

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6963168号  
(P6963168)

(45) 発行日 令和3年11月5日(2021.11.5)

(24) 登録日 令和3年10月19日(2021.10.19)

(51) Int.Cl. F I  
G O 6 F 13/00 (2006.01) G O 6 F 13/00 3 5 3 C

請求項の数 5 (全 21 頁)

|   |   |
|---|---|
| <p>(21) 出願番号 特願2017-121571 (P2017-121571)<br/>                 (22) 出願日 平成29年6月21日 (2017. 6. 21)<br/>                 (65) 公開番号 特開2019-8417 (P2019-8417A)<br/>                 (43) 公開日 平成31年1月17日 (2019. 1. 17)<br/>                 審査請求日 令和2年3月10日 (2020. 3. 10)</p> | <p>(73) 特許権者 000005223<br/>                 富士通株式会社<br/>                 神奈川県川崎市中原区上小田中4丁目1番<br/>                 1号<br/>                 (74) 代理人 110002918<br/>                 特許業務法人扶桑国際特許事務所<br/>                 (72) 発明者 野澤 康文<br/>                 神奈川県川崎市中原区上小田中4丁目1番<br/>                 1号 富士通株式会社内<br/> <br/>                 審査官 今川 悟</p> |
|---|---|

最終頁に続く

(54) 【発明の名称】 情報処理装置、メモリ制御方法およびメモリ制御プログラム

(57) 【特許請求の範囲】

【請求項 1】

他の情報処理装置からの要求に応じてコネクションを追加可能な通信部と、  
 維持されているコネクションの数および前記維持されているコネクションそれぞれにおいて最後に通信が行われたタイミングを示すコネクション情報を記憶する記憶部と、  
 単位時間当たりの前記維持されているコネクションの増加数が第1の閾値を超えておりかつ前記維持されているコネクションの数が第2の閾値を超えている場合、前記コネクション情報に基づいて、前記維持されているコネクションのうち前記最後に通信が行われたタイミングが古い方から優先的に少なくとも1つのコネクションを選択し、選択した前記少なくとも1つのコネクションを切断させる処理部と、  
 を有する情報処理装置。

【請求項 2】

前記コネクション情報は、前記最後に通信が行われたタイミングに応じた順序で前記維持されているコネクションの識別情報を記載したリスト情報を含む、  
 請求項 1 記載の情報処理装置。

【請求項 3】

前記少なくとも1つのコネクションとして、前記処理部は、前記維持されているコネクションの数と前記第2の閾値との差に相当する数のコネクションを選択する、  
 請求項 1 記載の情報処理装置。

【請求項 4】

情報処理装置が実行するメモリ制御方法であって、

前記情報処理装置で維持されているコネクションの数と、前記維持されているコネクションそれぞれにおいて最後に通信が行われたタイミングとを監視し、

単位時間当たりの前記維持されているコネクションの増加数が第1の閾値を超えておりかつ前記維持されているコネクションの数が第2の閾値を超えている場合、前記維持されているコネクションのうち前記最後に通信が行われたタイミングが古い方から優先的に少なくとも1つのコネクションを選択し、

選択した前記少なくとも1つのコネクションを切断させる、  
メモリ制御方法。

【請求項5】

10

コンピュータに、

前記コンピュータで維持されているコネクションの数と、前記維持されているコネクションそれぞれにおいて最後に通信が行われたタイミングとを監視し、

単位時間当たりの前記維持されているコネクションの増加数が第1の閾値を超えておりかつ前記維持されているコネクションの数が第2の閾値を超えている場合、前記維持されているコネクションのうち前記最後に通信が行われたタイミングが古い方から優先的に少なくとも1つのコネクションを選択し、

選択した前記少なくとも1つのコネクションを切断させる、  
処理を実行させるメモリ制御プログラム。

【発明の詳細な説明】

20

【技術分野】

【0001】

本発明は情報処理装置、メモリ制御方法およびメモリ制御プログラムに関する。

【背景技術】

【0002】

2つの情報処理装置がネットワークを介して通信する場合、当該2つの情報処理装置の間にTCP (Transmission Control Protocol) コネクションなどのコネクションを確立することがある。1つの情報処理装置が複数の他の情報処理装置と通信し得る場合、当該1つの情報処理装置が同時に複数のコネクションを維持することもある。例えば、1つの情報処理装置が複数の他の情報処理装置それぞれからコネクション確立要求を受信し、コネクション確立要求に応じて新規コネクションを確立することで、当該情報処理装置と複数の他の情報処理装置の間に複数のコネクションを維持することがある。

30

【0003】

新規コネクションを確立すると、通信相手から受信したデータを一時的に保存する受信バッファなどを含むメモリ領域をその新規コネクションに対して割り当てることになる。複数のコネクションを維持することがある情報処理装置は、最後に通信が行われてから所定時間経過してもまだ切断されていない古いコネクションを切断するなど、経過時間に応じて古いコネクションを整理することがある。

【0004】

例えば、アプリケーションサーバとデータサーバとの間のコネクションを管理するWebシステムが提案されている。提案のWebシステムは、予めアプリケーションサーバとデータサーバとの間に複数のコネクションを確立してコネクションプールに入れておき、アプリケーションサーバで実行されるサブレットに対して動的にコネクションを割り当てる。Webシステムは、サブレットに対して割り当てたコネクションのうち、使用されていない状態が所定時間続いているコネクションの割り当てを解除する。

40

【0005】

また、アプリケーションサーバが複数のデータベースサーバに接続する計算機システムが提案されている。提案の計算機システムは、アプリケーションサーバが複数のコネクションを確立してコネクションプールに入れておき、各コネクションに対して接続先のデータベースサーバを識別する識別子を付与しておく。計算機システムは、あるコネクション

50

で障害が検出されると、障害が検出されたコネクションと同じ識別子が付与されたコネクションのみをコネクションプールから削除する。

【0006】

また、サービス利用装置とサービス提供装置との間に維持するコネクションの数を制御するコンピュータシステムが提案されている。提案のコンピュータシステムは、予めサービス利用装置とサービス提供装置との間に複数のコネクションを確立する。コンピュータシステムは、サービス提供装置の負荷を監視し、負荷が増加した場合にはサービス利用装置が使用可能なコネクションの数を減少させる。

【先行技術文献】

【特許文献】

10

【0007】

【特許文献1】特開2000-29814号公報

【特許文献2】特開2007-226398号公報

【特許文献3】特開2009-181481号公報

【発明の概要】

【発明が解決しようとする課題】

【0008】

ところで、複数のノードに並列にプログラムを実行させる並列処理では、これら複数のノードが一斉に特定の情報処理装置と通信するデータ転送フェーズが発生することがある。例えば、プログラムの起動前に複数のノードそれぞれが、データサーバなどの特定の情報処理装置からローカル記憶装置に入力データをコピーする「ステージイン」が行われることがある。また、プログラムの終了後に複数のノードそれぞれが、ローカル記憶装置からデータサーバなどの特定の情報処理装置に出力データを転送する「ステージアウト」が行われることがある。データ転送フェーズでは、短時間のうちに特定の情報処理装置と複数のノードとの間に複数の新規コネクションが確立され得る。

20

【0009】

しかし、情報処理装置に既に比較的多数のコネクションが維持されているときに並列処理のデータ転送フェーズが発生すると、メモリ領域が不足して、一部のノードについて新規コネクションの確立およびデータ転送の開始が迅速に行われられない可能性がある。このとき、通常の切断契機に従って既存コネクションが切断されるのを待っている間は、データ転送フェーズの完了が遅延して並列処理の完了が遅れてしまうおそれがある。

30

【0010】

1つの側面では、本発明は、並列処理におけるデータ転送フェーズの遅延を抑制する情報処理装置、メモリ制御方法およびメモリ制御プログラムを提供することを目的とする。

【課題を解決するための手段】

【0011】

1つの態様では、通信部と記憶部と処理部とを有する情報処理装置が提供される。通信部は、他の情報処理装置からの要求に応じてコネクションを追加可能である。記憶部は、維持されているコネクションの数および維持されているコネクションそれぞれにおいて最後に通信が行われたタイミングを示すコネクション情報を記憶する。処理部は、単位時間当たりの維持されているコネクションの増加数が第1の閾値を超えておりかつ維持されているコネクションの数が第2の閾値を超えている場合、コネクション情報に基づいて、維持されているコネクションのうち最後に通信が行われたタイミングが古い方から優先的に少なくとも1つのコネクションを選択し、選択した少なくとも1つのコネクションを切断させる。

