

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6251417号  
(P6251417)

(45) 発行日 平成29年12月20日 (2017.12.20)

(24) 登録日 平成29年12月1日 (2017.12.1)

(51) Int.Cl.	F I				
<b>G06F 3/06 (2006.01)</b>	G06F	3/06	301U		
<b>G06F 13/38 (2006.01)</b>	G06F	3/06	301S		
<b>G06F 13/12 (2006.01)</b>	G06F	13/38	310A		
<b>G06F 9/52 (2006.01)</b>	G06F	13/12	340B		
<b>G06F 12/00 (2006.01)</b>	G06F	9/46	472Z		
請求項の数 10 (全 25 頁) 最終頁に続く					

(21) 出願番号	特願2016-556061 (P2016-556061)	(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
(86) (22) 出願日	平成26年10月27日 (2014.10.27)	(74) 代理人	110000279 特許業務法人ウィルフォート国際特許事務所
(86) 国際出願番号	PCT/JP2014/078497	(72) 発明者	高田 有時 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
(87) 国際公開番号	W02016/067339	審査官	田中 啓介
(87) 国際公開日	平成28年5月6日 (2016.5.6)		
審査請求日	平成29年3月1日 (2017.3.1)		
最終頁に続く			

(54) 【発明の名称】 ストレージシステム、及び、記憶制御方法

(57) 【特許請求の範囲】

【請求項1】

複数プロセッサを含んだコントローラと、  
記憶デバイスが接続されるインターフェイスデバイスと、  
前記インターフェイスデバイスに関連付けられており専有キューと共有キューとを含んだ複数のキューと、

を有し、

前記複数のキューの各々には、そのキューに割り当てられているプロセッサから前記インターフェイスデバイスへ送られるデータが格納され、前記複数のキューの各々から前記インターフェイスデバイスにデータが送られるようになっており、

前記専有キューは、1つのプロセッサだけが割り当てられているキューであり、データの格納時に排他処理が不要なキューであり、

前記共有キューは、複数の前記プロセッサが割り当てられているキューであり、データの格納時に排他処理が必要なキューであり、

前記コントローラは、

前記専有キューに割り当てられたプロセッサと、前記共有キューに割り当てられたプロセッサと、の割当てを切替えた場合に、前記共有キューにデータが格納される共有頻度が所定の閾値以上小さくなる場合に、前記割当ての切替えを実行する  
ストレージシステム。

【請求項2】

前記コントローラは、

前記複数のプロセッサがそれぞれ前記インターフェイスデバイスへデータを発行したデータ発行頻度に基づいて、前記専有キューに割り当てるプロセッサを決定する請求項 1 記載のストレージシステム。

【請求項 3】

複数プロセッサを含んだコントローラと、

記憶デバイスが接続されるインターフェイスデバイスと、

前記インターフェイスデバイスに関連付けられており専有キューと共有キューとを含んだ複数のキューと、

を有し、

前記複数のキューの各々には、そのキューに割り当てられているプロセッサから前記インターフェイスデバイスへ送られるデータが格納され、前記複数のキューの各々から前記インターフェイスデバイスにデータが送られるようになっており、

前記専有キューは、1つのプロセッサだけが割り当てられているキューであり、データの格納時に排他処理が不要なキューであり、

前記共有キューは、複数の前記プロセッサが割り当てられているキューであり、データの格納時に排他処理が必要なキューであり、

前記コントローラは、

前記複数のキューの内の利用負荷指標が最大のキューである第 1 のキューについて、そのキューのキュー種別を、共有から専有に変更するか否か判定するようになっており、

前記利用負荷指標は、キューにおける 1 のプロセッサあたりの利用負荷の大きさを示す指標であって、そのキューに割り当てられているプロセッサ数が増えると小さくなり、各プロセッサからそのキューにデータが格納された頻度が増えると大きくなる値である

ストレージシステム。

【請求項 4】

前記コントローラは、

前記第 1 のキューに割り当てられている複数のプロセッサの内、前記第 1 のキューにデータを格納する頻度が最小のプロセッサを、前記複数のキューの内の前記利用負荷指標が最小のキューである第 2 のキューに割り当てる割当て切替えを実行したと仮定した場合の前記第 1 のキューの利用負荷指標及び前記第 2 のキューの利用負荷指標をそれぞれ推定し、

前記推定された第 1 のキューの利用負荷指標が、前記推定された第 2 のキューの利用負荷指標よりも大きい場合、前記第 2 のキューに前記最小のプロセッサを割り当てる割当て切替えを実行する

請求項 3 記載のストレージシステム。

【請求項 5】

前記利用負荷指標は、キューに割り当てられているプロセッサ数の逆数と、各プロセッサからそのキューにデータが格納された頻度の合計との積に基づく値である

請求項 4 記載のストレージシステム。

【請求項 6】

複数プロセッサを含んだコントローラと、

記憶デバイスが接続されるインターフェイスデバイスと、

前記インターフェイスデバイスに関連付けられており専有キューと共有キューとを含んだ複数のキューと、

を有し、

前記複数のキューの各々には、そのキューに割り当てられているプロセッサから前記インターフェイスデバイスへ送られるデータが格納され、前記複数のキューの各々から前記インターフェイスデバイスにデータが送られるようになっており、

前記専有キューは、1つのプロセッサだけが割り当てられているキューであり、データ

10

20

30

40

50

の格納時に排他処理が不要なキューであり、

前記共有キューは、複数の前記プロセッサが割り当てられているキューであり、データの格納時に排他処理が必要なキューであり、

前記複数のプロセッサのそれぞれでプロセスを実行することによりデータが発行されるようになっており、

前記コントローラは、

前記共有キューにデータが格納される頻度である共有頻度及び前記専有キューに割り当てられたプロセッサの負荷に関する所定の条件が満たされた場合、前記共有キューに割り当てられたいずれかのプロセッサで実行されているプロセスを、前記専有キューに割り当てられたプロセッサに移動させるプロセス移動を実行する

ストレージシステム。

【請求項 7】

前記所定の条件とは、前記プロセスを移動させた場合に前記共有キューにデータが格納される共有頻度が所定の閾値以上小さくなり、且つ、前記専有キューに割り当てられたプロセッサの負荷が所定の閾値未満である

請求項 6 記載のストレージシステム。

【請求項 8】

ストレージシステムにおけるデータの記憶制御方法であって、

前記ストレージシステムは、複数のプロセッサを含んだコントローラと、記憶デバイスが接続されるインターフェイスデバイスと、前記インターフェイスデバイスに関連付けられており専有キューと共有キューを含んだ複数のキューと、を有し、

前記専有キューは、1つのプロセッサだけが割り当てられているキューであり、

前記共有キューは、複数の前記プロセッサが割り当てられているキューであり、

前記コントローラは、

前記インターフェイスデバイスへ送られるデータが前記専有キューに割り当てられたプロセッサから発行されたものである場合、排他処理を行わずに前記データを前記専有キューに格納し、

前記インターフェイスデバイスへ送られるデータが前記共有キューに割り当てられたプロセッサから発行されたものである場合、排他処理を行って前記データを前記共有キューに格納し、

前記専有キューに割り当てられたプロセッサと、前記共有キューに割り当てられたプロセッサと、の割当てを切替えた場合に、前記共有キューにデータが格納される共有頻度が所定の閾値以上小さくなる場合に、前記割当ての切替えを実行する  
記憶制御方法。

【請求項 9】

ストレージシステムにおけるデータの記憶制御方法であって、

前記ストレージシステムは、複数のプロセッサを含んだコントローラと、記憶デバイスが接続されるインターフェイスデバイスと、前記インターフェイスデバイスに関連付けられており専有キューと共有キューを含んだ複数のキューと、を有し、

前記専有キューは、1つのプロセッサだけが割り当てられているキューであり、

前記共有キューは、複数の前記プロセッサが割り当てられているキューであり、

前記コントローラは、

前記インターフェイスデバイスへ送られるデータが前記専有キューに割り当てられたプロセッサから発行されたものである場合、排他処理を行わずに前記データを前記専有キューに格納し、

前記インターフェイスデバイスへ送られるデータが前記共有キューに割り当てられたプロセッサから発行されたものである場合、排他処理を行って前記データを前記共有キューに格納し、

前記複数のキューの内の利用負荷指標が最大のキューである第 1 のキューについて、そのキューのキュー種別を、共有から専有に変更するか否かを判定するようになっており、

10

20

30

40

50

前記利用負荷指標は、キューにおける1のプロセッサあたりの利用負荷の大きさを示す指標であって、そのキューに割り当てられているプロセッサ数が増えると小さくなり、各プロセッサからそのキューにデータが格納された頻度が大きくなると大きくなる値である

記憶制御方法。

【請求項10】

ストレージシステムにおけるデータの記憶制御方法であって、

前記ストレージシステムは、複数のプロセッサを含んだコントローラと、記憶デバイスが接続されるインターフェイスデバイスと、前記インターフェイスデバイスに関連付けられており専有キューと共有キューを含んだ複数のキューと、を有し、

前記専有キューは、1つのプロセッサだけが割り当てられているキューであり、

前記共有キューは、複数の前記プロセッサが割り当てられているキューであり、

前記複数のプロセッサのそれぞれでプロセスを実行することによりデータが発行されるようになっており、

前記コントローラは、

前記インターフェイスデバイスへ送られるデータが前記専有キューに割り当てられたプロセッサから発行されたものである場合、排他処理を行わずに前記データを前記専有キューに格納し、

前記インターフェイスデバイスへ送られるデータが前記共有キューに割り当てられたプロセッサから発行されたものである場合、排他処理を行って前記データを前記共有キューに格納し、

前記共有キューにデータが格納される頻度である共有頻度及び前記専有キューに割り当てられたプロセッサの負荷に関する所定の条件が満たされた場合、前記共有キューに割り当てられたいずれかのプロセッサで実行されているプロセスを、前記専有キューに割り当てられたプロセッサに移動させるプロセス移動を実行する

記憶制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージシステム及び記憶制御方法の技術に関する。

【背景技術】

【0002】

ストレージシステムにおいて、プロセッサは、キューを介してI/F（インターフェイス）デバイスとデータを遣り取りする。複数のプロセッサが1つのキューを利用する構成において、プロセッサは、キューにデータを格納する際、他のプロセッサからそのキューが利用されないように、排他処理を行う必要がある。排他処理は、プロセッサの処理負荷を高めるので、できるだけ実行されないことが望ましい。そこで、複数のプロセッサを備える計算機において、I/Fデバイスに対して、複数のプロセッサの各々が専用に利用するキューを設けることにより、排他処理を不要とする技術が知られている（特許文献1、2、3）。

