

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2023-507702

(P2023-507702A)

(43)公表日 令和5年2月27日(2023.2.27)

(51)国際特許分類		F I			
G 0 6 F	3/06 (2006.01)	G 0 6 F	3/06	3 0 1 S	
G 0 6 F	9/50 (2006.01)	G 0 6 F	9/50	1 2 0 Z	

審査請求 未請求 予備審査請求 未請求 (全56頁)

(21)出願番号	特願2022-514540(P2022-514540)	(71)出願人	591003943 インテル・コーポレーション アメリカ合衆国 9 5 0 5 4 カリフォルニア州・サンタクララ・ミッション カレッジ ブレーバード・2 2 0 0
(86)(22)出願日	令和2年9月18日(2020.9.18)	(74)代理人	110000877 弁理士法人 R Y U K A 国際特許事務所
(85)翻訳文提出日	令和4年3月30日(2022.3.30)	(72)発明者	バンヤイ、クリストファー ジェイ . アメリカ合衆国 9 5 0 5 4 カリフォルニア州・サンタクララ・ミッション カレッジ ブレーバード・2 2 0 0
(86)国際出願番号	PCT/US2020/051560	(72)発明者	コーエン、デイビッド イー . アメリカ合衆国 9 5 0 5 4 カリフォルニア州・サンタクララ・ミッション カ
(87)国際公開番号	WO2021/133443		
(87)国際公開日	令和3年7月1日(2021.7.1)		
(31)優先権主張番号	16/729,075		
(32)優先日	令和1年12月27日(2019.12.27)		
(33)優先権主張国・地域又は機関	米国(US)		
(81)指定国・地域	AP(BW,GH,GM,KE,LR,LS,MW,MZ,NA,RW,SD,SL,ST,SZ,TZ,UG,ZM,ZW),EA(AM,AZ,BY,KG,KZ,RU,TJ,TM),EP(AL,A T,BE,BG,CH,CY,CZ,DE,DK,EE,ES,FI,FR,GB,GR,HR,HU,IE,IS,IT,LT,LU,LV,MC, 最終頁に続く		最終頁に続く

(54)【発明の名称】 クラウド・ネイティブなワークロードのためのデータ管理プラットフォームにおけるストレージ管理

(57)【要約】

計算サーバおよびストレージ・サーバを含むデータ管理プラットフォームが提供される。ストレージ・サーバは、ストレージ・サーバに通信可能に結合された複数のストレージ・デバイスを管理する。計算サーバおよびストレージ・サーバは、ネットワークを介して通信可能に結合されている。ストレージ・サーバによって管理される複数のストレージ・デバイスは、複数のストレージ・デバイスのストレージ容量を計算サーバから独立してスケールアップすることができるように、計算サーバから分離されている。

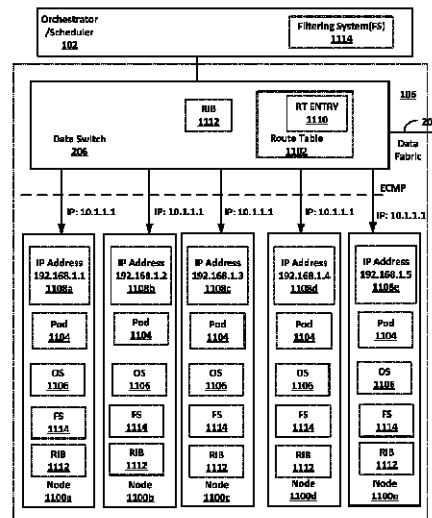


FIG. 11

【特許請求の範囲】**【請求項 1】**

計算サーバと、

ストレージ・サーバであって、前記ストレージ・サーバに通信可能に結合された複数のストレージ・デバイスを管理する、ストレージ・サーバと、を備え、前記計算サーバおよび前記ストレージ・サーバがネットワークを介して通信可能に結合され、前記ストレージ・サーバによって管理される前記複数のストレージ・デバイスが、前記複数のストレージ・デバイスのストレージ容量を前記計算サーバから独立してスケールリングすることができるように、前記計算サーバから分離されている、
装置。

10

【請求項 2】

前記ストレージ・サーバが、前記ネットワークに通信可能に結合されたネットワーク・インターフェース・コントローラと、
複数のコアおよびラスト・レベル・キャッシュ・メモリを含むシステム・オン・チップであって、前記複数のコアが前記ラスト・レベル・キャッシュ・メモリに通信可能に結合され、前記ラスト・レベル・キャッシュ・メモリが複数のキャッシュ・ウェイを含み、前記複数のキャッシュ・ウェイの一部が、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、システム・オン・チップと、
をさらに備える、請求項 1 に記載の装置。

20

【請求項 3】

前記複数のキャッシュ・ウェイの前記一部が初期化中に割り当てられる、請求項 2 に記載の装置。

【請求項 4】

前記論理ボリュームが前記計算サーバによる使用のためのデータを記憶する、請求項 3 に記載の装置。

【請求項 5】

前記論理ボリュームと前記ラスト・レベル・キャッシュ・メモリ内の前記複数のキャッシュ・ウェイとの間で転送される前記データが、前記ネットワーク・インターフェース・コントローラに通信可能に結合されたネットワークを介して前記ストレージ・サーバと前記計算サーバとの間で転送される、請求項 4 に記載の装置。

30

【請求項 6】

前記複数のコアのうちの少なくとも 1 つが、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、請求項 2 に記載の装置。

【請求項 7】

前記システム・オン・チップに結合された外部メモリであって、転送されるデータを記憶するために、ラスト・レベル・キャッシュの前記一部における前記複数のキャッシュ・ウェイのすべてが、前記論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられると、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間で転送されるデータを一時的に記憶する、外部メモリをさらに備える、請求項 2 に記載の装置。

40

【請求項 8】

ストレージ・サーバによって、前記ストレージ・サーバに通信可能に結合された複数のストレージ・デバイスを管理するステップと、
ネットワークを介して、前記ストレージ・サーバを計算サーバと通信可能に結合するステップであって、前記ストレージ・サーバによって管理される前記複数のストレージ・デバイスが、前記複数のストレージ・デバイスのストレージ容量を前記計算サーバから独立

50

してスケーリングすることができるように、前記計算サーバから分離されている、ステップと、

を含む方法。

【請求項 9】

前記ストレージ・サーバが、前記ネットワークに通信可能に結合されたネットワーク・インターフェース・コントローラと、

複数のコアおよびラスト・レベル・キャッシュ・メモリを含むシステム・オン・チップであって、前記複数のコアが前記ラスト・レベル・キャッシュ・メモリに通信可能に結合され、前記ラスト・レベル・キャッシュ・メモリが複数のキャッシュ・ウェイを含み、前記複数のキャッシュ・ウェイの一部が、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、システム・オン・チップと、

10

を有する、請求項 8 に記載の方法。

【請求項 10】

前記複数のキャッシュ・ウェイの前記一部が初期化中に割り当てられる、請求項 9 に記載の方法。

【請求項 11】

前記論理ボリュームが前記計算サーバによる使用のためのデータを記憶する、請求項 10 に記載の方法。

20

【請求項 12】

前記論理ボリュームと前記ラスト・レベル・キャッシュ・メモリ内の前記複数のキャッシュ・ウェイとの間で転送される前記データが、前記ネットワーク・インターフェース・コントローラに通信可能に結合されたネットワークを介して前記ストレージ・サーバと前記計算サーバとの間で転送される、請求項 11 に記載の方法。

【請求項 13】

前記複数のコアのうちの少なくとも 1 つが、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、請求項 9 に記載の方法。

30

【請求項 14】

前記システム・オン・チップに結合された外部メモリであって、転送されるデータを記憶するために、ラスト・レベル・キャッシュの前記一部における前記複数のキャッシュ・ウェイのすべてが、前記論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられると、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間で転送されるデータを一時的に記憶する、外部メモリをさらに備える、請求項 9 に記載の方法。

【請求項 15】

請求項 8 から 14 のいずれか一項に記載の方法を実行するための手段を備える装置。

【請求項 16】

システムによって実行されることに応答して、前記システムに、
ストレージ・サーバによって、前記ストレージ・サーバに通信可能に結合された複数のストレージ・デバイスを管理させ、

40

ネットワークを介して、前記ストレージ・サーバを計算サーバと通信可能に結合させ、
前記ストレージ・サーバによって管理された前記複数のストレージ・デバイスが、前記複数のストレージ・デバイスのストレージ容量を前記計算サーバから独立してスケーリングすることができるように、前記計算サーバから分離されている、
複数の命令を含むコンピュータプログラム。

【請求項 17】

前記ストレージ・サーバが、前記ネットワークに通信可能に結合されたネットワーク・

50

インターフェース・コントローラと、

複数のコアおよびラスト・レベル・キャッシュ・メモリを含むシステム・オン・チップであって、前記複数のコアが前記ラスト・レベル・キャッシュ・メモリに通信可能に結合され、前記ラスト・レベル・キャッシュ・メモリが複数のキャッシュ・ウェイを含み、前記複数のキャッシュ・ウェイの一部が、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、システム・オン・チップと、

を有する、請求項 16 に記載のコンピュータプログラム。

【請求項 18】

前記複数のキャッシュ・ウェイの前記一部が初期化中に割り当てられる、請求項 17 に記載のコンピュータプログラム。

【請求項 19】

前記論理ボリュームが前記計算サーバによる使用のためのデータを記憶する、請求項 18 に記載のコンピュータプログラム。

【請求項 20】

前記論理ボリュームと前記ラスト・レベル・キャッシュ・メモリ内の前記複数のキャッシュ・ウェイとの間で転送される前記データが、前記ネットワーク・インターフェース・コントローラに通信可能に結合されたネットワークを介して前記ストレージ・サーバと前記計算サーバとの間で転送される、請求項 19 に記載のコンピュータプログラム。

【請求項 21】

前記複数のコアのうちの少なくとも 1 つが、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、請求項 17 に記載のコンピュータプログラム。

【請求項 22】

前記システム・オン・チップに結合された外部メモリであって、転送されるデータを記憶するために、ラスト・レベル・キャッシュの前記一部における前記複数のキャッシュ・ウェイのすべてが、前記論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられると、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間で転送されるデータを一時的に記憶する、外部メモリをさらに備える、請求項 17 に記載のコンピュータプログラム。

【請求項 23】

請求項 16 から 22 のいずれか一項に記載のコンピュータプログラムを格納する、少なくとも 1 つのコンピュータ可読記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

優先権の主張

本出願は、米国特許法第 365 条 (c) の下で、2019 年 12 月 27 日に出願された「STORAGE MANAGEMENT IN A DATA MANAGEMENT PLATFORM FOR CLOUD-NATIVE WORKLOADS」と題する米国特許出願第 16 / 729075 号の優先権を主張し、その全体が本明細書に組み込まれる。

【背景技術】

【0002】

クラウド・コンピューティングは、インターネットを介してサーバ、ストレージ、データベース、および広範なアプリケーション・サービスへのアクセスを提供する。クラウド・サービス・プロバイダは、インターネットを介して会社または個人がアクセスすることができる 1 つまたは複数のデータセンタ内のサーバにホストされる、ネットワーク・サー

10

20

30

40

50

ビスおよびビジネス・アプリケーションなどのクラウド・サービスを提供する。ハイパースケールのクラウド・サービス・プロバイダは、典型的には数十万のサーバを有する。ハイパースケールのクラウド内の各サーバは、ユーザデータ、例えば、ビジネスインテリジェンス、データマイニング、分析、ソーシャルメディア、およびマイクロサービスのためのユーザデータを記憶するためのストレージ・デバイスを含む。クラウド・サービス・プロバイダは、クラウド・サービスを利用する企業および個人（テナントとも呼ばれる）から収益を生み出す。例えば、テナントは、テナントに割り当てられたストレージの量に基づいて、クラウド・サービス・プロバイダに月額料金を支払うことでデータを記憶することができる。

【0003】

10

今日の大多数の企業のデータセンタは、規模および性能においてペタバイトのデータを効果的に管理および処理する能力を有していない。人工知能（AI）推論および分析などのデータ集約型アプリケーションおよびツールは、より安全で、より高速で、拡張性の高い方法で移動、記憶、処理する必要のある爆発的な量のデータおよび遠隔測定データを生成し、消費する。ハイパースケール・データセンタでは、典型的には、データセンタに追加のサーバを追加することによってこれを実行する。しかしながら、データセンタで実行されているワークロードに応じて、これらのサーバの1つのタイプのコンポーネントが過剰に利用されることがある一方で、別のタイプのコンポーネントが十分に活用されていない場合があり、これは、顧客およびサービス・プロバイダが投資の使用を最適化していないことを意味する。

20

【発明の概要】

【0004】

特許請求される主題の実施形態の特徴は、以下の詳細な説明が進むにつれて、および同様の数字が同様の部分を示す図面を参照することによって、明らかになるであろう。

【図面の簡単な説明】

【0005】

【図1】データ管理プラットフォーム（DMP）の一実施形態の概念図である。

【図2】図1に示すデータ管理プラットフォームの物理クラスタにおける実施形態のブロック図である。

【図3】図2に示すデータ管理プラットフォームにおける計算サーバのうちの1つの実施形態のブロック図である。

30

【図4】図2に示すデータ管理プラットフォームにおけるアクセラレータ・サーバのうちの1つの実施形態のブロック図である。

【図5】図2に示すデータ管理プラットフォームにおける計算サーバからソリッド・ステート・ドライブへのアクセスを示す論理図である。

【図6】ストレージ・サービスのリソース制御を実行する図4に示すアクセラレータ・サーバのブロック図である。

【図7】アクセラレータ・サーバ内のネットワーク・インターフェース・コントローラとソリッド・ステート・ドライブとの間のデータ転送を示す図である。

【図8】L3キャッシュおよびネットワーク・インターフェース・コントローラを介してソリッド・ステート・ドライブからデータプレーンにデータを移動する方法を示す流れ図である。

40

【図9】ソリッド・ステート・ドライブおよびネットワーク・インターフェース・コントローラによって共有されるラスト・レベル・キャッシュのキャッシュ・ウェイのN個のセットを分離するようにラスト・レベル・キャッシュを構成するための実施形態を示す図である。

【図10】ソリッド・ステート・ドライブおよびネットワーク・インターフェース・コントローラによって共有されるラスト・レベル・キャッシュのキャッシュ・ウェイのセットを構成する方法の実施形態を示す流れ図である。

【図11】健全なシステムで通常動作している場合の、図2に示す物理クラスタにおける

50

データ管理プラットフォーム内のラックの実施形態のブロック図である。

【図 1 2】故障したシステムにおける縮退動作の場合の、図 2 に示す物理クラスタにおけるデータ管理プラットフォーム内のラックの実施形態のブロック図である。

【図 1 3】ノード状態/障害を検出するためにデータ管理プラットフォームによって使用される、ノード内のポッドにあるコンテナ内のメトリクス・エクスポートを示すブロック図である。

【図 1 4】物理クラスタにおけるデータ管理プラットフォームのハードウェア障害を管理する方法を示す流れ図である。

【図 1 5】性能しきい値を監視および管理してノード状態および障害を検出するためにデータ管理プラットフォームによって使用されるノード内のハードウェア・イベントおよび測定値を示すブロック図である。

【図 1 6】ストレージ・ノードの性能を監視するためにラック内のストレージ・ノードで実施される方法を示す流れ図である。

【図 1 7】計算ノードの実施形態のブロック図である。

【図 1 8】計算ノードの別の実施形態のブロック図である。

【図 1 9】論理リソースを自動的に追加および除去するリソース・マネージャを含む、図 2 に示す物理クラスタにおけるデータ管理プラットフォーム内のラックの実施形態のブロック図である。

【図 2 0】図 1 9 に示すデータ管理プラットフォーム内のラックにおける圧迫の検出に応じて論理リソースを自動的に追加または除去する方法の流れ図である。

【図 2 1】ストレージ自己修復メカニズムを含む、データ管理プラットフォームにおける物理クラスタの一実施形態のブロック図である。

【図 2 2】図 2 1 に示すクラスタにおけるワークロードのマッピングの一実施形態を示す図である。

【発明を実施するための形態】

【0006】

以下の詳細な説明は、特許請求される主題の例示的な実施形態を参照して進められるが、その多くの代替形態、修正形態、および変形形態が当業者には明らかであろう。したがって、特許請求される主題は、広くとらえられ、添付の特許請求の範囲に記載されるように定義されることが意図されている。

【0007】

データ管理プラットフォームは、アクセラレータ・サーバおよび計算サーバを含む。アクセラレータ・サーバによって管理されるストレージ・デバイスは、ストレージ容量が計算とは無関係にスケールできるように計算サーバから分離されている。

【0008】

本発明の様々な実施形態および態様は、以下で論じる詳細を参照して説明され、添付の図面は様々な実施形態を示す。以下の説明および図面は、本発明を例示するものであり、本発明を限定するものとして解釈されるべきではない。本発明の様々な実施形態の完全な理解を提供するために、多数の具体的な詳細が説明される。しかしながら、ある場合には、本発明の実施形態の簡潔な議論を提供するために、よく知られたまたは従来の詳細については説明されない。

【0009】

本明細書における「一実施形態」または「実施形態」への言及は、実施形態と併せて説明される特定の特徵、構造、または特性が、本発明の少なくとも1つの実施形態に含まれることを意味する。本明細書の様々な箇所における「一実施形態では」という語句の出現は、必ずしもすべてが同じ実施形態を指すわけではない。

【0010】

図 1 は、データ管理プラットフォーム (DMP) 100 の一実施形態の概念図である。図 1 に示す実施形態では、データ管理プラットフォーム 100 は、ルーティング・インターコネクト 110 を介して相互接続されたラック 106 を有するラックを中心とした物理

10

20

30

40

50

クラスタである。ルーティング・インターコネクト 110 は、多段 Clost ポロジに配列されたイーサネット（登録商標）ファブリック、または任意の他の Open Systems Interconnect (OSI) レイヤ 3 ルーティング・インターコネクトとすることができる。

【0011】

データセンタ内のラック 106 は、サーバ、ネットワーキング・デバイス、ケーブル、および他のデータセンタ・コンピューティング機器を収納するように設計された一種の物理的な鋼鉄および電子フレームワークである。各ラック 106 は、ルーティング・インターコネクト 110 に接続し、1つまたは複数の計算サーバ 112、アクセラレータ・サーバ 114、ユーティリティ・サーバ 118、およびインフラストラクチャ・サーバ 116 を含むことができる。サーバは、ノードと呼ぶこともできる。

10

【0012】

ユーティリティ・サーバ 118 は、物理クラスタを初期化するために使用される。初期化の間、ユーティリティ・サーバ 118 は、オーケストレーションおよびスケジューリング機能を実行する。一実施形態では、オーケストレータ/スケジューラ 102 のための機能を実行するために Kubernetes (K8) が使用される。Kubernetes は、アプリケーション・デプロイメント、スケーリング、および管理を自動化するためのオープン・ソースのコンテナ・オーケストレーション・システムである。Kubernetes 制御プレーンは、インフラストラクチャ・サーバ 116 上でホストされている。Kubernetes Host Agent は、すべての計算サーバ 112 およびアクセラレータ・サーバ 114 上で実行される。

20

【0013】

アプリケーション・デプロイメントは、仮想マシンを使用することによって自動化することもできる。オーケストレータ/スケジューラ 102 の他の例として、OpenShift (Docker および Kubernetes 上に構築される Red Hat 社の PaaS (platform as a service)) および Pivotal Container Service (PKS) が挙げられる。

【0014】

制御プレーンマネージャ 104 は、仮想マシンなどのインフラストラクチャリソースを作成、管理、および更新するための機能を実行することができる。制御プレーンマネージャ 104 は、物理マシンおよびネットワークスイッチを初期化することもできる。制御プレーンマネージャ 104 の例としては、Fleet、Red Hat Satellite、Teraform および Metal As A Service (MaaS) が挙げられる。

30

【0015】

計算サーバ 112、アクセラレータ・サーバ 114、ユーティリティ・サーバ 118、およびインフラストラクチャ・サーバ 116 のそれぞれは、Baseboard Management Controller (BMC) 120 を含む。BMC 120 は、計算サーバ 112、アクセラレータ・サーバ 114、ユーティリティ・サーバ 118、およびインフラストラクチャ・サーバ 116 の物理的状态を監視し、Management API 108 を介して動作を監視および制御するサービスを提供する専用のサービスプロセッサである。Management API 108 の例としては、IPMI (Intelligent Platform Management Interface)、Redfish (登録商標) (Distributed Management Task Force (DMTF) 標準) および Dell (登録商標) Open Manage Enterprise (OME) が挙げられる。

40

【0016】

図 2 は、物理クラスタ 200 における図 1 に示すデータ管理プラットフォーム (DMP) 100 の実施形態のブロック図である。物理クラスタ 200 は、N 個のラック 106、106-1、... 106-N を有する。一実施形態では、N は 16 である。各ラック 10

50

6 は、計算サーバ 1 1 2 と、少なくとも 1 つのアクセラレータ・サーバ 1 1 4 と、を含む。各計算サーバ 1 1 2 およびアクセラレータ・サーバ 1 1 4 は、データスイッチ 2 0 6 および管理スイッチ 2 0 8 に通信可能に結合されている。各ラック 1 0 6 のデータスイッチ 2 0 6 は、同じラック 1 0 6 内の計算サーバ 1 1 2 およびアクセラレータ・サーバ 1 1 4 と、他のラック 1 0 6 内の計算サーバ 1 1 2 およびアクセラレータ・サーバ 1 1 4 と、複数のラック 1 0 6 によって共有されるインフラストラクチャ・サーバ 1 1 6 との間のデータプレーン 2 0 4 (データ・ファブリックとも呼ばれる) を提供する。各ラック 1 0 6 の管理スイッチ 2 0 8 は、ラック 1 0 6 と複数のラック 1 0 6 によって共有されるユーティリティ・サーバ 1 1 8 との間に制御プレーン 2 0 2 (管理ネットワークとも呼ばれる) を提供する。

10

【 0 0 1 7 】

図 3 は、図 2 に示す物理クラスタ 2 0 0 における計算サーバ 1 1 2 のうちの 1 つの実施形態のブロック図である。図示する実施形態では、計算サーバ 1 1 2 は、システム・オン・チップ 3 0 6 と、ネットワーク・インターフェース・コントローラ 3 0 2 と、計算サーバ制御ロジック 3 0 4 と、を含む。ネットワーク・インターフェース・コントローラ 3 0 2 は、図 2 に示すデータプレーン 2 0 4 に通信可能に結合されている。システム・オン・チップ 3 0 6 内の組み込みネットワーク・インターフェース・コントローラ 3 0 8 は、図 2 に示す制御プレーン 2 0 2 に通信可能に結合されている。

【 0 0 1 8 】

図 4 は、図 2 に示す物理クラスタ 2 0 0 におけるアクセラレータ・サーバ 1 1 4 のうちの 1 つの実施形態のブロック図である。図示する実施形態では、アクセラレータ・サーバ 1 1 4 は、ストレージ処理タスクを実行し、ストレージ・サーバ 4 0 0 と呼ぶことができる。

20

【 0 0 1 9 】

ストレージ・サーバ 4 0 0 は、システム・オン・チップ 3 0 6 と、ネットワーク・インターフェース・コントローラ 3 0 2 と、1 つまたは複数のソリッド・ステート・ドライブ 4 0 4 に通信可能に結合されたストレージ・サーバ制御ロジック 4 0 2 と、を含む。一実施形態では、ストレージ・サーバ制御ロジック 4 0 2 は、P C I (P e r i p h e r a l C o m p o n e n t I n t e r c o n n e c t) E x p r e s s (P C I e) プロトコルを用いて、ソリッド・ステート・ドライブ 4 0 4 およびネットワーク・インターフェース・コントローラ 3 0 2 に通信可能に結合されている。システム・オン・チップ 3 0 6 内の組み込みネットワーク・インターフェース・コントローラ 3 0 8 は、図 2 に示す制御プレーン 2 0 2 に通信可能に結合されている。

