

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関  
国際事務局

(43) 国際公開日  
2012年10月4日(04.10.2012)



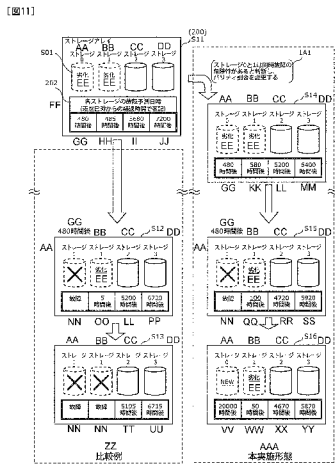
(10) 国際公開番号  
WO 2012/132408 A1

- (51) 国際特許分類:  
G06F 13/10 (2006.01) G06F 3/06 (2006.01)
- (21) 国際出願番号: PCT/JP2012/002116
- (22) 国際出願日: 2012年3月27日(27.03.2012)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:  
特願 2011-081205 2011年3月31日(31.03.2011) JP
- (71) 出願人(米国を除く全ての指定国について): パナソニック株式会社(PANASONIC CORPORATION) [JP/JP]; 〒5718501 大阪府門真市大字門真1006番地 Osaka (JP).
- (72) 発明者; および
- (75) 発明者/出願人(米国についてのみ): 大坪 紹二(OHTSUBO, Shohji). 阿部 敏久(ABE, Toshihisa). 幸 裕弘(YUKI, Yasuhiro). 寺田 吉希(TERADA, Yoshiki). 廣瀬 勝彦(HIROSE, Katsuhiko).
- (74) 代理人: 新居 広守(NII, Hiromori); 〒5320011 大阪府大阪市淀川区西中島5丁目3番10号タナカ・イトーピア新大阪ビル6階新居国際特許事務所内 Osaka (JP).
- (81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

[続葉有]

(54) Title: ARRAY MANAGEMENT APPARATUS, ARRAY MANAGEMENT METHOD, COMPUTER PROGRAM, AND INTEGRATED CIRCUIT

(54) 発明の名称: アレイ管理装置、アレイ管理方法、コンピュータプログラム、集積回路



200... TABLEAU DE MÉMOIRES  
1A1... LES MÉMOIRES 0 ET 1 SONT DÉTERMINÉES COMME PRÉSENTANT UN RISQUE DE DYSFONCTIONNEMENT EN MÊME TEMPS, ET LES RAPPORTS DE PARITÉ DE CELLES-CI SONT CHANGÉS

AA... MÉMOIRE 0  
BB... MÉMOIRE 1  
CC... MÉMOIRE 2  
DD... MÉMOIRE 3  
EE... DÉTERIORATION  
FF... DATE/HEURE PRÉDITE DE DYSFONCTIONNEMENT DE CHAQUE DES MÉMOIRES (EXPRIMÉE SOUS LA FORME D'UN TEMPS ÉCOULÉ À PARTIR DE LA DATE/HEURE ACTUELLE)  
GG... 480 HEURES PLUS TARD  
HH... 485 HEURES PLUS TARD  
II... 5850 HEURES PLUS TARD  
JJ... 7200 HEURES PLUS TARD  
KK... 590 HEURES PLUS TARD  
LL... 5200 HEURES PLUS TARD  
MM... 6400 HEURES PLUS TARD  
NN... DYSFONCTIONNEMENT  
OO... 5 HEURES PLUS TARD  
PP... 6720 HEURES PLUS TARD  
QQ... 100 HEURES PLUS TARD  
RR... 4720 HEURES PLUS TARD  
SS... 8920 HEURES PLUS TARD  
TT... 5160 HEURES PLUS TARD  
UU... 6715 HEURES PLUS TARD  
VV... 20000 HEURES PLUS TARD  
WW... 50 HEURES PLUS TARD  
XX... 4670 HEURES PLUS TARD  
YY... 5670 HEURES PLUS TARD  
ZZ... EXEMPLE COMPARATIF  
AAA... PRÉSENT MODE DE RÉALISATION

(57) Abstract: An array management apparatus (10) that manages an array (200) composed of a plurality of storages is provided with: a determining unit (100x) that determines whether a possibility of data stored in the array (200) being lost, due to at least two storages among the plurality of storages malfunctioning at the same time, exists; and a changing control unit (100y) that changes, when the possibility of data stored in the array (200) being lost is determined as existing, the malfunctioning timing of at least one storage of the aforementioned two storages. This array management apparatus (10) is thereby able to prevent data stored in the array from getting lost.

(57) 要約: 複数のストレージが構成するアレイ(200)を管理するアレイ管理装置(10)が、前記複数のストレージのうち少なくとも2台のストレージが同時期に故障することによって、アレイ(200)が保持するデータが消失する可能性があるか否かを判断する判断部(100x)と、アレイ(200)が保持するデータが消失する可能性があると判断された場合には、前記少なくとも2台のストレージのうち少なくとも一つの故障時期を変更する、変更制御部(100y)とを備え、アレイが保持するデータの消失を防ぐことができる。

WO 2012/132408 A1



(84) 指定国 (表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK,

SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

添付公開書類:

— 国際調査報告 (条約第 21 条(3))

## 明 細 書

発明の名称：

アレイ管理装置、アレイ管理方法、コンピュータプログラム、集積回路  
技術分野

[0001] 本発明は、複数のストレージが構成するアレイを管理するアレイ管理装置に関する。

### 背景技術

[0002] 従来技術として、複数のストレージが構成するアレイを管理するアレイ管理装置がある（特許文献1などを参照）。

[0003] アレイでは、例えば、RAID（Redundant Arrays of Inexpensive Disks）5、RAID6などの技術により、データが保持される。

### 先行技術文献

#### 特許文献

[0004] 特許文献1：特開平8-137633号公報

### 発明の概要

#### 発明が解決しようとする課題

[0005] しかしながら、上記従来技術では、アレイの状態が、将来、複数のストレージが同時期に故障することによって、アレイが保持するデータが消失してしまう状態（パリティを使用しても、アレイが保持するデータを復旧することができない状態のこと。RAID5の場合には、2台以上のストレージが同時故障する状態。RAID6の場合には、3台以上のストレージが同時故障する状態）になる可能性がある場合が生じる。この場合に、同時故障する可能性がある複数のストレージの故障時期をずらすことができず、アレイが保持するデータが消失する可能性が高くなってしまっていたという課題がある。

[0006] 本発明は、上記従来技術の課題を解決するものであり、アレイが保持するデータの消失を防ぐことを目的とする。

## 課題を解決するための手段

[0007] 上記目的を達成するために、本発明のアレイ管理は、複数のストレージが構成するアレイを管理するアレイ管理装置であって、前記複数のストレージのうち少なくとも2台のストレージが同時期に故障することによって、前記アレイが保持するデータが消失する可能性があるか否かを判断する判断部と、前記アレイが保持する前記データが消失する可能性があると判断された場合には、前記少なくとも2台のストレージのうちの少なくとも一方の故障時期を変更する、故障時期変更部とを備えるアレイ管理装置である。

[0008] つまり、例えば、前記故障時期変更手段は、前記各ストレージに対して書き込む、パリティとデータとの割合を変更することによって、前記少なくとも2台のストレージのうちの少なくとも一つの故障時期を変更する。

[0009] これにより、将来、複数のストレージが同時期に故障することによって、アレイが保持するデータが消失する可能性があると判断した場合には、同時故障する可能性がある複数のストレージの故障時期をずらされ、データの消失を防ぐことができる。

[0010] なお、これらの全般的または具体的な態様は、システム、方法、集積回路、コンピュータプログラムまたは記録媒体で実現されてもよく、システム、方法、集積回路、コンピュータプログラムおよび記録媒体の任意な組み合わせで実現されてもよい。

## 発明の効果

[0011] 将来、複数のストレージが同時期に故障することによって、アレイが保持するデータが消失する可能性があると判断した場合には、同時故障する可能性がある複数のストレージの故障時期をずらす。このことで、アレイが保持するデータ消失を防ぐことができる。

## 図面の簡単な説明

[0012] [図1]図1は、本システムを示す図である。

[図2]図2は、本システムを示す図である。

[図3]図3は、ストレージのアドレス空間の概念図である。

- [図4]図4は、4つのストレージなどを示す図である。
- [図5]図5は、4つのストレージについての表を示す図である。
- [図6]図6は、ストレージなどを示す図である
- [図7]図7は、アドレス空間などを示す図である。
- [図8]図8は、本システムを示す図である。
- [図9]図9は、同時故障回避の処理のフローチャートである。
- [図10A]図10Aは、S1402での動作の詳細を示すフローチャートである。
- 。
- [図10B]図10Bは、2つの時刻の間の差1F1、1F2などを示す図である。
- 。
- [図10C]図10Cは、2つの時刻の間の差1F3、1F4などを示す図である。
- 。
- [図11]図11は、故障が生じるまでの時間などを示す図である。
- [図12]図12は、アレイ管理部の構成を示す図である。
- [図13]図13は、本システムの動作のフローチャートである。
- [図14]図14は、本システムの動作のフローチャートである。
- [図15]図15は、処理のそれぞれの段階での表を示す図である。
- [図16]図16は、S1102における動作の詳細を示すフローチャートである。
- [図17]図17は、S1203における動作の詳細を示すフローチャートである。
- [図18]図18は、2つのグラフなどを示す図である。
- [図19]図19は、S1401での動作の詳細を示すフローチャートである。
- [図20]図20は、S1501での動作の詳細を示すフローチャートである。
- [図21]図21は、S1503での動作の詳細を示すフローチャートである。
- [図22]図22は、データ量のグラフを示す図である。
- [図23]図23は、3つのグラフなどを示す図である。
- [図24]図24は、S1404での動作の詳細を示すフローチャートである。

- [図25]図25は、S1901での動作の詳細を示すフローチャートである。
- [図26]図26は、S2002での動作の詳細を示すフローチャートである。
- [図27]図27は、S2003での動作の詳細を示すフローチャートである。
- [図28]図28は、上欄、中央欄、下欄の3つの欄の表を示す図である。
- [図29]図29は、保存可能なデータ量などを示す表である。
- [図30]図30は、発生頻度などを示す表である。
- [図31]図31は、保存可能なデータ量などを示す図である。
- [図32]図32は、処理効率などを示す図である。
- [図33]図33は、ブロック書き換え量のグラフを示す図である。
- [図34]図34は、ブロック書き換え量などを示す図である。
- [図35A]図35Aは、パリティ割合などを示す図である。
- [図35B]図35Bは、4つのストレージなどを示す図である。
- [図35C]図35Cは、書き換え量などを示す表を示す図である。
- [図36]図36は、パリティ割合などを示す表を示す図である。
- [図37]図37は、3つのグラフなどを示す図である。
- [図38]図38は、4つのグラフなどを示す図である。
- [図39]図39は、S2004での動作の詳細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

- [0013] 以下、本発明の実施形態におけるシステム1について、図面を参照しながら説明する。
- [0014] 実施形態のレイ管理装置10は、複数のストレージ（ストレージ21、22など）が構成するレイ200を管理するレイ管理装置10であって、前記複数のストレージのうち少なくとも2台のストレージ（ $\alpha$ 、 $\beta$ ）が同時期に故障することによって、前記レイ200が保持するデータが消失する可能性があるか否かを判断する判断部100xと（閾値1F2、差1F1などを参照）、レイ200が保持するデータが消失する可能性があるかと判断された場合には、前記少なくとも2台のストレージのうちの少なくとも一つ（例えば $\alpha$ ）の故障時期（日時1tを参照）を変更する故障時期変更部1

00yとを備えるアレイ管理装置10である。

[0015] そして、故障時期変更部100yは、各ストレージに対して書き込む、パリティ（領域201p参照）とデータ（領域201p参照）との割合を変更することによって、前記少なくとも2台のストレージのうちの少なくとも一つ（例えば $\alpha$ ）の故障時期を変更する。

[0016] 例えば、前記判断部100xは、前記複数のストレージのうち少なくとも2台のストレージが同時期に故障する可能性があるか否かを判断する同時故障判断部（第1部分100x1）と、前記少なくとも2台のストレージが同時期に故障すると判断された場合に、前記少なくとも2台のストレージが同時期に故障することによって、前記アレイ200が保持するデータが消失する可能性があるか否かを判断するアレイ故障判断部（第2部分100x2）とを含んでもよい。

[0017] 例えば、前記アレイ200により保持される前記データの消失は、第1の前記ストレージが故障する日時が第1の日時である場合には、発生し難く、第2の日時である場合には発生し易く、第1の前記ストレージは、記憶されるデータとパリティとの全体に占める、前記パリティの割合が、第1の割合である場合と、第2の割合である場合とがあるストレージであり、第1の前記ストレージの前記日時は、前記第1の割合である場合には、消失が発生し難い前記第1の日時であり、前記第2の割合である場合には、発生し易い前記第2の日時であり、前記判断部100xは、2つの割合のうち一方の前記割合を前記第1の割合と特定し、他方の前記割合を前記第2の割合と特定し、前記故障時期変更部100yは、第1の前記ストレージでの前記割合を、特定された一方の前記割合にさせ、特定された他方の前記割合にはさせなくてもよい。

[0018] 例えば、消失が発生し難い前記第1の日時は、第2の前記ストレージが故障する日時から閾値以上に離れた日時であり、消失が発生し易い前記第2の日時は、第2の前記ストレージが故障する日時から閾値以内の日時でもよい。

- [0019] なお、これらの全般的または具体的な態様は、システム、方法、集積回路、コンピュータプログラムまたは記録媒体で実現されてもよく、システム、方法、集積回路、コンピュータプログラムまたは記録媒体の任意な組み合わせで実現されてもよい。
- [0020] 以下、詳しく説明される。
- [0021] (同時故障回避)
- 図1は、システム1の実施イメージ図である。
- [0022] 図2は、システム1の構成図である。
- [0023] システム1は、アレイ管理装置10と、アレイ200(図2参照)とを含む。
- [0024] アレイ管理装置10は、アレイ200の動作を制御する。
- [0025] なお、アレイ管理装置10は、図1に示されるように、例えば、DVD(Digital Video Disc)レコーダ、BD(Blu-ray Disc)レコーダの全体または一部などでもよい。つまり、例えば、アレイ管理装置10は、放送波から受信される動画のデータなどのデータを記録するストレージ(図1のストレージ25、26などを参照)備えるストレージ装置の全体または一部などでもよい。
- [0026] アレイ(ストレージアレイ)200は、データを保持する。
- [0027] なお、例えば、アレイ管理装置10は、アレイ200により保持されるデータを、アレイ200と、アレイ管理装置10との間で、転送させてもよい。
- [0028] なお、アレイ管理装置10は、例えば、この転送をさせる下位側転送制御部55(図2)を備えてもよい。
- [0029] そして、例えば、アレイ管理装置10は、アレイ200により保持されるデータを、アレイ管理装置10と、保持されるデータの処理をする処理部との間で、転送させてもよい。
- [0030] アレイ管理装置10は、例えば、この転送をさせる上位側転送制御部54(図2)を備えてもよい。

- [0031] そして、例えば、上述の処理部は、アレイ管理装置10の外部にあるパーソナルコンピュータなどの装置などで動作するオペレーションシステムにおけるファイルシステムなどでもよい。
- [0032] なお、アレイ200により保持されるデータは、例えば、アレイ200と、このオペレーションシステムの上で動作するアプリケーションとの間を、上述ファイルシステムなどを介して転送されて、このアプリケーションにより利用されてもよい。
- [0033] アレイ200は、ストレージ201（図2）を複数個、含んでなり、これら複数のストレージ201を用いて、データの保持を行う。
- [0034] 含まれる複数のストレージ201のうちの少なくとも1つは、例えば、図1のストレージ25、26などで示されるような、アレイ管理装置10に設けられたストレージなどでもよい。
- [0035] また、少なくとも1つは、例えば、図1のストレージ21～24で示されるような、アレイ管理装置10の外部にあるストレージでもよい。
- [0036] 具体的には、少なくとも1つは、ストレージ21～23で示されるような、アレイ管理装置10が設けられた住宅と同じ住宅に設けられたストレージでもよい。
- [0037] また、少なくとも1つは、ストレージ24で示されるような、他の住宅に設けられたストレージでもよい。なお、例えば、他の住宅は、親夫婦の住む住宅である一方で、アレイ管理装置10が設けられた住宅は、子供夫婦が住む住宅でもよい、なお、他の住宅にあるストレージ24は、例えば、インターネットなどであるネットワーク1DLを介して、アレイ管理装置10からの、動作の制御を受けてもよい。
- [0038] なお、図1は、単なる一例である。システム1においては、図1に示される通りでもよいし、図1に示される通りでなくてもよい。
- [0039] 図3は、ストレージのアドレス空間の概念図である。
- [0040] それぞれのストレージ201（例えば、図3のストレージ3）の記憶領域は、複数のブロック $1 \times B2$ （ $i$ 個のブロック $Bk03 \sim Bki3$ などを参

照)に分かれる。

[0041] そして、アレイ200は、それぞれのストレージ201(図3のストレージ0~3)におけるブロック(bk00、bk01、bk02、bk03参照)が集まってなる記憶領域であるストライプ1xB1(図3のストライプ0など)を有する。

[0042] 図4は、ストレージ201に対するデータ・パリティ格納の概念図である。

[0043] ストライプ1xB1(例えば、図4のストライプi)に含まれるブロック1xB2(図3)としては、データのブロック1xB2d(ブロックBDi0、BDi1、BDi2などを参照)と、パリティのブロック1xB2p(ブロックBPi0などを参照)とがある。

[0044] ブロック1xB2pに記憶されるパリティは、そのブロック1xB2pのあるストレージ201以外の他のストレージ201(例えば図4のストレージ2)で故障が生じた場合に用いられる。

[0045] つまり、故障が生じた場合に、そのブロック1xB2pに記憶されるパリティと、そのブロック1xB2pのストライプ1xB1における、故障が生じてないそれぞれのストレージ(ストレージ0、1)でのデータ(ブロックBdi0、Bdi1参照)とが用いられる。

[0046] すなわち、これらから、そのストライプ1xB1における、故障が生じたストレージ(ストレージ2)におけるデータ(ブロックBDi2のデータ)が生成される。

[0047] このため、ストライプ1xB1におけるブロック1xB2dへとデータが書き込まれる際には、そのブロック1xB2dにデータが書き込まれるのと共に、そのストライプ1xB1のブロック1xB2pに対して、パリティが書き込まれる。

[0048] そして、ストライプ1xB1における、データが記憶されるブロック1xB2dの個数は、例えば、図4における3個などのように、2個以上である。

- [0049] このため、パリティのブロック  $1 \times B_2 p$  への書き込みが行われる回数は、その書き込みが、何れの、データのブロック  $1 \times B_2 d$  への書き込みがされる際にも行われて、比較的多い回数である。
- [0050] つまり、それぞれの、データのブロック  $1 \times B_2 d$  への書き込みの回数は、比較的少ない。
- [0051] パリティのブロック  $1 \times B_2 p$  (例えば  $B_{P i 0}$ ) への書き込みの回数は、例えば、それぞれの、データのブロック  $1 \times B_2 d$  ( $B_{D i 0} \sim B_{D i 2}$ ) への書き込みの回数が合計された回数などである。つまり、例えば、図4の例では、データのブロック  $1 \times B_2 d$  の個数が3個であり、データのブロック  $1 \times B_2 d$  への書き込みの回数の3倍などである。
- [0052] 図5は、それぞれのストレージ  $201$  でのパリティ割合  $201 r$  (第7行) などを示す表の図である。
- [0053] 図5の表におけるそれぞれの列では、その列のストレージ  $201$  のことが示される。
- [0054] 前述の通り、ストレージ  $201$  におけるブロック  $1 \times B_2$  としては、データのブロック  $1 \times B_2 d$  と、パリティのブロック  $1 \times B_2 p$  とがある。
- [0055] 図6は、ストレージ  $201$  における、データの領域  $201 d$  と、パリティの領域  $201 p$  などを示す図である。
- [0056] つまり、ストレージ  $201$  の記憶領域の全体である領域  $201 A$  は、データの領域  $201 d$  と、パリティの領域  $201 p$  とに分かれる。データの領域  $201 d$  は、それぞれの、データのブロック  $1 \times B_2 d$  の領域が集まってなる領域である(図4参照)。そして、パリティの領域  $201 p$  は、それぞれの、パリティのブロック  $1 \times B_2 p$  の領域が集まってなる領域である。
- [0057] パリティ割合  $201 r$  (図5) は、全体の領域  $201 A$  の大きさに対する、パリティの領域  $201 p$  の割合をいう。つまり、パリティ割合  $201 r$  は、例えば、ストレージ  $201$  に含まれる、データのブロック  $1 \times B_2 d$  の個数と、パリティのブロック  $1 \times B_2 p$  の個数との合計に占める、パリティのブロック  $1 \times B_2 p$  の個数の割合などである。

- [0058] ストレージ201のパリティ割合201rが、比較的低い割合である場合には、パリティのブロック $1 \times B2p$ が比較的少ない。パリティブロックは、データブロックと比較して、アクセス頻度が高いため、パリティ割合が少なければ、パリティ割合が多いストレージと比較して、ストレージの劣化が起こりにくく、ストレージ201が故障する日時 $1t$ （図10B）が、比較的遅い（図示される日時 $1t2$ を参照）。
- [0059] 一方で、パリティ割合201rが、比較的高い方の割合である場合には、逆に、故障する日時 $1t$ が、比較的早い（日時 $1t1$ を参照）。
- [0060] こうして、ストレージ201のパリティ割合201rが一方の割合である場合における、ストレージ201の、故障の日時 $1t$ は、他方の割合である場合における日時 $1t$ とは異なる。
- [0061] 図7は、アレイ200におけるアドレス空間の概念図である。
- [0062] なお、上述の通りであるので、アレイ200におけるアドレス空間200asは、それぞれの、ストライプ $1 \times B1$ （図3：例えば図4のストライプ0、1など）での部分に分かれる。そして、それぞれのストライプ $1 \times B1$ （例えばストライプ1）での部分は、さらに、そのストライプ $1 \times B1$ における、それぞれのブロック $1 \times B2$ （図7におけるブロックBD10、BD11、BD12を参照）での部分に分かれる。
- [0063] 図8は、システム1を示す図である。
- [0064] アレイ管理装置10は、アレイ管理部100（図8）を備える。
- [0065] アレイ管理部100は、パリティ割合管理部101と、パリティ割合算出部102と、パリティ割合変更部103と、故障予測部104とを備える。
- [0066] なお、アレイ管理装置10は、図2に示されるように、例えば、CPU52、RAM53、ROM51などを含んでなるコンピュータを含んでもよい。
- [0067] そして、アレイ管理部100（アレイ管理部100が備える、パリティ割合管理部101などのそれぞれの機能ブロック）は、このコンピュータにより、コンピュータプログラム101pg（図2）が実行されることで、アレ

イ管理装置 10 に実現される機能の機能ブロックでもよい。

[0068] なお、このコンピュータプログラム 101 pg は、例えば、ROM 51 により記憶され、記憶されたコンピュータプログラム 101 pg が、上記のコンピュータにより実行されてもよい。