40

【0012】

また、1つの態様では、情報処理装置が実行するメモリ制御方法が提供される。また、1つの態様では、コンピュータに実行させるメモリ制御プログラムが提供される。

【発明の効果】

【0013】

50

1つの側面では、並列処理におけるデータ転送フェーズの遅延を抑制できる。

【図面の簡単な説明】

【0014】

【図1】第1の実施の形態の情報処理装置を説明する図である。

【図2】第2の実施の形態の並列処理システムの例を示す図である。

【図3】データサーバのハードウェア例を示すブロック図である。

【図4】ノードとデータサーバのソフトウェア例を示すブロック図である。

【図5】コネクション管理リストの例を示す図である。

【図6】コネクション数テーブルの例を示す図である。

【図7】ジョブ実行の手順例を示すフローチャートである。

【図8】コネクション確立の手順例を示すフローチャートである。

【図9】ファイルアクセスの手順例を示すフローチャートである。

【図10】コネクション切断の手順例を示すフローチャートである。

【図11】コネクション切断判定の手順例を示すフローチャートである。

【発明を実施するための形態】

【0015】

以下、本実施の形態を図面を参照して説明する。

[第1の実施の形態]

第1の実施の形態を説明する。

【0016】

図1は、第1の実施の形態の情報処理装置を説明する図である。

第1の実施の形態の情報処理装置10は、ネットワークを介して他の情報処理装置と通信する。情報処理装置10は、サーバコンピュータでもよく、他の情報処理装置が使用するデータを管理するデータサーバでもよい。例えば、情報処理装置10は、プログラムを並列に実行する複数のノードを含む並列処理システムに用いられる。

【0017】

情報処理装置10は、通信部11、記憶部12および処理部13を有する。

通信部11は、ルータやスイッチなどの有線通信装置に接続される有線通信インタフェースでもよいし、基地局やアクセスポイントなどの無線通信装置に接続される無線通信インタフェースでもよい。記憶部12は、RAM(Random Access Memory)などの揮発性の半導体メモリでもよいし、HDD(Hard Disk Drive)やフラッシュメモリなどの不揮発性のストレージでもよい。処理部13は、例えば、CPU(Central Processing Unit)やDSP(Digital Signal Processor)などのプロセッサである。ただし、処理部13は、ASIC(Application Specific Integrated Circuit)やFPGA(Field Programmable Gate Array)などの特定用途の電子回路を含んでもよい。プロセッサは、RAMなどのメモリ(記憶部12でもよい)に記憶されたプログラムを実行する。複数のプロセッサの集合を「マルチプロセッサ」または単に「プロセッサ」と言うこともある。

【0018】

通信部11は、他の情報処理装置からの要求に応じてコネクションを確立することができる。コネクションは、例えば、LAN(Local Area Network)のTCPコネクションやInfiniBandのコネクションなどである。通信部11は、例えば、プロトコルに規定される所定の手順に従って、2以上のメッセージの交換を通じてコネクションを確立する。通信部11は、確立されたコネクションを用いて他の情報処理装置との間で信頼性の高いデータ通信を行うことができる。例えば、通信部11は、あるコネクションにおいてデータ読み出し要求を受信し、要求されたデータを当該コネクションで返信する。また、例えば、通信部11は、あるコネクションにおいてデータ書き込み要求を受信し、書き込み結果を当該コネクションで返信する。

【0019】

情報処理装置10は、同時に複数のコネクションを維持することができ、複数の他の情報処理装置との間でコネクションを維持することもできる。例えば、情報処理装置10は

10

20

30

40

50

、並列処理システムに含まれる複数のノードからの要求に応じて複数のコネクションを確立し、それら複数のコネクションを同時に維持してもよい。情報処理装置10は、維持しているコネクションそれぞれに対して、RAM領域などのメモリ領域を消費することになる。例えば、情報処理装置10は、維持しているコネクションそれぞれについて、通信部11が受信するメッセージを一時的に保持する受信バッファや、通信部11が送信するメッセージを一時的に保持する送信バッファを用意する。

#### 【0020】

維持しているコネクションが多いほど、情報処理装置10のメモリ消費量も多くなる。情報処理装置10が維持可能なコネクションの数は、情報処理装置10が有するメモリ領域の総量に依存することがある。例えば、通信部11が他の情報処理装置からコネクション確立要求を受信したとき、追加しようとするコネクションに対応するメモリ領域を確保できなければコネクション確立に失敗することがある。その場合、通信部11は当該他の情報処理装置に対してエラーメッセージを返信することがある。

10

#### 【0021】

確立されたコネクションは切断されることがある。通信部11は、例えば、プロトコルに規定される所定の手順に従って、2以上のメッセージの交換を通じてコネクションを切断する。コネクションの切断は、他の情報処理装置からの要求に応じて実行されてもよいし、情報処理装置10の判断によって実行されてもよい。例えば、情報処理装置10は、維持しているコネクションのうち最後にデータ通信が行われてから所定時間（例えば、数十分程度）経過したコネクション、すなわち、無通信時間が閾値を超えたコネクションを切断してもよい。コネクションを切断すると、情報処理装置10は当該切断したコネクションに対応するメモリ領域を解放することができる。

20

#### 【0022】

記憶部12は、コネクション情報15を記憶する。コネクション情報15は、情報処理装置10に維持されているコネクションおよびその使用状況を示す。コネクション情報15は、維持されているコネクションの数を特定するために用いることができる。また、コネクション情報15は、維持されているコネクションそれぞれにおいて最後にデータ通信が行われたタイミングを特定するために用いることができる。例えば、コネクション情報15は、維持されているコネクションの識別情報を、最後にデータ通信が行われた時刻の古い順に並べた連結リストを含んでもよい。また、コネクション情報15は、現在のコネクションの数を示すコネクション数情報を更に含んでもよい。ただし、上記の連結リストを用いて、維持されているコネクションの数をカウントするようにしてもよい。

30

#### 【0023】

情報処理装置10は、コネクションの追加状況や各コネクションにおける通信状況を監視し、コネクション情報15を継続的に更新する。例えば、情報処理装置10は、コネクションが追加されると、追加されたコネクションの識別情報を連結リストの末尾に追加する。また、例えば、情報処理装置10は、あるコネクションにおいてデータ通信が行われると、当該コネクションの識別情報を連結リストの末尾に移動する。この場合、連結リストの先頭の識別情報が示すコネクションは、最後にデータ通信が行われたタイミングが最も古いコネクション、すなわち、無通信時間が最も長いコネクションとなる。

40

#### 【0024】

処理部13は、記憶部12に記憶されたコネクション情報15を参照して、コネクションを切断するタイミングおよび切断すべきコネクションを判定する。ここで、処理部13は、単位時間当たり（例えば、1秒間当たり）のコネクションの増加数を算出し、算出した増加数が第1の閾値（例えば、数百程度）を超えているか否かを判定する。また、処理部13は、維持されているコネクションの総数が第2の閾値（例えば、数千～数万程度）を超えているか否かを判定する。単位時間当たりの増加数が第1の閾値を超えており、かつ、コネクション総数が第2の閾値を超えている場合、処理部13は、メモリ領域の不足を解消するために早急に既存コネクションの切断を行うことを決定する。

#### 【0025】

50

単位時間当たりの増加数が第1の閾値を超えること、すなわち、接続の急激な増加は、並列処理システムに含まれる複数のノードがデータ転送フェーズを実行するときが発生し得る。データ転送フェーズでは、ジョブの実行に使用される複数のノードがジョブの途中で一斉にデータの送信または受信を行う。データ転送フェーズには、複数のノードでプログラムの実行が終了した後、それら複数のノードが結果データを情報処理装置10に一斉に保存する「ステージアウト」が含まれ得る。また、データ転送フェーズには、複数のノードでプログラムが起動する前、それら複数のノードが入力データを情報処理装置10から一斉にコピーする「ステージイン」が含まれ得る。

#### 【0026】

単位時間当たりの増加数が第1の閾値を超えた場合、処理部13は、並列処理システムにおいてデータ転送フェーズが発生したと推定することが考えられる。データ転送フェーズでは、少なくとも一部のノードのデータ転送が正常に終了しないと、そのデータ転送フェーズ全体が完了せず、ジョブの完了が遅延することになる。データ転送フェーズの遅延は、ジョブの投入から完了までの所要時間に大きな影響を与える。データ転送フェーズが発生したときに接続総数が第2の閾値を超えている場合、すなわち、接続超過が生じている場合、情報処理装置10のメモリ領域の不足によって接続の追加や接続上のデータ送信が阻害されるおそれがある。そこで、データ転送フェーズが発生したときには接続超過を早急に解消することが好ましい。

