

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2024-510230

(P2024-510230A)

(43)公表日 令和6年3月6日(2024.3.6)

(51)国際特許分類	F I	テーマコード(参考)
G 0 6 T 7/00 (2017.01)	G 0 6 T 7/00 3 5 0 C	5 L 0 9 6
G 0 6 V 10/82 (2022.01)	G 0 6 T 7/00 6 6 0 B	
G 0 6 N 3/02 (2006.01)	G 0 6 V 10/82	
G 0 6 N 3/0464(2023.01)	G 0 6 N 3/02	
	G 0 6 N 3/0464	

審査請求 有 予備審査請求 未請求 (全17頁)

(21)出願番号 特願2023-556536(P2023-556536)
 (86)(22)出願日 令和4年3月31日(2022.3.31)
 (85)翻訳文提出日 令和5年9月13日(2023.9.13)
 (86)国際出願番号 PCT/IB2022/053034
 (87)国際公開番号 WO2022/208440
 (87)国際公開日 令和4年10月6日(2022.10.6)
 (31)優先権主張番号 63/168,467
 (32)優先日 令和3年3月31日(2021.3.31)
 (33)優先権主張国・地域又は機関 米国(US)
 (31)優先権主張番号 63/279,916
 (32)優先日 令和3年11月16日(2021.11.16)
 (33)優先権主張国・地域又は機関 米国(US)
 (31)優先権主張番号 17/701,991

最終頁に続く

(71)出願人 000002185
 ソニーグループ株式会社
 東京都港区港南1丁目7番1号
 (71)出願人 504257564
 ソニー コーポレーション オブ アメリカ
 アメリカ合衆国 ニューヨーク 1 0 0 1
 0 , ニューヨーク , マディソン アベ
 ニュー 2 5
 (74)代理人 100092093
 弁理士 辻居 幸一
 (74)代理人 100109070
 弁理士 須田 洋之
 (74)代理人 100067013
 弁理士 大塚 文昭
 (74)代理人

最終頁に続く

(54)【発明の名称】 顔表情、身体ポーズ形状及び衣服パフォーマンスキャプチャのための暗黙的微分可能レンダラーを用いたマルチビューニューラル人間予測

(57)【要約】

ニューラルヒューマンパフォーマンスキャプチャフレームワーク(MVS-Perf)が、校正されたマルチビュー画像セットから、人物の骨格、体形及び衣服の変位、並びに外観を取り込む。MVS-Perfは、単眼人体メッシュ復元(monocular human mesh recovery)において絶対位置を予測する曖昧さに対処し、NeRFからのボリューム表現をアニメーションに適したパフォーマンスキャプチャに仲介する。MVS-Perfは、マルチビュー画像から特徴マップを抽出して特徴量に融合するモジュール、特徴量を裸の人間パラメータベクトルに回帰して、骨格ポーズ、体形及び表情を含むSMPL-Xスキントイト人体メッシュ(SMPL-X skin-tight body mesh)を生成するモジュール、ニューラル放射場及び変形場を活用し、微分可能レンダリングを使用して衣服を裸体上の変位として推測するモジュールという3つのモジュールを含む。SMPL-Xスキントイト人体メッシュ頂点に補間された変位ベクトルを加算することによって、着衣姿の人体メッシュを取得する。取得

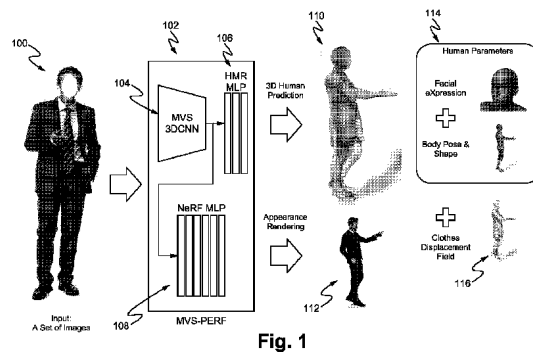


Fig. 1

【特許請求の範囲】**【請求項 1】**

装置の非一時的なものにプログラムされた方法であって、
画像セットを入力として取得することと、
ニューラルネットワークを使用して前記画像セットを処理することと、
を含み、前記処理は、
前記画像セットを 1 又は 2 以上の特徴に符号化することと、
前記特徴を人間パラメータに回帰させることと、
前記ニューラルネットワークを微調整することと、
クエリ 3 D 光線を、前記画像セットに基づく R G B カラー及び衣服 - 身体変位に復号
することと、
を含む、
ことを特徴とする方法。 10

【請求項 2】

前記画像セットは、サイズ $N \times w \times h \times c$ の 4 D テンソルを含み、ここで、 N はビュー
の数、 w は画像の幅、 h は画像の高さ、 c は画像のチャンネルである、
請求項 1 に記載の方法。

【請求項 3】

前記ニューラルネットワークは、前記画像セットから正面ビューを基準ビューとして選
択し、特徴量を抽出する、 20
請求項 1 に記載の方法。

【請求項 4】

前記ニューラルネットワークは、全ての特徴量を人間のポーズ、形状、表情パラメータ
に回帰させる、
請求項 3 に記載の方法。

【請求項 5】

前記ニューラルネットワークは、前記パラメータに従って人間の裸体メッシュを生成す
る、
請求項 4 に記載の方法。

【請求項 6】

前記裸体メッシュは、バウンディングボックス内の占有フィールドに変換される、 30
請求項 5 に記載の方法。

【請求項 7】

前記ニューラルネットワークは、ビューの各中心からの光線方向に関連する身体メッシ
ュの近くのいずれかの 3 D 点について、前記 R G B カラーと、裸体の表面を示す 3 D 変位
ベクトルとを生成する、
請求項 6 に記載の方法。

【請求項 8】

カメラビューの全ての画素から放たれる全ての光線を問い合わせることにより、着衣姿
の人体の外観が R G B 画像としてレンダリングされ、サンプリングされた点から前記 3 D
変位ベクトルを使用して裸体を変形させることにより、着衣姿の身体メッシュが取得され
る、 40
請求項 7 に記載の方法。