30

【 0 0 2 0 】

ストレージ・サーバ制御ロジック 4 0 2 は、システム・オン・チップ 3 0 6 によってオフロードされたストレージ処理タスクを実行し、計算とストレージを独立したスケラブルなリソースに分解することを可能にする。

【 0 0 2 1 】

図 5 は、図 2 に示す物理クラスタ 2 0 0 における計算サーバ 1 1 2 からソリッド・ステート・ドライブ 4 0 4 へのアクセスを示す論理図である。オペレーティング・システムのカーネル空間 5 0 2 の仮想ルーティング機能 5 0 8 は、ユーザ空間 5 0 0 のリレーショナル・データベース管理システム 5 0 6 に、データプレーン 2 0 4 を介してストレージ・サーバ 4 0 0 (図 4) のソリッド・ステート・ドライブ 4 0 4 に記憶されたデータへのアクセスを提供する。仮想ルーティング機能 5 0 8 は、F I B (F o r w a r d i n g I n f o r m a t i o n B a s e) 5 1 0 と、ルートおよびポリシーを記憶するフローテーブル 5 1 2 と、を含む。

40

【 0 0 2 2 】

ルータ 5 0 4 は、仮想マネージャおよびコンテナに安全なネットワーク接続性を提供する。ルータ 5 0 4 の一例は、C a l i c o である。C a l i c o は、コンテナおよび仮想マシンワークロードのための安全なネットワーク接続性を提供する。C a l i c o は、ル

50

ーティングテーブルを構築するために、OSI (Open System Interconnection) モデルのレイヤ3 (ネットワークレイヤ) およびBGP (Border Gateway Protocol) を使用する。Calicoは、フラットなレイヤ3ネットワークを作成し、すべてのラック106に完全にルーティング可能なインターネット・プロトコル (IP) アドレスを割り当てる。ワークロードは、ベアメタル性能のためにIPカプセル化またはネットワークアドレス変換なしで通信することができる。Calicoは、Felix (ノードごとのドメインデーモン) を使用して、ルートを構成し、ネットワークポリシーを実施する。

【0023】

ラスト・レベル・キャッシュ (LLC) およびメインメモリ帯域幅などの共有リソースは、データ管理プラットフォーム (DMP) におけるワークロード性能に大きな影響を及ぼす。これらのリソースをより厳密に監視および管理することにより、ますます厳しくなるパフォーマンスSLA (performance service-level agreement) を含む、より厳しいワークロード要求を満たすデプロイメントが可能になる。

【0024】

図6は、ストレージ・サービスのリソース制御を行う図4に示すストレージ・サーバ400のブロック図である。

【0025】

ストレージ・サーバ400は、プロセッサと、メモリと、入出力 (I/O) 制御ロジックを1つのSoCパッケージに統合したシステム・オン・チップ (SOCまたはSoC) 604と、を含む。SoC604は、少なくとも1つの中央処理装置 (CPU) モジュール608およびメモリコントローラ614を含む。他の実施形態では、メモリコントローラ614は、SoC604の外部にあってもよい。CPUモジュール608は、レベル1 (L1) およびレベル2 (L2) キャッシュ620を含む少なくとも1つのプロセッサコア602と、CPUモジュール608内の他のプロセッサコア602と共有されるレベル3 (L3) キャッシュ606と、を含む。

【0026】

図示されていないが、プロセッサコア602のそれぞれは、内部に、実行ユニット、プリフェッチバッファ、命令キュー、分岐アドレス計算ユニット、命令デコーダ、浮動小数点ユニット、リタイアメントユニットなどを含むことができる。CPUモジュール608は、一実施形態によれば、Intel (登録商標) Corporationによって提供されるものなどのシングルコアまたはマルチコア汎用プロセッサに対応することができる。

【0027】

I/Oサブシステム612内には、1つまたは複数のI/Oインターフェース616が存在し、プロセッサコア602内で利用されるホスト通信プロトコルを特定のI/Oデバイスと互換性のあるプロトコルに変換する。変換のためにI/Oインターフェースを利用することができるプロトコルの一部には、PCI (Peripheral Component Interconnect) Express (PCIe)、USB (Universal Serial Bus)、SATA (Serial Advanced Technology Attachment)、およびIEEE (Institute of Electrical and Electronics Engineers) 1594「Firewire (登録商標)」が含まれる。

【0028】

I/Oインターフェース616は、メモリ630および/またはL3キャッシュ606を介して、1つまたは複数のソリッド・ステート・ドライブ404およびネットワーク・インターフェース・コントローラ302と通信することができる。ソリッド・ステート・ドライブ404は、SAS (Serial Attached SCSI (Small Computer System Interface))、PCIe (Periphe

10

20

30

40

50

ral Component Interconnect Express)、NVMe (NVM Express) over PCIe (Peripheral Component Interconnect Express)、およびSATA (Serial ATA (Advanced Technology Attachment))を含むがこれらに限定されない様々なプロトコルのうちの1つまたは複数を使用して、1つまたは複数のバスを介して互いに通信可能におよび/または物理的に結合することができる。他の実施形態では、他のストレージ・デバイス、例えば、ハード・ディスク・ドライブ(HDD)などの他のストレージ・デバイスをソリッド・ステート・ドライブ404の代わりに使用することができ、ハード・ディスク・ドライブおよび/またはソリッド・ステート・ドライブは、RAID (Redundant Array of Independent Disks)として構成することができる。

【0029】

NVMe (Non-Volatile Memory Express)規格は、ホストソフトウェアが、高速シリアルコンピュータ拡張バスであるPCIe (Peripheral Component Interconnect Express)を介して不揮発性メモリサブシステム(例えば、ソリッド・ステート・ドライブ404)と通信するためのレジスタレベルのインターフェースを定義する。NVM Express規格は、www.nvmexpress.orgで入手可能である。PCIe規格は、www.pcisig.comで入手可能である。

【0030】

一実施形態では、メモリ630は揮発性メモリであり、メモリコントローラ614は揮発性メモリコントローラである。揮発性メモリは、デバイスへの電力が遮断された場合に、その状態(したがって、内部に記憶されたデータ)が不確定になるメモリである。動的揮発性メモリは、状態を維持するためにデバイスに記憶されたデータをリフレッシュする必要がある。動的揮発性メモリの一例には、DRAM (ダイナミック・ランダム・アクセス・メモリ)、またはSDRAM (同期式DRAM)などの何らかの変形が含まれる。本明細書に記載されるメモリサブシステムは、DDR3 (ダブル・データ・レートversion 3、2007年6月27日のJEDEC (Joint Electronic Device Engineering Council)によるオリジナルリリース)などのいくつかのメモリ技術と互換性があり得る。DDR4 (DDRバージョン4、JEDECによって2012年9月に公開された初期仕様)、DDR4E (DDR version 4)、LPDDR3 (低電力DDR version 3、JESD209-3B、JEDECによって2013年8月に公開された)、LPDDR4 (LPDDR version 4、JESD209-4、2014年8月にJEDECによって最初に公開された)、WIO2 (Wide Input/Output version 2、JESD229-2、2014年10月にJEDECによって最初に公開された)、HBM (広帯域メモリ、JESD325、2013年10月にJEDECによって最初に公開された)、DDR5 (DDR version 5、JEDECによって現在審議中)、LPDDR5、JEDECによって現在審議中のHBM2 (HBM version 2)など、またはその他、またはメモリ技術の組合せ、およびこのような仕様の派生版もしくは拡張版に基づく技術である。JEDEC規格は、www.jedec.orgで入手可能である。

【0031】

別の実施形態では、メモリ630は不揮発性メモリ(NVM)であり、メモリコントローラ614は不揮発性メモリコントローラである。不揮発性メモリデバイスは、デバイスへの電力が遮断されてもその状態が確定しているメモリである。不揮発性メモリデバイスとしては、シングルレベルもしくはマルチレベル相変化メモリ(PCM)またはスイッチ付き相変化メモリ(PCMS)などの、バイトアドレス指定可能なライト・イン・プレイス(write-in-place)3次元クロス・ポイント・メモリ・デバイスあるいは他のバイトアドレス指定可能なライト・イン・プレイスNVMデバイス(永続メモリとも呼ばれる)、カルコゲナイド相変化材料(例えば、カルコゲナイドガラス)を使用する

NVMデバイス、金属酸化物ベース、酸素空孔ベースおよび導電性ブリッジ・ランダム・アクセス・メモリ(CB-RAM)を含む抵抗メモリ、ナノワイヤメモリ、強誘電体ランダム・アクセス・メモリ(FeRAM、FRAM(登録商標))、メモリスタ技術を組み込んだ磁気抵抗ランダム・アクセス・メモリ(MRAM)、スピン・トランスファ・トルク(STT)MRAM、スピントロニクス磁気接合メモリベースのデバイス、磁気トンネル接合(MTJ)ベースのデバイス、DW(磁壁)およびSOT(スピン軌道トランスファ)ベースのデバイス、サイリスタベースのメモリデバイス、または上記のいずれかの組合せ、または他のメモリを挙げることができる。

【0032】

さらに別の実施形態では、メモリ630は、1つまたは複数のメモリモジュールに含まれ得るバイトアドレス指定可能なライト・イン・プレイスNVMデバイスおよび揮発性メモリデバイスの両方を含む。

10

【0033】

ワークロード性能に大きな影響を及ぼす共有リソースには、プロセッサキャッシュおよびメモリ帯域幅リソースが含まれ、これらは、アプリケーション性能およびランタイム決定論に大きな影響を及ぼす可能性がある。これらのリソースをより綿密に監視および管理することにより、ネットワーク機能仮想化(NFV)などの新しいワークロードをサポートするための、ますます厳しくなるパフォーマンスSLA(performance service-level agreement)を含む、より厳しいワークロード要求を満たすデプロイメントが可能になる。

20

【0034】

図7は、ネットワーク・インターフェース・コントローラ302とストレージ・サーバ400のソリッド・ステート・ドライブ404との間のデータ転送を示す。I/Oアダプタ616は、ソリッド・ステート・ドライブ404からの第1のPCIeインターフェース708を介した通信を管理するための第1のPCIeアダプタ702と、ネットワーク・インターフェース・コントローラ302への第1のPCIeインターフェース708を介した通信を管理するための第2のPCIeインターフェース704と、を含む。ネットワーク・インターフェース・コントローラ302は、リモート・ダイレクト・メモリ・アクセス(RDMA)、例えば、L3キャッシュ606および/またはメモリ630から、データプレーン204に通信可能に結合された計算サーバ112または別のアクセラレータ・サーバ114のメモリへのダイレクト・メモリ・アクセスを使用して、データを交換することができる。

30

【0035】

ソリッド・ステート・ドライブ404とネットワーク・インターフェース・コントローラ302は、L3キャッシュ606および/またはメモリ630を介してデータを交換する。L3キャッシュ606は、ラスト・レベル・キャッシュ(LLC)と呼ぶこともできる。レベル3(L3)キャッシュ606は、CPUモジュール608内の他のプロセッサコア602と共有されることに加えて、第1のPCIeインターフェース708および第2のPCIeインターフェース710とも共有される。

【0036】

複数のエージェント(プロセッサコア602、第1のPCIeインターフェース708、および第2のPCIeインターフェース710)がすべて同じL3キャッシュ606に競合的にアクセスすると、L3キャッシュ606におけるキャッシュミス、L3キャッシュ606からメモリ630へのキャッシュ・エビクション、およびエージェントのトランザクションにおけるレイテンシの変動が大きくなる可能性がある。ネットワーク・インターフェース・コントローラ302の帯域幅がソリッド・ステート・ドライブ404の帯域幅とよく一致し、L3キャッシュ606が十分なサイズのものである場合、ソリッド・ステート・ドライブ404とネットワーク・インターフェース・コントローラ302との間のデータ転送の大部分は、パス714を介したメモリ630への退避(「流出(spill)」)なしにL3キャッシュ606を介して行われる。

40

50

【 0 0 3 7 】

L 3 キャッシュ 6 0 6 からメモリ 6 3 0 への流出を最小限に抑えるために、L 3 キャッシュ 6 0 6 の設定可能な部分（キャッシュ・ウェイのサブセット）は、ソリッド・ステート・ドライブ 4 0 4 とネットワーク・インターフェース・コントローラ 3 0 2 との間で転送されるデータを専ら記憶する。図 6 に示す実施形態では、N 個のプロセッサコア 6 0 2 - 1、. . . 6 0 2 - N が存在する。L 3 キャッシュ 6 0 6 のキャッシュ・ウェイ 7 1 2 の第 1 のサブセットは、プロセッサコア 6 0 2 - 1 およびプロセッサコア 6 0 2 - 2 の両方に専用である。L 3 キャッシュ 6 0 6 のキャッシュ・ウェイ 7 0 6 の第 2 のサブセットは、プロセッサコア 6 0 2 - 3、ソリッド・ステート・ドライブ 4 0 4、およびネットワーク・インターフェース・コントローラ 3 0 2 に専用である。

10

【 0 0 3 8 】

一実施形態では、Intel（登録商標）キャッシュ・アロケーション・テクノロジー（CAT）を使用して、L 3 キャッシュ 6 0 6 のキャッシュ・ウェイのサブセットを特定のプロセッサコア 6 0 2 - 1、. . . 6 0 2 - N および / または I / O メモリ空間（PCIe）に専用にすることができ、どのエージェントが L 3 キャッシュ 6 0 6 のキャッシュ・ウェイ（または部分）の特定のサブセットを共有 / 競合するかを制御することができる。他のすべてのエージェントは、ソリッド・ステート・ドライブ 4 0 4 とネットワーク・インターフェース・コントローラ 3 0 2 との間で転送されるデータを専ら記憶する L 3 キャッシュ 6 0 6 の第 2 のセットのキャッシュ・ウェイ 7 0 6 を使用することから除外される。

【 0 0 3 9 】

L 3 キャッシュ 6 0 6 のキャッシュ・ウェイ 7 0 6 の設定可能な第 2 のサブセットの使用により、ワークロードの変動が低減され、他の同じ場所に配置されたワークロードに対するストレージ・サービスのより正確で予測可能なリソース割り当てが提供され、他のワークロードと同じ場所に配置されたストレージ・サービスのより正確なサービスレベルを予測することが可能になる。3 つのサービス品質メカニズム（キャッシュ、コア / 入力 / 出力メモリ、および論理ボリューム帯域幅スロットリング）を組み合わせ、調整可能なリソース共有、分離、および変動低減を提供する。

20

【 0 0 4 0 】

ストレージ・サービスおよびネットワーキングに関連するコンテナ（例えば、Kubernetes コンテナまたは仮想マシンコンテナ）およびスレッドには、キャッシュ・ウェイまたはバッファ空間のサブセットおよびメモリ帯域幅のサブセット（一実施形態ではメモリトランザクションクレジット - メモリ帯域幅施行）への有効な割り当てが割り当てられて、ネットワーク / ストレージ機能を、キャッシュ / バッファウェイおよび / またはメモリ帯域幅のサブセットに制約する。ストレージ・サーバ・ネットワーク割り当てと L 3 キャッシュ 6 0 6 の帯域幅 / サイズ割り当てとの適切な帯域幅マッチングにより、ストレージ・サーバ 4 0 0 は、メモリ 6 3 0 へのデータの流出をほとんど、または全くせずに、L 3 キャッシュ 6 0 6 のキャッシュ・ウェイ 7 0 6 の第 2 の部分を全体がまたはほぼ全体が通過するネットワーク・データ・フローとの間のストレージをサポートすることができる。

30

【 0 0 4 1 】

加えて、単一のストレージ・デバイス（NVMe ソリッド・ステート・ドライブなどであるが、これに限定されない）上の論理ボリューム（LVM）のアクセス帯域幅を設定した帯域幅、例えば、200 メガバイト / 秒（MB / s）に分割するためのオペレーティング・システムのメカニズム（例えば、Linux（登録商標）オペレーティング・システムにおけるデバイス・マップ）が存在する。論理ボリューム・レート・サービスの品質制御を上述のキャッシュ・ウェイの設定可能なサブセットと組み合わせることにより、より調整可能かつ予測可能なやり方で他のワークロードと共存するストレージ・サービス・ワークロードを提供する全体的なシステムソリューションが提供される。

40

【 0 0 4 2 】

図 8 は、L 3 キャッシュ 6 0 6 およびネットワーク・インターフェース・コントローラ

50

302を介してソリッド・ステート・ドライブ404からデータプレーン204にデータを移動させる方法を示す流れ図である。データは、ネットワーク・インターフェース・コントローラ302からソリッド・ステート・ドライブ404へ、L3キャッシュ606を介してデータプレーン204へ反対方向に移動させることもできる。

【0043】

ブロック800において、固定数のキャッシュ・ウェイ（例えば、キャッシュ・ウェイ706の第2のサブセット）がL3キャッシュ606内に割り当てられ、ソリッド・ステート・ドライブ404およびネットワーク・インターフェース・コントローラ302によって共有されるデータを記憶する。キャッシュ・ウェイの固定数は、システム性能要件に基づいて調整可能である。ソリッド・ステート・ドライブ404とネットワーク・インターフェース・コントローラ302との間で転送されるデータを記憶するためだけに使用するためにキャッシュ・ウェイ706の第2のサブセット内に割り当てられるキャッシュ・ウェイの数は、他のコアが使用するために利用可能なL3キャッシュ内のキャッシュ・ウェイの数を減らし、他のコアの性能を低下させることになる。

10

【0044】

一実施形態では、選択される固定数は、動作中に動的に変更されない。L3キャッシュ606の設定可能な部分のN個のキャッシュ・ウェイは、1つまたは複数のプロセッサコア602（例えば、602-2）によっても共有される。L3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットは、他のエージェントから分離されている。L3キャッシュ606内の残りのキャッシュ・ウェイ（例えば、キャッシュ・ウェイ712の第1のセット）は、他のエージェントによって使用/共有され得る。

20

【0045】

ブロック802において、ソリッド・ステート・ドライブ404は、L3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットに（ダイレクト・メモリ・アクセスを介して）データを直接書き込むように構成され、ネットワーク・インターフェース・コントローラ302は、L3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットからデータを直接（ダイレクト・メモリ・アクセスを介して）読み出すように構成されている。ソリッド・ステート・ドライブ404は、ネットワーク・インターフェース・コントローラ302がL3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットからデータを読み出している間に、L3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットにデータを書き込む。

30

【0046】

ブロック804において、ソリッド・ステート・ドライブ404がL3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットに書き込んでいるレートと、ネットワーク・インターフェース・コントローラ302がL3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットからデータを読み出しているレートとが一致しない場合、処理はブロック806に進む。ソリッド・ステート・ドライブ404がL3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットに書き込んでいるレートと、ネットワーク・インターフェース・コントローラ302がL3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットからデータを読み出しているレートとが一致する場合、処理はブロック802に進む。

40

【0047】

ブロック806において、データは、ソリッド・ステート・ドライブ404がL3キャッシュ606のキャッシュ・ウェイ706の第2のサブセットにデータを書き込み続けることを可能にするために、L3キャッシュ606からメモリ630に退避される。L3キャッシュからのデータの退避は、キャッシュ流出（cache spill）と呼ぶことができる。

【0048】

図9は、ソリッド・ステート・ドライブ404およびネットワーク・インターフェース・コントローラ302によって共有されるL3キャッシュ606（ラスト・レベル・キャ

50

ッシュとも呼ばれる)のキャッシュ・ウェイ906のN個のセットを分離するようにラスト・レベル・キャッシュを構成する実施形態を示す。Intel(登録商標)キャッシュ・アロケーション・テクノロジー(CAT)は、スレッド/アプリ/仮想メモリ(VM)/コンテナをグループ化することができるリソース制御タグとして機能するサービスクラス(CLOS)を含む。各サービスクラスは、所与のサービスクラスによってラスト・レベル・キャッシュがどれだけ使用され得るかを示す関連付けられたリソース容量ビットマスク(CBM)を有する。図9に示す実施形態では、サービス・クラス・テーブル902は、CLOS1~CLOS4とラベル付けされた4つのサービスクラスを有する。他の実施形態では、4つより多いまたは少ないサービスクラスが存在し得る。各CLOSレジスタは、プロセッサコア602ごとにビットを有し、ビットの状態は、コアが特定のサービス 10
 クラスの一部であるかどうかを示す。強制マスク904は、キャッシュマッピングをそれぞれのCLOSレジスタにおいて符号化されたキャッシュ・ウェイ・ビットマップに制限する回路/ロジックである。

【0049】

図示する実施形態では、各サービスクラスCLOS1~CLOS4は、マスク1~マスク4とラベル付けされた関連付けられた容量ビットマスクを有する。容量ビットマスクの各ビットの値は、サービスクラスに利用可能なL3キャッシュ606の量を示し、キャッシュ・ウェイ906のセット内のキャッシュ・ウェイのどれだけがサービスクラスCLOS1~CLOS4によって共有されているかどうかを示す。

【0050】

一実施形態では、サービスクラスに関連付けられた容量ビットマスクは、ネットワーク・インターフェース・コントローラによって共有されるデータを専ら記憶し、キャッシュ・ウェイ906のセットの一部は、ソリッド・ステート・ドライブ404およびネットワーク・インターフェース・コントローラ302によって共有される。

【0051】

図10は、ソリッド・ステート・ドライブ404およびネットワーク・インターフェース・コントローラ302によって共有されるラスト・レベル・キャッシュ内のキャッシュ・ウェイのセットを構成および使用方法の一実施形態を示す流れ図である。

【0052】

キャッシュ・アロケーション・テクノロジーは、アプリケーションの優先度またはサービス 30
 クラス(COSまたはCLOS)に基づいたリソース割り当てを可能にする。プロセッサは、アプリケーション(または個々のスレッド)を割り当てることができるサービスクラスのセットを公開する。それぞれのアプリケーションまたはスレッドに対するキャッシュの割り当ては、それらが関連付けられているクラスに基づいて制限される。各サービスクラスは、容量を表す容量ビットマスクを使用して設定することができ、クラス間の重複および分離の度合いを示すことができる。各論理プロセッサには、アプリケーション、スレッド、または仮想マシンがスケジュールされるときに、オペレーティング・システム/仮想マシンマネージャがサービスクラスを指定することができるように、公開されたレジスタが存在する。サービスクラスの使用は、リソース間で一貫しており、サービスクラスは、複数のリソース制御属性が付けられてもよく、これにより、コンテキストスワップ時 40
 のソフトウェアのオーバーヘッドが低減される。リソースごとに新しいタイプのサービスクラス・タグを追加するのではなく、サービスクラス管理のオーバーヘッドは一定である。指示されたアプリケーション/スレッド/コンテナ/VMに対するキャッシュ割り当ては、クラスおよびそのクラスに関連付けられたビットマスクに基づいてハードウェアによって自動的に制御される。ビットマスクは、L3キャッシュ用のモード・ステータス・レジスタを介して設定することができる。

【0053】

ブロック1000において、キャッシュ・アロケーション・テクノロジーにより、オペレーティング・システム(OS)、ハイパーバイザ/仮想マシンマネージャ(VMM)、または同様のシステムサービス管理エージェントが、アプリケーションが埋めることができ 50

るキャッシュ空間の量を指定することができるようになる。キャッシュ階層のどのレベルがサポートされているか、および最大割り当てビットマスクサイズなどの特定のキャッシュ・アプリケーション・テクノロジーの能力を照会するために、列挙型サポートが提供されている。

【 0 0 5 4 】

ブロック 1 0 0 2 において、オペレーティング・システムまたはハイパーバイザは、割り当てビットマスクのリストを介して特定のサービスクラスに利用可能なリソースの量を設定する。利用可能な容量マスクのビット長は、L 3 キャッシュの構成に依存する。