#### 【0027】

既存接続の切断を行うことを決定した場合、処理部13は、接続情報15に基づいて、維持されている接続のうち最後に通信が行われたタイミングが古い方から優先的に少なくとも1つの接続を選択する。選択する接続の数は、例えば、現在の接続総数と第2の閾値との差とする。ただし、1回に選択する接続の数を所定数としてもよい。接続情報が上記の連結リストを含む場合、処理部13は、連結リストの先頭から順に少なくとも1つの識別情報を抽出するようにしてもよい。選択される接続は、例えば、データ転送フェーズが発生したジョブとは異なる別のジョブに関する接続である。

#### 【0028】

処理部13は、選択した接続を通信部11に切断させる。処理部13は、切断した接続に対応するメモリ領域を解放することが可能となる。例えば、接続14a, 14b, 14cを含むn個の接続が維持されており、最後に通信が行われたタイミングは接続14b, 14a, 14cの順に古いとする。この場合、処理部13は、接続14bを最も優先して切断し、その次に接続14aを優先して切断し、その次に接続14cを優先して切断することになる。

#### 【0029】

情報処理装置10は、単位時間当たりの増加数に基づく切断判定とは別に、各接続の無通信時間に基づく切断判定を併せて行ってもよい。その場合、維持されている接続は、無通信時間がまだ所定時間に達していない接続となる。その場合であっても、単位時間当たりの増加数が第1の閾値を超えた場合には、接続超過を放置していると並列処理のデータ転送フェーズが遅延するおそれがあるため、無通信時間が所定時間に達するのを待たずに接続を切断することが好ましい。

#### 【0030】

第1の実施の形態の情報処理装置10によれば、維持されている接続の数と、維持されている接続それぞれについて最後に通信が行われたタイミングとが監視される。単位時間当たりの維持されている接続の増加数が第1の閾値を超えておりかつ維持されている接続の数が第2の閾値を超えている場合、最後に通信が行われたタイミングが古い方から優先的に少なくとも1つの接続が選択される。そして、選択された少なくとも1つの接続が切断される。これにより、情報処理装置10のメモリ領域の不足によって並列処理におけるデータ転送フェーズが遅延することを抑制でき、並列処理を行うジョブの完了が遅延することを抑制できる。

10

20

30

40

50

## 【 0 0 3 1 】

## [ 第 2 の実施の形態 ]

次に、第 2 の実施の形態を説明する。

図 2 は、第 2 の実施の形態の並列処理システムの例を示す図である。

## 【 0 0 3 2 】

第 2 の実施の形態の並列処理システムは、複数のノードを用いてプログラムを並列実行することができる。並列処理システムは、ネットワーク 3 0、データサーバ 1 0 0、1 0 0 a、1 0 0 b を含む複数のデータサーバ、ノード 2 0 0、2 0 0 a、2 0 0 b、2 0 0 c を含む複数のノードおよびジョブ管理サーバ 3 0 0 を有する。

## 【 0 0 3 3 】

この並列処理システムは、例えば、データセンタなどの情報処理施設に設置される。ネットワーク 3 0 は、例えば、情報処理施設の内部でデータ通信を行う有線ローカルネットワークである。ネットワーク 3 0 は、LAN や InfiniBand を含んでもよい。複数のデータサーバ、複数のノードおよびジョブ管理サーバ 3 0 0 は、それぞれサーバコンピュータであり、ネットワーク 3 0 に接続されている。

## 【 0 0 3 4 】

データサーバ 1 0 0、1 0 0 a、1 0 0 b は、不揮発性の記憶装置を利用して、並列処理を行うジョブに関するファイルを記憶する。2 以上のデータサーバが共通の分散ファイルシステムを形成してもよい。あるジョブに関する 2 以上のファイルが異なるデータサーバに分散して記憶されてもよい。データサーバ 1 0 0、1 0 0 a、1 0 0 b は、あるジョブのプログラムが読み込む入力データを含む入力ファイルを記憶することがある。また、データサーバ 1 0 0、1 0 0 a、1 0 0 b は、あるジョブのプログラムが出力した結果データを含む出力ファイルを記憶することがある。データサーバ 1 0 0、1 0 0 a、1 0 0 b は、あるジョブで実行されるプログラムを含むプログラムファイルを記憶してもよい。

## 【 0 0 3 5 】

入力ファイルやプログラムファイルは、例えば、ジョブの開始前に予めユーザによってデータサーバ 1 0 0、1 0 0 a、1 0 0 b に格納される。出力ファイルは、例えば、プログラムの実行が終了してジョブを完了させる際に、ノード 2 0 0、2 0 0 a、2 0 0 b、2 0 0 c によってデータサーバ 1 0 0、1 0 0 a、1 0 0 b に格納される。出力ファイルは、その後、ユーザによって読み出されることがある。

## 【 0 0 3 6 】

ノード 2 0 0、2 0 0 a、2 0 0 b、2 0 0 c は、それぞれ CPU、RAM、HDD などのローカルの計算リソースを有し、これらローカルの計算リソースを用いてプログラムを並列に実行する。あるジョブの実行に使用されるノードは、ジョブ管理サーバ 3 0 0 によって割り当てられる。1 つのジョブに対して 2 以上のノードが割り当てられ得る。ノード 2 0 0、2 0 0 a、2 0 0 b、2 0 0 c は、データの保存にデータサーバ 1 0 0、1 0 0 a、1 0 0 b を利用する。データサーバ 1 0 0、1 0 0 a、1 0 0 b から見て、ノード 2 0 0、2 0 0 a、2 0 0 b、2 0 0 c はクライアントとすることができる。

## 【 0 0 3 7 】

ノード 2 0 0、2 0 0 a、2 0 0 b、2 0 0 c は、ジョブのプログラムを起動する前に、データサーバ 1 0 0、1 0 0 a、1 0 0 b から HDD などのローカルストレージに入力ファイルを一括してコピーする「ステージイン」を行うことがある。プログラムの起動前に入力ファイルをコピーしておくことで、プログラム実行途中の通信を抑制してプログラムの実行効率を向上させることができる。また、ノード 2 0 0、2 0 0 a、2 0 0 b、2 0 0 c は、ジョブのプログラムが終了した後、ローカルストレージからデータサーバ 1 0 0、1 0 0 a、1 0 0 b に出力ファイルを一括して転送する「ステージアウト」を行うことがある。転送した出力ファイルはローカルストレージから削除してよい。これにより、ジョブを完了させて当該ジョブに対するノードの割り当てを解放することができる。

## 【 0 0 3 8 】

ここで、あるジョブについてステージインやステージアウトなどの「ステージング」が

10

20

30

40

50

発生すると、当該ジョブに割り当てられた複数のノードが1以上のデータサーバに一齐にアクセスすることになる。データサーバ100, 100a, 100bとノード200, 200a, 200b, 200cとの間のファイル転送は、信頼性確保のためコネクション型通信によって行われる。このため、ステージングが発生すると、データサーバ100, 100a, 100bが維持するコネクションが急増する。

#### 【0039】

ジョブ管理サーバ300は、複数のノードを利用したジョブの実行を管理する。ジョブ管理サーバ300は、ユーザが使用する端末装置(図示せず)からジョブの要求を受け付けることがある。ジョブ管理サーバ300は、受け付けた複数のジョブのスケジューリングを行い、それら複数のジョブそれぞれに対してノードを割り当てる。

10

#### 【0040】

ジョブの要求はジョブスクリプトファイルに記載されることがある。ジョブスクリプトファイルは、ステージインコマンド、プログラム起動コマンド、ステージアウトコマンドなどを含むことがある。ステージインコマンドでは入力ファイルのパスが指定される。プログラム起動コマンドではプログラムファイルのパスが指定される。ステージアウトコマンドでは出力ファイルのパスが指定される。ジョブ管理サーバ300は、例えば、ジョブスクリプトファイルに従って、ノード200, 200a, 200b, 200cにステージインやプログラム起動やステージアウトを指示する。

#### 【0041】

図3は、データサーバのハードウェア例を示すブロック図である。

20

データサーバ100は、CPU101、RAM102、HDD103、画像信号処理部104、入力信号処理部105、媒体リーダー106および通信インタフェース107を有する。データサーバ100a, 100b、ノード200, 200a, 200b, 200cおよびジョブ管理サーバ300も、データサーバ100と同様のハードウェアを用いて実装できる。なお、CPU101は、第1の実施の形態の処理部13に対応する。RAM102またはHDD103は、第1の実施の形態の記憶部12に対応する。通信インタフェース107は、第1の実施の形態の通信部11に対応する。

