリに存在するディレクトリのメタデータ、及びルートディレクトリに存在するファイルのメタデータのみが復元される。そしてその後は第 1 サーバ装置 3 a に対してクライアント装置 2 からまだ復元されていないファイルに対してデータ I / O 要求が発生する度に、段階的に第 1 サーバ装置 3 a（第 1 ストレージ装置 1 0 a）にディレクトリイメージが復元されていく。

10

【 0 1 7 1 】

このように障害回復後、データ I / O 要求の受け付けを開始する前にディレクトリイメージの全体を復元してしまうのではなく、ディレクトリイメージを段階的に復元するようにすることで、サービスが再開される前にディレクトリイメージの全体を復元しておくようにしておく場合に比べ、障害発生からサービス再開までに要する時間を短縮することができ、ユーザの業務等への影響を防ぐことができる。

【 0 1 7 2 】

またディレクトリイメージが完全に復元されるまでの間は、第 1 ストレージ装置 1 0 a の資源を節約することができる。さらにディレクトリイメージが完全に復元されるまでの間は記憶容量の消費が抑えられるため、例えば、記憶容量の小さなストレージ装置を障害が発生した第 1 ストレージ装置 1 0 a の代替装置として用いるようなこともできる。

20

【 0 1 7 3 】

< 復元対象の追加 >

ところで、例えば、第 1 サーバ装置 3 a や第 1 ストレージ装置 1 0 a が十分な性能や記憶容量を備えている場合や、ユーザがサービスの迅速に完全復旧されることを望んでいるような場合には、図 3 5 に示したオンデマンド復元処理 S 3 5 0 0 によって障害発生前における第 1 ストレージ装置 1 0 a のディレクトリイメージを迅速に復元することが好ましい。

【 0 1 7 4 】

しかし前述したオンデマンド復元処理 S 3 5 0 0 によるディレクトリイメージの復元速度は、クライアント装置 2 からのデータ I / O 要求の発生頻度に依存しているため、データ I / O 要求の発生頻度が少ない場合はディレクトリイメージが完全に復元されるまでに長時間を要してしまうことになる。

30

【 0 1 7 5 】

そこで本実施形態の情報処理システム 1 は、このようなディレクトリイメージの復旧速度の低下を防ぐべく、オンデマンド復元処理 S 3 5 0 0 において、第 2 サーバ装置 3 b が、第 1 サーバ装置 3 a から復元のためのディレクトリイメージを要求された際、クライアント装置 2 から受け付けたデータ I / O 要求が所定の条件を満たすことを条件として、第 1 サーバ装置 3 a に送信するディレクトリイメージを追加し、ディレクトリイメージの復元を自動的に促進させる仕組みを備えている。

40

【 0 1 7 6 】

尚、上記の所定の条件としては、例えば次のようなものが考えられる。

【 0 1 7 7 】

（条件 1）アクセス対象ファイルのデータサイズが、現在から遡って所定時間内に発生したデータ I / O 要求のアクセス対象ファイルのデータサイズの平均値よりも小さい。

【 0 1 7 8 】

（条件 2）アクセス対象ファイルのデータサイズが予め設定された閾値よりも小さい。

【 0 1 7 9 】

また上記の仕組みにおいて追加するディレクトリイメージの選出方法としては、例えば次のようなものが考えられる。

50

## 【 0 1 8 0 】

( 選出方法 1 ) 既に復元されているディレクトリの配下に存在するファイルのメタデータ及び / 又は実体を選出する。

## 【 0 1 8 1 】

ここで一般に既に復元されているディレクトリの配下に存在するファイルはその後アクセスされる可能性が高く、選出方法 1 に従いそのようなディレクトリの配下に存在するファイルのディレクトリイメージを復元することで、クライアント装置 2 に対する応答性能の向上が期待できる。

## 【 0 1 8 2 】

( 選出方法 2 ) 既に復元されているディレクトリの配下に存在するディレクトリのメタデータを選出する。

10

## 【 0 1 8 3 】

ここで既に復元されているディレクトリの配下に存在するディレクトリはその後アクセスされる可能性が高いので、選出方法 2 に従い既に復元されているディレクトリの配下に存在するディレクトリのメタデータを復元しておくことでクライアント装置 2 に対する応答性能の向上が期待できる。

## 【 0 1 8 4 】

( 選出方法 3 ) 第 1 サーバ装置 3 a に障害が発生する前に第 1 ストレージ装置 1 0 a に実体が格納されていたファイルの実体 ( 障害が発生する前にスタブ化フラグが O F F になっていたファイル ) を選出する。

20

## 【 0 1 8 5 】

第 1 サーバ装置 3 a に障害が発生する前に第 1 ストレージ装置 1 0 a に実体が格納されていたファイルは元々アクセス頻度が高いファイルであった可能性が高い。そこでそのようなファイルの実体を優先して第 1 ストレージ装置 1 0 a に復元しておくようにすればクライアント装置 2 に対する応答性能の向上を期待することができる。

## 【 0 1 8 6 】

尚、第 1 サーバ装置 3 a は、障害が発生する前に第 1 ストレージ装置 1 0 a に実体が格納されていたファイルであるか否かを、例えば、第 2 サーバ装置 3 b にそのファイルのスタブ化フラグ 2 3 1 1 を問い合わせることにより把握する ( そのファイルのスタブ化フラグ 2 3 1 1 が O F F であれば障害が発生する前に第 1 ストレージ装置 1 0 a に実体が格納されていたことになる ) 。

30

## 【 0 1 8 7 】

( 選出方法 4 ) アクセス対象ファイルよりも高い重要度が設定されているファイルのメタデータ及び / 又は実体を選出する。

## 【 0 1 8 8 】

一般に高い重要度が設定されているファイルはクライアント装置 2 からアクセスされる可能性が高いファイルであることが多い。そこでそのようなファイルのメタデータ及び / 又は実体を復元しておくようにすることで、クライアント装置 2 に対する応答性能の向上を期待することができる。

40

## 【 0 1 8 9 】

尚、第 1 サーバ装置 3 a は、第 1 サーバ装置 3 a から第 2 サーバ装置 3 b に問い合わせを行うことにより、第 1 ストレージ装置 1 0 a にメタデータが未だ復元されていないファイルの重要度 ( i n o d e 管理テーブル 1 9 1 2 の重要度 2 3 1 6 の内容 ) を取得する。

## 【 0 1 9 0 】

( 選出方法 5 ) 障害発生時点から遡って所定時間内におけるアクセス頻度がアクセス対象ファイルよりも高いファイルを選出する。

## 【 0 1 9 1 】

障害発生時点から遡って所定時間内におけるアクセス頻度が高いファイルは、クライアント装置 2 からアクセスされる可能性が高いファイルであると考えられる。そこでそのよ

50

うなファイルのメタデータ及びノ又は実体を復元しておくようにすることで、クライアント装置2に対する応答性能の向上を期待することができる。

【0192】

尚、第1サーバ装置3aは、第2サーバ装置3bにファイルアクセスログ368の内容を問い合わせることにより、障害発生時点から遡って所定時間内におけるファイルのアクセス頻度を取得する。

【0193】

尚、以上に列挙した方法は選出方法の一例に過ぎず、選出方法は以上に示したものに限定されない。例えば以上に列挙した選出方法の二つ以上を組み合わせることで復元するディレクトリイメージを選出するようにしてもよい。例えば単独の選出方法だと選出される復元対象が多すぎてしまう場合には、複数の選出方法を組み合わせることで復元対象を絞り込むことができる。

10

【0194】

図37は、前述したオンデマンド復元処理S3500において、データI/O要求が前述した所定の条件を満たす場合に前述した所定の選出方法に従い復元するディレクトリイメージを追加するようにした処理（以下、オンデマンド復元処理（復元対象追加有）S3700と称する。）を説明する図である。以下、同図とともにオンデマンド復元処理（復元対象追加有）S3700について説明する。

【0195】

第1サーバ装置3aは、クライアント装置2からデータI/O要求を受け付けると（S3711）、当該データI/O要求のアクセス対象ファイルのメタデータが第1ストレージ装置10aに存在するか否か（既に復元されているか否か）を判断する（S3712）。

20

【0196】

アクセス対象ファイルのメタデータが既に復元されている場合には、第1サーバ装置3aは、受け付けたデータI/O要求の対象や種類、管理方式、スタブ化の有無等に応じて、受け付けたデータI/O要求に対応する処理を行い、クライアント装置2に応答を返す（S3718）。

【0197】

一方、アクセス対象ファイルのメタデータが復元されていない場合には、第1サーバ装置3aは、ルートディレクトリを起点としてアクセス対象ファイルが存在するディレクトリレベル（ディレクトリ階層）に至るまでのディレクトリイメージを復元するためのデータを第2サーバ装置3bに要求する（ここまでの処理は図35に示したオンデマンド復元処理S3500と同様である。）。

30

【0198】

第2サーバ装置3bは、上記要求を受信すると、データI/O要求が前述した所定の条件を満たすか否かを判断し、所定の条件を満たす場合には、前述した所定の選出方法に従い追加するディレクトリイメージをさらに選出する。そしてサーバ装置3bは、上記要求に指定されているディレクトリイメージを復元するためのデータと、上記選出したディレクトリイメージを復元するためのデータとを、第2ストレージ装置10bから取得し、第1サーバ装置3aに送信する（S3713～S3715）。

40

【0199】

第1サーバ装置3aは、第2サーバ装置3bから上記データを受信すると、受信したデータを用いて、第1ストレージ装置10aにディレクトリイメージを復元する（S3716）。

【0200】

次に第1サーバ装置3aは、アクセス対象ファイルのスタブ化フラグ2311をONに、レプリケーションフラグ2314をOFFに、メタデータ同期要フラグ2312をONに、夫々設定する（S3717）。

【0201】

50

そして第1サーバ装置3aは、受け付けたデータI/O要求の対象や種類、管理方式、スタブ化の有無等に応じて、受け付けたデータI/O要求に対応する処理を行い、クライアント装置2に応答を返す(S3718)。

【0202】

以上に説明したオンデマンド復元処理(復元対象追加有)S3700によれば、データI/O要求が所定の条件を満たす場合は、復元されるディレクトリイメージが自動的に追加される。このため、ディレクトリイメージの復旧速度を自動的に速めることができ、第1ストレージ装置10aのディレクトリイメージを迅速に障害発生前の状態に復元することができる。

【0203】

また以上に説明したオンデマンド復元処理(復元対象追加有)S3700では、復元するディレクトリイメージを追加するか否かの判断やディレクトリイメージを復元するためのデータの取得に関する処理を、もっぱら第2サーバ装置3b側で行うようにしている。そのため、第1サーバ装置3a側に特別な仕組みを設ける必要がなく、また第1サーバ装置3aの代替装置の選出に際し、機種や製造元(ベンダー)を一致させる必要がない等、情報処理システム1の柔軟な運用形態が可能になる。

【0204】

<再スタブ化の回避>

前述したスタブ化候補選出処理S2500(図25)では、ファイル格納領域の残容量がスタブ化閾値未満になったことを条件としてスタブ化の候補となるファイルが選出され、選出されたファイルは前述したスタブ化処理S2600(図26)において実際にスタブ化(第1ストレージ装置10aから実体が削除)されるが、スタブ化候補選出処理S2500及びスタブ化処理S2600は、図35に説明したオンデマンド復元処理S3500(又は図37に示したオンデマンド復元処理(復元対象追加有)S3700。以下、オンデマンド復元処理S3500のみ表記する。)の実行中においても実行されてしまうことがある。