【先行技術文献】

【特許文献】

【0003】

【特許文献1】US 8 595 385

【特許文献2】US 8 578 106

【特許文献3】特開 2010 - 128664号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、I/Fデバイスに対して、複数のプロセッサの各々が専用に利用するキ

10

20

30

40

50

ュー（「専有キュー」という）を設ける構成の場合、プロセッサの数が多いと、その分、専有キューの数も多く、延いては、I/Fデバイスにおけるリソース（レジスタなど）の使用量も大きい。したがって、プロセッサ数が多いと、I/Fデバイスのリソースが不足し、結果として、プロセッサ毎に専有キューを設けることができない場合もある。それを避けるために、I/Fデバイスに大きなリソースを載せることが考えられるが、I/Fデバイスに大きなリソースを載せると、I/Fデバイスのコストが増加してしまう。

【0005】

そこで、本発明の目的は、I/Fデバイスのリソース使用量と排他処理の負荷の両方を抑制したストレージシステム及び記憶制御方法を提供することにある。

【課題を解決するための手段】

10

【0006】

一実施形態に係るストレージシステムは、複数のプロセッサを含んだコントローラと、記憶デバイスが接続されるインターフェイスデバイスと、インターフェイスデバイスに関連付けられている複数のキューとを有する。複数のキューの各々には、そのキューに割り当てられているプロセッサからインターフェイスデバイスに送られるデータが格納され、複数のキューの各々からインターフェイスデバイスにデータが送られるようになっている。複数のキューは、専有キューと共有キューとを含む。専有キューは、複数のプロセッサの内の1つのプロセッサである第1プロセッサだけが割り当てられているキューであり、データの格納時に排他処理が不要なキューである。共有キューは、複数のプロセッサの内の2以上の第2プロセッサが割り当てられるキューであり、データの格納時に排他処理が必要なキューである。第2プロセッサは、第1プロセッサ以外のいずれかのプロセッサである。

20

【発明の効果】

【0007】

本発明によれば、I/Fデバイスのリソース使用量と排他処理の負荷の両方を抑制できる。

【図面の簡単な説明】

【0008】

【図1】第1の実施形態に係るストレージシステムの動作概要を示す。

【図2】第1の実施形態に係るストレージシステムの構成例を示す。

30

【図3】キュー管理テーブルの構成例を示す。

【図4】プロセッサ別I/O（Input/Output）数テーブルの構成例を示す。

【図5】I/O要求をキューに格納する処理の例を示すフローチャートである。

【図6】専有キューを何れのプロセッサに割り当てるかを決定する処理の例を示すフローチャートである。

【図7】I/Fデバイスに関する複数のキューの内、何れを専有キューとするかを決定する処理の例を示すフローチャートである。

【図8】キュー割り当て変更処理の例を示すフローチャートである。本処理は、図7のS310に相当する。

【図9】第2の実施形態に係るストレージシステムの動作概要を示す。

40

【図10】第2の実施形態に係るストレージシステムの構成例を示す。

【図11】プロセス別I/O数テーブルの構成例を示す。

【図12】プロセスを何れのプロセッサで実行するかを決定する処理の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】

【0009】

以下、実施形態を説明する。以下の説明では、「xxxテーブル」の表現にて情報を説明することがあるが、情報は、どのようなデータ構造で表現されていてもよい。すなわち、情報がデータ構造に依存しないことを示すために、「xxxテーブル」を「xxx情報」と呼ぶことができる。

50

## 【0010】

また、以下の説明では、「xxx部」を主語として処理を説明する場合があるが、「xxx部」は、コンピュータプログラム（「プログラム」という）の一種であってもよい。プログラムは、プロセッサによって実行されることで、定められた処理を、適宜に記憶資源（例えばメモリ）及びネットワークインターフェイスデバイスの内の少なくとも1つを用いながら行うため、処理の主語が、プロセッサ、そのプロセッサを有する装置とされてもよい。プロセッサが行う処理の一部又は全部が、ハードウェア回路（例えば、ASIC（Application Specific Integrated Circuit）など）で行われてもよい。プログラムは、プログラムソースからインストールされてよい。プログラムソースは、プログラム配布サーバ又は記憶メディア（例えば可搬型の記憶メディア）であってもよい。また、プロセッサ及びメモリをまとめてコントローラと呼んでもよい。

10

## 【0011】

また、以下の説明では、同種の要素を区別して説明する場合には、「xxx1a」、「xxx1b」のように、参照符号を使用し、同種の要素を区別しないで説明する場合には、「xxx1」のように参照符号の内の共通番号のみを使用することがある。

<第1の実施形態>

## 【0012】

図1は、第1の実施形態に係るストレージシステム1aの動作概要を示す。

## 【0013】

20

第1の実施形態に係るストレージシステム1aは、記憶装置14と、記憶装置14へのデータの入出力を制御するコントローラ2aとを有する。コントローラ2aは、複数のプロセッサ11（例えば11a~11c）と、メモリ12aと、I/Fデバイス13（例えば13a）とを備える。

## 【0014】

メモリ12aは、プロセッサ11がI/Fデバイス13へ発行するデータ（I/O要求など）が一時的に格納されるキュー20（例えばキュー20a及び20b）を有する。

## 【0015】

各キュー20は、I/Fデバイス13に関連付けられている。プロセッサ11からI/Fデバイス13へ発行されたデータは、発行先であるI/Fデバイス13に関連付けられているキュー20に格納される。キュー20に格納されたデータは、順番に（例えばFIFOで）、I/Fデバイス13に渡されて処理される。

30

## 【0016】

コントローラ2aは、I/Fデバイス13に複数のキュー20が関連付けられている場合、何れかのキュー20にデータを格納する。ここで、1つのI/Fデバイス13に関連付けられている複数のキューの内、一部のキュー（1又は2以上のキュー）は、それぞれ、専有キュー20aであり、残りのキュー（1又は2以上のキュー）は、それぞれ、共有キュー20bである。つまり、1つのI/Fデバイス13につき、少なくとも1つの専有キュー20aと少なくとも1つの共有キュー20bとが存在する。専用キュー20aは、複数のプロセッサ11のうちのいずれか1つのプロセッサだけが割り当てられたキュー20である。以下、専有キュー20aに割り当てられているプロセッサを、「第1のプロセッサ」と言い、第1のプロセッサ以外のプロセッサを「第2のプロセッサ」と言うことがある。共有キュー20bは、全ての第2のプロセッサのうちの2以上の第2のプロセッサが割り当てられているキュー20である。共有キュー20bと第2のプロセッサの対応関係については、様々な対応関係を採用可能である。例えば、全ての共有キュー20bの各々に、全ての第2のプロセッサの各々が割り当てられてもよいし、第1の共有キュー20bに、2以上の第2のプロセッサが割り当てられ、第2の共有キュー20bに、2以上の第2のプロセッサが割り当てられてもよい。第1のプロセッサも第2のプロセッサも、それぞれ、そのプロセッサが割り当てられているキューにデータを格納することになる。

40

## 【0017】

50

共有キュー 20 b は、割り当てられている 2 以上の第 2 のプロセッサ 11 から利用され得る。よって、コントローラ 2 a は、第 2 のプロセッサがデータを共有キュー 20 b に格納する際、共有キュー 20 b について排他処理を行う必要がある。

【0018】

専有キュー 20 a は、第 1 のプロセッサ 11 以外のプロセッサ 11 から利用される可能性が無い。よって、コントローラ 2 a は、データを専有キュー 20 a に格納する際、専有キュー 20 a について排他処理を行う必要が無い。

【0019】

排他処理を行う場合は、排他処理を行わない場合と比較して、キュー 20 に対してデータを格納する処理に要する時間が長くなり、又、プロセッサ 11 の処理負荷も高くなる。なお、排他処理は、コントローラ 2 a において、データを出力するプロセッサにより行われてもよいし、そのプロセッサに代えて又は加えて、所定のロジック（例えば、I/F デバイス 13 との間でデータをやり取りする複数のプロセッサとは別に設けられたプロセッサであって排他処理等の所定処理を行うためのプロセッサ）により行われてもよい。

10

【0020】

ストレージシステムが備えるプロセッサ 11 と同じ数の分、専有キュー 20 a を設けることが難しい場合もある。なぜなら、キュー 20 の数が多いと、その分 I/F デバイス 13 で使用されるリソース量（例えばレジスタ 21 の数）が大きく、故に、I/F デバイス 13 のリソース量が不足してしまうことがあるからである。

【0021】

20

そこで、本実施形態では、上述のとおり、複数のキュー 20 の内の一部のキュー 20 を、それぞれ、1 つのプロセッサにのみ利用される専有キュー 20 a とし、残りのキューを、それぞれ、2 以上のプロセッサに利用される共有キュー 20 b とする。

【0022】

これにより、ストレージシステムが備えるプロセッサ 11 と同じ数の専有キュー 20 a を設けること無く（言い換えれば、1 つの I/F デバイス 13 に関連付けられるキュー 20 の数を、プロセッサ 11 の数より少なくでき、以って、I/F デバイスのリソース使用量を抑え）、且つ、排他処理が発生する回数を抑制する（つまり、排他処理の負荷を抑制する）ことができる。このため、例えば、排他処理がプロセッサ 11 により行われるのであれば、プロセッサ 11 の処理負荷を低減することになる。

30

【0023】

図 1 において、専有キュー 20 a はプロセッサ 11 c に割り当てられていて、共有キュー 20 b は、プロセッサ 11 a 及び 11 b に割り当てられているとする（図 1 の細い破線を参照）。

【0024】

その後、プロセッサ 11 a、11 b 及び 11 c の中で、プロセッサ 11 b の I/F デバイス 13 に対するデータの発行頻度が最大になった場合、次の処理が行われて良い。すなわち、コントローラ 2 a は、専有キュー 20 a に、プロセッサ 11 c に代えてプロセッサ 11 b を割り当て（図 1 の太い破線矢印を参照）、プロセッサ 11 c を、専有キュー 20 a に代えて共有キュー 20 b に割り当てる（図 1 の太い実線矢印を参照）。これにより、共有キュー 20 b にデータが格納される頻度が減少するため、排他処理が実行される回数も減少する。以下、第 1 の実施形態を詳細に説明する。

40

【0025】

図 2 は、第 1 の実施形態に係るストレージシステム 1 a の構成例を示す。

【0026】

ストレージシステム 1 a は、1 以上の記憶装置 14 と、それらの記憶装置 14 に接続されたコントローラ 2 a とを有する。コントローラ 2 a は、複数のプロセッサ（例えば 11 a ~ 11 c）11 と、メモリ 12 a と、I/F デバイス 13（例えば 13 a 及び 13 b）とを備え、これらの要素 11、12 a、13 は、内部バス 15 で接続されている。内部バス 15 は、双方向にデータの送信及び受信が可能である。内部バス 15 の一例としては、

