

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7250709号  
(P7250709)

(45)発行日 令和5年4月3日(2023.4.3)

(24)登録日 令和5年3月24日(2023.3.24)

(51)国際特許分類	F I		
G 0 6 T 7/00 (2017.01)	G 0 6 T 7/00	3 5 0 C	
G 0 6 T 7/593(2017.01)	G 0 6 T 7/593		
G 0 6 T 7/70 (2017.01)	G 0 6 T 7/70	Z	
G 0 6 N 3/04 (2023.01)	G 0 6 N 3/04		
G 0 6 N 3/08 (2023.01)	G 0 6 N 3/08		

請求項の数 14 (全23頁)

(21)出願番号	特願2019-571451(P2019-571451)	(73)特許権者	514108838 マジック リープ, インコーポレイテッド Magic Leap, Inc. アメリカ合衆国 フロリダ 33322, プランテーション, ウエスト サンライズ ブルバード 7500 7500 W SUNRISE BLVD , PLANTATION, FL 333 22 USA
(86)(22)出願日	平成30年6月27日(2018.6.27)	(74)代理人	100078282 弁理士 山本 秀策
(65)公表番号	特表2020-526818(P2020-526818 A)	(74)代理人	100113413 弁理士 森下 夏樹
(43)公表日	令和2年8月31日(2020.8.31)	(74)代理人	100181674 弁理士 飯田 貴敏
(86)国際出願番号	PCT/US2018/039804		
(87)国際公開番号	WO2019/005999		
(87)国際公開日	平成31年1月3日(2019.1.3)		
審査請求日	令和3年6月16日(2021.6.16)		
(31)優先権主張番号	62/526,203		
(32)優先日	平成29年6月28日(2017.6.28)		
(33)優先権主張国・地域又は機関	米国(US)		

最終頁に続く

(54)【発明の名称】 畳み込み画像変換を使用して同時位置特定およびマッピングを実施する方法およびシステム

(57)【特許請求の範囲】

【請求項1】

2つの画像に基づいてホモグラフィを算出する方法であって、前記方法は、  
 第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することと、  
 前記第1の画像に基づく第1の2D点群と、前記第2の画像に基づく第2の2D点群とを生成することと、  
 前記第1の2D点群および前記第2の2D点群をニューラルネットワークに提供することと、  
 前記ニューラルネットワークによって、前記ニューラルネットワークへの入力として前記第1の2D点群および前記第2の2D点群を提供したことに応答して、前記ホモグラフィを生成することと  
 を含み、  
 前記ニューラルネットワークは、  
 複数の点を含む3D点群を生成することと、  
 2つのカメラ姿勢を取得することであって、前記2つのカメラ姿勢から前記複数の点が少なくとも部分的に視認可能であり、前記2つのカメラ姿勢は、前記第1のカメラ姿勢および前記第2のカメラ姿勢とは異なる、ことと、  
 前記2つのカメラ姿勢を使用して前記複数の点を2つの2D平面上に投影することにより、2つの2D点群を生成することと、

前記ニューラルネットワークによって、前記ニューラルネットワークへの入力として前記2つの2D点群を提供したことに応答して、特定のホモグラフィを生成することと、  
前記2つのカメラ姿勢に基づいて、グラウンドトゥールースホモグラフィを決定することと、

前記特定のホモグラフィおよび前記グラウンドトゥールースホモグラフィを使用して前記ニューラルネットワークを修正することと

によって以前に訓練されたものである、方法。

【請求項2】

前記第1の画像は、第1の瞬間において第1のカメラによって捕捉され、前記第2の画像は、前記第1の瞬間後の第2の瞬間において前記第1のカメラによって捕捉されている、請求項1に記載の方法。

10

【請求項3】

前記第1の2D点群および前記第2の2D点群は、第1のニューラルネットワークを使用して生成され、前記ニューラルネットワークは、第2のニューラルネットワークである、請求項1に記載の方法。

【請求項4】

前記3D点群は、1つ以上の幾何学形状をサンプリングすることによって生成される、請求項1に記載の方法。

【請求項5】

前記2つのカメラ姿勢は、少なくとも30%の視覚的重複を有する、請求項1に記載の方法。

20

【請求項6】

拡張現実（AR）デバイスであって、前記ARデバイスは、  
 カメラと、  
 前記カメラに通信可能に結合されたプロセッサと  
 を備え、

前記プロセッサは、  
 前記カメラから、第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することと、

前記第1の画像に基づく第1の2D点群と、前記第2の画像に基づく第2の2D点群とを生成することと、

30

前記第1の2D点群および前記第2の2D点群をニューラルネットワークに提供することと、

前記ニューラルネットワークによって、前記ニューラルネットワークへの入力として前記第1の2D点群および前記第2の2D点群を提供したことに応答して、ホモグラフィを生成することと

を含む動作を実施するように構成されており、

前記ニューラルネットワークは、  
複数の点を含む3D点群を生成することと、

2つのカメラ姿勢を取得することであって、前記2つのカメラ姿勢から前記複数の点が少なくとも部分的に視認可能であり、前記2つのカメラ姿勢は、前記第1のカメラ姿勢および前記第2のカメラ姿勢とは異なる、ことと、

40

前記2つのカメラ姿勢を使用して前記複数の点を2つの2D平面上に投影することにより、2つの2D点群を生成することと、

前記ニューラルネットワークによって、前記ニューラルネットワークへの入力として前記2つの2D点群を提供したことに応答して、特定のホモグラフィを生成することと、

前記2つのカメラ姿勢に基づいて、グラウンドトゥールースホモグラフィを決定することと、

前記特定のホモグラフィおよび前記グラウンドトゥールースホモグラフィを使用して前記ニューラルネットワークを修正することと

50

によって以前に訓練されたものである、ARデバイス。

【請求項 7】

前記第 1 の 2D 点群 および前記第 2 の 2D 点群 は、第 1 のニューラルネットワークを使用して生成され、前記ニューラルネットワークは、第 2 のニューラルネットワークである、請求項 6 に記載の AR デバイス。

【請求項 8】

前記 3D 点群 は、1 つ以上の幾何学形状をサンプリングすることによって生成される、請求項 6 に記載の AR デバイス。

【請求項 9】

前記 2 つのカメラ姿勢 は、少なくとも 30 % の視覚的重複 を有する、請求項 6 に記載の AR デバイス。 10

【請求項 10】

命令を備えている非一過性コンピュータ読み取り可能な媒体であって、前記命令は、プロセッサによって実行されると、

第 1 のカメラ姿勢に基づく第 1 の画像と、第 2 のカメラ姿勢に基づく第 2 の画像とを受信することと、

前記第 1 の画像に基づく第 1 の 2D 点群 と、前記第 2 の画像に基づく第 2 の 2D 点群 とを生成することと、

前記第 1 の 2D 点群 および前記第 2 の 2D 点群 をニューラルネットワークに提供することと、 20

前記ニューラルネットワークによって、前記ニューラルネットワークへの入力として前記第 1 の 2D 点群 および前記第 2 の 2D 点群を提供したことに応答して、ホモグラフィを生成することと

を含む動作を前記プロセッサに実施させ、

前記ニューラルネットワークは、

複数の点を含む 3D 点群を生成することと、

2 つのカメラ姿勢を取得することであって、前記 2 つのカメラ姿勢から前記複数の点が少なくとも部分的に視認可能であり、前記 2 つのカメラ姿勢は、前記第 1 のカメラ姿勢および前記第 2 のカメラ姿勢とは異なる、ことと、

前記 2 つのカメラ姿勢を使用して前記複数の点を 2 つの 2D 平面上に投影することにより、2 つの 2D 点群を生成することと、 30

前記ニューラルネットワークによって、前記ニューラルネットワークへの入力として前記 2 つの 2D 点群を提供したことに応答して、特定のホモグラフィを生成することと、

前記 2 つのカメラ姿勢に基づいて、グラウンドトゥールスホモグラフィを決定することと、

前記特定のホモグラフィおよび前記グラウンドトゥールスホモグラフィを使用して前記ニューラルネットワークを修正することと

によって以前に訓練されたものである、非一過性コンピュータ読み取り可能な媒体。

【請求項 11】

前記第 1 の画像は、第 1 の瞬間において第 1 のカメラによって捕捉され、前記第 2 の画像は、前記第 1 の瞬間後の第 2 の瞬間において前記第 1 のカメラによって捕捉されている、請求項 10 に記載の非一過性コンピュータ読み取り可能な媒体。 40

【請求項 12】

前記第 1 の 2D 点群 および前記第 2 の 2D 点群 は、第 1 のニューラルネットワークを使用して生成され、前記ニューラルネットワークは、第 2 のニューラルネットワークである、請求項 10 に記載の非一過性コンピュータ読み取り可能な媒体。

【請求項 13】

前記 3D 点群 は、1 つ以上の幾何学形状をサンプリングすることによって生成される、請求項 10 に記載の非一過性コンピュータ読み取り可能な媒体。

【請求項 14】

前記2つのカメラ姿勢は、少なくとも30%の視覚的重複を有する、請求項10に記載の非一過性コンピュータ読み取り可能な媒体。

【発明の詳細な説明】

【技術分野】

【0001】

(関連出願の相互参照)