【0205】

そして例えばスタブ化閾値が高めに設定されていた場合や、障害が発生した第1ストレージ装置10aの代替として用意されたストレージ装置が十分な記憶容量を備えていない場合には、オンデマンド復元処理S3500によって第1ストレージ装置10aに実体が復元されたファイルが、オンデマンド復元処理S3500(又は図37に示したオンデマンド復元処理(復元対象追加有)S3700)によって直ぐに再びスタブ化の候補として選出されスタブ化されてしまう(以下、この現象のことを再スタブ化と称する。)

【0206】

そしてこのような再スタブ化が頻発すると、情報処理システム1の資源が浪費され、情報処理システム1の運用効率が低下してしまうことになる。

【0207】

そこで本実施形態の情報処理システム1は、再スタブ化の発生を抑制すべく、再スタブ化の発生を随時監視し、再スタブ化の発生状況に応じてディレクトリイメージの復元を自動的に抑制する仕組みを備えている。

【0208】

図38は、上記仕組みに関して第2サーバ装置3bによって行われる処理(以下、再スタブ化回避処理S3800と称する。)を説明する図である。以下、同図とともに再スタブ化回避処理S3800について説明する。

【0209】

第2サーバ装置3bは、前述したオンデマンド復元処理S3500の実行中、単位時間当たりの再スタブ化の発生頻度が予め設定された閾値(以下、再スタブ化頻度閾値と称する。)以上になっているか否か、もしくは再スタブ化の発生時間間隔が予め設定された閾値(以下、再スタブ化発生時間間隔閾値と称する。)未満になっているか否かを監視する(S3811~S3813)。

10

20

30

40

50

## 【0210】

ここで再スタブ化が発生したか否かの判断は、例えば、リストアログ365の内容と、メタデータの同期処理S2900における第1サーバ装置3aから第2サーバ装置3bへのスタブ化フラグ2311の更新通知(スタブ化フラグ2311をOFFからONにする通知)とに基づいて行う。

## 【0211】

例えば、第2サーバ装置3bは、第1ストレージ装置10aにディレクトリイメージが復元されてから予め設定された所定時間内に、そのディレクトリイメージのデータ(メタデータ又は実体)のスタブ化フラグ2311がONにされたことをもって再スタブ化が発生したと判断する。

10

## 【0212】

同図に示すように、第2サーバ装置3bは、上記監視において再スタブ化の発生頻度が再スタブ化頻度閾値以上になっていること、もしくは再スタブ化の発生時間間隔が再スタブ化発生時間間隔閾値未満になっていることを検知すると、第1サーバ装置3aに送信するディレクトリイメージ(図37に示したオンデマンド復元処理(復元対象追加有)S3700において追加されるディレクトリイメージを含む)の量を抑制(減少)させる。尚、この抑制には、第1サーバ装置3aへのディレクトリイメージの送信を中止する場合も含まれる(S3814)。

## 【0213】

ここで上記抑制の具体的な方法としては、例えば、次のようなものが考えられる。

20

## 【0214】

(抑制方法1) データI/O要求がメタデータのみを対象としている場合に当該ファイルの実体を復元しないようにする。

## 【0215】

これによれば実体を復元するのに要する負荷を軽減することができる。またデータI/O要求がメタデータのみを対象としている場合には、当該ファイルの実体まで復元する必要がないため、実体を復元しなくてもデータI/O要求の処理に影響を与えることもない。

## 【0216】

(抑制方法2) 前述した(選出方法1)~(選出方法5)の一つ以上を用いたディレクトリイメージの選出を行っている場合に、更に別の選出方法を重複して適用するようにする。

30

## 【0217】

選出方法を重複して適用するようにすることで、段階的に再スタブ化の発生を抑制することができ、再スタブ化の発生状況に併せて、第1サーバ装置3aに送信するディレクトリイメージの量を適切に抑制することができる。

## 【0218】

(抑制方法3) 前述した(選出方法4)において判断に用いている重要度の閾値を更に高く設定する。

## 【0219】

重要度の閾値を更に高く設定することで、再スタブ化の抑制を容易に実施することができる。また重要度の閾値を段階的に高く設定していくようにすれば、再スタブ化の発生状況に併せて、第1サーバ装置3aに送信するディレクトリイメージの量を適切に抑制することができる。

40

## 【0220】

(抑制方法4) 前述した(選出方法5)においてアクセス頻度の判断に用いているアクセス頻度の閾値を更に高く設定する。

## 【0221】

アクセス頻度の閾値を更に高く設定することで、再スタブ化の抑制を容易に実施することができる。またアクセス頻度の閾値を段階的に高く設定していくようにすれば、再スタ

50

ブ化の発生状況に併せて、第1サーバ装置3aに送信するディレクトリイメージの量を適切に抑制することができる。

【0222】

尚、第2サーバ装置3bは上記監視を継続して行い、再スタブ化の発生頻度が再スタブ化頻度閾値未満になるとともに、再スタブ化の発生時間間隔が再スタブ化発生時間間隔閾値以上になると、自動的に上記抑制を解除する。ここで解除には、抑制を一度に全て解除してしまう場合のほか、少しずつディレクトリイメージを追加していく等、抑制を段階的に解除していく場合も含む(S3814)。

【0223】

以上に説明したように、再スタブ化回避処理S3800によれば、再スタブ化が頻繁に発生すると第2サーバ装置3bから第1サーバ装置3aに送信されるディレクトリイメージの量が自動的に抑制されるので、再スタブ化の発生を抑制することができる。そのため、再スタブ化による情報処理システム1の資源の浪費を防ぐことができ、再スタブ化に起因する情報処理システム1を運用効率の低下を防ぐことができる。

10

【0224】

また以上に説明した再スタブ化回避処理S3800は、第2サーバ装置3bが主体となっていくため、第1サーバ装置3a側に特別な仕組みを設ける必要がない。そのため、再スタブ化を抑制するための仕組みを情報処理システム1に容易に実現することができる。また特別な性能や仕様が要求されないことで第1ストレージ装置10aの選択肢が拡がり、ベンダーや機種等に依存することなく、ハードウェア及びソフトウェアを自由に選択することができる。

20

【0225】

<処理詳細>

次に情報処理システム1において行われる処理の詳細について説明する。

【0226】

図39は、図24に示したレプリケーション開始処理S2400の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0227】

第1サーバ装置3aは、クライアント装置2等からレプリケーション開始要求を受け付けたか否かをリアルタイムに監視している(S3911)。第1サーバ装置3aは、クライアント装置2等からレプリケーション開始要求を受け付けると(S3911:YES)(図24のS2411)、受け付けたレプリケーション開始要求に指定されているファイルのデータ(メタデータ及び実体)の格納先(RAIDグループの識別子、ブロックアドレス等)を第2サーバ装置3bに問い合わせる(S3912)。

30

【0228】

第2サーバ装置3bは、上記の問い合わせがあると(S3921)、第2ストレージ装置10bの空き領域を探索してファイルのデータの格納先を決定し、決定した格納先を第1サーバ装置3aに通知する(S3922)。

【0229】

第1サーバ装置3aは、上記通知を受信すると(S3913)、受け付けたレプリケーション開始要求に指定されているファイルのデータ(メタデータ及び実体)を、第1ストレージ装置10aから読み出し(S3914)(図24のS2412)、読み出したファイルのデータをS3922で通知された格納先とともに第2サーバ装置3bに送信する(S3915)(図24のS2413)。

40

【0230】

また第1サーバ装置3aは、当該ファイルのメタデータ(第1ストレージ装置10aに格納されている当該ファイルのメタデータ)のレプリケーションフラグ2314をON、メタデータ同期要フラグ2312をONに夫々設定する(S3916)(図24のS2414)。

【0231】

50

尚、メタデータ同期要フラグ2312をONに設定しておくことで、前述した同期処理S2900により第1ストレージ装置10aに格納されているファイルのメタデータと、その複製として第2ストレージ装置10bに格納されているファイルのメタデータの一致性が、同期又は非同期に確保(保証)される。

【0232】

一方、第2サーバ装置3bは、第1サーバ装置3aからファイルのデータを受信すると(S3923)、受信したファイルのデータを、当該ファイルとともに受信した格納先で特定される第2ストレージ装置10bの位置に格納する(S3924)。

【0233】

図40は、図25に示したスタブ化候補選出処理S2500の詳細を説明するフローチャートである。以下、同図とともに説明する。

10

【0234】

第1サーバ装置3aは、ファイル格納領域の残容量がスタブ化閾値未満になっているかを随時監視し(S4011、S4012)、ファイル格納領域の残容量がスタブ化閾値未満になっていることを検知すると、前述した所定の選出基準に従い第1ストレージ装置10aに格納されているレプリケーションファイルの中からスタブ化の候補を選出する(S4012)(図25のS2511)。

【0235】

そして第1サーバ装置3aは、スタブ化の候補を選出すると(S4013)、選出したレプリケーションファイルのスタブ化フラグ2311をON、レプリケーションフラグ2314をOFF、メタデータ同期要フラグ2312をONに夫々設定する(S4014)(図25のS2512)。

20

【0236】

図41は、図26に示したスタブ化処理S2600の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0237】

第1サーバ装置3aは、第1ストレージ装置10aのファイル格納領域に格納されているファイルの中から、スタブ化候補として選出されているファイル(スタブ化フラグ2311がONに設定されているファイル)を随時抽出する(S4111、S4112)。

【0238】

30

そして第1サーバ装置3aは、抽出したファイルの実体を第1ストレージ装置10aから削除するとともに(S4113)、抽出したファイルのメタデータからそのファイルの第1ストレージ装置10aの格納先を示す情報に無効な値を設定し(例えば、メタデータの当該ファイルの格納先を設定する欄(例えばブロックアドレス2218)にNULL値やゼロを設定する。)(S4114)、メタデータ同期要フラグ2312をONに設定する(S4115)(図26のS2611)。

【0239】

図42は、図27に示したレプリケーションファイル更新処理S2700の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0240】

40

第1サーバ装置3aは、クライアント装置2からレプリケーションファイルに対する更新要求を受け付けたか否かをリアルタイムに監視している(S4211)。第1サーバ装置3aは、更新要求を受け付けると(S4211:YES)(図27のS2711)、第1ストレージ装置10aに格納されている、当該更新要求の対象になっているレプリケーションファイルのデータ(メタデータ又は実体)を、受け付けた更新要求に従って更新する(S4212)(図27のS2712)。

【0241】

また第1サーバ装置3aは、メタデータを更新した場合はそのレプリケーションファイルのメタデータ同期要フラグ2312をONに設定し(S4213)、レプリケーションファイルの実体を更新した場合はそのレプリケーションファイルの実体同期要フラグ23

50

13をONに設定する(S4214)(図27のS2713)。

【0242】

図43は、図28に示したレプリケーションファイル参照処理S2800の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0243】

第1サーバ装置3aは、クライアント装置2からレプリケーションファイルに対する参照要求を受け付けたか否かをリアルタイムに監視している(S4311)。第1サーバ装置3aは、参照要求を受け付けると(S4311:YES)(図28のS2811)、そのレプリケーションファイルのデータ(メタデータ又は実体)を第1ストレージ装置10aから読み出し(S4312)(図28のS2812)、読み出したデータに基づきクライアント装置2に対して応答する情報を生成し、生成した応答情報をクライアント装置2に送信する(S4313)(図28のS2813)。

