

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7650983号  
(P7650983)

(45)発行日 令和7年3月25日(2025.3.25)

(24)登録日 令和7年3月14日(2025.3.14)

(51)国際特許分類	F I
G 0 6 T 15/00 (2011.01)	G 0 6 T 15/00 5 0 1
G 0 6 F 12/0875(2016.01)	G 0 6 F 12/0875 1 0 6
G 0 6 F 12/0895(2016.01)	G 0 6 F 12/0895 1 1 8
G 0 6 F 9/38 (2018.01)	G 0 6 F 9/38 3 1 0 J

請求項の数 15 (全15頁)

(21)出願番号	特願2023-539265(P2023-539265)	(73)特許権者	591016172
(86)(22)出願日	令和3年12月22日(2021.12.22)		アドバンスト・マイクロ・デバイス ・インコーポレイテッド
(65)公表番号	特表2024-501015(P2024-501015 A)		ADVANCED MICRO DEVI CES INCORPORATED
(43)公表日	令和6年1月10日(2024.1.10)		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、 オーガスティ ン ドライブ 2 4 8 5
(86)国際出願番号	PCT/US2021/064797	(73)特許権者	503447036
(87)国際公開番号	WO2022/146810		サムスン エレクトロニクス カンパニー リミテッド
(87)国際公開日	令和4年7月7日(2022.7.7)		大韓民国・1 6 6 7 7・キョンギ-ド・ スウォン-シ・ヨントン-ク・サムスン -ロ・1 2 9
審査請求日	令和6年12月13日(2024.12.13)	(74)代理人	100108833
(31)優先権主張番号	17/134,790		
(32)優先日	令和2年12月28日(2020.12.28)		
(33)優先権主張国・地域又は機関	米国(US)		
早期審査対象出願			

最終頁に続く

(54)【発明の名称】 キャッシュラインに対するミス要求の選択的生成

(57)【特許請求の範囲】

【請求項1】

装置であって、

複数のサブセットに分割されたキャッシュラインを含むテクスチャキャッシュと、  
グラフィックスパイプライン内の少なくとも一つの計算ユニットと、を備え、

前記計算ユニットは、前記テクスチャキャッシュ内のキャッシュラインの複数のサブセ  
ットのうち第1のサブセットに関連付けられたアドレスへのメモリアクセス要求に対する  
キャッシュミスと、色圧縮又は深度圧縮の何れが有効であることを示す、前記キャッシュラ  
インに記憶されたデータの特性と、に応じて、前記第1のサブセットに対するミス要求を  
選択的に生成するように構成されている、

装置。

【請求項2】

前記少なくとも一つの計算ユニットは、メモリアクセス要求に関連付けられたキャッシ  
ュミスが、前記複数のサブセットのうち前記第1のサブセットのみにマッピングするか、  
又は、前記複数のサブセットのうち前記第1のサブセットに追加された又は前記第1のサ  
ブセット以外の一つ以上のサブセットにマッピングするかを判定するように構成されてい  
る、

請求項1の装置。

【請求項3】

前記少なくとも一つの計算ユニットは、前記キャッシュミスが、前記複数のサブセット

のうち前記第 1 のサブセットに追加された又は前記第 1 のサブセット以外のサブセットにマッピングすることに応じて、フルキャッシュラインに対するミス要求を生成するように構成されている、

請求項 2 の装置。

【請求項 4】

前記少なくとも 1 つの計算ユニットは、前記メモリアクセス要求が、前記第 1 のサブセットのみにマッピングすることに応じて、色圧縮及び深度圧縮のうち少なくとも 1 つがテクスチャデータに対して有効にされているか否かを判定するように構成されている、

請求項 3 の装置。

【請求項 5】

前記少なくとも 1 つの計算ユニットは、色圧縮及び深度圧縮のうち少なくとも 1 つが前記テクスチャデータに対して有効にされていることに応じて、前記フルキャッシュラインに対するミス要求を生成するように構成されている、

請求項 4 の装置。

【請求項 6】

前記少なくとも 1 つの計算ユニットは、色圧縮及び深度圧縮のうち少なくとも 1 つが前記テクスチャデータに対して有効にされていないことに応じて、前記キャッシュラインの前記第 1 のサブセットに対するミス要求を生成するように構成されている、

請求項 4 の装置。

【請求項 7】

前記少なくとも 1 つの計算ユニットは、前記メモリアクセス要求の時間的局所性及び空間的局所性のうち少なくとも 1 つに基づいて、前記第 1 のサブセット又は前記複数のサブセットに対するミス要求を選択的に生成するように構成されている、

請求項 1 の装置。

【請求項 8】

前記少なくとも 1 つの計算ユニットは、メモリアクセス要求シーケンスが、前記複数のサブセットにアクセスすることが予想されることに応じて、前記第 1 のサブセットにおけるキャッシュミスに応じて、前記複数のサブセットに対するミス要求を生成するように構成されている、

請求項 7 の装置。

【請求項 9】

前記少なくとも 1 つの計算ユニットは、前記第 1 のサブセット内のキャッシュミスに応じて、且つ、メモリアクセス要求シーケンスが、閾値を上回る空間的局所性を有することに応じて、前記複数のサブセットに対するミス要求を生成するように構成されている、

請求項 7 の装置。

【請求項 10】

前記少なくとも 1 つの計算ユニットは、前記メモリアクセス要求が前記閾値を下回る空間的局所性を有することに応じて、前記第 1 のサブセットに対するミス要求を生成するように構成されている、

請求項 9 の装置。

【請求項 11】

方法であって、

複数のサブセットに分割されたキャッシュラインを含むテクスチャキャッシュ内のキャッシュラインに対するミス要求を検出することと、

前記キャッシュラインの第 1 のサブセットに関連付けられたアドレスに対するキャッシュミスと、色圧縮又は深度圧縮の何れが有効であることを示す、前記キャッシュラインに記憶されたデータの特性と、に応じて、前記第 1 のサブセットに対するミス要求を選択的に生成することと、を含む、

方法。

【請求項 12】

10

20

30

40

50

メモリアクセス要求に関連付けられたキャッシュミスが、前記複数のサブセットのうち前記第 1 のサブセットのみにマッピングするか、又は、前記複数のサブセットのうち前記第 1 のサブセットに追加された又は前記第 1 のサブセット以外の 1 つ以上のサブセットにマッピングするかを判定することを含む、

請求項 1 1 の方法。

【請求項 1 3】

前記キャッシュミスが前記複数のサブセットにマッピングすることに応じて、フルキャッシュラインに対するミス要求を生成することを含む、

請求項 1 2 の方法。

【請求項 1 4】

色圧縮及び深度圧縮のうち少なくとも 1 つがテクスチャデータに対して有効にされているか否かを判定することを含む、

請求項 1 3 の方法。

【請求項 1 5】

前記ミス要求を選択的に生成することは、色圧縮及び深度圧縮のうち少なくとも 1 つが前記テクスチャデータに対して有効にされていることに応じて、前記複数のサブセットに対するミス要求を生成することを含む、