本願は、その内容が参照することによってその全体として本明細書に組み込まれる、2017年6月28日に出願され、「METHOD AND SYSTEM FOR PERFORMING SIMULTANEOUS LOCALIZATION AND MAPPING USING CONVOLUTIONAL IMAGE TRANSFORMATION」と題された、米国仮特許出願第62/526,203号の非仮出願であり、その優先権の利益を主張する。

10

【背景技術】

【0002】

現代のコンピューティングおよびディスプレイ技術は、いわゆる「仮想現実」または「拡張現実」体験のためのシステムの開発を促進しており、デジタル的に再現された画像またはその一部が、現実であるように見える、もしくはそのように知覚され得る様式でユーザに提示される。仮想現実または「VR」シナリオは、典型的に、他の実際の実世界の視覚的入力に対する透過性を伴わずに、デジタルまたは仮想画像情報の提示を伴い、拡張現実または「AR」シナリオは、典型的に、ユーザの周囲の実際の世界の可視化への拡張として、デジタルまたは仮想画像情報の提示を伴う。

20

【0003】

これらのディスプレイ技術において行われる進歩にもかかわらず、当技術分野において、拡張現実システム、特に、ディスプレイシステムに関連する改良された方法、システム、およびデバイスの必要がある。

【発明の概要】

【課題を解決するための手段】

【0004】

本開示は、概して、同時位置特定およびマッピング(SLAM)を実施するためのシステムおよび方法に関する。より具体的に、本開示の実施形態は、頭部搭載型仮想現実(VR)、複合現実(MR)、および/または拡張現実(AR)デバイスにおいて、畳み込み画像変換を使用して、SLAMを実施するためのシステムおよび方法を提供する。本発明の実施形態は、ユーザによって装着されるデバイスによって捕捉された画像を分析し、それによって、表示される仮想コンテンツの正確度を改良することによって、ユーザ/デバイス移動の正確な検出を可能にする。本発明は、ARデバイスを参照して説明され得るが、本開示は、コンピュータビジョンおよび画像ディスプレイシステムにおける種々の用途にも適用可能である。

30

【0005】

本発明の第1の側面では、2つの画像に基づいてホモグラフィを算出する方法が、提供される。方法は、第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することを含み得る。方法は、第1の画像に基づく第1の点群と、第2の画像に基づく第2の点群とを生成することも含み得る。方法は、第1の点群および第2の点群をニューラルネットワークに提供することをさらに含み得る。方法は、ニューラルネットワークによって、第1の点群および第2の点群に基づいて、ホモグラフィを生成することをさらに含み得る。いくつかの実施形態では、第1の点群および第2の点群は、2次元(2D)点群である。いくつかの実施形態では、第1の画像は、第1の瞬間において第1のカメラによって捕捉されている。いくつかの実施形態では、第2の画像は、第1の瞬間後の第2の瞬間において第1のカメラによって捕捉されている。いくつかの実施形態では、第1の点群および第2の点群は、第1のニューラルネットワークを使用して生成され、ニューラルネットワークは、第2のニューラルネットワークである。

40

50

## 【 0 0 0 6 】

いくつかの実施形態では、ニューラルネットワークは、複数の点を含む1つ以上の3D点群のうちの各3次元(3D)点群のために、複数の点の閾値距離内の3D軌道を決定することと、3D軌道をサンプリングし、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢を取得することとであって、複数の点は、少なくとも部分的に特定の第1のカメラ姿勢および特定の第2のカメラ姿勢から視認可能である、ことと、特定の第1のカメラ姿勢に基づいて、複数の点を第1の2D平面上に投影することと、第1の2D点群を生成し、特定の第2のカメラ姿勢に基づいて、複数の点を第2の2D平面上に投影することと、第2の2D点群を生成し、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢に基づいて、第1の2D点群と第2の2D点群との間のグラウンドトゥールスホモグラフィを決定することと、ニューラルネットワークによって、第1の2D点群および第2の2D点群に基づいて、特定のホモグラフィを生成することと、特定のホモグラフィをグラウンドトゥールスホモグラフィと比較することと、比較に基づいて、ニューラルネットワークを修正することとによって事前に訓練されている。いくつかの実施形態では、複数の3D点群は、1つ以上の幾何学形状をサンプリングすることによって生成される。いくつかの実施形態では、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢は、少なくとも30%重複を有する。

10

## 【 0 0 0 7 】

本発明の第2の側面では、ARデバイスが、提供される。ARデバイスは、カメラを含み得る。ARデバイスは、カメラに通信可能に結合されたプロセッサも含み、プロセッサは、カメラから、第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することと、第1の画像に基づく第1の点群と、第2の画像に基づく第2の点群とを生成することと、第1の点群および第2の点群をニューラルネットワークに提供することと、ニューラルネットワークによって、第1の点群および第2の点群に基づいて、ホモグラフィを生成することとを含む動作を実施するように構成され得る。いくつかの実施形態では、第1の点群および第2の点群は、2D点群である。いくつかの実施形態では、第1の点群および第2の点群は、第1のニューラルネットワークを使用して生成され、ニューラルネットワークは、第2のニューラルネットワークである。

20

## 【 0 0 0 8 】

本発明の第3の側面では、非一過性コンピュータ読み取り可能な媒体が、提供される。非一過性コンピュータ読み取り可能な媒体は、命令を含み得、命令は、プロセッサによって実行されると、第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することと、第1の画像に基づく第1の点群と、第2の画像に基づく第2の点群とを生成することと、第1の点群および第2の点群をニューラルネットワークに提供することと、ニューラルネットワークによって、第1の点群および第2の点群に基づいて、ホモグラフィを生成することとを含む動作をプロセッサに実施させる。いくつかの実施形態では、第1の点群および第2の点群は、2D点群である。いくつかの実施形態では、第1の画像は、第1の瞬間において第1のカメラによって捕捉され、第2の画像は、第1の瞬間後の第2の瞬間において第1のカメラによって捕捉されている。いくつかの実施形態では、第1の点群および第2の点群は、第1のニューラルネットワークを使用して生成され、ニューラルネットワークは、第2のニューラルネットワークである。

30

40

本明細書は、例えば、以下の項目も提供する。

(項目1)

2つの画像に基づいてホモグラフィを算出する方法であって、前記方法は、

第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することと、

前記第1の画像に基づく第1の点群と、前記第2の画像に基づく第2の点群とを生成することと、

前記第1の点群および前記第2の点群をニューラルネットワークに提供することと、

前記ニューラルネットワークによって、前記第1の点群および前記第2の点群に基づい

50

て、前記ホモグラフィを生成することと

を含む、方法。

(項目2)

前記第1の点群および前記第2の点群は、2次元(2D)点群である、項目1に記載の方法。

(項目3)

前記第1の画像は、第1の瞬間において第1のカメラによって捕捉され、前記第2の画像は、前記第1の瞬間後の第2の瞬間において前記第1のカメラによって捕捉されている、項目2に記載の方法。

(項目4)

前記第1の点群および前記第2の点群は、第1のニューラルネットワークを使用して生成され、前記ニューラルネットワークは、第2のニューラルネットワークである、項目2に記載の方法。

(項目5)

前記ニューラルネットワークは、  
複数の点を含む1つ以上の3D点群のうちの各3次元(3D)点群のために、  
前記複数の点の閾値距離内の3D軌道を決定することと、  
前記3D軌道をサンプリングし、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢を取得することであって、前記複数の点は、少なくとも部分的に前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢から視認可能である、ことと、

前記特定の第1のカメラ姿勢に基づいて、前記複数の点を第1の2D平面上に投影し、第1の2D点群を生成することと、

前記特定の第2のカメラ姿勢に基づいて、前記複数の点を第2の2D平面上に投影し、第2の2D点群を生成することと、

前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢に基づいて、前記第1の2D点群と前記第2の2D点群との間のグラウンドトゥルスホモグラフィを決定することと、

前記ニューラルネットワークによって、前記第1の2D点群および前記第2の2D点群に基づいて、特定のホモグラフィを生成することと、

前記特定のホモグラフィを前記グラウンドトゥルスホモグラフィと比較することと、  
前記比較に基づいて、前記ニューラルネットワークを修正することと

によって以前に訓練されている、項目2に記載の方法。

(項目6)

前記複数の3D点群は、1つ以上の幾何学形状をサンプリングすることによって生成される、項目5に記載の方法。

(項目7)

前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢は、少なくとも30%重複を有する、項目5に記載の方法。

(項目8)

拡張現実(AR)デバイスであって、前記ARデバイスは、  
カメラと、

前記カメラに通信可能に結合されたプロセッサと  
を備え、

前記プロセッサは、

前記カメラから、第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することと、

前記第1の画像に基づく第1の点群と、前記第2の画像に基づく第2の点群とを生成することと、

前記第1の点群および前記第2の点群をニューラルネットワークに提供することと、  
前記ニューラルネットワークによって、前記第1の点群および前記第2の点群に基づい

10

20

30

40

50

て、ホモグラフィを生成することと

を含む動作を実施するように構成されている、ARデバイス。

(項目9)

前記第1の点群および前記第2の点群は、2次元(2D)点群である、項目8に記載のARデバイス。

(項目10)

前記第1の点群および前記第2の点群は、第1のニューラルネットワークを使用して生成され、前記ニューラルネットワークは、第2のニューラルネットワークである、項目9に記載のARデバイス。

(項目11)