[0069] なお、図 2 では、説明の便宜上、パリティ割合管理部 101 などが、この ROM 51 の位置に図示される。

[0070] 図 9 は、同時故障回避の処理のフローチャートである。

[0071] S 1401 では、それぞれのストレージ 201 で故障が生じる日時である故障予測日時が故障予測部 104 により算出される。S 1402 では、同時故障が発生する可能性があるか否かが、パリティ割合算出部 102 により算出される。

[0072] 図 10A は、同時故障発生危険性判断の処理のフローチャートである。図 9 の S 1402 では、例えば、この図 10A の処理がされてもよい。

[0073] 図 10B は、第 1 の日時  $t_1$ 、第 2 の日時  $t_2$  などを示す図である。

[0074] S 1801 では、最も早い故障予測日時（図 10B の日時  $t$  を参照）が算出されたストレージ 201 が、第 1 のストレージ  $\alpha$ （図 10B 参照）として特定される。

[0075] S 1802（図 10A）では、第 2 のストレージ  $\beta$  が特定される。第 2 のストレージ  $\beta$  は、特定された、複数のストレージ 201 の故障予測日時のうちで、 $M+1$  番目に早い故障予測日時（図 10B の日時  $t_B$ ）が算出されたストレージ 201 である。

[0076] この  $M$  は、アレイ 200 における故障許容台数（最大故障許容台数）である。故障許容台数は、その台数以下の台数のストレージ 201 で故障が生じるだけであれば、アレイ 200 により保持されるデータの消失が生じず、その台数よりも多い台数での故障が生じた場合に、消失が生じる台数である。つまり、故障許容台数は、例えば、アレイ 200 が、RAID 5 での動作をする際には、1 台であり、RAID 6 での動作をする際には、2 台である。

[0077] つまり、第 2 のストレージ  $\beta$  は、その第 2 のストレージ  $\beta$  が故障するまで

に生じた、それぞれのストレージ201での故障が継続するまで、その第2のストレージ $\beta$ も故障すると、データの消失が発生する1以上のストレージ201のうちで、故障予測日時が最も早いストレージ201である。

[0078] S1803 (図10A)では、第1のストレージ $\alpha$ の故障予測日時(日時1t)と、第2のストレージ $\beta$ の故障予測日時(日時1B)との間の差1F1が、閾値1F2以内であるか否かが判定される。

[0079] つまり、第1のストレージ $\alpha$ の故障予測日時(日時1t)が、第2のストレージ $\beta$ の故障予測日時(日時1B)から、閾値1F2以内で、第2のストレージ $\beta$ の故障予測日時(日時1B)と同時期であるか否かが判定される。

[0080] 閾値1F2は、上述の差1F1が、その閾値1F2以内である場合には(日時1t2参照)、第1のストレージ $\alpha$ の故障予測時刻(日時1t)以後で、かつ、第2のストレージ $\beta$ が故障する故障予測時刻(日時1B)よりも前に、この第1のストレージ $\alpha$ などの、故障したストレージ201の交換、修理などがされず、データの消失が発生する(と推測される)値の閾値である。

[0081] つまり、閾値1F2は、上述の差1F1が、その閾値1F2より大きい場合には、その大きな差1F1の期間のうちに、故障したストレージ201の交換などがされて、データの消失が発生しない(と推測される)値の閾値である。

[0082] S1804 (図10A)では、上述の差1F1が、閾値1F2以内と判定される場合に(S1803: YES)、データの消失が発生する(と推測される)ことが判定される。他方、差1F1が、閾値1F2以内ではないと判定される場合には(S1803: NO)、データの消失が発生しない(と推測される)ことが判定される。

[0083] 図10Cは、差1F3、差1F4などを示す図である。

[0084] なお、例えば、図10Cで示される通りでもよい。

[0085] つまり、上述の第1のストレージ $\alpha$ は、例えば、図10Cで書かれる第1のストレージ $\gamma$ 1で、第2のストレージ $\beta$ は、第2のストレージ $\gamma$ 2でもよ

い。そして、第3のストレージ $\gamma$ は、その第3のストレージ $\gamma$ の故障予測日時が、第1のストレージ $\gamma$ 1の故障予測日時以後で、かつ、第2のストレージ $\gamma$ 2以前であるストレージ201である。

[0086] そして、例えば、第1のストレージ $\gamma$ 1と、第3のストレージ $\gamma$ との間の差1F3が、所定の閾値より大きい場合にのみ、消失が発生しないとの判定がされ(S1803:NO)、閾値より小さい場合には、消失が発生するとの判定がされてもよい(S1803:YES)。

[0087] 同様に、例えば、第3のストレージ $\gamma$ と、第2のストレージ $\gamma$ 2との間の差1F4が、所定の閾値より大きい場合にのみ、消失が発生しないとの判定がされ(S1803:NO)、閾値より小さい場合には、消失が発生するとの判定がされてもよい(S1803:YES)。

[0088] S1403(図9)では、データの消失が発生する可能性があるか否かの判定(S1402)に基づいて、行われる処理が選択される。

[0089] つまり、消失が発生しないと判定される場合には(S1403:NO)、現在の直前に行われた、アレイ200への制御と同じ制御が、パリティ割合変更部103により行われてもよい。

[0090] 他方、S1404では、消失が発生すると判定される場合に(S1403:YES)、現在の直前に行われた、アレイ200への制御とは別の制御が、パリティ割合変更部103により行われる。別の制御は、その制御がされれば、第1のストレージ $\alpha$ が故障する日時1tが、第2のストレージ $\beta$ が故障する日時1Bから閾値1F2より離れていて、同時期ではなく、データの消失が発生しない(と推測される)制御である。つまり、別の制御は、例えば、その制御がされれば、アレイ200に含まれるストレージ201におけるパリティ割合201r1が、現在の直前に行われた制御がされる際での、そのストレージ201におけるパリティ割合201r2と違っていて、上述の差1F1が比較的大きく、故障の時期が同時期ではない制御ではある。

[0091] 図11は、同時故障回避のイメージ図である。

[0092] 図11における、本実施形態の欄で示されるように、上述の技術によれば

、データの消失が生じると推測されることが判定される場合には（S 1 4 0 3 : Y E S、図 1 1 の処理 1 A 1）、別の制御がされて（S 1 4 0 4）、行われる制御が変更されることにより、データの消失が発生してしまうことが回避できる（図 1 1 の S 1 4 ~ S 1 6）。

[0093] これに対して、図 1 1 に示される比較例は、このような本実施形態とは別の、想定されるその他の技術におけるケースを示す。この比較例では、上述された、同時期か否かの判定（図 9 : S 1 4 0 2）がされず、行われる制御が変更されなくて、データの消失が発生してしまう（図 1 1 の S 1 2 ~ S 1 3）。

[0094] こうして、比較例では、データの消失が発生してしまう一方で、本実施形態では、データの消失の発生が回避できる。

[0095] 図 1 2 は、アレイ管理装置 1 0 を示す図である。

[0096] なお、アレイ管理装置 1 0 は、例えば、判断部 1 0 0 x と、変更制御部（故障時期偏向部） 1 0 0 y とを備えてもよい。

[0097] そして、判断部 1 0 0 x は、上述された、パリティ割合算出部 1 0 2（図 8 など）と、故障予測部 1 0 4 とを含み、変更制御部 1 0 0 y は、パリティ割合管理部 1 0 1 と、パリティ割合変更部 1 0 3 とを含んでもよい。

[0098] つまり、例えば、パリティ割合算出部 1 0 2 により、パリティ割合算出部 1 0 2 の上述の動作がされることにより、その動作が、パリティ割合算出部 1 0 2 を用いて、判断部 1 0 0 x により行われてもよい。なお、故障予測部 1 0 4 などについても、それぞれ、この例と同様である。

[0099] （受付データ書き込み）

図 1 3 は、受付データ書き込みの処理のフローチャートである。

[0100] 本システム 1 においては、例えば、図 1 3 で示される処理が行われてもよい。

[0101] S 1 0 0 1 では、アレイ 2 0 0 に対して書き込むことが要求されるデータが、アレイ管理装置 1 0 により受け付けられる。S 1 0 0 2 では、受け付けられたデータが記録されるストライプ 1 x B 1（図 3 参照）が特定される。

S 1 0 0 3では、そのデータのパリティが、アレイ管理装置10により生成される。S 1 0 0 4では、パリティ割合管理部101により記憶される情報1011（図8）に基づく処理が行われる。

[0102] ここで、情報1011は、特定されたストライプ $1 \times B 1$ における、受け付けられたデータが記録される（1以上の）ブロック $1 \times B 2$ と、生成された、そのデータのパリティが記録されるブロック $1 \times B 2$ とを、それぞれ特定する情報である。

[0103] なお、情報1011は、例えば、図5の表の第2行～第5行におけるデータ構造を有しても良い。情報1011は、図11に示される表などでもよい。

[0104] S 1 0 0 5では、受け付けられたデータが、特定される、そのデータが記録されるべき（1以上の）ブロック $1 \times B 2 d$ に記録され、生成されたパリティが、特定される、そのパリティが記録されるべきブロック $1 \times B 2 p$ に記録される。

[0105] （初期化）

図14は、パリティ割合の初期化の処理のフローチャートである。

[0106] 図15は、データ・パリティ対応情報作成作業における対応表（先述の図5参照）の遷移の例を示す図である。

[0107] S 1 1 0 1では、それぞれのストレージ201におけるパリティ割合201rの目標値（図15のそれぞれの表における第7行）が、アレイ管理装置10により特定される。なお、各ストレージのパリティ割合の目標値は、各ストレージにおいて非均等であっても良い。S 1 1 0 2では、特定されたそれぞれの目標値に基づいて、情報1011が、アレイ管理装置10により作成される。

[0108] 図16は、パリティ割合の初期化の処理のフローチャートである。先述された図14のS 1 1 0 2では、例えば、この図16の処理がされてもよい。

[0109] S 1 2 0 1では、情報1011の表（図5）における、それぞれのストレージ（各列）のパリティ割合201r（第7行）が0に初期化される。S 1

202では、最初のストライプ（例えば、図5のストライプ0）が、選択ストライプとして選択される。S1203では、情報1011の表における、選択された選択ストライプの行（例えば、ストライプ0の行）の、それぞれのストレージ（各列）の欄の値が特定される。

[0110] S1204では、全てのストライプの処理が終了したか否かが判定される。S1205では、終了していないと判定された場合に（S1204：NO）、選択ストライプとして、直前のS1203の処理における選択ストライプの次のストライプが選択される（図15の処理1C1～1C2などを参照）。一方で、全てのストライプの処理が終了したと判定される場合には（S1204：YES）、図16の処理が終了する。

[0111] 図17は、ストライプに対するデータ・パリティ対応情報割り当ての処理のフローチャートである。図16のS1203では、例えば、この図17の処理がされてもよい。

[0112] S1301では、所定の台数のストレージ201が選択される。

[0113] つまり、先述された、パリティ割合201rの目標値（S1101）から、現在の割合（図15の表に第6行）が減じられた値（第8行）の順序が、一番大きい順序から、上述の所定の台数番目の順序までの、それぞれの順序のストレージ201が選択される。

[0114] なお、上述の所定の台数は、アレイ200がRAID5で動作するならば、1台であり、RAID6で動作するならば、2台である。

[0115] S1302では、この図17の処理がされる対象のストライプにおける、選択された、所定の台数のストレージ201のうちのそれぞれのストレージ201のブロック1×B2が、パリティのブロック1×B2pとして特定される（図3、図4など参照）。

[0116] そして、このS1302では、他のそれぞれのストレージ201のブロック1×B2が、データのブロック1×B2dとして特定される（図15の処理1C1、1C2などを参照）。

[0117] S1303では、上述のS1301、S1302での特定の結果が反映さ

れた後における、それぞれのストレージ 201 での、現在の割合（図 15 の表の第 6 行）が算出される。

[0118] （故障予測日時算出の具体例）

図 18 は、故障予測日時算出の処理のイメージ図である。

[0119] 図 19 は、故障予測日時算出の処理のフローチャートである。先述された図 9 の S 1401 では、例えば、この図 19 の処理がされてもよい。

[0120] S 1501 では、図 18 の下欄に示される 2 つのグラフのうちの、左のグラフで示される関数  $f$  が特定される。

[0121] 図 18 で示されるように、関数  $f$  は、日時と、その日時における、アレイ 200 への書き込み量の累積量との関係の関数である。

[0122] S 1502 では、最初のストレージが、選択ストレージとして選択される。

[0123] S 1503 では、選択された選択ストレージについての、図 18 の下欄の右のグラフで示される関数  $g_j$  が特定される。図 18 で示されるように、関数  $g_j$  は、アレイ 200 に対する書き込みの累積量（下欄のグラフの横軸）と、アレイ 200 に対して、その累積量の書き込みがされた日時における、選択ストレージへの書き込み量の累積量との関係を特定する関数である。

[0124] S 1504 では、選択ストレージにおける保証量が特定される。保証量は、その量のデータ量までの書き込みがされるまでは、その選択ストレージで故障が生じないことが、その選択ストレージのメーカーなどにより保証された量である。つまり、保証量は、この量のデータ量だけの書き込みがされた日時（に近い日時）に、故障が生じることが推定される量である。

[0125] S 1505 では、選択ストレージに対して、保証量のデータ量だけの書き込みがされる際における、アレイ 200 に対する書き込みの累積量が特定される。

[0126] S 1506 では、アレイ 200 に対する書き込みの累積量が、特定された上述の累積量になる日時が、その選択ストレージの故障予測日時として特定される。

- [0127] S 1 5 0 7では、全てのストレージの故障予測日時が、上述のS 1 5 0 6で特定されたか否かが判定される。そして、されていると判定された場合には（S 1 5 0 7 : Y E S）、図 1 9の処理が終了する。他方、S 1 5 0 8では、全てのストレージについての処理が終わってないと判定される場合に（S 1 5 0 7 : Y E S）、直前のS 1 5 0 3～S 1 5 0 7での処理での選択ストレージの次のストレージ2 0 1が、次のS 1 5 0 4～S 1 5 0 7での処理での選択ストレージとして特定される。
- [0128] 図 2 0は、累積書き込みデータ量の増加度合い算出の処理のフローチャートである。上述された図 1 9のS 1 5 0 1では、例えば、この図 2 0の処理がされてもよい。
- [0129] S 1 6 0 1では、現在の日時が取得される。なお、例えば、アレイ管理装置 1 0に設けられた、現在の日時を計時する時計から、計時された、現在の日時が取得されてもよい。
- [0130] S 1 6 0 2では、アレイ 2 0 0に対する書き込みの、現在における累積量が取得される。なお、例えば、アレイ 2 0 0などの、アレイ管理装置 1 0の外部に設けられた所定の記憶部により、この累積量が記憶されて、記憶された累積量が、アレイ管理装置 1 0へと取得されてもよい。
- [0131] S 1 6 0 3では、アレイ管理装置 1 0の記憶領域に、取得された、現在の日時と、アレイ 2 0 0への書き込みの、現在の累積量とが、互いに対応付けて記憶される。
- [0132] 図 1 8の下欄の2つのグラフのうちの左側のグラフにおける、それぞれの丸印は、アレイ管理装置 1 0の記憶領域において対応付けて記憶された、日時、および、その日時における、アレイ 2 0 0に対する書き込みの累積量を示す。
- [0133] S 1 6 0 4では、それぞれの丸印での日時および累積量に基づいて、先述の関数  $f$ （図 1 8の左のグラフを参照）が特定される。なお、この特定の処理としては、従来技術における、この種の処理における技術（回帰分析の処理など）が用いられてもよい。

- [0134] 図21は、ストレージに対する累積書き換え量の増加度合い算出の処理のフローチャートである。先述された図19のS1503では、この図21の処理がされてもよい。
- [0135] S1701では、アレイ200に対する書き込みの、現在における累積量が取得される。S1702では、選択ストレージ（先述）に対する書き込みの、現在の累積量が取得される。なお、例えば、S1702で取得されるこの累積量が、選択ストレージに設けられた所定の記憶部により記憶されて、この記憶部により記憶された累積量が、アレイ管理装置10へと取得されてもよい。
- [0136] S1703では、S1701で取得された累積量と、S1702で取得された累積量とが、それぞれ、アレイ管理装置10の記憶領域に記憶される。
- [0137] S1704では、選択ストレージの関数 $g_j$ （図18の右のグラフ参照）が特定される。なお、この特定の処理では、例えば、S1604での、関数 $f$ の特定の処理と同様に、回帰分析の技術などが用いられてもよい。
- [0138] 図22は、3つの故障予測関数のデータのグラフを示す図である。
- [0139] 図示されるように、ストレージ $j$ の故障予測関数は、「関数 $f \times$ 関数 $g_j$ 」（関数 $f$ と関数 $g_j$ の合成関数）である。そして、互いに異なるストレージ201での故障予測関数は、通常は、互いに異なる関数である。また、図示されるように、例えば、互いに異なるストレージ201の保証量（先述）は、互いに異なってもよい。
- [0140] 図23は、同時故障発生のイメージ図である。
- [0141] 図23では、同時故障が発生する第1のケースを示す上段の欄と、第2のケースを示す中断の欄と、第3のケースを示す下段の欄とが示される。
- [0142] （変更処理の具体例）
- 図24は、同時故障回避の処理のフローチャートである。図9のS1404では、例えば、この図24の処理がされてもよい。
- [0143] S1901では、それぞれのストレージ201における、適切（最適）なパリティ割合201 $r$ が算出される。

- [0144] S 1 9 0 2では、情報1 0 1 1が作成される。作成される情報1 0 1 1は、その情報での、図5の表における、それぞれのストレージ2 0 1でのパリティ割合が、算出された、そのストレージ2 0 1での適切（最適）なパリティ割合2 0 1 rである情報である。
- [0145] S 1 9 0 3では、それぞれのストライプについて、適宜、スワップの処理がされる。これにより、それぞれのストレージ2 0 1のそれぞれのブロック1 × B 2について、そのブロック1 × B 2が、データのブロック1 × B 2 dであるか、パリティのブロック1 × B 2 pであるかが、作成された情報1 0 1 1（図5参照）に沿ったものにされて、それぞれのストレージ2 0 1のデータ構造が、作成された情報1 0 1 1でのデータ構造にされる。
- [0146] S 1 9 0 4では、作成された情報1 0 1 1が、この図2 4の処理がされる以前に、パリティ割合管理部1 0 1により記憶される情報1 0 1 1（図8参照）に対して上書きされて、この上書きがされて以後は、パリティ割合管理部1 0 1により、作成された新しい情報1 0 1 1が記憶される。
- [0147] 図2 5は、最適パリティ割合算出の処理のフローチャートである。図2 4のS 1 9 0 1では、例えば、この図2 5の処理がされてもよい。
- [0148] S 2 0 0 1では、関数f（先述の図1 8の左のグラフ参照）が特定される。S 2 0 0 2では、それぞれのストレージ2 0 1の許容残書き換え量が特定される。
- [0149] 図2 6は、各ストレージの許容残書き換え量算出の処理のフローチャートである。図2 5のS 2 0 0 2では、例えば、この図2 6の処理がされてもよい。
- [0150] S 2 1 0 1では、最初のストレージが、選択ストレージとして選択される。S 2 1 0 2では、選択ストレージの保証量が取得される。S 2 1 0 3では、選択ストレージに対する書き込みの、現在における累積量が取得される。S 2 1 0 4では、取得された保証量から、現在の累積量が減じられた量が、許容残書き換え量として特定される。
- [0151] S 2 1 0 5では、全ストレージ2 0 1の処理が終了したか否かが判定され

る。終了したと判定される場合には（S 2 1 0 5 : Y E S）、図 2 6 の処理が終了する。一方で、S 2 1 0 6 では、終了していないと判定される場合に（S 2 1 0 6 : N O）、選択ストレージとして、次のストレージ 2 0 1 が選択される。

[0152] S 2 0 0 3（図 2 6）では、関数が特定される。特定される関数は、ストレージ 2 0 1 における、パリティ割合 2 0 1  $r$  と、そのパリティ割合 2 0 1  $r$  での、そのストレージ 2 0 1 への書き込み量との間の関係を特定する関数である（図 3 3 参照）。

[0153] 図 2 7 は、ストレージのパリティ割合とブロック書き換え量との処理のフローチャートである。図 2 5 の S 2 0 0 3 では、例えば、この図 2 7 の処理がされてもよい。

[0154] 図 2 8 は、データブロック同時書き換えの説明図である。

[0155] 図 2 8 の表では、アレイ 2 0 0 への書き込みの際における第 1 のパターンが上欄で示され、第 2 のパターンが中欄で示され、第 3 のパターンが下欄で示される。第 1 のパターンでは、1 個のストレージ 2 0 1（ストレージ 3）に、パリティの書き込みがされると共に、1 個のストレージ 2 0 1（ストレージ 0）に、データの書き込みがされる。第 2 のパターンでは、1 個のストレージ 2 0 1（ストレージ 3）に、パリティの書き込みがされると共に、2 個のストレージ 2 0 1（ストレージ 0、1）に、データの書き込みがされる。第 3 のパターンでは、1 個のストレージ 2 0 1（ストレージ 3）に、パリティの書き込みがされると共に、3 個のストレージ 2 0 1（ストレージ 0～2）に、データの書き込みがされる。

[0156] 図 2 9 は、同時書き換えデータブロック数と、保存データ量との関係図である。

[0157] それぞれのパターンでの、アレイ 2 0 0 への書き込みでの、書き込みがされるデータ量は、例えば、図 2 9 で示される通りである。S 2 2 0 1（図 2 7）では、頻度の情報が取得される。取得される、頻度の情報は、アレイ 2 0 0 への書き込みとして、第 1 のパターンの書き込みが生じる第 1 の頻度（

確率)と、第2のパターンの書き込みが生じる第2の頻度(確率)と、第3のパターンの書き込みが生じる第3の頻度(確率)とをそれぞれ特定する。

[0158] 図30は、データブロック同時書き換え発生頻度情報管理テーブルの例である。

[0159] 図30では、取得される情報により、第1の頻度として25%が特定され、第2の頻度として40%が特定され、第3の頻度として35%が特定されるケースが示される。

[0160] 図31は、データブロック同時書き換え発生頻度に応じたデータの割り当ての図である。

[0161] アレイ200に対して書き込みがされるのに際して、第1、第2、第3のパターンでの書き込みが、上述の第1、第2、第3の頻度(25%、40%、35%)で行われる。そして、それぞれのパターンでの書き込みでは、1MB、2MB、3MBの書き込みがされる。

[0162] このため、第1、第2、第3のパターンで書き込みがされるデータ量の比は、「 $25\% \times 1\text{MB} : 40\% \times 2\text{MB} : 35\% \times 3\text{MB}$ 」=「 $25\text{MB} : 80\text{MB} : 105\text{MB}$ 」である。

[0163] このため、アレイ200に対して、あるデータ量(2520MB)の書き込みがされる際における、第1、第2、第3のパターンでの書き込みのデータ量は、 $2520\text{MB} \times (25\text{MB} / (25\text{MB} + 80\text{MB} + 105\text{MB})) = 2520\text{MB} \times (25\text{MB} / 210\text{MB}) = 300\text{MB}$ 、 $2520\text{MB} \times (80\text{MB} / 210\text{MB}) = 960\text{MB}$ 、 $2520\text{MB} \times (105\text{MB} / 210\text{MB}) = 1260\text{MB}$ である(図31の最下部を参照)。