請求項 1 4 の方法。

【発明の詳細な説明】

【背景技術】

【0001】

グラフィックス処理ユニット (Graphics Processing Unit、GPU) は、プログラマブルシェーダ及び固定機能ハードウェアブロックシーケンスで形成されるグラフィックスパイプラインを使用して三次元 (three-dimensional、3D) グラフィックスを処理する。例えば、フレーム内で見えるオブジェクトの 3D モデルは、三角形、他の多角形又はパッチのセットによって表すことができ、これらはグラフィックスパイプラインで処理され、ユーザに表示するためのピクセルの値を生成する。三角形、他の多角形又はパッチは、まとめてプリミティブと呼ばれる。レンダリングプロセスは、テクスチャをプリミティブにマッピングして、プリミティブの解像度よりも高い解像度を有する視覚的詳細を組み込むことを含む。GPU は、グラフィックスパイプラインにおいて処理されているプリミティブにマッピングするためにテクスチャ値が利用可能であるように、テクスチャ値を記憶するために使用される専用メモリを含む。テクスチャは、ディスク上に記憶することができ、又は、グラフィックスパイプラインによって必要とされる場合に手続き的に生成することができる。専用 GPU メモリに記憶されたテクスチャデータは、ディスクからテクスチャをロードすることによって又はデータを手続き的に生成することによってポピュレートされる。頻繁に使用されるテクスチャデータは、シェーダ又は固定機能ハードウェアブロックによってアクセスされる 1 つ以上のテクスチャキャッシュにキャッシュされる。

【0002】

本開示は、添付の図面を参照することによってより良好に理解され、その多くの特徴及び利点が当業者に明らかになる。異なる図面における同じ符号の使用は、類似又は同一のアイテムを示す。

【図面の簡単な説明】

【0003】

【図 1】いくつかの実施形態による、キャッシュラインの部分に対するミス要求を選択的に生成する処理システムのブロック図である。

【図 2】いくつかの実施形態による、高次ジオメトリプリミティブを処理して、所定の解像度で三次元 (3D) シーンのラスタ化された画像を生成するように構成されたグラフィックスパイプラインを示す図である。

【図 3】いくつかの実施形態による、第 1 の読み取りサイクルにおいて複数のセクタにわたって分散された要求と、第 2 の読み取りサイクルにおいて単一のセクタに制約された要

10

20

30

40

50

求と、を有するキャッシュラインのブロック図である。

【図4】いくつかの実施形態による、第1の読み取りサイクル及び第2の読み取りサイクル中に高度の時間的局所性を示さない要求を有するキャッシュラインのブロック図である。

【図5】いくつかの実施形態による、キャッシュラインの部分に対するミス要求を選択的に生成する方法のフロー図である。

【発明を実施するための形態】

【0004】

テクスチャキャッシュ内のキャッシュラインは、通常、大量のデータを保持するように構成され、例えば、テクスチャキャッシュラインの幅は、128バイト又は1024(1K)ビット程度とすることができる。広いキャッシュラインは、グラフィックス処理の特性である大きな及び/又は可変サイズのデータブロックのキャッシュを容易にする。テクスチャデータは、4×4ピクセルフットプリント又は8×8ピクセルフットプリントを有するタイル等のタイルに記憶される。タイルのサイズは、テクスチャフォーマットにも依存し、テクスチャフォーマットは、8ビットフォーマット、32ビットフォーマット、128ビットフォーマット等のように、各ピクセルを表すために使用されるビット数を示す。したがって、8×8ピクセルフットプリントを有するタイルは、テクスチャフォーマットに応じて、526ビット、2048ビット、8192ビット又は他のビット数で表すことができる。動作中、テクスチャキャッシュは、サイクルごとに最大N個のメモリアクセス要求(例えば、読み取り要求又は書き込み要求)を受信し、ここで、Nはベクトルのサイズであり(例えば、ベクトルサイズは64、32又は16とすることができる)、各キャッシュミスは、より高レベルのキャッシュ又はメモリからキャッシュラインを取り出す要求を生成する。キャッシュラインのサイズが大きいと仮定すると、要求されたデータが複数のキャッシュラインにわたって分散される場合、キャッシュミス要求は、元のアクセス要求内のデータ量にかかわらず、かなりのメモリ帯域幅を消費する。更に、全ての要求サイクルについてフルキャッシュラインを有効にすることは、グラフィックスパイプラインによって使用されているデータを記憶するために必要とされないキャッシュの部分が無効にすることによって、電力節約の機会を制限する。

【0005】

図1～図5は、キャッシュラインのサブセットに関連付けられたアドレスへのメモリアクセス要求に対するキャッシュミスに応じて、テクスチャキャッシュ内のキャッシュラインのサブセットに対するミス要求を選択的に生成することによって、テクスチャキャッシュとシステムメモリ(又はより高レベルのキャッシュ)との間のメモリ帯域幅を節約しながら、電力消費を潜在的に低減するためのシステム及び技術を開示する。いくつかの実施形態では、キャッシュラインは2つ以上のセクタに分割される。フルキャッシュラインに対するミス要求は、例えば、メモリアクセス要求内のアドレスに基づいて、キャッシュライン内の全てのセクタにマッピングするメモリアクセス要求(読み取り要求等)によるキャッシュミスに応じて生成される。メモリアクセス要求がキャッシュラインの単一のセクタにマッピングされる場合、ミス要求は、テクスチャデータの1つ以上のヒューリスティック又は特性の評価に基づいて、フルキャッシュライン又はキャッシュラインのセクタのうち何れかに対して選択的に生成される。例えば、テクスチャデータに対して色圧縮又は深度圧縮が有効にされている場合、ミス要求がフルキャッシュラインに対して生成される。テクスチャデータに対して圧縮が有効にされていない場合、ミス要求は、メモリアクセス要求によって示されるキャッシュラインのセクタに対してのみ生成される。また、ミス要求は、メモリアクセス要求の時間的局所性に基づいてキャッシュラインのサブセットに対して選択的に生成される。例えば、メモリアクセス要求シーケンスがキャッシュラインの異なるセクタにアクセスすると予想される場合、何れかのセクタにおけるキャッシュミスに応じて、フルキャッシュラインに対するミス要求が生成される。また、ミス要求は、メモリアクセス要求の空間的局所性に基づいてキャッシュラインのサブセットに対して選択的に生成される。例えば、メモリアクセス要求シーケンスが、隣接する、近接する又は近くのアドレスにアクセスすることが予想される場合、ミス要求は、何れかのセクタにお

10

20

30

40

50

けるミスに応じて、フルキャッシュラインについて生成される。対照的に、メモリアクセス要求のアドレスが分散しており、低い空間的局所性を有する場合、ミス要求は、キャッシュミスを含むキャッシュラインのセクタに対してのみ生成される。