前記ニューラルネットワークは、

複数の点を含む1つ以上の3D点群のうちの各3次元(3D)点群のために、

前記複数の点の閾値距離内の3D軌道を決定することと、

前記3D軌道をサンプリングし、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢を取得することであって、前記複数の点は、少なくとも部分的に前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢から視認可能である、ことと、

前記特定の第1のカメラ姿勢に基づいて、前記複数の点を第1の2D平面上に投影し、第1の2D点群を生成することと、

前記特定の第2のカメラ姿勢に基づいて、前記複数の点を第2の2D平面上に投影し、第2の2D点群を生成することと、

前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢に基づいて、前記第1の2D点群と前記第2の2D点群との間のグラウンドトゥルスホモグラフィを決定することと、

前記ニューラルネットワークによって、前記第1の2D点群および前記第2の2D点群に基づいて、特定のホモグラフィを生成することと、

前記特定のホモグラフィを前記グラウンドトゥルスホモグラフィと比較することと、

前記比較に基づいて、前記ニューラルネットワークを修正することと

によって以前に訓練されている、項目9に記載のARデバイス。

(項目12)

前記複数の3D点群は、1つ以上の幾何学形状をサンプリングすることによって生成される、項目11に記載のARデバイス。

(項目13)

前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢は、少なくとも30%重複を有する、項目11に記載のARデバイス。

(項目14)

命令を備えている非一過性コンピュータ読み取り可能な媒体であって、前記命令は、プロセッサによって実行されると、

第1のカメラ姿勢に基づく第1の画像と、第2のカメラ姿勢に基づく第2の画像とを受信することと、

前記第1の画像に基づく第1の点群と、前記第2の画像に基づく第2の点群とを生成することと、

前記第1の点群および前記第2の点群をニューラルネットワークに提供することと、

前記ニューラルネットワークによって、前記第1の点群および前記第2の点群に基づいて、ホモグラフィを生成することと

を含む動作を前記プロセッサに実施させる、非一過性コンピュータ読み取り可能な媒体。

(項目15)

前記第1の点群および前記第2の点群は、2次元(2D)点群である、項目14に記載の非一過性コンピュータ読み取り可能な媒体。

(項目16)

前記第1の画像は、第1の瞬間において第1のカメラによって捕捉され、前記第2の画

10

20

30

40

50

像は、前記第1の瞬間後の第2の瞬間において前記第1のカメラによって捕捉されている、項目15に記載の非一過性コンピュータ読み取り可能な媒体。

(項目17)

前記第1の点群および前記第2の点群は、第1のニューラルネットワークを使用して生成され、前記ニューラルネットワークは、第2のニューラルネットワークである、項目15に記載の非一過性コンピュータ読み取り可能な媒体。

(項目18)

前記ニューラルネットワークは、複数の点を含む1つ以上の3D点群のうちの各3次元(3D)点群のために、

前記複数の点の閾値距離内の3D軌道を決定することと、

前記3D軌道をサンプリングし、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢を取得することであって、前記複数の点は、少なくとも部分的に前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢から視認可能である、ことと、

前記特定の第1のカメラ姿勢に基づいて、前記複数の点を第1の2D平面上に投影し、第1の2D点群を生成することと、

前記特定の第2のカメラ姿勢に基づいて、前記複数の点を第2の2D平面上に投影し、第2の2D点群を生成することと、

前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢に基づいて、前記第1の2D点群と前記第2の2D点群との間のグラウンドトゥールースホモグラフィを決定することと、

前記ニューラルネットワークによって、前記第1の2D点群および前記第2の2D点群に基づいて、特定のホモグラフィを生成することと、

前記特定のホモグラフィを前記グラウンドトゥールースホモグラフィと比較することと、

前記比較に基づいて、前記ニューラルネットワークを修正することと

によって以前に訓練されている、項目15に記載の非一過性コンピュータ読み取り可能な媒体。

(項目19)

前記複数の3D点群は、1つ以上の幾何学形状をサンプリングすることによって生成される、項目18に記載の非一過性コンピュータ読み取り可能な媒体。

(項目20)

前記特定の第1のカメラ姿勢および前記特定の第2のカメラ姿勢は、少なくとも30%重複を有する、項目18に記載の非一過性コンピュータ読み取り可能な媒体。

【図面の簡単な説明】

【0009】

種々の実施形態の性質および利点のさらなる理解は、以下の図を参照することによって実現され得る。添付の図では、類似コンポーネントまたは特徴は、同一参照標識を有し得る。さらに、同一タイプの種々のコンポーネントは、参照標識の後、ダッシュと、類似コンポーネント間で区別する第2の標識とが続くことによって区別され得る。第1の参照標識のみが、本明細書で使用される場合、説明は、第2の参照標識に関係なく、同一の第1の参照標識を有する類似コンポーネントのうちの任意の1つに適用可能である。

【0010】

【図1】図1は、本発明のいくつかの実施形態による同時位置特定およびマッピング(SLAM)アプローチの3つのサブタスクの視覚的概要を図示する。

【0011】

【図2】図2は、本発明のいくつかの実施形態によるDH-SLAM追跡システムを図示する。

【0012】

【図3】図3は、本発明のいくつかの実施形態によるMagicPointNetアーキテクチャを図示する。

【0013】

10

20

30

40

50

【図4】図4は、本発明のいくつかの実施形態による Magic Point Net 合成訓練データを図示する。

【0014】

【図5】図5は、本発明のいくつかの実施形態による点ベースの Homography Net アーキテクチャを図示する。

【0015】

【図6】図6は、本発明のいくつかの実施形態による点ベースの Homography Net データ生成を図示する。

【0016】

【図7】図7は、本発明のいくつかの実施形態によるスマートアンカ SLAM システムを図示する。

10

【0017】

【図8】図8は、本発明のいくつかの実施形態による2つの画像に基づいてホモグラフィを算出する方法を図示する。

【0018】

【図9】図9は、本発明のいくつかの実施形態によるニューラルネットワークを訓練する方法を図示する。

【0019】

【図10】図10は、本発明のいくつかの実施形態によるウェアラブル拡張現実 (AR) デバイスの概略図を図示する。

20

【発明を実施するための形態】

【0020】

画像カテゴリ化およびオブジェクト検出等のコンピュータビジョンタスクにおける深層学習成功の多くは、ImageNet および MS-COCO のような大規模な注釈が付けられたデータベースの利用可能性から生じる。しかしながら、同時位置特定およびマッピング (SLAM) のような姿勢追跡および再構成問題に対して、代わりに、Microsoft Kinect に基づく Freiburg-TUM RGB-D データセット、立体視カメラおよび IMU に基づく EuroC ドローン/MAV データセット、ならびに4つのカメラ、GPS データ、および Velodyne レーザスキャナを装備する車両の KITTI 運転データセット等のより小規模なデバイス特有のデータセットのフラグメント化されたエコシステムが存在する。

30

【0021】

SLAM のための ImageNet は、現在、存在しない。実際、多数の環境およびシナリオのための正確なグラウンドトゥルース (ground-truth) 姿勢測定を取得することは、困難である。グラウンドトゥルースセンサと視覚的 SLAM システムとの間の正確な整列を得ることは、かなりの努力を要求し、異なるカメラにおける変動にわたってスケールリングすることは容易ではない。フォトリアリスティックレンダリングは、SLAM タスクのための全ての関連幾何学的変数が 100% 正確度で記録され得るので、有用であり得る。フォトリアリスティックシーケンス上のベンチマーキング SLAM は、利点を有し得るが、そのようなレンダリングされた画像に基づく訓練は、多くの場合、ドメイン適合問題に悩まされ、多くの深層ネットは、過剰適合すると考えられる。

40

【0022】

実際は、SLAM モデルが、ロボットおよび複合現実におけるそれら等の内蔵プラットフォーム上で大規模に起動するように効率的である場合、好ましい。本発明の実施形態は、そのようなシステムが内蔵プラットフォーム上で展開され得るように、完全フレーム予測とは対照的に、幾何学的一貫性に焦点を当てる。深層学習システムによって行われる完全フレーム予測は、ある利点を有するが、多くの場合では、点を予測すること/整列させることは、メトリックレベル姿勢復元のために十分である。深層ネットの展開は、通常、ネットを可能な限り小規模にするためのカスタムオフライン最適化プロシージャを伴い、したがって、本発明のいくつかの場合では、完全視覚的フレーム予測問題は、完全に省略さ

50

れる。

#### 【0023】

本発明の実施形態は、未知の環境における6自由度単眼カメラ位置特定のために、DH-SLAMと称されるシステムを含み得、それは、主に、高センサ雑音、低照明、および主要なオクルージョンの存在下でホモグラフィをロバストに推定するように訓練される畳み込みニューラルネットワークを装備している。システムは、3つの深層ConvNetsを含み得る。第1のネットワークであるMagicPointNetは、単一画像に動作し、画像内の顕著な点を抽出し（顕著な点は、設計によって、画像内で分離され、十分に分散させられている）、非最大抑制等の任意の追加の後処理を利用しないこともある。第2のネットワークであるPointHomographyNetは、MagicPointNetからの対の点応答マップに動作し、2つの点画像を関連させるホモグラフィを推定し得る。そして、ホモグラフィ推定は、標準的MVG技法を使用して再推定されるか、または、場面が高度に非平面である場合、基本行列にアップグレードされ得る。第3のネットワークであるRelocalizationNetは、単一画像を信頼性があり、かつ高速な再位置特定のために使用される高速画像比較を可能にする低次元埋め込みベクトルに変換し得る。ネットワークは、単純合成データで訓練され、単純合成データは、それらが高価な外部カメラグラウンドトゥールズ機器または高度なグラフィックレンダリングパイプラインに依拠しないので、ネットワークを訓練しやすくする。システムは、高速かつ効率的であり、CPU上で30+FPSで起動することが可能である。