【請求項 9】

前記ニューラルネットワークは、教師ありモード又は自己教師ありモードで実装される
、
請求項 1 に記載の方法。

【請求項 10】

アプリケーションを記憶するように構成された非一時的メモリと、
前記アプリケーションを処理するように構成されたプロセッサと、 50

を備えた装置であって、前記アプリケーションは、
画像セットを入力として取得し、
ニューラルネットワークを使用して前記画像セットを処理する、
ように構成され、前記処理は、
前記画像セットを1又は2以上の特徴に符号化することと、
前記特徴を人間パラメータに回帰させることと、
前記ニューラルネットワークを微調整することと、
クエリ3D光線を、前記画像セットに基づくRGBカラー及び衣服 - 身体変位に復号
することと、
を含む、
ことを特徴とする装置。

10

【請求項11】

前記画像セットは、サイズ $N \times w \times h \times c$ の4Dテンソルを含み、ここで、 N はビュー
の数、 w は画像の幅、 h は画像の高さ、 c は画像のチャンネルである、
請求項10に記載の装置。

【請求項12】

前記ニューラルネットワークは、前記画像セットから正面ビューを基準ビューとして選
択し、特徴量を抽出する、
請求項10に記載の装置。

【請求項13】

前記ニューラルネットワークは、全ての特徴量を人間のポーズ、形状、表情パラメータ
に回帰させる、
請求項12に記載の装置。

20

【請求項14】

前記ニューラルネットワークは、前記パラメータに従って人間の裸体メッシュを生成す
る、
請求項13に記載の装置。

【請求項15】

前記裸体メッシュは、バウンディングボックス内の占有フィールドに変換される、
請求項14に記載の装置。

30

【請求項16】

前記ニューラルネットワークは、ビューの各中心からの光線方向に関連する身体メッシ
ュの近くのいずれかの3D点について、前記RGBカラーと、裸体の表面を示す3D変位
ベクトルとを生成する、
請求項15に記載の装置。

【請求項17】

カメラビューの全ての画素から放たれる全ての光線を問い合わせることにより、着衣姿
の人体の外観がRGB画像としてレンダリングされ、サンプリングされた点から前記3D
変位ベクトルを使用して裸体を変形させることにより、着衣姿の身体メッシュが取得され
る、
請求項16に記載の装置。

40

【請求項18】

前記ニューラルネットワークは、教師ありモード又は自己教師ありモードで実装される
、
請求項10に記載の装置。

【請求項19】

アプリケーションを記憶するように構成された非一時的メモリと、
前記アプリケーションを処理するように構成されたプロセッサと、
を備えた装置であって、前記アプリケーションは、

入力画像を特徴に符号化するように構成されたマルチビューステレオ3D畳み込みニ

50

ューラルネットワーク (M V S - 3 D C N N) と、

前記特徴を人間パラメータに回帰させるように構成された人間メッシュ復元多層パーセプトロン (H M R M L P) と、

前記 M V S - 3 D C N N を微調整するように構成され、クエリ 3 D 光線 (3 D 位置及び方向) を R G B カラー及び衣服 - 身体変位に復号するニューラル輝度場多層パーセプトロン (N e R F M L P) と、

を含む、

ことを特徴とする装置。

【請求項 20】

前記画像セットは、サイズ $N \times w \times h \times c$ の 4 D テンソルを含み、ここで、 N はビューの数、 w は画像の幅、 h は画像の高さ、 c は画像のチャンネルである、請求項 19 に記載の装置。

【請求項 21】

前記 M V S - 3 D C N N は、前記画像セットから正面ビューを基準ビューとして選択し、特徴量を抽出する、請求項 20 に記載の装置。

【請求項 22】

前記 H M R M L P は、全ての特徴量を人間のポーズ、形状、表情パラメータに回帰させる、請求項 21 に記載の装置。

【請求項 23】

前記パラメータに従って人間の裸体メッシュを生成するように構成されたモデルをさらに備える、請求項 22 に記載の装置。

【請求項 24】

前記裸体メッシュは、バウンディングボックス内の占有フィールドに変換される、請求項 23 に記載の装置。

【請求項 25】

前記 N e R F M L P は、ビューの各中心からの光線方向に関連する身体メッシュの近くのいずれかの 3 D 点について、前記 R G B カラーと、裸体の表面を示す 3 D 変位ベクトルとを生成する、請求項 24 に記載の装置。

【請求項 26】

カメラビューの全ての画素から放たれる全ての光線を問い合わせることにより、着衣姿の人体の外観が R G B 画像としてレンダリングされ、サンプリングされた点から前記 3 D 変位ベクトルを使用して裸体を変形させることにより、着衣姿の身体メッシュが取得される、請求項 25 に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

〔関連出願との相互参照〕

本出願は、2021年11月16日に出願された「顔表情、身体ポーズ形状及び衣服パフォーマンスキャプチャのための暗黙的微分可能レンダラーを用いたマルチビューニューラル人間予測 (MULTIVIEW NEURAL HUMAN PREDICTION USING IMPLICIT DIFFERENTIABLE RENDER FOR FACIAL EXPRESSION, BODY POSE SHAPE AND CLOTHES PERFORMANCE CAPTURE)」という名称の米国仮特許出願シリアル番号第 63 / 279, 916 号、及び 2021 年 3 月 31 日に出願された「顔表情、身体ポーズ形状及び衣服変位のための暗黙的微分可能レンダラーを用いたマルチビ

10

20

30

40

50

ューニューラル人間予測 (MULTIVIEW NEURAL HUMAN PREDICTION USING IMPLICIT DIFFERENTIABLE RENDER FOR FACIAL EXPRESSION, BODY POSE SHAPE AND CLOTHES DISPLACEMENT)」という名称の米国仮特許出願シリアル番号第 63 / 168, 467 号の米国特許法第 119 条に基づく優先権の利益を主張するものであり、これらの両文献はその全体が全ての目的で引用により本明細書に組み入れる。

【0002】

本発明は、娯楽産業のための 3 次元コンピュータビジョン及びグラフィックスに関する。具体的には、本発明は、映画、TV、音楽及びゲームコンテンツ制作のための 3 次元コンピュータビジョン及びグラフィックスを取得して処理することに関する。