50

PCI Express などがある。

【0027】

I/Fデバイス13は、プロセッサ11と記憶装置14との間のI/Oを制御するデバイスである。I/Fデバイス13は、プロセッサ11からライトコマンド、リードコマンド及び消去コマンドなどのデータ(例えばI/O要求)を受領し、そのデータを記憶装置14へ渡す。また、I/Fデバイス13は、記憶装置14からコマンド結果(例えばI/O要求の結果を含んだ応答)などのデータを受領し、そのデータをプロセッサ11へ渡す。I/Fデバイス13の一例としては、ホストバスアダプタ及びSATAコントローラなどがある。I/Fデバイス13は、I/Oデバイスと呼ばれてもよい。I/Fデバイス13に接続されるデバイスは、記憶装置14に限られない。I/Fデバイス13(例えば13a)は、リソースを有し、リソースには、そのI/Fデバイス13に関連付けられている複数のキュー(例えば20a及び20b)にそれぞれ対応付けられた複数のレジスタ21(例えば21a及び21b)が設けられる。レジスタ21は、I/Fデバイス13のリソースに設けられた記憶領域の一例である。

10

【0028】

プロセッサ11は、メモリ12aに記憶されているコンピュータプログラムを実行する。プロセッサ11は、I/Fデバイス13へコマンドなどのデータを送信したり、I/Fデバイス13からコマンド結果などのデータを受信したりする。プロセッサの一例としては、CPU(Central Processing Unit)及びLSI(Large-Scale Integrated circuit)などがある。

20

【0029】

メモリ12aには、プロセッサ11及びI/Fデバイス13などからアクセスされるコンピュータプログラム及びデータなどが格納される。メモリ12aの一例としては、DRAM(Dynamite Random Access Memory)、MRAM(Magnetic Random Access Memory)及びFeRAM(Ferroelectric Random Access Memory)などがある。

【0030】

メモリ12aは、I/Fデバイス13毎に、I/Fデバイス13に関連付けられたキューセット(複数のキュー20)を有する。図2の例によれば、I/Fデバイス13aに対して、キュー20a及び20bが関連付けられており、I/Fデバイス13bに対して、キュー20c及び20dが関連付けられている。いずれのI/Fデバイス13についても、関連付けられているキュー20の数は、いずれかのキュー20に割り当てられ得るプロセッサ11の数より少なくてもよい。また、キュー20の数は、全てのI/Fデバイス13について同じでもよいし異なってもよい。

30

【0031】

メモリ12aは、複数のプロセス31(例えば31a~31d)と、プロセススケジューラ32と、I/O制御部33と、キュー制御部34と、キュー割当部35とを記憶する。これらの要素31~35は、プログラムであって、プロセッサ11によって読み出されて実行されてよい。I/O制御部33と、キュー制御部34及びキュー割当部35とが、それぞれ、複数のプロセッサ11の各々で実行されてもよいし、I/O制御部33が、複数のプロセッサ11の各々で実行され、キュー制御部34及びキュー割当部35のうちの少なくとも一方が、複数のプロセッサ11とは別のプロセッサ(図示せず)により実行されてもよい。

40

【0032】

メモリ12aは、キュー管理テーブル100と、プロセッサ別I/O数テーブル120とを記憶する。これらの要素100、120は、データであって、プロセッサ11で実行されるプログラム31~35によって読み出されたり、書き換えられたりしてよい。

【0033】

プロセス31は、プロセッサ11で実行されるプログラムである。1つのプロセス31は、複数のプロセッサ11の内の1つのプロセッサ11で実行される。

50

## 【0034】

プロセススケジューラ32は、各プロセッサ11における各プロセス31の処理時間をスケジューリングする。プロセススケジューラ32は、OS(Operating System)の機能として提供され、プロセス31は、そのOS上で実行されてよい。

## 【0035】

I/O制御部33は、プロセス31のI/Oを制御する。例えば、I/O制御部33は、プロセス31から発行されたコマンドを解析して、所定のキュー20に格納したり、I/Fデバイス13から発行されたコマンド結果を、プロセス31に通知したりする。

## 【0036】

キュー制御部34は、キュー20を制御する。キュー制御部34は、プロセッサ11がI/Fデバイス13へデータを送信する際、そのI/Fデバイス13に関連付いているキュー20にデータを格納する。キュー制御部34は、共有キュー20にデータを格納する際、排他処理を行う。例えば、排他処理は、(1)共有キュー20が他のプロセッサ11からロックされていないことを確認すること、(2)ロックがされていないことの確認が取れたら共有キュー20を他のプロセッサ11から書き込みできないようにロックすること、(3)共有キュー20へのデータの書き込みが完了した場合にロックを解除すること、を含んだ処理である。

10

## 【0037】

キュー制御部34は、専有キュー20にデータを格納する際、上述の排他処理を行わない。すなわち、キュー制御部34は、キュー20のロック及びロック解除を行うことなく、専有キュー20にデータを格納する。

20

## 【0038】

I/Fデバイス13は、自己のレジスタ21を参照し、そのレジスタ21に対応付けられているキュー20からFIFOでデータを取得する。そして、I/Fデバイス13は、その取得したデータを、例えば、記憶装置14に送信する。

## 【0039】

キュー割当部35は、キュー20とプロセッサ11との対応関係を制御する。具体的には、キュー割当部35は、専有キュー20に割り当てられる1つのプロセッサ11を決定したり、共有キュー20に割り当てられる2以上のプロセッサ11を決定したりする。又は、キュー割当部35は、複数のキュー20の内、何れを専有キューとするか(何れを共有キューとするか)を決定してもよい。

30

## 【0040】

キュー管理テーブル100は、I/Fデバイス13とキュー20との対応関係と、キュー20とプロセッサ11との対応関係と、キュー20のモード(専有キュー又は共有キュー)と、を管理するためのテーブルである。キュー管理テーブル100の詳細については後述する(図3参照)。

## 【0041】

プロセッサ別I/O数テーブル120は、プロセッサ11がI/Fデバイス13に対して発行したI/O要求の数を管理するためのテーブルである。プロセッサ別I/O数テーブル120の詳細については後述する(図4参照)。

40

## 【0042】

図3は、キュー管理テーブル100の構成例を示す。

## 【0043】

キュー管理テーブル100は、I/Fデバイス13とキュー20との対応関係と、キュー20とプロセッサ11との対応関係と、キュー20のモード(専有キュー又は共有キュー)と、を管理する。キュー管理テーブル100は、例えば、キュー20毎にレコードを有し、各レコードが、フィールド値として、デバイスID101と、キューID102と、プロセッサID103と、モード104とを有する。

## 【0044】

デバイスID101は、キュー20が関連付けられているI/Fデバイス13を識別す

50

るための情報である。キューID102は、キュー20を識別するための情報である。プロセスID103は、キュー20に割り当てられているプロセス11を識別するための情報である。

【0045】

モード104は、キュー20が「専有キュー」であるか「共有キュー」であるかを示す情報である。モード104は、「専有」の場合は「ON」、「共有」の場合は「OFF」とするフラグであってもよい。

【0046】

図3に示すキュー管理テーブル100の例によれば、次のことがわかる。すなわち、デバイスID101「D1」のI/Fデバイス13aには、キューID102「Q1」及び「Q2」のキュー20a、20bが関連付けられている。そして、キューID102「Q1」のキュー20aは、「専有キュー」(104)であり、プロセスID103「U1」に割り当てられている。キューID102「Q2」のキュー20bは、「共有キュー」(104)であり、プロセスID103「U2」及び「U3」のプロセス11b、11cに割り当てられている。

10

【0047】

図4は、プロセス別I/O数テーブル120の構成例を示す。

【0048】

プロセス別I/O数テーブル120は、プロセス11がI/Fデバイス13に対して発行したI/O要求の数を管理するためのテーブルである。

20

【0049】

プロセス別I/O数テーブル120は、プロセス11毎にレコードを有し、各レコードが、フィールド値として、プロセスID121と、デバイスID122と、I/O数123とを有する。プロセスID121及びデバイスID122は、図3で説明したとおりである。

【0050】

I/O数123は、プロセス11がI/Fデバイス13に対して所定期間に発行したI/O要求の数である。I/O数123は、IOPS(Input/Output Per Second)であってもよい。また、I/O数123において、シーケンシャルアクセスに係るI/O要求と、ランダムアクセスに係るI/O要求とが区別されてもよい。

30

【0051】

図4に示すプロセス別I/O数テーブル120の例によれば、次のことがわかる。すなわち、プロセスID121「U1」のプロセス11aは、デバイスID122「D1」のI/Fデバイス13aに対して、所定期間に「431回」のI/O要求を発行した。プロセスID121「U1」のプロセス11aは、デバイスID122「D2」のI/Fデバイス13bに対して、所定期間に「5回」のI/O要求を発行した。

【0052】

図5は、I/O要求をキュー20に格納する処理の例を示すフローチャートである。

【0053】

I/O制御部33は、I/O要求を発行したプロセス11のプロセスIDを特定する(S101)。そして、I/O制御部33は、I/O要求が指定するI/Fデバイス13のデバイスIDを特定する(S102)。

40

【0054】

次に、I/O制御部33は、キュー管理テーブル100から、その特定したデバイスID101に対応するキューID102を特定する。そして、I/O制御部33は、キュー管理テーブル100から、その取得したキューID102に対応するモード104を特定する(S103)。

【0055】

次に、I/O制御部33は、その特定したモード104が「専有」又は「共有」の何れであるかを判定する(S104)。

50

## 【 0 0 5 6 】

S 1 0 4 で特定したモード 1 0 4 が「共有」の場合 ( S 1 0 4 : 共有 )、I / O 制御部 3 3 は、次の処理を行う。すなわち、I / O 制御部 3 3 は、S 1 0 3 において特定したキュー ID 1 0 2 の共有キュー 2 0 をロックし ( S 1 0 5 )、その後、その共有キュー 2 0 に I / O 要求を格納する ( S 1 0 6 )。そして、I / O 制御部 3 3 は、共有キュー 2 0 のロックを解除し ( S 1 0 7 )、本処理を終了する。すなわち、I / O 制御部 3 3 は、共有キュー 2 0 に対して排他処理を行う。

## 【 0 0 5 7 】

S 1 0 4 で取得したモードが「専有」の場合 ( S 1 0 4 : 専有 )、I / O 制御部 3 3 は、その専有キュー 2 0 に I / O 要求を格納し ( S 1 1 1 )、本処理を終了する。すなわち、I / O 制御部 3 3 は、専有キュー 2 0 に対して排他処理を行わない。よって、専有キュー 2 0 が利用される方が、共有キュー 2 0 が利用される場合よりも、プロセッサ 1 1 の処理負荷が小さくなる。

