

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6922005号  
(P6922005)

(45) 発行日 令和3年8月18日(2021.8.18)

(24) 登録日 令和3年7月30日(2021.7.30)

(51) Int.Cl.		F I			
<b>G06T</b>	<b>7/00</b>	<b>(2017.01)</b>	G06T	7/00	350C
<b>G06T</b>	<b>7/11</b>	<b>(2017.01)</b>	G06T	7/11	
<b>G06N</b>	<b>3/04</b>	<b>(2006.01)</b>	G06N	3/04	

請求項の数 19 外国語出願 (全 27 頁)

(21) 出願番号	特願2020-7450 (P2020-7450)	(73) 特許権者	000003078
(22) 出願日	令和2年1月21日(2020.1.21)		株式会社東芝
(65) 公開番号	特開2020-119568 (P2020-119568A)		東京都港区芝浦一丁目1番1号
(43) 公開日	令和2年8月6日(2020.8.6)	(74) 代理人	100108855
審査請求日	令和2年3月2日(2020.3.2)		弁理士 蔵田 昌俊
(31) 優先権主張番号	1900883.8	(74) 代理人	100103034
(32) 優先日	平成31年1月22日(2019.1.22)		弁理士 野河 信久
(33) 優先権主張国・地域又は機関	英国 (GB)	(74) 代理人	100179062
			弁理士 井上 正
		(74) 代理人	100075672
			弁理士 峰 隆司
		(74) 代理人	100153051
			弁理士 河野 直樹
		(74) 代理人	100162570
			弁理士 金子 早苗

最終頁に続く

(54) 【発明の名称】 コンピュータビジョンのシステムおよび方法

(57) 【特許請求の範囲】

【請求項1】

画像を受け取ることと、

第1の特徴マップを作るために共通の処理ステージを用いて前記画像を処理することと

、  
前記第1の特徴マップを受け取る第1および第2の並列分岐を備える並列処理ステージに前記第1の特徴マップを入力することと、

セマンティックセグメント化された画像を作るために前記第1および第2の分岐の出力を融合ステージで結合することと、

を備え、

前記融合ステージは、前記第1の分岐の前記出力をアップサンプリングすることと、前記第1の分岐の前記アップサンプリングされた出力を前記第2の分岐の前記出力に加算することと、を備え、前記第1の分岐の前記アップサンプリングされた出力は、前記第2の分岐の前記出力への加算の前に深さ方向畳み込みを受け、

前記共通の処理ステージは、ニューラルネットワークを備え、前記ニューラルネットワークは、第1の特徴マップを作るために分離可能な畳み込みを実行し、前記画像をダウンサンプリングするように構成された少なくとも1つの分離可能な畳み込みモジュールを有し、前記第1の分岐は、分離可能な畳み込みを実行するように構成された少なくとも1つの分離可能な畳み込みモジュールを備えるニューラルネットワークを備える、

画像をセグメント化する画像処理方法。

## 【請求項 2】

第2の分岐は、前記第1の特徴マップが前記第1の分岐の前記出力と結合されることを可能にするためにスキップ接続として機能する、請求項1に記載の画像処理方法。

## 【請求項 3】

前記第1および第2の分岐からの前記出力は、1つのステージのみで結合される、請求項1または2に記載の画像処理方法。

## 【請求項 4】

前記第1の分岐内の層の数は、前記共通の処理ステージ内の層の数より多い、請求項1から3のいずれかに記載の画像処理方法。

## 【請求項 5】

前記加算することは、アップサンプリングされ、深さ方向畳み込みされた前記第1の分岐の出力と前記第2の分岐との対応する値を加算することを備える、請求項1から4のいずれかに記載の画像処理方法。

10

## 【請求項 6】

前記融合ステージにおける前記深さ方向畳み込みは、1より大きい膨張係数を用いて実行される、請求項1に記載の画像処理方法。

## 【請求項 7】

前記第1の分岐のアップサンプリングされた出力は、1より大きい膨張係数を有する深さ方向畳み込みを受け、前記第2の分岐の前記出力は、加算の前に2次元畳み込みを受ける、請求項6に記載の画像処理方法。