【 0 0 5 5 】

ブロック 1 0 0 4 において、コンテキスト・スイッチがある場合、処理はブロック 1 0 0 6 に進む。そうでない場合は、処理はブロック 1 0 0 8 に進む。

【 0 0 5 6 】

ブロック 1 0 0 6 では、現在実行中のアプリケーションのサービスクラスが実行環境（オペレーティング・システム / 仮想）に伝達される。新しいスレッドのサービスクラスが現在実行中のアプリケーション・サービスクラスと異なる場合、異なるサービスクラスをロードすることができる。処理はブロック 1 0 0 8 に進む。

【 0 0 5 7 】

ブロック 1 0 0 8 において、メモリ要求がある場合、処理はブロック 1 0 1 0 に進む。そうでない場合は、処理はブロック 1 0 0 4 に進む。

【 0 0 5 8 】

ブロック 1 0 1 0 では、メモリアクセスに関連付けられたサービスクラスを使用して、キャッシュ割り当てを実施する。処理はブロック 1 0 0 4 に進む。

【 0 0 5 9 】

図 2 に戻ると、物理クラスタ 2 0 0 におけるデータ管理プラットフォーム 1 0 0 のハードウェア障害の復旧は、ヘルスチェックとともにハードウェアまたはソフトウェアのロードバランスによって実行することができる。ハードウェア障害としては、計算サーバ 1 1 2、アクセラレータ・サーバ 1 1 4、データスイッチ 2 0 6、管理スイッチ 2 0 8、インフラストラクチャ・サーバ 1 1 6 およびユーティリティ・サーバ 1 1 8 のハードウェア障害を挙げることができる。サーバ（計算サーバ 1 1 2、アクセラレータ・サーバ 1 1 4、インフラストラクチャ・サーバ 1 1 6、およびユーティリティ・サーバ 1 1 8）のそれぞれは、ノードと呼ぶこともできる。典型的には、ロギングシステムを使用してイベントにフラグを立て、特定のイベントがログに記録されたときにオペレータが手動で介入して故障したハードウェアを除去または交換する。

【 0 0 6 0 】

しかしながら、現在のロードバランスは、データ管理プラットフォーム 1 0 0 のハードウェアの故障または劣化に基づいた、ハードウェア・コンポーネントの劣化および / または性能の劣化を考慮していない。加えて、現在のロードバランスは、企業のデータセンタにおいてスケーリングすることが困難である。

【 0 0 6 1 】

図 1 1 は、ハードウェアの故障または劣化がないシステムにおける通常動作の場合の、図 2 に示す物理クラスタ 2 0 0 におけるデータ管理プラットフォーム 1 0 0 内のラック 1 0 6 の実施形態のブロック図である。ハードウェア障害およびハードウェア劣化の透過的なシステムサービス修復のための方法およびシステムにより、プロセッサのハードウェア・イベントおよび測定値を直接的かつ効率的に公開することができる。テレメトリは、ルーティング情報ベース（RIB）1 1 1 2、転送情報ベース（FIB）5 1 0（図 5）、フィルタリングシステム（FS）1 1 1 4、およびインターネット・プロトコル・エニキャスト（Internet Protocol Anycast）への統合インターフェースと連携して、ハードウェア・イベントに基づいて、サーバ（計算サーバ 1 1 2、アクセラレータ・サーバ 1 1 4、インフラストラクチャ・サーバ 1 1 6、およびユーティリティ・サーバ 1 1 8）からの動的ルートを選択的に許可または抑制する。ルーティング情

報ベース (R I B) 1 1 1 2 は、特定のネットワークの宛先へのルートを記憶するデータテーブルである。

【 0 0 6 2 】

動的ルーティングは、最適なデータルーティングを提供するネットワーク技術である。動的ルーティングにより、ルータは、リアルタイムに変化する論理的なネットワークレイアウトに応じてパスを選択することができる。動的ルーティングでは、ルータ上で動作するルーティングプロトコルが、動的ルートテーブルの作成、維持、および更新を担当する。動的ルートとは、あるエンドポイントへのネットワーク・トラフィックを、環境に基づいて異なるルートを經由して転送することができるプロセスである。

【 0 0 6 3 】

サーバ (ノード) 1 1 0 0 a ~ e における故障したコンポーネントは、サーバ (ノード) 1 1 0 0 a ~ e 上で実行されている1つまたは複数のアプリケーションの機能および性能に影響を及ぼす可能性がある。サーバ (ノード) 1 1 0 0 a ~ e における故障したコンポーネントの例としては、ソリッド・ステート・ドライブ、メモリモジュール、または配電ユニットが挙げられる。データ管理プラットフォーム 1 0 0 内のサーバ (ノード) 1 1 0 0 a ~ e が劣化している場合、フィルタリングシステム (F S) 1 1 1 4 内のノード障害検出器によって故障コンポーネント・イベントが検出され、サービスに関連付けられたルートがルートテーブル 1 1 0 2 から撤回される。一実施形態では、ルートは、ルータ 5 0 4 (図 5)、例えば、K u b e r n e t e s の C a l i c o によって撤回される。

【 0 0 6 4 】

一実施形態では、オペレーティング・システム (O S) 1 1 0 6 は、L i n u x (登録商標) オペレーティング・システムである。サーバ (ノード) 1 1 0 0 a ~ e 上の B G P (B o r d e r G a t e w a y P r o t o c o l) クライアントは、F I B 5 1 0 からルーティング状態を読み取り、それを他のサーバ (ノード) 1 1 0 0 a ~ e 上で実行される他の B G P クライアントに配信する。F I B 5 1 0 内のルートは、要求に応答してエージェントによってセットアップされ、特定のワークロードのための接続性をプロビジョンする。B G P クライアントは、F I B 5 1 0 の更新に応答して、更新されたルートを他のサーバ (ノード) 1 1 0 0 a ~ e 上で実行されている B G P クライアントに配信する。

【 0 0 6 5 】

一実施形態では、F I B 5 1 0 にルートをセットアップするエージェントは F e l i x であり、B G P クライアントは B I R D である。B I R D は、U n i x (登録商標) ライクなオペレーティング・システム上でインターネット・プロトコル・パケットをルーティングするためのオープン・ソース実装である。F e l i x は、ルートを設定し、ネットワークポリシーを実施するためのノードごとのドメインデーモンである。

【 0 0 6 6 】

フィルタリングシステム 1 1 1 4 内のノード障害検出器は、ノードのハードウェアメトリクスを監視し、アラート (例えば、故障コンポーネント・イベント) を生成する。一実施形態では、障害または劣化は、プラットフォーム・テレメトリを介して検出され、障害 / 劣化イベントは、オープン・コレクタ、例えば、「c o l l e c t e d」に伝達され、その後、是正措置をとるイベント・ハンドラに伝達される。是正措置の一例は、ルートを除去することである。

【 0 0 6 7 】

インターネット・プロトコル・エニーキャストの実装では、故障したサーバ・サービスに関連付けられたルートが撤回され、接続されているピアからのルートの除去がトリガされる。接続されているピアは、データ管理プラットフォーム 1 0 0 のデータプレーン・ネットワークにあるすべてのネットワークデバイス (サーバおよびスイッチの両方) である。フローは、健全なまたは利用可能なサーバノード 1 1 0 0 a ~ e に透過的にリダイレクトされる。

【 0 0 6 8 】

図示する実施形態では、複数のノード (N 個) 1 1 1 0 を含むラック 1 0 6 が示されて

10

20

30

40

50

おり、N個のノード1110a~1110eのうちの5個が示されている。各ノードは、計算サーバ112、アクセラレータ・サーバ114、インフラストラクチャ・サーバ116、またはユーティリティ・サーバ118とすることができる物理サーバである。ユーティリティ・サーバ118は、データ管理プラットフォームにおいて管理タスクを実行する制御プレーン・サーバ・ノードと呼ぶこともできる。

【0069】

一実施形態では、物理クラスタ200には、最大16個のラック106、およびラックあたり最大20個のノード1110が存在する。他の実施形態では、16個を超えるラックおよびラックあたり20個のノードが存在することができる。最初の3つのラックにはノードあたり1台のユーティリティ・サーバ118があり、第2および第3のラックにはラックあたり1台のインフラストラクチャ・サーバ116があり、最初の3つのラック106にはラック106あたり最大14台の計算サーバ112があり、次の13のラック106には最大16台の計算サーバがあり、ラック106ごとに最大4台のアクセラレータ・サーバ114がある。アクセラレータ・サーバ114は、ストレージ処理タスクを実行し、ストレージ・サーバ400(図4)と呼ぶことができる。

10

【0070】

一実施形態では、各ノード1110a~1110eは、ポッド1104およびオペレーティング・システム(OS)1106(例えば、Red Hat Enterprise Linux(登録商標)(RHEL)オペレーティング・システム)を含む。ポッド1104は、Kubernetesアプリケーションの基本的な実行単位であり、作成またはデプロイすることができるKubernetesオブジェクトモデルにおける最小かつ最も単純な単位である。ポッド1104は、デプロイメントの単位、すなわち、単一のコンテナ、または密に結合されリソースを共有する少数のコンテナのいずれかを含むことができる、Kubernetesにおけるアプリケーションの単一のインスタンスを表す。

20

【0071】

ポッド1104は、ストレージ/ネットワークが共有された1つまたは複数のコンテナのグループである。ポッド1104内のコンテナは、インターネット・プロトコル(IP)アドレスおよびポート空間を共有し、標準的なプロセス間通信を使用して他のポッド1104と通信することができる。異なるポッド1104内のコンテナは、別個のインターネット・プロトコル・アドレスを有し、ポッド1104のIPアドレスを使用して互いに通信する。

30

【0072】

エニーキャストは、単一の宛先アドレスが2つ以上のエンドポイント宛先への複数のルーティングパスを有するネットワークアドレス指定およびルーティング方法である。ルータ504は、ホップ数、距離、最低コスト、レイテンシ測定値に基づいて、または最も混雑していないルートに基づいて、ノード1110a~e間のパスを選択する。通常の動作では、ラック106内の各ノード1110a~eは、分散共通サービスのために同じIP(Internet Protocol)アドレス(Anycastアドレス)をアドバタイズする。

【0073】

図11に示す例を参照すると、サービスは、ノード1110a~eのそれぞれからアドバタイズされ、そのサービスに関連付けられたエニーキャストアドレス(IPアドレス)は、6つのノード1100a~eすべてにわたって同じである。この例では、IPアドレス「10.1.1.1」である。各ノード1110a~eは、データスイッチ206のルートテーブル1102に記憶された固有のイーサネット(登録商標)・アドレスを有する。ラック106内のデータスイッチ206は、トップ・オブ・ラック(TOR)スイッチと呼ぶこともできる。

40

【0074】

例えば、ノード1100aのIPアドレスは「10.1.1.1」であり、ノード1100aのイーサネット(登録商標)・アドレスは192.168.1.1である。ハード

50

ウェア障害または劣化イベントがない場合、オーケストレータ/スケジューラ102（例えば、Kubernetes）によって管理されるルートテーブル1102により、（ラック106内のすべてのノード1100a～eを経由する）すべてのルートをアドバタイズすることが可能である。データスイッチ206は、宛先に到達するために、単一のIPアドレス（10.1.1.1）および6つのパス（ノード1100a～eのうちの1つを経由する）を見る。宛先は、アプリケーション・インスタンスである。一実施形態では、アプリケーション・インスタンスはKubernetesサービスである。アプリケーションは、データ管理プラットフォーム100内のネットワーク・トラフィックを負荷分散し、アプリケーションおよびデータへのアクセスを提供するために、複数のアプリケーション・インスタンスとして生成され得る。

10

【0075】

データスイッチ206は、内蔵の負荷分散方法、例えば、等価コスト・マルチパス・ルーティング（ECMP）を使用して、ノード1100a～eへのパスのうちの1つを選択することができる。等価コスト・マルチパス・ルーティング（ECMP）は、単一の宛先へのネクスト・ホップ・パケット転送が、ルーティング・メトリクス計算において最上位に位置（tie）する複数の「最良パス」を介して行うことができるルーティング戦略である。マルチパス・ルーティングは、単一のルータに限定されたホップごとの判断であるため、ほとんどのルーティングプロトコルと併せて使用することができる。

【0076】

図12は、故障したシステムにおける縮退動作について、図2に示す物理クラスタ200におけるデータ管理プラットフォーム100内のラック106の実施形態のブロック図である。

20

【0077】

障害または劣化イベント（例えば、計算サーバ112における故障したネットワーク・インターフェース・コントローラ302、ストレージ・ノード400における故障したソリッド・ステート・ドライブ404、またはノード1100a～eにおける不安定なオペレーティング・システム1106）中に、影響を受けたノード1100a～eは、影響を受けているアプリケーション（サービスとも呼ばれる）に関連付けられたルートのアドバタイズメントを抑制する。

【0078】

オーケストレータ/スケジューラ102がKubernetesである実施形態では、障害または劣化イベントがネットワーク接続性に関連している場合、イベントは、Kubernetes\OpenShiftおよびKubernetesネットワークコンポーネントによって処理される。Kubernetesは、ノード1100a～eがネットワーク上で利用可能でないことを検出する。ネットワークコンポーネントは、物理クラスタ200にわたってルートテーブル1102を更新する。

30

【0079】

障害または劣化イベントがネットワーク接続性に関連していない場合、イベントは、ポッド1104内のLMA（Logging Monitoring Alerting）スタックによって処理される。各ノード1100a～e上で生成されるエクスポートは、LMAスタックにメトリクスデータを定期的に提供する。メトリクスに基づいて、ノード1100a～eは追加のラベルでマークされ、潜在的に追加のアクションが行われる可能性がある。例えば、故障または劣化したノード1100a～e上で実行されているコンテナを、別のノード1100a～e上で再スケジュールすることができる。

40

【0080】

図12に示すように、ノード1100aは、劣化または障害イベントを有し、ルートテーブル1102を介して、影響を受けたサービスに関連付けられたルートを抑制する（決してアドバタイズしない）か、または撤回する。本例では、エニーキャスト・アドレスは、ラック106内の各ノード1100a～eに対して10.1.1.1である。ノード1100aへのルートが撤回され、例えば、ノード1100aはそのルートのアドバタイズ

50

を停止し、データスイッチ 206 はそのルートをルートテーブル 1102 から除去する。

【0081】

エニーキャスト・アドレス (10.1.1.1) への 6 つの利用可能なパスのうち、ノード 1100a へのパス (宛先 10.1.1.1、ネクスト・ホップ 192.168.1.1) がデータスイッチ 206 のルートテーブル 1102 から除去される。ネクスト・ホップは、それぞれのノード 1100a ~ e に関連付けられた固有のインターネット・プロトコル (IP) アドレスである。エニーキャスト・IP アドレスは、それぞれのノード 1100a ~ e 上のアプリケーション・インスタンスまたはポッド 1104 の IP アドレスである。ネットワーク・トラフィックは、ノード 1100a ~ e 上のポッド 1104 に転送され、次いで、ノード 1100a ~ e 内のアプリケーション・インスタンスに転送される。残りのすべてのトラフィックフローは、データスイッチ 206 のルートテーブル 1102 を介して利用可能な残りのパスに分配される。

10

【0082】

Linux (登録商標) オペレーティング・システムを使用するシステムの実施形態では、生のセンサデータから得られる値 (例えば、「sysfs」インターフェースを介した Linux (登録商標) 「libsensors」ライブラリによる「臨界最大電力一致」) などの障害または劣化イベント中に、影響を受けたノード 1100a ~ e のポッド 1104 内のイベントアクション検出器およびマネージャは、影響を受けた Kubernetes サービスに関連付けられたルートのアダプタイズメントを抑制する。

【0083】

Kubernetes サービスは、クラスタ 200 内で実行されるポッド 1104 の論理セットを定義する抽象化であり、これらはすべて同じ機能を提供する。作成されると、各 Kubernetes サービスには、ルートとなる固有のインターネット・プロトコル (IP アドレス) (cluster IP と呼ばれる) が割り当てられる。割り当てられた IP アドレスは、Kubernetes サービスが活着している間に変更されない。

20

【0084】

ポッド 1104 は、オーケストレータ/スケジューラ 102 内の Kubernetes サービスと通信するように構成され得る。Kubernetes サービスへの通信は、Kubernetes サービスのメンバであるポッド 1104 に自動的に負荷分散される。複数のノードが同じサービス IP をアダプタイズすることができ、これは「エニーキャスト (Anycast)」と呼ばれる。Kubernetes サービスの一例は、NGINX (オープン・ソース、高性能 HTTP サーバおよびリバースプロキシ) および IMA P / POP3 プロキシサーバ、DNS (Domain Name System) および Apache (オープン・ソースのウェブサーバ) などのアプリケーションをホストするポッド 1104 またはポッド 1104 のセットをバックアップする Cluster IP である。

30

【0085】

イベントが検出された後、影響を受けたノード 1100a のポッド 1104 内のイベントアクション検出器およびマネージャは、影響を受けたノード 1100a 上のサービスに関連付けられたルートをブラックホール化する (別のパスを抑制し、アダプタイズしない) スクリプトを開始する。

40

【0086】

故障したノード (この場合、ノード 1100a) が再び正常に機能するようになると、故障したノードは、以前に故障したサービスに関連付けられたエニーキャスト IP アドレス (10.1.1.1) を自動的にアダプタイズし、透過的に物理クラスタ 200 に再挿入される。データスイッチ 206 は、ノード 1100a の別のパス (宛先 (ポッドの IP アドレス) 10.1.1.1、ネクスト・ホップ (ノード 1100a の IP アドレス) 192.168.1.1) を検出し、これをマルチパス宛先として既存のルートテーブル 1102 に追加する。

【0087】

50

図13は、ノード状態および障害を検出するためにデータ管理プラットフォーム100によって使用される、ノード1100a~eのポッド1104(図11)のコンテナ内のメトリクス・エクスポートを示すブロック図である。図示する実施形態では、4つのメトリクス・エクスポートおよび他のエクスポート1314が存在する。各メトリクス・エクスポートは、ポッド1104(図11)の別のコンテナ内にある。

【0088】

デバイス・マップ・エクスポート1310は、デバイス・マップ・ボリュームから低レベルのメトリクスを収集する。デバイス・マップ・ボリュームから収集される低レベルのメトリクスの例には、平均読み取り/書き込み時間、平均待ち時間、利用率、キューサイズ、1秒あたりの書き込み/読み取り数、1秒あたりの読み取り/書き込みサイズ、1秒あたりのマージされた読み取り/書き込みが含まれる。

10

【0089】

ストレージ・エクスポート1312は、ソリッド・ステート・ドライブから低レベルのメトリクスを収集する。収集される低レベルのメトリクスの例には、ソリッド・ステート・ドライブ内の不揮発性メモリに対するプログラムおよび消去が失敗した回数、ならびにエンド・ツー・エンドのエラー検出回数、巡回冗長検査(CRC)エラー回数、時限ワークロード・タイマ、熱スロットル状態、再試行バッファ・オーバーフロー回数、ウェア・レベリング回数、時限ワークロード媒体摩耗、時限ワークロードホスト読み取り/書き込み比率、停電切迫(pli)ロック損失回数、ソリッド・ステート・ドライブ内の不揮発性メモリに書き込まれるバイト、ホストによってソリッド・ステート・ドライブに書き込まれるバイト、および残りのシステム領域の残存寿命が含まれる。

20

【0090】

メモリ帯域幅エクスポート1308は、メモリ帯域幅モニタに基づいて低レベルのメトリクスを収集する。プロセッサ・カウンタ・モニタ(PCM)は、アプリケーション・プログラミング・インターフェース(API)であり、Intel(登録商標)プロセッサの性能およびエネルギー・メトリクスを監視するためのAPIに基づくツールのセットである。メモリ帯域幅エクスポート1308は、プロセッサ・カウンタ・モニタを使用して、低レベルのメトリクスを収集する。収集されるメモリ帯域幅に関連する低レベルのメトリクスの例には、チャンネル読み出し/書き込み、メモリ読み出し/書き込みMega Byte/秒、読み出し/書き込み、メモリ、およびMemory Mega Byte/秒が含まれる。

30

【0091】

ネットワーク・インターフェース・コントローラ1306のエクスポートは、ネットワーク・インターフェース・コントローラから低レベルのメトリクスを収集する。収集される低レベルのメトリクスの例には、送信キューのドロップ、送信キューの停止、バッファ外受信、送信エラー、バッファしきい値通過受信、および受信/送信シグナル・インテグリティが含まれる。

【0092】

ポッド1104に含まれる他のエクスポートには、サーバ・シャーシ・エクスポート1316、ノード・エクスポート1318、およびブラックボックス・エクスポート1320が含まれる。サーバ・シャーシ・エクスポート1316は、サーバ・シャーシから低レベルのメトリクスを収集する。ノード・エクスポート1318は、オペレーティング・システム・レベルのメトリクスを収集する。ブラックボックス・エクスポート1320は、HTTP(Hyper Text Transfer Protocol)/TCP(Transmission Control Protocol)エンドポイントに関連するメトリクスを収集する。

40

【0093】

一部のエクスポート(デバイス・マップ1310およびストレージ1312)は、ソリッド・ステート・ドライブ404上のメトリクスのみを監視しているため、ストレージ・ノード1302内でのみ使用される。図13に示すように、非ストレージ・ノード130

50

4 (例えば、計算サーバ112、ユーティリティ・サーバ118、インフラストラクチャ・サーバ116、および非ストレージ・アクセラレータ・サーバ114)は、デバイス・マップ1310およびストレージ1312のエクスポートを含まない。データ管理プラットフォーム100は、エクスポート(デバイス・マップ、ソリッド・ステート・ドライブ、メモリ帯域幅、およびネットワーク・インターフェース・コントローラ)からのメトリクスに基づいて、トラフィックを健全なノード上のアプリケーション・インスタンスにリダイレクトすることによって、このようなイベントを検出し、反応することができる。

【0094】

データ管理プラットフォーム・クラスタ200内のすべてのノードがすべて正常に動作している場合、複数のノード上でアプリケーション・インスタンスを生成することに制限はない。その場合、ネットワークからのトラフィックは、図12に関連して説明したように動作している。エクスポートのうちの一つが、ノードにハードウェア問題があることを示すメトリクスを検出した場合、ポッド1104のLMA(Logging Monitoring Alerting)は、図13に関連して説明したように、不健全なノードを除外し、ネットワーク・トラフィックをブロックするアクションを実行する。

10

【0095】

図14は、物理クラスタ200におけるデータ管理プラットフォーム100のハードウェア障害を管理するための方法を示す流れ図である。

【0096】

ブロック1400では、図13に関連して説明したエクスポート(ネットワーク・インターフェース・エクスポート1306、メモリ帯域幅エクスポート1308、デバイス・マップ・エクスポート1310、ストレージ・エクスポート1312、および他のエクスポート1314)は、ノード1100a~eのメトリクスを継続的に監視する。ノード1100a~eのポッド1104内のLMAは、エクスポートからメトリクスを収集する。

20

【0097】

ブロック1402において、エクスポートから収集されたすべてのメトリクスが良好である場合、処理はブロック1404に進む。そうでない場合は、処理はブロック1410に進む。

【0098】

ブロック1404において、ノード1100a~eが動作可能である場合、ノードは動作可能とマークされ、すべてのメトリクスは良好であり、ノード1100a~eがエラーなしで動作していることを示す。ノードが非動作状態から回復し、以前に非動作とマークされていた場合、ノードは動作可能とマークされる。

30

【0099】

ブロック1406において、ラック106内のすべてのノード1100a~eが動作可能である。RIB1112は、以前は非動作であったノード1100a~eへのルートを復元し、回復したノード1100a~e上のアプリケーション・インスタンスへのトラフィックを復元するように更新される。

【0100】

ブロック1408において、ラック106内の動作ノード1100a~e上のすべてのアプリケーション・インスタンスに対してトラフィックが再開される。処理はブロック1400に進む。