10

#### 【0024】

本発明の種々の実施形態では、3つのカスタムSLAM畳み込みニューラルネットワークが、3つの別個の訓練ドメインとともに提示される。しかしながら、この特定の実装は、要求されず、ニューラルネットワークのうちの1つ以上のものは、いくつかの実施形態では、組み合わせられ得る。ホモグラフィが各システムの重要なコンポーネントであり得るので、ホモグラフィ的に導かれる単眼SLAMシステムが、提示される。本発明の実施形態は、どんな種類のローカル特徴点記述子も使用しないこともあり、手動データ注釈または高価なオフライングラフィックレンダリングパイプラインを要求しないこともある。

20

#### 【0025】

図1は、本発明のいくつかの実施形態によるSLAMアプローチの3つのサブタスクの視覚的概要を図示する。サブタスク102では、幾何学的点プリミティブが、結像効果（センサ雑音、照明、テクスチャ、およびモーションぼけ等）を除去するために、画像から抽出される。サブタスク104では、対の点画像が、比較され、画像を関連させる、グローバル姿勢情報を抽出する。サブタスク106では、画像は、高速画像マッチングのために、低次元ベクトルに圧縮される。

30

#### 【0026】

サブタスク102は、画像を点のような幾何学的エンティティの組にマッピングすることに関する。サブタスク102の1つの目標は、照明、陰影、および全体的グローバル照明変動のような迷惑変数の損傷を元に戻すことである。角検出に類似した技法を使用して、入力画像内の信頼性がある場所を抽出することも望ましくあり得る。しかしながら、画像の角を分析的に定義する代わりに、カスタム2D形状レンダラが、角検出器を訓練するために使用される。結果として生じる2D画像場所が、場面点の準稠密被覆がロバストな変換推定のために重要であり得るので、システムから準稠密的に出力される。

40

#### 【0027】

サブタスク104は、一对の画像間の相対的姿勢を見出すことに関する。いくつかの場合では、ホモグラフィは、ワーピングされた自然画像の大規模データセットからサンプリングされた複数対のグレースケール画像に基づいて訓練されたCNNから算出され得る。本発明のいくつかの実施形態では、CNNは、点画像、すなわち、MagicPointNetによって出力された画像の種類的空間を使用して訓練される。準稠密点画像の空間は、完全RGB画像の空間よりかなり小さくあり得るので、かなりより高い性能が、完全RGB画像を取り扱うために必要であるものより少ない数のニューロンを使用して取得さ

50

れ得る。本明細書に提示される姿勢推定ネットワークは、点応答画像に動作し得、ホモグラフィモードおよび基本モードの両方において起動し得る。2つの異なる方法において3D世界を取り扱う能力は、複数視点幾何学形状の状況では、場面平面性の仮定が適用できることも、できないこともあるので、重要であり得る。平面性仮定が適用できる場合、ホモグラフィが、推定され、後に、個々の点深度を配慮せずに、 $(R, t)$ 推定値にアップグレードされることができる。代わりに、ソルバが、基本行列推定に基づく場合、場面平面性は、退化E推定値を提供し、全ての他の推定は、失敗し得る。故に、視認可能場面幾何学形状が高度に非平面であるとき、E、R、t、および点深度は、直接、対処され得る。

#### 【0028】

サブタスク106は、画像の視覚的埋め込みを作成することに関する。他のタスクのために、姿勢に敏感であることは、重要であり得るが、埋め込みのために、姿勢に敏感でないことが、望ましくあり得る。実世界画像を使用することは、2D形状または点応答画像のようなエンジニアリングされた表現の上に、埋め込みを学習することが補助となる可能性が低くあり得るので、このタスクのために重要であり得る。代わりに、画像の大規模な自由に利用可能なデータセット（すなわち、ImageNet）が、使用され、ホモグラフィが、画像をワーピングさせるために使用される。学習中、2つの画像は、それらがホモグラフィ的に関連する場合、類似埋め込みベクトルを有するように強制され得る。

#### 【0029】

図2は、本発明のいくつかの実施形態によるDH-SLAM追跡システム200の一般的アーキテクチャを図示する。いくつかの場合では、対の画像（例えば、画像202および画像204）が、MagicPointNetによって処理され、MagicPointNetは、画像内の顕著な点を検出し、一对の点画像（例えば、点画像206および点画像208）を生成するように訓練される。点画像は、次いで、PointHomographyNetによって、一緒に処理され、時間Tにおける点画像と時間T+1における点画像を関連させるホモグラフィHを算出する。ネットワークからのH推定値は、ほぼ平面の場面に対して、多くの場合、点を互いの数ピクセル以内に投影し得、したがって、単純最近傍対応が、より精密なH'を再推定するために十分であるか、または、Hは、複数視点幾何学形状技法を使用して非平面場面をより詳細に説明する基本行列Fにアップグレードされ得る。そして、H'およびFの両方は、カメラ行列Kを使用して、カメラの相対的姿勢に分解されることができる。

#### 【0030】

いくつかの場合では、SLAMシステムは、3つのサブシステムに分解される：画像を2D点場所の規準的サブ空間の中にもたらすMagicPoint検出器、グローバル変換推定点ベースのHomographyNet、RelocNetと呼ばれる視覚的画像埋め込みエンジン。いくつかの実施形態では、各サブシステムは、別個の畳み込みニューラルネットワークである。MagicPointNetおよびRelocNetの両方が、単一画像に動作する一方、PointHomographyNetは、一对の画像に動作する。いくつかの実施形態では、ネットワークは、以下の問題を解決することが可能である：雑音の多い画像においてロバストな2D場所を検出すること、2つの画像間の相対的姿勢を算出すること、および、再位置特定。エンジニアリングされた特徴記述子（ORBまたはSIFT）の両方に大きく依拠する従来のアプローチと異なり、本発明の実施形態は、記述子を画像内の個々の点に関連付けないこともある。これらの従来の特徴ベースのSLAMシステムと異なり、相対的姿勢推定は、無記述子方式で実施され得る。グローバル画像全体記述子に類似し得る埋め込みも、使用され得る。埋め込みは、擬似ホモグラフィ不変量であるようにエンジニアリングされ得る。設計によって、ホモグラフィによって関連付けられる2つの画像は、所与の多様体上で近接し得る。

#### 【0031】

いくつかの場合では、第1のステップは、画像内で顕著かつ位置特定可能な2D場所を検出することを含み得る。このステップは、HarrisまたはFAST等の角様応答マップを算出し、極大値を検出し、非最大抑制（non-maximal suppress

10

20

30

40

50

sion)を採用することによって実施され得る。追加のステップが、これらの極大値を画像全体を通して分散させるために実施され得る。このプロセスは、高度な専門領域の知識および手動エンジニアリングを伴い得、それは、一般化およびロバスト性を限定し得る。SLAM設定における有用性を増加させるために、点検出器によって検出された点は、フレームにわたる対応が容易であるように、画像全体を通して広く分散させられ、互いから分離され得る。システムが高センサ雑音シナリオおよび低光量において点を検出することも望ましくあり得る。いくつかの場合では、信頼度スコアが、検出された各点のために取得され得、それは、スプリアス点を排除することに役立つように使用され得る。これらの像点が画像内のローカル高勾配縁に対応する必要はないが、代わりに、プロブの中心等の他の低レベルキューに対応し得、それが、従来の方検出器より大きい受け入れ可能野を

10

#### 【0032】

図3は、本発明のいくつかの実施形態によるMagicPointNetアーキテクチャ300を図示する。いくつかの実施形態では、MagicPointNetは、グレースケール画像に動作し、入力各ピクセルのための「点性(point-ness)」確率を出力する。これは、明示的デコーダと組み合わせられたVGG式エンコーダを用いて、行われ得る。最終 $15 \times 20 \times 65$ テンソルにおける各空間場所は、ローカル $8 \times 8$ 領域に関する確率分布+単一ダストピンチャンネルを表し得、それは、点が検出されないことを表し得る( $8 \times 8 + 1 = 65$ )。ネットワークは、2D形状レンダラからの点監視を使用した標準クロスエントロピー損失を使用して訓練される。

20

#### 【0033】

いくつかの場合では、MagicPointNetは、カスタムConvNetアーキテクチャおよび訓練データパイプラインを使用して実装される。いくつかの場合では、画像Iを等価分解能を伴う点応答画像Pにマッピングすることが重要であり得、出力各ピクセルは、入力におけるそのピクセルのための「点性」の確率に対応する。稠密予測のためのネットワーク設計は、エンコーダ-デコーダペアを伴い得、空間分解能は、プーリングまたはストライド畳み込みを介して減らされ、そして、上方畳み込み演算を介して、完全分解能に戻るようアップサンプリングされる。アップサンプリング層は、より高い算出負担を追加し得るので、MagicPointNetは、明示的デコーダを用いて実装され、モデルの算出を低減させ得る。畳み込みニューラルネットワークは、VGG式エンコーダを使用して、画像の寸法を $120 \times 160$ から $15 \times 20$ セルグリッドに低減させ、各空間位置のための65チャンネルを伴い得る。いくつかの実施形態では、QQVGA分解能は、算出負担を減少させるために、 $120 \times 160$ であり得る。65チャンネルは、ピクセルのローカル非重複 $8 \times 8$ グリッド領域+余剰ダストピンチャンネルに対応し得、それは、その $8 \times 8$ 領域内で点が検出されないことに対応する。ネットワークは、 $3 \times 3$ 畳み込み後、BatchNorm正規化およびReLU非線形性を使用して、完全に畳み込まれ得る。最終畳み込み層(convlayer)は、 $1 \times 1$ 畳み込みであり得る。