【0006】

図1は、いくつかの実施形態による、キャッシュラインの部分に対するミス要求を選択的に生成する処理システム100のブロック図である。処理システム100は、ダイナミックランダムアクセスメモリ(Dynamic Random-Access Memory、DRAM)等の非一時的なコンピュータ可読記憶媒体を使用して実装されるメモリ105又は他の記憶コンポーネントを含むか又はそれらへのアクセスを有する。しかしながら、場合によっては、メモリ105は、スタティックランダムアクセスメモリ(Static Random-Access Memory、SRAM)、不揮発性RAM等を含む他のタイプのメモリを使用して実装することもできる。メモリ105は、処理システム100において実装される処理ユニットの外部に実装されるために外部メモリと呼ばれる。また、処理システム100は、メモリ105等のように、処理システム100において実装されるエンティティ間の通信をサポートするためのバス110を含む。処理システム100のいくつかの実施形態は、他のバス、ブリッジ、スイッチ、ルータ等を含むが、これらは明確にするために図1には示されていない。

【0007】

本明細書で説明される技術は、様々な実施形態では、様々な並列プロセッサ、例えば、ベクトルプロセッサ、グラフィックス処理ユニット(GPU)、汎用GPU(GPGPU)、非スカルプロセッサ、高並列プロセッサ、人工知能(AI)プロセッサ、推論エンジン、機械学習プロセッサ、他のマルチスレッド処理ユニット等の何れかで利用される。図1は、いくつかの実施形態による、並列プロセッサ、特に、グラフィックス処理ユニット(GPU)115の一例を示す。グラフィックス処理ユニット(GPU)115は、ディスプレイ120上に提示するための画像をレンダリングする。例えば、GPU115は、オブジェクトをレンダリングして、ディスプレイ120に提供されるピクセルの値を生成し、ディスプレイ120は、ピクセル値を使用して、レンダリングされたオブジェクトを表す画像を表示する。GPU115は、命令を同時に又は並列に実行する複数の計算ユニット(CU)121、122、123(本明細書ではまとめて「計算ユニット121~123」と呼ぶ)を実装する。いくつかの実施形態では、計算ユニット121~123は、1つ以上の単一命令複数データ(SIMD)ユニットを含み、計算ユニット121~123は、ワークグループプロセッサ、シェーダアレイ、シェーダエンジン等に集約される。GPU115において実装される計算ユニット121~123の数は、設計上の選択の問題であり、GPU115のいくつかの実施形態は、図1に示されるよりも多い又は少ない計算ユニットを含む。計算ユニット121~123は、本明細書で説明するように、グラフィックスパイプラインを実装するために使用することができる。GPU115のいくつかの実施形態は、汎用コンピューティングのために使用される。GPU115は、メモリ105に記憶されたプログラムコード125等の命令を実行し、GPU115は、実行された命令の結果等の情報をメモリ105に記憶する。

【0008】

また、処理システム100は、バス110に接続され、したがってバス110を介してGPU115及びメモリ105と通信する中央処理装置(Central Processing Unit、CPU)130を含む。CPU130は、命令を同時に又は並列に実行する複数のプロセッサコア131、132、133(本明細書ではまとめて「プロセッサコア131~133」と呼ぶ)を実装する。CPU130において実装されるプロセッサコア131~133の数は、設計上の選択の問題であり、いくつかの実施形態は、図1に示されるよりも多い又は少ないプロセッサコアを含む。プロセッサコア131~133は、メモリ105に記憶されたプログラムコード135等の命令を実行し、CPU130は、実行された命令の結果等の情報をメモリ105に記憶する。また、CPU130は、GPU115にドローコールを発行することによって、グラフィックス処理を開始することができる。CPU130のいくつかの実施形態は、同時に又は並列に命令を独立して実行する複数のプロセッ

10

20

30

40

50

サコア（明確化のために図 1 には示さず）を含む。

【 0 0 0 9 】

入力/出力（Input/Output、I/O）エンジン 1 4 5 は、ディスプレイ 1 2 0 と関連付けられた入力又は出力動作、及び、キーボード、マウス、プリンタ、外部ディスク等のような処理システム 1 0 0 の他の要素を扱う。I/O エンジン 1 4 5 は、I/O エンジン 1 4 5 がメモリ 1 0 5、GPU 1 1 5 又は CPU 1 3 0 と通信するようにバス 1 1 0 に結合される。図示される実施形態では、I/O エンジン 1 4 5 は、コンパクトディスク（Compact Disk、CD）、デジタルビデオディスク（Digital Video Disc、DVD）等の非一時的なコンピュータ可読記憶媒体を使用して実装される、外部記憶コンポーネント 1 5 0 上に記憶される情報を読み取る。また、I/O エンジン 1 4 5 は、GPU 1 1 5 又は CPU 1 3 0 による処理の結果等の情報を外部記憶コンポーネント 1 5 0 に書き込むことができる。

10

【 0 0 1 0 】

図示した実施形態では、GPU 1 1 5 内の計算ユニット 1 2 1 ~ 1 2 3 は、本明細書では集合的に「キャッシュ 1 5 1 ~ 1 5 2」と呼ばれる 1 つ以上のキャッシュ 1 5 1、1 5 3、1 5 3 を含む（又はそれらに関連付けられる）。キャッシュ 1 5 1 ~ 1 5 3 は、L 1 キャッシュ、L 2 キャッシュ、L 3 キャッシュ、又は、キャッシュ階層内の他のキャッシュを含むことができる。キャッシュ 1 5 1 ~ 1 5 3 の部分は、計算ユニット 1 2 1 ~ 1 2 3 上で実行されるグラフィックスパイプラインのためのテクスチャキャッシュを実装するために使用される。キャッシュ 1 5 1 ~ 1 5 3 内のキャッシュラインは、キャッシュラインの 1 つ以上のセクタ等のサブセットに分割される。グラフィックスパイプラインは、キャッシュラインのサブセットに関連付けられたアドレスへのメモリアクセス要求に対するキャッシュミスに応じて、テクスチャキャッシュ内のキャッシュラインのサブセットに対するミス要求を選択的に生成する。いくつかの実施形態では、キャッシュラインは、第 1 のセクタ及び第 2 のセクタに分割される。本明細書で説明するように、要求サイクル中に受信されたメモリアクセス要求に対するキャッシュミスが第 1 のセクタ内に（排他的に又は主に）あることに応じて、ミス要求が第 1 のセクタに対して生成され、第 2 のセクタに対するミス要求の生成がバイパスされる。

20

【 0 0 1 1 】

図 2 は、いくつかの実施形態による、高次ジオメトリプリミティブを処理して、所定の解像度で三次元（3D）シーンのラスタ化された画像を生成するように構成されたグラフィックスパイプライン 2 0 0 を示す。グラフィックスパイプライン 2 0 0 は、図 1 に示される処理システム 1 0 0 のいくつかの実施形態で実施される。グラフィックスパイプライン 2 0 0 の図示された実施形態は、DX 1 1 仕様に従って実装される。グラフィックスパイプライン 2 0 0 の他の実施形態は、Vulkan、Metal、DX 1 2 等の他のアプリケーションプログラミングインターフェース（Application Programming Interfaces、API）に従って実装される。グラフィックスパイプライン 2 0 0 は、ラスタ化前のグラフィックスパイプライン 2 0 0 の部分を含むジオメトリ部 2 0 1 と、ラスタ化後のグラフィックスパイプライン 2 0 0 の部分を含むピクセル処理部 2 0 2 と、に細分される。