#### 【0042】

CPU101は、プログラムの命令を実行するプロセッサである。CPU101は、HDD103に記憶されたプログラムやデータの少なくとも一部をRAM102にロードし、プログラムを実行する。なお、CPU101は複数のプロセッサコアを含んでもよく、データサーバ100は複数のプロセッサを有してもよく、以下で説明する処理を複数のプロセッサまたはプロセッサコアを用いて並列に実行してもよい。また、複数のプロセッサの集合を「マルチプロセッサ」または単に「プロセッサ」と言うことがある。

30

#### 【0043】

RAM102は、CPU101が実行するプログラムやCPU101が演算に用いるデータを一時的に記憶する揮発性の半導体メモリである。なお、データサーバ100は、RAM以外の種類のメモリを備えてもよく、複数個のメモリを備えてもよい。

#### 【0044】

HDD103は、OS(Operating System)やミドルウェアなどのソフトウェアのプログラム、および、データを記憶する不揮発性の記憶装置である。なお、データサーバ100は、フラッシュメモリやSSD(Solid State Drive)などの他の種類の記憶装置を備えてもよく、複数の不揮発性の記憶装置を備えてもよい。

40

#### 【0045】

画像信号処理部104は、CPU101からの命令に従って、データサーバ100に接続されたディスプレイ111に画像を出力する。ディスプレイ111としては、CRT(Cathode Ray Tube)ディスプレイ、液晶ディスプレイ(LCD:Liquid Crystal Display)、プラズマディスプレイ、有機EL(OEL:Organic Electro-Luminescence)ディスプレイなど、任意の種類ディスプレイを用いることができる。

#### 【0046】

50

入力信号処理部 105 は、データサーバ 100 に接続された入力デバイス 112 から入力信号を取得し、CPU 101 に出力する。入力デバイス 112 としては、マウス・タッチパネル・タッチパッド・トラックボールなどのポインティングデバイス、キーボード、リモートコントローラ、ボタンスイッチなどを用いることができる。また、データサーバ 100 に、複数の種類の入力デバイスが接続されていてもよい。

【0047】

媒体リーダ 106 は、記録媒体 113 に記録されたプログラムやデータを読み取る読み取り装置である。記録媒体 113 として、例えば、磁気ディスク、光ディスク、光磁気ディスク (MO : Magneto-Optical disk)、半導体メモリなどを使用できる。磁気ディスクには、フレキシブルディスク (FD : Flexible Disk) や HDD が含まれる。光ディスクには、CD (Compact Disc) や DVD (Digital Versatile Disc) が含まれる。

10

【0048】

媒体リーダ 106 は、例えば、記録媒体 113 から読み取ったプログラムやデータを、RAM 102 や HDD 103 などの他の記録媒体にコピーする。読み取られたプログラムは、例えば、CPU 101 によって実行される。なお、記録媒体 113 は可搬型記録媒体であってもよく、プログラムやデータの配布に用いられることがある。また、記録媒体 113 や HDD 103 を、コンピュータ読み取り可能な記録媒体とすることがある。

【0049】

通信インタフェース 107 は、ネットワーク 30 に接続され、ネットワーク 30 を介してノード 200, 200a, 200b, 200c と通信を行うインタフェースである。通信インタフェース 107 は、スイッチやルータなどの通信装置とケーブルで接続される。

20

【0050】

図 4 は、ノードとデータサーバのソフトウェア例を示すブロック図である。

データサーバ 100 は、共有ファイル記憶部 121、管理情報記憶部 122、接続メモリ領域 123、分散ファイルシステム (分散 FS) データ処理部 124 および通信制御部 125 を有する。共有ファイル記憶部 121 は、例えば、HDD 103 を用いて実装される。管理情報記憶部 122 は、例えば、RAM 102 または HDD 103 を用いて実装される。接続メモリ領域 123 は、例えば、RAM 102 を用いて実装される。分散 FS データ処理部 124 は、例えば、CPU 101 が実行するミドルウェアプログラムを用いて実装される。通信制御部 125 は、例えば、CPU 101 が実行するデバイスドライバを用いて実装される。データサーバ 100a, 100b などの他のデータサーバも、データサーバ 100 と同様のユニットによって実現できる。

30

【0051】

共有ファイル記憶部 121 は、分散ファイルシステムのディレクトリ体系のもとで管理されるファイルを記憶する。共有ファイル記憶部 121 に記憶されるファイルは、分散ファイルシステムのファイルパスによって識別される。共有ファイル記憶部 121 には、ジョブに割り当てられた複数のノードによって使用される 1 以上の入力ファイルや、これら複数のノードが生成した複数の出力ファイルが格納され得る。

【0052】

管理情報記憶部 122 は、複数のノードとの間に確立された接続の管理に用いられる管理情報を記憶する。後述するように、あるジョブのステージインやステージアウトの際に、データサーバ 100 は、当該ジョブに使用される複数のノードからの要求に応じて一斉に複数の新たな接続を確立することがある。その際、データサーバ 100 は、管理情報を参照して既存の古い接続を切断することがある。

40

【0053】

接続メモリ領域 123 は、データサーバ 100 が維持する接続それぞれに対して割り当てられたメモリ領域を含む。ある接続に対応するメモリ領域は、当該接続を用いて受信したメッセージを一時的に保持する受信バッファを含む。また、ある接続に対応するメモリ領域は、当該接続を用いて送信しようとするメッセージを一時的に保持する送信バッファを含む。また、メモリ領域は、メッ

50

ページの順序管理に用いられるシーケンス番号などの各種の情報を記憶する。

【 0 0 5 4 】

分散 F S データ処理部 1 2 4 は、ファイルアクセス要求を処理する。分散 F S データ処理部 1 2 4 は、ファイルパスを含む読み出し要求に対して、ファイルパスが示すファイルデータを共有ファイル記憶部 1 2 1 から読み出し、読み出したファイルデータをファイルアクセス結果として返信する。また、分散 F S データ処理部 1 2 4 は、ファイルパスおよびファイルデータを含む書き込み要求に対して、ファイルパスに対応付けてファイルデータを共有ファイル記憶部 1 2 1 に書き込む。分散 F S データ処理部 1 2 4 は、書き込み成否を示す書き込み結果をファイルアクセス結果として返信する。

【 0 0 5 5 】

また、分散 F S データ処理部 1 2 4 は、新たなコネクションの確立やコネクションを用いたファイルアクセス要求の受信を監視し、管理情報記憶部 1 2 2 に記憶された管理情報を更新する。分散 F S データ処理部 1 2 4 は、管理情報に基づいて、データサーバ 1 0 0 に維持されているコネクションの数が所定の条件を満たすか判定し、所定の条件を満たす場合には古いコネクションを切断するよう通信制御部 1 2 5 に指示する。コネクションを切断する契機および切断すべきコネクションの選択については後述する。

【 0 0 5 6 】

通信制御部 1 2 5 は、データサーバ 1 0 0 と複数のノードとの間のコネクション型通信を制御する。通信制御部 1 2 5 は、あるノードからコネクション確立要求を受信すると、新たなコネクションに対応するメモリ領域をコネクションメモリ領域 1 2 3 に確保して、コネクション確立応答を返信する。また、通信制御部 1 2 5 は、新たなコネクションにコネクション識別子を付与し、当該コネクションを用いた通信を以降行えるようにコネクション識別子を分散 F S データ処理部 1 2 4 に通知する。

【 0 0 5 7 】

また、通信制御部 1 2 5 は、あるコネクションの受信バッファにファイルアクセス要求が到着したとき、分散 F S データ処理部 1 2 4 が当該ファイルアクセス要求を処理できるように制御する。また、通信制御部 1 2 5 は、あるコネクションについて分散 F S データ処理部 1 2 4 がファイルアクセス結果を生成したとき、当該コネクションの送信バッファからファイルアクセス結果が送信されるように制御する。

【 0 0 5 8 】

また、通信制御部 1 2 5 は、分散 F S データ処理部 1 2 4 からコネクション識別子を指定してコネクション切断要求を受け付けると、コネクション識別子が示すコネクションの相手ノードに対してコネクション切断要求を送信する。通信制御部 1 2 5 は、相手ノードからコネクション切断応答を受信すると、そのコネクションに対応するメモリ領域をコネクションメモリ領域 1 2 3 から解放する。ただし、通信制御部 1 2 5 は、コネクション切断要求を送信してから一定時間以内にコネクション切断応答を受信しない場合、コネクションが切断されたものとみなしてメモリ領域を解放してもよい。なお、コネクション確立要求、コネクション確立応答、コネクション切断要求およびコネクション切断応答の伝送は、例えば、コネクションレス型通信によって行われる。