10

## 【 0 0 5 8 】

図 6 は、専有キュー 2 0 を何れのプロセッサ 1 1 に割り当てるかを決定する処理の例を示すフローチャートである。

## 【 0 0 5 9 】

本処理は、I / F デバイス 1 3 に関する複数のキュー 2 0 の内、専有キューとして利用されるキューの数が決まっている場合の例である。また、本処理は、1 つの I / F デバイス 1 3 に対する処理の例である。ストレージシステム 1 a が複数の I / F デバイス 1 3 を備える場合、I / F デバイス 1 3 毎に本処理が実行される。以下、1 つの I / F デバイス 1 3 を例に取り、図 6 の説明において、その I / F デバイス 1 3 を「対象 I / F デバイス」と言う。

20

## 【 0 0 6 0 】

キュー割当部 3 5 は、専有キュー 2 0 毎にループ A の処理を実行する。以下、1 つのループ A の処理を例に取り、そのループ A の処理の対象となる専有キュー 2 0 を「対象専有キュー」という。対象 I / F デバイス 1 3 に関連付けられている複数のキュー 2 0 の内、何れのキューが専有キューであるかは、キュー管理テーブル 1 0 0 のデバイス ID 1 0 1 及びモード 1 0 4 を参照することによりわかる。

## 【 0 0 6 1 】

次に、キュー割当部 3 5 は、ループ A の処理において、プロセッサ 1 1 のそれぞれについて、ループ B の処理を実行する。以下、1 つのループ B の処理を例に取り、そのループ B の処理の対象となるプロセッサ 1 1 を「対象プロセッサ」という。

30

## 【 0 0 6 2 】

次に、キュー割当部 3 5 は、プロセッサ別 I / O 数テーブル 1 2 0 を参照し、現状における対象 I / F デバイス 1 3 に関する排他数を算出する ( S 2 0 3 )。ここで算出した排他数を「前排他数」という。

## 【 0 0 6 3 】

排他数とは、対象 I / F デバイス 1 3 に関連付けられているキュー 2 0 に、プロセッサ 1 1 から発行された I / O 要求を格納するにあたって、排他処理が実行された回数を示す。すなわち、専有キュー 2 0 に I / O 要求が格納された場合、排他数はカウントされず、共有キュー 2 0 に I / O 要求が格納された場合、排他数はカウントされる。現状における I / F デバイス 1 3 に関する排他数は、対象 I / F デバイス 1 3 に関連付けられている全ての共有キュー 2 0 にそれぞれ対応した I / O 数 1 2 3 の合計であってもよい。また、排他数は、共有キュー 2 0 に I / O 要求が格納される頻度という観点から「共有頻度」と呼ばれてもよい。

40

## 【 0 0 6 4 】

前排他数は、次のように算出される。すなわち、例えば、図 3 のキュー管理テーブル 1 0 0 に示すように、デバイス ID 1 0 1 「D 1」の対象 I / F デバイス 1 3 a には、キュー ID 1 0 2 「Q 1」の専有キュー 2 0 a と、キュー ID 1 0 2 「Q 2」の共有キュー 2

50

0 b とが関連付けられている。そして、専有キュー 20 a には、プロセッサ ID 103 「U 1」のプロセッサ 11 a が割り当てられており、共有キュー 20 b には、プロセッサ ID 103 「U 2」及び「U 3」のプロセッサ 11 b 及び 11 c が割り当てられている。そして、図 4 のプロセッサ別 I/O 数テーブル 120 に示すように、プロセッサ ID 121 「U 1」のプロセッサ 11 a から、デバイス ID 122 「D 1」の I/F デバイス 13 a に対する I/O 数 123 は「431 回」である。プロセッサ ID 121 「U 2」のプロセッサ 11 b から、I/F デバイス 13 a に対する I/O 数 123 は「19 回」である。プロセッサ ID 121 「U 3」のプロセッサ 11 c から、I/F デバイス 13 a に対する I/O 数 123 は「2 回」である。この場合、対象 I/F デバイス 13 a に関する前排他数は、共有キュー 20 b に割り当てられているプロセッサ 11 b、11 c からの I/O 数 123 の合計「19 + 2 = 21」となる。

10

## 【0065】

次に、キュー割当部 35 は、対象専有キューに対象プロセッサ 11 を割り当てたと仮定した場合における、対象 I/F デバイスに関する排他数を算出（推定）する（S204）。ここで算出した排他数を「後排他数」という。

## 【0066】

前排他数は、次のように算出される。すなわち、例えば、仮に、専有キュー 20 a にプロセッサ 11 b を割り当てた場合、元々専有キュー 20 a に割り当てられていたプロセッサ 11 a は、共有キュー 20 b に割り当てられることになる。したがって、I/F デバイス 13 a（デバイス ID 「D 1」）に関する後排他数は、共有キュー 20 b に割り当てられることになるプロセッサ 11 a（プロセッサ ID 「U 1」）及びプロセッサ 11 c（プロセッサ ID 「U 3」）の I/O 数 123 の合計「431 + 2 = 433」となる。

20

## 【0067】

次に、キュー割当部 35 は、後排他数と前排他数との差 Z1 が所定の閾値 C1（C1 は 0 以下の値）未満であるか否かを判定する（S205）。つまり、キュー割当部 35 は、「Z1（= 後排他数 - 前排他数）< 閾値 C 0」であるか否かを判定する。

## 【0068】

S205 の判定が肯定的な場合（S205：YES）、キュー割当部 35 は、対象専有キューに、対象専用キューに割り当てられていたプロセッサに代えて対象プロセッサ 11 を割り当てる（S210）。すなわち、キュー割当部 35 は、キュー管理テーブル 100 において、キュー ID 102 が対象専有キューと一致するレコードにおけるプロセッサ ID 103 を、対象プロセッサ 11 に更新する処理を行う。なぜなら、対象専有キューに対象プロセッサ 11 を割り当てた方が、I/F デバイス 13 に関する排他数が減少する可能性が高いと考えられるからである。そして、キュー割当部 35 は、ループ B の処理を抜ける。

30

## 【0069】

S205 の判定が否定的な場合（S205：NO）、キュー割当部 35 は、次の処理を行う。キュー割当部 35 は、全てのプロセッサ 11 についてループ B の処理を完了した場合、ループ B の処理を抜け、まだ未処理のプロセッサ 11 が残っている場合、ループ B の処理を繰り返す。すなわち、キュー割当部 35 は、対象専有キューに割り当てるプロセッサ 11 を現状のままとする。なぜなら、対象専有キューに割り当てるプロセッサ 11 を変更しても、I/F デバイス 13 に関する排他数が減少する可能性が低いと考えられるからである。

40

## 【0070】

キュー割当部 35 は、ループ B の処理を抜けた後、次の処理を行う。キュー割当部 35 は、I/F デバイスに関する全ての専有キューについてループ A の処理を完了した場合、ループ A の処理を抜けて本処理を終了し、まだ未処理の専有キューが残っている場合、ループ A の処理を繰り返す。

## 【0071】

図 7 は、I/F デバイス 13 に関連付けられている複数のキューの内、何れを専有ク

50

ーとするかを決定する処理の例を示すフローチャートである。

【0072】

本処理は、I/Fデバイス13に関連付けられている複数のキューの内、専有キューとして利用されるキューの数が決まっていなかった場合の例である。すなわち、本処理の場合、必要に応じて専有キューの数が増減し得る。また、本処理は、1つのI/Fデバイス13に対する処理の例である。ストレージシステム1aが複数のI/Fデバイス13を備える場合、I/Fデバイス13毎に本処理が実行される。キュー割当部35は、図6に示す処理と図7に示す処理の何れかを採用しても良いし、所定のタイミングで図6に示す処理と図7に示す処理とを切り替えても良い。以下、1つのI/Fデバイス13を例に取り、図7の説明において、そのI/Fデバイス13を「対象I/Fデバイス」という。

10

【0073】

キュー割当部35は、キュー管理テーブル100及びプロセッサ別I/O数テーブル120を参照し、対象I/Fデバイスに関連付いている各キュー20の利用指標を算出する(S301)。ここで算出した利用指標を、「前利用指標」という。利用指標とは、キュー20がプロセッサ11にどのように利用されているか(つまり、キューの利用態様)を示す指標である。利用指標は、キュー20に割り当てられているプロセッサ11の数と、各プロセッサ11からそのキュー20に発行されたI/O数とに基づいて算出されてよい。利用指標は、キュー20に発行されたI/O数が多くなると、大きくなる値であってよい。反対に、利用指標は、キュー20に割り当てられているプロセッサ11の数が多くなると、小さくなる値であってよい。例えば、利用指標は、キュー20に割り当てられている各プロセッサ11からそのキュー20に発行されたI/O数の合計が、そのキュー20に割り当てられているプロセッサ数の逆数によって重み付けされた値であってよい。この観点から、「利用指標」は、1のキュー20における1のプロセッサ11あたりの利用負荷指標と定義されてもよい。

20

【0074】

例えば、キュー20に割り当てられているプロセッサ11の数が1つの場合に対する係数を「2.0」、プロセッサ11の数が2つの場合に対する係数を「1.5」、プロセッサ11の数が3つの場合に対する係数を「1.0」とする。そして、図3のキュー管理テーブル100に示すように、デバイスID101「D2」のI/Fデバイス13bに関連付いているキューID102「Q3」のキュー20cには、プロセッサID103「U3」のプロセッサ11cだけが割り当てられている。そして、図4のプロセッサ別I/O数テーブル120に示すように、プロセッサ11c(プロセッサID121「U3」)からデバイスID122「D2」のI/Fデバイス13bに発行されたI/O数123は「11回」である。この場合、キュー20cの前利用指標は、「 $2.0 \times 11 = 22$ 」と算出されてよい。

30

【0075】

また、デバイスID「D2」のI/Fデバイス13bに関連付いているキューID102「Q4」のキュー20dには、プロセッサID103「U1」及び「U2」のプロセッサ11a及び11bが割り当てられている。そして、プロセッサ11a及び11bからキュー20dに発行されたI/O数123は、それぞれ「5回」及び「189回」である。この場合、キュー20dの前利用指標は、「 $1.5 \times (5 + 189) = 291$ 」と算出されてよい。