30

【 0 0 1 2 】

グラフィックスパイプライン 2 0 0 は、バッファを実装し、頂点データ、テクスチャデータ等を記憶するために使用される 1 つ以上のメモリ又はキャッシュの階層等のストレージリソース 2 0 5 へのアクセスを有する。図示される実施形態では、ストレージリソース 2 0 5 は、データを記憶するために使用されるローカルデータストア（LDS）2 0 6 回路と、グラフィックスパイプライン 2 0 0 によるレンダリング中に頻繁に使用されるデータをキャッシュするために使用されるキャッシュ 2 0 7 と、を含む。ストレージリソース 2 0 5 は、図 1 に示されるメモリ 1 0 5 のいくつかの実施形態を使用して実装され得る。

40

【 0 0 1 3 】

入力アセンブラ 2 1 0 は、シーンのモデルの部分を表すオブジェクトを定義するために使用される、ストレージリソース 2 0 5 から情報にアクセスする。プリミティブの一例が

50

三角形 2 1 1 として図 2 に示されているが、グラフィックスパイプライン 2 0 0 のいくつかの実施形態では、他のタイプのプリミティブが処理される。三角形 2 1 1 は、1 つ以上の辺 2 1 4 によって接続された 1 つ以上の頂点 2 1 2 を含む（明確にするために、図 2 には各々の 1 つのみが示されている）。頂点 2 1 2 は、グラフィックスパイプライン 2 0 0 のジオメトリ処理部 2 0 1 中にシェーディングされる。

【 0 0 1 4 】

頂点シェーダ 2 1 5 は、図示される実施形態ではソフトウェアで実装されており、プリミティブの単一の頂点 2 1 2 を入力として論理的に受信し、単一の頂点を出力する。頂点シェーダ 2 1 5 等のシェーダのいくつかの実施形態は、複数の頂点が同時に処理されるように、単一命令 - 複数データ (SIMD) 処理を実装する。グラフィックスパイプライン 2 0 0 は、グラフィックスパイプライン 2 0 0 に含まれる全てのシェーダが、共有大規模 SIMD 計算ユニット上に同じ実行プラットフォームを有するように、統一されたシェーダモデルを実装する。したがって、頂点シェーダ 2 1 5 を含むシェーダは、本明細書では統一されたシェーダプール 2 1 6 と呼ばれるリソースの共通セットを使用して実装される。

【 0 0 1 5 】

ハルシェーダ 2 1 8 は、入力パッチを定義するために使用される入力高次パッチ又は制御ポイント上で動作する。ハルシェーダ 2 1 8 は、テッセレーション係数及び他のパッチデータを出力する。いくつかの実施形態では、ハルシェーダ 2 1 8 によって生成されたプリミティブは、テッセレータ 2 2 0 に提供される。テッセレータ 2 2 0 は、ハルシェーダ 2 1 8 からオブジェクト（パッチ等）を受信し、例えば、ハルシェーダ 2 1 8 によってテッセレータ 2 2 0 に提供されたテッセレーション係数に基づいて、入力オブジェクトをテッセレーションすることにより、入力オブジェクトに対応するプリミティブを識別する情報を生成する。テッセレーションは、例えば、テッセレーションプロセスによって生成されたプリミティブの粒度を指定するテッセレーション係数によって示されるように、パッチ等の入力高次プリミティブを、より細かいレベルの詳細を表す低次出力プリミティブのセットに細分する。したがって、シーンのモデルは、（メモリ又は帯域幅を節約するため）より少数の高次プリミティブによって表され、追加の詳細は、高次プリミティブをテッセレーションすることによって追加される。

【 0 0 1 6 】

ドメインシェーダ 2 2 4 は、ドメインの場所及び（任意選択的に）他のパッチデータを入力する。ドメインシェーダ 2 2 4 は、提供された情報で動作し、入力ドメインの場所及び他の情報に基づいて、出力のための単一の頂点を生成する。図示した実施形態では、ドメインシェーダ 2 2 4 は、三角形 2 1 1 及びテッセレーション係数に基づいてプリミティブ 2 2 2 を生成する。ジオメトリシェーダ 2 2 6 は、入力プリミティブを受信し、入力プリミティブに基づいてジオメトリシェーダ 2 2 6 によって生成される最大 4 つのプリミティブを出力する。図示した実施形態では、ジオメトリシェーダ 2 2 6 は、テッセレートされたプリミティブ 2 2 2 に基づいて出力プリミティブ 2 2 8 を生成する。

【 0 0 1 7 】

プリミティブの 1 つのストリームが 1 つ以上のスキャンコンバータ 2 3 0 に提供され、いくつかの実施形態では、プリミティブの最大 4 つのストリームは、ストレージリソース 2 0 5 内のバッファに連結される。スキャンコンバータ 2 3 0 は、シェーディング動作、クリッピング、透視分割、切断及びビューポート選択等の他の動作を実行する。スキャンコンバータ 2 3 0 は、グラフィックスパイプライン 2 0 0 のピクセル処理部 2 0 2 において後で処理されるピクセルのセット 2 3 2 を生成する。

【 0 0 1 8 】

図示された実施形態では、ピクセルシェーダ 2 3 4 は、ピクセルフロー（例えば、ピクセルのセット 2 3 2 を含む）を入力し、入力ピクセルフローに応じて 0 又は別のピクセルフローを出力する。出力マージブロック 2 3 6 は、ピクセルシェーダ 2 3 4 から受信したピクセルに対してブレンド、深度、ステンシル又は他の動作を実行する。

【 0 0 1 9 】

10

20

30

40

50

グラフィックスパイプライン 200 内のシェーダの一部又は全部は、ストレージリソース 205 に記憶されたテクスチャデータを使用してテクスチャマッピングを実行する。例えば、ピクセルシェーダ 234 は、ストレージリソース 205 からテクスチャデータを読み取り、テクスチャデータを使用して 1 つ以上のピクセルをシェーディングすることができる。次いで、シェーディングされたピクセルは、ユーザに提示するためにディスプレイに提供される。本明細書で説明するように、グラフィックスパイプライン 200 内のシェーダによって使用されるテクスチャデータは、キャッシュ 207 を使用してキャッシュされる。ミス要求は、キャッシュ 207 におけるキャッシュミスに応じて、例えばキャッシュ 207 のキャッシュラインの部分におけるアドレスの位置、要求又はデータのヒューリスティック又は特性、時間的局所性、空間的局所性等に基づいて、選択的に生成される。