10

【背景技術】

【0003】

例えば Facebook FrankMocap などの従来のシステムは、単一画像から裸体の形状及びポーズのみを予測する。このようなシステムは、衣服表面を予測することができない。このようなシステムは 2D 画像変換法であり、マルチビュー入力に対処することができない。

【0004】

暗黙的パーツネットワーク (Implicit Part Network) は、スキャン又は再構成された点群から身体及び衣服の両方を予測するが、3D スキャンを必要とし、入力としての RGB 画像にも、顔表情及び外観にも対処することができない。また、暗黙的パーツネットワークは、ボクセルを身体又は衣服として識別するラベルのみを予測した後人間事前モデル (human prior model) を明示的にフィットさせ、低速である。Neural Body 及び Animatable NeRF は、ニューラル輝度場 (Neural Radiance Field: NeRF) を使用して、顔表情を含まない衣服人体 (clothes human body) を予測する。しかしながら、これらは低解像度に制限される高密度の潜在コードボリューム (dense latent code volume) の作成を必要とし、従って人体形状が粗くなってしまう。また、これらは、メッシュ頂点の対応関係を含まないポリメトリックな人体モデルしか復元することができない。

20

30

【発明の概要】

【課題を解決するための手段】

【0005】

マルチビューニューラル人間予測 (Multiview neural human prediction) が、カメラ校正を与えられたマルチビュー画像セットから、骨格、体形、並びに衣服の変位及び外観を含む 3D 人間モデルを予測することを含む。

【0006】

1 つの態様では、ニューラルネットワークが、異なるビューからの単一画像又は複数画像であることができる入力画像セットを受け取って、層状 3D 人間モデル (layered 3D human model) を予測する。画像セットは、 $N \times w \times h \times c$ のサイズの 4D テンソルを含み、ここで、 N はビューの数であり、 w は画像の幅であり、 h は画像の高さであり、 c は画像のチャンネルである。画像セットのためのカメラ情報は既知である。出力モデルは、内側から外側に向かって、予測されたポーズの骨格、顔表情を含む予測された形状の裸の 3D 身体 (例えば、ブレンドシェイプ (blendshapes) 及び関節回転によってパラメータ化された SMPL-X モデル)、及び入力画像から推測される衣服変位及び外観 RGB 色の 3D 場という 3 つの層を含む。裸の 3D 人体メッシュを衣服変位場 (clothes displacement field) に従って変形させることによって着衣姿の人体メッシュ (clothed body mesh) が取得される。

40

【0007】

50

別の態様では、ニューラルネットワークが、入力画像セットを特徴に符号化するマルチビューステレオ3D畳み込みニューラルネットワーク(MVS-3DCNN)、特徴を人間パラメータに回帰させる人間メッシュ復元多層パーセプトロン(human mesh recovery multilayer perceptron: HMR MLP)、及びMVS-3DCNNを微調整してクエリ3D光線(3D位置及び方向)をRGBカラー及び衣服-身体変位に復号するニューラル輝度場多層パーセプトロン(neural radiance field multilayer perceptron: NeRF MLP)という3つのサブネットワークで構成される。

【0008】

別の態様では、テスト/推論モードにおいて、層状3D人間モデルの予測が、訓練データ内のカメラのビュー範囲内で、明示的な数値最適化を伴わずに、小さな入力セットについて、装置に依存せず、完全に自動であり、リアルタイムである。訓練済みニューラルネットワークを用いて予測する際には、MVS-3DCNNが、マルチビュー画像セットを入力として受け取り、正面ビューを基準ビューとして選択し、特徴量を抽出する。HMR MLPは、全ての特徴量を人間のポーズ、形状、顔表情パラメータに回帰させる。SMP L-Xモデルは、パラメータに従って人間の裸体メッシュを生成する。その後、裸体メッシュは、バウンディングボックス内の占有フィールドに変換される。訓練済みNeRF MLPは、ビューの各中心からの光線方向に関連する身体メッシュの近くのいずれかの3D点について、RGBカラーと、裸体の表面を示す3D変位ベクトルとを生成する。カメラビュー(入力ビューと同じビュー、又はいずれかの新規ビュー)の全ての画素から放たれる全ての光線を問い合わせることにより、着衣姿の人体の外観をRGB画像としてレンダリングすることができる。サンプリングされた点から3D変位ベクトルを使用して裸体を変形させることにより、SMP L-Xモデルと同じ頂点对応のSMP L-X+Dなどの着衣姿の人体メッシュを取得することができる。

【0009】

別の態様では、ニューラルネットワークの訓練が、教師あり及び自己教師ありという2つの事例を含む。教師ありの事例では、例えばH36Mデータセットなどの、既知の人間パラメータを有するラベル付きデータセットが与えられる。グランドトゥルース(GT)のパラメータ及び形状を、CNN回帰されたパラメータ及び形状と比較する。その差分を形状損失として計算する。一方で、入力画像セット内のサンプリングされた画素から光線を投げ、NeRF MLPが光線をレンダリングして、パラメータを裸体の密度及び3D衣服変位の関数である色及び密度に回帰させる。色損失は、サンプリングされた画素色とレンダリングされた色との差分の合計によって計算される。一方で、モーションキャプチャデータセットなどの、GT人間パラメータが未知である既存のデータセットでは、自己教師あり/自己改善訓練(self-improving training)が利用される。各訓練反復では、MVS-3DCNNからパラメータを回帰させた後に、これらをSMP Lify Xなどの最適化ベースの人間予測アルゴリズムに送り、明示的な数値最適化法(explicit numerical optimization approaches)によって最適化する。最適化されたパラメータは、CNN回帰されたパラメータと比較されて形状損失になる。残りのステップは教師あり訓練と同じであるが、自己改善訓練は教師ありの事例よりも多くのエポック及び長い時間を要する。全体的なニューラルネットワークの訓練は、形状損失及び色損失の両方を最小化するAdamなどの並列最適化アルゴリズムによって実行され、最適化されたネットワークの重みが出力される。

【図面の簡単な説明】

【0010】

【図1】いくつかの実施形態によるニューラル人間予測のフローチャートを示す図である。

【図2】いくつかの実施形態による、全てのネットワークMVS-3DCNN、HMR MLP及びNeRF MLPの重みが既知である、テンソル表記によって表される前方予測のワークフローを示す図である。