20

## 【請求項 8】

前記第1および第2の分岐の前記結合された出力は、分類器によって処理される、請求項1から7のいずれかに記載の画像処理方法。

## 【請求項 9】

前記第1の分岐のステージ内の前記分離可能な畳み込みモジュールは、深さ方向畳み込みモジュールである、請求項1から8のいずれかに記載の画像処理方法。

## 【請求項 10】

前記第1の分岐内の前記分離可能な畳み込みモジュールは、深さ方向の分離可能な畳み込みモジュールである、請求項1から9のいずれかに記載の画像処理方法。

## 【請求項 11】

前記第1の分岐内の前記分離可能な畳み込みモジュールは、ボトルネックアーキテクチャモジュールである、請求項1から10のいずれかに記載の画像処理方法。

30

## 【請求項 12】

複数のボトルネック残差アーキテクチャモジュールが存在する、請求項11に記載の画像処理方法。

## 【請求項 13】

前記複数のボトルネック残差アーキテクチャモジュールの後にピラミッドプーリングモジュールが設けられる、請求項12に記載の画像処理方法。

## 【請求項 14】

モデルをトレーニングする方法であって、前記モデルは、画像をセマンティックセグメント化するモデルであり、前記モデルは、

40

第1の特徴マップを作るための共通の処理ステージと、

前記第1の特徴マップを受け取る第1および第2の並列分岐を備える並列処理ステージと、

セマンティックセグメント化された画像を作るために前記第1および第2の分岐の出力を結合する融合ステージと、

を備え、

前記融合ステージは、前記第1の分岐の前記出力をアップサンプリングすることと、前記第1の分岐の前記アップサンプリングされた出力を前記第2の分岐の前記出力に加算することと、を備え、前記第1の分岐の前記アップサンプリングされた出力は、前記第2の

50

分岐の前記出力への加算の前に深さ方向畳み込みを受け、

前記共通の処理ステージはニューラルネットワークを備え、前記ニューラルネットワークは、第1の特徴マップを作るために、分離可能な畳み込みを実行し、前記画像をダウンサンプリングするように構成された少なくとも1つの分離可能な畳み込みモジュールを有し、前記第1の分岐は、分離可能な畳み込みを実行するように構成された少なくとも1つの分離可能な畳み込みモジュールを備えるニューラルネットワークを備え、

前記トレーニングする方法は、

トレーニングデータを提供することと、ここで前記トレーニングデータは、画像および前記画像に関するセマンティックセグメント化された情報を備え、

入力として前記画像を使用し、出力として前記セマンティックセグメント化された情報を使用して前記モデルをトレーニングすることと、ここで前記共通の処理ステージおよび前記並列処理ステージは一緒にトレーニングされる、

を含む、方法。

【請求項15】

トレーニング中にフィルタの個数を適応させることと、より重要性の低いフィルタを破棄するためにその個数を減らすこととをさらに備える、請求項14に記載の方法。

【請求項16】

前記共通の処理ステージおよび/または第1の処理分岐ステージに対して少なくとも1つの追加出力を加えるためにトレーニング中に前記モデルを適応させることをさらに備え、前記方法は、入力として前記画像を使用してトレーニングすることと、前記出力と前記少なくとも1つの追加出力との両方での前記セマンティックセグメント化された情報の両方との比較によって損失を決定することと、両方の出力からの前記決定された損失を使用することによってトレーニング中に重みを更新することとをさらに備える、請求項14に記載の方法。

【請求項17】

インターフェースとプロセッサと

を備え、

前記インターフェースは、画像入力を有し、画像を受け取るように適応され、

前記プロセッサは、

第1の特徴マップを作るために共通の処理ステージを用いて前記画像を処理し、

前記第1の特徴マップを受け取る第1および第2の並列分岐を備える並列処理ステージに前記第1の特徴マップを入力し、

セマンティックセグメント化された画像を作るために前記第1および第2の分岐の出力を融合ステージで結合するように適応され、

前記融合ステージは、前記第1の分岐の前記出力をアップサンプリングすることと、前記第1の分岐の前記アップサンプリングされた出力を前記第2の分岐の前記出力に加算することと、を備え、前記第1の分岐の前記アップサンプリングされた出力は、前記第2の分岐の前記出力への加算の前に深さ方向畳み込みを受け、