40

【0101】

ブロック1410において、エクスポートから収集されたメトリクスのすべてが良好ではなく、ノード1100a~eにおける障害または劣化イベントを示す。ノード1100a~eは、非動作とマークされる。エニークキャスト・サービス・アドバイズメントおよびECMPは、他のノード1100b~eの他のアプリケーション・インスタンスを介してアクセスをアプリケーションに転送する。

【0102】

ブロック1412において、データネットワークは非動作ノード1100aに利用可能

50

ではなく、非動作ノード1100a上のアプリケーション・インスタンスへのアクセスは利用可能ではない。クラスタ内のすべてのノードのRIBが更新される。

【0103】

ブロック1414において、動作ノード1100b~eのアプリケーション・インスタンスへのトラフィックが、ラック106内のすべての動作ノード1100b~eに対して再開される。アプリケーション・インスタンスへのトラフィックは、非動作ノード1100aに送信されない。処理はブロック1400に進む。

【0104】

現在のロードバランサ(ソフトウェアまたはハードウェア)は、ヘルスチェック、スクリプティングまたは監視システムとともに、性能しきい値を超える(例えば、中央処理装置(CPU)の利用率が90%を超える)と動的に反応しない。これらの性能しきい値をより綿密に監視および管理することにより、ネットワーク機能仮想化(NFV)などの新しいワークロードをサポートするための、ますます厳しくなるパフォーマンスSLA(performance service-level agreement)を含む、より厳しいワークロード要求を満たすデプロイメントが可能になる。

10

【0105】

一実施形態では、圧迫状態および性能劣化の指標を提供する性能しきい値に応答する動的かつ透過的なスケーリングは、定義された性能しきい値に基づいてKubernetesサービスレベルごとに行われる。これにより、トリガされた性能しきい値に基づく動的検出および透過的なサービススケーリングが可能になり、より最適化されたスケラブルなKubernetes実装が可能になる。

20

【0106】

前述したように、Kubernetes制御プレーンは、インフラストラクチャ・サーバ116上でホストされ、Kubernetesホスト・エージェントは、すべての計算サーバ112およびアクセラレータ・サーバ114上で実行される。ハードウェア・イベントおよび測定値の直接的かつ効率的な公開が、ルーティングおよび情報ベース(RIB)1112への統合インターフェースと併せて提供される。ハードウェア・イベントおよび測定値の例には、「sysfs」インターフェースを介してLinux(登録商標)「libsensors」ライブラリを通じて公開される生のセンサデータなどのテレメトリが含まれる。ノードは、以前はミニオン(minion)として知られていたKubernetesのワーカーマシンである。ノードは、クラスタに応じて、仮想マシンまたは物理マシン(サーバ)である場合がある。各ノードは、ポッド1104を実行するのに必要なサービスを含む。ノード上のKubernetesサービスには、コンテナランタイム(コンテナを実行し、ノード上のコンテナイメージを管理するソフトウェア)、kubelelet(pod1104を実行する)、およびkube-proxy(クラスタ内の各ノード上で実行され、ノード上のネットワークルールを維持するKubernetesサービスの一部を実装するネットワークプロキシ)が含まれる。

30

【0107】

図15は、性能しきい値1502を監視および管理してノード状態および障害を検出するためにデータ管理プラットフォーム100によって使用される、ノード1500のポッド1104(図11)におけるハードウェア・イベントおよび測定値を示すブロック図である。

40

【0108】

ハードウェア・イベントおよび測定値の公開は、ポッド1104およびIPエニーキャストのイベント検出器およびモニタ1504を通して提供される。ハードウェア・イベントおよび測定値の公開により、CPU利用率などの定義された性能しきい値に基づいて、サーバ(データ管理プラットフォーム100のノード)からの動的ルートを選択的に許可または抑制することができる。これらのしきい値は、ランタイム前またはランタイム中に設定することができる。

【0109】

50

性能しきい値の例には、CPU利用率、ソリッド・ステート・ドライブ404の1秒当たりのInput/Output Operation IOPS、または帯域幅利用率が含まれる。性能しきい値は、特定のKubernetesサービスに関連付けられ、Kubernetesサービスごとの粒度を提供する。性能しきい値イベントが検出されるか、またはメトリクスが一致した後、イベント検出器およびモニタ1504は、影響を受けたノード上のサービスに関連付けられたルートをブラックホール化する。

【0110】

ネットワークにおいて、ブラックホールとは、着信または発信トラフィックが、その意図された受信者にデータが到達しなかったことを送信元に知らせることなく、静かに廃棄（または「ドロップ」）されるネットワーク内の場所を指す。ネットワークのトポロジを調べる際に、ブラックホール自体は見ることはできず、失われたトラフィックを監視することによってのみ検出され得る。

10

【0111】

利用率しきい値は、設定された期間にわたるリソースの割合を指定する。例えば、リソースがソリッド・ステート・ドライブ404へのNVMeインターフェースの帯域幅である場合、帯域幅の利用率しきい値は、ソリッド・ステート・ドライブ404のNVMeインターフェースの最大帯域幅（例えば、毎秒25ギガビット）の75%とすることができる。利用率しきい値がフィルタリングシステム1114によって満たされた場合、その所与のサービスに関連付けられたルートは撤回される。IPエニキャストを使用する実施形態では、故障ノードサービスに関連付けられたルートが撤回され、接続されたピアからのルートの除去がトリガされる。フローは、「動作範囲」内にあるノードに透過的にリダイレクトされる。

20

【0112】

トリガされた利用率しきい値は、設定された期間を超えた場合に、しきい値通知をトリガするリソースの割合を指定する。各ノードは、ノードがアダタイズするルートのセットを有する。ルートのセットは、ノードのルートテーブル1102で見ることができる。

【0113】

トリガされる性能しきい値がない場合、すべてのアクティブなルートがアダタイズされる（デフォルトの動作モード）。接続されたデータスイッチ206には1つのエニキャストIPアドレスが見えており、例えば、エニキャストIPアドレス（10.1.1.1）と、5つのノード1100a～eのうちの1つを経由して、この場合Kubernetesサービスである宛先に到達する5つのパスとが利用可能である。負荷分散方法、例えばECMPを使用して、パスのうちの1つを選択することができる。

30

【0114】

トリガされた性能しきい値イベントの間、影響を受けたノードは、影響を受けているサービスに関連付けられたルートのアダタイズメントを抑制する。ノードは、影響を受けたノードのサービスに関連付けられたルートをブラックホール化することによって、ルートのアダタイズメントを抑制する。

【0115】

図12を参照すると、ノード1100aは、マッチング性能しきい値イベントを有し、影響を受けたサービスに関連付けられたルートを抑制する（決してアダタイズしない）か、または撤回する。本例では、アドレスは、エニキャストIPアドレス10.1.1.1である。ルート（パス）は撤回され、すなわち、ノード1100aは、エニキャストIPアドレス10.1.1.1のパスのアダタイズを停止し、データスイッチ206は、エニキャストIPアドレス10.1.1.1のノード1100aを経由するパスをルートテーブル1102から除去する。そのアドレスへの5つの利用可能なパスのうち、ノード1100aへのパス（宛先10.1.1.1、ネクスト・ホップ192.168.1.1）がデータスイッチ206のルートテーブル1102から除去される。これにより、影響を受けたノードへの接続が除去される。残りのトラフィックフローはすべて、残りのパスに分散される。

40

50

【0116】

「性能が影響を受けた」ノード（この場合、ノード1100a）が正常に機能した後、ノード1100aは、自動的にネクスト・ホップ（以前に障害を受けたサービスに関連付けられたイーサネット（登録商標）・アドレス（宛先10.1.1.1、ネクスト・ホップ192.168.1.1））をアドバタイズする。ノード1100aは、ノード1100a上のサービスに関連付けられた以前にインストールされたブラックホールルートを除去することによって、透過的にクラスタに再挿入される。データスイッチ206は、ブラックホール化された（すなわち、そのアドレス（宛先10.1.1.1、ネクスト・ホップ192.168.1.1）に対して抑制され、別のパスをアドバタイズしない）ルートを検出し、これをマルチパス宛先としてルートテーブル1102に追加する。

10

【0117】

図16は、ストレージ・ノード400の性能を監視するためにラック106内のストレージ・ノード400において実施される方法を示す流れ図である。

【0118】

ブロック1600において、図15に関連して説明したイベント検出器およびモニタ1504は、ノード内の性能しきい値1502を継続的に監視する。

【0119】

ブロック1602において、性能しきい値が所定のしきい値最大値と一致しない場合、処理はブロック1604に進む。それらが一致する場合、処理はブロック1610に進む。

20

【0120】

ブロック1604において、ノードがサービス・レベル・アグリーメント・パラメータを満たしている場合、ノードは適合しているとマークされる。ノードが以前に非適合とマークされていた場合、ノードは適合とマークされる。

【0121】

ブロック1606において、クラスタ内のすべてのノードが適合しており、ルーティング情報ベース（RIB）が更新されて、以前に適合していなかったノードへのルートを復元し、このようなノード上のアプリケーション・インスタンスへのトラフィックを復元する。

【0122】

ブロック1608において、クラスタ内の適合しているノード上のすべてのアプリケーション・インスタンスに対してトラフィックが再開される。処理はブロック1600に進む。

30

【0123】

ブロック1610において、ノードがSLAパラメータを満たさない場合、ノードは非適合とマークされる。

【0124】

ブロック1612において、データネットワークは非適合のノードには利用できず、アプリケーション・インスタンスへのアクセスも利用できない。クラスタ内のすべてのノードのRIBが更新される。

40

【0125】

ブロック1614において、クラスタ内のすべての適合しているノードに対してトラフィックが再開される。トラフィックは適合していないノードには送信されない。処理はブロック1600に進む。

【0126】

図2に戻ると、前述したように、サーバ（計算サーバ112、アクセラレータ・サーバ114、インフラストラクチャ・サーバ116、およびユーティリティ・サーバ118）のそれぞれは、ノードと呼ぶこともできる。オーケストレータ/スケジューラ102は、固定数のノードを管理する。ノードの数は、データセンタ内のトラフィックのピークに対応するように選択され、典型的にはオーバープロビジョニングされる。現在のデータセン

50

タでは、ワークロードが圧迫下にある場合、オーケストレータ/スケジューラ 102 は、ワークロードを抑制するか、またはワークロードがデータセンタの性能を低下させる圧迫下にあるノード上の追加のワークロードのスケジューリングを防止することができる。

【0127】

典型的には、データセンタ内の負荷がCPU、メモリ、またはストレージの点で容量に達すると、手動でデータセンタのサイズ変更が実行される。データセンタのサイズ変更には、新しいノードの追加、プロビジョニング、およびコンフィギュレーションが含まれる。負荷が減少すると、データセンタはさらにオーバープロビジョニングされる。

【0128】

一実施形態では、データセンタの総所有コスト(TCO)は、データセンタにおけるリソースの過剰サブスクリプションを減少させることによって低下させることができる。オーケストレータが管理するデータセンタ内の様々な圧迫状態を監視し、追加の論理リソースで既存のノードのサイズ変更を要求することによって、総所有コスト(TCO)が下げられる。

【0129】

図17は、計算ノード1704の実施形態のブロック図である。計算ノード1704は、プロセッサ、メモリ、および入出力(I/O)制御ロジックを1つのSoCパッケージに統合したシステム・オン・チップ(SoCまたはSoC)604を含む。SoC604は、少なくとも1つの中央処理装置(CPU)モジュール608およびメモリコントローラ614を含む。

【0130】

図示する実施形態では、SoC604は、内部グラフィック処理装置(GPU)1700も含む。内部GPU1700は、1つまたは複数のGPUコアと、GPUコアのためのグラフィックス関連データを記憶することができるGPUキャッシュと、を含むことができる。GPUコアは、1つまたは複数の実行ユニットならびに1つまたは複数の命令キャッシュおよびデータキャッシュを内部に含むことができる。さらに、内部グラフィック処理装置(GPU)1700は、1つまたは複数の頂点処理ユニット、ラスタ化ユニット、メディア処理ユニット、およびコーデックなど、図17に示されていない他のグラフィックス論理ユニットを含むことができる。

【0131】

他の実施形態では、メモリコントローラ614は、SoC604の外部にあってもよい。CPUモジュール608は、レベル1(L1)およびレベル2(L2)キャッシュ620を含む少なくとも1つのプロセッサコア102と、CPUモジュール608内の他のプロセッサコア102と共有されるレベル3(L2)キャッシュ606と、を含む。

【0132】

一実施形態では、メモリ630は揮発性メモリである。さらに別の実施形態では、メモリ630は、1つまたは複数のメモリモジュールに含まれ得るバイトアドレス指定可能なライト・イン・プレイスNVMデバイスおよび揮発性メモリデバイスの両方を含む。リソース・マネージャ・エージェント1706およびワークロード1708は、メモリ630に記憶される。

【0133】

計算ノード1704は、永続メモリ1702も含む。永続メモリ1702は、バイトアドレス指定可能なライト・イン・プレイス3次元クロス・ポイント・メモリ・デバイス、または他のバイトアドレス指定可能なライト・イン・プレイス不揮発性メモリデバイス、または他のメモリを含むことができる。バイトアドレス指定可能なライト・イン・プレイス3次元クロス・ポイント・メモリ・デバイスの一例は、3DXPoint(例えば、Intel(登録商標)Optane(登録商標)およびMicron(登録商標)QuantX(登録商標))である。

【0134】

図18は、計算ノード1804の別の実施形態のブロック図である。計算ノード180

10

20

30

40

50

4 は、プロセッサ、メモリ、および入出力（I/O）制御ロジックを1つのSoCパッケージに統合したシステム・オン・チップ（SOCまたはSoC）604を含む。SoC604は、少なくとも1つの中央処理装置（CPU）モジュール608およびメモリコントローラ614を含む。

【0135】

計算ノード1804はまた、SoC604内の入力/出力（I/O）サブシステム612に通信可能に結合されたフィールド・プログラマブル・ゲート・アレイ（FPGA）1800およびアクセラレータ1802を含む。一実施形態では、FPGA1800は、Intel（登録商標）のAgilex（登録商標）FPGAデバイスである。

【0136】

図19は、論理リソースを自動的に追加および除去するためのリソース・マネージャ1950を含む、図2に示す物理クラスタにおけるデータ管理プラットフォーム100内のラック106の実施形態のブロック図である。ラック106は、複数の計算ノード1904-1、1904-2、1904-3およびアクセラレータ・ノード1902を含む。計算ノード1904-1、1904-2、1904-3は、図17に関連して説明した計算ノード1704または図18に関連して説明した計算ノード1804を含むことができる。

【0137】

リソース・マネージャ1950は、メトリクスを監視して、いつ論理リソースを自動的にアタッチし、構成するかを決定する。一実施形態では、リソース・マネージャ1950は、オーケストレータ/スケジューラ102内にある。他の実施形態では、リソース・マネージャ1950は、計算ノード1904-1、1904-2、1904-3のうちの1つ、またはデータ管理プラットフォーム100内の別のコンポーネントに含まれる。リソース・マネージャ1950は、アクセラレータ・ノード1902と、データ管理プラットフォーム100のすべてのメトリクスとにアクセスする。データ管理プラットフォーム100内のリソース・マネージャ1950は、ユーザの介入なしに、論理リソース（例えば、メモリ、ストレージボリューム、グラフィック処理装置（GPU）、およびフィールド・プログラマブル・ゲート・アレイ（FPGA）論理リソース）を自動的にアタッチ、デタッチ、および構成する。

【0138】

図19に示す特定の非限定的な例では、ラック106内に3つの計算ノード1904-1、1904-2、1904-3と1つのアクセラレータ・ノード1902がある。オーケストレータ/スケジューラ102は、ラック106内の計算ノード1904-1、1904-2、1904-3のそれぞれにおけるワークロードおよびプロセスを監視する。

【0139】

データ管理プラットフォーム100は、オーケストレータ・メトリクス1911、ノード・メトリクス1910、およびワークロード・メトリクス1912を含む。オーケストレーション・メトリクス1911は、オーケストレータ/スケジューラ102によって管理される。ワークロード・メトリクス1912は、リソース・マネージャ150によってアクセス可能である。ワークロード・メトリクス1912は、ワークロードによって公開され、リソース・マネージャ1950によって照会されるか、またはメトリクス・アグリゲータ1914によって照会され得て、メトリクス・アグリゲータ1914は、リソース・マネージャ1950によって照会される。ノード・メトリクス1910は、ノード・エクスポート（例えば、github.com/prometheus/node_exporter）によって公開される。ノード・メトリクス1910は、リソース・マネージャ1950によって照会され、またはメトリクス・アグリゲータ1914によって照会され得て、メトリクス・アグリゲータ1914は、リソース・マネージャ1950によって照会される。

【0140】

オーケストレータ/スケジューラ102は、基本的なノード・メトリクスをノード・メ

10

20

30

40

50

トリクス 1910 に記憶する。基本的なノード・メトリクスには、計算ノードごとのワークロードの数、計算ノードごとのプロセスの数、圧迫状態、計算ノードごとの CPU 利用率、および計算ノードごとのメモリ利用率が含まれる。圧迫状態は、計算ノード 1904 が圧迫下にあるかどうかを示す。

【0141】

計算ノード 1904 が、計算ノード 1904 上で実行されるワークロード 1708 の性能に影響を与える高いリソース利用率を経験している場合、計算ノード 1904 は圧迫下にある。計算ノード 1904 が圧迫下にあるかどうかを判定するために、追加のノード・メトリクスが監視され、ノード・メトリクス 1910 に記憶される。追加のノード・メトリクスには、プロセスごとの CPU 利用率、プロセスごとのメモリ帯域幅利用率、プロセスごとのメモリ利用率、プロセスごとのストレージレイテンシ、プロセスごとのストレージ利用率、プロセスごとの毎秒のストレージ入力/出力、プロセスごとの GPU および/または FPGA 利用率、ならびにプロセスごとの GPU および/または FPGA レイテンシが含まれる。

10

【0142】

オーケストレータ/スケジューラ 102 はまた、ワークロード・メトリクス 1912 を監視し、記憶する。ワークロード・メトリクス 1912 には、クライアントの数、平均応答レイテンシ、およびパーセンタイル・メトリクスが含まれる。パーセンタイル・メトリクスの例は、99 パーセンタイル・レイテンシまたは 99.9 パーセンタイル・レイテンシであり、これは、99% または 99.9% のワークロードに対する最大レイテンシである。

20

【0143】

リソース・マネージャ 1950 は、メトリクス (ノード・メトリクス 1910、ワークロード・メトリクス 1912、およびオーケストレータ・メトリクス (1911)) を集約して、圧迫状態が発生したときに圧迫状態を検出する。リソース・マネージャ 1950 はまた、メトリクスを集約して、圧迫状態が発生する前に圧迫状態を検出する。時系列分析アルゴリズム (Time Series Analysis algorithm) を使用することにより、圧迫状態の発生前に圧迫状態を検出することができる。時系列分析アルゴリズムとしては、マルコフチェーン/連鎖アルゴリズムまたは人工知能アルゴリズム (例えば、ニューラルネットワークもしくは遺伝的アルゴリズム) が挙げられる。さらに、リソース・マネージャ 1950 は、メトリクスを集約して、どのリソース (メモリ/ディスク/GPU/FPGA) が圧迫下にあるかを検出し、計算ノード 1904 - 1、1904 - 2、1904 - 3 のうちの 1 つまたは複数へのさらなるリソースの追加を要求する。

30

【0144】

図 20 は、図 19 に示すデータ管理プラットフォーム 100 におけるラック 106 内の圧迫の検出にตอบสนองして論理リソースを自動的に追加または除去する方法の流れ図である。

【0145】

一般に、圧迫検出は複数の入力源に依存する。圧迫検出は、事後 (post-factum) に発生した事象、例えば、50 パーセンタイル、99 パーセンタイル、または 99.9 パーセンタイルのレイテンシ・スパイク (ワークロード・メトリクス) に基づくことができる。圧迫検出は、事前 (pre-factum)、すなわち、計算ノード 1904 - 1、1904 - 2、1904 - 3 およびワークロード 1708 - 1、...、1708 - 9 におけるリソース利用率の増加に基づいて、50 パーセンタイル、99 パーセンタイル、または 99.9 パーセンタイルのレイテンシ・スパイクが検出される前に基づくこともできる。

40

【0146】

圧迫検出が事後である場合、オーケストレータ・メトリクス 1911 およびノード・メトリクス 1910 を用いて、計算ノード 1904 - 1、1904 - 2、1904 - 3 および圧迫検出に関連付けられたリソースを検出する。圧迫検出が事前である場合、ノード 1

50

904-1、1904-2、1904-3およびワークロード1708-1、...、1708-9におけるリソース利用率の増加に基づいて圧迫が存在するという予測が行われる。

【0147】

ブロック2000において、リソース・マネージャ1950はシステム・メトリクスを監視する。監視されるシステム・メトリクスには、オーケストレータ・メトリクス1911、ノード・メトリクス1910、およびワークロード・メトリクス1912が含まれる。

【0148】

ブロック2002において、計算ノード1904-1、1904-2、1904-3がストレス下にある場合、圧迫状態が発生する。システム・メトリクスを監視しながら、リソース・マネージャ1950は、圧迫状態のステータスがアクティブであるか非アクティブであるかを検出することができる。圧迫状態は、圧迫状態が発生しようとしている場合、圧迫状態が終了しようとしている場合、または圧迫状態が進行中である場合にアクティブである。リソース・マネージャ1950によってアクティブな圧迫状態が検出された場合、処理はブロック2004に進む。そうでない場合は、処理はブロック2000に進み、メトリクスを監視し続ける。

10

【0149】

ブロック2004では、アクティブな圧迫状態が検出される。圧迫状態の例は、99パーセントایل・レイテンシまたは99.9パーセントایل・レイテンシ・スパイク（ワークロード・メトリクス）である。アクティブな圧迫状態によって影響を受ける計算ノード1904-1、1904-2、1904-3上で実行されるアプリケーションが決定される。処理はブロック2006に進む。

20

【0150】

ブロック2006において、アクティブな圧迫状態によって影響を受ける計算ノード1904-1、1904-2、1904-3が決定される。

【0151】

ブロック2008において、検出されたアクティブな圧迫状態は、圧迫状態が発生しようとしているか、終了しようとしているか、または進行中であるかであり得る。圧迫状態が発生しようとしているか、または進行中である場合、処理はブロック2012に進む。圧迫状態が終了しようとしている場合、処理はブロック2010に進む。

30

【0152】

ブロック2010において、圧迫状態が終了しようとしている場合、論理リソースが計算ノード1904-1、1904-2、1904-3から除去される。処理はブロック2000に進み、メトリクスを監視し続ける。

【0153】

ブロック2012において、圧迫状態が発生しようとしているか、または進行中である場合、論理リソースが計算ノード1904-1、1904-2、1904-3に追加される。一実施形態では、圧迫状態を経験している計算ノード1904-1、1904-2、1904-3により多くの論理リソースが追加される。追加の論理リソースは、計算ノード1904-1、1904-2、1904-3上のすべてのワークロード1708によって使用することができ、または計算ノード1904-1、1904-2、1904-3上の特定のワークロード1708によってのみ使用されるように制限することができる。圧迫を緩和するために追加することができる論理リソースには、ストレージ、メモリ、アクセラレータ、およびフィールド・プログラマブル・ゲート・アレイ（FPGA）リソースが含まれる。

40

【0154】

ディスク（例えば、ソリッド・ステート・ドライブ404（図4））に対する圧迫状態は、ディスクの容量不足またはディスクへの入力/出力レイテンシの増加によるものであり得る。ディスクの（事後または事前の）圧迫状態を検出すると、リソース・マネージャ