10

20

30

40

50

【図3】いくつかの実施形態による、スーパービジョンを使用してネットワークを訓練するワークフローを示す図である。

【図4】いくつかの実施形態による、自己改善戦略においてネットワークを訓練するワークフローを示す図である。

【図5】いくつかの実施形態による、各ビューのMVS-3DCNNのNeRF-MLPへのアライメントを示す図である。

【発明を実施するための形態】

【0011】

ニューラル人間予測が、画像セット（単一の画像又はマルチビュー画像）から骨格のポーズ、体形、並びに衣服の変位及び外観を含む3D人間モデルを予測することを含む。ニューラル人間予測の実施形態は、ニューラルネットワークの使用方法について説明する。マルチビューニューラル人間予測は、単一画像ベースのモーションキャプチャ（mocap）及び人間リフティング（human lifting）を品質及びロバスト性において上回り、メモリコストの高いまばらな点群を入力として受け取って低速で実行する暗黙的パーツネットワークなどの身体衣服予測ネットワークのアーキテクチャを単純化し、3Dボリューム全体を符号化するNeural Bodyなどの潜在コードベースのネットワークの解像度制限を回避する。

10

【0012】

図1は、いくつかの実施形態によるニューラル人間予測のフローチャートである。ステップ100において、被写体の周囲で撮影された写真セットなどの、入力画像セットI、単一画像、又はマルチビュー画像を入力として取得する。入力Iは、 $N \times w \times h \times c$ のサイズの4Dテンソルとして表され、Nはビューの数であり、w、h、cはそれぞれ画像幅、画像高さ及び画像チャンネルである。カメラは既に校正済みであり、従ってカメラ情報（例えば、カメラパラメータ）は全て既知である。画像前処理として、Detectron2及びimage Grab-Cutなどの既存の手法を使用して被写体のバウンディングボックス及び前景マスクを抽出する。画像はバウンディングボックスによって切り取られ、同じアスペクト比で $w \times h$ のサイズにズームされる。画像境界は黒で塗りつぶされる。

20

【0013】

ニューラルネットワーク（MVS-PERF）102は、入力画像セットを特徴に符号化するマルチビューステレオ3D畳み込みニューラルネットワーク（MVS-3DCNN）104、特徴を人間パラメータに回帰させる人間メッシュ復元多層パーセプトロン（HM-MLP）106、及びMVS-3DCNNを微調整してクエリ3D光線（3D位置及び方向）をRGBカラー及び衣服-身体変位に復号するニューラル輝度場多層パーセプトロン（NeRF-MLP）108という3つのコンポーネントで構成される。

30

【0014】

ステップ104において、深層2DCNNが各ビューから画像特徴を抽出する。各畳み込み層の後には、最後の層を除いてバッチ正規化（BN）層及び整流化線形ユニット（rectified linear unit: ReLU）が続く。2つのダウンサンプリング層も配置される。2DCNNの出力は、 $w/4 \times h/4 \times 32$ のサイズの特徴マップである。

40

【0015】

その後、あるビューを基準ビューとして選択し、その視錐台（view frustum）を透視投影及び近遠面（near far planes）に従って被写体の作業空間全体をカバーするように設定する。この錐台を、近い面及び遠い面の両方に平行なd個の深度面によって近くから遠くにサンプリングする。全ての特徴マップを各深度面に変換してブレンドする。 $i = 1, 2, \dots, N$ であるいずれかのビューiについて、（1をインデックスとする）基準ビューに対する 3×3 のホモグラフィ画像ワーピング行列（homography image warping matrix）が以下の数式によって与えられる。

50

$$H_i(z) = K_i \left(R_i R_i^T + \frac{-R_1^T t_1 + R_i^T t_i}{z} n^T \right) K_i^{-1}$$

【0016】

ここで、 K 、 $[R, t]$ はカメラの固有パラメータ及び外部パラメータを表し、 z は深度面から基準ビューのカメラ中心までの距離であり、 n は深度面の法線方向である。

【0017】

全ての画像が深度面にワーブされた後に、全ての特徴の分散

10

$$\sum_{i=1}^W (V_i - \bar{V}_i)^2 / N$$

によって座標 (u, v, z) におけるコストを決定する。

$$\bar{V}_i$$

は、全てのビューの平均特徴値である。

コストボリュームのサイズは、 $d \times w / 4 \times h / 4$ である。

【0018】

20

ステップ106において、人間メッシュ復元多層パーセプトロン(HMR MLP)が、フラット化層(flatten layer)及びドロップアウト層(dropout layer)によって分離された3層の線形回帰を含む。HMR MLPは、MVS 3DCNNからの特徴量を人体パラメータ reg_{114} に回帰させる。

【0019】

人体パラメータ reg は、SMPL-Xなどの人体パラメトリックモデルを3D裸体メッシュ202に操作することができる。通常、SMPL-X表現 reg は、骨格ポーズ(各関節の3次元回転角)、身長及び体重などの体形を制御するボディブレンドシェイプパラメータ、並びに顔表情を制御するフェイシャルブレンドシェイプパラメータを含む。

reg は、ブレンドシェイプパラメータを使用してTポーズメッシュを構築し、これを線形スキニングモデルの骨格ポーズによってポーズメッシュに変形させる。

30

【0020】

一方では、ステップ108において、コストボリュームがニューラル輝度場(NeRF)などの微分可能なレンダリングMLPに送られる。NeRF MLPは、3D位置 x 及び方向 d によって表されるクエリ光線を4チャンネルカラーRGBにマッピングする関数 M として $c(x, d) = M(x, d, f)$ のように定式化される。 f は、錐台MVS 3DCNN104のコストボリュームからNeRFボリュームへの特徴マップであり、 c は、NeRF MLPネットワークの重みであり、 ρ は、3Dポイントがメッシュ内に存在する場合の確率の占有密度を表す。裸体の占有密度場 ρ_b は、錐台104のメッシュ202(図2)を変換することによって直接取得することができる。また、着衣姿の身体密度場 ρ_c は、3次元変位ベクトル場 D と特徴量マップ f との関数： $\rho_c(D, f)$ として表すことができる。3次元変位ベクトル場 D_{116} は、着衣姿の身体表面204上の点が裸体表面上の点とどのように関連しているかを表す。NeRF MLPを訓練すると、変位ベクトル場 D も最適化される。