前記共通の処理ステージは、ニューラルネットワークを備え、前記ニューラルネットワークは、第1の特徴マップを作るために分離可能な畳み込みを実行し、前記画像をダウンサンプリングするように構成された少なくとも1つの分離可能な畳み込みモジュールを有し、前記第1の分岐は、分離可能な畳み込みを実行するように構成された少なくとも1つの分離可能な畳み込みモジュールを備えるニューラルネットワークを備える、

画像をセグメント化する画像処理システム。

【請求項18】

車両の検出システムであって、前記検出システムは、画像を受け取り、前記画像をセグメント化することによって前記画像から物体を決定するように適応された、請求項17の画像処理システムを含む、検出システム。

【請求項19】

請求項1から16のいずれかの方法をコンピュータに実行させるように適応されたコン

10

20

30

40

50

コンピュータ可読命令を担持する非一時的キャリアメディア。

【発明の詳細な説明】

【技術分野】

【0001】

諸実施形態はコンピュータビジョンのシステムおよび方法に関する。

【背景技術】

【0002】

多数のコンピュータビジョンタスク、たとえば物体の認識およびレジストレーションは、画像の領域にラベルが与えられる、画像のセグメンテーションを必要とする。セマンティックセグメンテーションは、詳細な画素レベルの仕様を提供し、これは、障害物検出および正確な境界検出をしばしば必要とするアプリケーションに特に適する。そのようなアプリケーションは、自律車両および運転者支援、組み込みデバイス、ならびにウェアラブルデバイスを含むが、これに限定はされない。

10

【0003】

近代のセマンティックセグメンテーション方法は、非常に正確な結果を達成するが、しばしば、効率の低下という犠牲を払う。畳み込みニューラルネットワーク(CNN)の最近の開発は、これらのネットワークによって達成される結果の顕著な改善を示す。しかし、これらの有効性は、モデルに含まれるオペレーションおよびパラメータの数に大きく依存する。最近のセマンティックセグメンテーション方法は、処理がハイエンドグラフィックス処理ユニット(GPU)上で実行される場合であっても、単一の画像の物体分類を実行するのに1秒超を要する。これらの方法の複雑さは、リアルタイムアプリケーションでのそれらの展開を妨げる。

20

【0004】

自律運転は、複雑なタスクであり、物体の検出および分類は、他のタイムクリティカルなタスクの前処理ステップにすぎない。そのようなシステムは、しばしば、リアルタイムよりも高速の能力を有するシステムから利益を得る。したがって、物体の検出および分類システムが、物体分類の精度を損なうことなく、リアルタイムよりも高速の性能が可能である、セマンティックセグメンテーションの新しい手法が必要である。

【図面の簡単な説明】

【0005】

【図1】車両で実施される実施形態に係るシステムの図。

【図2】図1のシステム内で使用され得るニューラルネットワークの簡潔なフロー図。

【図3】物体分類用の高速ニューラルネットワークの詳細なフロー図。

【図4】一実施形態によるネットワークアーキテクチャの概略図。

【図5】一実施形態によるピラミッドプーリングモジュールアーキテクチャの概略図。

【図6】物体分類用のニューラルネットワークの例のトレーニングプロセスのフロー図。

【図7】パラメータ枝刈りプロセスを示すフロー図。

【図8】一実施形態による都市景観に対するトレーニング曲線の図。

【図9】都市景観妥当性検査セットに関する視覚的比較の図。第1列は入力RGB画像、第2の列はグラウンドトゥールラベル、最後の列はFast-SCNNの出力である。

40

【図10】Fast-SCNNのセグメンテーション結果の視覚的比較の図。第1列は入力RGB画像、第2の列はFast-SCNNの出力、最後の列はスキップ接続の寄与を0にした後のFast-SCNNの出力である。

【図11】一実施形態に従って使用され得るデバイスの概略図。

【発明を実施するための形態】

【0006】

一実施形態では、

画像を受け取ることと、

第1の特徴マップを作るために共通の処理ステージを用いて前記画像を処理することと

50

前記第1の特徴マップを並列処理ステージに入力することと、ここで前記第2の処理ステージは、第1の特徴マップを受け取る第1および第2の並列分岐を備え、

セマンティックセグメント化された画像を作るために第1および第2の分岐の出力を結合することとを備え、

共通の処理ステージは、ニューラルネットワークを備え、ニューラルネットワークは、第1の特徴マップを作るために分離可能(separable)な畳み込みを実行し、画像をダウンサンプリングするように構成された少なくとも1つの分離可能な畳み込みモジュールを有し、前記第1の分岐は、分離可能な畳み込みを実行するように構成された少なくとも1つの分離可能な畳み込みモジュールを備えるニューラルネットワークを備える、