10

【0244】

図44は、図29に示した同期処理S2900の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0245】

第1サーバ装置3aは、クライアント装置2からレプリケーションファイルの同期要求を受け付けたか否かをリアルタイムに監視している(S4411)。第1サーバ装置3aは、同期要求を受け付けると(S4411:YES)(図29のS2911)、第1ストレージ装置10aのファイル格納領域に格納されているレプリケーションファイルのうち、メタデータ同期要フラグ2312又は実体同期要フラグ2313の少なくともいずれかがONに設定されているファイルを取得する(S4412)(図29のS2912)。

20

【0246】

そして第1サーバ装置3aは、取得したファイルのメタデータ又は実体を第2サーバ装置3bに送信するとともに(S4413)、当該レプリケーションファイルのメタデータ同期要フラグ2312又は実体同期要フラグ2313をOFFに設定する(S4414)(図29のS2913)。

【0247】

一方、第2サーバ装置3bは、メタデータ又は実体を受信すると(S4421)(図29のS2913)、受信したメタデータ又は実体に対応する、第2ストレージ装置10bに格納されているファイルのメタデータ又は実体を、受信したメタデータ又は実体(もしくは更新差分)に基づき更新する(S4422)(図29のS2914)。

30

【0248】

図45は、図30に示したメタデータアクセス処理S3000の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0249】

第1サーバ装置3aは、クライアント装置2からスタブ化されているファイルのメタデータに対するアクセス要求(参照要求又は更新要求)を受け付けたか否かをリアルタイムに監視している(S4511)。

【0250】

第1サーバ装置3aは、スタブ化されているファイルのメタデータに対するアクセス要求を受け付けると(S4511:YES)(図30のS3011)、受け付けたアクセス要求の対象になっている第1ストレージ装置10aのメタデータを取得し(S4512)、受け付けたアクセス要求に従い(S4513)、メタデータの参照(読み出したメタデータに基づく応答情報のクライアント装置2への送信)(S1514)、又はメタデータの更新を行う(S4515)(図30のS3012)。またメタデータの内容を更新した場合は(S4515)、そのファイルのメタデータ同期要フラグ2312にONを設定する(図30のS3013)。

40

【0251】

図46は、図31に示したスタブ化ファイル実体参照処理S3100の詳細を説明する

50

フローチャートである。以下、同図とともに説明する。

【0252】

第1サーバ装置3aは、クライアント装置2からスタブ化ファイルの実体に対する参照要求を受け付けると(S4611: YES)(図31のS3111)、当該スタブ化ファイルの実体が第1ストレージ装置10aに格納されているか否かを判断する(S4612)(図31のS3112)。

【0253】

スタブ化ファイルの実体が第1ストレージ装置10aに格納されている場合には(S4612: YES)、第1サーバ装置3aは、第1ストレージ装置10aから当該スタブ化ファイルの実体を読み出し、読み出した実体に基づきクライアント装置2に应答する情報を生成し、生成した应答情報をクライアント装置2に送信する(S4613)(図31のS3113)。

10

【0254】

一方、スタブ化ファイルの実体が第1ストレージ装置10aに格納されていない場合には(S4612: NO)、第1サーバ装置3aは、第2サーバ装置3bに対してスタブ化ファイルの実体を要求する(リコール要求)(S4614)(図31のS3114)。

【0255】

第1サーバ装置3aは、上記取得要求に応じて第2サーバ装置3bから送られてくるスタブ化ファイルの実体を受信すると(S4621、S4622、S4615)(図31のS3115)、受信した実体に基づき应答情報を生成し、生成した应答情報をクライアント装置2に送信する(S4616)(図31のS3116)。

20

【0256】

また第1サーバ装置3aは、上記第2サーバ装置3bから受信した実体を第1ストレージ装置10aに格納し、当該スタブ化ファイルのメタデータの当該ファイルの実体の格納先を示す情報(例えばブロックアドレス2218)に、当該ファイルの第1ストレージ装置10aにおける格納先を示す内容を設定する(S4617)。

【0257】

また第1サーバ装置3aは、当該ファイルのスタブ化フラグ2311をOFFに、レプリケーションフラグ2314をONに、メタデータ同期要フラグ2312をONに、夫々設定する(S4618)(図31のS3117)。

30

【0258】

図47は、図32に示したスタブ化ファイル実体更新処理S3200の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0259】

第1サーバ装置3aは、クライアント装置2からスタブ化ファイルの実体に対する更新要求を受け付けると(S4711: YES)(図32のS3211)、当該スタブ化ファイルの実体が第1ストレージ装置10aに格納されているか否かを判断する(S4712)(図32のS3212)。

【0260】

スタブ化ファイルの実体が第1ストレージ装置10aに格納されている場合(S4712: YES)、第1サーバ装置3aは、第1ストレージ装置10aに格納されている当該スタブ化ファイルの実体を更新要求の内容に従って更新するとともに(S4713)、当該スタブ化ファイルの実体同期要フラグ2313をONに設定する(S4714)(図32のS3213)。

40

【0261】

一方、上記判断の結果、スタブ化ファイルの実体が第1ストレージ装置10aに格納されていない場合には(S4712: NO)、第1サーバ装置3aは、第2サーバ装置3bに当該スタブ化ファイルの実体の取得要求(リコール要求)を送信する(S4715)(図32のS3214)。

【0262】

50

第1サーバ装置3aは、上記要求に応じて第2サーバ装置3bから送られてくるファイルの実体を受信すると(S4721、S4722、S4716)(S3215)、受信した実体の内容を更新要求の内容に従って更新し(S4717)、更新後の実体を当該スタブ化ファイルの実体として第1ストレージ装置10aに格納する(S4718)(図32のS3216)。

【0263】

また第1サーバ装置3aは、当該スタブ化ファイルのスタブ化フラグ2311をOFF、レプリケーションフラグ2314をOFF、メタデータ同期要フラグ2312をONに夫々設定する(S4719)。

【0264】

図48は、図33に示した仮想マシン復旧処理S3300及び図34に示したディレクトリイメージ事前回復処理S3400の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0265】

まず第1サーバ装置3aにおいて、ブートルード等を用いて記録媒体に記録されているインストールプログラムを実行し、第1サーバ装置3aに仮想化制御部305をインストールし、仮想化制御部305の機能を開始させる(S4811)(図33のS3311、S3312)。

【0266】

次に機能を開始した仮想化制御部305が、第2サーバ装置3bに仮想マシンイメージの取得要求を送信する(S4812)(図33のS3313)。

【0267】

第2サーバ装置3bは、第1サーバ装置3aから上記取得要求を受信すると(S4821)、当該取得要求に指定されている仮想マシンイメージを第2ストレージ装置10bから取得し、取得した仮想マシンイメージを第1サーバ装置3aに送信する(S4822)(図33のS3314、S3315)。

【0268】

第1サーバ装置3aは、第2サーバ装置3bから仮想マシンイメージを受信すると(S4813)(図33のS3316)、受信した仮想マシンイメージを第1ストレージ装置10aに格納し(S4814)、受信した仮想マシンイメージに基づき仮想マシン310の動作を開始する(S4815)(図33のS3317)。

【0269】

次に第1サーバ装置3aは、第2サーバ装置3bに対し、仮想マシン復旧処理S3300によって再起動された仮想マシン310のファイルシステム312が障害発生前に構成していたディレクトリイメージのルートディレクトリに存在するディレクトリのメタデータ、及びルートディレクトリに存在するファイルのメタデータの取得要求を、第2サーバ装置3bに送信する(S4816)(図34のS3411)。

【0270】

第2サーバ装置3bは、上記取得要求を受信すると(S4823)、要求されているルートディレクトリに存在するディレクトリのメタデータ、及びルートディレクトリに存在するファイルのメタデータを第2ストレージ装置10bから取得し、取得したメタデータを第1ストレージ装置10aに送信する(S4824)(図34のS3412、S3413)。

【0271】

次に第1サーバ装置3aは、第2サーバ装置3bからメタデータを受信すると(S4817)(図34のS3413)、受信したメタデータに従ったディレクトリイメージを第1ストレージ装置10aに構成(復元)する(S4818)(図34のS3414)。またこのとき、第1サーバ装置3aは、メタデータ同期要フラグ2312をONに、実同期要フラグ2313をONに夫々設定する(S4819)。

【0272】

10

20

30

40

50

そして第1サーバ装置3aは、第1ストレージ装置10aに上記ディレクトリイメージが構成されると、第1サーバ装置3aはクライアント装置2へのサービスを開始する(S4820)(図34のS3415)。

【0273】

図49は、図35に示したオンデマンド復元処理S3500の詳細を説明するフローチャートである。以下、同図とともに説明する。

【0274】

第1サーバ装置3aは、クライアント装置2から、あるファイルについてデータI/O要求を受け付けると(S4911: YES)(図35のS3511)、受け付けたデータI/O要求の対象になっているファイル(アクセス対象ファイル)のメタデータが、第1ストレージ装置10aに存在するか否かを調べる(S4912)(図35のS3512)。

10

【0275】

そしてメタデータが第1ストレージ装置10aに復元されている場合(S4912: YES)、第1サーバ装置3aは、受け付けたデータI/O要求の対象や種類、管理方式、スタブ化の有無等に応じて、受け付けたデータI/O要求に対応する処理を行い、クライアント装置2に応答を返す(S4913)(図35のS3518)。

【0276】

一方、アクセス対象ファイルのメタデータが第1ストレージ装置10aに復元されていない場合には(S4912: NO)、第1サーバ装置3aは、ルートディレクトリを起点としてアクセス対象ファイルが存在するディレクトリレベルに至るまでのディレクトリイメージを復元するためのデータを第2サーバ装置3bに要求する(S4914)。

20

【0277】

第2サーバ装置3bは、要求されたデータを第2ストレージ装置10bから取得し、取得したデータを第1サーバ装置3aに送信する(S4921、S4922、S4915)。

【0278】

第1サーバ装置3aは、第2サーバ装置3bから送られてくるデータを受信すると(S4915)、そのデータを用いてディレクトリイメージを第1ストレージ装置10aに復元する(S4916)(図35のS3513~S3516)。

30

【0279】

また第1サーバ装置3aは、アクセス対象ファイルのスタブ化フラグ2311をONに、レプリケーションフラグ2314をOFFに、メタデータ同期要フラグ2312をONに、夫々設定する(S4917)(図35のS3517)。

【0280】

次に第1サーバ装置3aは、受け付けたデータI/O要求の対象や種類、管理方式、スタブ化の有無等に応じて、受け付けたデータI/O要求に対応する処理を行い、クライアント装置2に応答を返す(S4918)(図35のS3518)。

【0281】

図50及び図51は、図37に示したオンデマンド復元処理(復元対象追加有)S3700の詳細を説明するフローチャートである。以下、同図とともに説明する。

40

【0282】

第1サーバ装置3aは、クライアント装置2から、あるファイルについてデータI/O要求を受け付けると(S5011: YES)(図37のS3711)、受け付けたデータI/O要求の対象になっているアクセス対象ファイルのメタデータが、第1ストレージ装置10aに存在するか否かを調べる(S5012)(図37のS3712)。