10

#### 【0020】

図 3 は、いくつかの実施形態による、第 1 の読み取りサイクル 301 において複数のセクタにわたって分散された要求と、第 2 の読み取りサイクル 302 において単一のセクタに制約された要求と、を有するキャッシュラインのブロック図である。キャッシュライン 300、305 は、図 1 に示されるキャッシュ 151 ~ 153 のいくつかの実施形態及び図 2 に示されるキャッシュ 207 のいくつかの実施形態におけるキャッシュラインを表す。キャッシュライン 300、305 は、グラフィックス処理のためのテクスチャを記憶するために使用され、したがって、キャッシュラインは比較的大きい。例えば、キャッシュライン 300、305 の各々は、対応する計算ユニット又は他のプロセッサ、プロセッサコア、処理要素等によるアクセスのために 128 バイト、すなわち、1 K ビットのデータを記憶することができる。図示された実施形態では、キャッシュライン 300、305 は 2 つのセクタに分割される。しかしながら、いくつかの実施形態では、キャッシュライン 300、305 は、3 つ以上のセクタに分割される。

20

#### 【0021】

第 1 の読み取りサイクル 301 中に、キャッシュライン 300 は、キャッシュライン 300 によって記憶されたバイトのサブセットを保持する位置 310 を示すアドレスへの読み取り要求を受信する。位置 310 は、キャッシュライン 300 の第 1 のセクタ 311 内にある。また、キャッシュライン 300 は、キャッシュライン 300 によって記憶されたバイトの別のサブセットを保持する位置 315 を示すアドレスへの読み取り要求を受信する。位置 315 は、キャッシュライン 300 の第 2 のセクタ 312 内にある。図示した実施形態では、位置 315 及び位置 310 への読み取り要求は、キャッシュライン 300 においてミスする。

30

#### 【0022】

第 2 の読み取りサイクル 302 中に、キャッシュライン 305 は、キャッシュライン 305 によって記憶されたバイトのサブセットを保持する位置 320 を示すアドレスへの読み取り要求を受信する。位置 320 は、キャッシュライン 305 の第 1 のセクタ 321 内にあり、位置 320 の何れもキャッシュライン 305 の第 2 のセクタ 322 内にはない。図示した実施形態では、位置 320 への読み取り要求は、キャッシュライン 305 においてミスする。

#### 【0023】

ミス要求は、キャッシュライン 300、305 におけるキャッシュミスの位置に基づいて、第 1 のセクタ 311、321、第 2 のセクタ 312、322、又は、両方のセクタ 311、312、321、322 (例えば、フルキャッシュライン 300 及び 305) に対して選択的に生成される。いくつかの実施形態では、キャッシュミスの他のヒューリスティック又は特性も、本明細書で説明されるように、ミス要求が第 1 のセクタ 311、321、第 2 のセクタ 312、322、又は、両方のセクタ 311、312、321、322 に対して生成されるか否かを判定するために使用される。例えば、キャッシュミスの測定、予想又は予測された空間的局所性又は時間的局所性を使用して、ミス要求がどのように生成されるかを判定することができる。図示した実施形態では、第 1 のセクタ 311 内の位置 310 及び第 2 のセクタ 312 内の位置 315 を含む第 1 の読み取りサイクル 301

40

50

内のキャッシュミスに応じて、フルキャッシュライン 3 0 0（例えば、セクタ 3 1 1 及び 3 1 2）に対するミス要求が生成される。第 2 の読み取りサイクル 3 0 2 中に、第 2 の読み取りサイクル 3 0 2 内のキャッシュミスが第 1 のセクタ 3 2 1 内のみにある位置 3 2 0 に対するものであることに応じて、ミス要求が第 1 のセクタ 3 2 1 に対してのみ生成される（第 2 のセクタ 3 2 2 に対するミス要求の生成はバイパスされる）。

#### 【 0 0 2 4 】

図 4 は、いくつかの実施形態による、第 1 の読み取りサイクル 4 0 1 及び第 2 の読み取りサイクル 4 0 2 に高度の時間的局所性を示さない要求を有するキャッシュライン 4 0 0 のブロック図である。キャッシュライン 4 0 0 は、図 1 に示されるキャッシュ 1 5 1 ~ 1 5 3 のいくつかの実施形態及び図 2 に示されるキャッシュ 2 0 7 のいくつかの実施形態におけるキャッシュラインを表す。キャッシュライン 4 0 0 は、グラフィックス処理のためのテクスチャを記憶するために使用され、したがって、キャッシュラインは、比較的大きく、例えば、1 2 8 バイトのデータである。図示された実施形態では、キャッシュライン 4 0 0 はセクタ 4 1 1、4 1 2 に分割される。しかしながら、いくつかの実施形態では、キャッシュライン 4 0 0 は、3 つ以上のセクタに分割される。

10

#### 【 0 0 2 5 】

第 1 の読み取りサイクル 4 0 1 中に、キャッシュライン 4 0 0 は、キャッシュライン 4 0 0 によって記憶されたバイトのサブセットを保持する位置 4 1 5 を示すアドレスへの読み取り要求を受信する。位置 4 1 5 は、全て、第 1 のセクタ 4 1 1 内に見出される。第 2 の読み取りサイクル 4 0 2 中に、キャッシュライン 4 0 0 は、キャッシュライン 4 0 0 によって記憶されたバイトのサブセットを保持する位置 4 2 0 を示すアドレスへの読み取り要求を受信する。位置 4 2 0 は、全て、キャッシュライン 4 0 0 の第 2 のセクタ 4 1 2 内にある。図示した実施形態では、ミス要求は、キャッシュミスの実際の又は予測された時間的局所性に少なくとも部分的に基づいて、第 1 のセクタ 4 1 1 又は第 2 のセクタ 4 1 2 に対して選択的に生成される。したがって、連続する読み取りサイクル（例えば、第 1 の読み取りサイクル 4 0 1 及び第 2 の読み取りサイクル 4 0 2）におけるキャッシュミスが第 1 のセクタ 4 1 1 及び第 2 のセクタ 4 1 2 にわたって分散されるので、ミス要求がフルキャッシュライン（例えば、第 1 のセクタ 4 1 1 及び第 2 のセクタ 4 1 2）に対して生成される。対照的に、キャッシュライン 4 0 0 への読み取り要求が高度の時間的局所性を示す場合、例えば、複数のサイクル中の読み取り要求が、セクタ 4 1 1、4 1 2 のうち何れかのみ位置するアドレスに対するものであると予想又は予測される場合、ミス要求は、対応するセクタに対してのみ生成される。