40

【0076】

次に、キュー割当部35は、各キューの中から、前利用指標が最大のキュー(「第1のキュー」という)と、前利用指標が最小のキュー(「第2のキュー」という)を特定する(S302)。つまり、第1のキューは、少数のプロセッサ11から多数のI/O要求を受けているキューであるといえる。第2のキューは、多数のプロセッサ11から少数のI/O要求を受けているキューであるといえる。上記の例において、キュー20cの前利用指標は「22」、キュー20dの前利用指標は「291」であるので、キュー割当部35は、第1のキューとしてキュー20dを特定し、第2のキューとしてキュー20cを特定

50

する。

【 0 0 7 7 】

次に、キュー割当部 3 5 は、I / F デバイス 1 3 に I / O 要求を発行する複数のプロセッサ 1 1 の中から、キューの割当の変更について検討するプロセッサ 1 1 を選択する。例えば、キュー割当部 3 5 は、プロセッサ別 I / O 数テーブル 1 2 0 を参照し、第 1 のキューに割り当てられている複数のプロセッサ 1 1 の中で、I / O 数 1 2 3 が最小のプロセッサ 1 1 を選択してもよい ( S 3 0 3 )。上記の例において、キュー割当部 3 5 は、第 1 のキュー 2 0 d ( キュー I D 「 Q 4 」 ) に割り当てられているプロセッサ 1 1 a 及び 1 1 b ( プロセッサ I D 「 U 1 」 及び 「 U 2 」 ) の中から、I / O 数 1 2 3 が最小のプロセッサ 1 1 a ( プロセッサ I D 「 U 1 」 ) を選択する。

10

【 0 0 7 8 】

次に、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 及びプロセッサ別 I / O 数テーブル 1 2 0 を参照し、選択プロセッサ 1 1 ( S 3 0 3 で選択したプロセッサ ) を第 2 のキューに割り当てたと仮定した場合における各キューの利用指標を算出 ( 推定 ) する ( S 3 0 4 )。ここで算出した利用指標を、「後利用指標」という。

【 0 0 7 9 】

第 2 のキューに選択プロセッサ 1 1 を割り当てたと仮定した場合、第 1 のキューの後利用指標は、選択プロセッサ 1 1 の割り当てが解除される分、係数が大きくなり ( プロセッサ数が減少するため )、I / O 数が減少する。一方、第 2 のキューの後利用指標は、選択プロセッサ 1 1 が割り当てられる分、係数が小さくなり ( プロセッサ数が増加するため )、I / O 数が増加する。

20

【 0 0 8 0 】

上記の例の場合、第 2 のキュー 2 0 c に選択プロセッサ 1 1 a を割り当てた場合、第 2 のキュー 2 0 c に、プロセッサ 1 1 a、1 1 c ( プロセッサ I D 「 U 1 」、「U 3 」 ) が割り当てられることとなる。よって、第 2 のキュー 2 0 c の後利用指標は、「 $1.5 \times ( 5 + 11 ) = 24$ 」と算出される。一方、第 1 のキュー 2 0 d には、プロセッサ 1 1 b ( プロセッサ I D 「 U 2 」 ) だけが割り当てられることとなる。よって、第 1 のキュー 2 0 d の後利用指標は、「 $2.0 \times 189 = 378$ 」と算出される。

【 0 0 8 1 】

次に、キュー割当部 3 5 は、第 1 のキューの後利用指標が、第 2 のキューの後利用指標以上であるか否かを判定する ( S 3 0 5 )。すなわち、キュー割当部 3 5 は、仮に第 2 のキューに選択プロセッサ 1 1 を割り当てた場合において、第 1 のキューと第 2 のキューとの間における利用指標の大小関係が逆転するか否かを判定する。

30

【 0 0 8 2 】

第 1 のキューの後利用指標が第 2 のキューの後利用指標よりも小さい場合 ( S 3 0 5 : N O )、キュー割当部 3 5 は、本処理を終了する。すなわち、キュー割当部 3 5 は、第 1 のキューと第 2 のキューとの間において利用指標の大小関係が逆転する等により利用指標の差が一定以上減少しない場合、キューに対する選択プロセッサ 1 1 の割り当てを変更しない。なぜなら、変更したとしても、I / F デバイス 1 3 に関連付いているキュー 2 0 の全体における排他処理の発生回数が減少する可能性が低いからである。

40

【 0 0 8 3 】

第 1 のキューの後利用指標が第 2 のキューの後利用指標以上である場合 ( S 3 0 5 : Y E S )、キュー割当部 3 5 は、キュー割り当て変更処理を実行し ( S 3 1 0 )、本処理を終了する。すなわち、キュー割当部 3 5 は、第 1 のキューと第 2 のキューとの間において利用指標の大小関係が逆転等せずに利用指標の差が一定以上減少する場合、本処理を終了する。キュー割り当て変更処理 ( S 3 1 0 ) の詳細については後述する ( 図 8 参照 )。上記の例の場合、第 1 のキュー 2 0 d の後利用指標は「 $378$ 」、第 2 のキュー 2 0 c の後利用指標は「 $24$ 」と算出されるので、キュー割当部 3 5 は、キュー割り当て変更処理 ( S 3 1 0 ) を実行する。

【 0 0 8 4 】

50

図 8 は、キュー割り当て変更処理の例を示すフローチャートである。本処理は、図 7 の S 3 1 0 に相当する。

【 0 0 8 5 】

キュー割当部 3 5 は、第 2 のキューに選択プロセッサ 1 1 を割り当てる ( S 4 0 1 ) 。上記の例の場合、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 において、選択プロセッサ ID 1 0 3 「 U 1 」を、キュー ID 1 0 2 「 Q 4 」との対応関係から、キュー ID 1 0 2 「 Q 3 」 ( 第 2 のキュー ) との対応関係へ変更する。

【 0 0 8 6 】

次に、キュー割当部 3 5 は、第 1 のキューに割り当てられているプロセッサを特定する ( S 4 0 2 ) 。上記の例の場合、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 において、キュー ID 1 0 2 「 Q 4 」 ( 第 1 のキュー ) と対応関係を有するプロセッサ ID 1 0 3 「 U 2 」を特定する。

10

【 0 0 8 7 】

そして、キュー割当部 3 5 は、第 1 のキューに割り当てられているプロセッサの数が「 1 」であるか否かを判定する ( S 4 0 3 ) 。

【 0 0 8 8 】

第 1 のキューに割り当てられているプロセッサの数が「 1 」でない場合 ( S 4 0 3 : N O ) 、キュー割当部 3 5 は、そのまま S 4 1 0 へ進む。つまり、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 において、第 1 のキューのモード 1 0 4 を「共有」のままとする。

20

【 0 0 8 9 】

第 1 のキューに割り当てられているプロセッサの数が「 1 」である場合 ( S 4 0 3 : Y E S ) 、キュー割当部 3 5 は、第 1 のキューを専有キューに変更する ( S 4 0 4 ) 。つまり、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 において、第 1 のキューのモード 1 0 4 を「専有」に変更する。そして、キュー割当部 3 5 は、 S 4 1 0 へ進む。

【 0 0 9 0 】

上記の例の場合、第 1 のキュー 2 0 d ( キュー ID 1 0 2 「 Q 4 」 ) と対応関係を有するプロセッサの数は「 1 」 ( プロセッサ ID 「 U 2 」のみ ) となるので、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 において、キュー ID 1 0 2 「 Q 4 」のモード 1 0 4 を「専有」に変更する。

30

【 0 0 9 1 】

次に、キュー割当部 3 5 は、第 2 のキューに割り当てられているプロセッサを特定する ( S 4 1 0 ) 。上記の例の場合、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 において、第 2 のキュー 2 0 c ( キュー ID 1 0 2 「 Q 3 」 ) と対応関係を有するプロセッサ ID 1 0 3 「 U 1 」及び「 U 3 」を特定する。

【 0 0 9 2 】

そして、キュー割当部 3 5 は、第 2 のキューに割り当てられているプロセッサの数が「 2 以上」であるか否かを判定する ( S 4 1 1 ) 。

【 0 0 9 3 】

第 2 のキューに割り当てられているプロセッサの数が「 2 以上」でない場合 ( S 4 1 1 : N O ) 、キュー割当部 3 5 は、本処理を終了し、図 7 に示す処理へ戻る。つまり、キュー割当部 3 5 は、キュー管理テーブル 1 0 0 において、第 2 のキューのモード 1 0 4 を「専有」のままとする。

40

【 0 0 9 4 】

第 2 のキューに割り当てられているプロセッサの数が「 2 以上」である場合 ( S 4 1 1 : Y E S ) 、キュー割当部 3 5 は、第 2 のキューのモード 1 0 4 を「共有」に変更する ( S 4 1 2 ) 。そして、キュー割当部 3 5 は、本処理を終了し、図 7 に示す処理に戻る。

【 0 0 9 5 】

上記の例の場合、第 2 のキュー 2 0 c ( キュー ID 1 0 2 「 Q 3 」 ) と対応関係を有するプロセッサの数は「 2 」 ( プロセッサ ID 「 U 1 」及び「 U 3 」 ) となるので、キュー

50

割当部 35 は、キュー管理テーブル 100 において、キュー ID 102 「Q3」のモード 104 を「共有」に変更する。

【0096】

以上の第 1 の実施形態の説明によれば、例えば以下の説明をできる。

【0097】

コントローラ 2a (例えば、プロセッサ 11 で実行されるキュー制御部 34 (図 2 参照)) は、I/F デバイス 13 に関連付けられているキュー 20 に対する I/O を制御する。コントローラ 2a は、複数のプロセッサ 11 の各々の I/F デバイス 13 に対するデータの発行頻度に基づいて、専有キュー 20a に割り当てられるプロセッサ 11 (「第 1 のプロセッサ」) を決定してよい。第 1 のプロセッサとされるプロセッサ 11 は、複数のプロセッサ 11 の内、I/F デバイス 13 に対するデータの発行頻度が最大のプロセッサ 11 であってよい。これにより、I/F デバイス 13 に対するデータの発行頻度が最大のプロセッサ 11 から発行されるデータについて、排他処理が不要となる。

10

【0098】

コントローラ 2a は、共有キューにデータが格納される頻度である共有頻度に関する所定の条件が満たされた場合、共有キュー 20b に割り当てられているプロセッサ (この段落で「第 2 のプロセッサ」) 11 を、専有キュー 20a に割り当てられているプロセッサ (この段落で「第 1 のプロセッサ」) 11 に代えて専有キュー 20a に割り当て、且つ、第 2 のプロセッサに代えて第 1 のプロセッサを共有キュー 20b に割り当てる割当て切替えを実行してもよい。これにより、排他処理の発生状況が変化した場合に、専有キュー 20a に割り当てられるプロセッサを変更することができる。