画像をセグメント化する画像処理方法が提供される。

10

#### 【0007】

開示されるシステムは、コンピュータ技術に結び付けられ、コンピューティングの領域から生じる技術的問題、すなわち、リアルタイムセマンティックセグメンテーションの提供という技術的問題に対処する。開示されるシステムは、2つの並列分岐を有する並列処理ステージに先行する共通の処理ステージによって、この技術的問題を解決する。共通の処理ステージは、画像をダウンサンプリングし、低レベルの特徴を計算する。したがって、共通の処理ステージは2つの特徴をサービスし、それにより低レベルの特徴が抽出されることを可能にし、また、画像からコンテキスト情報を抽出するように構成された分岐に入力することを可能にするために、画像をダウンサンプリングする。

#### 【0008】

20

上記方法は、リアルタイムセマンティックセグメンテーションを可能にし、実際に、高解像度画像データ(1024×2048画素)に対してリアルタイムセマンティックセグメンテーションを超える結果を生み出すことができる。したがって、上記方法は、少メモリ組み込みデバイス上での効率的な計算に適する。上の実施形態は、高速セグメンテーション用の2分岐方法を使用する。第1の「共通の」処理ステージは、多重解像度分岐の浅い低レベル特徴(shallow low level feature)を同時に計算する「ダウンサンプリングを学習する」モジュールと考えられ得る。

#### 【0009】

このアーキテクチャは、高解像度での空間詳細とより低い解像度で抽出される深い特徴(deep feature)との組合せを可能にする。

30

#### 【0010】

上の実施形態では、第1の分岐は、画像またはシーンを構成する物体を分類するか他の形で識別することを目指し、シーン内に存在する物体にラベルを割り当てる。対して、第2の分岐は、主に、高解像度セグメンテーション結果の物体境界に関する情報を維持する処理を担う。

#### 【0011】

上の実施形態では、第1の分岐は、深層ネットワークを備える。

#### 【0012】

上記の実施形態のモデルは、以下の理由から計算コストを削減する。

i) コンテキスト処理では、より複雑で抽象的な特徴を学習する必要があり、したがってより深いネットワークが必要である。しかし、より高解像度の入力は不要であり、ゆえにコンテキスト分岐上でより低解像度の入力を使用することによってモデルコストが低減される。

40

ii) 境界処理では、高解像度入力の必要がある。しかし、大きい視野を見る必要はなく、ゆえに詳細分岐内で少数の層を使用することによってモデルコストが低減される。

iii) したがって、上で述べたように2つの異なる目的のために2つの分岐を動作させることは、全体的なモデルコストを下げる。

iv) 第2の分岐は、スキップ接続と考えられることができ、発明者らは、第1および第2の分岐が、初期処理ステージを共有できることを理解した。

v) 画像処理ネットワークが、少容量になるように設計された。

50

## 【 0 0 1 3 】

一実施形態では、第2の分岐は、第1の特徴マップが第1の分岐の出力と結合されることを可能にするためにスキップ接続として機能する。いくつかの実施形態では、特徴マップは、未変更で第2の分岐を通過される。他の実施形態では、第1の特徴マップは、たとえば、第1の分岐の出力と結合される前に第1の特徴マップの解像度を変更するために、1つまたは複数のニューラルネットワーク層を通過される。一実施形態では、第2の分岐は、1つまたは複数の2D畳み込み層を備える。

## 【 0 0 1 4 】

一実施形態では、第1および第2の分岐からの出力は、1つのステージのみで結合される。一実施形態では、2つの分岐だけがある。

10

## 【 0 0 1 5 】

第1の分岐は、コンテキスト情報が抽出されることを可能にする複数のモジュールを備え、第1の分岐は、深層ネットワークを備え、第1の分岐内の層の個数は、共通処理ステージ内の層の個数より多い。