【0283】

メタデータが第1ストレージ装置10aに復元されている場合(S5012: YES)、第1サーバ装置3aは、受け付けたデータI/O要求の対象や種類、管理方式、スタブ化の有無等に応じて、受け付けたデータI/O要求に対応する処理を行い、クライアント

50

装置 2 に応答を返す ( S 5 0 1 3 ) ( 図 3 7 の S 3 7 1 8 )。

【 0 2 8 4 】

一方、アクセス対象ファイルのメタデータが第 1 ストレージ装置 1 0 a に復元されていない場合には ( S 5 0 1 2 : N O )、第 1 サーバ装置 3 a は、ルートディレクトリを起点としてアクセス対象ファイルが存在するディレクトリレベルに至るまでのディレクトリイメージを復元するためのデータを、第 2 サーバ装置 3 b に要求する ( S 5 0 1 4 )。

【 0 2 8 5 】

第 2 サーバ装置 3 b は、上記要求を受信すると、データ I / O 要求が前述した所定の条件を満たすか否かを判断する ( S 5 0 2 2 )。

【 0 2 8 6 】

所定の条件を満たさない場合には ( S 5 0 2 2 : N O )、S 5 0 2 4 に進む。一方、所定の条件を満たす場合には ( S 5 0 2 2 : Y E S )、第 2 サーバ装置 3 b は、前述した所定の選出方法に従い追加するディレクトリイメージを選出する ( S 5 0 2 3 )。

【 0 2 8 7 】

S 5 0 2 4 では、第 2 サーバ装置 3 b は、S 5 0 2 1 で受信した要求に指定されているディレクトリイメージを復元するためのデータと、S 5 0 2 3 で選出したディレクトリイメージを復元するためのデータとを、第 2 ストレージ装置 1 0 b から取得し、取得したデータを第 1 サーバ装置 3 a に送信する ( 図 3 7 の S 3 7 1 3 ~ S 3 7 1 5 )。

【 0 2 8 8 】

第 1 サーバ装置 3 a は、上記データを受信すると ( S 5 0 1 5 )、受信したデータを用いて、第 1 ストレージ装置 1 0 a にディレクトリイメージを復元する ( S 5 0 1 6 ) ( 図 3 7 の S 3 7 1 6 )。

【 0 2 8 9 】

次に第 1 サーバ装置 3 a は、アクセス対象ファイルのスタブ化フラグ 2 3 1 1 を O N に、レプリケーションフラグ 2 3 1 4 を O F F に、メタデータ同期要フラグ 2 3 1 2 を O N に、夫々設定する ( S 5 0 1 7 ) ( 図 3 7 の S 3 7 1 7 )。

【 0 2 9 0 】

そして第 1 サーバ装置 3 a は、受け付けたデータ I / O 要求の対象や種類、管理方式、スタブ化の有無等に応じて、受け付けたデータ I / O 要求に対応する処理を行い、クライアント装置 2 に応答を返す ( S 5 0 1 8 ) ( 図 3 7 の S 3 7 1 8 )。

【 0 2 9 1 】

図 5 2 は、図 3 8 に示した再スタブ化回避処理 S 3 8 0 0 の詳細を説明するフローチャートである。以下、同図とともに説明する。

【 0 2 9 2 】

第 2 サーバ装置 3 b は、オンデマンド復元処理 S 3 5 0 0 ( 又は図 3 7 に示したオンデマンド復元処理 ( 復元対象追加有 ) S 3 7 0 0 ) の実行中、単位時間当たりの再スタブ化の発生頻度が予め設定された閾値 ( 以下、再スタブ化頻度閾値と称する。 ) 以上になっているか否か、もしくは再スタブ化の発生時間間隔が予め設定された閾値 ( 以下、再スタブ化発生時間間隔閾値と称する。 ) 未満になっているか否かを監視する ( S 5 2 1 1、S 5 2 1 2 ) ( 図 3 8 の S 3 8 1 1 ~ S 3 8 1 3 )。

【 0 2 9 3 】

第 2 サーバ装置 3 b は、上記監視において再スタブ化の発生頻度が再スタブ化頻度閾値以上になっていることを検知すると ( S 5 2 1 1 : Y E S )、抑制フラグ管理テーブル 3 6 6 に管理される抑制フラグ 3 6 6 1 を O N に設定する ( S 5 2 1 3 )。

【 0 2 9 4 】

また第 2 サーバ装置 3 b は、上記監視において、再スタブ化の発生時間間隔が再スタブ化発生時間間隔閾値未満になっていることを検知すると ( S 5 2 1 2 : Y E S )、抑制フラグ 3 6 6 1 を O N に設定する ( S 5 2 1 3 )。

【 0 2 9 5 】

また第 2 サーバ装置 3 b は、上記監視において、再スタブ化の発生頻度が再スタブ化頻

10

20

30

40

50

度閾値未満になっており（S5211：NO）、かつ、再スタブ化の発生時間間隔が再スタブ化発生時間間隔閾値以上になっている場合は（S5212：NO）、抑制フラグ3661をOFFに設定する（S5214）（図38のS3814）。

【0296】

S5215では、第2サーバ装置3bは、抑制フラグ3661がONかOFFかを判断する。抑制フラグ3661がONであれば（S5215：ON）、第2サーバ装置3bから第1サーバ装置3aに送信するディレクトリイメージの量を抑制する処理を開始する（S5216）。尚、既に抑制を開始している場合には抑制を継続する。

【0297】

一方、抑制フラグ3661がOFFであれば（S5215：OFF）、第2サーバ装置3bは、第1サーバ装置3aに送信するディレクトリイメージの量を抑制する処理を終了する（S5217）。尚、既に抑制を終了している場合は抑制なしの状態を維持する。

【0298】

以上詳細に説明したように、本実施形態の情報処理システム1にあつては、第1サーバ装置3aの障害復旧に際し、第2サーバ装置3bが、第1サーバ装置3aがデータI/O要求の受け付けを開始するのに先立ち、第2ストレージ装置10bに記憶しているファイルのデータのうち最上位のディレクトリから所定の下位階層までのディレクトリイメージを第1サーバ装置3aに送信し、第1サーバ装置3aが、第2サーバ装置3bから送られてくるディレクトリイメージを第1ストレージ装置10aに復元した後、データI/O要求の受け付けを再開する。

【0299】

このように、本実施形態の情報処理システム1にあつては、第1サーバ装置3aの障害復旧に際し、障害発生前に第1ストレージ装置10aに存在していたディレクトリイメージの全てを復元するのではなく、最上位のディレクトリから所定の下位階層までのディレクトリイメージのみを復元するので、障害発生前に第1ストレージ装置10aに存在していた全てのディレクトリイメージを復元する場合に比べて復元に要する時間を短縮することができ、早期にサービスを再開することができる。また全てのディレクトリイメージを復元する場合に比べて情報処理システム1に与える負荷が少なく済む。

【0300】

また第2サーバ装置3bは、第1サーバ装置3aから第1ストレージ装置10aに復元されていないディレクトリイメージを要求された場合に、要求の対象になっているディレクトリイメージを第2ストレージ装置10bから読み出して送信するとともに、所定の選出方法に従って選出したディレクトリイメージとは異なる追加のディレクトリイメージを第2ストレージ装置10bから読み出して送信し、第1サーバ装置3aは、第2サーバ装置3bから送られてくるディレクトリイメージに基づきデータI/O要求を処理するとともに、ディレクトリイメージ及び第2サーバ装置3bから送られてくる追加のディレクトリイメージを第1ストレージ装置10aに復元する。

【0301】

このように、本実施形態の情報処理システム1にあつては、第2サーバ装置3bが、第1サーバ装置3aから第1ストレージ装置10aに復元されていないディレクトリイメージを要求された場合に、要求の対象になっているディレクトリイメージに加えて、所定の選出方法に従って選出したディレクトリイメージとは異なる追加のディレクトリイメージを第2ストレージ装置10bから読み出して送信し、第1サーバ装置3aは、ディレクトリイメージ及び追加のディレクトリイメージの双方のディレクトリイメージを第1ストレージ装置10aに復元するので、ディレクトリイメージの復旧速度を自動的に速めることができる。

【0302】

また第2サーバ装置3bは、再スタブ化の発生頻度が予め設定された閾値以上になっているか、もしくは、再スタブ化の発生時間間隔が予め設定された閾値未満になっている場合に、第1サーバ装置3aへのディレクトリイメージ又は追加のディレクトリイメージの

10

20

30

40

50

送信を自動的に抑制するので、再スタブ化の発生を抑制することができ、再スタブ化によって情報処理システム1の資源が浪費されてしまうのを防ぐことができる。

【0303】

以上、本実施形態について説明したが、上記実施形態は本発明の理解を容易にするためのものであり、本発明を限定して解釈するためのものではない。本発明は、その趣旨を逸脱することなく、変更、改良され得ると共に、本発明にはその等価物も含まれる。

【0304】

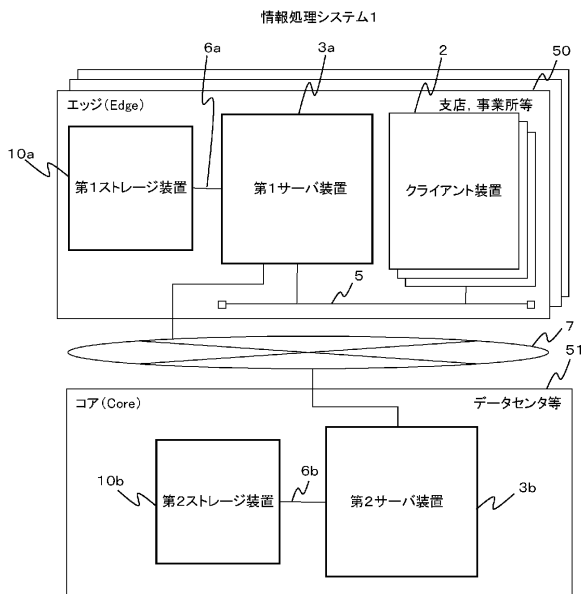
例えば、以上の説明では、ファイル共有処理部311、ファイルシステム312、データ操作要求受付部313、データ複製/移動処理部314、ファイルアクセスログ取得部317、及びカーネル/ドライバ318の各機能が仮想マシン310において実現されているものとして説明したが、これらの機能は必ずしも仮想マシン310において実現されるものでなくてもよい。

【0305】

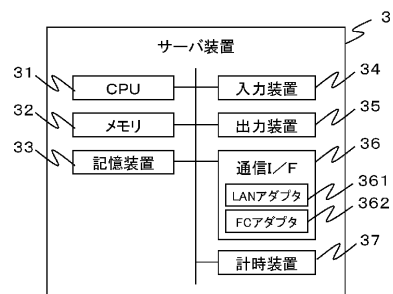
また前述したディレクトリイメージ事前復旧処理S3400では、ファイルシステム312が障害発生前に構成していたディレクトリイメージのルートディレクトリに存在するディレクトリのメタデータ、及びルートディレクトリに存在するファイルのメタデータを復旧させるようにしているが、第1サーバ装置3aの能力に余裕がある場合には、より下位階層までのディレクトリイメージを復旧させるようにしてもよい。