20

【0099】

所定の条件は、後共有頻度が前共有頻度よりも所定の閾値以上小さいことであってよい。ここで、後共有頻度とは、第 2 プロセッサを第 1 プロセッサに代えて専有キュー 20a に割り当て、且つ、第 2 プロセッサに代えて第 1 プロセッサを共有キュー 20b に割り当てる割当て切替えを実行したと仮定した場合の共有頻度 (共有キューにデータが格納される頻度) であってよい。前共有頻度とは、その割当て切替え前の共有頻度であってよい。これにより、排他処理の発生状況が変化した場合に、専有キュー 20a に、排他処理を発生させる可能性のより高いプロセッサ 11 を割り当てることができる。結果として、キュー 20 の全体における排他処理の発生回数を減少させることができる。

30

【0100】

コントローラ 2a は、複数のキュー 20 の各々について、専有及び共有の何れのキュー種別であるかを管理し、共有のキューについては、そのキューにロックをかけることを含んだ排他処理を行ってそのキューにデータを格納し、専有のキューについては、排他処理無しにそのキューにデータを格納してよい。また、コントローラ 2a は、複数のプロセッサと複数のキューとの割当関係を繰り返し更新してよい。そして、コントローラ 2a は、キューに割り当てられているプロセッサ数が 1 つだけになった場合に、そのキューのキュー種別を共有から専有に変更してよい。

【0101】

40

コントローラ 2a は、複数のキューの内の利用負荷指標が最大のキューである第 1 のキューについて、そのキューのキュー種別を、共有から専有に変更するか否か判定してよい。ここで、利用負荷指標とは、キューにおける 1 のプロセッサあたりの利用負荷の大きさを示す指標であって、そのキューに割り当てられているプロセッサ数が増えると小さくなり、各プロセッサからそのキューにデータが格納された頻度が大きくなると大きくなる値であってよい。

【0102】

コントローラ 2a は、上記の第 1 のキューに割り当てられている複数のプロセッサ 11 の内、その第 1 のキューにデータを格納する頻度が最小のプロセッサ 11 を、複数のキューの内の利用負荷指標が最小のキューである第 2 のキューに割り当てる割当て切替えを実

50

行したと仮定した場合の第1のキューの利用負荷指標及び第2のキューの利用負荷指標をそれぞれ推定し、その推定した第1のキューの利用負荷指標が、その推定された第2のキューの利用負荷指標よりも大きい場合に、第2のキューに最小のプロセッサを割り当てる割当て切替えを実行してもよい。

< 第2の実施形態 >

【0103】

以下、第2の実施形態を説明する。その際、第1の実施形態の相違点を主に説明し、第1の実施形態との共通点については説明を省略又は簡略する。

【0104】

図9は、第2の実施形態に係るストレージシステム1bの動作概要を示す。

10

【0105】

第2の実施形態に係るストレージシステム1bは、複数のプロセッサ11と、メモリ12bと、I/Fデバイス13(例えば13b)と、キュー20(例えば20c及び20d)とを有する。しかし、ストレージシステム1bは、専有キュー20とプロセッサ11との割当て関係を変更せず、I/O要求を発行するプロセス31を、適切なプロセッサ11に移動させる。

【0106】

コントローラ2bは、プロセス31のプロセッサ11への配置を制御する。プロセス31から発行されたデータは、そのプロセス31を実行するプロセッサ11が割り当てられているキュー20に格納される。そこで、コントローラ2bは、I/Fデバイス13に対するデータの発行頻度に基づいて、専有キュー20に割り当てられている第1のプロセッサ11で実行されるプロセス31を決定してもよい。これにより、排他処理を発生させる可能性のより高いプロセス31を、専有キュー20aに割り当てられているプロセッサ11で実行することができる。つまり、キュー20の全体における排他処理の発生回数が減少し得る。

20

【0107】

コントローラ2bは、専有キューに割り当てられていて負荷が所定の閾値未満であるプロセッサ(この段落で「第1のプロセッサ」)11に、共有キューに割り当てられているプロセッサ(この段落で「第2のプロセッサ」)11で実行されている1以上のプロセス31を移動させるプロセス移動を実行したと仮定した場合の排他数(共有キューで排他処理が発生する回数)を推定する。コントローラ2bは、その推定した排他数が、プロセス移動前の排他数よりも所定の閾値以上小さい場合、そのプロセス移動を実行してもよい。これにより、排他処理の発生状況が変化した場合、排他処理を発生させる可能性のより高いプロセス31を、専有キュー20aに割り当てられているプロセッサ11に移動させることができる。

30

【0108】

図9において、ストレージシステム1bは、I/Fデバイス13bに関連するキュー20c、20dを有する。キュー20cは「専有キュー」であり、キュー20dは「共有キュー」である。専有キュー20cには、プロセッサ11cが割り当てられている。共有キュー20dには、プロセッサ11a及び11cが割り当てられている。また、プロセッサ11aではプロセス31aが実行されており、プロセッサ11bではプロセス31bが実行されており、プロセッサ11cではプロセス31cが実行されている。

40

【0109】

その後、プロセス31a、31b及び31cの中で、プロセス31bのI/Fデバイス13bに対するデータの発行頻度が最大になった場合、次の処理が行われて良い。すなわち、コントローラ1bは、プロセス31bを、専有キュー20cに割り当てられているプロセッサ11cへ移動させ、プロセス31cを、共有キュー20dに割り当てられているプロセッサ11bへ移動させる。これにより、共有キュー20dにデータが格納される頻度が減少し、その結果、排他処理の実行される回数も減少する。

【0110】

50

図10は、第2の実施形態に係るストレージシステム1bの構成例を示す。

【0111】

コントローラ1bのメモリ12bにおいて、キュー割当部35に代えて、プロセス配置部37が記憶される。プロセス配置部37は、プロセス31を実行するプロセッサ11を決定する。プロセス配置部37は、プロセス31を実行するプロセッサ11を変更する。すなわち、プロセス配置部37は、或るプロセッサ11で実行されているプロセス31を、別のプロセッサ11へ移動させることができる。

【0112】

また、プロセッサ別I/O数テーブル140の構成は、第1の実施形態でのプロセッサ部I/O数テーブル120の構成と異なる。プロセス別I/O数テーブル140は、プロセス31がI/Fデバイス13に対して発行したI/O要求の数を管理するためのテーブルである。プロセス別I/O数テーブル140は、プロセス31が実行されているプロセッサ11と、プロセッサ11の使用率とを合わせて管理しても良い。プロセス別I/O数テーブル140の詳細については後述する(図11参照)。

10

【0113】

図11は、プロセス別I/O数テーブル140の構成例を示す。

【0114】

プロセス別I/O数テーブル140は、プロセス31がI/Fデバイス13に対して発行したI/O要求の数を管理する。さらに、プロセス別I/O数テーブル140は、プロセス31が処理されているプロセッサ11と、プロセッサ11の使用率も管理する。

20

【0115】

プロセス別I/O数テーブル140は、例えば、プロセス31毎にレコードを有し、各レコードが、フィールド値として、プロセスID141と、プロセッサID142と、プロセッサ使用率143と、デバイスID144と、I/O数145とを有する。

【0116】

プロセスID141は、プロセス31を識別するための情報である。プロセッサID142及びデバイスID144は、図3において説明したとおりである。I/O数145は、図4において説明したとおりである。

【0117】

図11に示すプロセス別I/O数テーブル140の例によれば、次のことがわかる。すなわち、プロセスID141「1001」のプロセス31aは、プロセッサID142「U1」のプロセッサ11aで実行されている。そのプロセッサID142「U1」のプロセッサ使用率143は「36%」である。プロセスID141「1001」のプロセス31は、デバイスID144「D1」のI/Fデバイス13aに対して、所定期間に「244回」(145)のI/O要求を発行し、デバイスID144「D2」のI/Fデバイス14bに対して、所定期間に「5回」(145)のI/O要求を発行した。

30

【0118】

図12は、プロセス31を何れのプロセッサ11で実行するかを決定する処理の例を示すフローチャートである。

【0119】

本処理は、I/Fデバイス13に関する複数のキュー20の内、専有キューとして利用されるキューの数が決まっている構成の例である。また、本処理は、1つのI/Fデバイス13に対する処理の例である。ストレージシステム1bが複数のI/Fデバイス13を備える場合、I/Fデバイス13毎に本処理が実行される。以下、1つのI/Fデバイス13を例に取り、図12の説明において、そのI/Fデバイス13を「対象I/Fデバイス」と言う。

40

【0120】

プロセス配置部37は、プロセス31毎にループAの処理を実行する。以下、1つのループAの処理を例に取り、そのループAの処理の対象となるプロセス31を「対象プロセス」という。どのようなプロセスが実行されているかについては、プロセススケジューラ

50

3 2 に問い合わせることによりわかる。

【 0 1 2 1 】

次に、プロセス配置部 3 7 は、ループ A の処理において、プロセッサ 1 1 のそれぞれについて、ループ B の処理を実行する。以下、1 つのループ B の処理を例に取り、そのループ B の処理の対象となるプロセッサ 1 1 を「対象プロセッサ」という。

【 0 1 2 2 】

次に、キュー割当部 3 5 は、プロセス別 I / O 数テーブル 1 4 0 を参照し、現状における対象 I / F デバイス 1 3 に関する排他数（前排他数）を算出する（S 6 0 3）。排他数については、図 6 において説明したとおりである。

【 0 1 2 3 】

前排他数は、次のように算出される。すなわち、例えば、図 3 のキュー管理テーブル 1 0 0 に示すように、デバイス ID 1 0 1 「D 2」の対象 I / F デバイス 1 3 b には、キュー ID 1 0 2 「Q 3」の専有キュー 2 0 c と、キュー ID 1 0 2 「Q 4」の共有キュー 2 0 d とが関連付けられている。そして、専有キュー 2 0 c には、プロセッサ ID 「U 3」のプロセッサ 1 1 c に割り当てられており、共有キュー 2 0 d には、プロセッサ ID 「U 1」及び「U 2」のプロセッサ 1 1 a 及び 1 1 b が割り当てられている。そして、図 1 1 のプロセス別 I / O 数テーブル 1 4 0 に示すように、プロセッサ 1 1 a（プロセッサ ID 1 4 2 「U 1」）で実行されているプロセス ID 1 4 1 「1 0 0 1」のプロセス 3 1 a から、デバイス ID 1 4 4 「D 2」の I / F デバイス 1 3 b に対する I / O 数 1 4 5 は「5 回」である。プロセッサ 1 1 b（プロセッサ ID 1 4 2 「U 2」）で実行されているプロセス ID 1 4 1 「1 0 0 2」のプロセス 3 1 b から、デバイス ID 1 4 4 「D 2」の I / F デバイス 1 3 b に対する I / O 数 1 4 5 は「1 8 9 回」である。プロセッサ 1 1 c（プロセッサ ID 1 4 2 「U 3」）で実行されているプロセス ID 1 4 1 「1 0 0 3」のプロセス 3 1 c から、デバイス ID 1 4 4 「D 2」の I / F デバイス 1 3 b に対する I / O 数 1 4 5 は「1 1 回」である。この場合、対象 I / F デバイス 1 3 b に関する前排他数は、共有キュー 2 0 d が割り当てられているプロセッサ 1 1 a 及び 1 1 b で実行されているプロセス 3 1 a 及び 3 1 b の I / O 数 1 4 5 の合計「5 + 1 8 9 = 1 9 4」となる。