## 【 0 0 1 6 】

一実施形態では、上記システムは、深さ方向 (depth-wise) の分離可能な畳み込みを使用して、標準的な畳み込み (Conv2d) を、空間畳み込みまたはチャネル方向 (channel-wise) 畳み込みとしても知られる深さ方向の畳み込み (DWConv) とそれに続く  $1 \times 1$  点方向 ( $1 \times 1$  point-wise) 畳み込み層とに因数分解する。したがって、チャネル間 (cross-channel) 相関および空間相関は、独立に計算され、これは、パラメータの数を劇的に低減し、より少数の浮動小数点演算と高速の実行時間とをもたらす。したがって、計算コストおよびメモリ要件が削減される。

20

## 【 0 0 1 7 】

一実施形態では、第1および第2の分岐の出力は、1つのステージのみで結合される。第1および第2の分岐の出力は、1回だけ結合される。

## 【 0 0 1 8 】

一実施形態では、第1および第2の分岐の出力は、融合ステージで結合され、前記融合ステージは、第1の分岐の出力をアップサンプリングすることと、第1の分岐のアップサンプリングされた出力を第2の分岐の出力に加算することとを備え、加算することは、第1の分岐のアップサンプリングされた出力と第2の分岐との対応する値を加算することを備える。

30

## 【 0 0 1 9 】

さらなる実施形態では、第1および第2の分岐の出力は、加算によって結合される。一実施形態では、サイズ  $a * b * c$  の第2の特徴マップ (アップサンプリングの後の低解像度から) は、サイズ  $a * b * c$  の結合された特徴マップCを作成するために、両方の数を加算することによってサイズ  $a * b * c$  の第2の特徴マップと結合される。連結 (concatenation) は不要であり、これは、メモリフットプリントを小さく保つことを可能にする。

## 【 0 0 2 0 】

第1の分岐からのアップサンプリングされた出力は、第2の分岐からの出力への加算の前に深さ方向畳み込みを受け、ここで深さ方向畳み込みは、1より大きい膨張係数 (dilation factor) を用いて実行される。

40

## 【 0 0 2 1 】

さらなる実施形態では、1より大きい膨張係数を有する深さ方向畳み込みを受けた第1の分岐からのアップサンプリングされた出力、および第2の分岐からの出力は、加算の前に2次元畳み込みを受ける。

## 【 0 0 2 2 】

一実施形態では、第1の分岐内の深さ方向畳み込みモジュールは、ボトルネックアーキテクチャモジュールである。ボトルネックアーキテクチャモジュールは、特徴拡張 (feature expansion) を実施し、その後に深さ方向畳み込みが行われ、その後に点畳み込みが

50

行われる。したがって、これらはまた、深さ方向の分離可能な畳み込みを実施しているとも考えられ得る。

【0023】

一実施形態では、第1の分岐の深さ方向畳み込みモジュールは、ボトルネック残差アーキテクチャモジュールである。複数の残差ボトルネックアーキテクチャモジュールが存在し得る。残差ボトルネックアーキテクチャモジュールからの出力チャンネルの個数は、複数の残差ボトルネックアーキテクチャモジュール内の後続モジュールでは増加し得る。

【0024】

第1の特徴マップは、前記深さ方向畳み込みモジュールによって処理される前に、標準的な畳み込みを受け得る。最終的な標準畳み込みモジュールは、残差ボトルネックアーキテクチャモジュールの後に設けられ得る。

10

【0025】

一実施形態では、ピラミッドプーリングモジュールが、最終的な標準畳み込みモジュールと残差ボトルネックアーキテクチャモジュールとの間に設けられる。

【0026】

共通の処理ステージは、複数の深さ方向の分離可能畳み込みモジュールを備えることができる。画像は、共通の処理ステージ内の前記深さ方向の分離可能畳み込みモジュールによって処理される前に、標準的な畳み込みを受け得る。

【0027】

一実施形態では、共通の処理ステージは、低レベル特徴の共有が有効であり、効率的に実施されることを保証するために、3つの層を備える。一実施形態では、第1の層は、標準的な畳み込み層 (Conv2D) であり、残りの2つの層は、深さ方向の分離可能な畳み込み層 (DSConv) である。DSConvは、計算的により効率的であるが、入力画像だけが3つのチャンネルを有し、DSConvの計算的な影響をこのステージで重大なものにするので、Conv2Dがここで使用される。