10

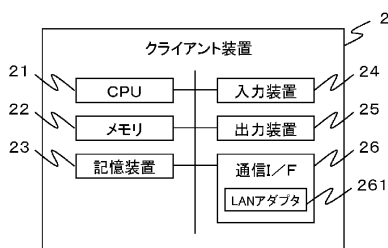
【図1】



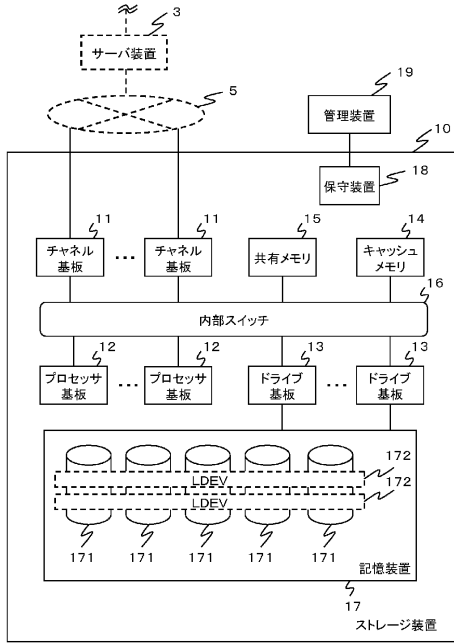
【図3】



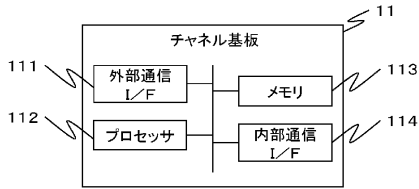
【図2】



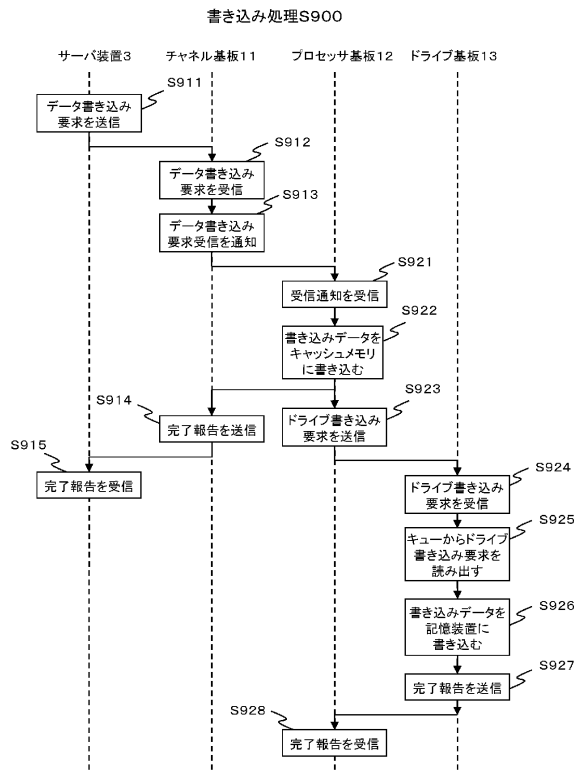
【図4】



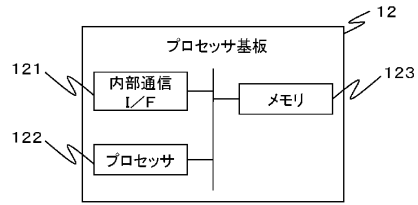
【図5】



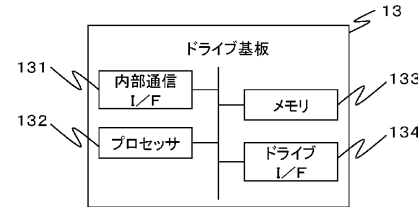
【図9】



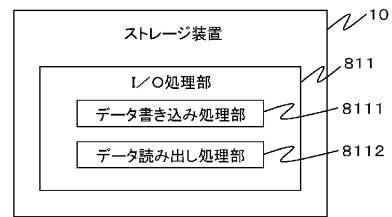
【図6】



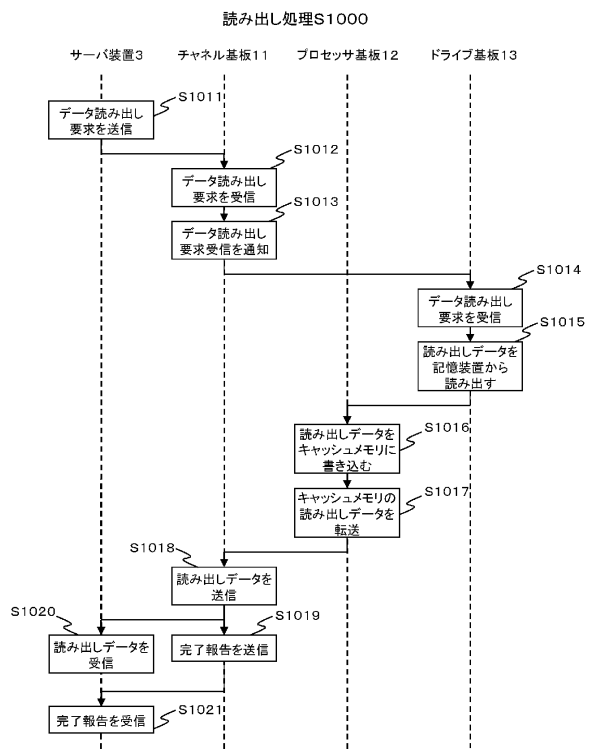
【図7】



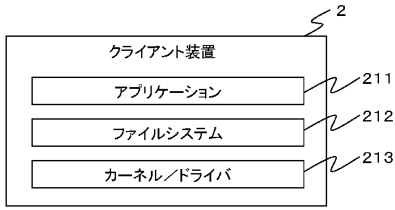
【図8】



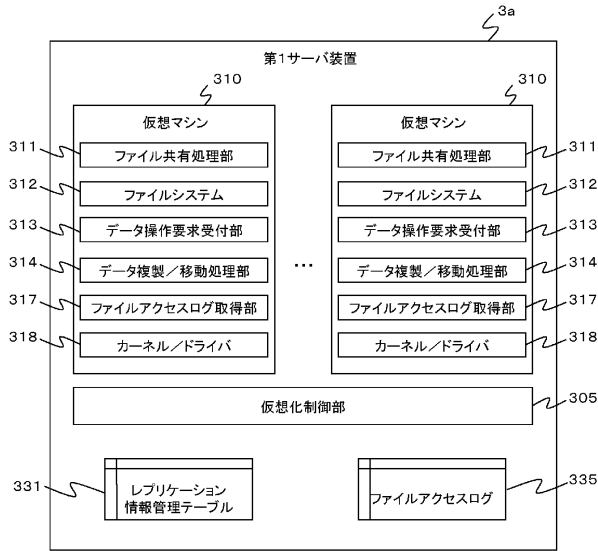
【図10】



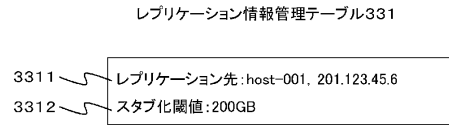
【図11】



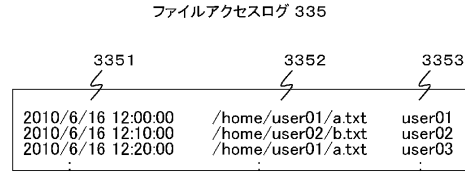
【図12】



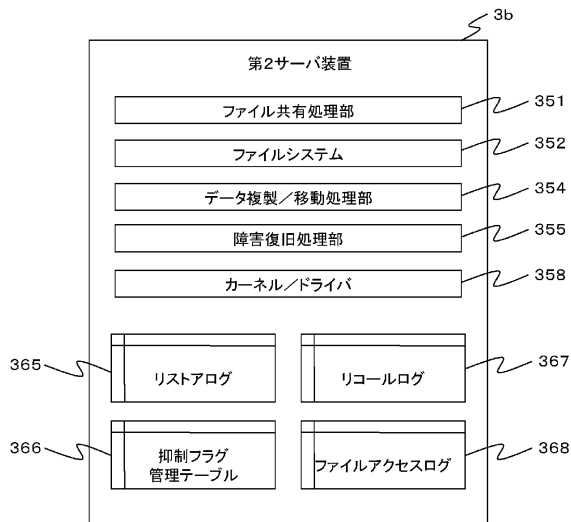
【図13】



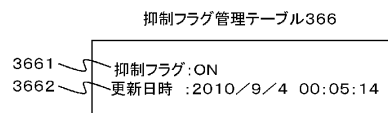
【図14】