40

【0021】

図2は、いくつかの実施形態による、全てのネットワークMVS 3DCNN、HMR MLP及びNeRF MLPの重みが訓練されて固定された、テンソル表記によって表される前方予測のワークフローである。透視投影画像からの画素の全ての光線200を問い合わせることによって、外観画像112がレンダリングされる。いくつかの実施形態では、3D人間予測110が実装される。人体の近くのサンプリングされた点を問い合わせる

50

ことによって、変位フィールド D_{116} が取得される。着衣姿の出力メッシュがテンプレートと同じトポロジーを有する人間パフォーマンスキャプチャタスクでは、各頂点に補間変位ベクトル (*interpolated displacement vector*) を追加することによって、裸体メッシュ $V_{b,202}$ を着衣姿の身体メッシュ $V_{c,204}$ に変形することができる。

【 0 0 2 2 】

図 3 は、いくつかの実施形態による、スーパービジョンを用いてネットワークを訓練するワークフローである。Human3.6Mなどの教師あり訓練データセットは、画像入力 I_{100} だけでなく、グランドトゥルス人間パラメータ gt_{300} 及び裸体メッシュ V_b 、 gt_{302} も含み、通常、これらはセンサ又は既存の手法によって取得される。この事例では、予測される裸体とグランドトゥルスとの差分を合計することによって、形状損失 304 が直接取得される。

10

$$\begin{aligned} \text{Shapeloss} = & w_{\theta} \left\| \theta_{reg} - \theta_{gt} \right\|^2 + w_v \sum \left\| V_b - V_{b,gt} \right\|^2 \\ & + w_j \sum \left\| J_b - J_{b,gt} \right\|^2 + w_{j,2D} \sum \left\| \Pi(J_b) - \Pi(J_{b,gt}) \right\|^2 \end{aligned}$$

ここで、 J は裸体の関節であり、 Π は各カメラビューの 3D 点の透視投影を表す。ネットワークを効果的に訓練するために、各訓練ステップでは、全てのビューが MVS 3DCNN の基準ビューとして順番に選択される。

20

【 0 0 2 3 】

一方で、典型的には画像顕著性 (*image saliency*) に比例する不均一なサンプリング戦略を使用して、入力画像セット 100 から光線 306 がサンプリングされる。高顕著性領域では多くの光線がサンプリングされ、平坦領域又は背景領域からは少ない光線がサンプリングされる。これらの光線は、MVS 3DCNN 104 からの特徴マップと共に NeRF MLP 106 に送られ、NeRF MLP 106 がサンプルの外観 RGB 色 308 をレンダリングする。入力画像内のサンプリングされた色とレンダリングされた色 308 との全ての差分を合計することによって色損失 310 が計算される。

【 0 0 2 4 】

Adamなどの並列化された確率的最適化アルゴリズム (*parallelized stochastic optimization algorithm*) を適用して、形状損失及び色損失の両方を最小化することによって全てのネットワーク MVS 3DCNN、HMR MLP、NeRF MLP の重みを訓練する。

30

【 0 0 2 5 】

図 4 は、いくつかの実施形態による、自己改善戦略においてネットワークを訓練するワークフローである。この事例では、訓練データセットが、注釈又は人間グランドトゥルスパラメータを含まない人間画像のみを提供する。入力セット 100 内の各画像について、回帰されたパラメータ reg_{114} を初期推測として選択することにより、SimplifyX アルゴリズムなどの最適化ベースの予測 400 を適用する。最適化ベースの予測は、最初に各画像上の人間の 2D キーポイントを検出し、非線形最適化を適用して 3D 人間にフィットさせる。

40

これらの 2D キーポイントに (opt_{402} によってパラメータ化された) メッシュ $V_{b,opt_{404}}$ を適用する。

$$\theta_{opt} = \arg \min \sum \left\| \Pi(V_{b,opt}) - K \right\|^2$$

【 0 0 2 6 】

ここで、 K は、キーポイントの検出された 2D 位置を示し、合計は全ての対応するキーポイント及び全てのビューを引き継ぐ。

50

【0027】

非線形最小二乗最適化は数値的に遅く、フィッティング精度は初期推測 reg に依存するが、信頼度は高い。十分なフィッティングの反復後には、 opt がグラントゥールスに近くなる。従って、自己改善訓練ワークフローは、以下に要約するように opt をグラントゥールスに向けて効率的に改善することができる。

自己改善訓練ワークフロー：

以下を実行

MVS - 3DCNN から reg を計算し、入力 I から HMR MLP を計算
 reg を初期推測、I を入力として、SMP Li fy X から opt を計算
 I から光線をサンプリングし、NeRF MLP からサンプリングされた色 c を計算
 Shape Loss 及び Color Loss を計算
 Shape Loss 及び Color Loss を最小化することによって MVS 3DCNN、HMR MLP 及び NeRF MLP のネットワークの重みを更新
 全ての訓練データについて重みが収束するまで反復

【0028】

図5に、いくつかの実施形態による、各ビューのMVS 3DCNNのNeRF MLPへのアライメントを示す。

【0029】

動作時には、例えばゲームスタジオにおけるマーカーレスモーションキャプチャ、又は人間3D表面再構成RGBカメラセットアップなどの、商業的及び/又は個人的マーカーレスパフォーマンスキャプチャ用途においてニューラル人間予測を直接適用することができる。マルチビューニューラル人間予測の実施形態の他の用途は、いずれかの拡張と組み合わせることができるリアルタイムバックボーン技術として、例えば深度センシングの入力、3Dモデリング、又は新規アニメーションを作成するための出力の使用を組み合わせることができる。マルチビューニューラル人間予測は、ゲーム用途、VR/AR用途、及びいずれかのリアルタイムヒューマンインタラクション用途において適用することもできる。マルチビューニューラル人間予測は、使用するハードウェア（例えば、GPUプロセッサの速度及びGPUメモリのサイズ）に応じて、予測のために少量のビューを処理する際にはリアルタイムとし、より多くのビュー（例えば、20）の場合には近リアルタイム処理及び予測を実装することができる。