[0164] つまり、第1、第2、第3のパターンでの書き込みがされるストライプの数は、 $300\text{MB} / 1\text{MB} = 300$ 個、 $960\text{MB} / 2\text{MB} = 480$ 個、 $1260\text{MB} / 3\text{MB} = 420$ 個である。

[0165] こうして、アレイ200に対する書き込みのデータ量(250MB)から、第1、第2、第3のパターンでの書き込みがされるストライプ数(300個、480個、420個)が特定される。

- [0166] S 2 2 0 2 (図 2 7) では、アレイ 2 0 0 に対して、単位データ量 (図 3 1 の 2 5 2 0 M B 参照) の書き込みがされる際における、第 1、第 2、第 3 のパターンでの書き込みがされる第 1、第 2、第 3 のストライプ数 (3 0 0 個、4 8 0 個、4 2 0 個を参照) が特定される。
- [0167] そして、第 1 のパターンで、第 1 のストライプ数 (3 0 0 個) の書き込みがされるのに際して、ストレージ 2 0 1 に対して行われる書き込みの回数は、以下の通りである。
- [0168] なお、以下では、説明の便宜上、ストレージ 2 0 1 のパリティ割合 2 0 1  $r$  が、 $x\%$  と略記される。
- [0169] つまり、第 1 のストライプ数 (3 0 0 個) のうちの、 $(100 - x)\%$  の割合の個数 ( $300 \text{ 個} \times (100 - x)\%$ ) のストライプについては、それらのストライプのうちのそれぞれのストライプにおける、そのストレージ 2 0 1 のブロック  $1 \times B 2$  は、データのブロックである。
- [0170] そして、第 1 のパターンでの書き込みでは、図 3 0 の上欄に示されるように、ストライプに含まれる、3 つの、データのブロックのうちの一つのブロックにのみ、書き込みがされる。
- [0171] このため、それらの、上述された、 $(300 \text{ 個} \times (100 - x)\%)$  の個数のストライプのうち、 $1/3$  の個数  $\{(300 \text{ 個} \times (100 - x)\%) / 3\}$  のストライプの書き込みにおいて、そのストレージ 2 0 1 への書き込みがされる。
- [0172] 図 3 2 は、データブロックにおける書き換え処理発生率の説明図である。
- [0173] なお、上述の  $1/3$  の点については、適宜、この図 3 2 も参照されたい。
- [0174] 一方で、上述された、第 1 のストライプ数 (3 0 0 個) のストライプのうち、 $x\%$  の割合の個数 ( $300 \text{ 個} \times x\%$ ) のストライプについては、そのストライプにおける、そのストレージ 2 0 1 のブロック  $1 \times B 2$  が、パリティのブロック  $1 \times B 2 p$  である。
- [0175] このため、これら、 $(300 \text{ 個} \times x\%)$  の個数のストライプの書き込みでは、必ず、そのストレージ 2 0 1 への書き込みがされる。

- [0176] このため、第1のストライプ数(300個)のストライプのうちの、「{(300個×(100-x)%} / 3} + (300個×x%)」の個数のストライプの書き込みにおいて、そのストレージ201への書き込みがされる。
- [0177] こうして、アレイ200に対して、単位データ量(2520MB参照)の書き込みがされる際に生じる、第1のパターンでの、ストレージ201への書き込みの回数 $W1(x)$ が、上述された、「{(300個×(100-x)%} / 3} + (300個×x%)」と特定される。
- [0178] 同様にして、第1のパターンでの、ストレージ201への書き込みの回数 $W2(x)$ 、および、第3のパターンでの、ストレージ201への書き込みの回数 $W3(x)$ も特定される。
- [0179] これにより、2520MBの書き込みがされる際における、ストレージ201への書き込みの回数 $J(x)$ が、 $J(x) = W1(x) + W2(x) + W3(x)$ と特定される。
- [0180] 図33は、パリティ割合201 $r$ (上述の $x$ )と、回数 $J(x)$ だけの書き込みで書き込まれるデータ量との関係を特定する関数 $K(x)$ のグラフを示す図である。先述の $J(x)$ が特定されることにより、この $K(x)$ が特定される。
- [0181] 図34は、同時書き換えデータブロック数を考慮したブロック書き換え量を示す図である。
- [0182] すなわち、上述された、第1のパターンでの回数である $W1(x)$ は、図34の上欄の左、中央、右の3つのグラフのうちの左のグラフでの関数である。そして、 $W2(x)$ は、中央のグラフでの関数であり、 $W3(x)$ は、右のグラフでの関数である。
- [0183] そして、下欄のグラフの関数は、上述された、 $J(x) = W1(x) + W2(x) + W3(x)$ から特定される $K(x)$ である。S2204~S2206(図27)では、上述の $W1(x)$ 、 $W2(x)$ 、 $W3(x)$ が特定される。S2207では、特定された、これら $W1(x)$ 、 $W2(x)$ 、 $W3$

(x) から、上述の  $K(x)$ 、つまり、図 34 の下欄のグラフの関数が特定される。

[0184] 図 35 A は、複数のストレージ 201 (ストレージ 0~3) でのパリティ割合 201 r の組み合わせを示す図である。

[0185] 図 35 B は、図 35 A の組み合わせの状態でのアレイ 200 を示す図である。

[0186] 図 35 C は、パリティ割合 201 r の組み合わせが、図 35 A の組み合わせである際における、それぞれのストレージ 201 での、書き込みがされるデータ量を示す図である。

[0187] 図 35 B の組み合わせでは、図 35 C で示されるように、ストレージ 0~4 での、書き込みのデータ量 (上述の  $K(x)$ ) の比が、「1 : 1 : 1 : 3」である。

[0188] 図 36 は、図 35 A での組み合わせとは異なる組み合わせを示す図である。

[0189] アレイ管理装置 10 により特定される適切な組み合わせは、通常は、図 35 A で示されるような組み合わせでなく、図 36 で示されるような組み合わせである。

[0190] 図 37 は、3つのグラフを示す図である。

[0191] 上段のグラフは、パリティ割合 201 r と、書き込みのデータ量との間の関係を示すグラフであり、上述の  $K(x)$  のグラフである (図 34 の下欄を参照)。中段のグラフは、パリティ割合 201 r と、アレイ 200 に対する書き込みのデータ量との間の関係を示すグラフ ( $L(x)$ ) である。上段のグラフ ( $J(x)$ ) から、この、中段のグラフ ( $L(x)$ ) が特定される。

[0192] 下段のグラフは、パリティ割合 201 r と、ストレージ 201 の故障予測日時との関係を示すグラフ ( $M(x)$ ) である。中段のグラフ ( $L(x)$ ) から、この、下段のグラフ ( $M(x)$ ) が特定される。

[0193] S2004 (図 25) では、特定された上述の  $J(x)$  から特定される、この  $M(x)$  に基づいて故障予測日時が特定される。つまり、例えば、それぞれのストレージ 201 での、そのストレージ 201 が、それぞれのパリティ

ィ割合  $201r$  のときにおける故障予測日時が特定される。

[0194] 図38は、4つのグラフを示す図である。すなわち、S2004では、図38の4つのグラフが特定される。

[0195] 図39は、各パリティ割合における各ストレージの故障予測日時算出の処理のフローチャートである。図25のS2004では、例えば、この図39の処理がされてもよい。

[0196] 図26の処理では、それぞれのストレージ201を対象ストレージとして、S3302の処理が行われる。

[0197] S3302（図39）では、それぞれのストレージ201（対象ストレージ）についての、それぞれのパリティ割合  $201r$  での故障予測日時が特定される。

[0198] S2005（図25）では、アレイ200に含まれる複数のストレージ201のパリティ割合  $201r$  の組合わせのそれぞれについて、その組合わせでの、それぞれのストレージ201の故障予測日時が特定される。つまり、S2005では、それぞれの組合わせについて、その組合わせにおける、複数の故障予測日時が特定される。

[0199] そして、S2005では、複数の組合わせ（上述）のうちから、特定される、複数の故障予測日時が、最も適切である組合わせを特定する（図38の処理1G、先述の図10Bなどを参照）。つまり、複数の組合わせから、最も適切な組合わせが特定される。なお、最も適切である組合わせでの、特定される複数の故障予測日時においては、例えば、先述された第1のストレージ  $\alpha$  の故障予測日時（図10Bの日時  $1t$ ）と、第2のストレージ  $\beta$  の故障予測日時  $1B$  との間の差  $1F1$  が、最大でもよい。

[0200] そして、図24のS1903～S1904では、特定された、最適である、パリティ割合の組合わせに基づいて、先述された処理が行われる。つまり、例えば、このS1903～S1904では、アレイ200の動作が制御されてもよい。つまり、制御がされた動作では、アレイ200における、それぞれのストレージ201でのパリティ割合が、特定された、最適な組合わせ

での、そのストレージ201のパリティ割合と同じ割合でもよい（先述された、S1903、S904の説明などを参照）。

[0201] （数式等による説明）

以下では、より詳細な説明がされる。

[0202] まず、記号の説明がされる。以下の説明で、記号「 $N_{\text{stripe}}$ 」は、アレイ200におけるストライプ $1 \times B1$ の数（ストライプ数）を示す。つまり、ストライプ番号 $i$ は、 $0 \sim N_{\text{stripe}} - 1$ の値と取りうる。また、記号「 $N_{\text{storage}}$ 」は、アレイ200におけるストレージ201の数（ストレージ数）を示す。ただし、スペアストレージは含まない。つまり、ストレージ番号 $j$ は、 $0 \sim N_{\text{storage}} - 1$ の値を取りうる。また、記号「 $N_{\text{data}}$ 」は、各ストライプにおける、データブロック（ブロック $1 \times B2d$ ）のストレージ数（図4等では3個）を示す。また、記号「 $N_{\text{parity}}$ 」は、各ストライプにおける、パリティブロック（ブロック $1 \times B2p$ ）のストレージ数を示す。つまり、RAID5では、 $N_{\text{parity}} = 1$ であり、また、RAID6では $N_{\text{parity}} = 2$ である。ここで、 $N_{\text{stripe}}$ 、 $N_{\text{storage}}$ 、 $N_{\text{data}}$ 、 $N_{\text{parity}}$ は、一般に、アレイの新規構築時に決定し、運用中は変化しない。例として、 $N_{\text{storage}} = 4$ 、 $N_{\text{data}} = 3$ 、 $N_{\text{parity}} = 1$ とする。

[0203] 図3などにおいて、それぞれのブロック $1 \times B2$ は、 $Bk \circ \Delta$ という形で表現される。ここで、 $\circ$ はストライプ番号であり、 $\Delta$ はストレージ番号である。

[0204] 図4、図7、図5などにおいて、データに対応するブロック（データのブロック $1 \times B2d$ ）は、 $BD \circ \#$ という形で表現される。なお、ここで、「 $\#$ 」は、「各ストライプにおけるデータブロックに対応するストレージ番号」である。一方、パリティに対応するブロック（ブロック $1 \times B2p$ ）は、 $BP \circ \$$ という形で表現される。ここで、「 $\$$ 」は、「各ストライプにおけるパリティブロックに対応するストレージ番号」である。

[0205] 次に、図18の2つのグラフのうちの、左のグラフについて述べる。

[0206] まず、記号の説明がされる。記号「 $T_{\text{now}}$ 」は、現在時刻を示す。記号「 $D_{\text{now}}$ 」は、アレイに対する、現在の累積データ書き込み量を示す。記号「 $T_j$ 」は

、アレイにおいて、ストレージ  $j$  が破損すると予測される日時を示す。記号「 $D_j$ 」は、日時  $T_j$  における、アレイに対する累積データ書き込み量を示す。記号「 $t_j$ 」は、現在時刻からストレージ  $j$  の故障予測日時までの所要時間（ $T_j = T_{\text{now}} + t_j$ ）を示す。記号「 $d_j$ 」は、日時  $T_j$  における、アレイに対する累積データ書き込み量（ $D_j = D_{\text{now}} + d_j$ ）を示す。

[0207] ここで、 $T_{\text{now}}$  は取得可能である。また、 $D_{\text{now}}$  は、予め、過去のデータ書き込み量を記録しておく仕組みがアレイ管理装置 10 に設けられ、それを利用し、取得可能とする。そして、以下の説明では、 $t_j$ 、 $d_j$  を求めることが述べられる。これら  $t_j$ 、 $d_j$  より、 $T_j$ 、 $D_j$  を決定することが可能である。

[0208] 故障予測関数  $f$  が、線形とすると、以下の式での表現がされる。ここで、 $a$  は係数である。

[0209] [数1]

$$a \cdot t_j = d_j \quad \dots (1)$$

[0210] 係数  $a$  については、「過去のデータ書き込み量」と、「その時点の日時」の記録を用いた回帰分析等により決定される。

[0211] 次に、図 18 の右のグラフについて説明される。

[0212] まず、記号の説明がされる。記号「 $W_{j\text{now}}$ 」は、ストレージ  $j$  に対する現在の累積ブロック書き換え量を示す。記号「 $W_{j\text{max}}$ 」は、ストレージ  $j$  に対するブロック書き換え量の保証値を示す。記号「 $w_j$ 」は、ストレージ  $j$  に対する許容残書き換え量（ $w_j = W_{j\text{max}} - W_{j\text{now}}$ ）を示す。ここで、 $W_{j\text{now}}$  は、予め、過去のブロック書き換え量を記録しておく仕組みをアレイ管理装置 10 に設けておき、それを利用し、取得可能とする。 $W_{j\text{max}}$  は、メーカー保証値、あるいは、加速度試験などにより、予め決定する（S.M.A.R.T. (Self-Monitoring, Analysis and Reporting Technology) を用いてもよい）。 $W_{j\text{now}}$ 、 $W_{j\text{max}}$  が決定可能なので、 $w_j$  も決定可能である。

[0213] 故障予測関数  $g_j$  が線形とすると、以下の式で表現される（ここで、 $b_j$  は係

数)。

[0214] [数2]

$$b_j \cdot d_j = w_j$$

・・・ (2)

[0215] 係数  $b_j$  は、ストレージごとに異なる。「過去のブロック書き換え量」と、「その時点のアレイに対するデータ書き込み量」の記録を用いた回帰分析等により決定される。

[0216]  $W_{j\text{now}}$ 、 $W_{j\text{max}}$ 、 $a$ 、 $b_j$  が決定されることで、上述の式(1)(2)を用いて、 $t_j$  を算出できる。

[0217] 例として、以下の場合を考える(データ量の単位をMB、時間の単位をhとする)。アレイに対する書き込みデータ量は、平均して1時間あたり、400MB ( $a=400$ ) である。アレイに対して、10MBのデータ書き込みが発生する際に、ストレージ  $j$  に対して、4MBのブロック書き換えが発生するとする ( $b_j=0.4$ )。ストレージ  $j$  における許容残書き換え量は40000MB ( $w_j=40000$ ) である。式(1)、(2)より、 $d_j=100000$ 、 $t_j=250$  であり、つまりストレージ  $j$  は、250時間後に破損するという予測が算出できる。全ストレージに対して、同様に、 $t_j$  を算出する。なお、故障予測関数  $f$ 、 $g_j$  は非線形でも良い。

[0218] 次に、図10A～図10Cについて述べられる。

[0219] まず、記号の説明がされる。記号「M」は、冗長化アレイにおける最大許容故障台数を示す。つまり、RAID5では $M=1$ であり、RAID6では $M=2$ である。RAID5、6の場合には、Mは、 $N_{\text{parity}}$ と一致する。記号「 $T_\alpha$ 」は、アレイにおいて、最初にストレージが破損すると予測される日時を示す。ここで、最初に破損すると予測されているストレージを $\alpha$ とする。記号「 $T_\beta$ 」は、アレイにおいて、 $M+1$ 番目にストレージが破損する(つまり、アレイが破損する)と予測される日時を示す。ここで、 $M+1$ 番目に破損すると予測されているストレージを $\beta$ とする。記号「 $\Delta_{\alpha\beta}$ 」は、ストレージ

$\alpha$ と、ストレージ $\beta$ の破損の間隔 ( $\Delta_{\alpha\beta} = T_{\beta} - T_{\alpha}$ ) を示す。

[0220] ここで、 $M$ は、一般に、アレイの新規構築時に決定し、運用中は変化しない。 $T_{\alpha}$ 、 $T_{\beta}$ は、全ストレージについて $T_j$ を算出後、並び替えることで、決定可能である。そして、 $T_{\alpha}$ 、 $T_{\beta}$ が決定可能なので、 $\Delta_{\alpha\beta}$ も決定可能である。そして、 $\Delta_{\alpha\beta}$ が、所定値以内である場合、同時故障発生の危険性があると判断する。

[0221] 続けて、図10Cについて述べられる。

[0222] まず、記号の説明がされる。記号「 $T_{\gamma}$ 」は、アレイにおいて、 $x$ 番目にストレージが破損すると予測される日時を示す。ここで、 $x$ 番目に破損すると予測されているストレージを $\gamma$ とする。 $x$ は、 $2 \sim M$ の値を取りうる。記号「 $T_{\gamma 1}$ 」は、アレイにおいて、 $x - 1$ 番目にストレージが破損すると予測される日時を示す。ここで、 $x - 1$ 番目に破損すると予測されているストレージを $\gamma 1$ とする。記号「 $T_{\gamma 2}$ 」は、アレイにおいて、 $x + 1$ 番目にストレージが破損すると予測される日時を示す。ここで、 $x - 2$ 番目に破損すると予測されているストレージを $\gamma 2$ とする。記号「 $\Delta_{\gamma\gamma 1}$ 」は、ストレージ $\gamma$ と、ストレージ $\gamma 1$ の破損の間隔の時間 ( $\Delta_{\gamma\gamma 1} = T_{\gamma} - T_{\gamma 1}$ ) を示す。記号「 $\Delta_{\gamma\gamma 2}$ 」は、ストレージ $\gamma$ と、ストレージ $\gamma 2$ の破損の間隔の時間 ( $\Delta_{\gamma\gamma 2} = T_{\gamma} - T_{\gamma 2}$ ) を示す。

[0223] ここで、 $x$ の値は $2 \sim M$ の整数値を取りうる。 $x$ が決定すれば、 $T_{\gamma}$ 、 $T_{\gamma 1}$ 、 $T_{\gamma 2}$ 、 $\Delta_{\gamma\gamma 1}$ 、 $\Delta_{\gamma\gamma 2}$ も決定可能である。任意の $x$ について、それぞれ $\Delta_{\gamma\gamma 1}$ 、 $\Delta_{\gamma\gamma 2}$ を算出する。そして、 $\Delta_{\gamma\gamma 1}$ または $\Delta_{\gamma\gamma 2}$ が、所定値以内である場合、同時故障発生の危険性があると判断する。

[0224] 次に、図28、図30について述べられる。

[0225] 記号の説明がされる。記号「 $v$ 」は、同時書き換えデータブロック数を示す。 $v$ は、 $1 \sim N_{\text{data}}$ の値を取りうる。

[0226] 以下の説明では、まず、 $v$ が、任意の値を取る場合について説明する。その後、 $v$ の発生頻度を用いて、全ての $v$ について算出した結果を統合して考える。

[0227] 次に、図29について述べられる。

[0228] まず、記号の説明がされる。記号「BK」は、ブロックサイズ（ここでは固定とするが、可変でもよい）を示す。記号「C<sub>v</sub>」は、同時書き換えデータブロック数vにおいて、1回のストライプ書き込みで保存可能なデータ量を示す。

[0229] C<sub>v</sub>は、以下の式により算出される。

[0230] [数3]

$$C_v = v \cdot BK \quad \dots (3)$$

[0231] 次に、図31について述べられる。

[0232] まず、記号の説明がされる。記号「R<sub>v</sub>」は、同時書き換えデータブロック数vにおける、書き込みストライプ数を示す。記号「P<sub>v</sub>」は、同時書き換えデータブロック数vにおける、同時書き換え発生頻度。（P<sub>v</sub>は0～1の値を取りうる。また、全同時書き換えデータブロック数における合計値はRAID5なら1、RAID6なら2となる）を示す。

[0233] ここで、媒介変数cを用いて、以下の式での表現ができる。

[0234] [数4]

$$R_v = c \cdot P_v \quad \dots (4)$$

[0235] 受付データのデータ量をdとすると、このdは、以下の式で算出される。

[0236] [数5]

$$d = \sum_{v=1}^{N_{data}} (R_v \cdot C_v) \quad \dots (5)$$

つまり、(3)、(4)より、下記の通りである。

[0237]

[数6]

$$d = c \cdot BK \cdot \sum_{v=1}^{N_{data}} (v \cdot P_v)$$

・・・ (6)

例として、以下の場合を考える（データ量の単位をMBとする）。つまり、 $N_{data} = 3$ 、 $P_1 = 0.25$ 、 $P_2 = 0.4$ 、 $P_3 = 0.35$ 、 $BK = 1$ 、 $d = 2520$ の場合が考えられる。

[0238] 式(6)を解くと、 $c = 1200$ となる。つまり、 $R_1 = 300$ 、 $R_2 = 480$ 、 $R_3 = 420$ と算出することができる。

[0239] 次に、図32について述べられる。まず、記号の説明がされる。記号「 $Q_v$ 」は、同時書き換えデータブロック数 $v$ における、データブロックにおける書き換え処理発生率を示す。

[0240]  $Q_v$ は、以下の式により算出可能である。

[0241] [数7]

$$Q_v = \frac{v}{N_{data}}$$

・・・ (7)

例として、 $N_{data} = 3$ 、 $v = 1$ なら、 $Q_1 = 1/3$ である。

[0242] 次に、図35A～図35Cについて述べられる。

[0243] 同時書き換えデータブロック数 $v$ において、 $R_v$ 個のストライプに対して書き込みを行うことを考える。ここで、各ストレージに対してブロック書き換えが発生する。

[0244] まず、各ストレージにおけるパリティ割合が0%または100%である場合を考える。

[0245] 最初に、記号の説明がされる。記号「 $w_{vdata}$ 」は、同時書き換えデータブロック数 $v$ において、 $R_v$ 個のストライプに対して書き込みを行う際に、「パリティ割合0%」のストレージに対して発生するブロック書き換え量を示す。記号「 $w_{vparity}$ 」は、同時書き換えデータブロック数 $v$ において、 $R_v$ 個のスト

ライブに対して書き込みを行う際に、「パリティ割合100%」のストレージに対して発生するブロック書き換え量を示す。

[0246]  $w_{vdata}$ 、 $w_{vparity}$ は、以下の式により算出される。

[0247] [数8]

$$W_{vdata} = Q_v \cdot R_v \cdot BK \quad \dots (8)$$

[0248] [数9]

$$W_{vparity} = R_v \cdot BK \quad \dots (9)$$

[0249] 例として、以下の場合を考える（データ量の単位をMB）。つまり、 $N_{data} = 3$ 、 $v = 1$ 、 $BK = 1$ 、 $R_v = 300$ である場合が考えられる。