20

30

#### 【 0 0 2 6 】

また、ミス要求は、読み取り要求の予測された空間的局所性に基づいて、キャッシュライン 4 0 0 の異なる部分に対して生成される。例えば、ピクセルシェーダがスクリーンにわたってスキャンしている場合、読み取り要求シーケンスは、隣接する又は近接するピクセル位置に関連付けられたローカルアドレスに対するものである可能性が高い。したがって、メモリアクセスシステムの効率及び性能は、後続の読み取り要求がキャッシュラインの他のセクタ内にあり得る近くのアドレスに対するものである可能性が高いので、現在の読み取りサイクル中のキャッシュミスが単一のセクタ内のみ（又は主に）にある場合でも、予測された空間的局所性が高い（例えば、閾値を上回る）場合、フルキャッシュラインをフェッチすることによって改善される可能性が高い。対照的に、予測された空間的局所性が低い（例えば、閾値を下回る）場合、現在の読み取りサイクル中のキャッシュミスが単一のセクタ内の位置に対するものである場合、ミス要求は単一のセクタに対してのみ生成され得る。いくつかの実施形態では、読み取り要求又はミス要求に関連付けられた情報は、ヒステリシスウィンドウのために保持され、ヒステリシスウィンドウ内の情報は、読み取り要求又はミス要求の時間的局所性又は空間的局所性を判定又は予測するために使用される。

40

#### 【 0 0 2 7 】

図 5 は、いくつかの実施形態による、キャッシュラインの部分に対するミス要求を選択

50

的に生成する方法 5 0 0 のフロー図である。方法 5 0 0 は、図 1 に示される処理システム 1 0 0 及び図 2 に示されるグラフィックスパイプライン 2 0 0 のいくつかの実施形態で実施される。キャッシュラインは、2 つのセクタ（又は部分若しくは半分）を含むが、いくつかの実施形態では、キャッシュラインは、より多くのセクタを含む。

#### 【 0 0 2 8 】

ブロック 5 0 5 において、キャッシュは、要求サイクル中にスレッドからキャッシュラインへの読み取り要求を受信する。読み取り要求は、キャッシュライン及び対応するメモリ内の位置を示すアドレスを含む。図示した実施形態では、読み取り要求はキャッシュライン内でミスし、これが、要求されたデータをバッキングメモリ又はより高レベルのキャッシュからフェッチするためのミス要求の選択的生成をトリガする。

10

#### 【 0 0 2 9 】

判定ブロック 5 1 0 において、キャッシュは、キャッシュミスを生じたスレッドが現在の要求サイクル中にキャッシュラインの両方のセクタ内の位置にマッピングする否かを判定する。マッピングする場合、方法 5 0 0 はブロック 5 2 0 に進む。マッピングしない場合、方法 5 0 0 は判定ブロック 5 1 5 に進む。

#### 【 0 0 3 0 】

判定ブロック 5 1 5 において、キャッシュは、特定のヒューリスティックに基づいて空間的局所性又は時間的局所性の尤度を判定する。いくつかの実施形態では、空間的局所性又は時間的局所性の尤度は、固定ヒューリスティック又はプログラマブルヒューリスティックとの一致があるか否かに基づいて判定される。例えば、キャッシュライン内の情報が色圧縮又は深度圧縮を使用して生成されるか否かにより、関連するデータが高度の空間的局所性を有し、後続の読み取り要求がキャッシュラインの両方のセクタ内のアドレスを含む可能性が高いことを示す。したがって、情報が高度の局所性を有すると予想される場合、方法 5 0 0 はブロック 5 2 0 に進む。そうでない場合、方法 5 0 0 はブロック 5 2 5 に進む。

20

#### 【 0 0 3 1 】

ブロック 5 2 0 において、キャッシュは、フルキャッシュラインに対するミス要求を生成する、すなわち、キャッシュは、キャッシュラインの全てのセクタを含むミス要求を生成する。ブロック 5 2 5 において、キャッシュは、要求サイクル中のスレッドに対するキャッシュミスにおけるアドレスに対応する位置を含むキャッシュラインのセクタ（又は部分若しくは半分）に対するミス要求を生成する。

30

#### 【 0 0 3 2 】

本明細書で開示するように、いくつかの実施形態では、装置は、複数のサブセットに分割されるキャッシュラインを含むテクスチャキャッシュと、グラフィックスパイプライン内の少なくとも 1 つの計算ユニットと、を含み、プロセッサは、キャッシュラインの第 1 のサブセットに関連付けられたアドレスへのメモリアクセス要求に対するキャッシュミスに応じて、テクスチャキャッシュ内のキャッシュラインの複数のサブセットのうち第 1 のサブセットに対するミス要求を選択的に生成するように構成されている。一態様では、少なくとも 1 つの計算ユニットは、メモリアクセス要求に関連付けられたキャッシュミスが、複数のサブセットのうち第 1 のサブセットのみにマッピングするの、又は、複数のサブセットのうち第 1 のサブセットに追加の又は第 1 のサブセット以外の 1 つ以上のサブセットにマッピングするかを判定するように構成されている。別の態様では、少なくとも 1 つの計算ユニットは、複数のサブセットのうち第 1 のサブセットに追加の又は第 1 のサブセット以外のサブセットへのキャッシュミスのマッピングに応じて、フルキャッシュラインに対するミス要求を生成するように構成されている。更に別の態様では、少なくとも 1 つの計算ユニットは、メモリアクセス要求が第 1 のサブセットのみにマッピングすることに応じて、色圧縮及び深度圧縮のうち少なくとも 1 つがテクスチャデータに対して有効にされているか否かを判定するように構成されている。

40

#### 【 0 0 3 3 】

一態様では、少なくとも 1 つの計算ユニットは、色圧縮及び深度圧縮のうち少なくとも

50

1つがテクスチャデータに対して有効にされていることに応じて、フルキャッシュラインに対するミス要求を生成するように構成されている。別の態様では、少なくとも1つの計算ユニットは、色圧縮及び深度圧縮のうち少なくとも1つがテクスチャデータに対して有効にされていないことに応じて、キャッシュラインの第1のサブセットに対するミス要求を生成するように構成されている。更に別の態様では、少なくとも1つの計算ユニットは、メモリアクセス要求の時間的局所性及び空間的局所性のうち少なくとも1つに基づいて、第1のサブセット又は複数のサブセットへのミス要求を選択的に生成するように構成されている。

#### 【0034】

一態様では、少なくとも1つの計算ユニットは、メモリアクセス要求シーケンスが複数のサブセットにアクセスすることが予想されることに応じて、第1のサブセットにおけるキャッシュミスに応じて、複数のサブセットに対するミス要求を生成するように構成されている。別の態様では、少なくとも1つの計算ユニットは、第1のサブセット内のキャッシュミスに応じて、且つ、閾値を上回る空間的局所性を有するメモリアクセス要求シーケンスに応じて、複数のサブセットに対するミス要求を生成するように構成されている。更に別の態様では、少なくとも1つの計算ユニットは、メモリアクセス要求が閾値を下回る空間的局所性を有することに応じて、第1のサブセットに対するミス要求を生成するように構成されている。