【0030】

本明細書で説明した方法は、いずれかのコンピュータ装置上に実装することができる。好適なコンピュータ装置の例としては、パーソナルコンピュータ、ラップトップコンピュータ、コンピュータワークステーション、サーバ、メインフレームコンピュータ、ハンドヘルドコンピュータ、携帯情報端末、セルラ/携帯電話機、スマート家電、ゲーム機、デジタルカメラ、デジタルカムコーダ、カメラ付き電話機、スマートホン、ポータブル音楽プレーヤ、タブレットコンピュータ、モバイル装置、ビデオプレーヤ、ビデオディスクライタ/プレーヤ（DVDライタ/プレーヤ、高精細ディスクライタ/プレーヤ、超高精細ディスクライタ/プレーヤなど）、テレビ、家庭用エンターテイメントシステム、拡張現実装置、仮想現実装置、スマートジュエリ（例えば、スマートウォッチ）、車両（例えば、自動走行車両）、又はその他のいずれかの好適なコンピュータ装置が挙げられる。

【0031】

顔表情、身体ポーズ形状及び衣服パフォーマンスキャプチャのための暗黙的微分可能レンダラーを用いたマルチビューニューラル人間予測のいくつかの実施形態

1. 装置の非一時的なものにプログラムされた方法であって、

画像セットを入力として取得することと、

ニューラルネットワークを使用して画像セットを処理することと、

を含み、処理は、

画像セットを1又は2以上の特徴に符号化することと、

特徴を人間パラメータに回帰させることと、

ニューラルネットワークを微調整することと、
クエリ3D光線を、画像セットに基づくRGBカラー及び衣服 - 身体変位に復号することと、
を含む、方法。

【0032】

2. 画像セットは、サイズ $N \times w \times h \times c$ の4Dテンソルを含み、ここで、 N はビューの数、 w は画像の幅、 h は画像の高さ、 c は画像のチャンネルである、条項1の方法。

【0033】

3. ニューラルネットワークは、画像セットから正面ビューを基準ビューとして選択し、特徴量を抽出する、条項1の方法。

【0034】

4. ニューラルネットワークは、全ての特徴量を人間のポーズ、形状、表情パラメータに回帰させる、条項3の方法。

【0035】

5. ニューラルネットワークは、パラメータに従って人間の裸体メッシュを生成する、条項4の方法。

【0036】

6. 裸体メッシュは、バウンディングボックス内の占有フィールドに変換される、条項5の方法。

【0037】

7. ニューラルネットワークは、ビューの各中心からの光線方向に関連する身体メッシュの近くのいずれかの3D点について、RGBカラーと、裸体の表面を示す3D変位ベクトルとを生成する、条項6の方法。

【0038】

8. カメラビューの全ての画素から放たれる全ての光線を問い合わせることにより、着衣姿の人体の外観がRGB画像としてレンダリングされ、サンプリングされた点から3D変位ベクトルを使用して裸体を変形させることにより、着衣姿の身体メッシュが取得される、条項7の方法。

【0039】

9. ニューラルネットワークは、教師ありモード又は自己教師ありモードで実装される、条項1の方法。

【0040】

10. アプリケーションを記憶するように構成された非一時的メモリと、
アプリケーションを処理するように構成されたプロセッサと、
を備えた装置であって、アプリケーションは、

画像セットを入力として取得し、

ニューラルネットワークを使用して画像セットを処理する、ように構成され、処理は、
画像セットを1又は2以上の特徴に符号化することと、

特徴を人間パラメータに回帰させることと、

ニューラルネットワークを微調整することと、

クエリ3D光線を、画像セットに基づくRGBカラー及び衣服 - 身体変位に復号することと、
を含む、装置。

【0041】

11. 画像セットは、サイズ $N \times w \times h \times c$ の4Dテンソルを含み、ここで、 N はビューの数、 w は画像の幅、 h は画像の高さ、 c は画像のチャンネルである、条項10の装置。

【0042】

12. ニューラルネットワークは、画像セットから正面ビューを基準ビューとして選択し、特徴量を抽出する、条項10の装置。

【0043】

10

20

30

40

50

- 13．ニューラルネットワークは、全ての特徴量を人間のポーズ、形状、表情パラメータに回帰させる、条項12の装置。
【0044】
- 14．ニューラルネットワークは、パラメータに従って人間の裸体メッシュを生成する、条項13の装置。
【0045】
- 15．裸体メッシュは、バウンディングボックス内の占有フィールドに変換される、条項14の装置。
【0046】
- 16．ニューラルネットワークは、ビューの各中心からの光線方向に関連する身体メッシュの近くのいずれかの3D点について、RGBカラーと、裸体の表面を示す3D変位ベクトルとを生成する、条項15の装置。
【0047】
- 17．カメラビューの全ての画素から放たれる全ての光線を問い合わせることにより、着衣姿の人体の外観がRGB画像としてレンダリングされ、サンプリングされた点から3D変位ベクトルを使用して裸体を変形させることにより、着衣姿の身体メッシュが取得される、条項16の装置。
【0048】
- 18．ニューラルネットワークは、教師ありモード又は自己教師ありモードで実装される、条項10の装置。
【0049】
- 19．アプリケーションを記憶するように構成された非一時的メモリと、アプリケーションを処理するように構成されたプロセッサと、を備えた装置であって、アプリケーションは、
入力画像を特徴に符号化するように構成されたマルチビューステレオ3D畳み込みニューラルネットワーク(MVS-3DCNN)と、
特徴を人間パラメータに回帰させるように構成された人間メッシュ復元多層パーセプトロン(HMR MLP)と、
MVS-3DCNNを微調整するように構成され、クエリ3D光線(3D位置及び方向)をRGBカラー及び衣服-身体変位に復号するニューラル輝度場多層パーセプトロン(NeRF MLP)と、
を含む、装置。
【0050】
- 20．画像セットは、サイズ $N \times w \times h \times c$ の4Dテンソルを含み、ここで、 N はビューの数、 w は画像の幅、 h は画像の高さ、 c は画像のチャンネルである、条項19の装置。
【0051】
- 21．MVS-3DCNNは、画像セットから正面ビューを基準ビューとして選択し、特徴量を抽出する、条項20の装置。
【0052】
- 22．HMR MLPは、全ての特徴量を人間のポーズ、形状、表情パラメータに回帰させる、条項21の装置。
【0053】
- 23．パラメータに従って人間の裸体メッシュを生成するように構成されたモデルをさらに備える、条項22の装置。
【0054】
- 24．裸体メッシュは、バウンディングボックス内の占有フィールドに変換される、条項23の装置。
【0055】
- 25．NeRF MLPは、ビューの各中心からの光線方向に関連する身体メッシュの近くのいずれかの3D点について、RGBカラーと、裸体の表面を示す3D変位ベクトルと