30

#### 【0034】

図4は、本発明のいくつかの実施形態によるMagicPointNet合成訓練データを図示する。いくつかの実施形態では、画像の着目点は、種々の視点、照明、および画像雑音変動にわたって安定している画像内の一意に識別可能な場所であり得る。スパーSSLAMシステムのための事前処理ステップとして使用されるとき、所与のSSLAMシステムのために良好に機能する点が、検出され得る。点検出アルゴリズムのハイパーパラメータを設計および選定することは、専門家および専門領域に特有の知識を利用し得る。

40

#### 【0035】

今日、着目点が標識された画像の大規模データベースは、存在しない。高価なデータ収集労力を回避するために、レンダラが、コンピュータビジョンライブラリを使用して実装され得る。いくつかの場合では、それらの各々のための角のグラウンドトゥール場所が既知である三角形、四辺形、市松模様、3D立方体、および楕円形等の単純幾何学的形状が、レンダリングされる。各オブジェクトの2D面の重心も、既知であり、それらは、追

50

加の着目点としての役割を果たす。2D面の中心を見出す単純かつロバストな検出器を設計することが困難であろうことに留意されたい。形状がレンダリングされた後、ホモグラフィワーピングが、各画像に適用され、訓練例の数を増大させ得る。随意に、雑音、例えば、照明変化、テクスチャエンジニアリングされる雑音、ガウス雑音、ごま塩雑音、それらの組み合わせ等の形態における大量の雑音が、画像のうちの1つ以上のもの（例えば、各画像）に適用され得る。データは、オンザフライで生成され得、それは、どんな例もネットワークによって2回経験されないことを意味する。ネットワークは、15×20グリッドにおける各セルのためのロジットがソフトマックス関数を通して送られた後、標準クロスエントロピー損失を使用して訓練され得る。

【0036】

いくつかの場合では、Point Homography Netは、Magic Pointによって生産されるような対の点画像を所与として、ホモグラフィを生産する。全ての画像の空間および相対的姿勢の空間（全ての画像の空間×相対的姿勢の空間）の代わりに、点画像の空間および相対的姿勢の空間（点画像の空間×相対的姿勢の空間）に動作するようにネットワークを設計することによって、照明、陰影、およびテクスチャ等のいくつかの考慮点は、重要性が低減させられるか、または無視され得る。さらに、適用できるために、測光一貫性仮定が、当てにされる必要はない。

【0037】

図5は、本発明のいくつかの実施形態による点ベースのHomography Netアーキテクチャ500を図示する。いくつかの場合では、対のバイナリ点画像が、連結され、そして、VGG式エンコーダを通してフィードされる。3×3ホモグラフィHが、完全結合層によって出力される。そして、ホモグラフィHは、その右下要素が1であるように、正規化され得る。損失が、第2の画像内に対応を有することが既知の1つの画像内の点を第2の画像に変換し、その位置を比較することによって算出される。Point Homography Netを訓練するための損失関数は、式(1)に示される。

【数10】

$$L_H = \sum_{n=1}^N \|Hx_n - x'_n\|_2 \quad (1)$$

【0038】

いくつかの場合では、Point Homography Netは、直接、Magic Point Netによって出力された点検出に動作するように設計される（但し、任意の従来の点検出器に動作することができる）。モデルは、対の準稠密15×20×65画像上で良好に機能し得る。この小空間分解能では、ネットワークは、非常にわずかな算出を使用する。入力チャンネル毎連結後、3×3畳み込み、最大プーリング、Batch Norm、およびReLUアクティブ化後、2つの完全結合層から成るVGG式エンコーダは、実装され得、それは、3×3ホモグラフィHの9つの値を出力する。

【0039】

図6は、本発明のいくつかの実施形態による点ベースのHomography Netデータ生成を図示する。いくつかの実施形態では、Point Homography Netを訓練するために、2つの仮想カメラの中にレンダリングされた点群の数百万の例が、生成され得る。点群は、平面、球体、および立方体を含む単純3D幾何学形状から生成され得る。2つの仮想カメラの位置は、図6に示されるように、区分線形平行移動およびランダム軸の周囲の回転から成るランダム軌道からサンプリングされる。いくつかの実施形態では、少なくとも30%視覚的重複を有するカメラ対が、ランダムにサンプリングされる。点が、2つのカメラフレームの中に投影されると、点入力ドロップアウトが、適用され、スプリアスおよび欠測点検出に対するネットワークのロバスト性を改良する。いくつかの場合では、性能は、独立して、合致の50%をランダムにドロップし、点の25%をラ

10

20

30

40

50

ンダムにドロップすることによって改良される。

【0040】

いくつかの考慮点が、典型的に、直接、 $3 \times 3$  行列を出力するようにネットワークを訓練するために考慮される。いくつかの場合では、訓練は、最終FC層バイアスが、単位行列を出力するように初期化されるとき、ホモグラフィHの座標が、範囲 $[-1, 1]$ に正規化されるとき、および、ホモグラフィHが8自由度を有し、かつ9つの要素を有するので、右下要素が1であるようにH数が正規化されるとき、最良に機能する。

【0041】

いくつかの実施形態では、埋め込みネットワークの1つの目標は、グローバル128次元記述子を入力画像に関連付けることであり得る。いくつかの実施形態では、埋め込みが、ホモグラフィ不変であることが望ましい。例えば、ホモグラフィによって関連する2つの画像は、同一埋め込みベクトルを有するべきであり、同一場面コンテンツを描写しない（したがって、同一平面ではない）2つの画像は、異なる埋め込みベクトルを有するべきである。

10

【0042】

埋め込みネットワークは、128 L2 - 正規化記述子を生産し得る。これは、VGGのようなエンコーダネットワークに加え、完全結合層によって行われ得る。いくつかの場合では、埋め込みネットワークは、ImageNetデータセットからの対のホモグラフィ的に関連する画像を使用して訓練され得る。全て同一画像からの画像パッチのトリプレット $(A, A', B)$ が、サンプリングされ得、 $(A, A')$ は、少なくとも30%重複を有し、ホモグラフィによって関連付けられ、 $(A, B)$ は、重複を有していない。ネットワークを訓練するために使用される、 $(A, A')$ 正対および $(A, B)$ 負対が、生成され得る。いくつかの実施形態では、当業者に明白であり得るように、シャムネットワークが、使用されることができ一方、他の実施形態では、2タワーアプローチまたはトリプレットネットワークが、使用され得る。

20

【0043】

MagicPointNetが、FAST角検出器およびHarris角検出器のような従来の角検出ベースラインに対して評価された。PointHomographyNetは、実センサからの合成データおよび画像シーケンスの両方に関して、古典的ORB+RANSACベースのホモグラフィ推定エンジンに対して評価された。合成ドット世界における評価の1つの利点は、点の組間のグラウンドトゥルス対応が既知であることである。雑音の量を変動させることが、2つのアルゴリズムが低下した程度を決定するために追加され得る。埋め込みネットワークを評価するために、最近傍のグリッドが、埋め込みメトリックを使用して算出される。これは、ベースResNetアクティブ化を使用して、アクティブ化空間内の最近傍を求めることと比較された。

30

【0044】

ホモグラフィSLAMシステム全体を評価するために、評価が、合致によって生産された最終 $(R, t)$ 推定値に関して実施された。追跡（最後のフレームに対する姿勢のみの推定）および埋め込み拡張追跡（最も近い埋め込みを伴うK枚の画像に対する姿勢の推定）の両方に関する数が、観察された。本発明の実施形態は、PTAM、ORB-SLAM、およびLSD-SLAMのようないくつかのオープンソースSLAMシステムと定質的に比較された。上で説明される評価の結果は、本発明の種々の実施形態が、従来のアプローチと比較して、より優れた性能を示すことを示した。

40

【0045】

本発明の実施形態は、PointHomographyNetおよびRelocNetとともに、少数の画像とともに、それらの関連付けられた点画像を使用する平面の周囲に設計される、小型拡張現実システムを含み得る。いくつかの場合では、データセット収集は、カメラ中心に向かって真っ直ぐに向いている法線 $[0, 0, 1]$ を伴う基準平面の頭部搭載構成から開始する短シーケンスを作成することを伴う。基準平面のそのような頭部搭載初期ビューを提供することは、ホモグラフィ分解からの2つの解の曖昧性解消ならび