【 0 0 5 9 】

ノード 2 0 0 は、ローカルファイル記憶部 2 2 1、コネクションメモリ領域 2 2 2、ステージング部 2 2 3、分散 F S クライアント部 2 2 4 および通信制御部 2 2 5 を有する。ローカルファイル記憶部 2 2 1 は、例えば、R A M または H D D を用いて実装される。コネクションメモリ領域 2 2 2 は、例えば、R A M を用いて実装される。ステージング部 2 2 3 および分散 F S クライアント部 2 2 4 は、例えば、C P U が実行するミドルウェアプログラムを用いて実装される。通信制御部 2 2 5 は、例えば、C P U が実行するデバイスドライバを用いて実装される。ノード 2 0 0 a、2 0 0 b、2 0 0 c などの他のノードも、ノード 2 0 0 と同様のユニットによって実現できる。

【 0 0 6 0 】

ローカルファイル記憶部 2 2 1 は、プログラムファイルのコピーや入力ファイルのコピ

10

20

30

40

50

ーを記憶する。入力ファイルは、ステージインによってデータサーバ100, 100a, 100bからコピーされる。ローカルファイル記憶部221に記憶された入力ファイルは、プログラムに対する入力として使用される。また、ローカルファイル記憶部221は、プログラムによって生成された出力ファイルを記憶する。出力ファイルは、ステージアウトによってデータサーバ100, 100a, 100bに転送される。ジョブが完了した後は、そのジョブに関するファイルはローカルファイル記憶部221から削除される。

#### 【0061】

コネクションメモリ領域222は、ノード200が維持するコネクションそれぞれに対して割り当てられたメモリ領域を含む。コネクション毎のメモリ領域は、受信バッファや送信バッファを含み、シーケンス番号などの各種の情報を記憶する。ノード200が維持するコネクションは、主にデータサーバ100, 100a, 100bとのコネクションであり、データサーバ100が維持するコネクションに比べて十分に少ない。

10

#### 【0062】

ステージング部223は、ジョブ管理サーバ300からの指示に応じて、ステージインおよびステージアウトを含むステージングを行う。ステージインでは、ステージング部223は、入力ファイルのパスを分散FSクライアント部224に指定して、入力ファイルをローカルファイル記憶部221にコピーさせる。ステージアウトでは、ステージング部223は、出力ファイルのパスを分散FSクライアント部224に指定して、ローカルファイル記憶部221に記憶された出力ファイルを転送させる。

#### 【0063】

20

分散FSクライアント部224は、分散ファイルシステムに対するファイルアクセス要求を生成する。分散FSクライアント部224は、ステージインにおいて、ステージング部223から指定されたファイルパスを含む読み出し要求を生成し、読み出し要求に対するファイルアクセス結果に含まれる入力ファイルデータをローカルファイル記憶部221に格納する。また、分散FSクライアント部224は、ステージアウトにおいて、ステージング部223から指定されたファイルパスとローカルファイル記憶部221に記憶された出力ファイルデータとを含む書き込み要求を生成する。

#### 【0064】

分散FSクライアント部224は、データサーバ100, 100a, 100bの通信アドレス（例えば、IP（Internet Protocol）アドレス）を予め知っている。分散FSクライアント部224は、コネクションが存在しないデータサーバに対してファイルアクセス要求を発行しようとする場合、宛先のデータサーバの通信アドレスを指定してコネクション確立要求を通信制御部225に通知する。コネクション確立が成功した場合、コネクション識別子が通信制御部225から分散FSクライアント部224に通知される。すると、分散FSクライアント部224は、宛先のデータサーバに対応するコネクション識別子を用いてファイルアクセス要求を発行する。

30

#### 【0065】

通信制御部225は、ノード200と1以上のデータサーバとの間のコネクション型通信を制御する。通信制御部225は、分散FSクライアント部224からコネクション確立要求を取得すると、新たなコネクションに対応するメモリ領域をコネクションメモリ領域222に確保する。また、通信制御部225は、分散FSクライアント部224から指定された通信アドレス宛てにコネクション確立要求を送信する。通信制御部225は、コネクション確立応答を受信すると、新たなコネクションにコネクション識別子を付与し、コネクション識別子を分散FSクライアント部224に通知する。なお、ここで付与されるコネクション識別子はノード200の中で一意であればよく、上記のように宛先のデータサーバで付与されるコネクション識別子とは異なってもよい。

40

#### 【0066】

また、通信制御部225は、あるコネクションについて分散FSクライアント部224がファイルアクセス要求を生成したとき、当該コネクションの送信バッファからファイルアクセス要求が送信されるように制御する。また、通信制御部225は、あるコネクショ

50

ンの受信バッファにファイルアクセス結果が到着したとき、分散FSクライアント部224が当該ファイルアクセス結果を処理できるように制御する。

【0067】

また、通信制御部225は、あるコネクションについてコネクション切断要求を受信すると、そのコネクションに対応するメモリ領域をコネクションメモリ領域222から解放し、相手データサーバに対してコネクション切断応答を返信する。

【0068】

図5は、コネクション管理リストの例を示す図である。

コネクション管理リスト131は、データサーバ100の管理情報記憶部122に記憶されている。コネクション管理リスト131は、複数の構造体をポインタで連結した連結リストである。各構造体は、コネクション識別子と1つ後ろの構造体を指す後ポインタと1つ前の構造体を指す前ポインタとを含む。ただし、先頭の構造体の前ポインタはNULL(空)であり、末尾の構造体の後ポインタはNULLである。

10

【0069】

4つのコネクションが維持されているとき、コネクション管理リスト131は構造体131a, 131b, 131c, 131dを含む。構造体131aは先頭の構造体であり、コネクション識別子#1と構造体131bを指す後ポインタ#1とを含む。構造体131bは前から2番目の構造体であり、コネクション識別子#2と構造体131cを指す後ポインタ#2と構造体131aを指す前ポインタ#2とを含む。構造体131cは前から3番目の構造体であり、コネクション識別子#3と構造体131dを指す後ポインタ#3と構造体131bを指す前ポインタ#3とを含む。構造体131dは末尾の構造体であり、コネクション識別子#4と構造体131cを指す前ポインタ#4とを含む。

20

【0070】

コネクション管理リスト131に含まれる複数の構造体は、先頭から末尾に向かってコネクションが「古い」順に並べられている。ここで言う「古い」とは、最後にファイルアクセス要求が送信された時刻が古いことであり、無通信時間が長いことに相当する。コネクション管理リスト131の先頭の構造体を示すコネクションは、維持されているコネクションのうち最も古いコネクションである。最も古いコネクションは、最後にファイルアクセス要求が送信された時刻が最も古く無通信時間が最も長いコネクションであり、LRU(Least Recently Used)コネクションとすることができる。また、コネクション管理リスト131の末尾の構造体を示すコネクションは、維持されているコネクションのうち最も新しいコネクションである。最も新しいコネクションは、最近にファイルアクセス要求が送信されており無通信時間が最も短いコネクションである。

30

【0071】

コネクション管理リスト131に含まれる複数の構造体の順序は、LRUアルゴリズムに従って管理される。新しいコネクションが確立されると、当該コネクションに対応する構造体がコネクション管理リスト131の末尾に追加される。あるコネクションでファイルアクセス要求が送信されると、当該コネクションに対応する構造体がコネクション管理リスト131の末尾に移動する。このとき、移動する構造体およびその前後の構造体のポインタが書き換えられる。例えば、構造体131bが末尾に移動する場合、構造体131bの後ポインタはNULLになり、構造体131bの前ポインタは構造体131dを指すように変更される。また、構造体131aの後ポインタは構造体131cを指すように変更され、構造体131cの前ポインタは構造体131aを指すように変更される。

40

【0072】

図6は、コネクション数テーブルの例を示す図である。

コネクション数テーブル132は、データサーバ100の管理情報記憶部122に記憶されている。コネクション数テーブル132は、現在コネクション数、前回コネクション数、増加速度、増加速度閾値、総数閾値および超過コネクション数の項目を有する。

【0073】

現在コネクション数は、データサーバ100が維持しているコネクションの数の最新値

50

である。現在コネクション数に比例する大きさのメモリ領域がコネクション維持のために消費されることになる。新しいコネクションが確立されると、確立されたコネクションの数だけ現在コネクション数が増加する。また、既存のコネクションが切断されると、切断されたコネクションの数だけ現在コネクション数が減少する。