を生成する、条項 2 4 の装置。

【 0 0 5 6 】

2 6 . カメラビューの全ての画素から放たれる全ての光線を問い合わせることにより、着衣姿の人体の外観が RGB 画像としてレンダリングされ、サンプリングされた点から 3 D 変位ベクトルを使用して裸体を変形させることにより、着衣姿の身体メッシュが取得される、条項 2 5 の装置。

【 0 0 5 7 】

本発明の構成及び動作の原理を容易に理解できるように、詳細を含む特定の実施形態に関して本発明を説明した。本明細書におけるこのような特定の実施形態及びこれらの実施形態の詳細についての言及は、本明細書に添付する特許請求の範囲を限定することを意図したものではない。当業者には、特許請求の範囲によって定められる本発明の趣旨及び範囲から逸脱することなく、例示のために選択した実施形態において他の様々な修正を行えることが容易に明らかになるであろう。

【符号の説明】

【 0 0 5 8 】

- 1 0 0 画像入力 I
- 1 0 2 ニューラルネットワーク (M V S - P E R F)
- 1 0 4 マルチビューステレオ 3 D 畳み込みニューラルネットワーク (M V S - 3 D C N N)
- 1 0 6 人間メッシュ復元多層パーセプトロン (H M R M L P)
- 1 0 8 ニューラル輝度場多層パーセプトロン (N e R F M L P)
- 1 1 0 3 D 人間予測
- 1 1 2 外観レンダリング
- 1 1 4 人体パラメータ θ_{reg}
- 1 1 6 3 次元変位ベクトル場 D