50

に各2D目印のための初期スケールを可能にする。

【0046】

図7は、本発明のいくつかの実施形態によるスマートアンカSLAMシステム700を図示する。いくつかの実施形態では、ユーザは、最初に、屋内環境内のほとんど平面の表面の写真を撮影することによって、「スマートアンカ」の組を構築する。これは、対の点画像のユーザのキーフレームインデックスおよび埋め込み( $E_0, E_1, E_2, \dots$ )を取り込む。インデックスが構築されると、システムは、追跡モードで起動されることができ、MagicPointNetおよびRelocNetは、時間Tにおける入力画像を処理し、点画像 $P_T$ および埋め込みベクトル $E_T$ を生産し得る。いくつかの実施形態では、 $E_T$ のドット積が、次いで、他の(例えば、1つおきの)埋め込み( $E_0, E_1, E_2, \dots$ )のうちの一つ以上のもので算出され、埋め込み多様体上の最近傍を見出す。図7に示される特定の実施形態では、 $E_2$ が、選択される。 $E_2$ に対応する点画像 $P_2$ が、次いで、 $P_T$ とともにPointHomographyNetの中にフィードされ、ホモグラフィを算出し得、それは、 $P_T$ における点を $P_2$ に変換する。ホモグラフィは、最後に、回転R、平行移動t、および主平面nに分解され得る。最後に、( $P_2, E_2$ )アンカに対応するコンテンツが、ワーピングされ、ARオーバーレイとして入力画像内に表示され得る。

10

【0047】

図8は、2つの画像に基づいてホモグラフィを算出する方法800を図示する。方法800のステップは、示されるものと異なる順序で実施され得、方法800の一つ以上のステップは、方法800の実施中、省略され得る。方法800の一つ以上のステップは、非一過性コンピュータ読み取り可能な媒体内に含まれる命令を実行するように構成されるプロセッサによって、実施および/または開始され得る。

20

【0048】

ステップ802では、第1の画像および第2の画像が、受信される。第1の画像は、第1のカメラ姿勢に基づき得、第2の画像は、第2のカメラ姿勢に基づき得、第2のカメラ姿勢は、第1のカメラ姿勢と異なる。いくつかの実施形態では、第1の画像および第2の画像は、同一カメラによって捕捉され得(それぞれ、第1の瞬間および第2の瞬間において、第2の瞬間は、第1の瞬間後に生じる)、および他の実施形態では、同時または2つの瞬間に、第1の画像が、第1のカメラによって捕捉され得、第2の画像が、第2のカメラによって捕捉され得る。

30

【0049】

ステップ804では、第1の点群が、第1の画像に基づいて生成され、第2の点群が、第2の画像に基づいて生成される。いくつかの実施形態では、第1のニューラルネットワークが、点群を生成するために使用され、すなわち、第1の画像は、第1のニューラルネットワークへの入力として提供され、第1の点群は、第1の画像に基づいて、第1のニューラルネットワークによって生成され、第2の画像は、第1のニューラルネットワークへの入力として提供され、第2の点群は、第2の画像に基づいて、第1のニューラルネットワークによって生成される。第1のニューラルネットワークは、本明細書に説明されるMagicPointNetであり得、それは、画像に基づいて、2D点群を生成し得る。

40

【0050】

ステップ806では、第1の点群および第2の点群は、第2のニューラルネットワークへの入力として提供される。第2のニューラルネットワークは、本明細書に説明されるPointHomographyNetであり得、それは、点群に基づいて、ホモグラフィを生成し得る。いくつかの実施形態では、第1のニューラルネットワークは、2つのネットワークの機能が単一システム内で組み合わせられ得るように、第2のニューラルネットワークと組み合わせられ得る。

【0051】

ステップ808では、ホモグラフィは、第2のニューラルネットワークを使用して、第1の点群および第2の点群に基づいて生成される。いくつかの実施形態では、生成された

50

ホモグラフィは、行列（例えば、 $3 \times 3$ ）を含み、それから、第1のカメラ姿勢と第2のカメラ姿勢との間の相対的回転および相対的平行移動（すなわち、相対的姿勢）が、抽出され得る。

【0052】

図9は、ニューラルネットワークを訓練する方法900を図示する。方法900のステップは、示されるものと異なる順序で実施され得、方法900の1つ以上のステップは、方法900の実施中、省略され得る。方法900を参照して説明されるニューラルネットワークは、方法800を参照して説明される第2のニューラルネットワークであり得、それは、本明細書に説明されるPoint Homography Netであり得る。方法900の1つ以上のステップは、非一過性コンピュータ読み取り可能な媒体内に含まれる命令を実行するように構成されるプロセッサによって、実施および/または開始され得る。

10

【0053】

ステップ902では、1つ以上の3D点群が、生成される。3D点群の各々は、複数の3D点を含み得る。いくつかの実施形態では、1つ以上の3D点群は、他の可能性の中でもとりわけ、平面、球体、立方体等の1つ以上の幾何学形状をランダムにサンプリングすることによって生成される。例えば、特定の幾何学形状（例えば、球体または立方体）の表面が、ランダムにサンプリングされ、複数の3D点を生産し得る。代替として、特定の幾何学形状の縁が、ランダムにサンプリングされ得るか、または、表面および縁の両方が、ランダムにサンプリングされ得る。いくつかの実施形態では、特定の幾何学形状の体積全体が、ランダムにサンプリングされ得る。

20

【0054】

いくつかの実施形態では、ステップ904 - 916の各々は、1つ以上の3D点群の各3D点群のために実施され得る。ステップ904では、3D軌道が、複数の点の近傍で決定され得る。例えば、3D軌道の全ては、複数の点の閾値距離内にあり得る。いくつかの実施形態では、3D軌道は、ランダム開始場所およびランダム終了場所を決定することによって形成される線形軌道である。他の実施形態では、または、同一実施形態では、3D軌道は、非線形軌道（例えば、湾曲）であるか、または、3D軌道は、ランダム開始場所、ランダム終了場所、および1つ以上の中間場所を決定することによって形成される一連の線形軌道である。

【0055】

ステップ906では、3D軌道が、サンプリングされ、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢を取得し得る。いくつかの実施形態では、複数の点は、少なくとも部分的に特定の第1のカメラ姿勢および特定の第2のカメラ姿勢から視認可能である。例えば、取得されるカメラ姿勢は、複数の点の少なくとも25%、50%、75%、または100%を視認するそれらのカメラ姿勢に制限され得る。カメラ姿勢が、所定の閾値（例えば、複数の点の少なくとも50%が視認可能である）を満たさない場合、カメラ姿勢は、破棄され、3D軌道は、再サンプリングされ、別のカメラ姿勢を取得する。いくつかの実施形態では、取得されるカメラ姿勢は、互いに視覚的重複の少なくともある閾値（例えば、30%）を有するように制限される。いくつかの実施形態では、視覚的重複は、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢の両方によって視認可能な、複数の点のうち点のパーセンテージに対応し得る。他の実施形態では、視覚的重複が、取得される姿勢間の共有視野に基づいて計算され得る。

30

40

【0056】

ステップ908では、複数の点は、特定の第1のカメラ姿勢に基づいて、第1の2D平面上に投影され、第1の2D点群を生成し、同様に、複数の点は、特定の第2のカメラ姿勢に基づいて、第2の2D平面上に投影され、第2の2D点群を生成する。いくつかの実施形態では、第1の2D平面は、特定の第1のカメラ姿勢の向きベクトルに直交し得、第2の2D平面は、特定の第2のカメラ姿勢の向きベクトルに直交し得る。

【0057】

ステップ910では、第1の2D点群と第2の2D点群との間のグラウンドトゥールース

50

ホモグラフィが、特定の第1のカメラ姿勢および特定の第2のカメラ姿勢に基づいて決定される。いくつかの実施形態では、グラウンドトゥルスホモグラフィは、最初に、特定の第1のカメラ姿勢と特定の第2のカメラ姿勢との間の相対的回転および相対的平行移動を決定し、そして、ニューラルネットワークによって生成されたホモグラフィと構造が一貫するホモグラフィ（行列）を形成することによって決定される。

【0058】

ステップ912では、第1の2D点群および第2の2D点群は、入力としてニューラルネットワークに提供され得、特定のホモグラフィが、第1の2D点群および第2の2D点群に基づいて、ニューラルネットワークによって生成され得る。

【0059】

ステップ914では、特定のホモグラフィは、グラウンドトゥルスホモグラフィと比較され、例えば、誤差信号を生成し得る。いくつかの実施形態では、誤差信号の大きさは、特定のホモグラフィとグラウンドトゥルスホモグラフィとの間の差異の大きさに比例し得る。1つの特定の实施形態では、誤差信号は、行列の対応する要素が互いから減算される標準行列減算方法を使用して計算される。他の実施形態では、または、同一実施形態では、誤差信号は、相対的回転における差異に対応する第1の成分と、相対的平行移動における差異に対応する第2の成分とを含む。いくつかの実施形態では、誤差信号は、相対的姿勢における差異に対応する単一成分を含む。

【0060】

ステップ916では、ニューラルネットワークは、ステップ914において実施される特定のホモグラフィとグラウンドトゥルスホモグラフィとの間の比較に基づいて、例えば、ニューラルネットワークの1つ以上の重みまたは係数を調節することによって、修正される。いくつかの実施形態では、ニューラルネットワークは、より大きい誤差信号がより大きい修正をニューラルネットワークに生じさせるように、ホモグラフィ間の計算された差異（すなわち、誤差信号）に基づいて、修正され得る。一般に、ニューラルネットワークを修正することは、ニューラルネットワークがより正確になるようにし、それによって、特定のホモグラフィとグラウンドトゥルスホモグラフィとの間の差異を減少させる。