**【 0 0 7 4 】**

前回コネクション数は、コネクション切断判定を前回行ったときの現在コネクション数である。コネクション切断判定は維持されているコネクションの一部を切断するか否かの判定であり、所定時間毎（例えば、1秒毎）に行われる。コネクション切断判定の際に、現在コネクション数が前回コネクション数として退避される。コネクション切断判定の後の所定時間の間に、現在コネクション数は前回コネクション数から変化していく。

10

**【 0 0 7 5 】**

増加速度は、所定時間におけるコネクションの増加数である。現在コネクション数から前回コネクション数を引いた値、すなわち、今回のコネクション切断判定時のコネクション数から前回のコネクション切断判定時のコネクション数を引いた値が、増加速度として算出される。ただし、現在コネクション数から前回コネクション数を引いた値が0未満である場合、増加速度を0とする。増加速度閾値は、増加速度と比較される閾値であり、予めデータサーバ100の管理者によって設定される。増加速度閾値は、例えば、数百程度に設定される。増加速度が増加速度閾値を超える場合、ステージインやステージアウトなどのステージングによってコネクションが急増していると推定される。

**【 0 0 7 6 】**

20

総数閾値は、現在コネクション数と比較される閾値であり、予めデータサーバ100の管理者によって設定される。総数閾値は、例えば、数千～数万程度に設定される。現在コネクション数が総数閾値を超える場合、メモリ領域の不足によってコネクション確立を拒否する可能性があるとして推定される。超過コネクション数は、総数閾値を超えるコネクション数であり、現在コネクション数から総数閾値を引いて算出される。ただし、現在コネクション数から総数閾値を引いた値が0未満である場合、超過コネクション数を0とする。

**【 0 0 7 7 】**

あるジョブでステージングが実行されているときにデータサーバ100でメモリ領域の不足が発生すると、そのジョブに属する一部のノードがコネクション確立を待たされてしまい、ステージングの完了が遅延するおそれがある。特に、ステージアウトの完了が遅延すると、プログラムの実行が終了しているにもかかわらずノードを解放することができず、他のジョブのスケジュールにも影響を与えてしまう。そこで、データサーバ100は、増加速度が増加速度閾値を超えており、かつ、現在コネクション数が総数閾値を超えている場合、コネクション超過を迅速に解消することとする。このとき、データサーバ100は、超過コネクション数だけ古いコネクションを切断する。

30

**【 0 0 7 8 】**

なお、データサーバ100は、後述するコネクション切断判定とは別に、無通信時間が所定時間（例えば、数十分程度）を超えるコネクションを自動的に切断するようにしてもよい。その場合、コネクション数テーブル132に記載される現在コネクション数は、無通信時間がまだ所定時間を超えていないコネクションの数となる。その場合であっても、データサーバ100は、ステージングの完了が遅延することを抑制するため、無通信時間が所定時間を超えるのを待たずに相対的に古いコネクションを切断する。無通信時間が所定時間を超える前に切断されたコネクションについては、相手ノードがまだ通信を行おうとしていた可能性もある。その場合には、相手ノードは改めてデータサーバ100に対してコネクション確立要求を送信することになる。このように、ステージングが発生した場合には、そのステージングの迅速な完了が優先される。

40

**【 0 0 7 9 】**

無通信時間が所定時間を超えるコネクションの切断は、分散FSデータ処理部124が行ってもよい。管理情報記憶部122に記憶される管理情報には、各コネクションについて最後に通信が行われた時刻またはその時刻からの経過時間が含まれてもよい。

50

## 【 0 0 8 0 】

次に、並列処理システムの処理手順について説明する。

図 7 は、ジョブ実行の手順例を示すフローチャートである。

( S 1 0 ) ジョブ管理サーバ 3 0 0 は、並列処理システムが有するノードの中から 1 以上の空きノードを選択してジョブに割り当て、ジョブを開始させる。

## 【 0 0 8 1 】

( S 1 1 ) ジョブ管理サーバ 3 0 0 は、ジョブスクリプトファイルなどに従って、割り当てた各ノードにステージインを実行させる。各ノードは、入力ファイルを 1 以上のデータサーバからノード内のローカルストレージにコピーする。ただし、ジョブによっては入力ファイルが存在せずステージインが実行されないこともある。

10

## 【 0 0 8 2 】

( S 1 2 ) ジョブ管理サーバ 3 0 0 は、ジョブスクリプトファイルなどに従って、割り当てた各ノードにプログラムを起動させる。各ノードはプログラムを起動する。

( S 1 3 ) 各ノードはプログラムの実行を終了する。ジョブ管理サーバ 3 0 0 は、ジョブに割り当てた全てのノードがプログラムの実行を終了したことを検出する。

## 【 0 0 8 3 】

( S 1 4 ) ジョブ管理サーバ 3 0 0 は、ジョブスクリプトファイルなどに従って、割り当てた各ノードにステージアウトを実行させる。各ノードは、出力ファイルをローカルストレージから 1 以上のデータサーバに移動させる。

## 【 0 0 8 4 】

( S 1 5 ) ジョブ管理サーバ 3 0 0 は、ジョブへのノードの割り当てを解放する。これによりジョブが完了し、解放されたノードは空きノードとなる。

20

図 8 は、コネクション確立の手順例を示すフローチャートである。

## 【 0 0 8 5 】

コネクション確立は、ステージインやステージアウトの開始時にノード毎に実行される。ここでは、ノード 2 0 0 がデータサーバ 1 0 0 に接続する場合を考える。

( S 2 0 ) 分散 F S クライアント部 2 2 4 は、データサーバ 1 0 0 の通信アドレスを指定してコネクション確立要求を通信制御部 2 2 5 に対して発行する。

## 【 0 0 8 6 】

( S 2 1 ) 通信制御部 2 2 5 は、新たなコネクションに対応するメモリ領域をコネクションメモリ領域 2 2 2 の中 ( R A M の中 ) に確保する。

30

( S 2 2 ) 通信制御部 2 2 5 は、データサーバ 1 0 0 にコネクション確立要求を送信する。コネクション確立要求は、例えば、コネクションレス型通信として送信される。

## 【 0 0 8 7 】

( S 2 3 ) 通信制御部 1 2 5 は、コネクション確立要求を受信する。

( S 2 4 ) 通信制御部 1 2 5 は、新たなコネクションに対応するメモリ領域をコネクションメモリ領域 1 2 3 の中 ( R A M 1 0 2 の中 ) に確保する。

## 【 0 0 8 8 】

( S 2 5 ) 通信制御部 1 2 5 は、ステップ S 2 4 でメモリ領域の確保に成功したか判断する。メモリ領域の確保に成功した場合、ステップ S 2 8 に処理が進む。メモリ領域不足などにより確保に失敗した場合、ステップ S 2 6 に処理が進む。

40

## 【 0 0 8 9 】

( S 2 6 ) 通信制御部 1 2 5 は、エラーメッセージをノード 2 0 0 に送信する。

( S 2 7 ) 通信制御部 2 2 5 は、エラーメッセージを受信する。そして、ステップ S 2 2 に処理が進む。通信制御部 2 2 5 は、一定時間待ってからコネクション確立要求をデータサーバ 1 0 0 に再送してもよい。また、通信制御部 2 2 5 が分散 F S クライアント部 2 2 4 にエラーを通知し、分散 F S クライアント部 2 2 4 が一定時間待ってからコネクション確立要求を通信制御部 2 2 5 に対して再発行してもよい。

## 【 0 0 9 0 】

( S 2 8 ) 通信制御部 1 2 5 は、新たなコネクションに対してコネクション識別子を付

50

与する。このコネクション識別子はデータサーバ100の中で重複していなければよく、ノード200とは独立に決めてよい。通信制御部125は、付与したコネクション識別子を分散FSデータ処理部124に通知する。

【0091】

(S29) 分散FSデータ処理部124は、通信制御部125から通知されたコネクション識別子を含む構造体を生成し、生成した構造体を管理情報記憶部122に記憶されたコネクション管理リスト131の末尾に追加する。

【0092】

(S30) 分散FSデータ処理部124は、管理情報記憶部122に記憶されたコネクション数テーブル132の現在コネクション数を1だけ増加させる。

10

(S31) 通信制御部125は、コネクション確立応答をノード200に送信する。コネクション確立応答は、例えば、コネクションレス型通信として送信される。ステップS29, S30のコネクション管理リスト131およびコネクション数テーブル132の更新とコネクション確立応答の送信とは、非同期に行ってもよい。