[0250] 式(8)、(9)より、 $w_{vdata} = 100$ 、 $w_{vparity} = 300$ 、つまりアレイに対して300MBのデータ書き込みを行う際に、パリティ割合0%のストレージに対するブロック書き換え量は100MBであるといえる。

[0251] 次に、図6、33、36について述べられる。

[0252] パリティ割合が0%および100%以外の場合について考える。

[0253] まず、記号の説明がされる。記号「 $w_{vr}$ 」は、同時書き換えデータブロック数 $v$ において、 $R_v$ 個のストライプに対して書き込みを行う際に、「パリティ割合 $r$ 」のストレージに対して発生するブロック書き換え量。以下の式により算出される（ただし $r$ は0~1の範囲を示す）。

[0254]  $w_{vr}$ は、以下の式により算出される。

[0255] [数10]

$$w_{vr} = (1 - r) \cdot w_{vdata} + r \cdot w_{vparity} \quad \dots (10)$$

[0256] 例として、 $w_{vdata} = 100$ 、 $w_{vparity} = 300$ 、 $R = 0.3$ とすると、 $w_{vr} =$

160と算出できる。

[0257] 次に、図34について述べられる。ここでは、前記で各 $v$ について算出した結果を統合して考える。まず、記号の説明がされる。記号「 $w_r$ 」は、「パリティ割合 $r$ 」のストレージに対して発生するブロック書き換え量を示す。

[0258] ここで、 $w_r$ は、各同時書き換えデータブロック数における合計値となる。つまり、 $w_r$ は、以下の式により算出される。

[0259] [数11]

$$w_r = \sum_{v=1}^{N_{data}} (w_{vr})$$

... (11)

(10)より、

[0260] [数12]

$$w_r = \sum_{v=1}^{N_{data}} ((1-r) \cdot w_{vdata} + r \cdot w_{vparity}) = \sum_{v=1}^{N_{data}} (w_{vdata} + (w_{vparity} - w_{vdata}) \cdot r)$$

... (12)

である。

[0261] (8)、(9)より、

[0262] [数13]

$$w_r = BK \cdot \sum_{v=1}^{N_{data}} (Q_v \cdot R_v + (1 - Q_v) \cdot R_v \cdot r)$$

... (13)

である。

[0263] (7)より、

[0264] [数14]

$$w_r = BK \cdot \sum_{v=1}^{N_{data}} \left( \frac{v}{N_{data}} \cdot R_v + \left(1 - \frac{v}{N_{data}}\right) \cdot R_v \cdot r \right)$$

... (14)

と表せる。

[0265] 例として、以下の場合を考える（データ量の単位をMB）。つまり、 $N_{data} = 3$ 、 $R_1 = 300$ 、 $R_2 = 480$ 、 $R_3 = 420$ 、 $BK = 1$ である場合が考えられる。

[0266] 式(14)への代入がされると、以下の式となる。

[0267] [数15]

$$w_r = 360r + 840 \quad \dots (15)$$

[0268] 上記の式に、例えば  $r = 0.3$  を代入すると、「アレイに対して2520 MBの書き込みを行う際の、パリティ割合30%のストレージに対するブロック書き換え量」が、 $w_r = 948$ として算出できる。

[0269] さらに、(14)に対して(4)を適用すると

[0270] [数16]

$$w_r = c \cdot BK \cdot \sum_{v=1}^{N_{data}} \left( \frac{v}{N_{data}} \cdot P_v + \left(1 - \frac{v}{N_{data}}\right) \cdot P_v \cdot r \right) \quad \dots (16)$$

である。(6)より、

[0271] [数17]

$$c \cdot BK = \frac{d}{\sum_{v=1}^{N_{data}} (v \cdot P_v)} \quad \dots (17)$$

なので、

[0272] [数18]

$$w_r = \frac{d}{\sum_{v=1}^{N_{data}} (v \cdot P_v)} \cdot \sum_{v=1}^{N_{data}} \left( \frac{v}{N_{data}} \cdot P_v + \left(1 - \frac{v}{N_{data}}\right) \cdot P_v \cdot r \right)$$

．．． (18)

である。整理すると

[0273] [数19]

$$w_r = d \cdot \left( \frac{1}{N_{data}} + \left( \frac{\sum_{v=1}^{N_{data}} (P_v)}{\sum_{v=1}^{N_{data}} (v \cdot P_v)} - \frac{1}{N_{data}} \right) \cdot r \right)$$

．．． (19)

である。ここで、

[0274] [数20]

$$\sum_{v=1}^{N_{data}} P_v = 1$$

．．． (20)

なので、

[0275] [数21]

$$w_r = d \cdot \left( \frac{1}{N_{data}} + \left( \frac{1}{\sum_{v=1}^{N_{data}} (v \cdot P_v)} - \frac{1}{N_{data}} \right) \cdot r \right)$$

．．． (21)

である。

[0276] つまり、式(21)より、 $r$ 、 $d$ 、 $N_{data}$ 、およびそれぞれの同時書き換えデータブロック数 $v$ における、同時書き換え発生頻度 $P_v$ が決定すれば、 $w_r$ を決定することができる。

[0277] ここで、

[数22]

$$\sum_{v=1}^{N_{data}} (v \cdot P_v)$$

．．． (22)

を $Y$ とおく。 $Y$ の値は、同時書き換え発生頻度 $P_v$ と $N_{data}$ により決定する。 $N_{data}$ は、アレイ新規構成時に決定し、 $P_v$ は、ストレージに対して書き込みを行うアプリケーションやシステムなど外部の特性により決定する（シーケンシャルライト、ランダムライトの割合、その時のデータ量など）。つまり、 $Y$ は、アレイを構成する各ストレージのパリティ割合の影響を受けない。

[0278] 次に、図37について述べられる。

[0279] ストレージ $j$ のパリティ割合を $r$ に再設定することを考える。まず、記号の説明がされる。記号「 $T_{jr}$ 」は、ストレージ $j$ のパリティ割合を $r$ に再設定した時に、アレイにおいて、ストレージ $j$ が破損すると予測される日時を示す。記号「 $D_{jr}$ 」は、日時 $T_{jr}$ における、アレイに対する累積データ書き込み量を示す。記号「 $t_{jr}$ 」は、現在時刻から、ストレージ $j$ のパリティ割合を $r$ に再設定した時のストレージ $j$ の故障予測日時までの所要時間（ $T_{jr} = T_{now} + t_{jr}$ ）を示す。記号「 $d_{jr}$ 」は、日時 $T_{jr}$ における、アレイに対する累積データ書き込み量を示す（ $D_{jr} = D_{now} + d_{jr}$ ）。

[0280] ストレージ $j$ のパリティ割合を $r$ に変更することで、故障予測日時を変化させることができる。なお、故障予測関数 $f$ については影響を受けない。ここで、式(20)において $w_r = w_j$ 、 $d = d_{jr}$ とすることで、 $d_{jr}$ を算出できる。

[0281] [数23]

$$w_j = d_{jr} \cdot \left( \frac{1}{N_{data}} + \left( \frac{1}{Y} - \frac{1}{N_{data}} \right) \cdot r \right) \quad \dots (23)$$

変形すると

[0282] [数24]

$$d_{jr} = \frac{N_{data} \cdot Y}{Y + (N_{data} - Y) \cdot r} \cdot w_j \quad \dots (24)$$

である。さらに、故障予測関数 $f$ を用いて故障予測日時を算出する。(1)より、

[0283] [数25]

$$t_{jr} = \frac{N_{data} \cdot Y}{Y + (N_{data} - Y) \cdot r \cdot a} \cdot w_j \quad \dots (25)$$

である。

[0284] 例として、以下の場合を考える（データ量の単位をMB）。つまり、 $N_{data} = 3$ 、 $P_1 = 0.25$ 、 $P_2 = 0.4$ 、 $P_3 = 0.35$ 、 $w_j = 40000$ である場合が考えられる。

[0285] 式(24)への代入がされると、以下の式が成り立つ。

[0286] [数26]

$$d_{jr} = \frac{840000}{7 + 3 \cdot r} \quad \dots (26)$$

さらに、 $a = 400$ とすると、

[0287] [数27]

$$t_{jr} = \frac{2100}{7 + 3 \cdot r} \quad \dots (27)$$

という式が成り立つ。

[0288] 上記の式に、例えば  $r = 0.3$  を代入すると、「ストレージ j のパリティ割合を 30% とした際の、故障予測日時（現在日時からの経過時間）が  $t_{jr} = 266$ 、つまり 266 時間後として算出できる。そして、例えば、パリティ割合を 0% に変更すれば  $t_{jr} = 300$ 、つまり故障予測日時を 34 時間変化させることができる。

[0289] 最後に、図 38 について述べられる。図 38 では、全てのストレージに対して、故障予測日時までの所要時間  $t_{jr}$  と、パリティ割合  $r$  の関係関数を導出することが示される。

[0290] （その他）

なお、例えば、次のようでもよい。

[0291] すなわち、アレイ200により保持されるデータの消失は、第1のストレージ $\alpha$ が故障する日時1tが第1の日時1t1である場合には、発生し難く、第2の日時1t2である場合には発生し易く、第1のストレージ $\alpha$ は、第1のストレージ $\alpha$ により記憶されるデータ（領域201d参照）とパリティ（領域201p参照）との全体（領域201A参照）に占める、パリティの割合であるパリティ割合201rが、第1の割合201r1である場合と、第2の割合201r2である場合とがあるストレージであり、第1のストレージ $\alpha$ の故障の日時1tは、第1の割合201r1である場合には、消失が発生し難い第1の日時1t1であり、第2の割合201r2である場合には、発生し易い第2の日時1t2であり、判断部100xは、予め定められた2つの割合のうち一方の割合を第1の割合201r1と特定し、他方の割合を第2の割合201r2と特定し、変更制御部100yは、第1のストレージ $\alpha$ での割合201rを、特定された一方の割合（第1の割合201r1）にさせ、特定された他方の割合（第2の割合201r2）にはさせない制御をしてもよい。

[0292] そして、消失が発生し難い第1の日時1t1は、アレイ200に含まれる、第1のストレージ $\alpha$ とは別の第2のストレージ $\beta$ が故障する日時1Bから閾値1F2以上に離れた日時であり、消失が発生し易い第2の日時1t2は、第2のストレージ $\beta$ が故障する日時1Bから閾値1F2以内の日時でもよい。

[0293] これにより、データの消失の発生が防げる。

[0294] しかも、単なる、パリティの割合に応じた処理がされるだけで済み、処理が簡単にできる。

[0295] なお、本実施例のアレイ管理装置は、典型的には半導体集積回路であるLSIとして実現される。これらは個別に1チップ化されてもよいし、一部またはすべてを含むように1チップ化されても良い。ここではLSIとしたが、集積度の違いにより、IC、システムLSI、スーパーLSI、ウルトラ

L S I と呼称されることもある。

[0296] また、集積回路化の手法はL S Iに限るものではなく、専用回路または汎用プロセッサで実現しても良い。L S I製造後に、プログラムすることが可能なFPGA (Field Programmable Gate Array) や、L S I内部の回路セルの接続や設定を再構成可能なリコンフィギュラブル・プロセッサを利用しても良い。

[0297] さらに、半導体技術の進歩または派生する別技術によりL S Iに置き換わる集積回路化の技術が登場すれば、当然、その技術を用いて機能ブロックの集積化を行っても良い。バイオ技術の適応などが可能性として有り得る。

[0298] さらに加えて、本実施例の～装置を集積化した半導体チップと、画像を描画するためのディスプレイとを組み合わせ、様々な用途に応じた描画機器を構成することができる。携帯電話やテレビ、デジタルビデオレコーダ、デジタルビデオカメラ、カーナビゲーション等における情報描画手段として、本発明を利用することが可能である。ディスプレイとしては、ブラウン管 (CRT) の他、液晶やPDP (プラズマディスプレイパネル)、有機ELなどのフラットディスプレイ、プロジェクターを代表とする投射型ディスプレイなどと組み合わせることが可能である。

[0299] なお、上記各実施の形態において、各構成要素は、専用のハードウェアで構成されるか、各構成要素に適したソフトウェアプログラムを実行することによって実現されてもよい。各構成要素は、CPUまたはプロセッサなどのプログラム実行部が、ハードディスクまたは半導体メモリなどの記録媒体に記録されたソフトウェアプログラムを読み出して実行することによって実現されてもよい。ここで、上記各実施の形態のレイ管理装置などを実現するソフトウェアは、次のようなプログラムである。

[0300] すなわち、このプログラムは、複数のストレージが構成するレイを管理するコンピュータに、前記複数のストレージのうち少なくとも2台のストレージが同時期に故障することによって、前記レイが保持するデータが消失する可能性があるか否かを判断する判断ステップと、前記レイが保持する

前記データが消失する可能性があるとは判断された場合には、前記少なくとも2台のストレージのうちの少なくとも一つの故障時期を変更する、故障時期変更ステップとを実行させるためのコンピュータプログラムなどである。

### 産業上の利用可能性

[0301] 本発明のアレイ管理装置は、様々な用途に利用可能である。例えば、携帯電話や携帯音楽プレーヤー、デジタルカメラ、デジタルビデオカメラ等の電池駆動の携帯表示端末や、テレビ、デジタルビデオレコーダ、カーナビゲーション等の高解像度の情報表示機器におけるメニュー表示やWebブラウザ、エディタ、EPG、地図表示等における情報表示手段として利用価値が高い。

[0302] 将来、複数のストレージが同時期に故障することによって、アレイが保持するデータが消失する可能性があるとは判断した場合には、同時故障する可能性がある複数のストレージの故障時期をずらす。このことで、アレイが保持するデータ消失を防ぐことができる。

### 符号の説明

- [0303] 1 x B 2、1 x B 2 d、1 x B 2 p ブロック
- 1 0 アレイ管理装置
  - 1 0 0 x 判断部
  - 1 0 0 y 変更制御部
  - 1 0 1 パリティ割合管理部
  - 1 0 2 パリティ割合算出部
  - 1 0 3 パリティ割合変更部
  - 1 0 4 故障予測部
  - 2 0 0 アレイ
  - 2 0 1 ストレージ
  - 2 0 1 d 領域
  - 2 0 1 p 領域

## 請求の範囲

- [請求項1] 複数のストレージが構成するアレイを管理するアレイ管理装置であって、
- 前記複数のストレージのうち少なくとも2台のストレージが同時期に故障することによって、前記アレイが保持するデータが消失する可能性があるか否かを判断する判断部と、
- 前記アレイが保持する前記データが消失する可能性があると判断された場合には、前記少なくとも2台のストレージのうちの少なくとも一つの故障時期を変更する、故障時期変更部と、
- を備えるアレイ管理装置。
- [請求項2] 前記故障時期変更部は、前記各ストレージに対して書き込む、パリティとデータとの割合を変更することによって、前記少なくとも2台のストレージのうちの少なくとも一つの故障時期を変更する、
- 請求項1記載のアレイ管理装置。
- [請求項3] 前記判断部は、
- 前記複数のストレージのうち少なくとも2台のストレージが同時期に故障する可能性があるか否かを判断する同時故障判断部と、
- 前記少なくとも2台のストレージが同時期に故障すると判断された場合に、前記少なくとも2台のストレージが同時期に故障することによって、前記アレイが保持するデータが消失する可能性があるか否かを判断するアレイ故障判断部とを含む
- 請求項1記載のアレイ管理装置。
- [請求項4] 前記アレイにより保持される前記データの消失は、第1の前記ストレージが故障する日時が第1の日時である場合には、発生し難く、第2の日時である場合には発生し易く、
- 第1の前記ストレージは、記憶されるデータとパリティとの全体に占める、前記パリティの割合が、第1の割合である場合と、第2の割合である場合とがあるストレージであり、

第1の前記ストレージの前記日時は、前記第1の割合である場合には、消失が発生し難い前記第1の日時であり、前記第2の割合である場合には、発生し易い前記第2の日時であり、

前記判断部は、2つの割合のうち一方の前記割合を前記第1の割合と特定し、他方の前記割合を前記第2の割合と特定し、

前記故障時期変更部は、第1の前記ストレージでの前記割合を、特定された一方の前記割合にさせ、特定された他方の前記割合にはさせない、

請求項1記載のアレイ管理装置。

[請求項5]

消失が発生し難い前記第1の日時は、第2の前記ストレージが故障する日時から閾値以上に離れた日時であり、

消失が発生し易い前記第2の日時は、第2の前記ストレージが故障する日時から閾値以内の日時である、

請求項4記載のアレイ管理装置。

[請求項6]

複数のストレージが構成するアレイを管理する集積回路であって、前記複数のストレージのうち少なくとも2台のストレージが同時期に故障することによって、前記アレイが保持するデータが消失する可能性があるか否かを判断する判断部と、

前記アレイが保持する前記データが消失する可能性があると判断された場合には、前記少なくとも2台のストレージのうちの少なくとも一つの故障時期を変更する、故障時期変更部と、

を備える集積回路。

[請求項7]

複数のストレージが構成するアレイを管理するアレイ管理方法であって、

前記複数のストレージのうち少なくとも2台のストレージが同時期に故障することによって、前記アレイが保持するデータが消失する可能性があるか否かを判断する判断ステップと、

前記アレイが保持する前記データが消失する可能性があるとして判断さ

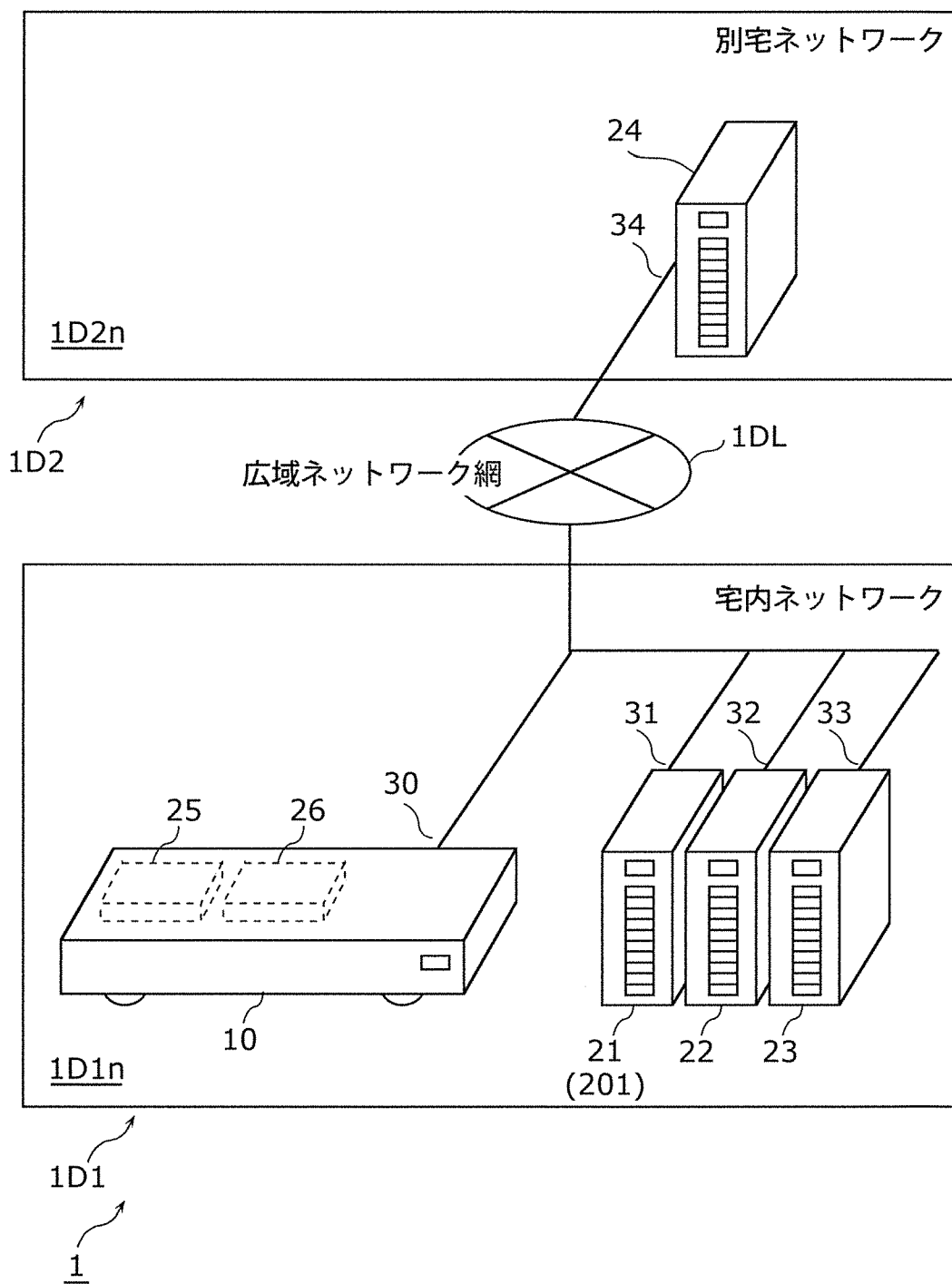
れた場合には、前記少なくとも2台のストレージのうちの少なくとも一つの故障時期を変更する、故障時期変更ステップと、  
を含むアレイ管理方法。

[請求項8]