【 0 1 2 4 】

次に、プロセス配置部 3 7 は、対象プロセッサ 1 1 に対象プロセス 3 1 を移動させた場合における、対象 I / F デバイス 1 3 に関する排他数（後排他数）を算出する（S 6 0 4）。

【 0 1 2 5 】

前排他数は、次のように算出される。すなわち、例えば、プロセス 3 1 b をプロセッサ 1 1 b からプロセッサ 1 1 c に移動させたと仮定した場合、プロセッサ 1 1 c が割り当てられているキュー 2 0 c は、プロセス 3 1 b 及び 3 1 c から利用されることになる。したがって、I / F デバイス 1 3 b（デバイス ID 「D 2」）に関する後排他数は、共有キュー 2 0 d に割り当てられているプロセッサ 1 1 a で実行されているプロセス 3 1 a（プロセス ID 「1 0 0 1」）の I / O 数 1 4 5 の「5」となる。

【 0 1 2 6 】

次に、プロセス配置部 3 7 は、後排他数と前排他数との差 Z 2 が所定の閾値 C 2（C 2 は 0 以下の値）未満であり、且つ、対象プロセッサ使用率 1 4 3 が閾値 C 3 未満であるか否かを判定する（S 6 0 5）。つまり、「Z 2（= 後排他数 - 前排他数）< 閾値 C 2 0」であり、且つ、「対象プロセッサ使用率 1 4 3 < 閾値 C 3」であるか否かを判定する。

【 0 1 2 7 】

S 6 0 5 の判定が肯定的な場合（S 6 0 5 : Y E S）、プロセス配置部 3 7 は、対象プロセス 3 1 を対象プロセッサ 1 1 に移動させる（S 6 1 0）。このとき、プロセス配置部 3 7 は、プロセス別 I / O 数テーブル 1 4 0 において、対象プロセス ID 1 4 1 に対象プロセッサ ID を対応付けるように更新する。なぜなら、プロセス 3 1 の移動先のプロセッサ 1 1 の使用率にも余裕があり、且つ、プロセス 3 1 を移動させた方が、I / F デバイス 1 3 に関連する排他数が減少する可能性が高いと考えられるからである。そして、プロセ

10

20

30

40

50

ス配置部 37 は、ループ B の処理を抜ける。

【 0 1 2 8 】

S 6 0 5 の判定が否定的な場合 ( S 6 0 5 : N O )、プロセス配置部 37 は、次の処理を行う。プロセス配置部 37 は、全てのプロセッサ 1 1 についてループ B の処理を完了した場合、ループ B の処理を抜け、まだ未処理のプロセッサ 1 1 が残っている場合、ループ B の処理を繰り返す。すなわち、プロセス配置部 37 は、対象プロセス 3 1 を移動させず、現状のままとする。なぜなら、プロセス 3 1 の移動先のプロセッサ 1 1 の使用率に余裕が無い、又は、プロセス 3 1 を移動しても I / F デバイス 1 3 に関する排他数が減少する可能性が低いと考えられるからである。

【 0 1 2 9 】

プロセス配置部 37 は、ループ B の処理を抜けた後、次の処理を行う。プロセス配置部 37 は、全てのプロセス 3 1 についてループ A の処理を完了した場合、ループ A の処理を抜けて本処理を終了し、まだ未処理のプロセス 3 1 が残っている場合、ループ A の処理を繰り返す。

【 0 1 3 0 】

複数のプロセッサのそれぞれでプロセスを実行することによりデータが発行されるようになっている第 2 の実施形態の説明によれば、例えば以下の説明をできる。

【 0 1 3 1 】

コントローラ 2 b は、複数のプロセスにそれぞれ対応した複数のデータ発行頻度に基づいて、専有キューに割り当てられている第 1 のプロセッサで実行するプロセスを決定する。ここで、プロセスのデータ発行頻度は、そのプロセスが I / F デバイスへデータを発行した頻度である。また、コントローラ 2 b は、その複数のデータ発行頻度のうち最大のデータ発行頻度に対応したプロセスを、専有キューに割り当てられている第 1 のプロセッサで実行するプロセスとしてよい。これにより、I / F デバイス 1 3 に対するデータの発行頻度が最大のプロセス 3 1 から発行されるデータについて、排他処理が不要となる。

【 0 1 3 2 】

コントローラ 2 b は、共有キューにデータが格納される頻度である共有頻度及び第 1 のプロセッサの負荷に関する所定の条件が満たされた場合、いずれかの第 2 のプロセッサである対象第 2 プロセッサで実行されているプロセスを、第 1 のプロセッサに移動させるプロセス移動を実行する。コントローラ 2 b は、所定の条件として、後共有頻度が前共有頻度よりも所定の閾値以上小さく、且つ、第 1 のプロセッサの負荷が所定の閾値未満であることを設定してもよい。ここで、後共有頻度とは、第 2 のプロセッサで実行されているプロセスを第 1 のプロセッサに移動させるプロセス移動を実行したと仮定した場合の共有頻度であってよい。そして、前共有頻度とは、プロセス移動前の共有頻度であってよい。これにより、排他処理の発生状況が変化した場合に、専有キュー 2 0 に割り当てられており、且つ、負荷が所定未満のプロセッサ 1 1 に、排他処理を発生させる可能性のより高いプロセス 3 1 を実行させることができる。結果として、プロセッサの負荷を考慮しつつ、キュー 2 0 の全体における排他処理の発生回数を減少させることができる。

【 0 1 3 3 】

上述した実施形態は、本発明の説明のための例示であり、本発明の範囲をそれらの実施形態にのみ限定する趣旨ではない。当業者は、本発明の要旨を逸脱することなしに、他の様々な態様で本発明を実施することができる。例えば、コントローラ 2 は、複数のプロセッサ 1 1 を含み、メモリ 1 2 及び I / F デバイス 1 3 は、コントローラ 2 の外に設けられていてもよい。

【 0 1 3 4 】

例えば、複数のプロセッサ 1 1 は、一つの集積回路に含まれる複数のコア ( マルチコア ) であってもよい。例えば、複数のプロセッサ 1 1 は、マルチスレッド C P U における複数のスレッドであってもよい。

【 符号の説明 】

【 0 1 3 5 】

10

20

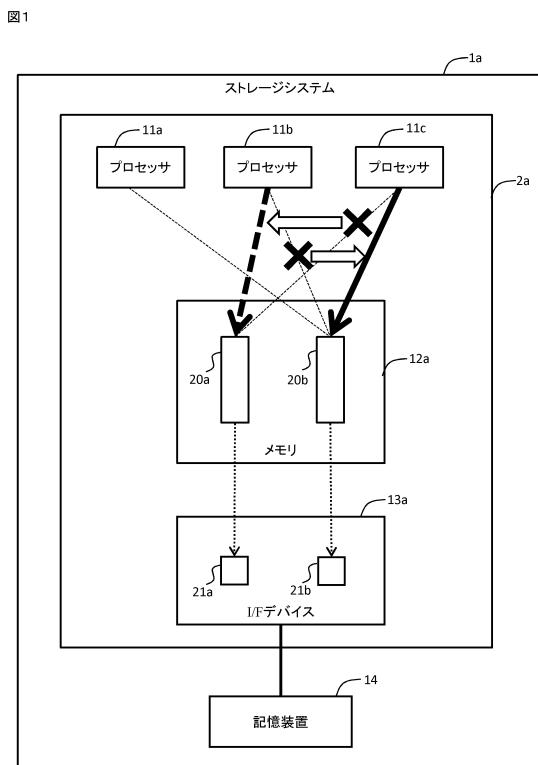
30

40

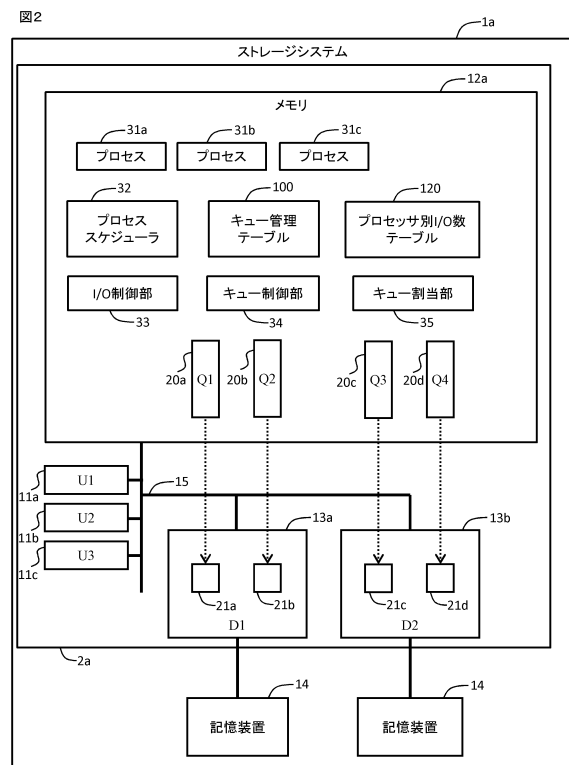
50

1 a、1 b : ストレージシステム 11 a、11 b、11 c : プロセッサ 12 a、12 b : メモリ 13 a、13 b : I/Fデバイス 20 a、20 b、20 c、20 d : キュー

【図1】



【図2】



【図3】

図3

キュー管理テーブル			
デバイスID	キューID	プロセッサID	モード
D1	Q1	U1	専有
	Q2	U2, U3	共有
D2	Q3	U3	専有
	Q4	U1, U2	共有

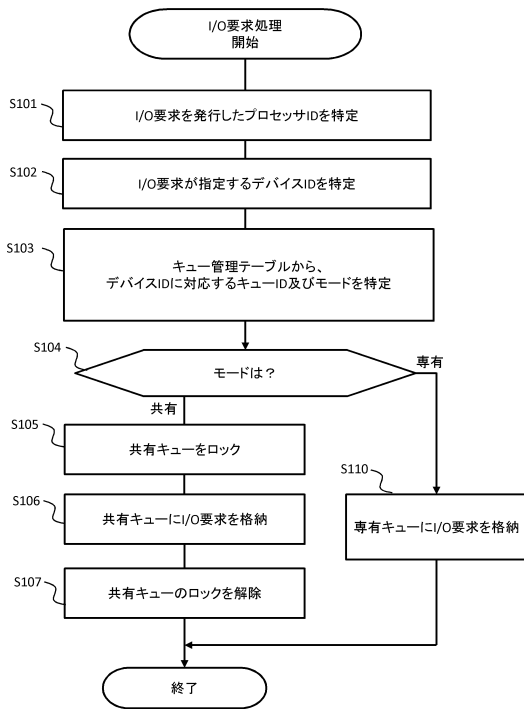
【図4】

図4

プロセッサ別I/O数テーブル		
プロセッサID	デバイスID	I/O数
U1	D1	431
	D2	5
U2	D1	19
	D2	189
U3	D1	2
	D2	11
...	...	...