10

#### 【0035】

いくつかの実施形態では、方法は、複数のサブセットに分割されるキャッシュラインを含むテクスチャキャッシュ内のキャッシュラインに対するミス要求を検出することと、キャッシュミスがキャッシュラインの第1のサブセットに関連付けられたアドレスに対するものであることに応じて、テクスチャキャッシュ内のキャッシュラインの複数のサブセットのうち第1のサブセットに対するミス要求を選択的に生成することと、を含む。一態様では、本方法は、メモリアクセス要求に関連付けられたキャッシュミスが、複数のサブセットのうち第1のサブセットのみにマッピングするのか、又は、複数のサブセットのうち第1のサブセットに追加の又は第1のサブセット以外の1つ以上のサブセットにマッピングするかを判定することを含む。別の態様では、本方法は、複数のサブセットへのキャッシュミスのマッピングに応じて、フルキャッシュラインに対するミス要求を生成することを含む。更に別の態様では、本方法は、色圧縮及び深度圧縮のうち少なくとも1つがテクスチャデータに対して有効にされているか否かを判定することを含む。

20

30

#### 【0036】

一態様では、ミス要求を選択的に生成することは、色圧縮及び深度圧縮のうち少なくとも1つがテクスチャデータに対して有効にされていることに応じて、複数のサブセットに対するミス要求を生成することを含む。別の態様では、ミス要求を選択的に生成することは、色圧縮及び深度圧縮のうち少なくとも1つがテクスチャデータに対して有効にされていないことに応じて、キャッシュラインの第1のサブセットに対するミス要求を生成することを含む。別の態様では、ミス要求を選択的に生成することは、第1のサブセットにおけるキャッシュミスに応じて、且つ、メモリアクセス要求シーケンスがキャッシュラインの異なるセクタにアクセスすることが予想されることに応じて、複数のサブセットに対するミス要求を生成することを含む。

40

#### 【0037】

一態様では、ミス要求を選択的に生成することは、第1のサブセットにおけるキャッシュミスに応じて、且つ、閾値を上回る空間的局所性を有するメモリアクセス要求シーケンスに応じて、複数のサブセットに対するミス要求を生成することを含む。別の態様では、ミス要求を選択的に生成することは、メモリアクセス要求が閾値を下回る空間的局所性を有することに応じて、第1のサブセットに対するミス要求を生成することを含む。

#### 【0038】

いくつかの実施形態では、装置は、第1のセクタ及び第2のセクタに分割されるキャッシュラインを含むテクスチャキャッシュと、グラフィックスパイプライン内の少なくとも

50

1つの計算ユニットと、を含み、少なくとも1つの計算ユニットは、要求サイクル中に受信されたメモリアクセス要求に対するキャッシュミスが第1のセクタ中にあることに応じて、第1のセクタに対するミス要求を生成し、第2のセクタに対するミス要求の生成をバイパスするように構成されている。一態様では、少なくとも1つの計算ユニットは、要求サイクル中に受信されたメモリアクセス要求に対するキャッシュミスが第1のセクタ及び第2のセクタ内にあることに応じて、第1のセクタ及び第2のセクタに対するミス要求を生成するように構成されている。

【0039】

コンピュータ可読記憶媒体は、命令及び/又はデータをコンピュータシステムに提供するために、使用中にコンピュータシステムによってアクセス可能な任意の非一時的な記憶媒体又は非一時的な記憶媒体の組み合わせを含む。このような記憶媒体には、限定されないが、光学媒体（例えば、コンパクトディスク（CD）、デジタル多用途ディスク（DVD）、ブルーレイ（登録商標）ディスク）、磁気媒体（例えば、フロッピー（登録商標）ディスク、磁気テープ、磁気ハードドライブ）、揮発性メモリ（例えば、ランダムアクセスメモリ（RAM）若しくはキャッシュ）、不揮発性メモリ（例えば、読取専用メモリ（ROM）若しくはフラッシュメモリ）、又は、微小電気機械システム（MEMS）ベースの記憶媒体が含まれ得る。コンピュータ可読記憶媒体（例えば、システムRAM又はROM）はコンピューティングシステムに内蔵されてもよいし、コンピュータ可読記憶媒体（例えば、磁気ハードドライブ）はコンピューティングシステムに固定的に取り付けられてもよいし、コンピュータ可読記憶媒体（例えば、光学ディスク又はユニバーサルシリアルバス（USB）ベースのフラッシュメモリ）はコンピューティングシステムに着脱可能に取り付けられてもよいし、コンピュータ可読記憶媒体（例えば、ネットワークアクセス可能ストレージ（NAS））は有線又は無線ネットワークを介してコンピュータシステムに結合されてもよい。

【0040】

いくつかの実施形態では、上述した技術の特定の態様は、ソフトウェアを実行する処理システムの1つ以上のプロセッサによって実装される。ソフトウェアは、非一時的なコンピュータ可読記憶媒体に記憶されるか、別の方法で明確に具体化された実行可能命令の1つ以上のセットを含む。ソフトウェアは、命令及び特定のデータを含んでもよく、当該命令及び特定のデータは、1つ以上のプロセッサによって実行されると、上述した技術の1つ以上の態様を実行するように1つ以上のプロセッサを操作する。非一時的なコンピュータ可読記憶媒体は、例えば、磁気又は光ディスク記憶デバイス、フラッシュメモリ等のソリッドステート記憶デバイス、キャッシュ、ランダムアクセスメモリ（RAM）、又は、他の不揮発性メモリデバイス（単数又は複数）等を含み得る。非一時的なコンピュータ可読記憶媒体に記憶された実行可能命令は、ソースコード、アセンブリ言語コード、オブジェクトコード、又は、1つ以上のプロセッサによって解釈され若しくは別の方法で実行可能な他の命令形式で実装可能である。

【0041】

上述したものに加えて、概要説明において説明した全てのアクティビティ又は要素が必要とされているわけではなく、特定のアクティビティ又はデバイスの一部が必要とされない場合があり、1つ以上のさらなるアクティビティが実行される場合があり、1つ以上のさらなる要素が含まれる場合があることに留意されたい。さらに、アクティビティが列挙された順序は、必ずしもそれらが実行される順序ではない。また、概念は、特定の実施形態を参照して説明された。しかしながら、当業者であれば、特許請求の範囲に記載されているような本発明の範囲から逸脱することなく、様々な変更及び変形を行うことができるのを理解するであろう。したがって、明細書及び図面は、限定的な意味ではなく例示的な意味で考慮されるべきであり、これらの変更形態の全ては、本発明の範囲内に含まれることが意図される。