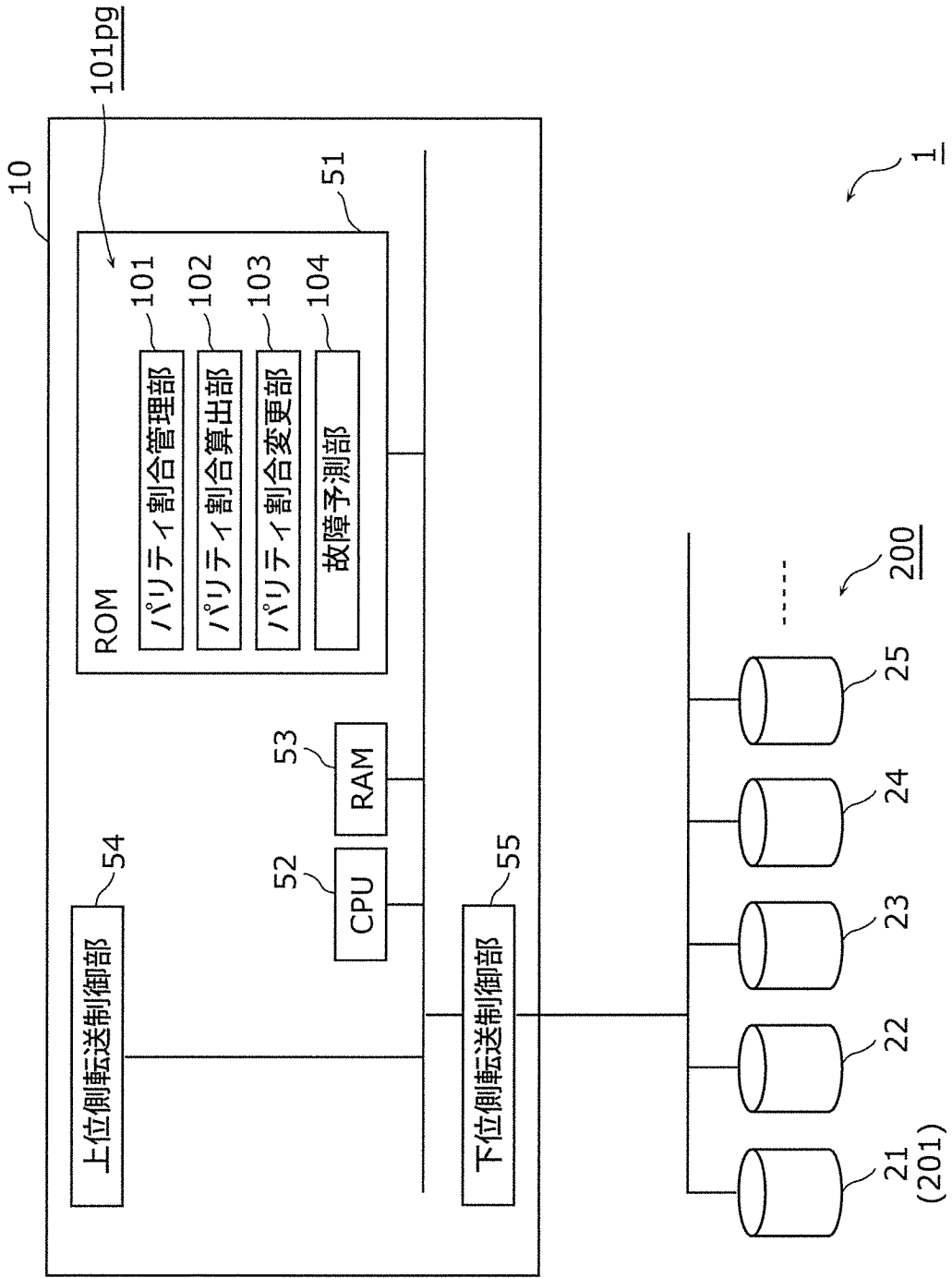
複数のストレージが構成するアレイを管理するコンピュータに、  
前記複数のストレージのうち少なくとも2台のストレージが同時期に故障することによって、前記アレイが保持するデータが消失する可能性があるか否かを判断する判断ステップと、

前記アレイが保持する前記データが消失する可能性があると判断された場合には、前記少なくとも2台のストレージのうちの少なくとも一つの故障時期を変更する、故障時期変更ステップと、  
を実行させるためのコンピュータプログラム。

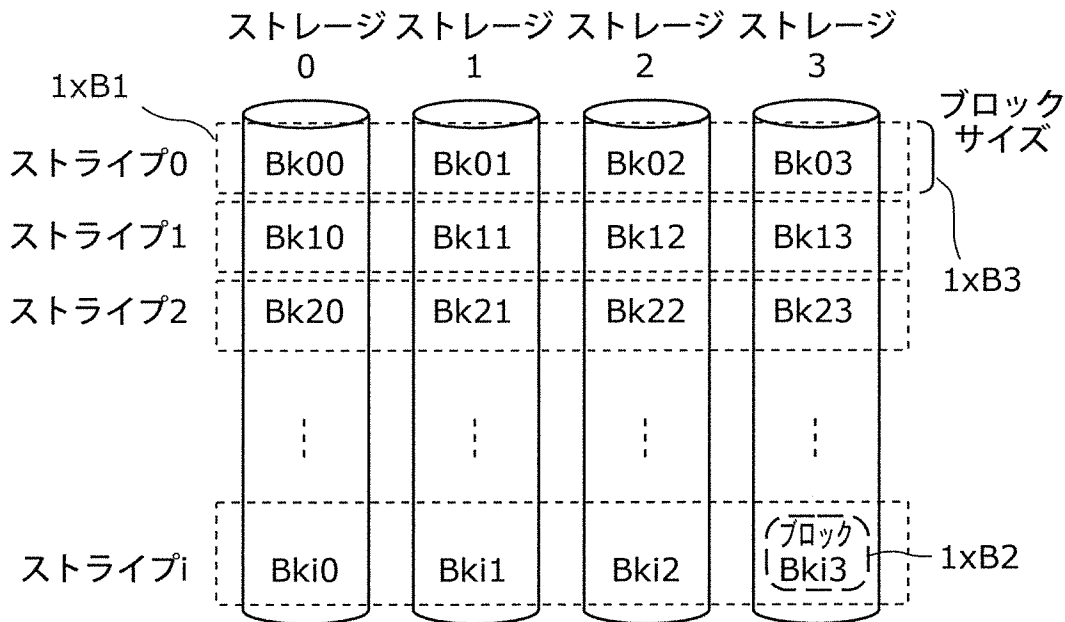
[図1]



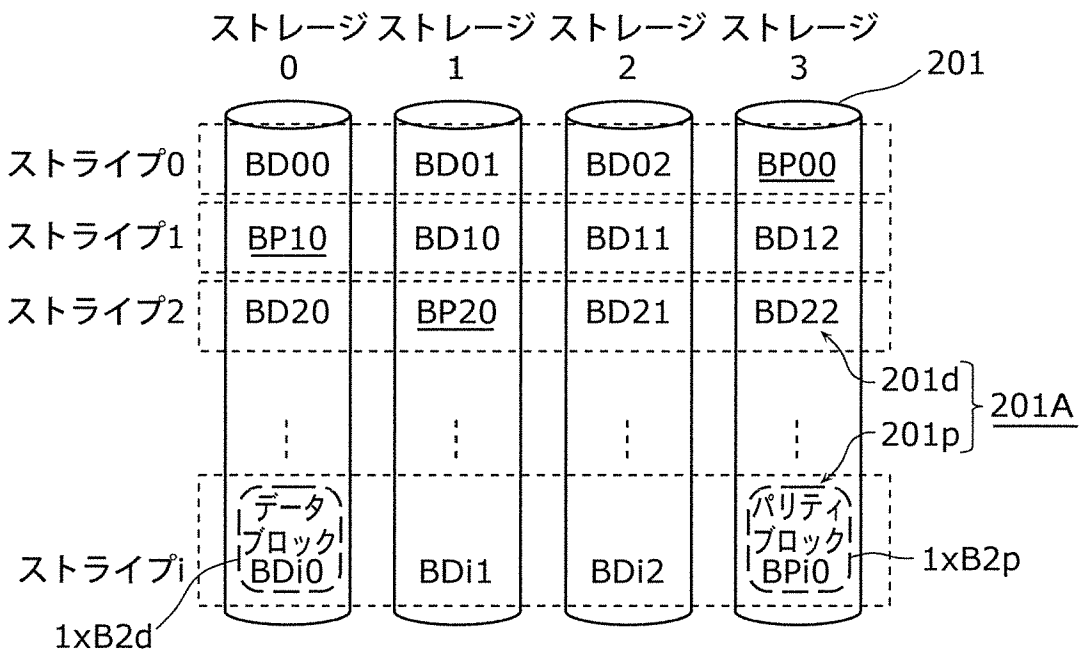
[図2]



[図3]



[図4]

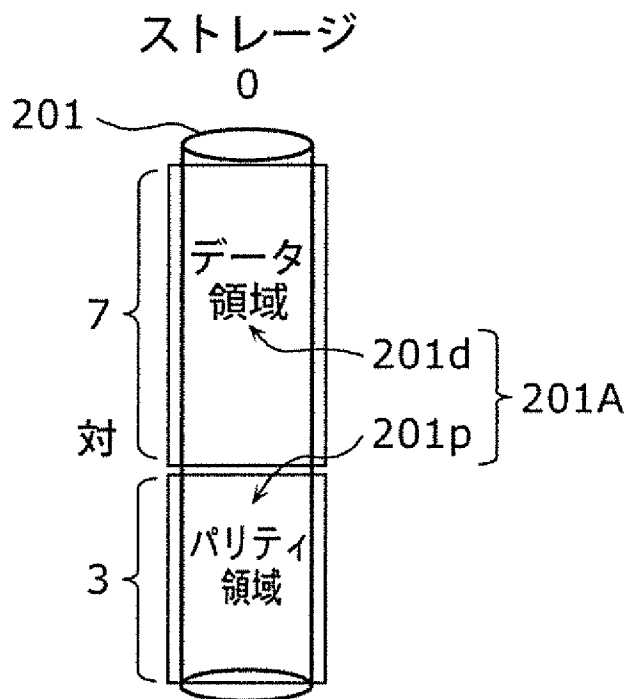


[図5]

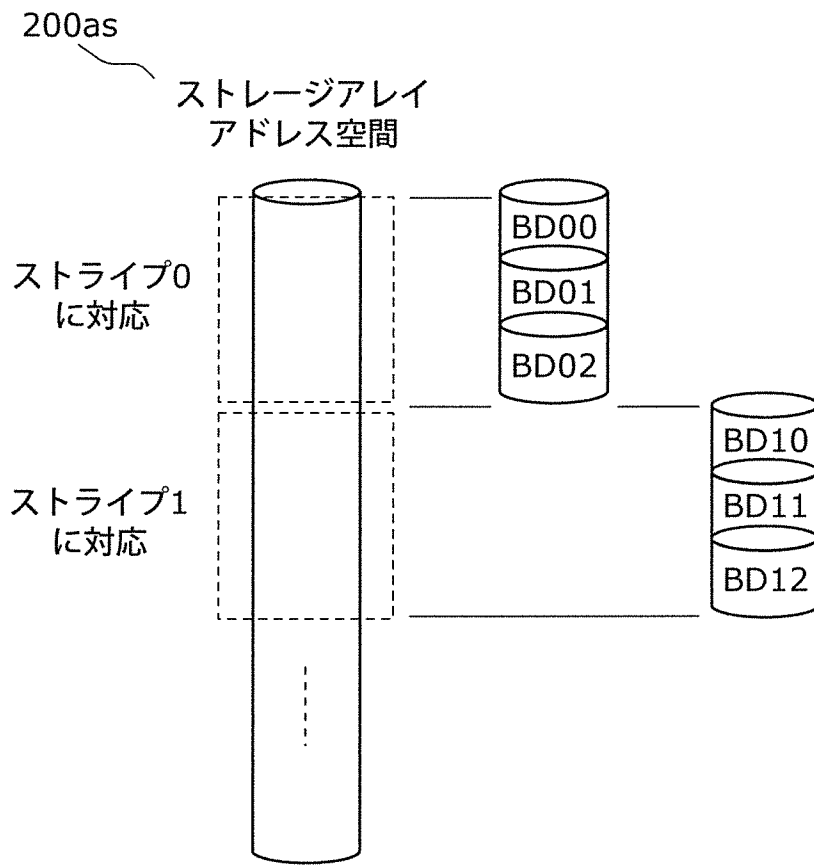
	ストレージ0	ストレージ1	ストレージ2	ストレージ3
ストライプ0	BD00	BD01	BD02	<u>BP00</u>
ストライプ1	<u>BP10</u>	BD10	BD11	BD12
ストライプ2	BD20	<u>BP20</u>	BD21	BD22
-----				
ストライプi	BDi0	BDi1	BDi2	BPi0
各ストレージに おけるパリティ割合	30%	20%	15%	35%

201r  
(201r1,201r2)

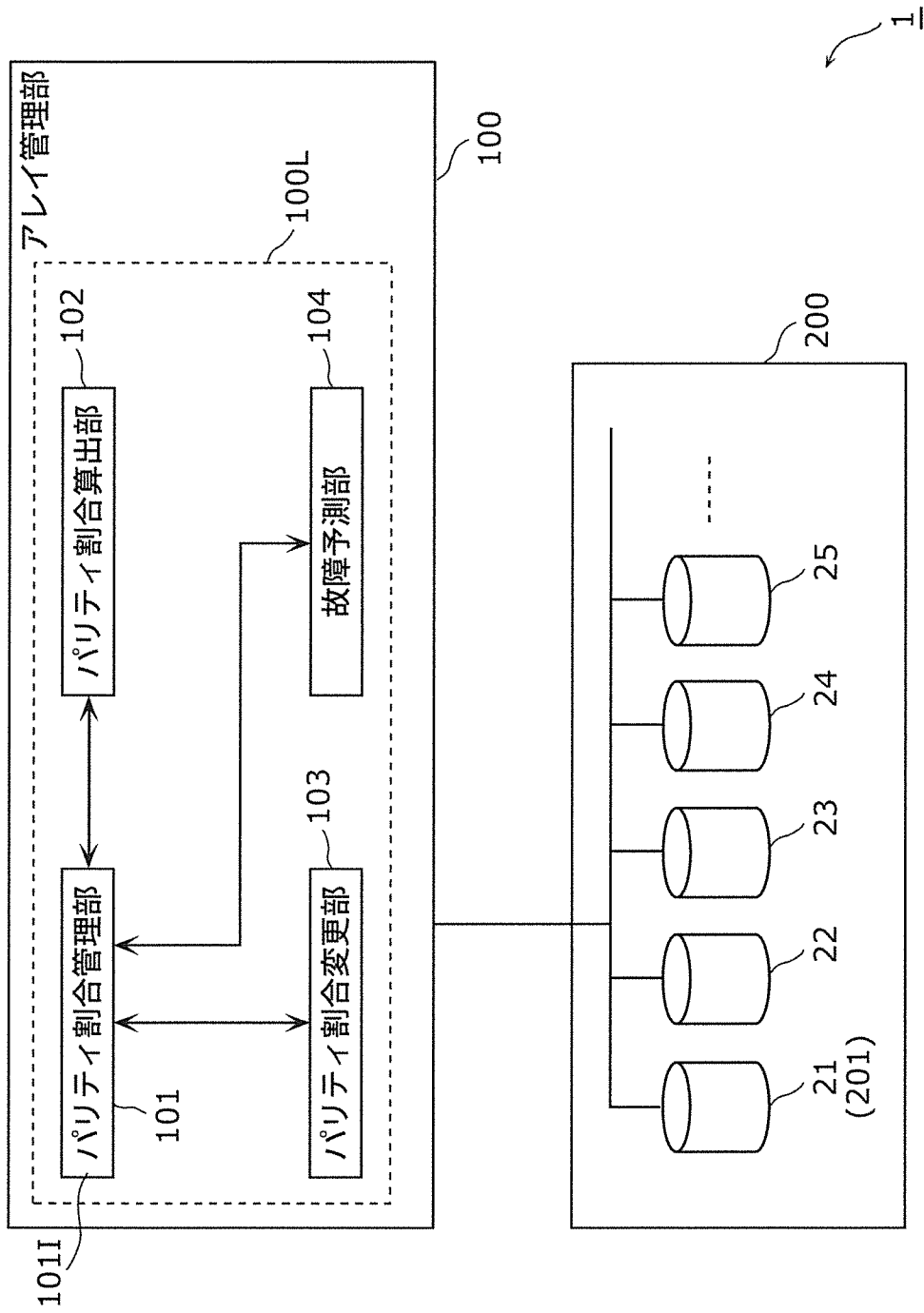
[図6]



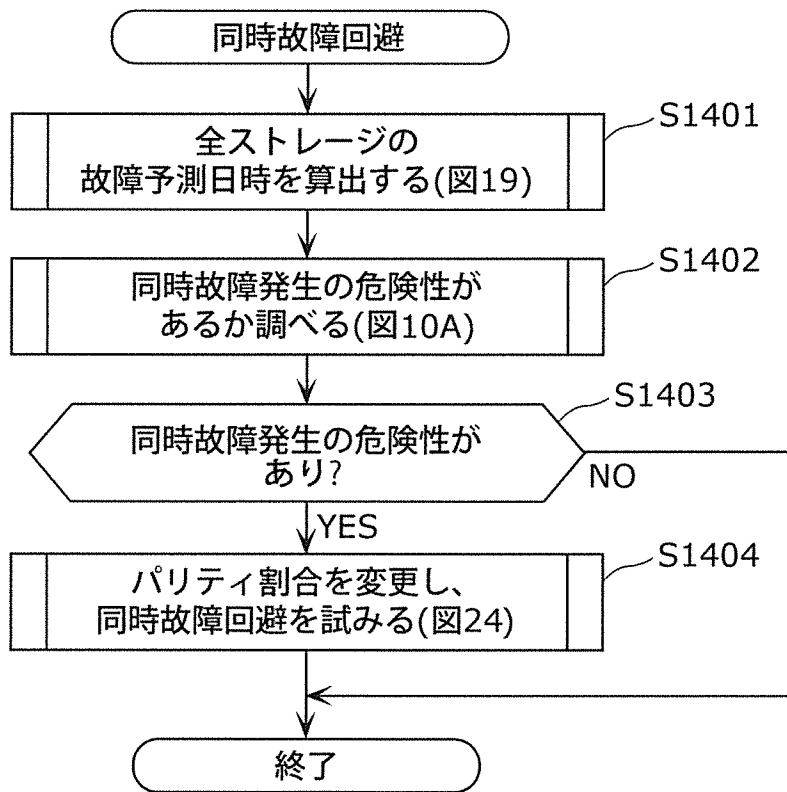
[図7]



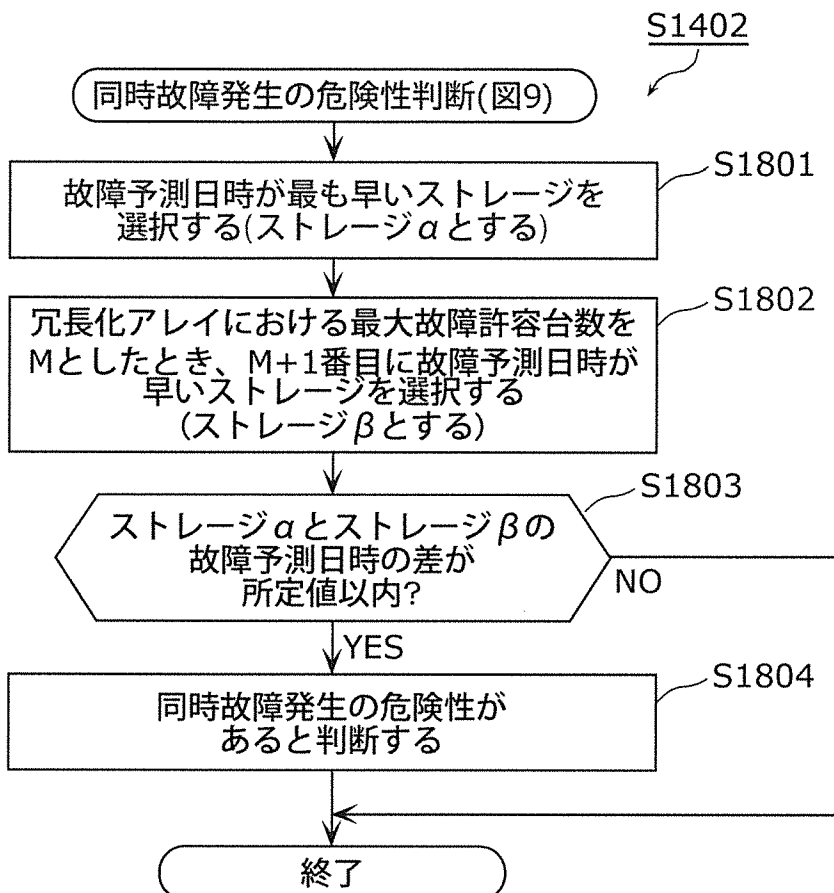
[図8]



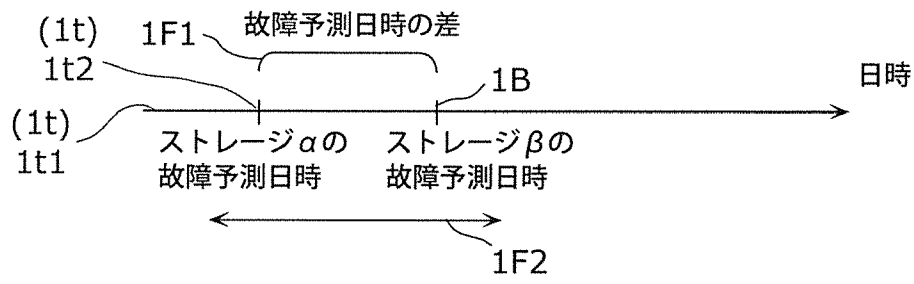
[図9]



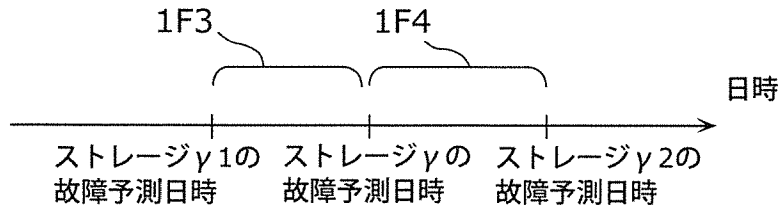
[図10A]



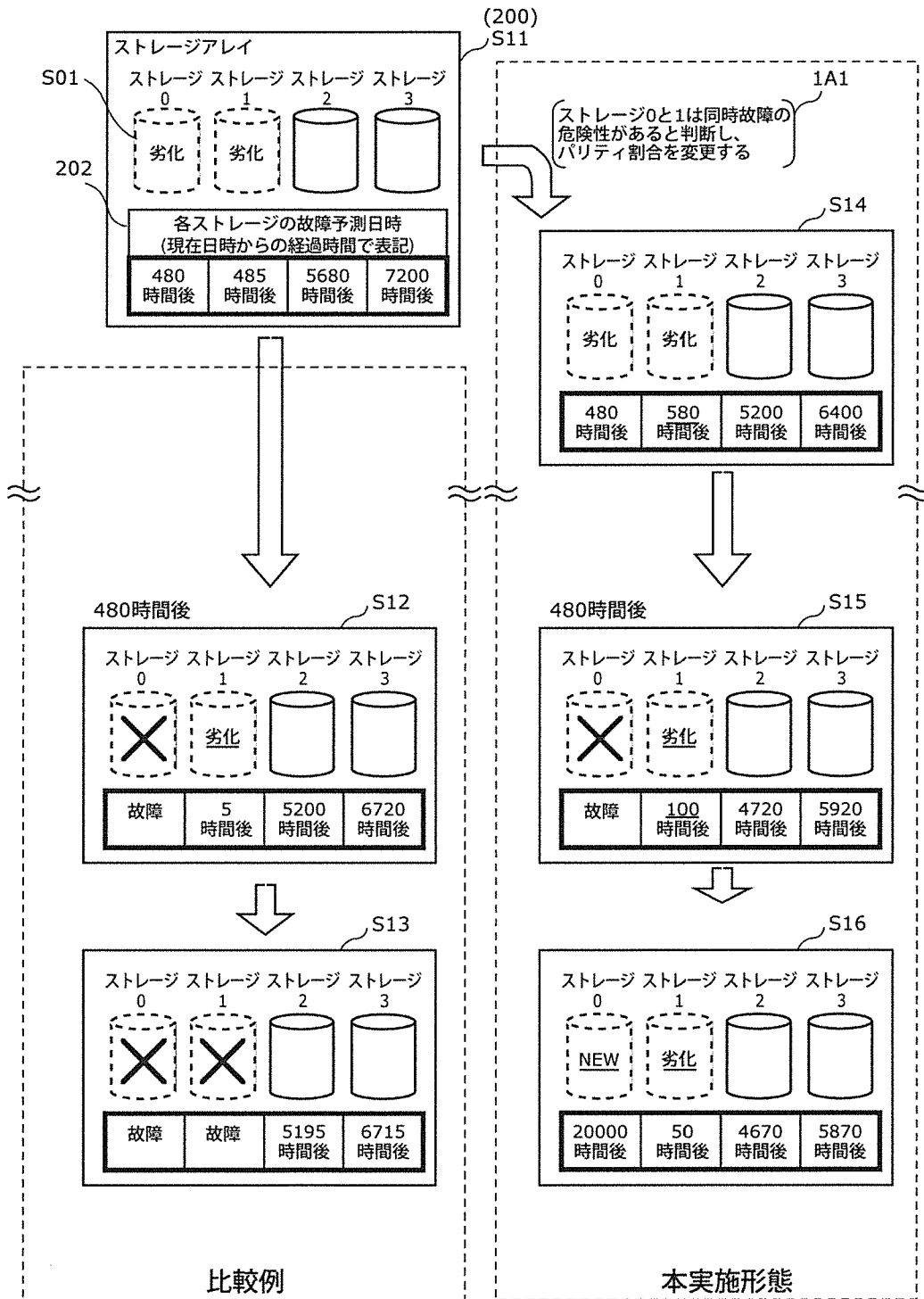
[図10B]