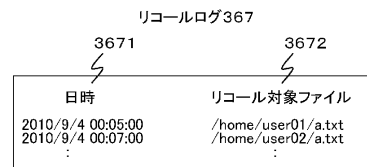
【図15】



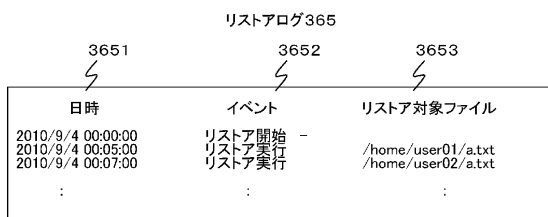
【図17】



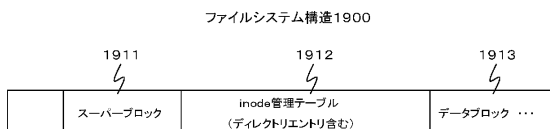
【図18】



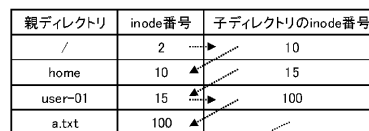
【図16】



【図19】

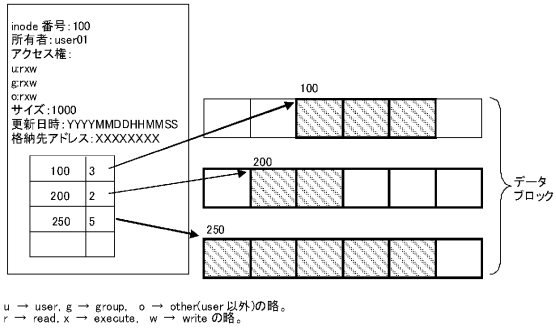


【図20】

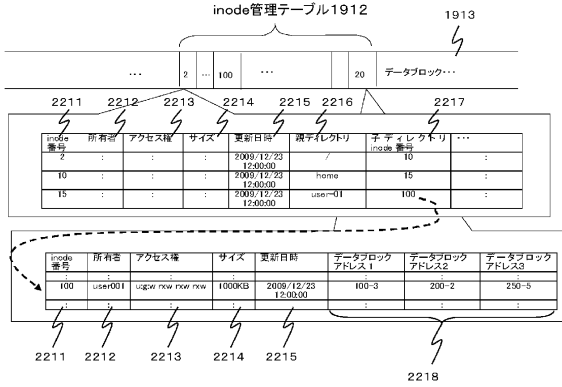


データブロック

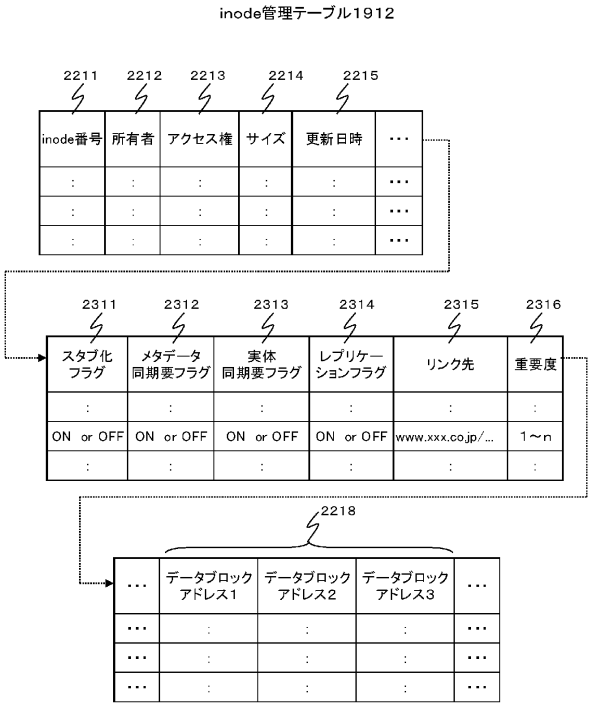
【図 2 1】



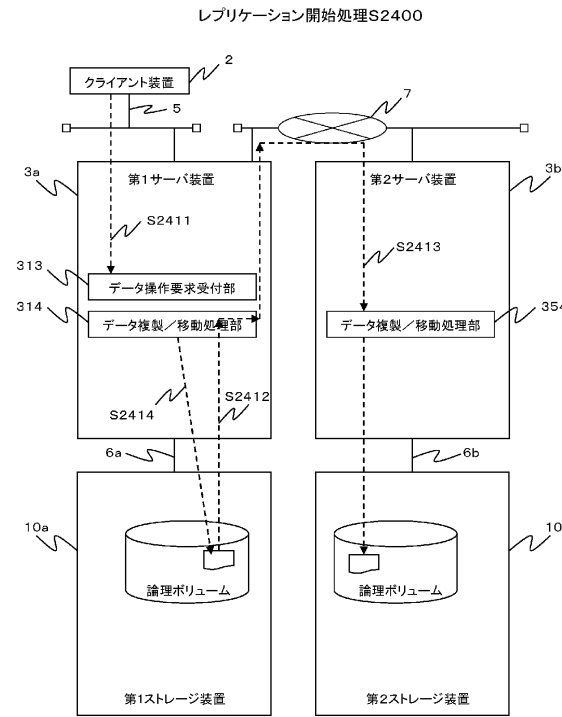
【図 2 2】



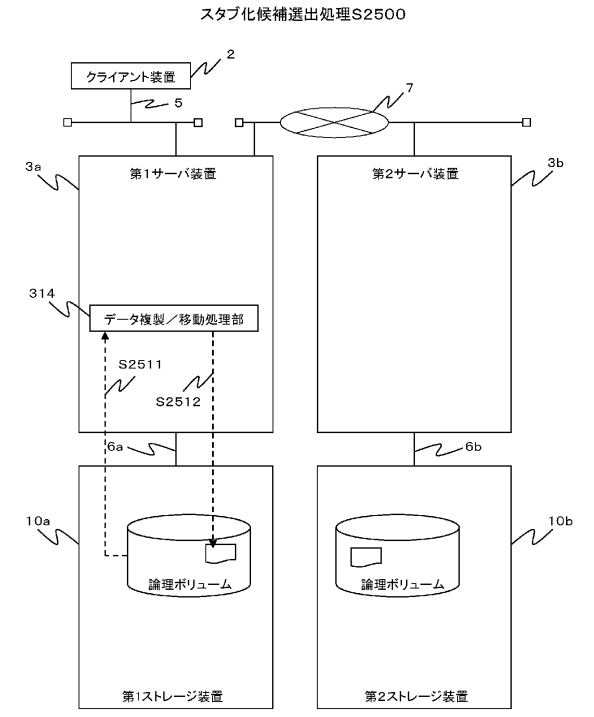
【図 2 3】



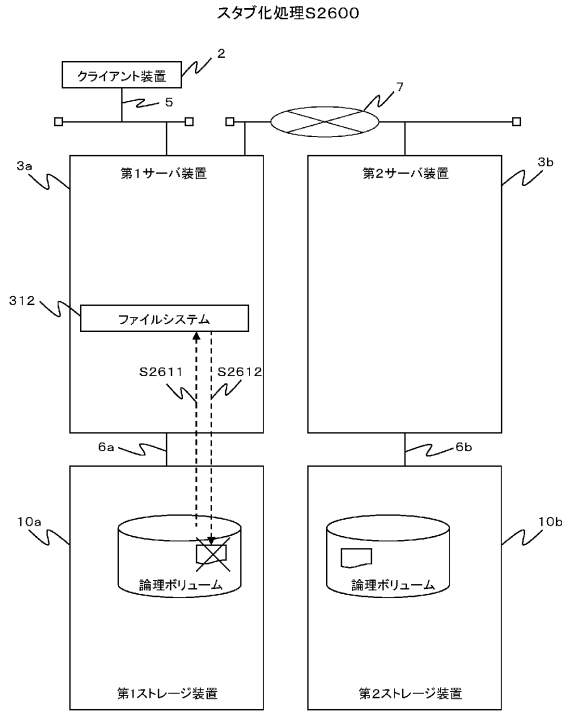
【図 2 4】



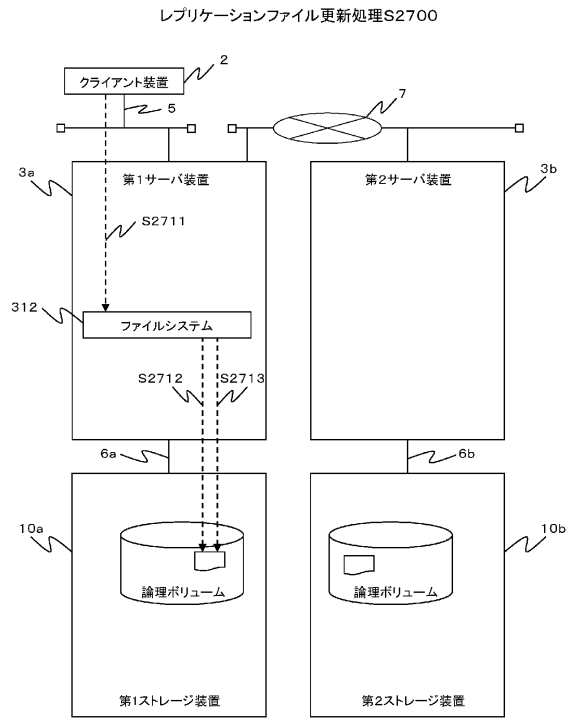
【図 2 5】



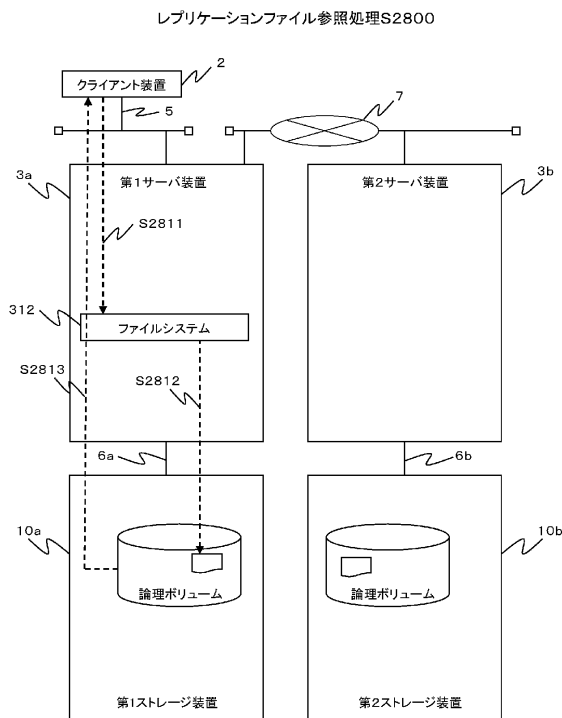
【図26】



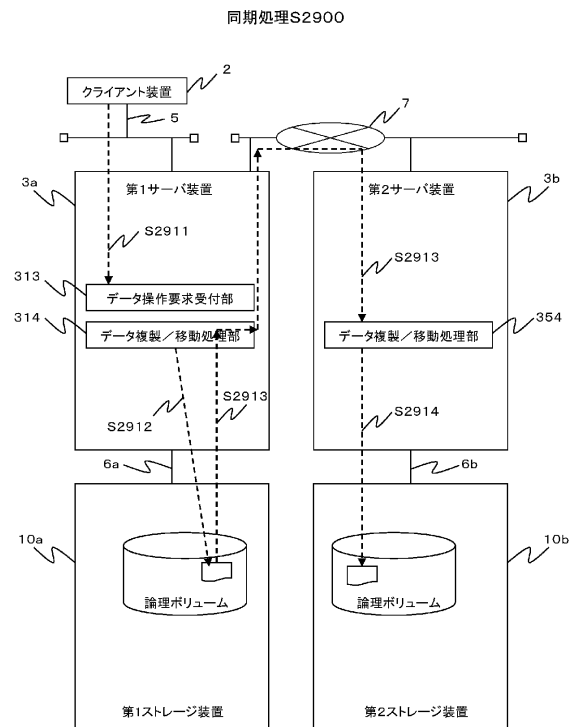
【図27】



【図28】