【0093】

(S32) 通信制御部225は、コネクション確立応答を受信する。

(S33) 通信制御部225は、新たなコネクションに対してコネクション識別子を付与する。このコネクション識別子はノード200の中で重複していなければよく、データサーバ100とは独立に決めてよい。通信制御部225は、付与したコネクション識別子を分散FSクライアント部224に通知する。

20

【0094】

図9は、ファイルアクセスの手順例を示すフローチャートである。

ファイルアクセスは、ステージインやステージアウトの中でノード毎に実行され得る。ここでは、ノード200がデータサーバ100にアクセスする場合を考える。

【0095】

(S40) 分散FSクライアント部224は、コネクション識別子を指定して、読み出し要求または書き込み要求を示すファイルアクセス要求を発行する。

(S41) 通信制御部225は、指定されたコネクション識別子が示すコネクションを用いて、ファイルアクセス要求をデータサーバ100に送信する。

【0096】

30

(S42) 通信制御部125は、ファイルアクセス要求を受信する。通信制御部125は、ファイルアクセス要求を受信したコネクションを示すコネクション識別子とファイルアクセス要求とを分散FSデータ処理部124に通知する。

【0097】

(S43) 分散FSデータ処理部124は、管理情報記憶部122に記憶されたコネクション管理リスト131の中から、通知されたコネクション識別子を含む構造体を検索し、検索した構造体をコネクション管理リスト131の末尾に移動させる。

【0098】

(S44) 分散FSデータ処理部124は、ファイルアクセス要求に対応するファイルアクセス処理を実行する。ファイルアクセス要求が読み出し要求である場合、分散FSデータ処理部124は、要求されたファイルデータを共有ファイル記憶部121から読み出し、読み出したファイルデータを含むファイルアクセス結果を生成する。ファイルアクセス要求が書き込み要求である場合、分散FSデータ処理部124は、書き込み要求に含まれるファイルデータを共有ファイル記憶部121に書き込み、書き込み成否を示すファイルアクセス結果を生成する。分散FSデータ処理部124は、ファイルアクセス要求の受信時のコネクション識別子を指定してファイルアクセス結果を出力する。

40

【0099】

(S45) 通信制御部125は、指定されたコネクション識別子が示すコネクションを用いて、ファイルアクセス結果をノード200に送信する。

(S46) 通信制御部225は、ファイルアクセス結果を受信する。通信制御部225

50

は、ファイルアクセス結果を受信したコネクションを示すコネクション識別子とファイルアクセス結果とを分散FSクライアント部224に通知する。

【0100】

図10は、コネクション切断の手順例を示すフローチャートである。

コネクション切断は、データサーバ側の判断で適宜実行される。ここでは、データサーバ100がノード200との間のコネクションを切断する場合を考える。

【0101】

(S50)分散FSデータ処理部124は、切断するコネクションを選択する。分散FSデータ処理部124は、選択したコネクションのコネクション識別子を指定してコネクション切断要求を通信制御部125に対して発行する。なお、切断するコネクションの選択を含むコネクション切断判定は、分散FSデータ処理部124によって定期的に行われる。コネクション切断判定については後述する。

10

【0102】

(S51)通信制御部125は、ノード200にコネクション切断要求を送信する。コネクション切断要求は、例えば、コネクションレス型通信として送信される。

(S52)通信制御部225は、コネクション切断要求を受信する。

【0103】

(S53)通信制御部225は、切断するコネクションに対応するメモリ領域をコネクションメモリ領域222(RAM)から解放する。

(S54)通信制御部225は、コネクション切断応答をデータサーバ100に送信する。コネクション切断応答は、例えば、コネクションレス型通信として送信される。

20

【0104】

(S55)通信制御部125は、コネクション切断応答を受信する。

(S56)通信制御部125は、切断するコネクションに対応するメモリ領域をコネクションメモリ領域123(RAM102)から解放する。

【0105】

図11は、コネクション切断判定の手順例を示すフローチャートである。

コネクション切断判定は、分散FSデータ処理部124によって反復実行される。

(S60)分散FSデータ処理部124は、単位時間(例えば、1秒)だけ待つ。

【0106】

(S61)分散FSデータ処理部124は、管理情報記憶部122に記憶されたコネクション数テーブル132から現在コネクション数と前回コネクション数を読み出す。分散FSデータ処理部124は、増加速度 = 現在コネクション数 - 前回コネクション数を算出し、算出した増加速度をコネクション数テーブル132に記録する。ただし、現在コネクション数 - 前回コネクション数が負の値である場合は増加速度を0とする。

30

【0107】

(S62)分散FSデータ処理部124は、コネクション数テーブル132に記録された現在コネクション数を前回コネクション数の項目にコピーする。

(S63)分散FSデータ処理部124は、コネクション数テーブル132から増加速度閾値を読み出し、ステップS61で算出した増加速度と増加速度閾値とを比較する。増加速度が増加速度閾値を超える場合、分散FSデータ処理部124は1以上のジョブでステージングが開始されたと推定し、ステップS64に処理が進む。増加速度が増加速度閾値以下である場合、今回のコネクション切断判定が終了する。

40

【0108】

(S64)分散FSデータ処理部124は、コネクション数テーブル132から総数閾値を読み出し、現在コネクション数と総数閾値とを比較する。現在コネクション数が総数閾値を超える場合、分散FSデータ処理部124はコネクション超過によりメモリ領域が不足するおそれがあると判定し、ステップS65に処理が進む。現在コネクション数が総数閾値以下である場合、今回のコネクション切断判定が終了する。

【0109】

50

(S 6 5) 分散 F S データ処理部 1 2 4 は、超過コネクション数 = 現在コネクション数 - 総数閾値を算出し、算出した超過コネクション数をコネクション数テーブル 1 3 2 に記録する。

【 0 1 1 0 】

(S 6 6) 分散 F S データ処理部 1 2 4 は、管理情報記憶部 1 2 2 に記憶されたコネクション管理リスト 1 3 1 の先頭から順に、すなわち、最後に通信が行われた時刻が古い順に、超過コネクション数だけコネクション識別子を抽出する。抽出したコネクション識別子を含む構造体はコネクション管理リスト 1 3 1 から削除してよい。

【 0 1 1 1 】

(S 6 7) 分散 F S データ処理部 1 2 4 は、コネクション数テーブル 1 3 2 に記録された現在コネクション数を超過コネクション数だけ減少させる。

10

(S 6 8) 分散 F S データ処理部 1 2 4 は、ステップ S 6 6 で抽出したコネクション識別子それぞれについてコネクション切断要求を発行する。発行されたコネクション切断要求それぞれに対して、図 1 0 のコネクション切断処理が実行される。

【 0 1 1 2 】

第 2 の実施の形態の並列処理システムによれば、データサーバの現在コネクション数とそれらコネクションの間で最後に通信が行われた時刻の相対的な古さとが管理される。コネクションの増加速度が閾値を超えた場合、1 以上のジョブでステージインやステージアウトなどのステージングが開始されたと推定される。ステージングが開始されたときにデータサーバでコネクション超過が生じていれば、最後に通信が行われた時刻が相対的に古い方から優先的にコネクションを切断することでコネクション超過が解消される。これにより、データサーバのメモリ領域不足によってステージングを行うノードがコネクション確立に失敗することを抑制できる。よって、ジョブのステージングが遅延することを抑制でき、ジョブの完了が遅延することを抑制できる。

20

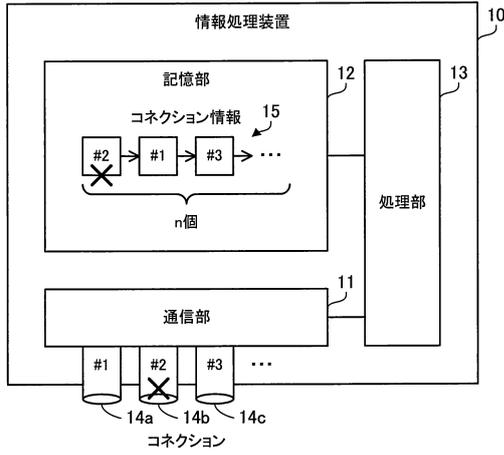
【 符号の説明 】

【 0 1 1 3 】

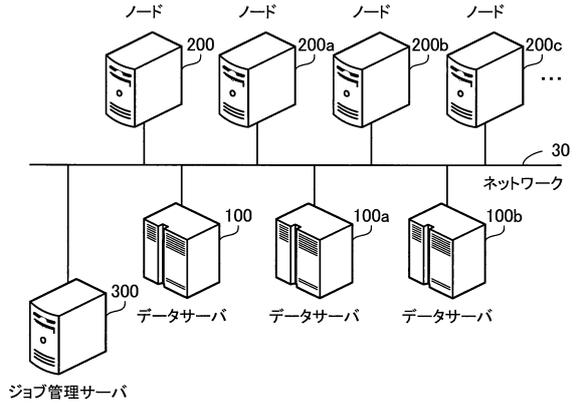
- 1 0 情報処理装置
- 1 1 通信部
- 1 2 記憶部
- 1 3 処理部
- 1 4 a , 1 4 b , 1 4 c コネクション
- 1 5 コネクション情報