50

1905は、アクセラレータ・ノード1902に新しいボリュームを作成し、新しく作成されたボリュームをそれぞれの計算ノード1904-1、1904-2、1904-3に論理的にアタッチするよう要求する。それぞれの計算ノード1904-1、1904-2、1904-3のリソース・マネージャ・エージェント1706は、新たに作成されたボリュームに対してファイルシステム拡張を実行し、計算ノード1904-1、1904-2、1904-3のうちの1つに実行中のワークロード1708のために新たに作成されたボリュームを直接マウントする。

【0155】

メモリ（例えば、メモリ630または永続メモリ1702（図17））に対する圧迫状態は、高いメモリ帯域幅使用量、計算ノード1904-1、1904-2、1904-3の低い空きメモリ、または計算ノード1904-1、1904-2、1904-3のワークロード1708のメモリ使用量のスパイクに起因し得る。リソース・マネージャ1950は、永続メモリ1702（図17）、シンプル・ストレージ・サービス（S3）エンドポイント、またはストレージ・ノード400内のリモート・ソリッド・ステート・ドライブ404を使用して、計算ノード1904-1、1904-2、1904-3に新しいメモリのプールを割り当てることができる。シンプル・ストレージ・サービスは、計算ノード1904-1、1904-2、1904-3のネットワーク・インターフェース・コントローラ302を介してアクセスすることができる。シンプル・ストレージ・サービスは、ネットワーク・インターフェース・コントローラ302を介してアクセス可能な複数のリモートドライブを使用して、1つのエンドポイントを提示する。シンプル・ストレージ・サービスは、ウェブ・サービス・インターフェースを介してオブジェクトストレージを提供するAPI（Application Programming Interface）である。Amazon（登録商標）シンプル・ストレージ・サービスは、オブジェクト・ストレージ・ソリューションにおけるデファクト・スタンダードである。Amazonシンプル・ストレージ・サービスと互換性のあるインターフェースの例には、Ceph RADOS Gateway、OpenIO、Scality、およびMinIOが含まれる。一実施形態では、シンプル・ストレージ・サービスは、MinIOによって提供される。リモート・ソリッド・ステート・ドライブ404は、計算ノード1904-1、1904-2、1904-3およびワークロード1708と同じラック106内のストレージ・ノード400にある。

【0156】

新しいメモリのプールは、リソース・マネージャ・エージェント1706を介して計算ノード1904-1、1904-2、1904-3が利用するためにアクセス可能である。リソース・マネージャ・エージェント1706は、新しいメモリのプールを、計算ノード1904-1、1904-2、1904-3に既に割り当てられているメモリの拡張としてマッピングする。永続メモリ1702に割り当てられた新しいメモリのプールは、ワークロード1708によって直接アクセス可能である。シンプル・ストレージ・サービスによって割り当てられた新しいメモリのプールは、ユーザ空間500からのオンデマンド・ページングを可能にするカーネル関数（例えば、「userfaultfd」関数）を介してワークロード1708に公開される。

【0157】

永続メモリ1702またはリモート・ソリッド・ステート・ドライブ404内の新たに割り当てられたメモリのプールは、論理メモリのウォーム階層（warm tier）として使用され、メモリ630は、揮発性メモリを含み、論理メモリのウォーム階層のためのキャッシュである。ローカルメモリ630は、低レイテンシおよび高帯域幅を有する論理メモリのホット階層である。永続メモリ1702は、メモリ630よりも容量が大きく、レイテンシが大きく、帯域幅が狭い。ソリッド・ステート・ドライブ302は、永続メモリ1702よりも容量が大きく、レイテンシが大きく、帯域幅が狭い。

【0158】

アクセラレータ1802またはFPGA1800に対する圧迫状態は、アクセラレータ

10

20

30

40

50

1802またはFPGA1800の十分なリソースの欠如によってワークロード1708が影響を受ける結果となる。1つまたは複数のワークロード1708によって使用されるアクセラレータ1802またはFPGA1800のリソースのためのジョブキューは、圧迫状態中に満たすことができる。

【0159】

一実施形態では、リソース・マネージャ1850は、CPU処理のストールに起因するパーセンタイル・レイテンシの増加を検出する。例えば、要求されたデータがCPUキャッシュにない場合、要求されたデータは、遠くのメモリまたはストレージからフェッチされなければならない。その結果、平均応答時間のばらつき（すなわち、平均からの偏差）が生じる。レイテンシ増加の検出に反応して、リソース・マネージャ1850は、RDMA (Remote Direct Memory Access) ベースの通信プロトコルを使用して、ストレージ・ノード404から新しいアクセラレータ・サーバまたはFPGAリソースを要求する。

10

【0160】

RDMAベースの通信プロトコルの例には、NVMeOF (NVMe over Fabric) またはFPGAリソース over Fabric (例えば、RDMA over Fabric) を使用してアクセス可能なFPGA) が含まれる。NVMe over Fabricは、ストレージ・ネットワーキング・ファブリック上でNVMeブロック・ストレージ・プロトコルのための一連のストレージ・ネットワーキング・ファブリックをサポートする共通のアーキテクチャを定義する。これには、ストレージ・システムへのフロントサイド・インターフェースを可能にすること、多数のNVMeデバイスにスケールアウトすること、およびNVMeデバイスおよびNVMeサブシステムにアクセスできるデータセンタ内の距離を拡大することが含まれる。

20

【0161】

新しい論理アクセラレータまたはFPGAは、計算ノード1904-1、1904-2、1904-3のうちの一つに接続され、ワークロード1708によって使用される。処理はブロック2000に進み、メトリクスを監視し続ける。

【0162】

ストレージレイヤの自己修復

シェアードナッシング・アーキテクチャ (SN) は、更新要求が単一のノードによって満たされる分散コンピューティング・アーキテクチャである。ノードは、計算ノード、メモリノードまたはストレージ・ノードとすることができる。その意図は、ノード間の競合を排除することである。各ノードは、メモリおよびストレージに独立してアクセスする。ノードは、メモリまたはストレージを共有しない。

30

【0163】

シェアードナッシング・アーキテクチャ・システムは、システムのボトルネックとなる中心的なリソースがないため、ノードを追加することによってスケールアップすることができる。シェアードナッシング・アーキテクチャの別の用語は、シャーディング (sharding) である。データベース・シャードは、データベースまたは検索エンジンのデータの水平パーティションである。個々のパーティションは、シャードまたはデータベース・シャードと呼ばれる。各シャードは、負荷を分散するために、別々のデータベース・サーバ・インスタンスに記憶される。

40

【0164】

シェアードナッシング・アーキテクチャを有するシステムで使用される分散アプリケーションは、データベース・サーバ・インスタンスに記憶されたそれらのシャードが永続的である必要がある。シェアードナッシング・アーキテクチャを有する分散アプリケーションの例には、ストラクチャード・クエリ・ランゲージ (SQL) データベース、シンプル・ストレージ・サービス (S3) オブジェクトストア、および時系列データベースが含まれる。ストラクチャード・クエリ・ランゲージは、プログラミングで使用されるドメイン固有言語であり、リレーショナル・データベース管理システム (RDBMS) に記憶され

50

たデータを管理するために、またはリレーショナル・データ・ストリーム管理システム (R D S M S) におけるストリーム処理のために設計されている。

【 0 1 6 5 】

データベース・サーバ・インスタンスまたはデータベース・サーバ・インスタンスのグループの障害は、データベース・サーバのユーザに影響を与える。障害は、データベース・サーバ・インスタンスに記憶されたデータに対するリクエストのレイテンシを増加させ、またはデータベース・サーバ・インスタンスに記憶されたデータに対するリクエストを失敗させることになる可能性がある。さらに、障害からの復旧は、障害が発生したデータベース・サーバ・インスタンスまたはデータベース・サーバ・インスタンスのグループに関連付けられたデータを復元しなければならないため、時間およびリソースがかかる。

10

【 0 1 6 6 】

アプリケーション・トポロジおよびアプリケーションがデプロイされている障害ドメインの知識を用いてデータベース・サーバ・インスタンスの復元を実行するためには、手動による介入が必要である。スケールアウト・アーキテクチャでは、アプリケーションはいくつかのプロセスから構成され、それぞれが K u b e r n e t e s のポッドで実行される。これらのポッドは、障害がアプリケーションの可用性またはアプリケーションが管理しているデータの耐久性に影響を与えないように、障害ドメイン、すなわちデータ管理プラットフォーム 1 0 0 内のラックにわたって分散されている。これらのポッドの分布がアプリケーション・トポロジである。

【 0 1 6 7 】

一実施形態では、ストレージ自己修復論理または回路とも呼ばれ得るストレージ自己修復メカニズムは、ストレージ・サブシステムを監視し、ストレージ・サブシステム (ストレージ・ノードおよびソリッド・ステート・ドライブ) を使用するワークロードを監視して、すべてのワークロードが利用可能な障害ドメインにわたって確実に分散されるようにする。

20

【 0 1 6 8 】

図 2 1 は、ストレージ自己修復メカニズム 2 1 0 8 を含むデータ管理プラットフォーム 1 0 0 における物理クラスタ 2 1 0 0 の一実施形態のブロック図である。物理クラスタ 2 1 0 0 は、オーケストレータ / スケジューラ 1 0 2 およびラック 1 0 6 を含む。一実施形態では、ストレージ自己修復メカニズム 2 1 0 8 は、オーケストレータ / スケジューラ 1 0 2 にある。他の実施形態では、ストレージ自己修復メカニズムは、データ管理プラットフォーム 1 0 0 の別の構成要素にあってもよい。

30

【 0 1 6 9 】

図 2 1 に示す特定の非限定的な例では、1つのデータスイッチ 2 0 6、3つの計算ノード 2 1 1 0 a ~ c、および2つのストレージ・ノード 2 1 0 2 a ~ b を有する1つのラック 1 0 6 がある。計算ノード 2 1 1 0 a ~ b およびストレージ・ノード 2 1 0 2 a ~ b は、データスイッチ 2 0 6 に通信可能に結合されている。

【 0 1 7 0 】

物理クラスタ 2 1 0 0 は、ストレージ・サブシステム (ストレージ・ノード 2 1 0 2 a ~ b およびソリッド・ステート・ドライブ 2 1 0 6 a ~ d) を使用するアプリケーションのための複数の障害ドメインを有する。第1の障害ドメインは、データスイッチ 2 0 6 であり、第2の障害ドメインは、ワークロード 2 1 0 4 a ~ c が実行される計算ノード 2 1 0 0 0 a ~ b であり、第3の障害ドメインは、ストレージ・ノード 2 1 0 2 a ~ b であり、第4の障害ドメインは、ソリッド・ステート・ドライブ 2 1 0 6 a ~ b である。

40

【 0 1 7 1 】

ストレージ自己修復メカニズム 2 1 0 8 は、複数の障害ドメインのそれぞれについてヘルスチェックを定期的に行う。データスイッチ 2 0 6 のためのストレージ自己修復メカニズムによって実行されるヘルスチェックの1つは、オーケストレータ / スケジューラ 1 0 2 がデータスイッチ 2 0 6 にアクセスできるかどうかを判定することである。一実施形態では、インターネット制御通知プロトコル (I C M P) を使用して、オーケストレー

50

タ102がデータスイッチにアクセスできるかどうかを判定することができる。例えば、ICMPエコー要求/応答または拡張エコー要求/応答メッセージを使用して、オーケストレータ/スケジューラ102がデータスイッチ206にアクセスできるかどうかを判定することができる。インターネット制御通知プロトコルは、エラー報告プロトコルであり、リクエスト・フォー・コメント(RFC)792によって定義されたインターネット・プロトコル(IP)の拡張である。

【0172】

データスイッチ206に対してストレージ自己修復メカニズム2108によって実行される別のヘルスチェックは、ルートがデータスイッチにおいて利用可能かどうかを判定することである。一実施形態では、「ip r g」コマンド(Linux(登録商標)ユーティリティコマンド)は、指定されたIPアドレスにバインドされている計算ノード2100a~cまたはストレージ・ノード2102a~bへのルートが見えるかどうか、およびポッドへのルートが見えるかどうかをチェックする。

10

【0173】

ストレージ自己修復メカニズム2108によって計算ノード2100a~cまたはストレージ・ノード2102a~bに対して実行されるヘルスチェックの1つは、計算ノード2110a~cまたはストレージ・ノード2102a~bがデータスイッチ206を介して到達可能であるかどうかを判定することである。一実施形態では、インターネット制御メッセージプロトコル(ICMP)を使用して、計算ノード2110a~cまたはストレージ・ノード2102a~bがデータスイッチにアクセスできるかどうかを判定することができる。

20

【0174】

ストレージ自己修復メカニズム2108によって実行される計算ノード2110a~cまたはストレージ・ノード2102a~bに対する別のヘルスチェックは、オーケストレータ102が計算ノード2110a~cまたはストレージ・ノード2102a~bを準備完了として報告するかどうかをチェックすることである。オーケストレータ102は、計算ノード2110a~cまたはストレージ・ノード2102a~bの健全性をチェックし、それぞれの計算ノード2110a~cまたはストレージ・ノードnがワークロードを受け入れる準備ができているかどうかを報告する。計算ノード2110a~cまたはストレージ・ノード2102a~bの健全性は、オペレーティング・システムの観点からのメモリおよびCPUチェック、ならびにオーケストレータとそれぞれの計算ノード2110a~cまたはストレージ・ノード2102a~bとの間のネットワーク接続性を含むことができる。

30

【0175】

ストレージ自己修復メカニズム2108によって実行されるソリッド・ステート・ドライブ2106a~dに対して実行されるヘルスチェックの1つは、ライト・アンプリフィケーション(write amplification)である。例えば、空きユーザ領域およびオーバープロビジョニングなどのライト・アンプリフィケーション率を使用して、ソリッド・ステート・ドライブ2106a~dにおける障害を予測することができる。ストレージ自己修復メカニズム2108によって実行されるソリッド・ステート・ドライブ2106a~dに対する別のヘルスチェックは、ソリッド・ステート・ドライブ2106a~dの健全性をチェックすることである。

40

【0176】

ソリッド・ステート・ドライブ2106a~dの健全性は、自己監視・分析・レポート技術(S.M.A.R.T.)を使用して監視することができる。S.M.A.R.Tとは、ソリッド・ステート・ドライブに含まれる監視システムであり、ソリッド・ステート・ドライブの信頼性の指標を監視して報告し、データ損失を防止するための予防措置を講じるために使用することができる。NANDベースのソリッド・ステート・ドライブ2106a~dに対するS.M.A.R.Tメトリクスの例には、プログラム失敗回数、消去失敗回数、ウェア・レベリング回数、エンド・ツー・エンド・エラー検出回数、巡回冗

50

長コード (CRC) エラー回数、時限ワークロード - 媒体摩耗、時限ワークロード - ホスト、読み取り / 書き込み比率、時限ワークロード・タイマ、熱スロットル・ステータス、再試行バッファ・オーバーフロー回数、PLI ロック損失回数、NAND バイト・ライト、ホスト・バイト・ライト、およびシステム・エリア・ライフ残量が含まれる。

【0177】

スケールアウトされたシェードナッシング・アーキテクチャでは、ワークロードは、複数のインスタンスを有する。受け入れられる失敗インスタンスの最小数は、ワークロード固有である。シンプル・ストレージ・サービスは、消去訂正符号およびチェックサムを使用して、ハードウェア障害およびサイレントデータ破損からデータを保護する。

【0178】

消去訂正符号は、失われたまたは破損したデータを再構築するための数学的アルゴリズムである。シンプル・ストレージ・サービスは、オブジェクトをデータブロックおよびパリティブロックに分割する。8個のデータブロックおよび4個のパリティブロックにより、最大4個のインスタンス障害が存在する場合にデータを回復することができる。3つのレプリカ (データベースのコピー) を有するデータベースでは、最大2つのインスタンスがデータを回復することができる。

【0179】

自己修復メカニズム 2108 は、入ってくる障害を検出することができ、故障したソリッド・ステート・ドライブ 2106 a ~ d に記憶されたデータの再生成をトリガし、スケジューリングすることができる。ストレージ自己修復メカニズム 2108 は、前述したヘルスチェックを介して得られたメトリクスを使用して、ストレージ・サブシステム (ストレージ・ノード 2102 a ~ b およびソリッド・ステート・ドライブ 2106 a ~ d) に対する自動アクションをトリガすることができる。

【0180】

図 22 は、図 21 に示すクラスタ 2100 におけるワークロードのマッピングの実施形態を示す。ストレージ自己修復メカニズム 2108 は、ストレージ・サブシステムへのワークロード 2104 a ~ c のマッピングを追跡する。例えば、マッピングは、論理ボリュームがマッピングされるソリッド・ステート・ドライブ 2106 a ~ d を追跡する。

【0181】

図 22 に示すように、ワークロード A (インスタンス 1) は、ラック 1 内の計算 1、アクセラレータ 1、ソリッド・ステート・ドライブ 1、ボリューム a にマッピングされる。

【0182】

ストレージ自己修復メカニズム 2108 を含むシステムでは、アプリケーションの復旧に必要な管理者 / オペレータの関与はない。また、クラスタ化されたアプリケーション・インスタンスの故障がシステム全体の性能に与えるレイテンシおよび帯域幅の影響も低減される。

【0183】

データスイッチ 206 の障害、またはラック 106 内のストレージ・ノード 2102 a ~ b および計算ノード 2110 a ~ c のすべてに関連する障害を検出すると、複数のワークロードが影響を受ける。ストレージ自己修復メカニズム 2108 は、障害によって影響を受けるワークロードを検出する。物理クラスタ 2100 内の利用可能なラックから別のラック 106 が選択される。オーケストレータ 102 内のストレージ自己修復メカニズム 2108 は、NVMe over Fabric インターフェースを介して、ソリッド・ステート・ドライブ上のボリュームを切断し、故障した計算ノードおよびストレージ・ノードからボリュームを除去する。

【0184】

故障したラックの計算ノード上で実行されていたワークロード用のリソースが、他のラックに作成される。ワークロードは、他のラックの計算ノード上で実行されるように再スケジューリングされる。再スケジューリングされたワークロード・インスタンスが再スケジューリングされた後、オーケストレータ 102 内のストレージ自己修復メカニズム 2108 は、スト

10

20

30

40

50

レージが他のラックに再作成された後に、他のラックのボリューム上のすべてのデータについて、ワークロードのワークロード「修復/回復」メカニズムをトリガし、それぞれのブロックにおけるエラーを検出するとブロックを修復する。

【0185】

ストレージ自己修復メカニズム2108が計算ノードの障害を検出するか、または計算ノードが障害を報告した場合、ストレージ自己修復メカニズム2108は、どのワークロードが故障した計算ノードで影響を受けたかを検出する。影響を受けたワークロードは、そのワークロードのインスタンスをまだホストしていない、同じラック内の別の計算ノード上で実行されるように再スケジュールされる。オーケストレータ102のストレージ自己修復メカニズム2108は、故障した計算ノードからのボリュームの切断と、他の計算ノードへのボリュームの接続とを要求する。ソリッド・ステート・ドライブとストレージ・ノードとの間の通信パスがNVMeOFを介する実施形態では、切断および接続要求は、NVMeOFインターフェースを介して送信される。ワークロード・インスタンスが他の計算ノード上で再開された後、オーケストレータ102は、ボリューム上のデータのすべてに対してワークロード「修復/回復」メカニズムをトリガし、エラー時に任意のブロックを修復する。

10

【0186】

複数のワークロードおよびソリッド・ステート・ドライブ上のデータに影響を与えるストレージ・ノード2102a~bの障害が回復できない場合、オーケストレータ102内のストレージ自己修復メカニズム2108は、影響を受けたワークロードおよび故障したストレージ・ノードで使用されるボリュームを決定する。ストレージ自己修復メカニズム2108は、影響を受けたすべてのワークロードをラック102内の異なる計算ノード2110a~c上に再スケジュールし、別のストレージ・ノード2100a~bのソリッド・ステート・ドライブ上に新しいボリュームを作成し、NVMeOFインターフェースを介してソリッド・ステート・ドライブ上の新しいボリュームを新しい計算ノード2110a~cに接続する。

20

【0187】

ラック106内の計算ノード2110a~cが既に同じタイプのワークロードをホストしている場合、オーケストレータ102は、物理クラスタ2100内の別のラック106を選択し、オーケストレータ102は、影響を受けたすべてのワークロードを、別のラック102内の計算ノード2100a~c上で実行するように再スケジュールする。オーケストレータ102は、他のラック106内の別のストレージ・ノード2100a~bのソリッド・ステート・ドライブ上に新しいボリュームを作成し、NVMeOFインターフェースを介して他のラック106内の新しい計算ノード2110a~cに接続する。ワークロード・インスタンスが他の計算ノード上で再開された後、オーケストレータ102は、ボリューム上のすべてのデータに対してワークロード「修復/回復」メカニズムをトリガし、エラー時に任意のブロックを修復する。

30

【0188】

ストレージ・ノード2102a~dの1つまたは複数のソリッド・ステート・ドライブが故障すると、複数のワークロードが影響を受ける。ストレージ自己修復メカニズム2108は、影響を受けたワークロード、すなわち、故障したソリッド・ステート・ドライブ上の論理ボリュームを使用しているワークロードを決定する。ストレージ自己修復メカニズム2108は、同じストレージ・ノード2102a~b内、または同じラック106内の別のストレージ・ノード2102a~b内の他の動作可能なソリッド・ステート・ドライブ上に新しいボリュームを作成する。新しいボリュームは、NVMeOFインターフェースを介して計算ノードに接続され、古いボリュームは切断される。ワークロード・インスタンスが他の計算ノード上で再開された後、オーケストレータ102は、ボリューム上のすべてのデータに対してワークロード「修復/回復」メカニズムをトリガし、エラー時に任意のブロックを修復する。

40

【0189】

50

本明細書に示す流れ図は、様々なプロセスアクションのシーケンスの例を提供する。流れ図は、ソフトウェアまたはファームウェアルーチンによって実行される動作、ならびに物理的動作を示すことができる。一実施形態では、流れ図は、ハードウェアおよび/またはソフトウェアで実装することができる有限状態機械(FSM)の状態を示すことができる。特定のシーケンスまたは順序で示されているが、別段の指定がない限り、アクションの順序は変更することができる。したがって、図示された実施形態は例として理解されるべきであり、プロセスは異なる順序で実行されてもよく、いくつかのアクションは並行して実行されてもよい。さらに、様々な実施形態において1つまたは複数のアクションを省略することができ、したがって、すべての実施形態においてすべてのアクションが必要なわけではない。他のプロセスフローも可能である。

10

【0190】

様々な動作または機能が本明細書で説明される限りにおいて、それらは、ソフトウェアコード、命令、設定、および/またはデータとして説明または定義することができる。コンテンツは、直接実行可能なもの(「オブジェクト」または「実行可能」形式)、ソースコード、または差分コード(「デルタ」または「パッチ」コード)とすることができる。本明細書で説明する実施形態のソフトウェアコンテンツは、コンテンツが格納された製造物を介して、または通信インターフェースを操作して通信インターフェースを介してデータを送信する方法を介して提供することができる。機械可読記憶媒体は、説明した機能または動作を機械に実行させることができ、記録可能/非記録可能媒体(例えば、読み出し専用メモリ(ROM)、ランダム・アクセス・メモリ(RAM)、磁気ディスク記憶媒体、光記憶媒体、フラッシュメモリデバイスなど)などの、機械(例えば、コンピューティングデバイス、電子システムなど)によってアクセス可能な形態で情報を記憶する任意の機構を含む。通信インターフェースには、ハードワイヤード媒体、ワイヤレス媒体、光媒体などのいずれかとインターフェースして別のデバイスと通信する任意の機構、例えば、メモリ・バス・インターフェース、プロセッサ・バス・インターフェース、インターネット接続、ディスクコントローラなどが含まれる。通信インターフェースは、ソフトウェアコンテンツを記述するデータ信号を提供するための通信インターフェースを準備するために、設定パラメータを提供し、および/または信号を送信することによって設定することができる。通信インターフェースは、通信インターフェースに送信される1つまたは複数のコマンドまたは信号を介してアクセスすることができる。