20

【0028】

一実施形態では、共通の処理ステージの層は、ストライド2を使用し、これに、バッチ正規化とReLUとが続く。畳み込み層および深さ方向層の空間カーネルサイズは、 $3 \times 3$  である。一実施形態では、深さ方向畳み込みと点方向畳み込みとの間の非線形性が、省略される。

30

【0029】

一実施形態では、第1および第2の分岐の結合された出力は、分類器によって処理される。一実施形態では、分類器は、少なくとも1つの分離可能な畳み込みモジュールを備える。分類器は、softmax層を備えることができる。少なくとも1つの分離可能な畳み込みモジュールは、深さ方向の分離可能な畳み込みモジュールである。

【0030】

一実施形態では、浮動小数点乗算が、整数演算または2進演算と比較して高コストなので、ランタイムは、DCNNフィルタおよびアクティブ化値に関する量子化技法を使用して、さらに削減され得る。

【0031】

一実施形態では、2進量子化技法が使用され得る。

40

【0032】

さらに、一実施形態では、事前にトレーニングされたネットワークのサイズを低減するために枝刈り (pruning) が適用され、その結果、より高速のランタイムと、より小さいパラメータセットと、より小さいメモリフットプリントとをもたらす。

【0033】

一実施形態では、モデルをトレーニングする方法が提供され、前記モデルは、画像をセマンティックセグメント化するモデルであり、モデルは、

第1の特徴マップを作るための共通の処理ステージと、

並列処理ステージと、ここで前記第2の処理ステージは、第1の特徴マップを受け取る

50

第 1 および第 2 の並列分岐を備え、

セマンティックセグメント化された画像を作るために第 1 および第 2 の分岐の出力を結合することと

を備え、共通の処理ステージはニューラルネットワークを備え、ニューラルネットワークは、第 1 の特徴マップを作るために、分離可能な畳み込みを実行し、画像をダウンサンプリングするように構成された少なくとも 1 つの分離可能な畳み込みモジュールを有し、前記第 1 の分岐は、分離可能な畳み込みを実行するように構成された少なくとも 1 つの分離可能な畳み込みモジュールを備えるニューラルネットワークを備え、

トレーニング方法は、

トレーニングデータを提供することと、ここでトレーニングデータは、画像および前記画像に関するセマンティックセグメント化された情報を備え、

入力として前記画像を使用し、出力としてセマンティックセグメント化された情報を使用して前記モデルをトレーニングすることとを備え、ここで 2 つのステージは一緒にトレーニングされる。

【 0 0 3 4 】

さらなる実施形態では、パラメータ枝刈りが、トレーニング中に実施され得る。

【 0 0 3 5 】

さらなる実施形態では、トレーニング方法は、前記第 1 の分岐に対して第 2 の出力を加えるためにトレーニング中にモデルを適応させることをさらに備え、方法は、入力として画像を使用してトレーニングすることと、出力と第 2 の出力との両方でのセマンティックセグメント化された情報の両方との比較によって損失を決定することと、両方の出力からの決定された損失を使用することによってトレーニング中に重みを更新することとをさらに備える。

【 0 0 3 6 】

—実施形態では、

インターフェースとプロセッサと

を備え、

前記インターフェースは、画像入力を有し、第 1 の画像を受け取るように適応され、前記プロセッサは、

第 1 の特徴マップを作るために共通の処理ステージを用いて前記画像を処理し、

並列処理ステージに前記第 1 の特徴マップを入力し、前記第 2 の処理ステージは、第 1 の特徴マップを受け取る第 1 および第 2 の並列分岐を備え、

セマンティックセグメント化された画像を作るために第 1 および第 2 の分岐の出力を結合する

ように適応され、ここで、共通の処理ステージは、ニューラルネットワークを備え、ニューラルネットワークは、第 1 の特徴マップを作るために分離可能な畳み込みを実行し、画像をダウンサンプリングするように構成された少なくとも 1 つの分離可能な畳み込みモジュールを有し、前記第 1 の分岐は、分離可能な畳み込みを実行するように構成された少なくとも 1 つの分離可能な畳み込みモジュールを備えるニューラルネットワークを備える