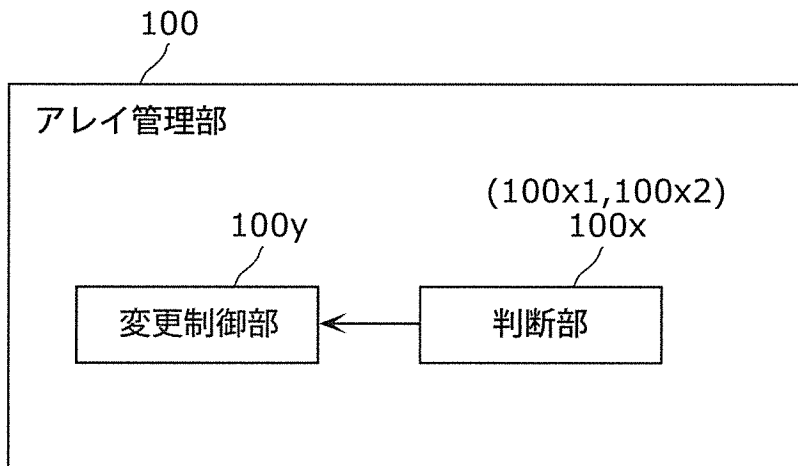
[図10C]



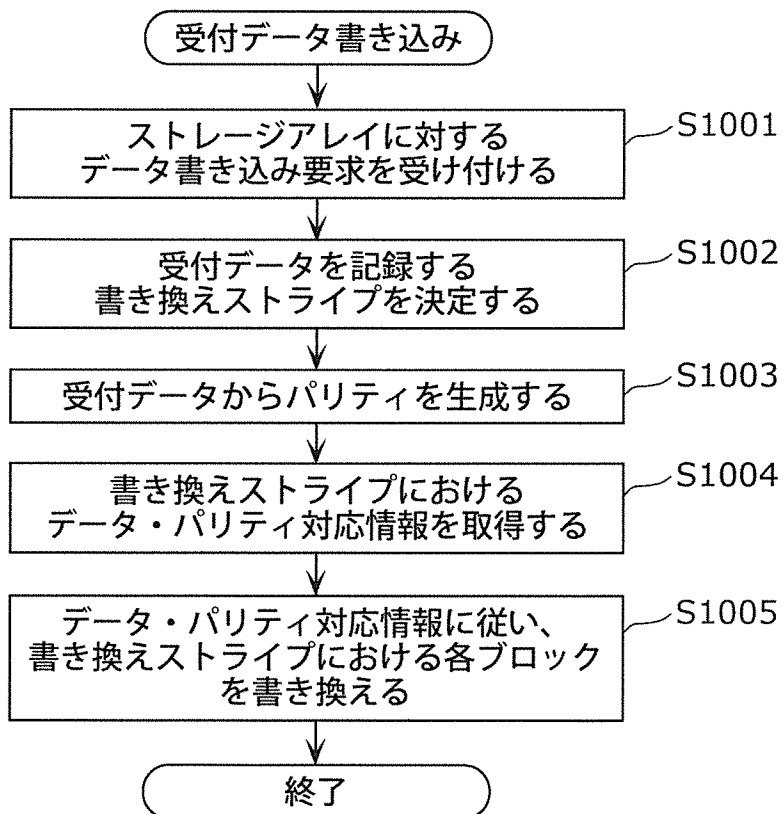
[図11]



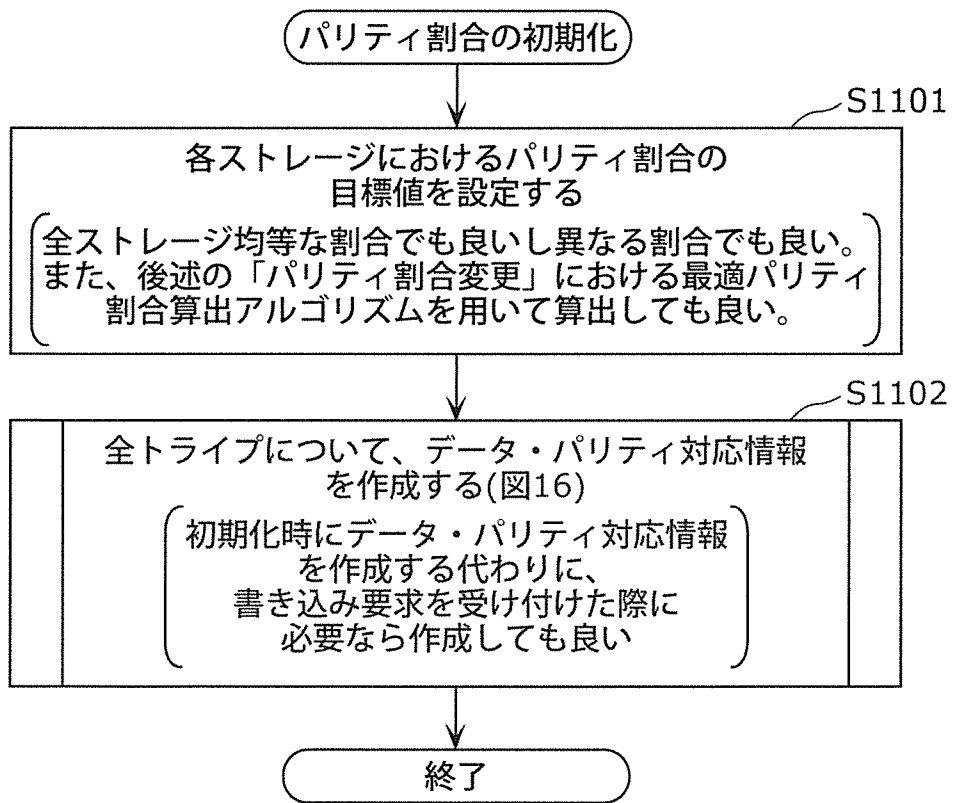
[図12]



[図13]



[図14]



[図15]

	ストレージ0	ストレージ1	ストレージ2	ストレージ3
ストライプ0	—	—	—	—
ストライプ1	—	—	—	—
ストライプ2	—	—	—	—
-----				
各ストレージにおけるパリティ割合	0%	0%	0%	0%
目標値	30%	20%	15%	35%
目標値から現在値を減じた値	30%	20%	15%	35%

↓ (ストライプ0において、ストレージ3に対応するブロックにパリティを割り当て、それ以外のブロックについてはデータを割り当てる) — 1C1

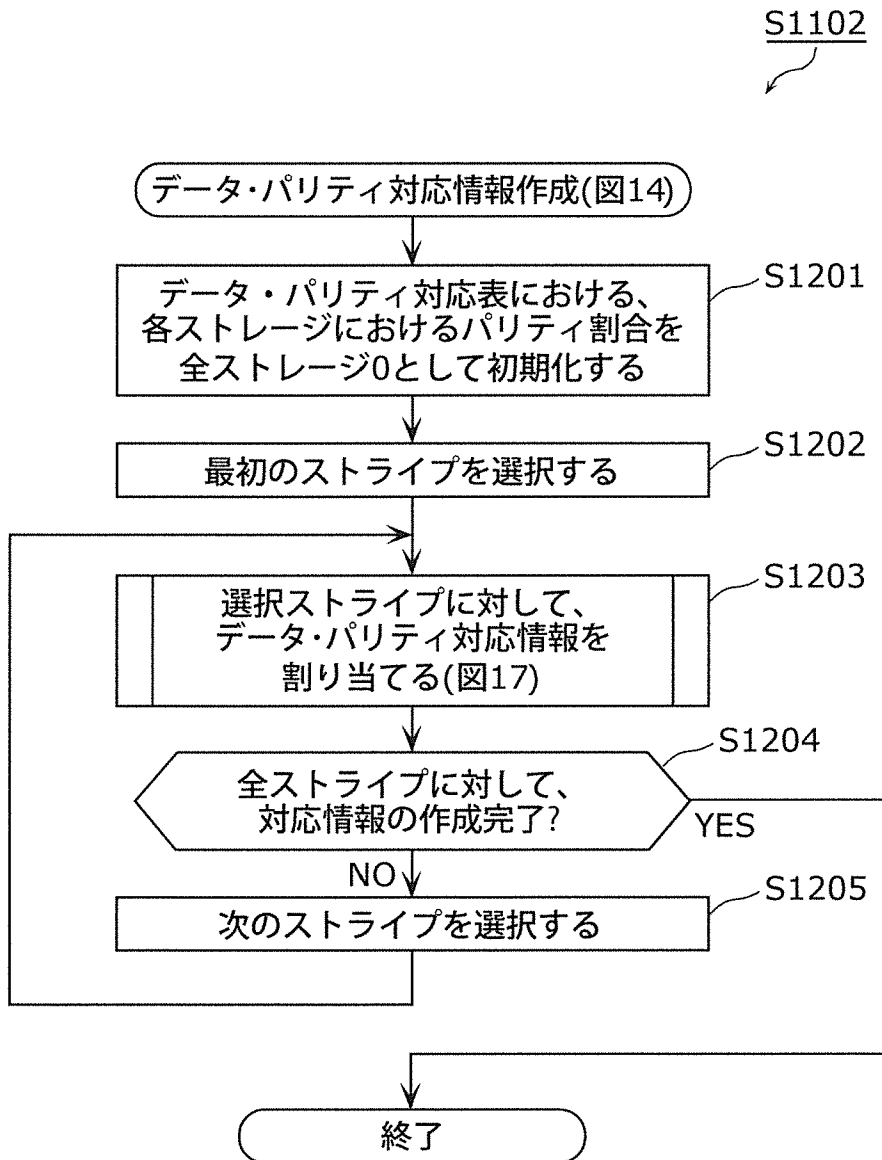
	ストレージ0	ストレージ1	ストレージ2	ストレージ3
ストライプ0	BD00	BD01	BD02	<u>BP00</u>
ストライプ1	—	—	—	—
ストライプ2	—	—	—	—
-----				
各ストレージにおけるパリティ割合	0%	0%	0%	0%
目標値	30%	20%	15%	35%
目標値から現在値を減じた値	30%	20%	15%	-65%

↓ (ストライプ1において、ストレージ0に対応するブロックにパリティを割り当て、それ以外のブロックについてはデータを割り当てる) — 1C2

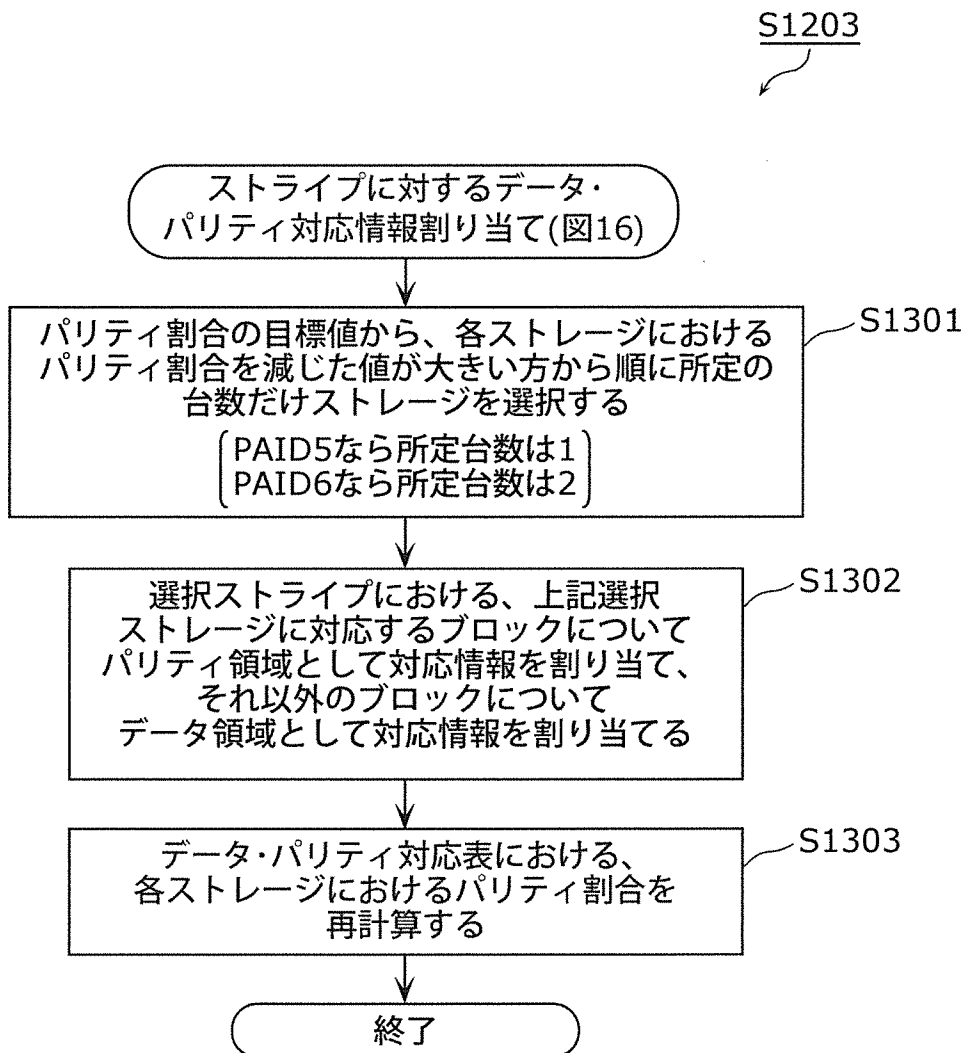
	ストレージ0	ストレージ1	ストレージ2	ストレージ3
ストライプ0	BD00	BD01	BD02	BP00
ストライプ1	BP10	BD10	BD11	BD12
ストライプ2	—	—	—	—
-----				
各ストレージにおけるパリティ割合	50%	0%	0%	0%
目標値	30%	20%	15%	35%
目標値から現在値を減じた値	-20%	20%	15%	-15%

以下、同様に繰り返す

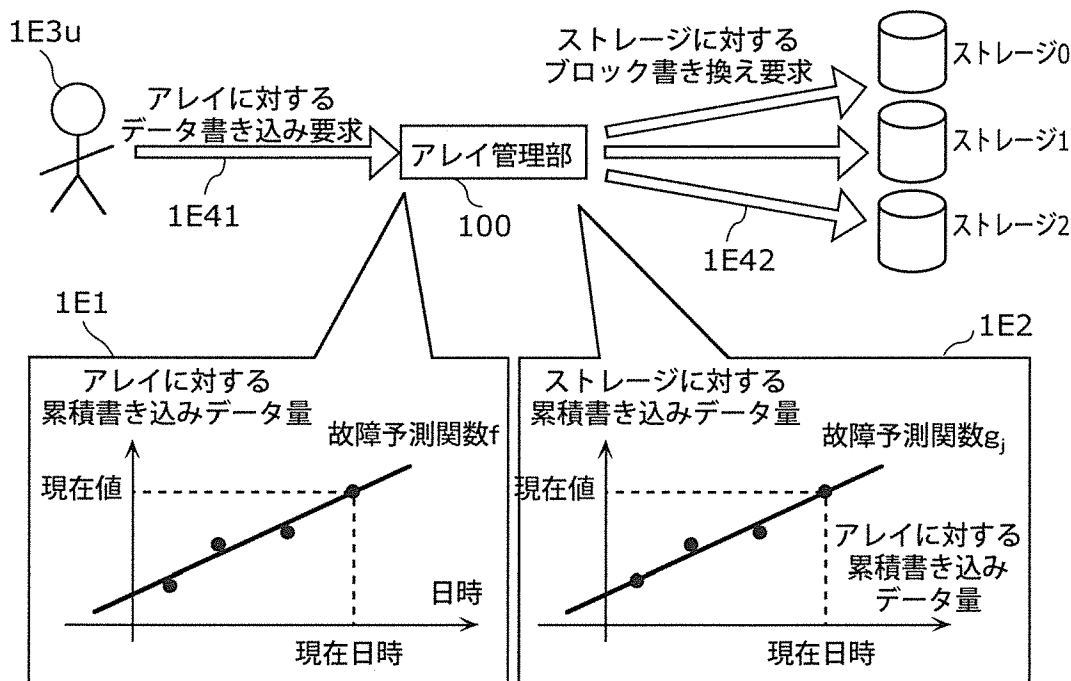
[図16]



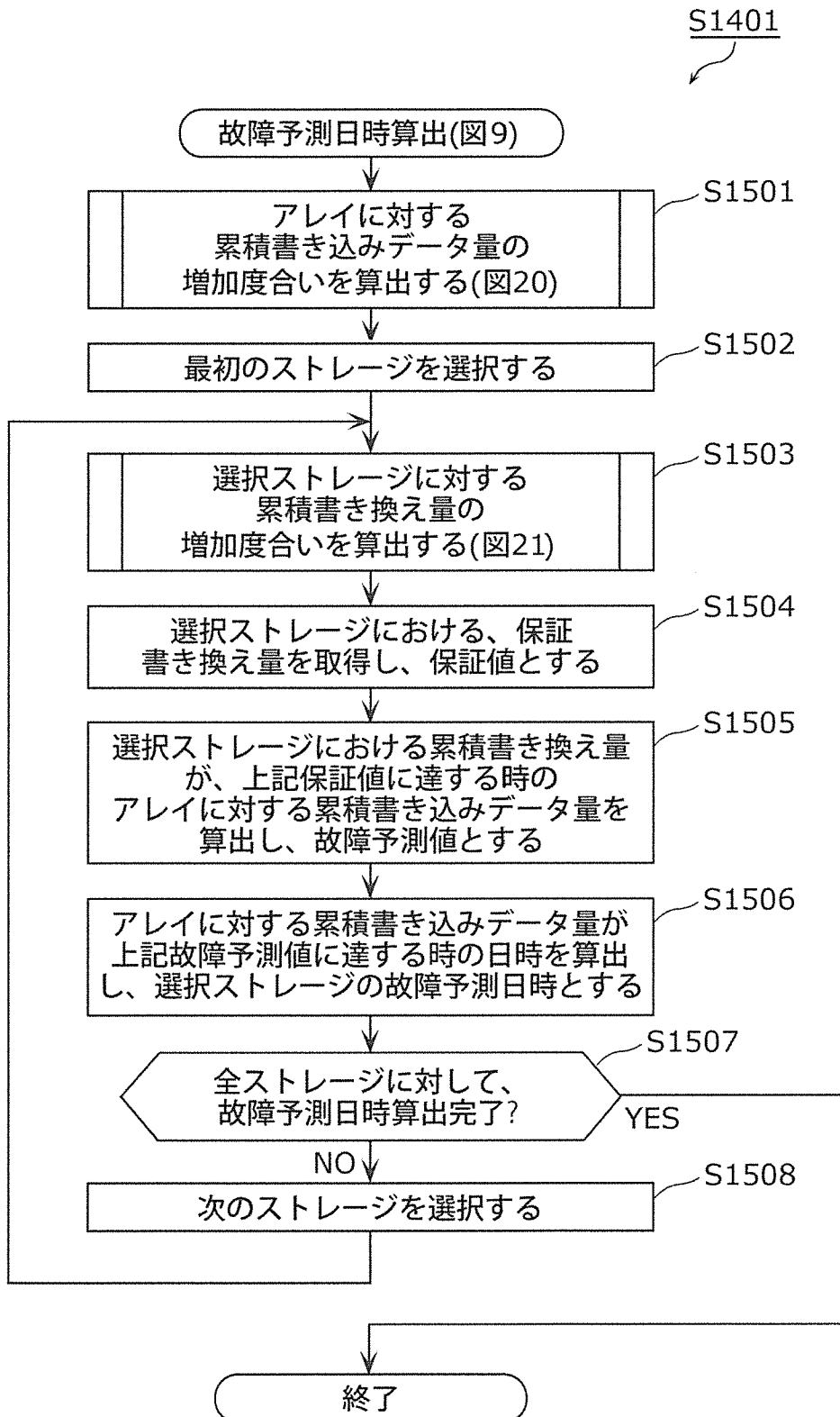
[図17]



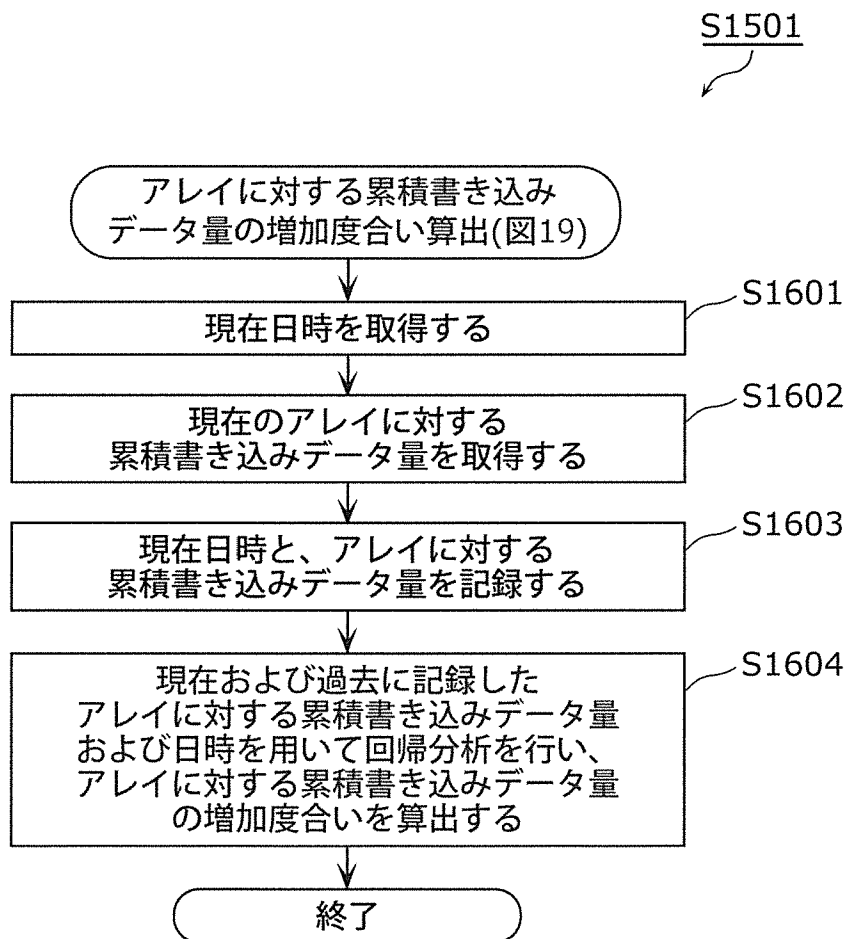
[図18]