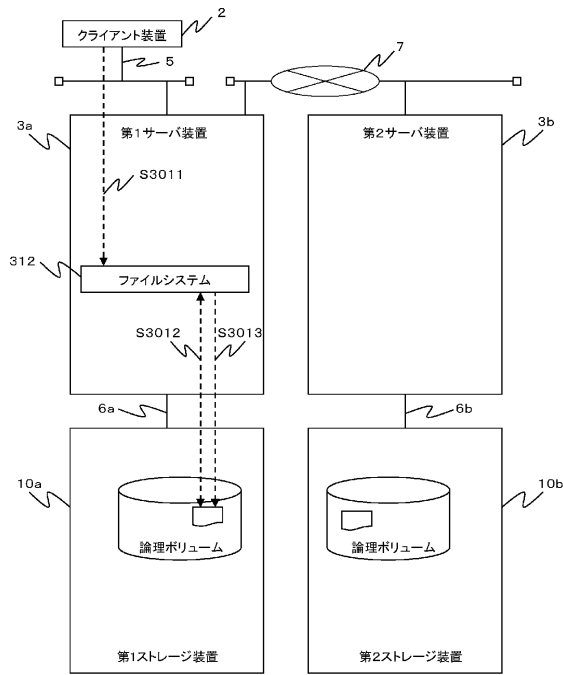


【図29】



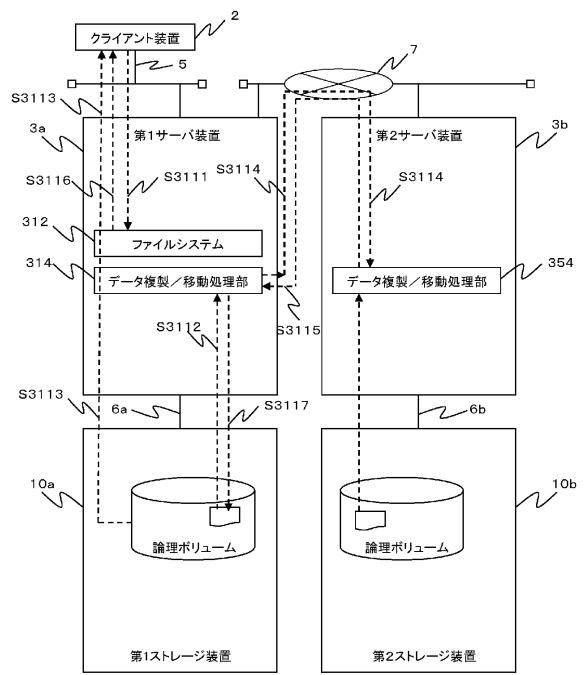
【図30】

メタデータアクセス処理S3000



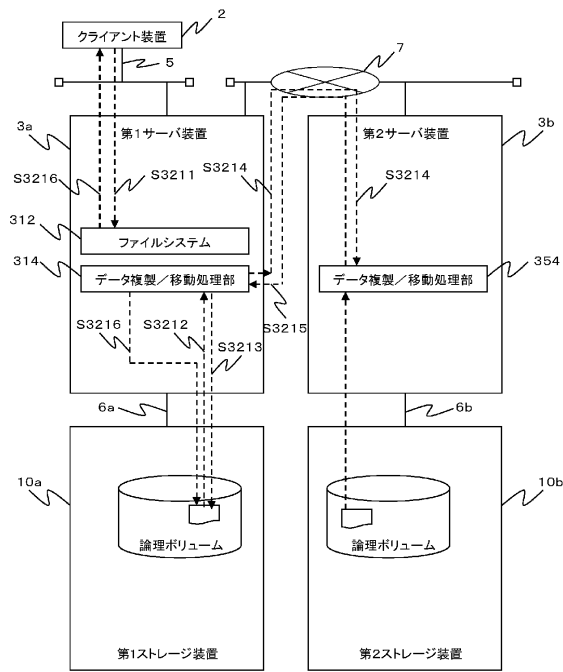
【図31】

スタブ化ファイル実参照処理S3100



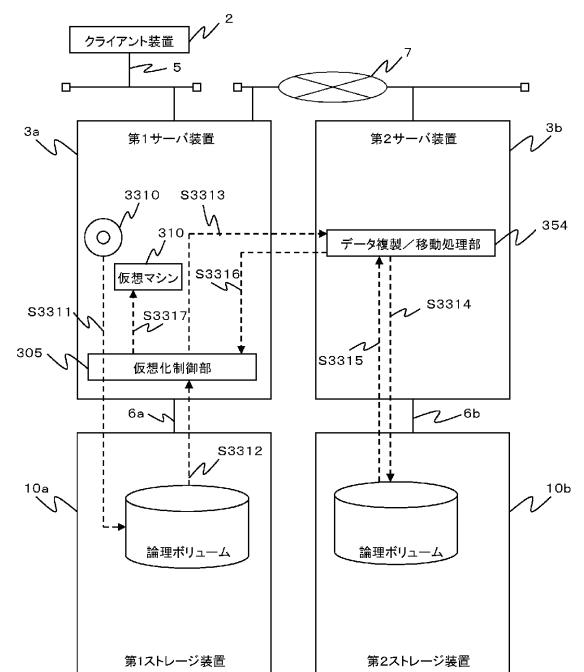
【図32】

スタブ化ファイル実更新処理S3200



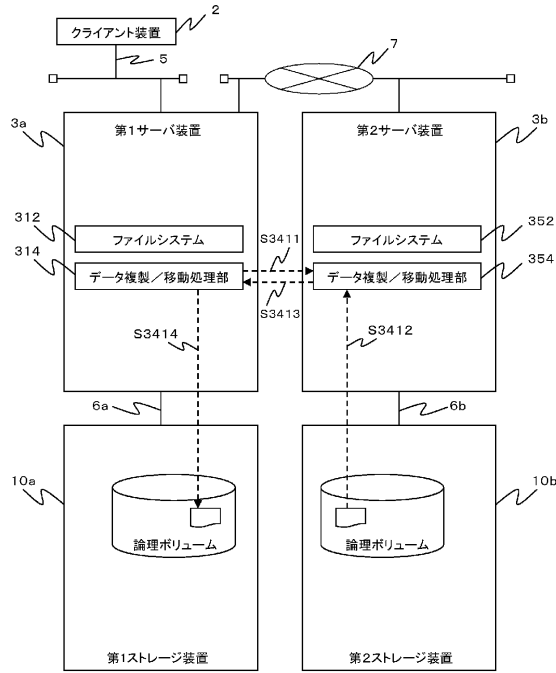
【図33】

仮想マシン復旧処理S3300



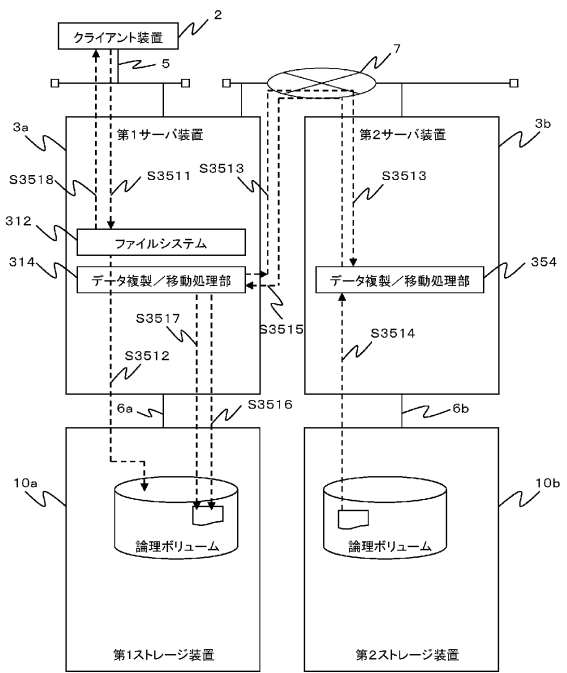
【図34】

ディレクトリイメージ事前回復処理S3400



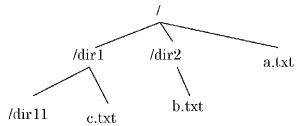
【図35】

オンデマンド復元処理S3500

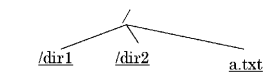


【図36】

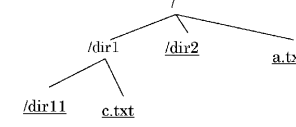
(0)



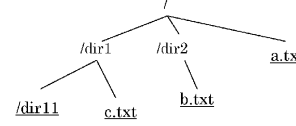
(A)



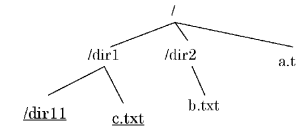
(B)



(C)

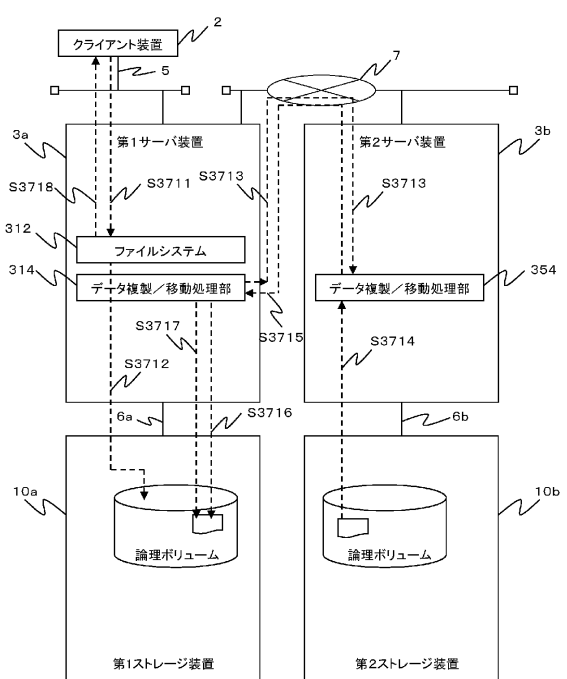


(D)