【0061】

図10は、本明細書に説明される実施形態のうちの1つ以上のものを採用し得るウェアラブルARデバイス1000の概略図を図示する。ARデバイス1000は、左接眼レンズ1002Aと、右接眼レンズ1002Bと、直接、左接眼レンズ1002A上またはその近傍に取り付けられる左正面に面した世界カメラ1006Aと、直接、右接眼レンズ1002B上またはその近傍に取り付けられる右正面に面した世界カメラ1006Bと、左側に面した世界カメラ1006Cと、右側に面した世界カメラ1006Dと、処理モジュール1050とを含み得る。ARデバイス1000のコンポーネントの一部または全部は、投影された画像がユーザによって視認され得るように、頭部搭載型であり得る。1つの特定の实装では、図10に示されるARデバイス1000のコンポーネントの全ては、ユーザによって装着可能な単一デバイス（例えば、単一ヘッドセット）ウェアラブル上に搭載される。別の実装では、処理モジュール1050は、ARデバイス1000の他のコンポーネントと物理的に別個であり、有線または無線接続性によって、それに通信可能に結合される。例えば、処理モジュール1050は、フレームに固定して取り付けられる構成、ユーザによって装着されるヘルメットまたは帽子に固定して取り付けられる構成、ヘッドホンに内蔵される構成、または別様に、ユーザに除去可能に取り付けられる構成（例えばリュック式構成、ベルト結合式構成等において）等、種々の構成において搭載され得る。

【0062】

処理モジュール1050は、プロセッサ1052と、不揮発性メモリ（例えば、フラッシュメモリ）等のデジタルメモリとを備え得、両方は、データの処理、キャッシュ、および記憶を補助するために利用され得る。データは、画像捕捉デバイス（例えば、カメラ1006）、マイクロホン、慣性測定ユニット、加速度計、コンパス、GPSユニット、無線デバイス、および/またはジャイロスコープから捕捉されたデータを含み得る。例えば

10

20

30

40

50

、処理モジュール1050は、カメラ1006からの画像1020、より具体的に、左正面に面した世界カメラ1006Aからの左正面画像1020A、右正面に面した世界カメラ1006Bからの右正面画像1020B、左側に面した世界カメラ1006Cからの左側画像1020C、および右側に面した世界カメラ1006Dからの右側画像1020Dを受信し得る。いくつかの実施形態では、画像1020は、単一画像、一对の画像、画像のストリームを備えているビデオ、ペアリングされた画像のストリームを備えているビデオ等を含み得る。画像1020は、ARデバイス1000が電源オンである間、周期的に、生成され、処理モジュール1050に送信され得るか、または、処理モジュール1050によってカメラのうちの1つ以上のものに送信される命令に応答して生成され得る。

#### 【0063】

接眼レンズ1002Aおよび1002Bは、プロジェクタ1014Aおよび1014Bからの光を向けるように構成される透明または半透明導波管を備え得る。具体的に、処理モジュール1050は、左プロジェクタ1014Aに、左投影画像1022Aを左接眼レンズ1002Aの中に出力させ得、右プロジェクタ1014Bに、右投影画像1022Bを右接眼レンズ1002Bの中に出力させ得る。いくつかの実施形態では、接眼レンズ1002の各々は、各々が異なる色および/または異なる深度平面に対応する複数の導波管を備え得る。

#### 【0064】

カメラ1006Aおよび1006Bは、それぞれ、ユーザの左および右眼の視野と実質的に重複する画像を捕捉するように位置付けられ得る。故に、カメラ1006Aおよび1006Bの場所は、ユーザの眼の近傍であり得るが、ユーザの視野を曖昧にするほど近傍ではない。代替として、または加えて、カメラ1006Aおよび1006Bは、それぞれ、投影された画像1022Aおよび1022Bの内部結合場所と整合するように位置付けられ得る。カメラ1006Cおよび1006Dは、ユーザの側面、例えば、ユーザの周辺視覚内またはユーザの周辺視覚外の画像を捕捉するように位置付けられ得る。カメラ1006Cおよび1006Dを使用して捕捉された画像1020Cおよび1020Dは、必ずしも、カメラ1006Aおよび1006Bを使用して捕捉された画像1020Aおよび1020Bと重複する必要はない。

#### 【0065】

ARデバイス1000の動作中、処理モジュール1050は、訓練されたネットワーク1056を使用して、カメラ1006の任意のものによる2つの捕捉された画像に基づいて、ホモグラフィを算出し得る。推定されたホモグラフィは、プロセッサ1052によって使用され、ユーザの移動に起因するユーザの視野の変化をより正確に反映する仮想コンテンツをレンダリングし得る。ネットワーク1056は、人工ニューラルネットワーク、畳み込みニューラルネットワーク、深層ネットワーク、または例を処理することによって徐々に「学習」し得る任意のタイプのネットワークもしくはシステムであり得る。いくつかの実施形態では、ネットワーク1056は、信号を1つのものから別のものに伝送することが可能である接続されたノードの集合を備えている。プロセッサ1052は、単一ネットワーク1056と通信し得るか、またはいくつかの実施形態では、プロセッサ1052は、第1のネットワーク(例えば、MagicPointNetに対応する)、第2のネットワーク(例えば、PointHomographyNetに対応する)、および第3のネットワーク(例えば、RelocNetに対応する)等の複数のニューラルネットワークと通信し得る。

#### 【0066】

いくつかの例示的構成が説明されたが、種々の修正、代替構造、および均等物が、本開示の精神から逸脱することなく、使用され得る。例えば、前述の要素は、より大きいシステムのコンポーネントであり得、他のルールが、本技術の用途に優先するか、または別様にそれを修正し得る。いくつかのステップは、前述の要素が検討される前、間、または後にも行われ得る。故に、前述の説明は、請求項の範囲を束縛するものではない。

#### 【0067】

10

20

30

40

50

本明細書および添付の請求項で使用されるように、単数形「a」、「an」、および「the」は、文脈によって明確に別様に示されない限り、複数参照を含む。したがって、例えば、「ユーザ」の言及は、複数のそのようなユーザを含み、「プロセッサ」の言及は、1つ以上のプロセッサおよび当業者に公知のその均等物等の言及を含む。

【0068】

単語「comprise（～を備えている）」、「comprising（～を備えている）」、「contains（～を含む）」、「containing（～を含む）」、「include（～を含む）」、「including（～を含む）」、および「includes（～を含む）」も、本明細書および以下の請求項で使用されるとき、述べられた特徴、整数、コンポーネント、またはステップの存在を規定するために意図されるが、それらは、1つ以上の他の特徴、整数、コンポーネント、ステップ、行為、またはグループの存在または追加を除外するものではない。

10

【0069】

本明細書に説明される例および実施形態が、例証目的のためだけのものであり、それに照らして、種々の修正または変更が、当業者に示唆され、本願の精神および権限ならびに添付される請求項の範囲内に含まれるものであることも理解されたい。