20

30

【0191】

本明細書で説明する様々な構成要素は、説明する動作または機能を実行するための手段とすることができる。本明細書で説明される各構成要素は、ソフトウェア、ハードウェア、またはこれらの組合せを含む。構成要素は、ソフトウェアモジュール、ハードウェアモジュール、専用ハードウェア(例えば、特定用途向けハードウェア、特定用途向け集積回路(ASIC)、デジタル信号プロセッサ(DSP)など)、組込みコントローラ、ハードワイヤード回路などとして実装され得る。

【0192】

本明細書で説明されるものに加えて、本発明の範囲から逸脱することなく、本発明の開示された実施形態および実施態様に対して様々な修正を行うことができる。

40

【0193】

したがって、本明細書における例示および実施例は、限定的な意味ではなく、例示的な意味で解釈されるべきである。本発明の範囲は、以下の特許請求の範囲を参照することによってのみ判断されるべきである。

【0194】

一般に、本明細書の説明に関して、一例では、装置は、計算サーバと、ストレージ・サーバに通信可能に結合された複数のストレージ・デバイスを管理するストレージ・サーバと、を含み、計算サーバおよびストレージ・サーバは、ネットワークを介して通信可能に結合され、ストレージ・サーバによって管理される複数のストレージ・デバイスは、複数のストレージ・デバイスのストレージ容量を計算サーバから独立してスケールリングするこ

50

とができるように、計算サーバから分離されている。

【0195】

一例では、ネットワーク・インターフェース・コントローラは、ネットワークに通信可能に結合され、システム・オン・チップは、複数のコアおよびラスト・レベル・メモリを備え、複数のコアは、ラスト・レベル・キャッシュ・メモリに通信可能に結合され、ラスト・レベル・キャッシュ・メモリは、複数のキャッシュ・ウェイを備え、複数のキャッシュ・ウェイの一部は、複数のストレージ・デバイス内の論理ボリュームおよびネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、論理ボリュームとネットワーク・インターフェース・コントローラとの間でデータを転送する。

【0196】

一例では、複数のキャッシュ・ウェイの一部は、初期化中に割り当てられる。

【0197】

一例では、論理ボリュームは、計算サーバによる使用のためのデータを記憶する。

【0198】

一例では、論理ボリュームとラスト・レベル・キャッシュ内の複数のキャッシュ・ウェイとの間で転送されるデータは、ネットワーク・インターフェース・コントローラに通信可能に結合されたネットワークを介してストレージ・サーバと計算サーバとの間で転送される。

【0199】

一例では、複数のコアのうちの少なくとも1つは、論理ボリュームとネットワーク・インターフェース・コントローラとの間でデータを転送するために、複数のストレージ・デバイス内の論理ボリュームおよびネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられる。

【0200】

一例では、システム・オン・チップに結合された外部メモリは、転送されるデータを記憶するために、ラスト・レベル・キャッシュの一部における複数のキャッシュ・ウェイのすべてが論理ボリュームおよびネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられると、論理ボリュームとネットワーク・インターフェースとの間で転送されるデータを一時的に記憶する。

【0201】

一般に、本明細書の説明に関して、一例では、装置は、それぞれがサーバにおけるメトリクスを監視する複数のサーバと、複数のサーバが通信可能に結合されたデータスイッチであって、複数のサーバにおけるサービスへの複数の動的ルートを含むルートテーブルを備える、データスイッチと、監視されたメトリクスに基づいて複数のサーバのうちの1つにおけるサービスへのルートを動的に許可または抑制するフィルタリングシステムと、を含む。

【0202】

一例では、サービスへのルートが、ルートを抑制するためにルートテーブルから除去される。

【0203】

一例では、サービスへのルートが、ルートを許可するためにルートテーブルに追加される。

【0204】

一例では、監視されるメトリクスは、複数のサーバにおける圧迫状態または性能劣化に関連するメトリクスを含む。

【0205】

一例では、圧迫状態に関連するメトリクスは、ワークロード・メトリクスを含む。

【0206】

一例では、性能劣化に関連するメトリクスは、メモリ帯域幅に関連するメトリクスを含む。

10

20

30

40

50

【0207】

一例では、監視されるメトリクスは、複数のサーバ内の劣化ハードウェア・コンポーネントに関連するメトリクスを含む。

【0208】

一例では、劣化ハードウェア・コンポーネントに関連するメトリクスは、ノード・メトリクス、オーケストレータ・メトリクスおよびワークロード・メトリクスを含む。

【0209】

一例では、サーバはストレージ・サーバであり、劣化ハードウェア・コンポーネントに関連するメトリクスは、ソリッド・ステート・ドライブ・メトリクスを含む。

【0210】

一般に、本明細書の説明に関して、一例では、方法は、複数のサーバにおけるメトリクスを監視するステップと、データスイッチのルートテーブルに、複数のサーバにおけるサービスへの複数の動的ルートを記憶するステップと、監視されたメトリクスに基づいて、複数のサーバのうちの1つにおけるサービスへのルートを動的に許可または抑制するステップと、を含む。

10

【0211】

一例では、サービスへのルートが、ルートを抑制するためにルートテーブルから除去される。

【0212】

一例では、サービスへのルートが、ルートを許可するためにルートテーブルに追加される。

20

【0213】

一例では、監視されるメトリクスは、複数のサーバにおける圧迫状態または性能劣化に関連するメトリクスを含む。

【0214】

一例では、圧迫状態に関連するメトリクスは、ワークロード・メトリクスを含む。

【0215】

一例では、性能劣化に関連するメトリクスは、メモリ帯域幅に関連するメトリクスを含む。

【0216】

一例では、監視されるメトリクスは、複数のサーバ内の劣化ハードウェア・コンポーネントに関連するメトリクスを含む。

30

【0217】

一般に、本明細書の説明に関して、一例では、データ管理プラットフォームは、それぞれがサーバにおけるメトリクスを監視する複数のサーバと、複数のサーバが通信可能に結合されたデータスイッチであって、複数のサーバにおけるサービスへの複数の動的ルートを含むルートテーブルを備える、データスイッチと、監視されたメトリクスに基づいて、複数のサーバのうちの1つにおけるサービスへのルートを動的に許可または抑制するフィルタリングシステムと、を備えるラック、を含む。

【0218】

一例では、サービスへのルートが、ルートを抑制するためにルートテーブルから除去される。

40

【0219】

一例では、サービスへのルートが、ルートを許可するためにルートテーブルに追加される。

【0220】

一例では、監視されるメトリクスは、複数のサーバにおける圧迫状態または性能劣化に関連するメトリクスを含む。

【0221】

一般に、本明細書の説明に関して、一例では、装置は、それぞれが論理リソースを使用

50

してワークロードを実行する複数の計算ノードと、複数の計算ノードに通信可能に結合されたオーケストレータであって、メトリクスを監視して計算ノードにおけるアクティブな圧迫状態を検出し、アクティブな圧迫状態の検出に応答して、圧迫状態がアクティブである間に新しい論理リソースを計算ノードにアタッチする、オーケストレータと、を含む。

【0222】

一例では、メトリクスは、オーケストレータ・メトリクス、ノード・メトリクスおよびワークロード・メトリクスを含む。

【0223】

一例では、ノード・メトリクスは、計算ノードごとのCPU利用率および計算ノードごとのメモリ利用率を含む。

【0224】

一例では、ワークロード・メトリクスは、クライアントの数、平均応答レイテンシ、およびパーセンタイル・メトリクスを含む。

【0225】

一例では、オーケストレータは、圧迫状態がアクティブでないときに、新しい論理リソースを計算ノードからデタッチする。

【0226】

一例では、新しいリソースは、圧迫状態がアクティブである間、計算ノード上の特定のワークロードによってのみ使用が制限される。

【0227】

一例では、新しいリソースは、圧迫状態がアクティブである間、計算ノード上のすべてのワークロードによって使用される。

【0228】

一般に、本明細書の説明に関して、一例では、装置は、それぞれが論理リソースを使用してワークロードを実行する複数の計算ノードと、複数のストレージ・デバイスに通信可能に結合された複数のストレージ・ノードを備えるストレージ・サブシステムと、複数の計算ノードおよびストレージ・サブシステムに通信可能に結合されたストレージ自己修復メカニズムであって、ストレージ・サブシステム上でヘルスチェックを実行し、ヘルスチェックを介して得られたメトリクスを使用してストレージ・サブシステム内の障害を検出し、ストレージ・サブシステム内の障害の修復を管理する、ストレージ自己修復メカニズムと、を備える、ラックを備えるデータ管理プラットフォームを含む。

【0229】

一例では、障害は計算ノードにあり、ストレージ自己修復メカニズムは、故障した計算ノード上で実行されるワークロードをラック内の別の計算ノードに再スケジュールする。

【0230】

一例では、障害は、計算ノードにあり、ラック内の他の計算ノードは、同じタイプのワークロードをホストしており、ストレージ自己修復メカニズムは、ワークロードを再スケジュールして、別のラックの別の計算ノード上で実行させる。

【0231】

一例では、障害はストレージ・デバイスにあり、ストレージ自己修復メカニズムは、ストレージ・ノード内の別のストレージ・デバイス上にボリュームを作成するために、ストレージ・デバイスに関連付けられたワークロードを決定する。

【0232】

一例では、ラックはデータスイッチをさらに備え、障害はデータスイッチにあり、ストレージ自己修復メカニズムは、ラックにおけるワークロードを決定し、別のラックを選択し、別のラックにおけるワークロードを再スケジュールする。

[他の可能な項目]

[項目1]

それぞれがサーバにおけるメトリクスを監視する複数のサーバと、

前記複数のサーバに通信可能に結合されたデータスイッチであって、前記複数のサーバ

10

20

30

40

50

におけるサービスへの複数の動的ルートを含むルートテーブルを有する、データスイッチと、

前記監視されたメトリクスに基づいて、前記複数のサーバのうちの1つにおける前記サービスへのルートを動的に許可または抑制するフィルタリングシステムと、
を備える装置。

[項目 2]

前記サービスへの前記ルートが前記ルートを抑制するために前記ルートテーブルから除去される、項目 1 に記載の装置。

[項目 3]

前記サービスへの前記ルートが前記ルートを許可するために前記ルートテーブルに追加される、項目 1 に記載の装置。

10

[項目 4]

前記監視されるメトリクスが前記複数のサーバにおける圧迫状態または性能劣化に関連するメトリクスを含む、項目 1 に記載の装置。

[項目 5]

前記圧迫状態に関連する前記メトリクスがワークロード・メトリクスを含む、項目 4 に記載の装置。

[項目 6]

前記性能劣化に関連する前記メトリクスがメモリ帯域幅に関連するメトリクスを含む、項目 4 に記載の装置。

20

[項目 7]

前記監視されるメトリクスが前記複数のサーバ内の劣化ハードウェア・コンポーネントに関連するメトリクスを含む、項目 1 に記載の装置。

[項目 8]

前記劣化ハードウェア・コンポーネントに関連する前記メトリクスがノード・メトリクス、オーケストレータ・メトリクス、およびワークロード・メトリクスを含む、項目 7 に記載の装置。

[項目 9]

前記複数のサーバのうちの1つがストレージ・サーバであり、前記劣化ハードウェア・コンポーネントに関連するメトリクスがソリッド・ステート・ドライブ・メトリクスを含む、項目 7 に記載の装置。

30

[項目 10]

複数のサーバにおけるメトリクスを監視するステップと、
データスイッチのルートテーブルに、前記複数のサーバにおけるサービスへの複数の動的ルートを記憶するステップと、

前記監視されたメトリクスに基づいて、前記複数のサーバのうちの1つにおける前記サービスへのルートを動的に許可または抑制するステップと
を含む方法。

[項目 11]

前記サービスへの前記ルートが前記ルートを抑制するために前記ルートテーブルから除去される、項目 10 に記載の方法。

40

[項目 12]

前記サービスへの前記ルートが前記ルートを許可するために前記ルートテーブルに追加される、項目 10 に記載の方法。

[項目 13]

前記監視されるメトリクスが前記複数のサーバにおける圧迫状態または性能劣化に関連するメトリクスを含む、項目 10 に記載の方法。

[項目 14]

前記圧迫状態に関連する前記メトリクスがワークロード・メトリクスを含む、項目 13 に記載の方法。

50

[項目 1 5]

前記性能劣化に関連する前記メトリクスがメモリ帯域幅に関連するメトリクスを含む、項目 1 3 に記載の方法。

[項目 1 6]

前記監視されるメトリクスが前記複数のサーバ内の劣化ハードウェア・コンポーネントに関連するメトリクスを含む、項目 1 0 に記載の方法。

[項目 1 7]

それぞれがサーバにおけるメトリクスを監視する複数のサーバと、前記複数のサーバに通信可能に結合されたデータスイッチであって、前記複数のサーバにおけるサービスへの複数の動的ルートを含むルートテーブルを含む、データスイッチと

、前記監視されたメトリクスに基づいて、前記複数のサーバのうちの 1 つにおける前記サービスへのルートを動的に許可または抑制するフィルタリングシステムと、

を有するラック、

を備えるデータ管理プラットフォーム。

[項目 1 8]

前記サービスへの前記ルートが前記ルートを抑制するために前記ルートテーブルから除去される、項目 1 7 に記載のデータ管理プラットフォーム。

[項目 1 9]

前記サービスへの前記ルートが前記ルートを許可するために前記ルートテーブルに追加される、項目 1 7 に記載のデータ管理プラットフォーム。

[項目 2 0]

前記監視されたメトリクスが前記複数のサーバにおける圧迫状態または性能劣化に関連するメトリクスを含む、項目 1 7 に記載のデータ管理プラットフォーム。

[項目 2 1]

計算サーバと、

ストレージ・サーバであって、前記ストレージ・サーバに通信可能に結合された複数のストレージ・デバイスを管理する、ストレージ・サーバと、を備え、前記計算サーバおよび前記ストレージ・サーバがネットワークを介して通信可能に結合され、前記ストレージ・サーバによって管理される前記複数のストレージ・デバイスが、前記複数のストレージ・デバイスのストレージ容量を前記計算サーバから独立してスケールリングすることができるように、前記計算サーバから分離されている、装置。

[項目 2 2]

前記ストレージ・サーバが、

前記ネットワークに通信可能に結合されたネットワーク・インターフェース・コントローラと、

複数のコアおよびラスト・レベル・メモリを含むシステム・オン・チップであって、前記複数のコアが前記ラスト・レベル・キャッシュ・メモリに通信可能に結合され、前記ラスト・レベル・キャッシュ・メモリが複数のキャッシュ・ウェイを含み、前記複数のキャッシュ・ウェイの一部が、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、システム・オン・チップと、

をさらに有する、項目 2 1 に記載の装置。

[項目 2 3]

前記複数のキャッシュ・ウェイの前記一部が初期化中に割り当てられる、項目 2 2 に記載の装置。

[項目 2 4]

前記論理ボリュームが前記計算サーバによる使用のためのデータを記憶する、項目 2 3 に記載の装置。

[項目 2 5]

前記論理ボリュームと前記ラスト・レベル・キャッシュ内の前記複数のキャッシュ・ウェイとの間で転送される前記データが、前記ネットワーク・インターフェース・コントローラに通信可能に結合されたネットワークを介して前記ストレージ・サーバと前記計算サーバとの間で転送される、項目 2 4 に記載の装置。

[項目 2 6]

前記複数のコアのうちの少なくとも 1 つが、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、項目 2 2 に記載の装置。

10

[項目 2 7]

前記システム・オン・チップに結合された外部メモリであって、転送されるデータを記憶するために、ラスト・レベル・キャッシュの一部における前記複数のキャッシュ・ウェイのすべてが、前記論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられると、前記論理ボリュームと前記ネットワーク・インターフェースとの間で転送されるデータを一時的に記憶する、外部メモリをさらに備える、項目 2 2 に記載の装置。

[項目 2 8]

ストレージ・サーバによって、前記ストレージ・サーバに通信可能に結合された複数のストレージ・デバイスを管理するステップと、

20

ネットワークを介して、前記ストレージ・サーバを計算サーバと通信可能に結合するステップであって、前記ストレージ・サーバによって管理される前記複数のストレージ・デバイスが、前記複数のストレージ・デバイスのストレージ容量を前記計算サーバから独立してスケールアップすることができるように、前記計算サーバから分離されている、ステップと、

を含む方法。

[項目 2 9]

前記ストレージ・サーバが、

前記ネットワークに通信可能に結合されたネットワーク・インターフェース・コントローラと、

30

複数のコアおよびラスト・レベル・メモリを含むシステム・オン・チップであって、前記複数のコアが前記ラスト・レベル・キャッシュ・メモリに通信可能に結合され、前記ラスト・レベル・キャッシュ・メモリが複数のキャッシュ・ウェイを含み、前記複数のキャッシュ・ウェイの一部が、前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、システム・オン・チップと、

をさらに有する、項目 2 8 に記載の方法。

[項目 3 0]

前記複数のキャッシュ・ウェイの前記一部が初期化中に割り当てられる、項目 2 9 に記載の方法。

40

[項目 3 1]

前記論理ボリュームが前記計算サーバによる使用のためのデータを記憶する、項目 3 0 に記載の方法。

[項目 3 2]

前記論理ボリュームと前記ラスト・レベル・キャッシュ内の前記複数のキャッシュ・ウェイとの間で転送される前記データが、前記ネットワーク・インターフェース・コントローラに通信可能に結合されたネットワークを介して前記ストレージ・サーバと前記計算サーバとの間で転送される、項目 3 1 に記載の方法。

[項目 3 3]

50

前記複数のコアのうちの少なくとも1つが前記複数のストレージ・デバイス内の論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられて、前記論理ボリュームと前記ネットワーク・インターフェース・コントローラとの間でデータを転送する、項目29に記載の方法。

[項目34]

前記システム・オン・チップに結合された外部メモリであって、転送されるデータを記憶するために、ラスト・レベル・キャッシュの一部における前記複数のキャッシュ・ウェイのすべてが、前記論理ボリュームおよび前記ネットワーク・インターフェース・コントローラによる排他的使用のために割り当てられると、前記論理ボリュームと前記ネットワーク・インターフェースとの間で転送されるデータを一時的に記憶する、外部メモリをさらに備える、項目29に記載の方法。

10

[項目35]

それぞれが論理リソースを使用してワークロードを実行する複数の計算ノードと、前記複数の計算ノードに通信可能に結合されたオーケストレータであって、メトリクスを監視して計算ノードにおけるアクティブな圧迫状態を検出し、前記アクティブな圧迫状態の検出に応答して、前記圧迫状態がアクティブである間に新しい論理リソースを前記計算ノードにアタッチする、オーケストレータと、

を備える装置。

[項目36]

前記メトリクスがオーケストレータ・メトリクス、ノード・メトリクス、およびワークロード・メトリクスを含む、項目35に記載の装置。

20

[項目37]

前記ノード・メトリクスが計算ノードごとのCPU利用率および計算ノードごとのメモリ利用率を含む、項目36に記載の装置。

[項目38]

前記ワークロード・メトリクスがクライアントの数、平均応答レイテンシ、およびパーセンタイル・メトリクスを含む、項目36に記載の装置。

[項目39]

前記オーケストレータが、前記圧迫状態がアクティブでないときに、前記新しい論理リソースを前記計算ノードからデタッチする、項目35に記載の装置。

30

[項目40]

前記新しい論理リソースが、前記圧迫状態がアクティブである間、前記計算ノード上の特定のワークロードによってのみ使用が制限される、項目36に記載の装置。

[項目41]

前記新しい論理リソースが、前記圧迫状態がアクティブである間、前記計算ノード上のすべてのワークロードによって使用される、項目36に記載の装置。

[項目42]

複数の計算ノードの各計算ノードによって論理リソースを使用してワークロードを実行するステップと、

前記複数の計算ノードに通信可能に結合されたオーケストレータによってメトリクスを監視して、計算ノードにおけるアクティブな圧迫状態を検出するステップと、

40

前記アクティブな圧迫状態の検出に応答して、前記オーケストレータが、前記圧迫状態がアクティブである間に新しい論理リソースを前記計算ノードにアタッチするステップと

、を含む方法。

[項目43]

前記メトリクスがオーケストレータ・メトリクス、ノード・メトリクス、およびワークロード・メトリクスを含む、項目42に記載の方法。

[項目44]

前記ノード・メトリクスが計算ノードごとのCPU利用率および計算ノードごとのメモ

50

リ利用率を含む、項目 4 3 に記載の方法。

[項目 4 5]

前記ワークロード・メトリクスがクライアントの数、平均応答レイテンシ、およびパーセンタイル・メトリクスを含む、項目 4 3 に記載の方法。

[項目 4 6]

前記オーケストレータが、前記圧迫状態がアクティブでないときに、前記新しい論理リソースを前記計算ノードからデタッチする、項目 4 2 に記載の方法。

[項目 4 7]

前記新しい論理リソースが、前記圧迫状態がアクティブである間、前記計算ノード上の特定のワークロードによってのみ使用が制限される、項目 4 2 に記載の方法。

10

[項目 4 8]

前記新しい論理リソースが、前記圧迫状態がアクティブである間、前記計算ノード上のすべてのワークロードによって使用される、項目 4 2 に記載の方法。

[項目 4 9]

それぞれが論理リソースを使用してワークロードを実行する複数の計算ノードと、複数のストレージ・デバイスに通信可能に結合された複数のストレージ・ノードを含むストレージ・サブシステムと、

前記複数の計算ノードおよび前記ストレージ・サブシステムに通信可能に結合されたストレージ自己修復メカニズムであって、前記ストレージ・サブシステム上でヘルスチェックを実行し、前記ヘルスチェックを介して得られたメトリクスを使用して前記ストレージ・サブシステム内の障害を検出し、前記ストレージ・サブシステム内の前記障害の修復を管理する、ストレージ自己修復メカニズムと、

20

を有するラック、

を備えるデータ管理プラットフォーム。

[項目 5 0]

前記障害が計算ノードにあり、前記ストレージ自己修復メカニズムが、故障した計算ノード上で実行されるワークロードを前記ラック内の別の計算ノードに再スケジュールする、項目 4 9 に記載のデータ管理プラットフォーム。

[項目 5 1]

前記障害が計算ノードにあり、前記ラック内の他の計算ノードが同じタイプのワークロードをホストしており、前記ストレージ自己修復メカニズムが前記ワークロードを再スケジュールして、別のラックの別の計算ノード上で実行させる、項目 4 9 に記載のデータ管理プラットフォーム。

30

[項目 5 2]

前記障害がストレージ・デバイスにあり、前記ストレージ自己修復メカニズムが、前記ストレージ・ノード内の別のストレージ・デバイス上にボリュームを作成するために、前記ストレージ・デバイスに関連付けられた前記ワークロードを決定する、項目 4 9 に記載のデータ管理プラットフォーム。

[項目 5 3]

前記ラックがデータスイッチをさらに有し、前記障害が前記データスイッチにあり、前記ストレージ自己修復メカニズムが、前記ラックにおける前記ワークロードを決定し、別のラックを選択し、前記別のラックにおける前記ワークロードを再スケジュールする、項目 4 9 に記載のデータ管理プラットフォーム。

40

[項目 5 4]

データ管理プラットフォーム内のラックにおいて実行される方法であって、

前記ラック内の複数の計算ノードのそれぞれによって、論理リソースを使用してワークロードを実行するステップと、

前記複数の計算ノードに通信可能に結合されたストレージ自己修復メカニズム、および前記ストレージ・サブシステム内の複数のストレージ・デバイスに通信可能に結合された複数のストレージ・ノードを含む前記ラック内のストレージ・サブシステムによって、前

50

記ストレージ・サブシステムに対するヘルスチェックを実行するステップと、

前記ストレージ自己修復メカニズムによって、前記ヘルスチェックを介して得られたメトリクスを使用して、前記ストレージ・サブシステムにおける障害を検出し、前記ストレージ・サブシステムにおける前記障害の修復を管理するステップと、