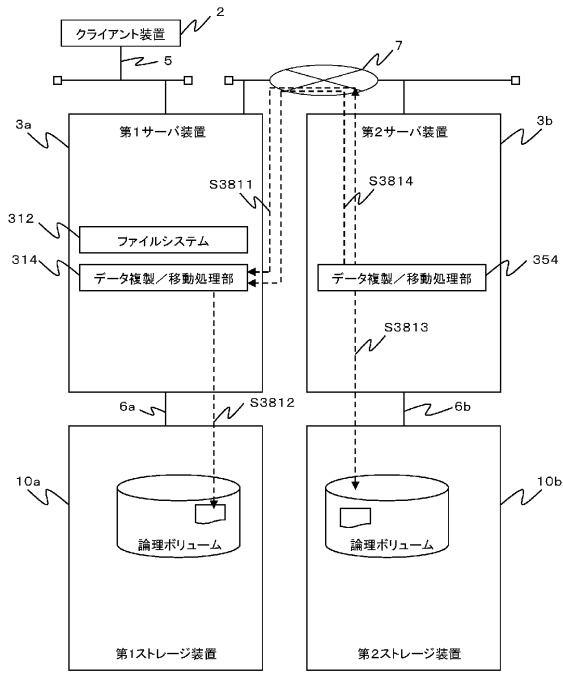
【図37】

オンデマンド復元処理(復元対象追加有)S3700



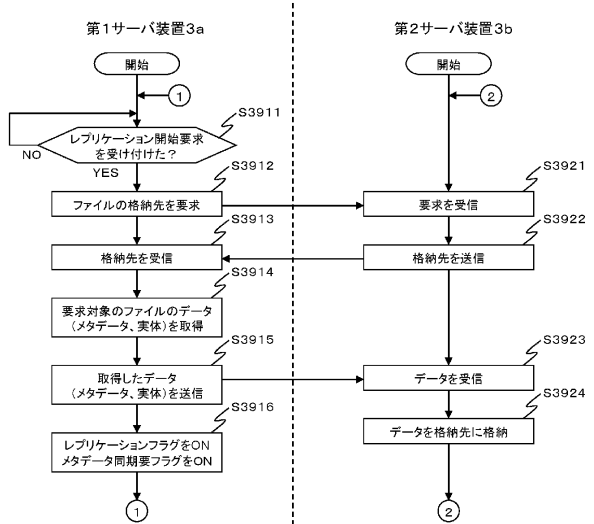
【図38】

再スタブ化回避処理S3800



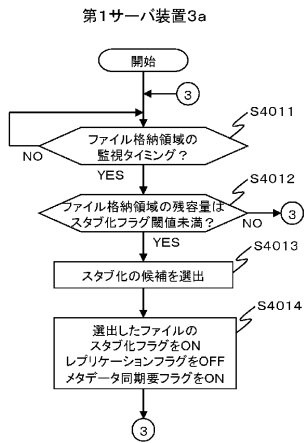
【図39】

レプリケーション開始処理S2400



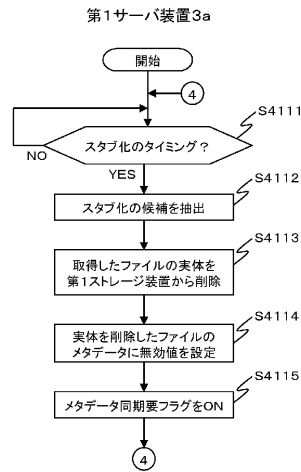
【図40】

スタブ化候補選出処理S2500



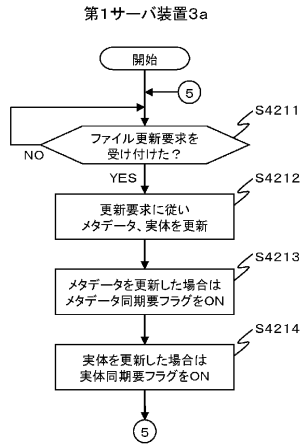
【図41】

スタブ化処理S2600



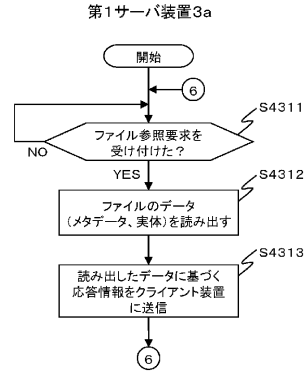
【図42】

レプリケーションファイル更新処理S2700



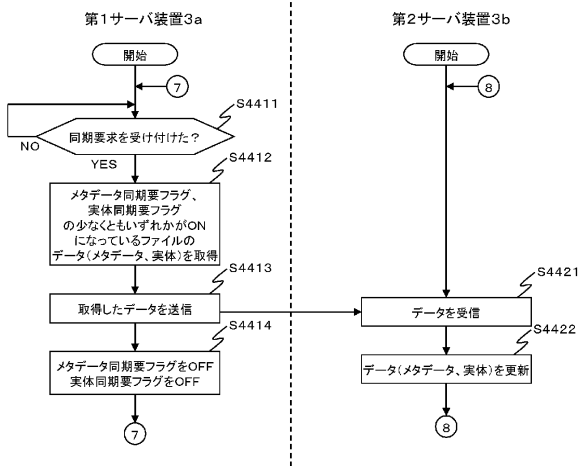
【図43】

レプリケーションファイル参照処理S2800



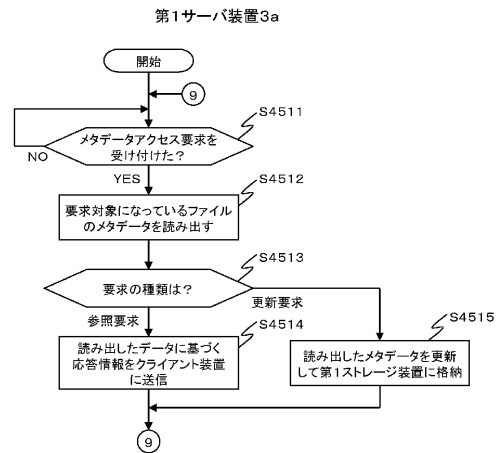
【図44】

同期処理S2900



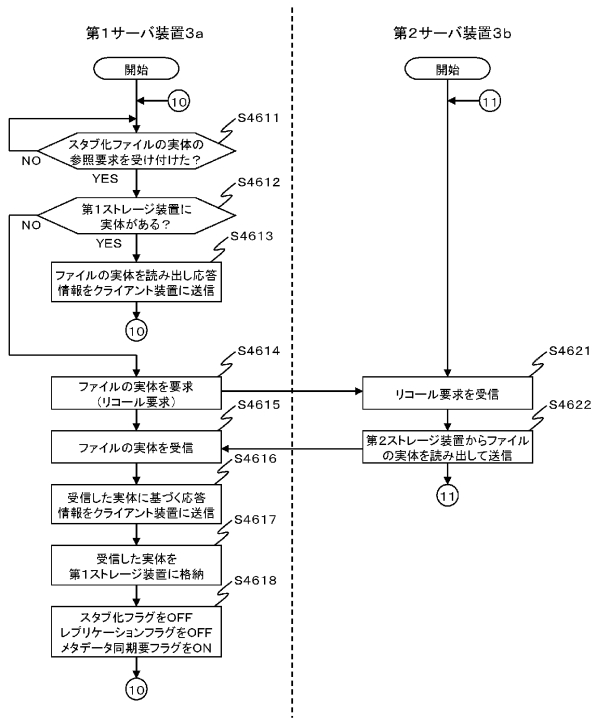
【図45】

メタデータアクセス処理S3000



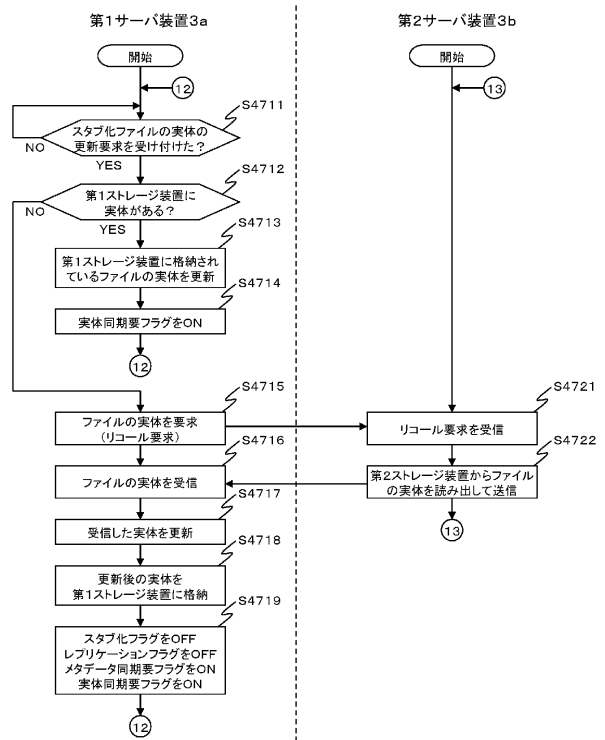
【図46】

スタブ化ファイル実参照処理S3100



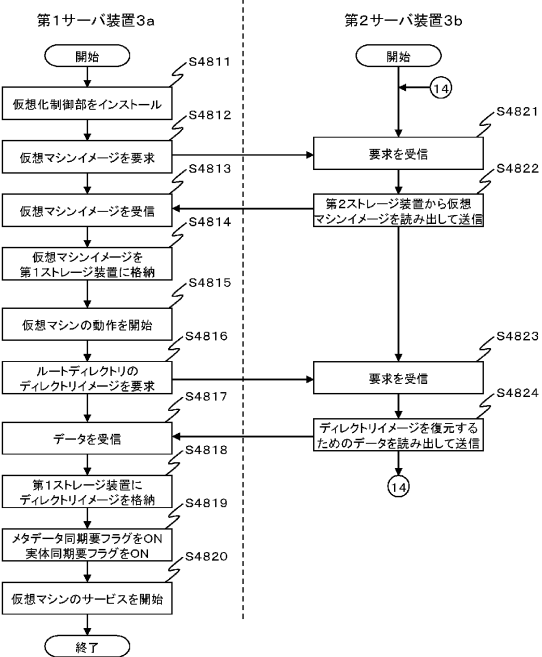
【図47】

スタブ化ファイル実更新処理S3200



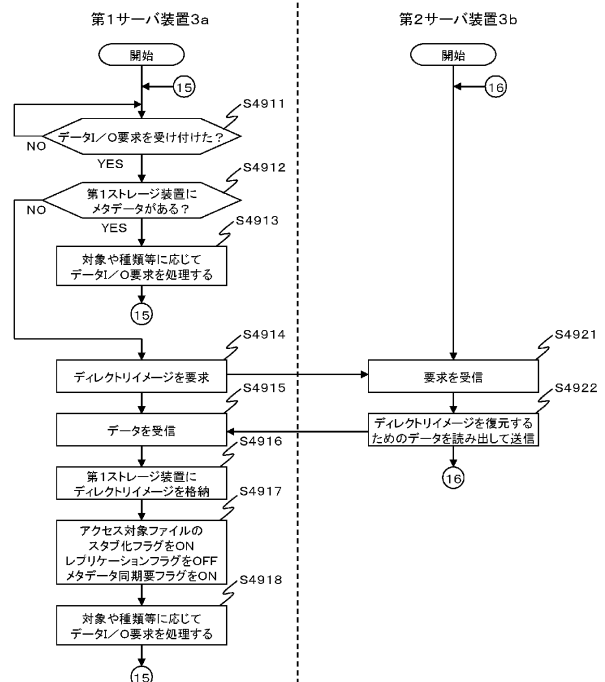
【図48】

仮想マシン復旧処理S3300、及び  
ディレクトリイメージ事前回復処理S3400



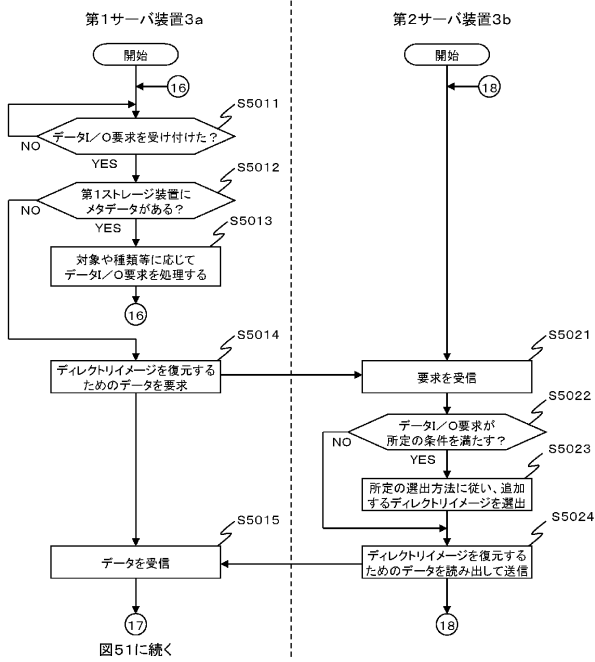
【図49】

オンデマンド復元処理S3500



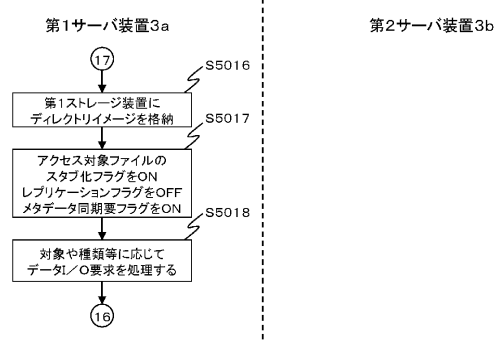
【図50】

オンデマンド復元処理(復元対象追加)S3700



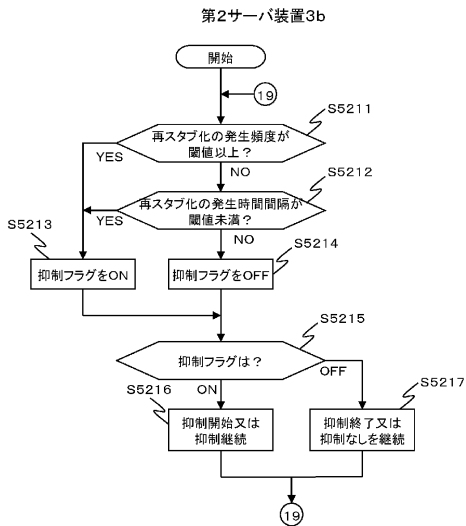
【図51】

オンデマンド復元処理(復元対象追加)S3700



【図52】

再スタブ化回避処理S3800



---

フロントページの続き

- (56)参考文献 特開2005-55947(JP,A)  
特開2005-141555(JP,A)  
特表2006-508473(JP,A)  
特開2008-33519(JP,A)  
特開2007-280099(JP,A)  
米国特許出願公開第2003/0158862(US,A1)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00  
G06F 3/06