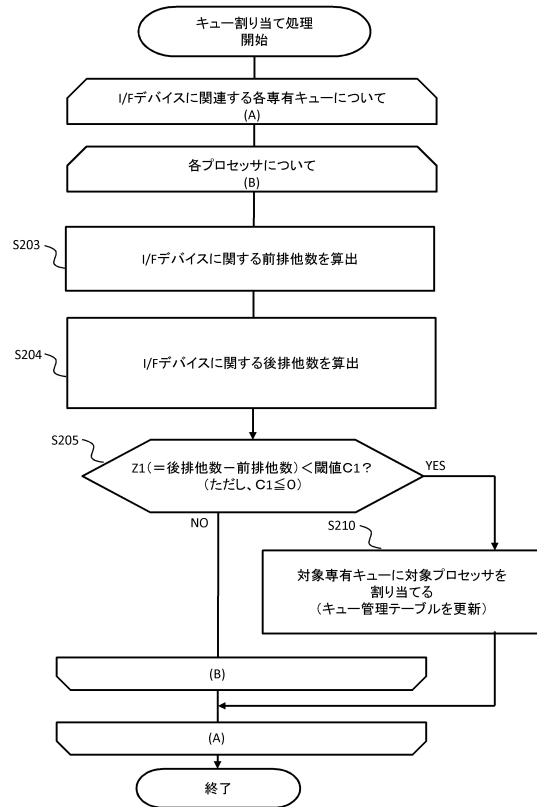
【図5】

図5

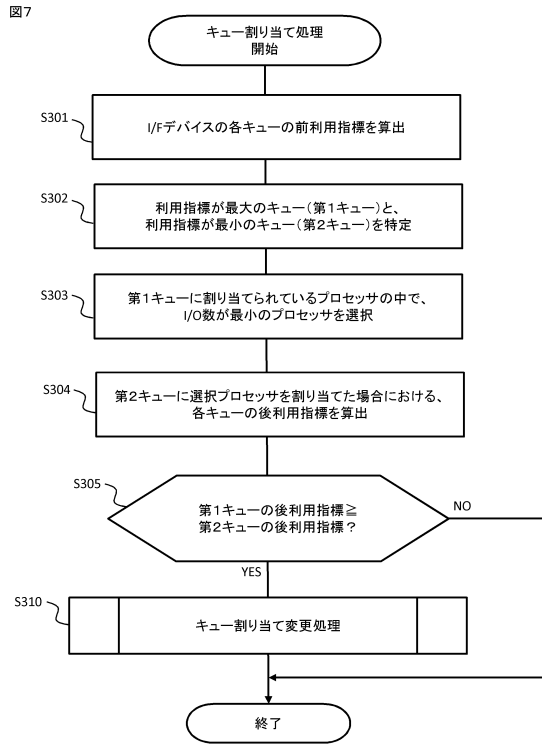


【図6】

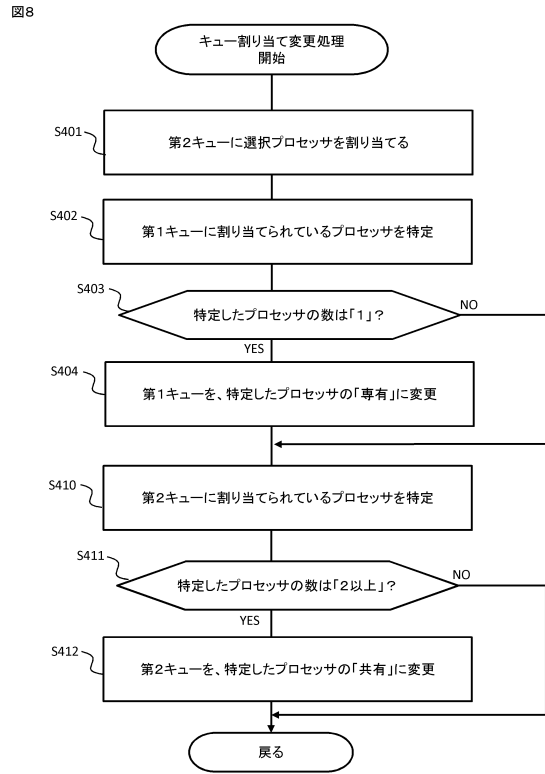
図6



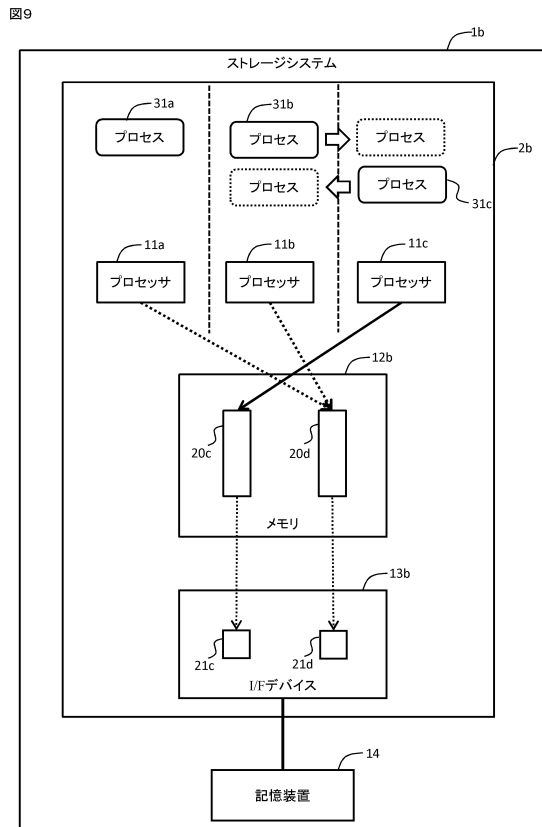
【図7】



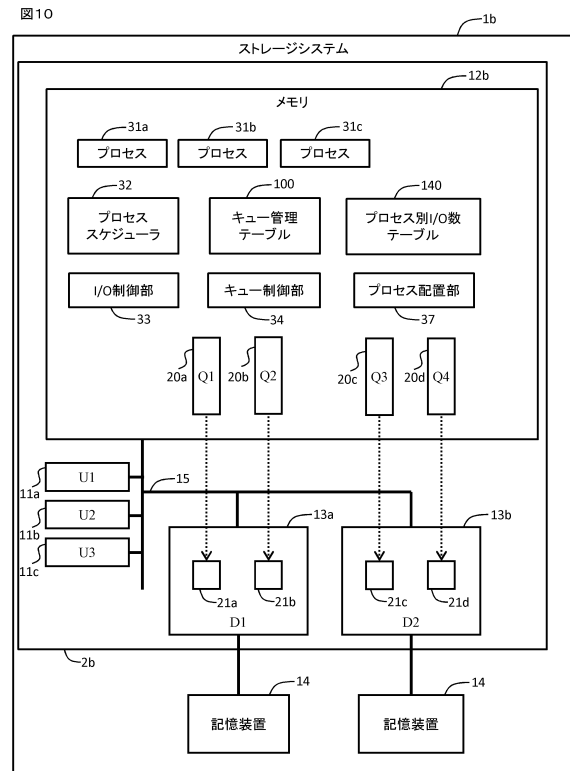
【図8】



【図9】



【図10】



【図11】

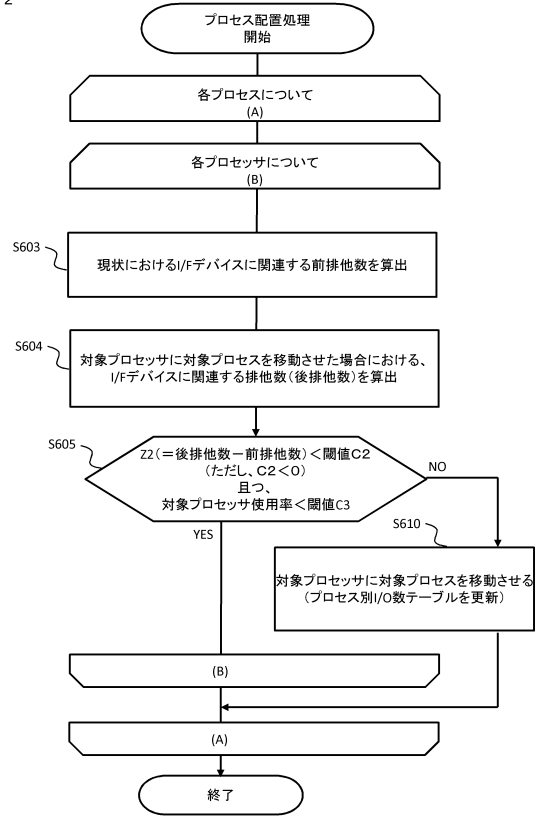
図11

140  
プロセス別I/O数テーブル

141 プロセスID	142 プロセッサID	143 プロセッサ使用率	144 デバイスID	145 I/O数
1001	U1	36	D1	244
			D2	5
1002	U2	15	D1	19
			D2	189
1003	U3	44	D1	2
			D2	11

【図12】

図12



---

フロントページの続き

(51)Int.Cl. F I  
G 0 6 F 12/00 5 1 4 A

(56)参考文献 特開2004-126694(JP,A)  
特開2002-063148(JP,A)  
特開2008-276326(JP,A)  
米国特許出願公開第2007/0248110(US,A1)

(58)調査した分野(Int.Cl., DB名)  
G 0 6 F 3 / 0 6 - 3 / 0 8、  
G 0 6 F 9 / 4 6 - 5 4、1 2 / 0 0  
G 0 6 F 1 3 / 1 0 - 1 3 / 1 4  
G 0 6 F 1 3 / 3 8 - 1 3 / 4 2