を含む方法。

[項目 5 5]

前記障害が計算ノードにあり、前記ストレージ自己修復メカニズムが、故障した計算ノード上で実行されるワークロードを前記ラック内の別の計算ノードに再スケジュールする、項目 5 4 に記載の方法。

[項目 5 6]

前記障害が計算ノードにあり、前記ラック内の他の計算ノードが同じタイプのワークロードをホストしており、前記ストレージ自己修復メカニズムが前記ワークロードを再スケジュールして、別のラックの別の計算ノード上で実行させる、項目 5 4 に記載の方法。

[項目 5 7]

前記障害がストレージ・デバイスにあり、前記ストレージ自己修復メカニズムが、前記ストレージ・ノード内の別のストレージ・デバイス上にボリュームを作成するために、前記ストレージ・デバイスに関連付けられた前記ワークロードを決定する、項目 5 4 に記載の方法。

[項目 5 8]

前記ラック内のデータスイッチにおける障害の検出にตอบสนองして、前記ストレージ自己修復メカニズムによって前記ラック内の前記ワークロードを決定するステップと、

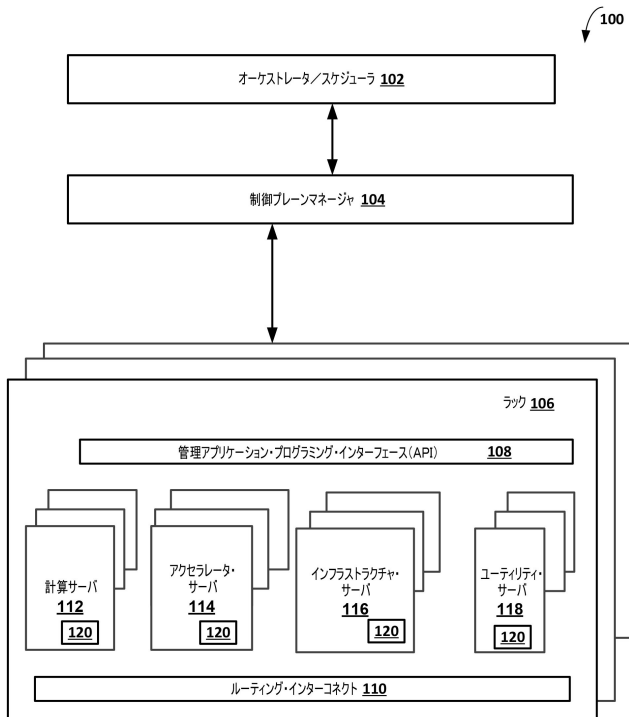
前記ストレージ自己修復メカニズムによって別のラックを選択するステップと、

前記ストレージ自己修復メカニズムによって前記別のラックにおける前記ワークロードを再スケジュールするステップと、

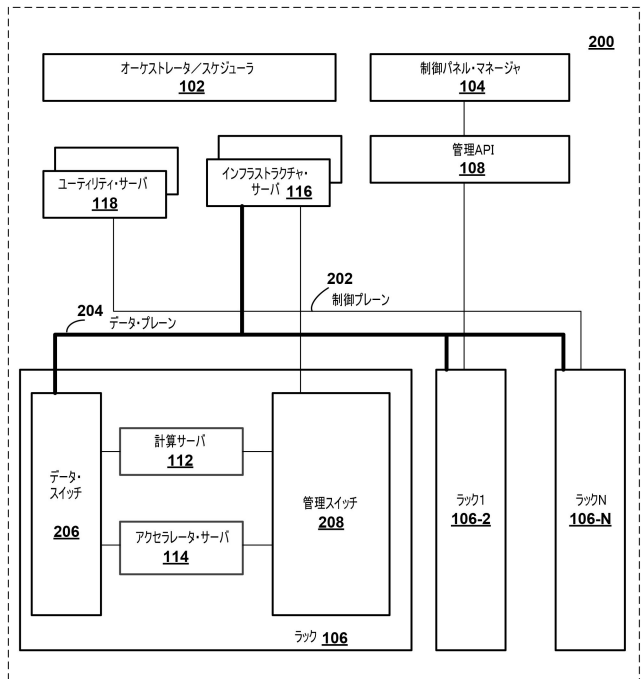
をさらに含む、項目 5 4 に記載の方法。

【 図面 】

【 図 1 】



【 図 2 】



10

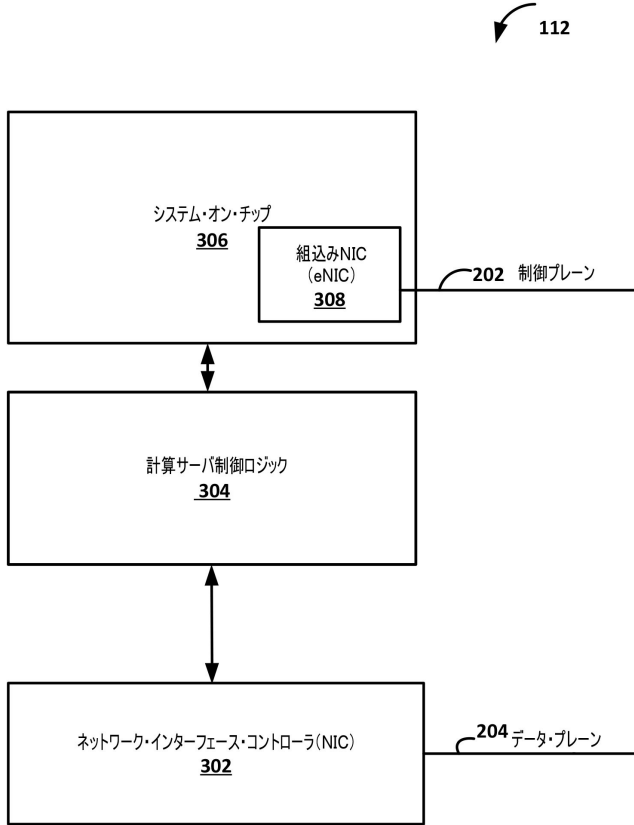
20

30

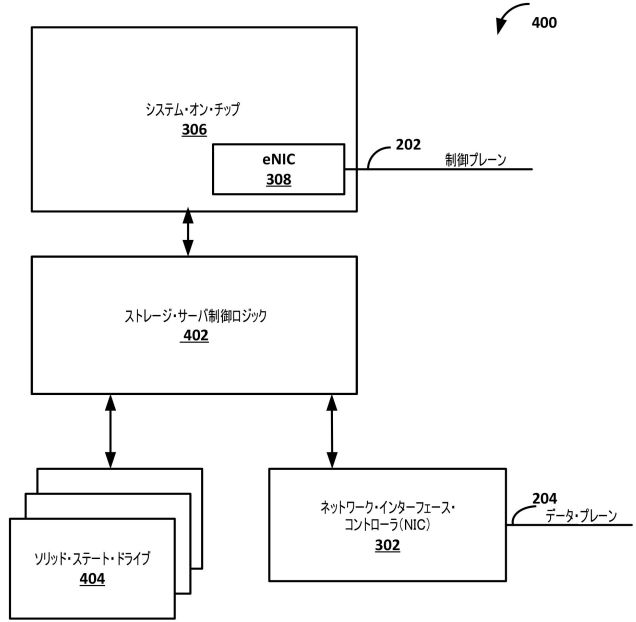
40

50

【 図 3 】



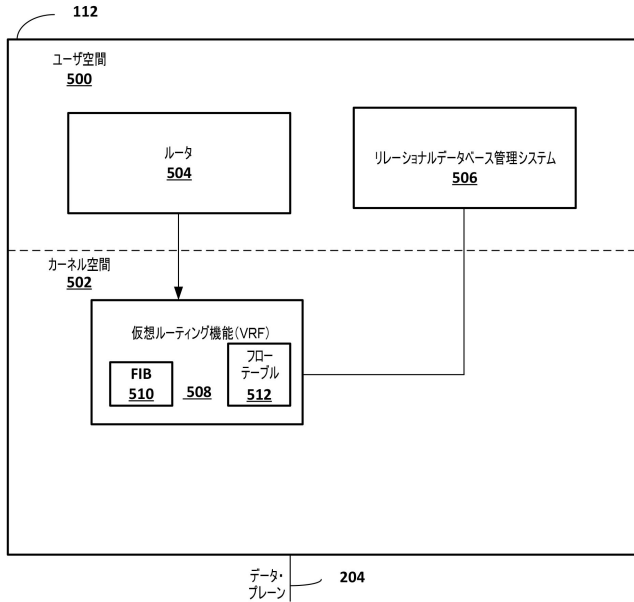
【 図 4 】



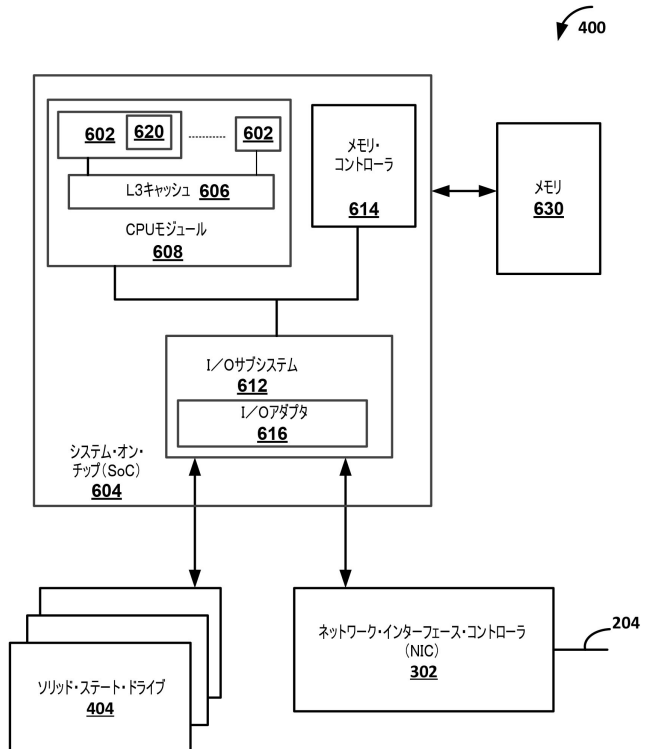
10

20

【 図 5 】



【 図 6 】

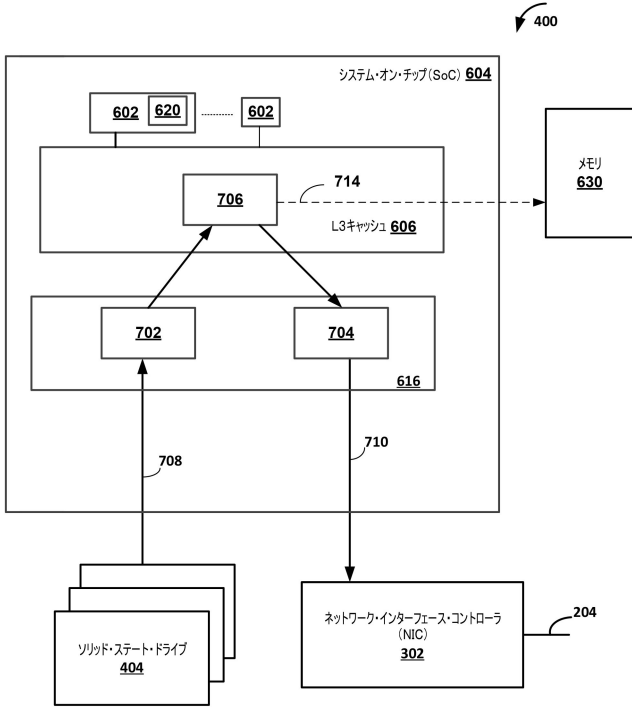


30

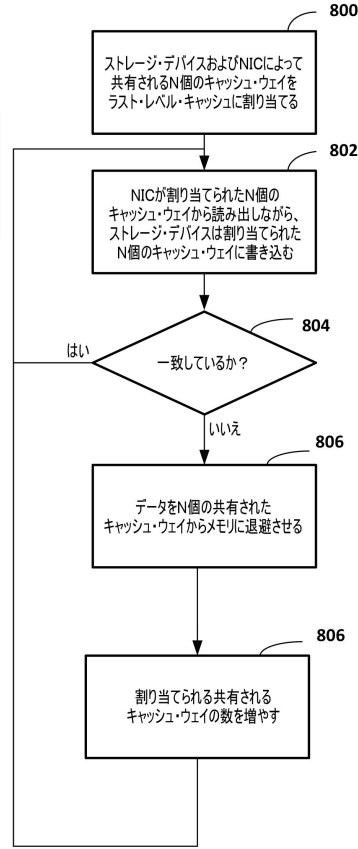
40

50

【 図 7 】



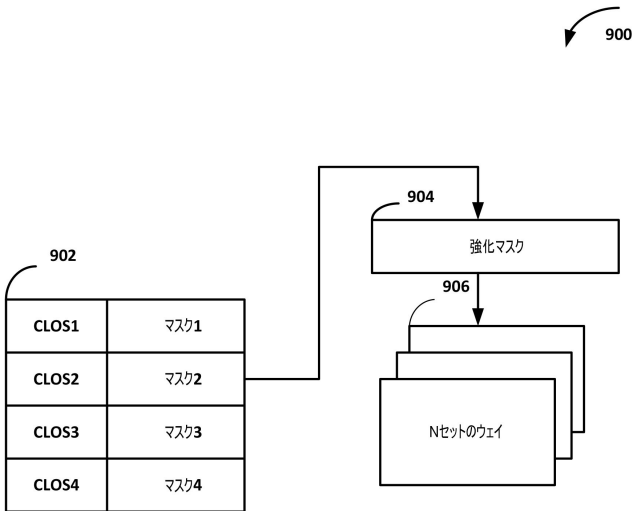
【 図 8 】



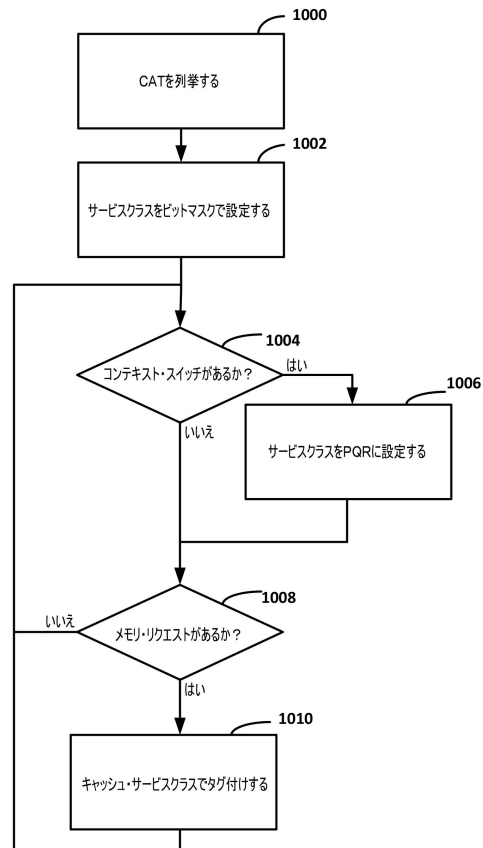
10

20

【 図 9 】



【 図 10 】

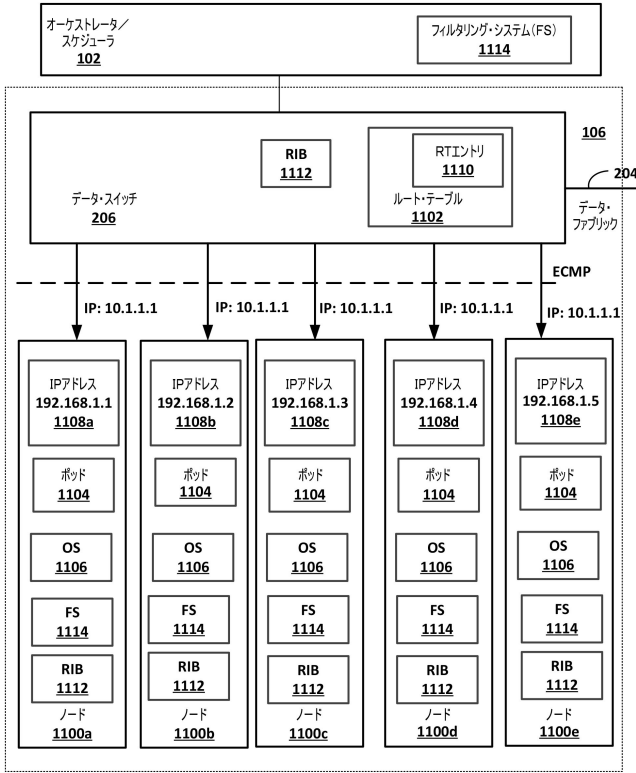


30

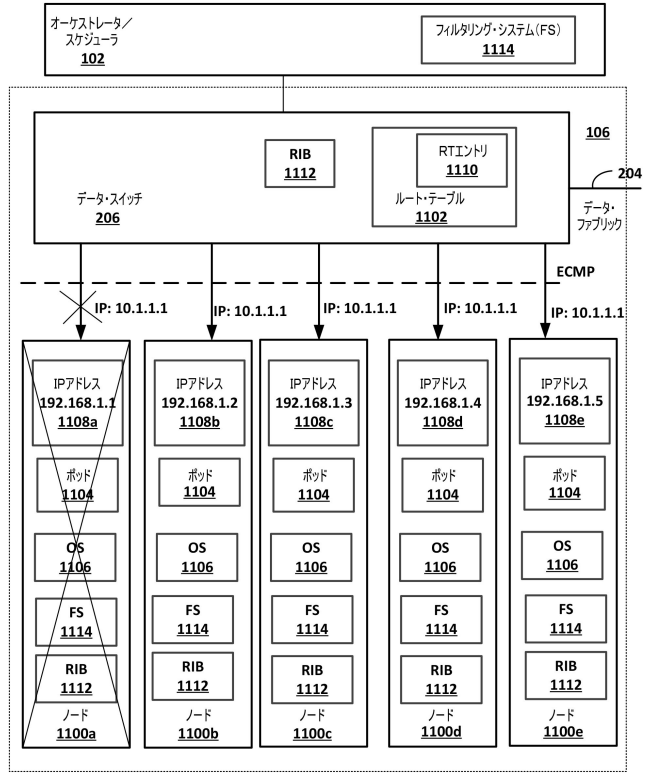
40

50

【図 1 1】



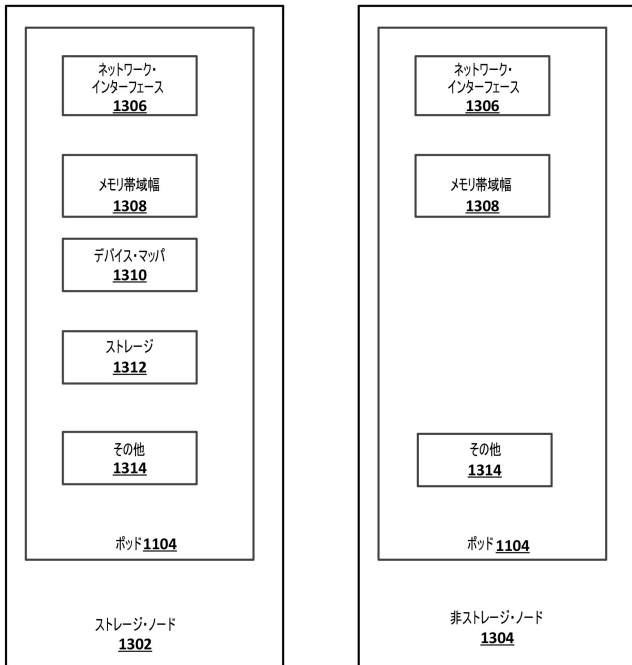
【図 1 2】



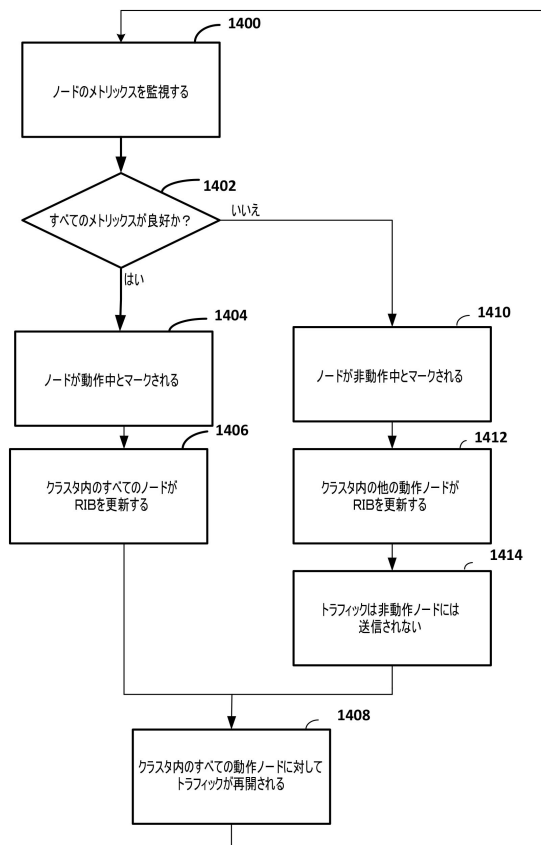
10

20

【図 1 3】



【図 1 4】

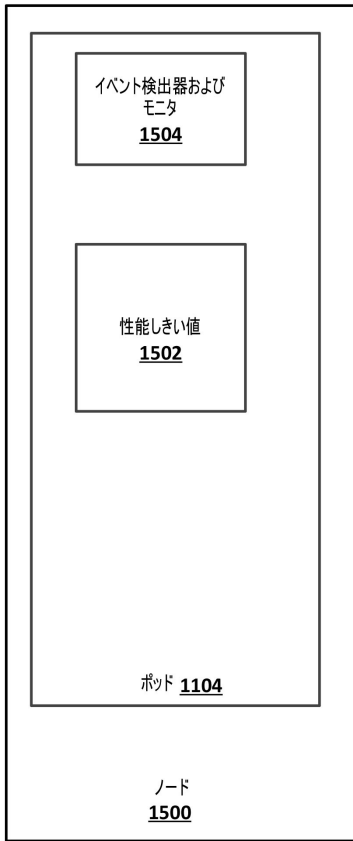


30

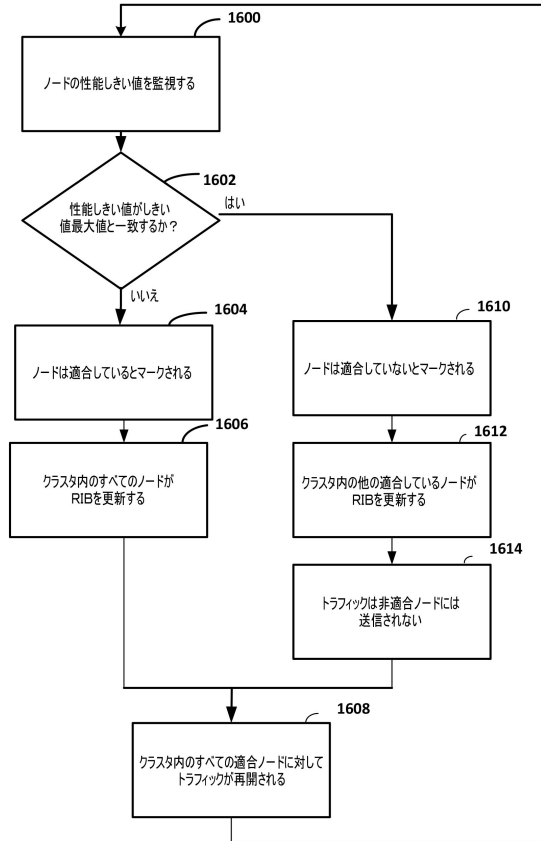
40

50

【 図 1 5 】

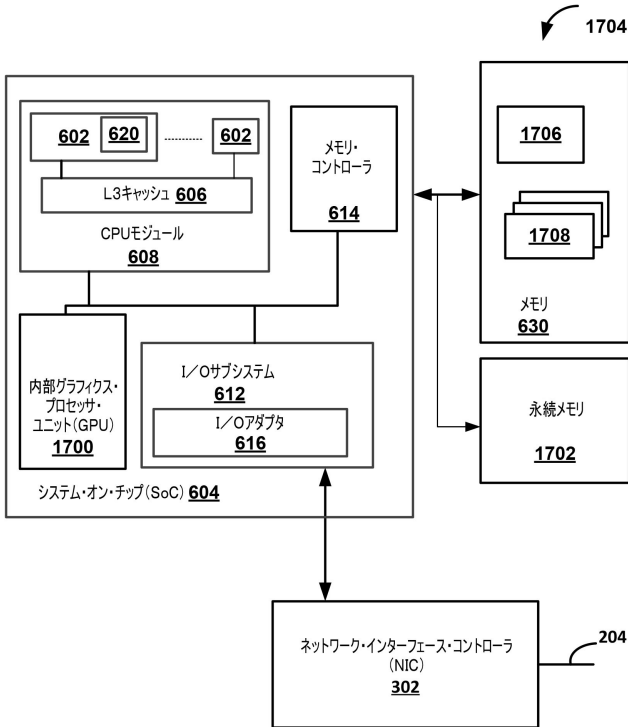


【 図 1 6 】

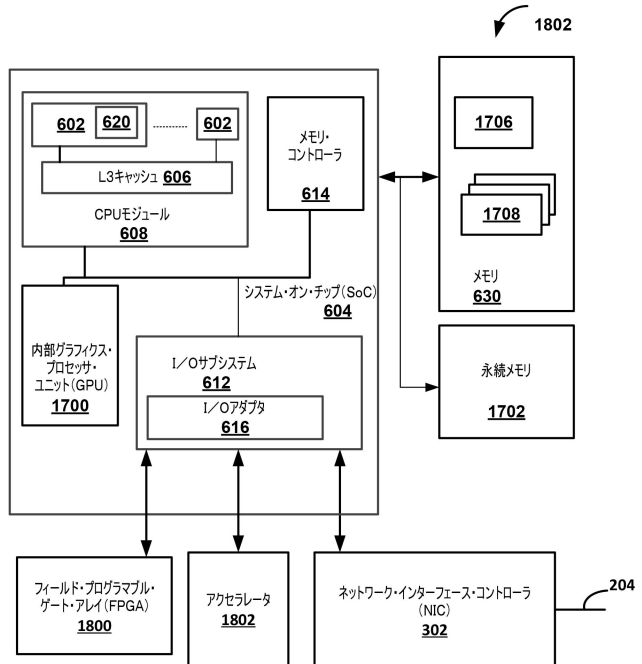


10
20

【 図 1 7 】

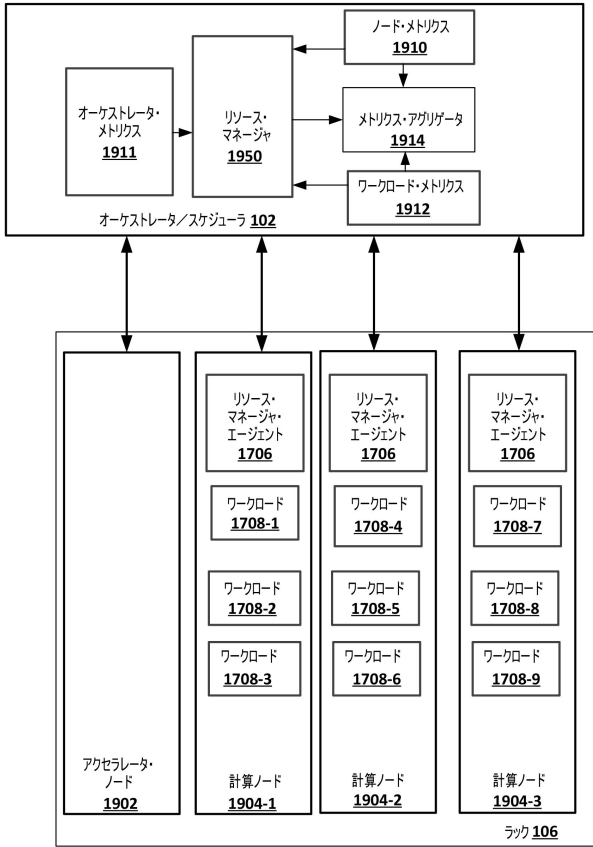


【 図 1 8 】

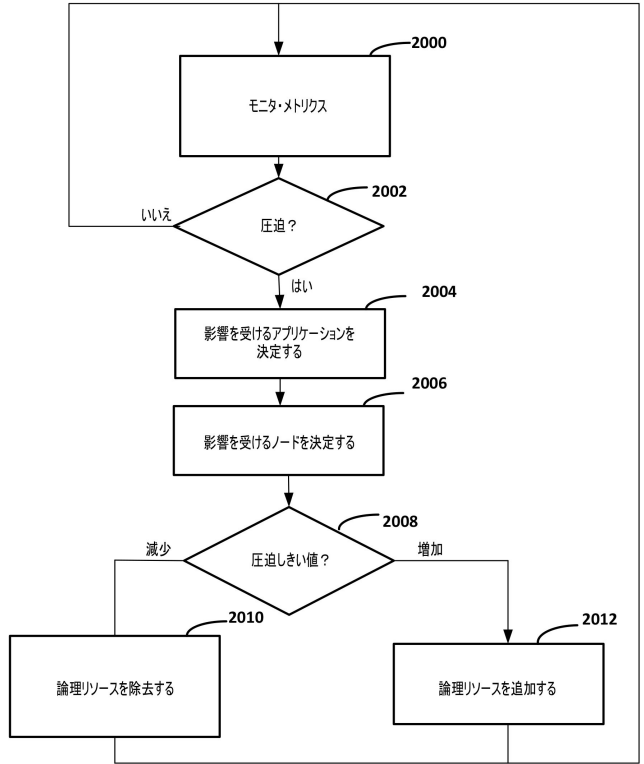


30
40

【図 19】



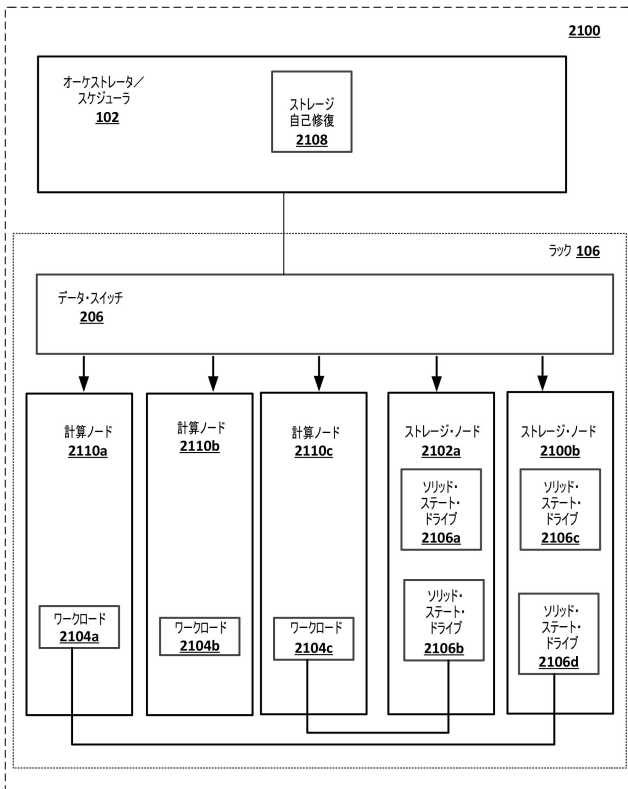
【図 20】



10

20

【図 21】



【図 22】

ワークロード	計算ノード	アクセラレータ・ノード	ソリッド・ステート・ドライブ	ボリューム	ラック
A (インスタンス1)	1	1	1	a	1
A (インスタンス2)	2	2	5	a	1
A (インスタンス3)	1	1	5	b	2

30

40

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US2020/051560

A. CLASSIFICATION OF SUBJECT MATTER		
H04L 12/24(2006.01); H04L 12/26(2006.01); H04L 12/935(2013.01); H04L 12/947(2013.01); H04L 12/861(2013.01); H04L 29/08(2006.01)		
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) H04L 12/24(2006.01); G06F 11/00(2006.01); G06F 11/07(2006.01); G06F 12/08(2006.01); G06F 12/10(2006.01); G06F 13/42(2006.01); G06F 15/173(2006.01); G06F 3/06(2006.01); G06F 9/50(2006.01); H04L 12/26(2006.01); H04L 12/851(2013.01); H04L 29/08(2006.01)		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Korean utility models and applications for utility models Japanese utility models and applications for utility models		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) eKOMPASS(KIPO internal) & Keywords: compute server, storage server, monitoring metric, performance degradation, attaching resource, detecting failure, moving workload		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 7275103 B1 (RUSSELL THRASHER et al.) 25 September 2007 (2007-09-25) column 1, lines 34-52; column 5, line 61 - column 6, line 11; column 7, lines 22-63; column 10, lines 58-67; column 23, line 65 - column 24, line 3; and figures 1-2	1-3,7-8,10-12,16-19
Y		4-6,9,13-15,20
Y	US 2018-0293023 A1 (NUTANIX, INC.) 11 October 2018 (2018-10-11) paragraphs [0014]-[0017]; and figure 1	4-6,9,13-15,20
A	US 2009-0150542 A1 (SATOMI YAHIRO et al.) 11 June 2009 (2009-06-11) paragraphs [0064]-[0435]; and figures 1-29	1-20
X	US 2019-0245924 A1 (ALIBABA GROUP HOLDING LIMITED) 08 August 2019 (2019-08-08) paragraphs [0005], [0027]-[0033]; and figure 1	21,28
A		22-27,29-34
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> 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 "D" document cited by the applicant in the international application "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 "&" document member of the same patent family	
Date of the actual completion of the international search 11 March 2021	Date of mailing of the international search report 11 March 2021	
Name and mailing address of the ISA/KR Korean Intellectual Property Office 189 Cheongsa-ro, Seo-gu, Daejeon 35208, Republic of Korea Facsimile No. +82-42-481-8578	Authorized officer YANG, Jeong Rok Telephone No. +82-42-481-5709	

Form PCT/ISA/210 (second sheet) (July 2019)

10

20

30

40

50

INTERNATIONAL SEARCH REPORT

International application No.

PCT/US2020/051560

C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	US 2017-0091108 A1 (INTEL CORPORATION) 30 March 2017 (2017-03-30) paragraphs [0028]-[0061]; and figures 1-6	21-34
A	KR 10-2016-0074648 A (MARVELL WORLD TRADE LTD.) 28 June 2016 (2016-06-28) paragraphs [0020]-[0091]; and figures 2-4	21-34
X	US 2018-0109610 A1 (AMAZON TECHNOLOGIES, INC.) 19 April 2018 (2018-04-19) paragraph [0049]; claim 30; and figure 3	35,39,42,46-48
Y		36-38,40-41,43-45
Y	US 2019-0273672 A1 (AT&T INTELLECTUAL PROPERTY I, L.P.) 05 September 2019 (2019-09-05) paragraph [0036]; and claim 7	36-38,40-41,43-45
A	US 9946577 B1 (10X GENOMICS, INC.) 17 April 2018 (2018-04-17) column 14, line 7 - column 38, line 45; and figures 1-7	35-48
X	US 10048996 B1 (AMAZON TECHNOLOGIES, INC.) 14 August 2018 (2018-08-14) column 1, lines 7-22; column 3, line 53 - column 4, line 12; column 16, lines 43-51; and figure 1	49,54
Y		50-53,55-58
Y	US 2016-0036924 A1 (MICROSOFT TECHNOLOGY LICENSING, LLC.) 04 February 2016 (2016-02-04) claims 5, 10, 14	50-53,55-58
A	US 2018-0300075 A1 (PURE STORAGE, INC.) 18 October 2018 (2018-10-18) paragraphs [0027]-[0193]; and figures 1A-15	49-58
PX	US 2020-0136943 A1 (INTEL CORPORATION) 30 April 2020 (2020-04-30) paragraphs [0030]-[0249]; and figures 1-22 The above document is a publication of the earlier application whose priority has been claimed in this international application.	1-58

10

20

30

40

50

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US2020/051560

Box No. III Observations where unity of invention is lacking (Continuation of item 3 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

Group I, Claims 1-20 relate to dynamically allowing or suppressing a route to a service in one of a plurality of servers based on monitored metrics.

Group II, Claims 21-34 relate to disaggregating a plurality of storage devices managed by a storage server from a compute server to enable storage capacity of the plurality of storage devices to scale independent of the compute server.

Group III, Claims 35-48 relate to attaching a new logical resource to a compute node while a pressure condition is active, in response to detection of the active pressure condition.

Group IV, Claims 49-58 relate to using metrics obtained via a health check to detect a failure in a storage subsystem and manage repair of the failure in the storage subsystem.

The invention listed as Groups I, II, III and IV do not relate to a single general inventive concept under PCT Rule 13.1, because under PCT Rule 13.2 they lack the same or corresponding special technical features for the following reasons; they are separate inventions with distinct fields of search.

- 1. As all required additional search fees were timely paid by the applicant, this international search report covers all searchable claims. 10
- 2. As all searchable claims could be searched without effort justifying additional fees, this Authority did not invite payment of additional fees. 20
- 3. As only some of the required additional search fees were timely paid by the applicant, this international search report covers only those claims for which fees were paid, specifically claims Nos.:
- 4. No required additional search fees were timely paid by the applicant. Consequently, this international search report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

- Remark on Protest**
- The additional search fees were accompanied by the applicant's protest and, where applicable, the payment of a protest fee. 30
 - The additional search fees were accompanied by the applicant's protest but the applicable protest fee was not paid within the time limit specified in the invitation.
 - No protest accompanied the payment of additional search fees.

40

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.

PCT/US2020/051560

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)			Publication date (day/month/year)
US	7275103	B1	25 September 2007	None			
US	2018-0293023	A1	11 October 2018	None			
US	2009-0150542	A1	11 June 2009	JP	2009-145962	A	02 July 2009
				JP	4782100	B2	28 September 2011
				US	7653725	B2	26 January 2010
US	2019-0245924	A1	08 August 2019	CN	110120915	A	13 August 2019
US	2017-0091108	A1	30 March 2017	CN	107924380	A	17 April 2018
				DE	112016004300	T5	21 June 2018
				US	10120809	B2	06 November 2018
				WO	2017-052909	A1	30 March 2017
KR	10-2016-0074648	A	28 June 2016	CN	106462504	A	22 February 2017
				CN	106463145	A	22 February 2017
				CN	106463145	B	30 August 2019
				EP	3060993	A1	31 August 2016
				EP	3138099	A1	08 March 2017
				JP	2016-541046	A	28 December 2016
				JP	2019-067417	A	25 April 2019
				JP	6431536	B2	28 November 2018
				JP	6796304	B2	09 December 2020
				US	10067687	B2	04 September 2018
				US	10097204	B1	09 October 2018
				US	10684949	B2	16 June 2020
				US	10761737	B2	01 September 2020
				US	2015-0113214	A1	23 April 2015
				US	2015-0242137	A1	27 August 2015
				US	2015-0318022	A1	05 November 2015
				US	2016-0062906	A1	03 March 2016
				US	2016-0239429	A1	18 August 2016
				US	2016-0320981	A1	03 November 2016
				US	2017-0177481	A1	22 June 2017
				US	2017-0344276	A1	30 November 2017
				US	2018-0293167	A1	11 October 2018
				US	2018-0373442	A1	27 December 2018
				US	2020-0301836	A1	24 September 2020
				US	9182915	B2	10 November 2015
				US	9323688	B2	26 April 2016
				US	9454991	B2	27 September 2016
				US	9477611	B2	25 October 2016
				US	9559722	B1	31 January 2017
				US	9594693	B2	14 March 2017
				US	9733841	B2	15 August 2017
				US	9928172	B2	27 March 2018
				WO	2015-061337	A1	30 April 2015
				WO	2015-168609	A1	05 November 2015
US	2018-0109610	A1	19 April 2018	CA	2984142	A1	10 November 2016
				CN	107567696	A	09 January 2018
				CN	107567696	B	12 January 2021
				EP	3289459	A1	07 March 2018
				JP	2018-518744	A	12 July 2018

Form PCT/ISA/210 (patent family annex) (July 2019)

10

20

30

40

50

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.

PCT/US2020/051560

Patent document cited in search report			Publication date (day/month/year)		Patent family member(s)			Publication date (day/month/year)	
					JP	6732798	B2	29 July 2020	
					US	10581964	B2	03 March 2020	
					US	2016-0323377	A1	03 November 2016	
					US	2020-0204623	A1	25 June 2020	
					US	9848041	B2	19 December 2017	
					WO	2016-178951	A1	10 November 2016	
US	2019-0273672	A1	05 September 2019	US	10348590	B2	09 July 2019		
				US	2017-0366428	A1	21 December 2017		
				WO	2017-0218417	A1	21 December 2017		
US	9946577	B1	17 April 2018	US	10162678	B1	25 December 2018		
				US	10452448	B2	22 October 2019		
				US	10795731	B2	06 October 2020		
				US	2019-0108064	A1	11 April 2019		
				US	2020-0159588	A1	21 May 2020		
US	10048996	B1	14 August 2018	None					
US	2016-0036924	A1	04 February 2016	US	10609159	B2	31 March 2020		
				WO	2016-022405	A1	11 February 2016		
US	2018-0300075	A1	18 October 2018	CN	111133409	A	08 May 2020		
				EP	3376361	A2	19 September 2018		
				EP	3376361	A3	10 October 2018		
				EP	3485362	A1	22 May 2019		
				EP	3485363	A1	22 May 2019		
				EP	3485364	A1	22 May 2019		
				EP	3485365	A1	22 May 2019		
				EP	3491551	A1	05 June 2019		
				EP	3491551	A4	18 March 2020		
				US	10007459	B2	26 June 2018		
				US	10019201	B1	10 July 2018		
				US	10146585	B2	04 December 2018		
				US	10162523	B2	25 December 2018		
				US	10191662	B2	29 January 2019		
				US	10203903	B2	12 February 2019		
				US	10275176	B1	30 April 2019		
				US	10275285	B1	30 April 2019		
				US	10331370	B2	25 June 2019		
				US	10331588	B2	25 June 2019		
				US	10353743	B1	16 July 2019		
				US	10360214	B2	23 July 2019		
				US	10366004	B2	30 July 2019		
				US	10452444	B1	22 October 2019		
				US	10459652	B2	29 October 2019		
				US	10481798	B2	19 November 2019		
				US	10534648	B2	14 January 2020		
				US	10545861	B2	28 January 2020		
				US	10585711	B2	10 March 2020		
				US	10613974	B2	07 April 2020		
				US	10649988	B1	12 May 2020		
				US	10671434	B1	02 June 2020		
				US	10671435	B1	02 June 2020		

Form PCT/ISA/210 (patent family annex) (July 2019)

10

20

30

40

50

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.

PCT/US2020/051560

Patent document cited in search report	Publication date (day/month/year)	Patent family member(s)	Publication date (day/month/year)
		US 10671439 B1	02 June 2020
		US 10776034 B2	15 September 2020
		US 10853281 B1	01 December 2020
		US 2018-0032279 A1	01 February 2018
		US 2018-0032280 A1	01 February 2018
		US 2018-0067771 A1	08 March 2018
		US 2018-0067772 A1	08 March 2018
		US 2018-0067775 A1	08 March 2018
		US 2018-0067881 A1	08 March 2018
		US 2018-0074951 A1	15 March 2018
		US 2018-0095662 A1	05 April 2018
		US 2018-0095667 A1	05 April 2018
		US 2018-0095788 A1	05 April 2018
		US 2018-0095871 A1	05 April 2018
		US 2018-0095872 A1	05 April 2018
		US 2018-0113640 A1	26 April 2018
		US 2018-0121088 A1	03 May 2018
		US 2018-0285024 A1	04 October 2018
		US 2018-0357263 A1	13 December 2018
		US 2019-0087094 A1	21 March 2019
		US 2019-0121542 A1	25 April 2019
		US 2019-0121566 A1	25 April 2019
		US 2019-0121673 A1	25 April 2019
		US 2019-0121889 A1	25 April 2019
		US 2019-0171388 A1	06 June 2019
		US 2019-0347195 A1	14 November 2019
		US 2020-0050361 A1	13 February 2020
		US 2020-0125941 A1	23 April 2020
		US 2020-0183827 A1	11 June 2020
		US 2020-0293378 A1	17 September 2020
		US 2020-0394304 A1	17 December 2020
		US 9740408 B1	22 August 2017
		US 9747039 B1	29 August 2017
		US 9892147 B1	13 February 2018
		WO 2018-022136 A1	01 February 2018
		WO 2018-022143 A1	01 February 2018
		WO 2018-022779 A1	01 February 2018
		WO 2018-048592 A1	15 March 2018
		WO 2018-067742 A1	12 April 2018
		WO 2018-067743 A1	12 April 2018
		WO 2018-067744 A1	12 April 2018
		WO 2018-067745 A1	12 April 2018
		WO 2018-067746 A1	12 April 2018
		WO 2018-075676 A1	26 April 2018
		WO 2018-075790 A1	26 April 2018
		WO 2018-218259 A1	29 November 2018
US	2020-0136943	A1	30 April 2020
		None	

Form PCT/ISA/210 (patent family annex) (July 2019)

10

20

30

40

50

フロントページの続き

MK,MT,NL,NO,PL,PT,RO,RS,SE,SI,SK,SM,TR),OA(BF,BJ,CF,CG,CI,CM,GA,GN,GQ,GW,KM,ML,MR,NE,SN,TD,TG),AE,AG,AL,AM,AO,AT,AU,AZ,BA,BB,BG,BH,BN,BR,BW,BY,BZ,CA,CH,CL,CN,CO,CR,CU,CZ,DE,DJ,DK,DM,DO,DZ,EC,EE,EG,ES,FI,GB,GD,GE,GH,GM,GT,HN,HR,HU,ID,IL,IN,IR,IS,IT,JO,JP,KE,KG,KH,KN,KP,KR,KW,KZ,LA,LC,LK,LR,LS,LU,LY,MA,MD,ME,MG,MK,MN,MW,MX,MY,MZ,NA,NG,NI,NO,NZ,OM,PA,PE,PG,PH,PL,PT,QA,RO,RS,RU,RW,SA,SC,SD,SE,SG,SK,SL,ST,SV,SY,TH,TJ,TM,TN,TR,TT,TZ,UA,UG,US,UZ,VC,VN,WS,ZA,ZM,ZW

レッジ ブレーバード・2200 インテル・コーポレーション内

(72)発明者 カルバリオ、ジョー

アメリカ合衆国 95054 カリフォルニア州・サンタクララ・ミッション カレッジ ブレーバード・2200 インテル・コーポレーション内

(72)発明者 スタハーフスキ、ミカル

アメリカ合衆国 95054 カリフォルニア州・サンタクララ・ミッション カレッジ ブレーバード・2200 インテル・コーポレーション内

(72)発明者 アロウリ、ブラサド

アメリカ合衆国 95054 カリフォルニア州・サンタクララ・ミッション カレッジ ブレーバード・2200 インテル・コーポレーション内

(72)発明者 シャルマッハ、ザイモン トマシュ

アメリカ合衆国 95054 カリフォルニア州・サンタクララ・ミッション カレッジ ブレーバード・2200 インテル・コーポレーション内