【 図 面 】

【 図 1 】

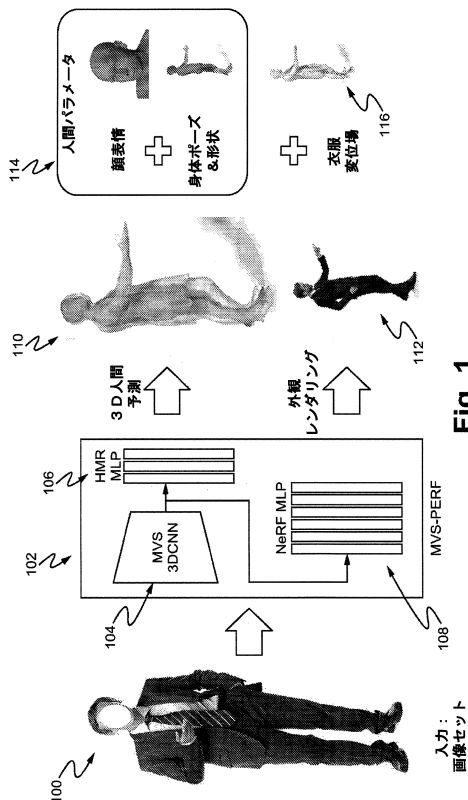


Fig. 1

【 図 2 】

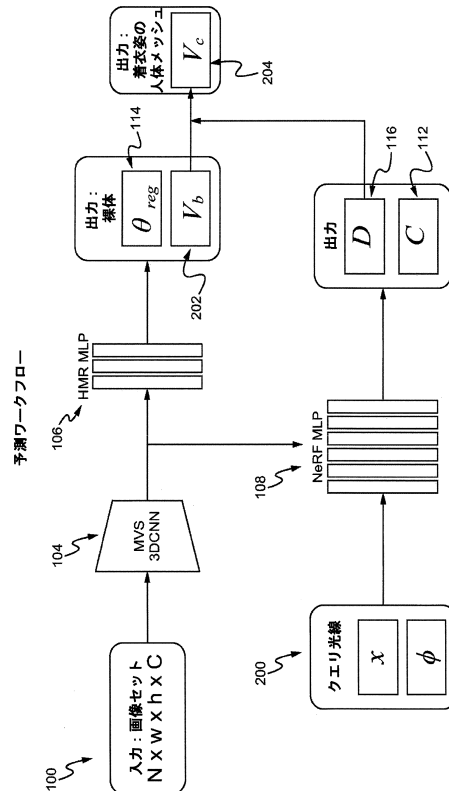


Fig. 2

10

20

30

40

50

【 図 3 】

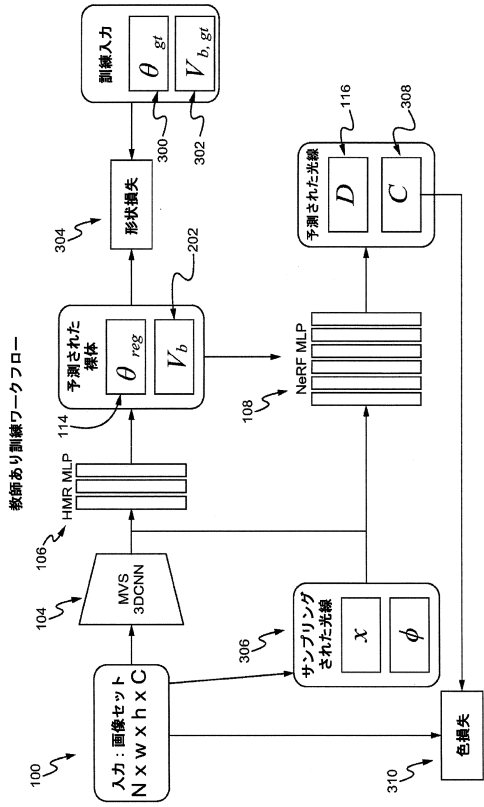


Fig. 3

【 図 4 】

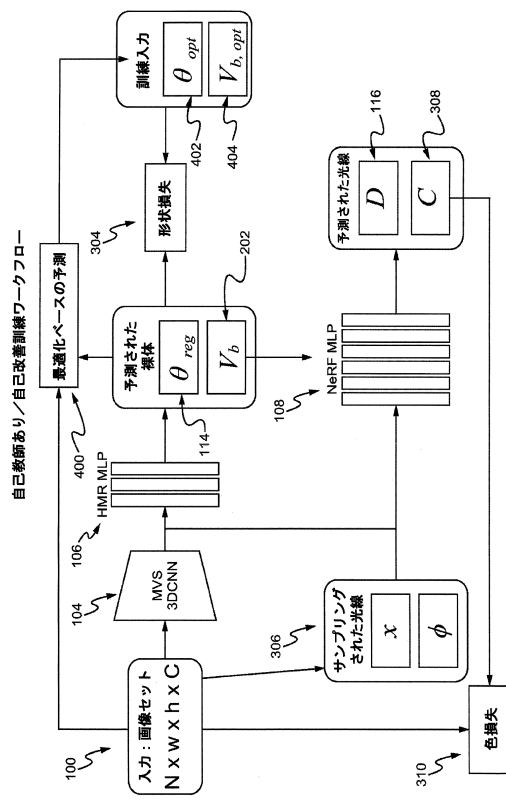


Fig. 4

【 図 5 】

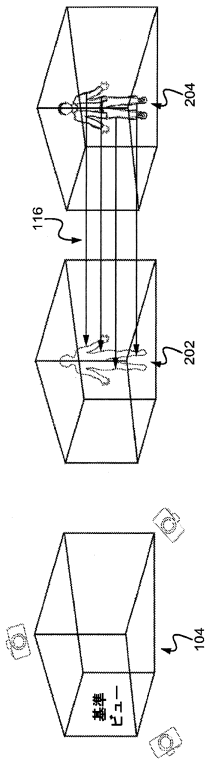


Fig. 5

10

20

30

40

50

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No
PCT/IB2022/053034

A. CLASSIFICATION OF SUBJECT MATTER		
INV.	G06V40/10	G06V10/82
		G06T15/08
		G06T19/20
ADD.		
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)		
G06V G06T		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
EPO-Internal		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	AMIT RAJ ET AL: "PVA: Pixel-aligned Volumetric Avatars", ARKIV.ORG, CORNELL UNIVERSITY LIBRARY, 201 OLIN LIBRARY CORNELL UNIVERSITY ITHACA, NY 14853, 7 January 2021 (2021-01-07), XP081854400, abstract page 5, section 4.1 page 3, section 3, first paragraph page 5, left-hand column, first paragraph page 4, section 3.2 page 5, section 3.6, first line to section 4.1, last line figure 2 ----- -/--	1-26
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input 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	"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	
"E" earlier application or patent but published on or after the international filing date	"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	
"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)	"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	
"O" document referring to an oral disclosure, use, exhibition or other means	"&" document member of the same patent family	
"P" document published prior to the international filing date but later than the priority date claimed		
Date of the actual completion of the international search	Date of mailing of the international search report	
17 June 2022	28/06/2022	
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016	Authorized officer Sagrebín-Mitzel, M	

10

20

30

40

1

50

INTERNATIONAL SEARCH REPORT

International application No
PCT/IB2022/053034

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>SIDA PENG ET AL: "Neural Body: Implicit Neural Representations with Structured Latent Codes for Novel View Synthesis of Dynamic Humans", ARXIV.ORG, CORNELL UNIVERSITY LIBRARY, 201 OLIN LIBRARY CORNELL UNIVERSITY ITHACA, NY 14853, 29 March 2021 (2021-03-29), XP081901285, the whole document</p> <p>-----</p>	1-26
A	<p>LI ZHONGGUO ET AL: "Learning to Implicitly Represent 3D Human Body From Multi-scale Features and Multi-view Images", 2020 25TH INTERNATIONAL CONFERENCE ON PATTERN RECOGNITION (ICPR), IEEE, 10 January 2021 (2021-01-10), pages 8968-8975, XP033909971, DOI: 10.1109/ICPR48806.2021.9412556 [retrieved on 2021-04-22] the whole document</p> <p>-----</p>	1-26
A	<p>SHIH-YANG SU ET AL: "A-NeRF: Surface-free Human 3D Pose Refinement via Neural Rendering", ARXIV.ORG, CORNELL UNIVERSITY LIBRARY, 201 OLIN LIBRARY CORNELL UNIVERSITY ITHACA, NY 14853, 11 February 2021 (2021-02-11), XP081881269, the whole document</p> <p>-----</p>	1-26

10

20

30

40

1

50

フロントページの続き

(32)優先日 令和4年3月23日(2022.3.23)

(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,AT,BE,BG,CH,CY,CZ,DE,DK,EE,ES,FI,FR,GB,GR,HR,HU,IE,IS,IT,LT,LU,LV,MC,MK,MT,NL,NO,PL,PT,RO,RS,SE,SI,SK,SM,TR),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,JM,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

(特許庁注：以下のものは登録商標)

1 . F A C E B O O K

上杉 浩

(74)代理人 100141553

弁理士 鈴木 信彦

(72)発明者 ジャン チン

アメリカ合衆国 カリフォルニア州 9 5 1 1 2 サンノゼ ノース ファースト ストリート 1 7 3
0 エムエス 3ダブリュ

(72)発明者 シャオ ハンユエン

アメリカ合衆国 カリフォルニア州 9 0 0 0 7 ロサンゼルス ウェスト サーティース ストリート
1 2 4 7 アpartment 3 1 3

F ターム(参考) 5L096 FA16 FA67 FA69 HA11 JA11 KA04

【要約の続き】

された輝度場は、入力された被写体のフリービューボリュメトリックレンダリング(free-view volumetric rendering)に使用される。

【選択図】 図1