【0042】

利益、他の利点及び問題に対する解決手段を、特定の実施形態に関して上述した。しか

10

20

30

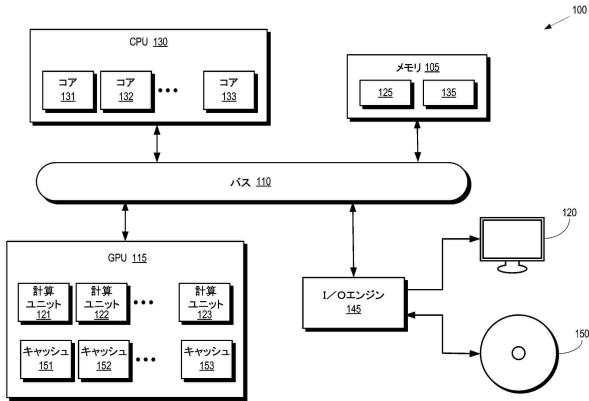
40

50

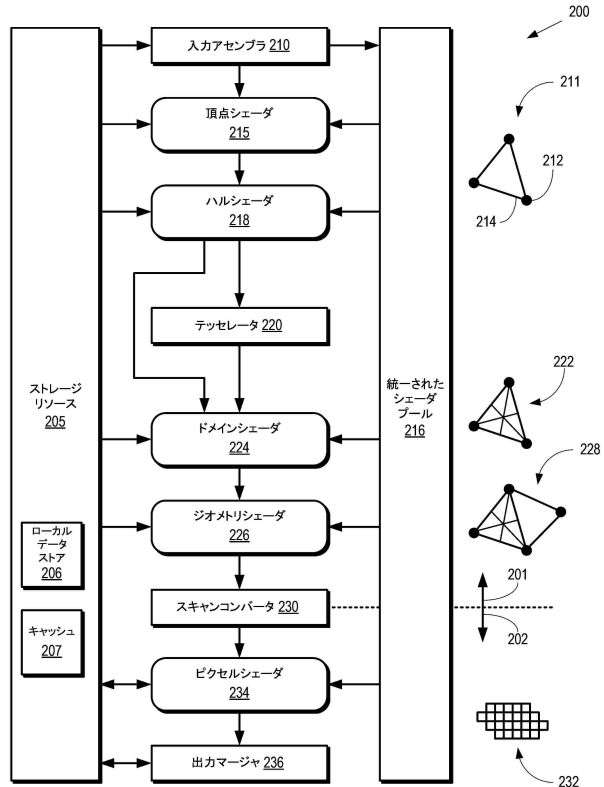
し、利益、利点、問題に対する解決手段、及び、何かしらの利益、利点若しくは解決手段が発生又は顕在化する可能性のある特徴は、何れか若しくは全ての請求項に重要な、必須の、又は、不可欠な特徴と解釈されない。さらに、開示された発明は、本明細書の教示の利益を有する当業者には明らかな方法であって、異なっているが同様の方法で修正され実施され得ることから、上述した特定の実施形態は例示にすぎない。添付の特許請求の範囲に記載されている以外に本明細書に示されている構成又は設計の詳細については限定がない。したがって、上述した特定の実施形態は、変更又は修正されてもよく、かかる変更形態の全ては、開示された発明の範囲内にあると考えられることが明らかである。したがって、ここで要求される保護は、添付の特許請求の範囲に記載されている。

【図面】

【図 1】



【図 2】



10

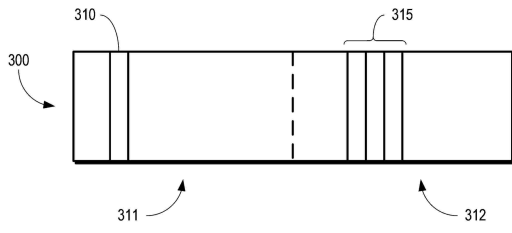
20

30

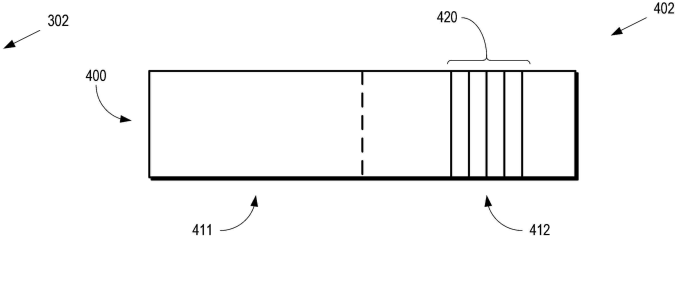
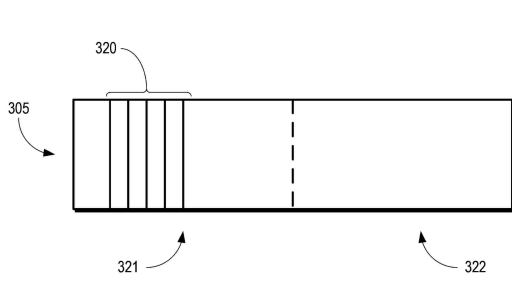
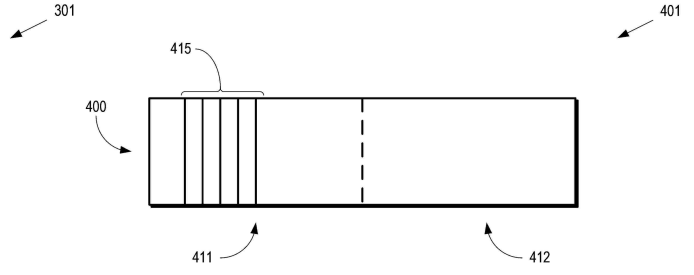
40

50

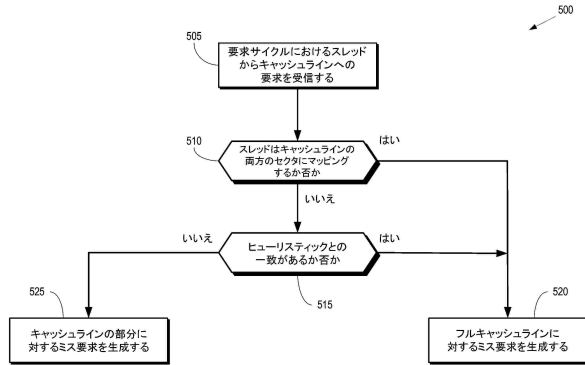
【 図 3 】



【 図 4 】



【 図 5 】



10

20

30

40

50

## フロントページの続き

- 弁理士 早川 裕司  
(74)代理人 100111615  
弁理士 佐野 良太  
(74)代理人 100162156  
弁理士 村雨 圭介  
(72)発明者 ファタネー エフ . ゴッドラット  
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、オーガスティン ドライブ 2 4 8 5  
(72)発明者 スティーブン ダブリュ . ソモギ  
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、オーガスティン ドライブ 2 4 8 5  
(72)発明者 チェンホン リウ  
大韓民国 1 6 6 7 7 ギョンギ - ド、スウォン - シ、ヨントン - ク、サムスン - 口、1 2 9  
審査官 渡部 幸和  
(56)参考文献 米国特許出願公開第 2 0 0 6 / 0 2 5 0 4 0 8 ( U S , A 1 )  
米国特許第 0 6 7 2 4 3 9 1 ( U S , B 1 )  
(58)調査した分野 (Int.Cl. , D B 名)  
G 0 6 T 1 5 / 0 0  
G 0 6 F 1 2 / 0 0  
G 0 6 F 9 / 3 8