30

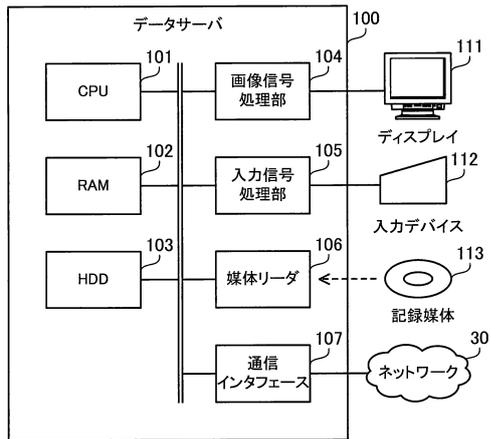
【図1】



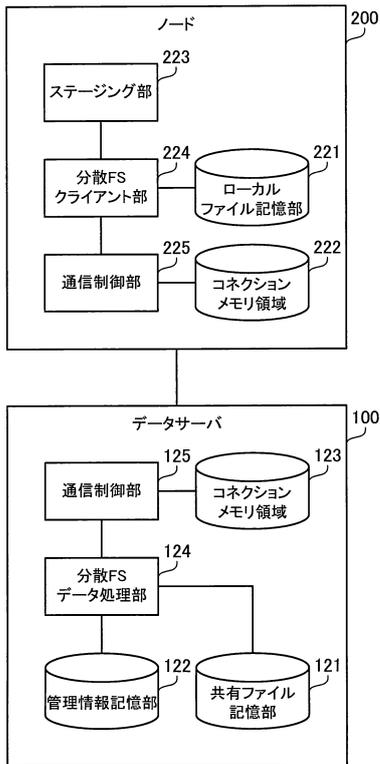
【図2】



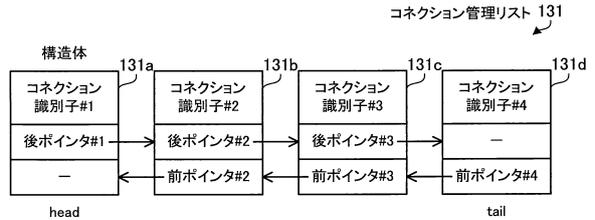
【図3】



【図4】



【図5】

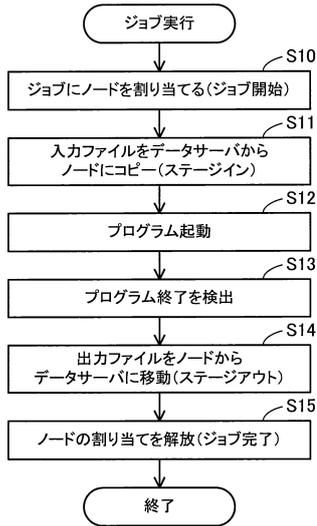


【図6】

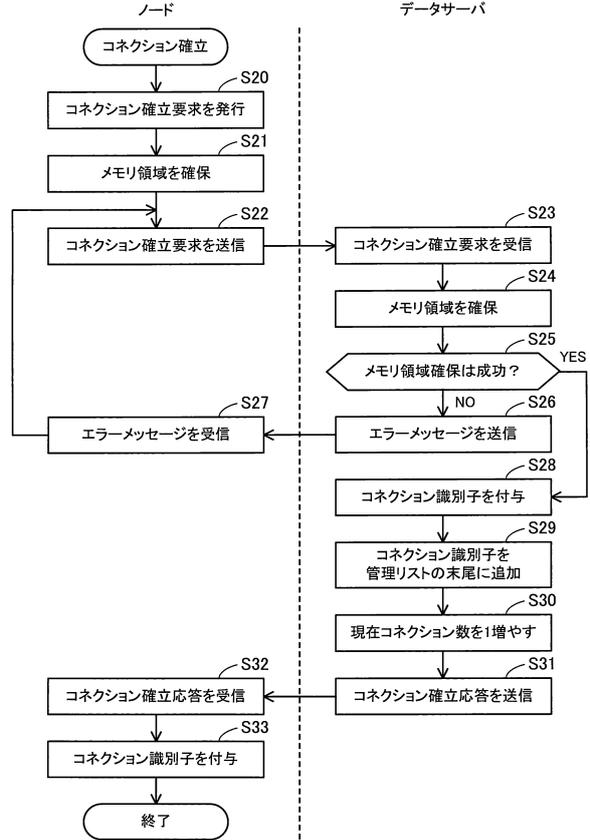
コネクション数テーブル

|           |      |
|-----------|------|
| 現在コネクション数 | 5050 |
| 前回コネクション数 | 4900 |
| 増加速度      | 150  |
| 増加速度閾値    | 100  |
| 総数閾値      | 5000 |
| 超過コネクション数 | 50   |

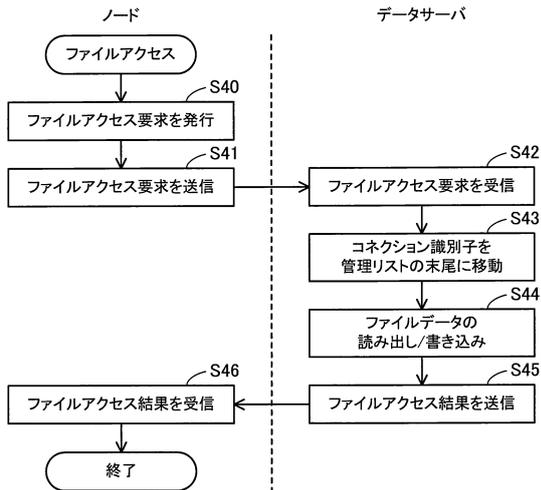
【図7】



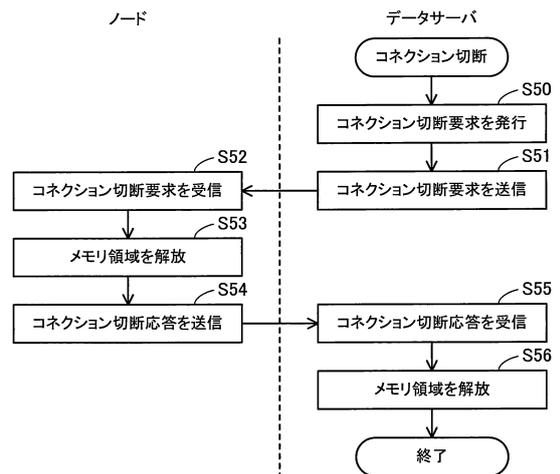
【図8】



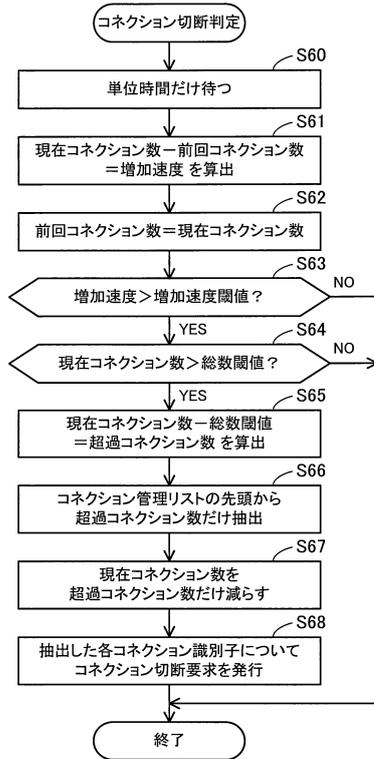
【図9】



【図10】



【図 11】



---

フロントページの続き

- (56)参考文献 特開2008-250505(JP,A)  
特開2002-158707(JP,A)  
特開平08-328986(JP,A)  
米国特許出願公開第2002/0062374(US,A1)  
中国特許出願公開第1909516(CN,A)  
特開2007-043281(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G06F 13/00