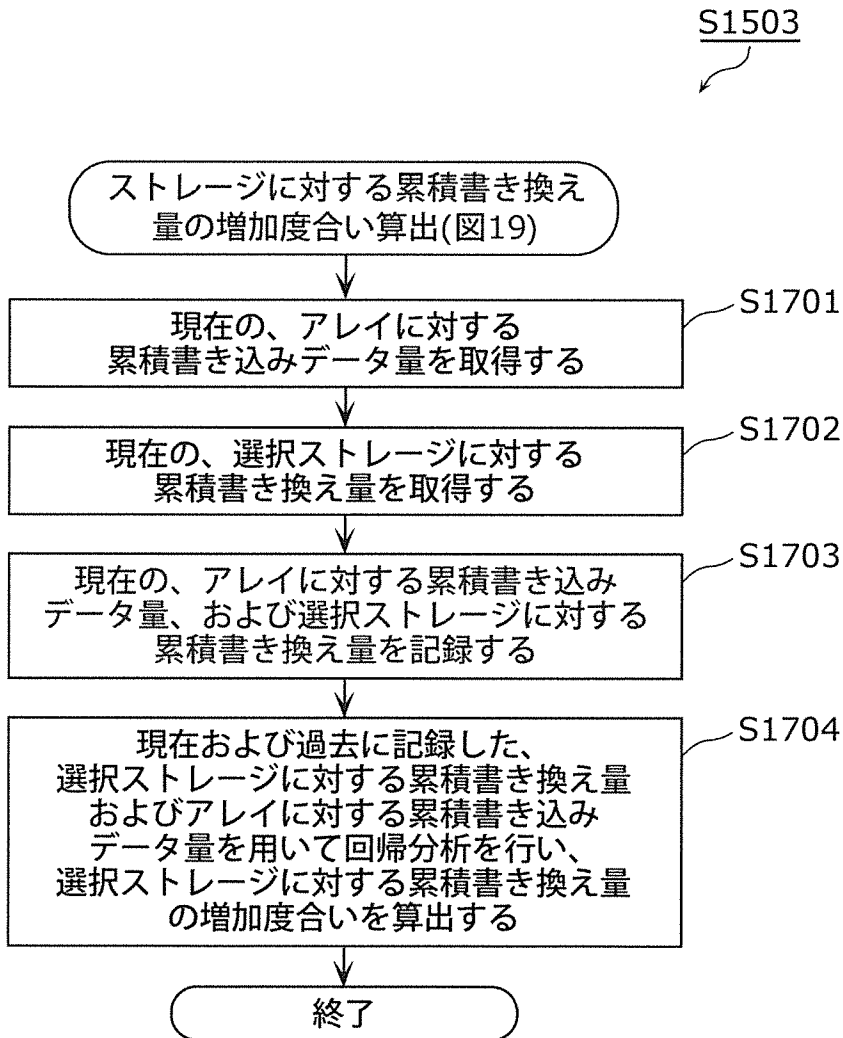
[図19]



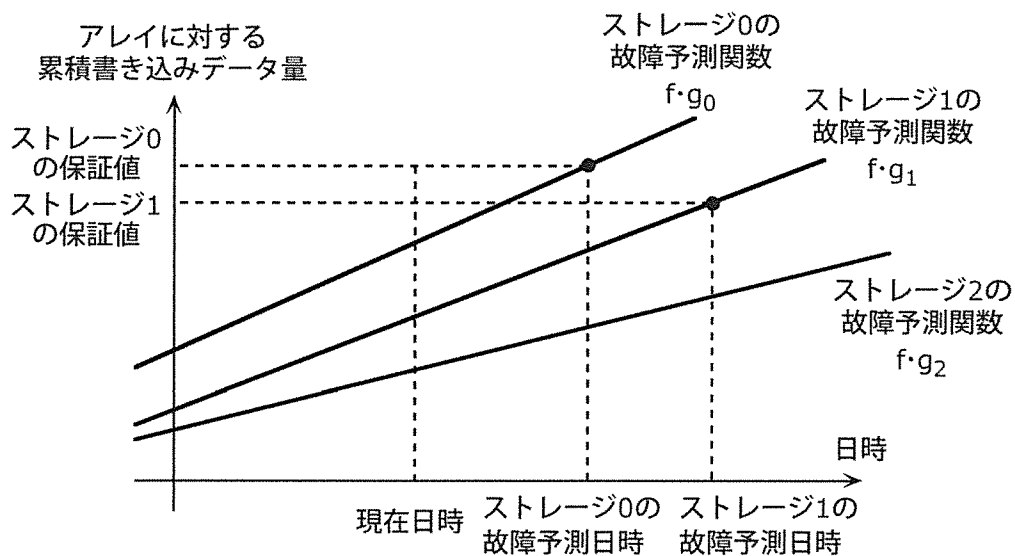
[図20]



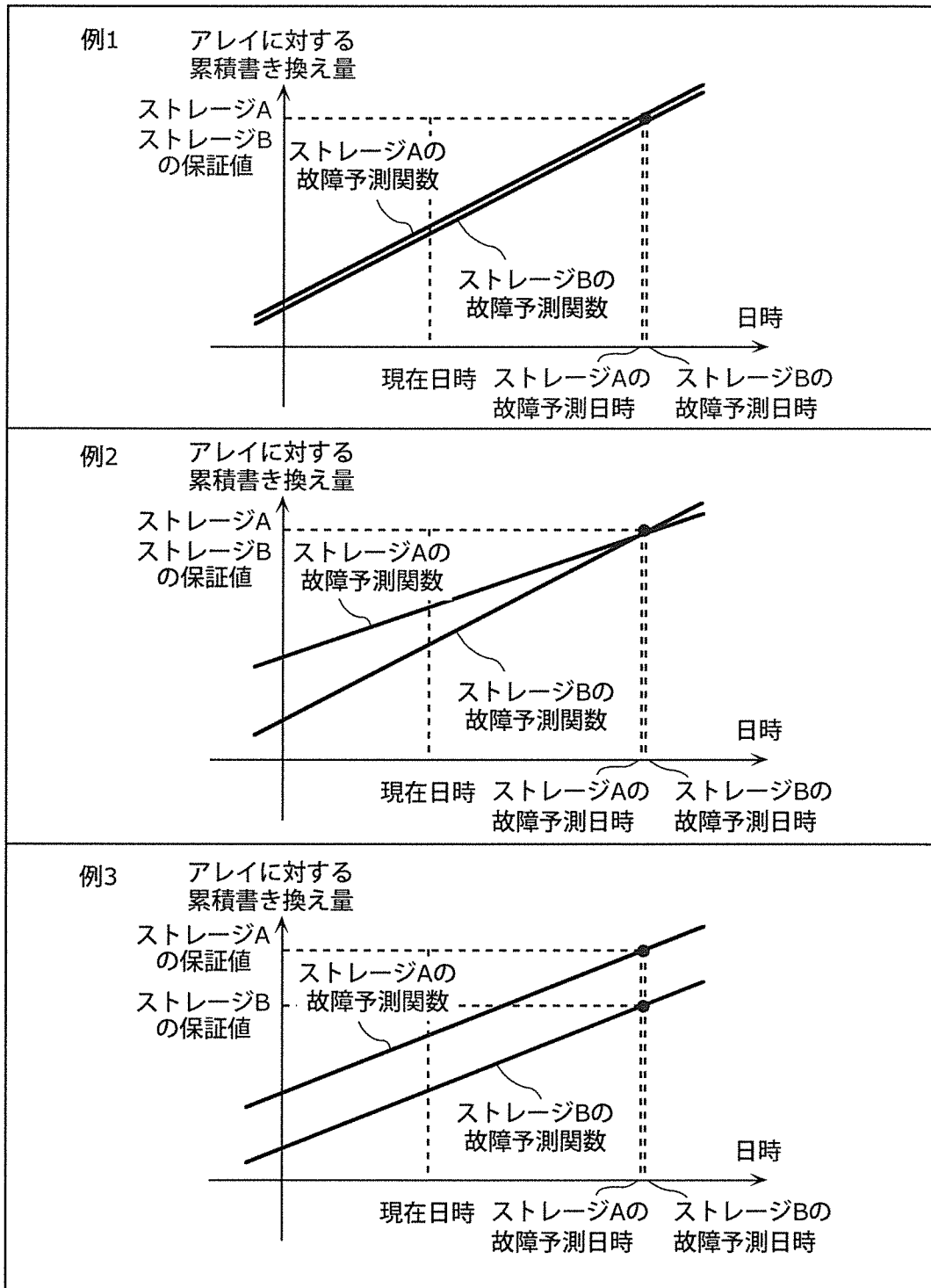
[図21]



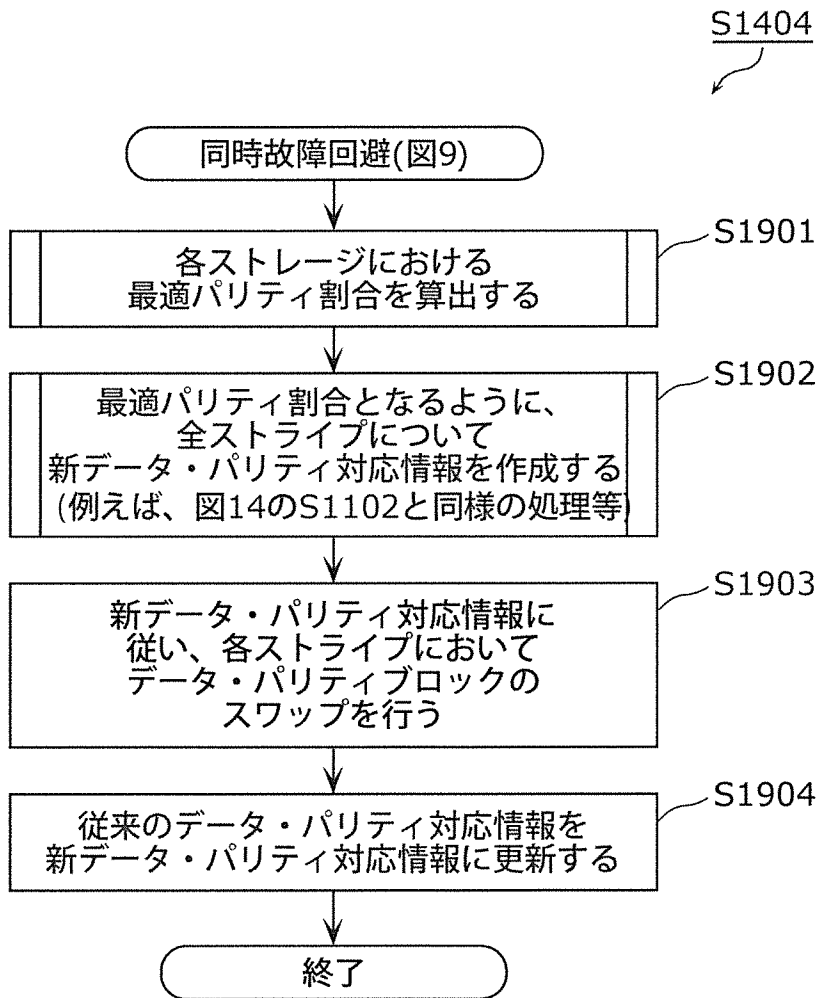
[図22]



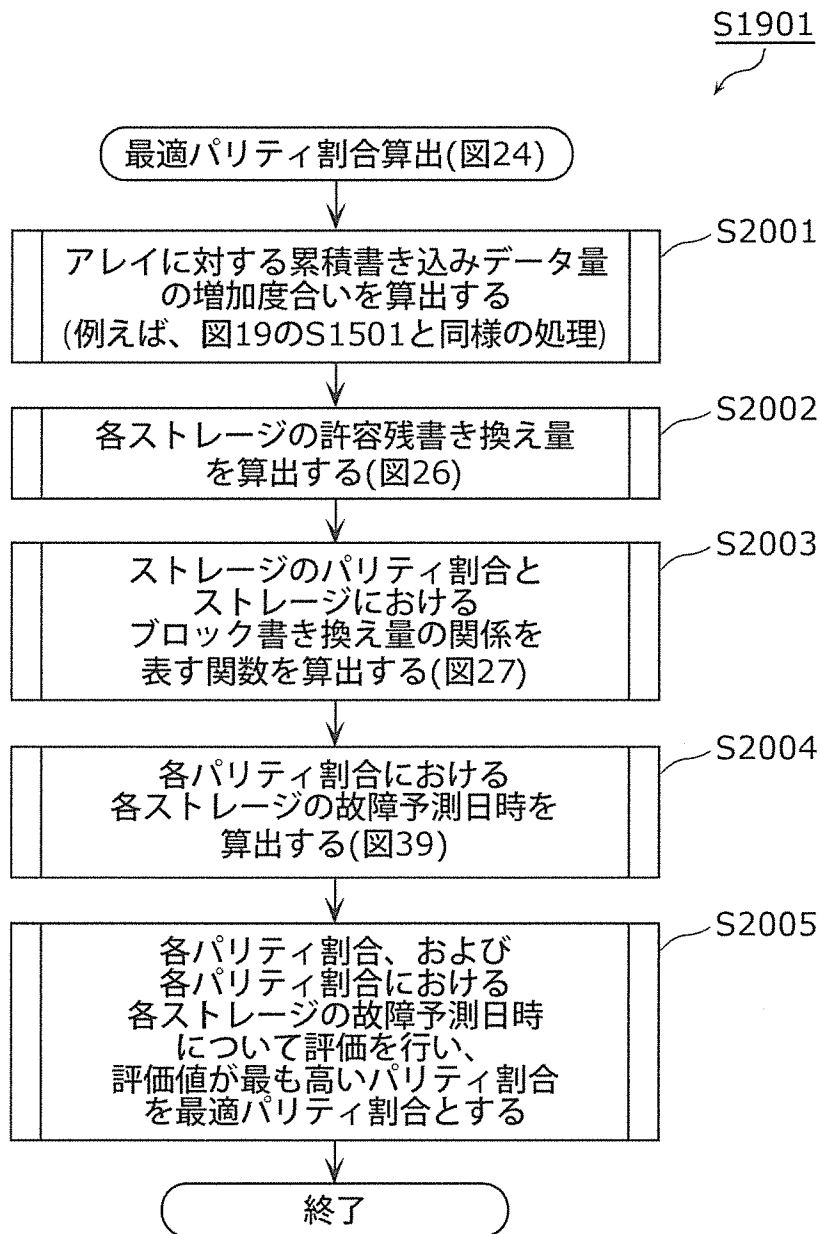
[図23]



[図24]

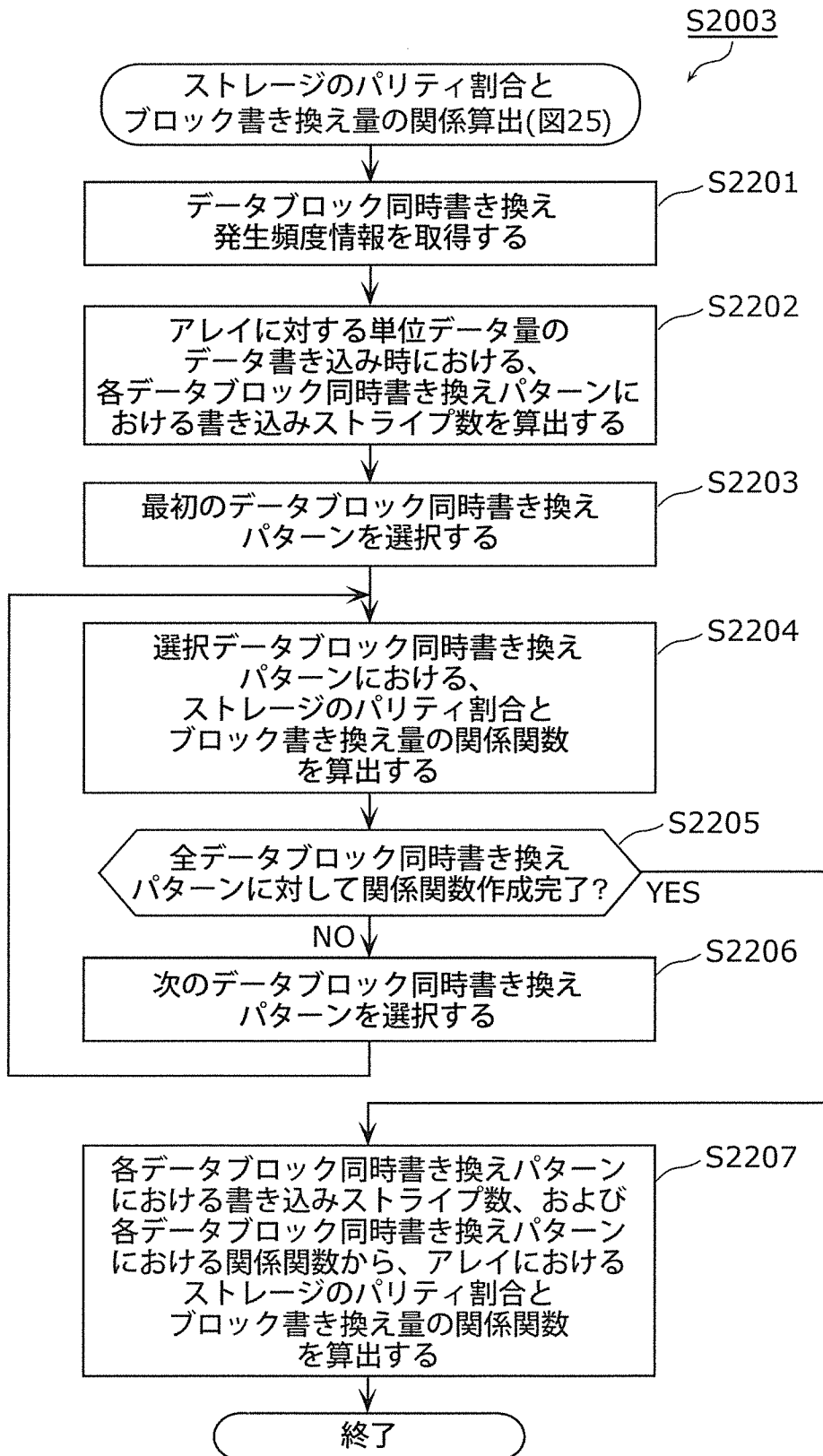


[図25]

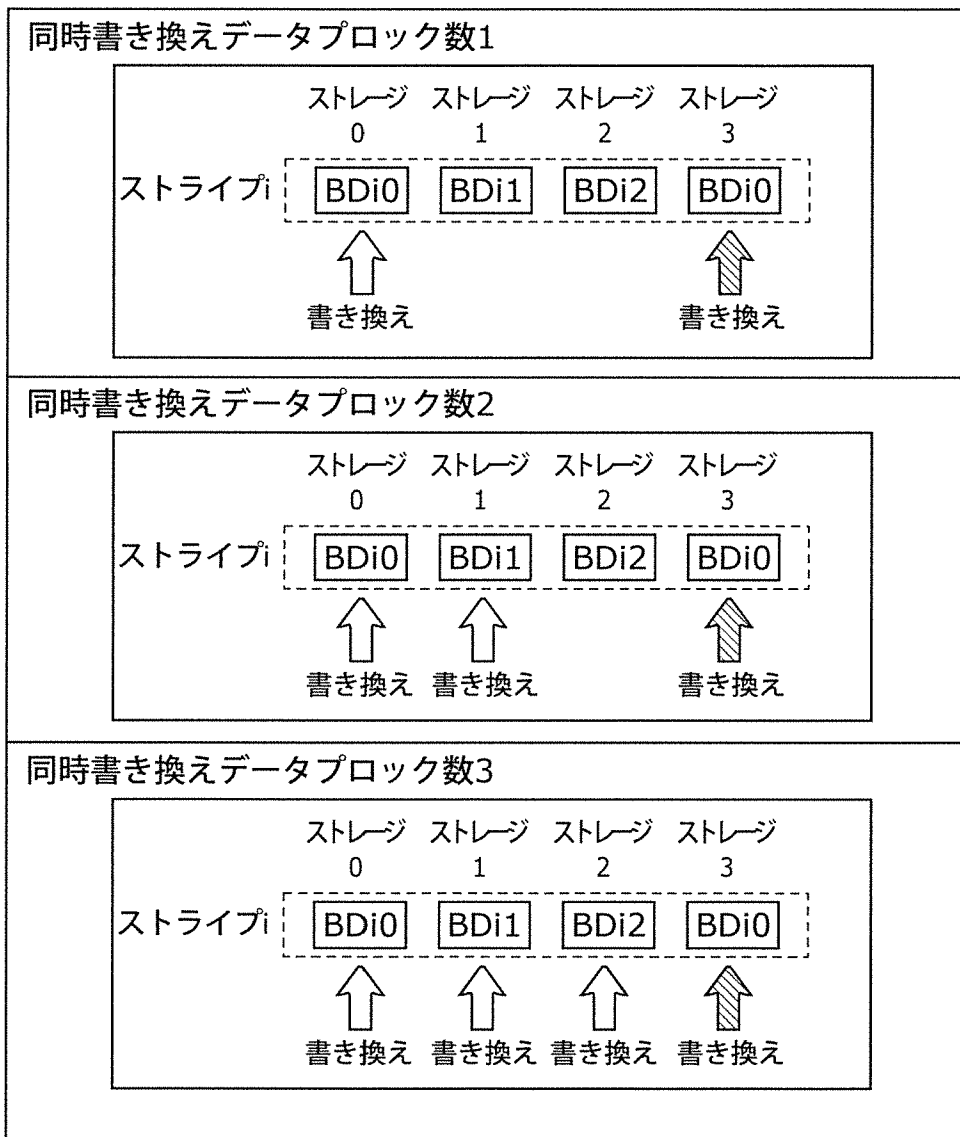




[図27]



[図28]