20

30

40

50

【図面】  
【図 1】

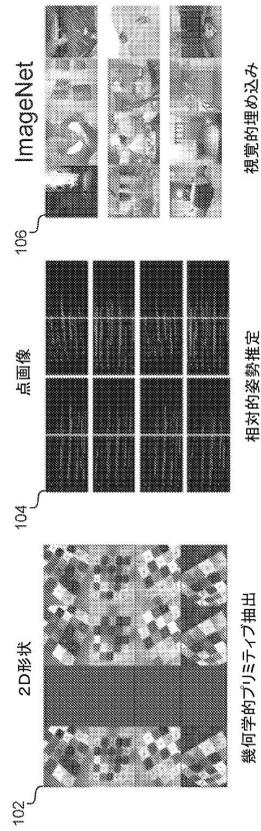


FIG. 1

【図 2】

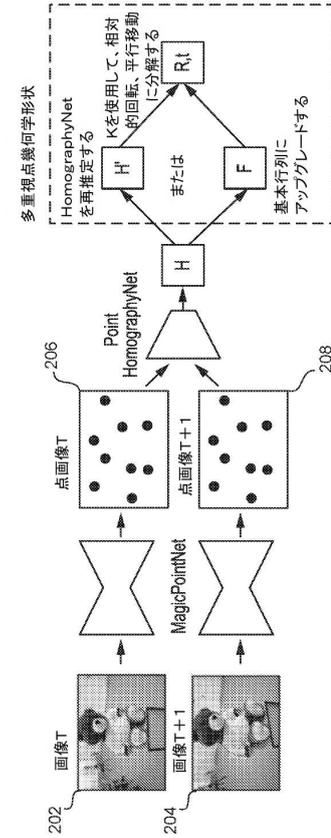


FIG. 2

【図 3】

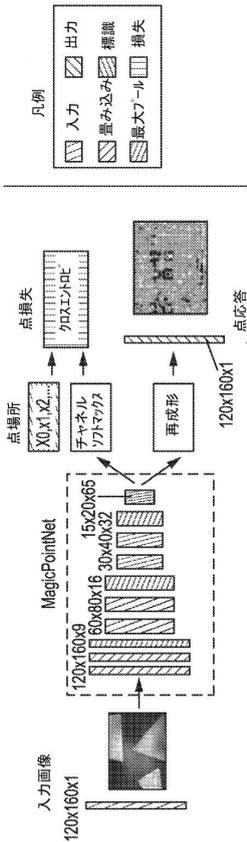


FIG. 3

【図 4】

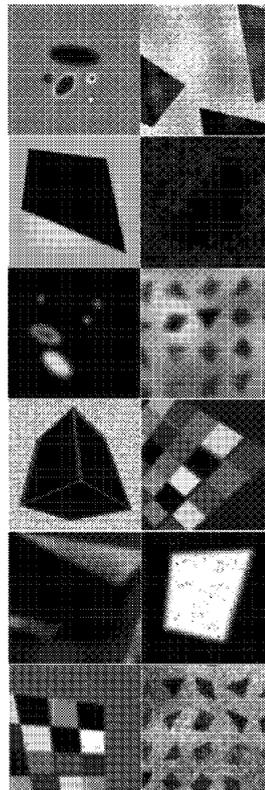


FIG. 4

【図5】

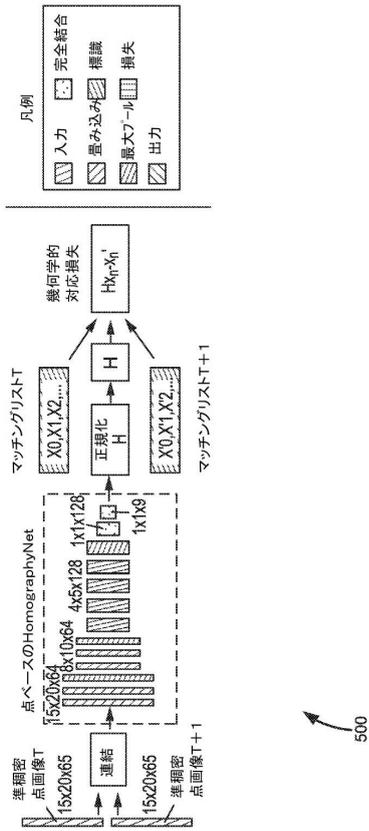


FIG. 5

【図6】

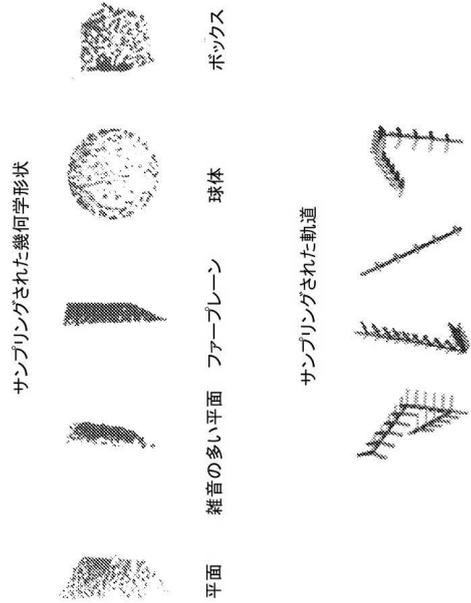


FIG. 6

【図7】

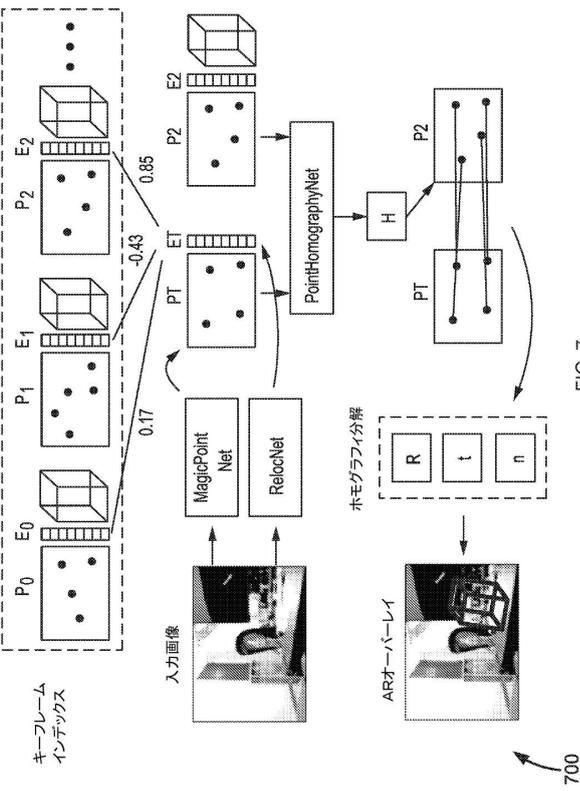


FIG. 7

【図8】

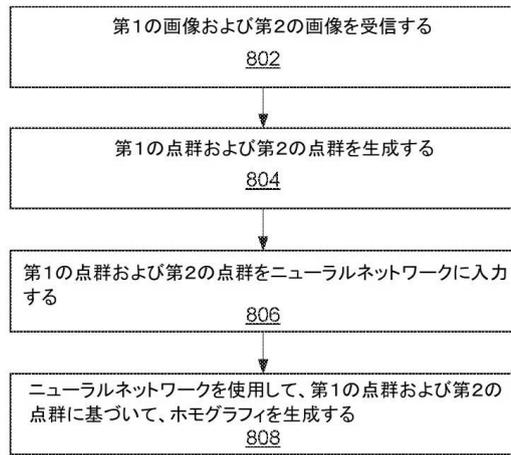


FIG. 8

10

20

30

40

50

【 図 9 】

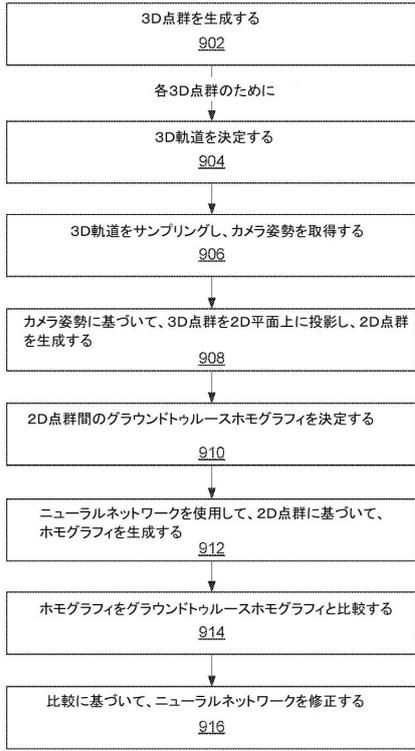


FIG. 9

900

【 図 10 】

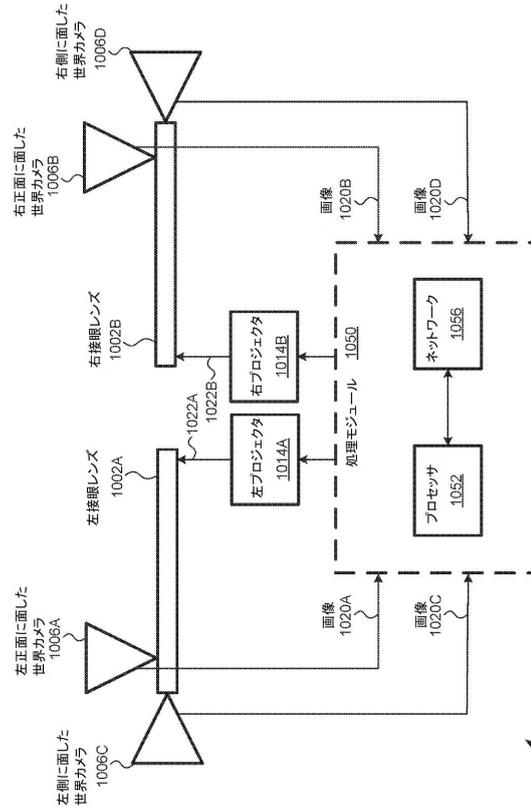


FIG. 10

1000

10

20

30

40

50

---

 フロントページの続き

- (74)代理人 100181641  
弁理士 石川 大輔
- (74)代理人 230113332  
弁護士 山本 健策
- (72)発明者 デトン, ダニエル エル.  
アメリカ合衆国 フロリダ 33322, プランテーション, ダブリュー. サンライズ プール  
バード 7500
- (72)発明者 マリシーウィッツ, トマシュ ジャン  
アメリカ合衆国 フロリダ 33322, プランテーション, ダブリュー. サンライズ プール  
バード 7500
- (72)発明者 ラビノビッチ, アンドリュウ  
アメリカ合衆国 フロリダ 33322, プランテーション, ダブリュー. サンライズ プール  
バード 7500
- 審査官 小池 正彦
- (56)参考文献 Hani Altwaijry, et al., Learning to Detect and Match Keypoints with Deep Architectures, vi  
sion.cornell.edu, 米国, 2016年08月01日, [https://vision.cornell.edu/se3/wp-content/upl  
oads/2016/08/learning-detect-match.pdf](https://vision.cornell.edu/se3/wp-content/uploads/2016/08/learning-detect-match.pdf)  
Daniel DeTone, et al., Deep Image Homography Estimation, arxiv.org, 米国, CORNELL U  
NIVERSIT, 2016年06月13日, <https://arxiv.org/pdf/1606.03798.pdf>
- (58)調査した分野 (Int.Cl., D B 名)
- |         |           |
|---------|-----------|
| G 0 6 T | 7 / 0 0   |
| G 0 6 T | 7 / 5 9 3 |
| G 0 6 T | 7 / 7 0   |
| G 0 6 N | 3 / 0 4   |
| G 0 6 N | 3 / 0 8   |