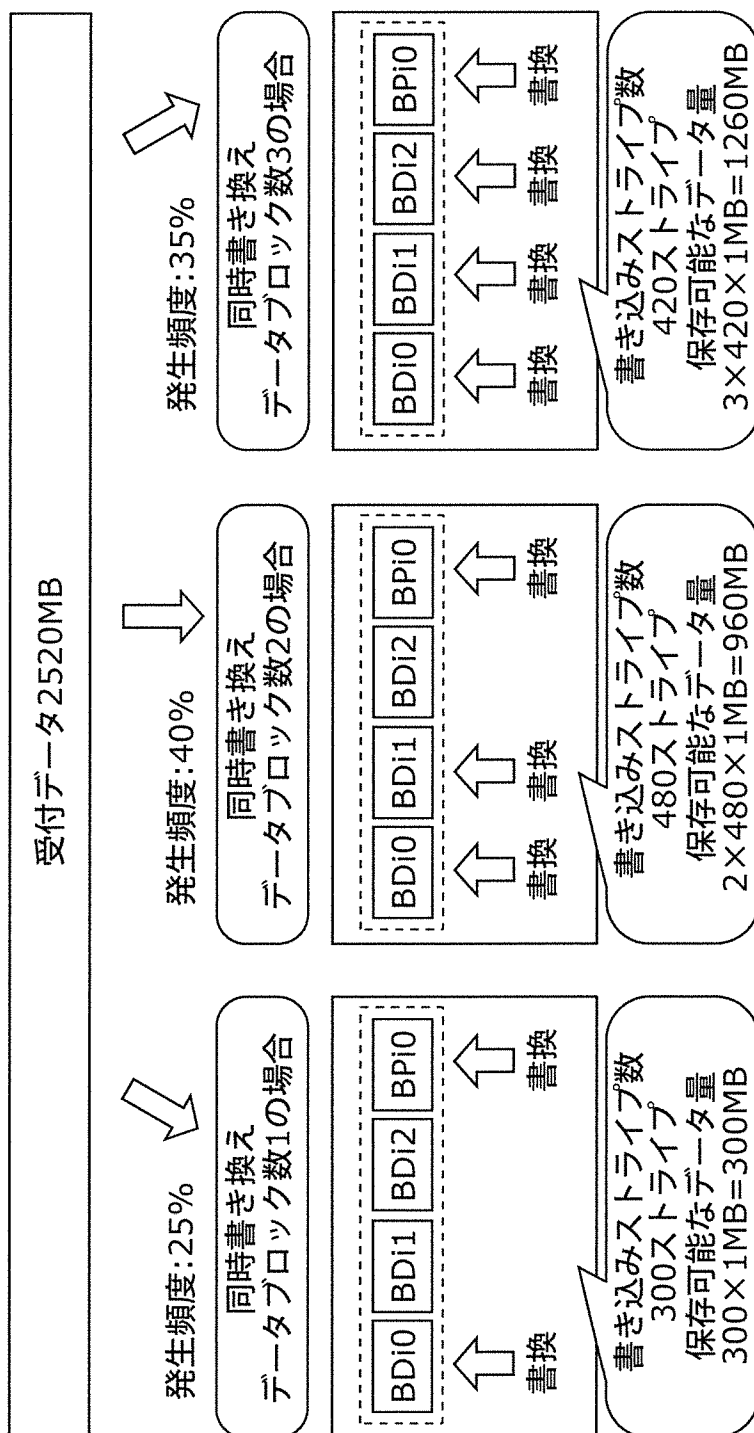
[図29]

同時書き換えデータブロック数	1回のストライプ書き込みで保存可能なデータ量
1	1MB
2	2MB
3	3MB

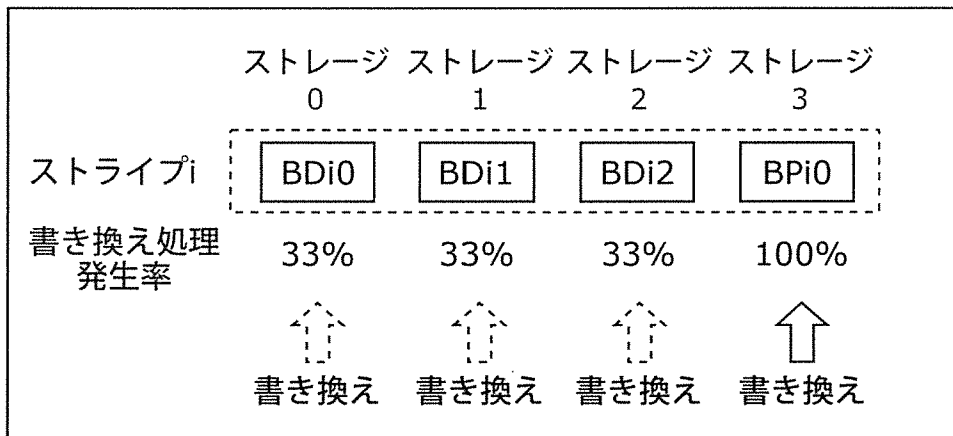
[図30]

同時書き換えデータブロック数	累積発生ストライプ数	発生頻度
1	60000	25%
2	96000	40%
3	84000	35%
(合計)	240000	

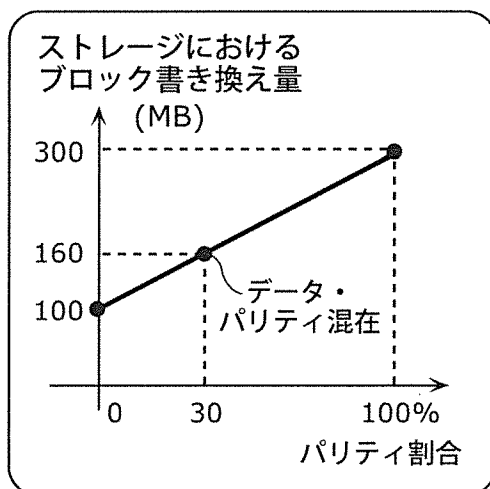
[図31]



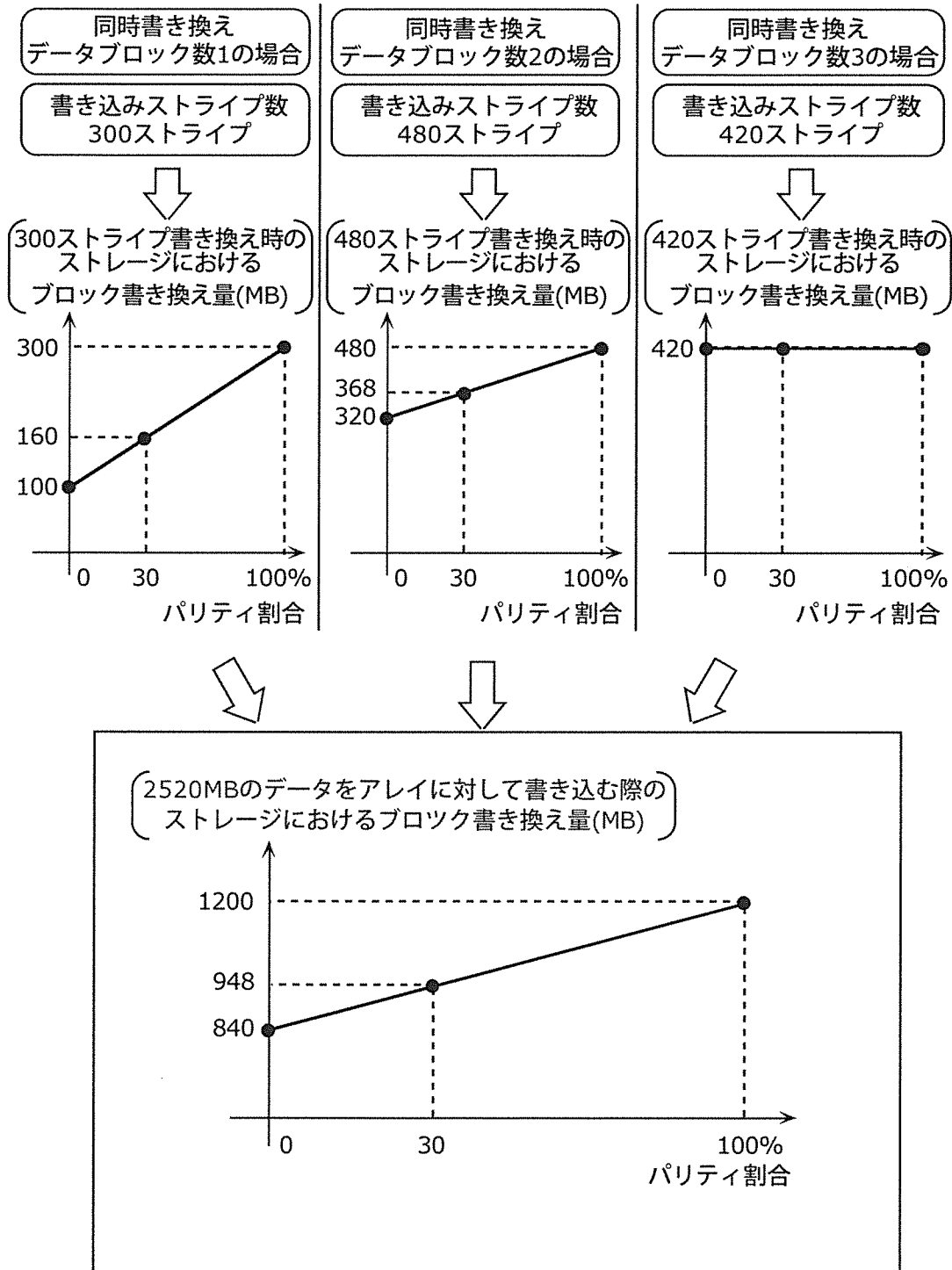
[図32]



[図33]



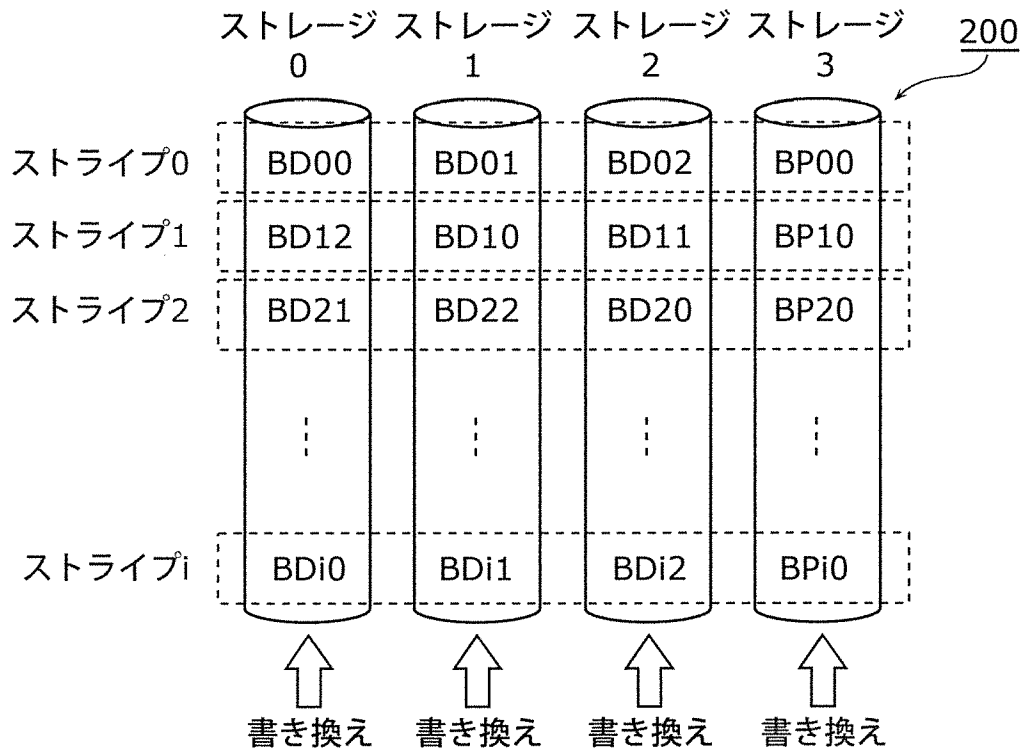
[図34]



[図35A]

	ストレージ				合計
	0	1	2	3	
パリティ割合	0%	0%	0%	100%	100%

[図35B]



[図35C]

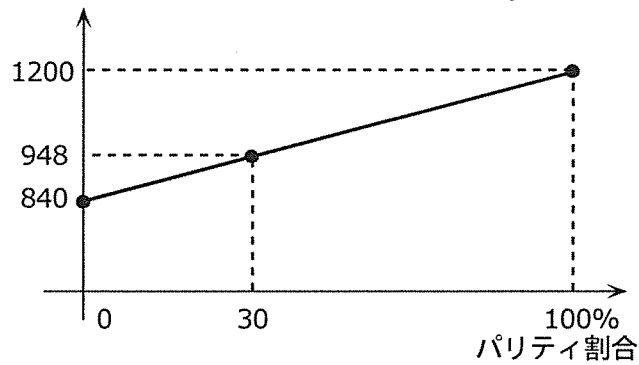
	ストレージ				合計
	0	1	2	3	
書き換え量	100MB	100MB	100MB	300MB	600MB

[図36]

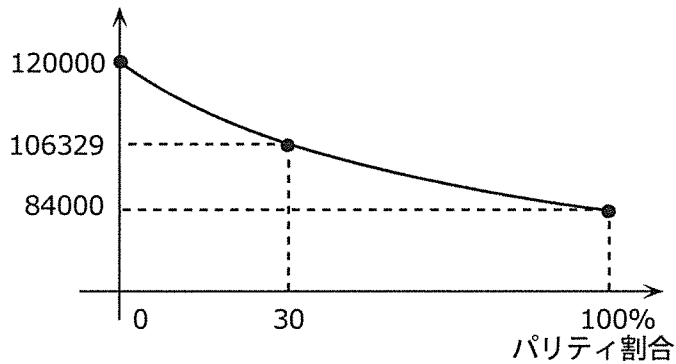
	ストレージ				合計
	0	1	2	3	
パリティ割合	30%	20%	15%	35%	100%

[図37]

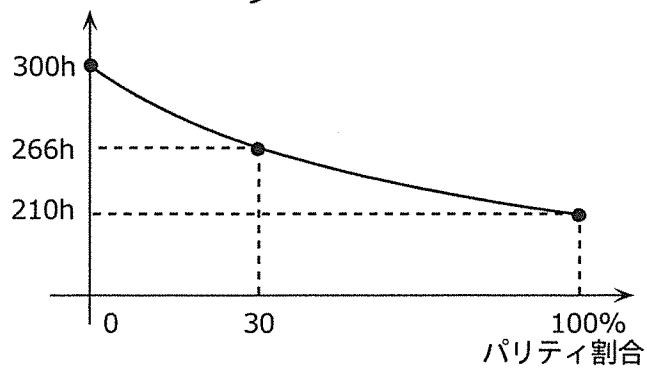
(2520MBのデータをアレイに対して書き込む際の  
ストレージにおけるブロック書き換え量(MB))



(ストレージに対して、40000MBのブロック書き換えを行う際の、  
アレイに対するデータ書き込み量(MB))

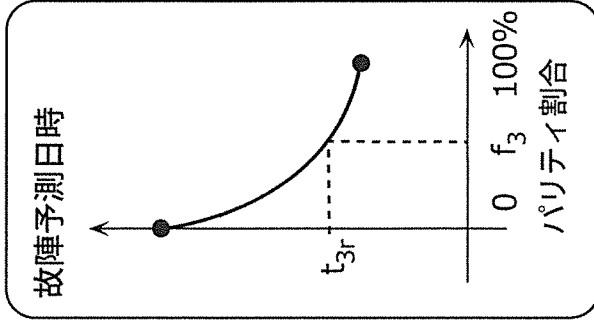


(故障予測日時  
(現在日時からの経過時間))

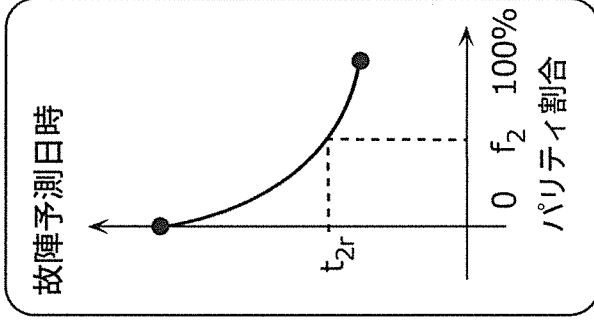


[図38]

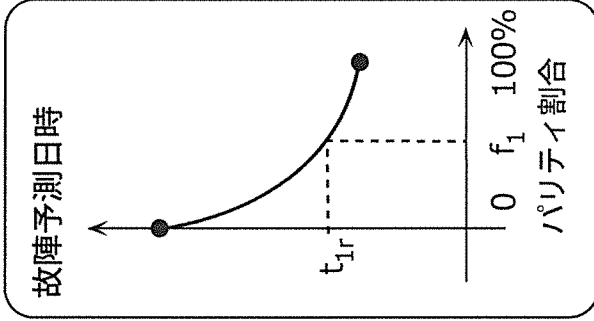
ストレージ3



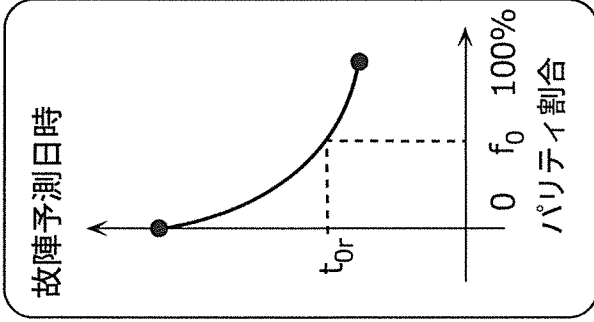
ストレージ2



ストレージ1

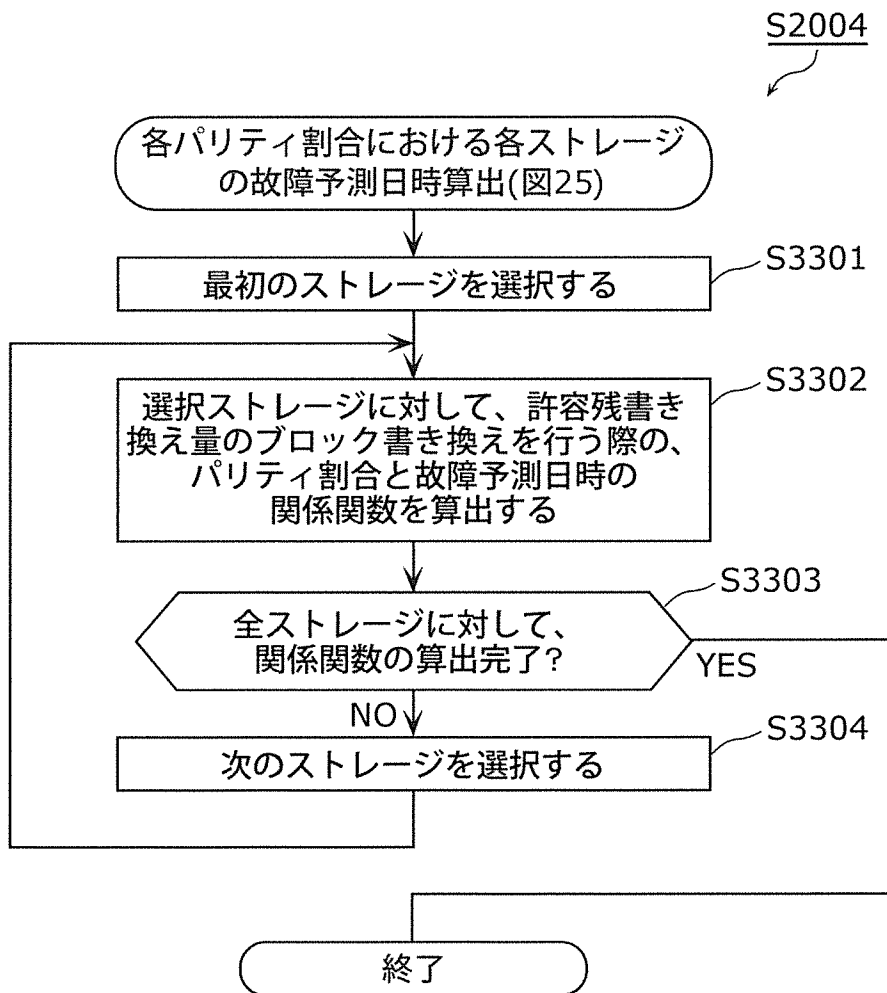


ストレージ0



1G  
 (最適なパリテイ割合を算出する)

[図39]



## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2012/002116

## A. CLASSIFICATION OF SUBJECT MATTER

G06F13/10(2006.01) i, G06F3/06(2006.01) i

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

G06F13/10, G06F3/06

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Jitsuyo Shinan Koho	1922-1996	Jitsuyo Shinan Toroku Koho	1996-2012
Kokai Jitsuyo Shinan Koho	1971-2012	Toroku Jitsuyo Shinan Koho	1994-2012

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	JP 2010-015516 A (Toshiba Corp.), 21 January 2010 (21.01.2010), paragraphs [0006] to [0016], [0045], [0059] to [0070] & US 2010/0005228 A1	1-8
Y	JP 2009-037304 A (Hitachi, Ltd.), 19 February 2009 (19.02.2009), paragraphs [0064] to [0066] & US 2009/0037656 A1	1-3, 6-8
Y	JP 2007-206993 A (Fujitsu Ltd.), 16 August 2007 (16.08.2007), paragraphs [0027], [0039] to [0042], [0063]; fig. 14, 15 & US 2007/0180294 A1	1-3, 6-8

 Further documents are listed in the continuation of Box C. See patent family annex.

\* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&amp;" document member of the same patent family

Date of the actual completion of the international search  
21 June, 2012 (21.06.12)Date of mailing of the international search report  
03 July, 2012 (03.07.12)Name and mailing address of the ISA/  
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.

A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int.Cl. G06F13/10(2006.01)i, G06F3/06(2006.01)i

B. 調査を行った分野

調査を行った最小限資料 (国際特許分類 (IPC))

Int.Cl. G06F13/10, G06F3/06

最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公報	1922-1996年
日本国公開実用新案公報	1971-2012年
日本国実用新案登録公報	1996-2012年
日本国登録実用新案公報	1994-2012年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
X	JP 2010-015516 A (株式会社東芝) 2010.01.21, 段落【0006】-【0016】、【0045】、 【0059】-【0070】 & US 2010/0005228 A1	1-8
Y	JP 2009-037304 A (株式会社日立製作所) 2009.02.19, 段落【0064】-【0066】 & US 2009/0037656 A1	1-3, 6-8
Y	JP 2007-206993 A (富士通株式会社) 2007.08.16, 段落【0027】、【0039】-【0042】、【0063】、 図14, 15 & US 2007/0180294 A1	1-3, 6-8

☐ C欄の続きにも文献が列挙されている。

☐ パテントファミリーに関する別紙を参照。

\* 引用文献のカテゴリー

「A」特に関連のある文献ではなく、一般的技術水準を示すもの  
 「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの  
 「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)  
 「O」口頭による開示、使用、展示等に言及する文献  
 「P」国際出願日前で、かつ優先権の主張の基礎となる出願

の日の後に公表された文献  
 「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの  
 「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの  
 「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの  
 「&」同一パテントファミリー文献

国際調査を完了した日

21.06.2012

国際調査報告の発送日

03.07.2012

国際調査機関の名称及びあて先

日本国特許庁 (ISA/J P)  
 郵便番号100-8915  
 東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)

横山 佳弘

5 T 3565

電話番号 03-3581-1101 内線 3568