画像をセグメント化する画像処理システムが提供される。

【 0 0 3 7 】

—実施形態では、このシステムは車両上で実施され、このシステムがその上で実施され得る車両は、このシステムに入力高解像度画像を供給する 1 つまたは複数のスチールデジタルカメラおよびまたはビデオデジタルカメラを備えた自律車両と半自律車両とを含むがこれに限定されない。—実施形態では、このシステムは、車両上に配置されたグラフィックス処理ユニットまたは中央処理装置上で実現される。このシステムの目的は、車両の周囲の物体を分類し、車両の最終目的地に向かう車両の安全なナビゲーションを容易にすることである。

【 0 0 3 8 】

したがって、さらなる実施形態では、車両の検出システムであって、前記検出システム

10

20

30

40

50

すなわち、上で説明した画像処理システムを備える検出システムは、前記画像を受信し、前記画像をセグメント化することによって前記画像から物体を決定するように適応される、検出システムが提供される。

【 0 0 3 9 】

図 1 は、システムが車両とともに移動するように車両に搭載されて提供された物体分類ネットワークの実施形態を示すのに使用される。図 1 は、自動車 1 の概略を示し、自動車 1 は、衝突回避システムを備える。衝突回避システムは、4 つのカメラ 3、5、7、9 を備え、カメラ 3、5、7、9 は、自動車 1 の各コーナーに設けられる。カメラのそれぞれは、観察可能な世界のより広い広がりを見ることが可能にする、広い視野 ( F O V ) を有する。一実施形態では、各カメラ 3、5、7、9 は、幅広いパノラマ画像を作る、非常に広角の魚眼レンズを与えられ得る。各カメラからの画像は、別々に作られるが、前部カメラ 3 および 5 の F O V は、視界が遮られる自動車 1 の前の部分内のエリアを残さないように、オーバーラップしてもよい。後部カメラ 7 および 9 の F O V も、視界が遮られる自動車 1 の背後の部分内のエリアを残さないように、オーバーラップしてもよい。

10

【 0 0 4 0 】

このシステムの一実施形態では、各カメラ 3、5、7、9 からの画像は、単一の中央処理装置 ( C P U ) または G P U によって別々に処理される。このシステムのさらなる実施形態では、各カメラ 3、5、7、9 は、別々の C P U または G P U を与えられ、この別々の C P U または G P U は、画像を処理し、処理された画像を自動車 1 の中央 C P U に転送する。

20

【 0 0 4 1 】

上記実施形態は、運転のための自律システムに関係する。しかし、この画像処理方法は、画像のセマンティックセグメンテーションを必要とするすべてのシステム、たとえばウェアラブル技術などにも適用され得る。

【 0 0 4 2 】

図 2 は、物体の識別と分類とに関する、一実施形態によるシステムのフロー図を示す。

【 0 0 4 3 】

一実施形態では、物体分類システムは、多重分岐アーキテクチャを有する畳み込みニューラルネットワークとエンコーダ - デコーダフレームワークとの組合せを備え、異なる解像度レベルでの初期畳み込みは、分岐によって共有される。

30

【 0 0 4 4 】

図 2 には、ニューラルネットワークが 2 つの分岐を備える実施形態による、多重分岐畳み込みニューラルネットワークのアーキテクチャが示されている。第 1 の分岐は、シーンの大域 ( global ) コンテキストを取り込む処理を担い、第 2 の分岐は、空間詳細を用いて大域コンテキストを洗練させる処理を担う。シーンの局所 ( local ) コンテキストは、フル解像度画像から抽出され、したがって、分岐は、少数の畳み込み層を備える。大域コンテキストは、前記分岐内のより多数の畳み込み層を可能にする、より低い画像解像度で取り込まれる。さらに、システムのこの実施形態では、エンコーダ - デコーダフレームワークも使用される。ニューラルネットワークのアーキテクチャ内にスキップ接続を統合することは、浅い低レベルの特徴が 2 つの分岐について同時に抽出されることを可能にする、2 つの分岐の初期層が共有されることを可能にする。スキップ接続は、ランタイム効率のために 1 回だけ使用される。さらに、高解像度分岐での ( すなわち、共有される初期層での ) プーリング畳み込み動作の使用は、ネットワークの低解像度分岐のために画像をより低解像度にダウンサンプリングする必要を回避する。スキップ接続は、ネットワークの初期層内で抽出されるシーンの局所コンテキストが、低解像度サブネットワーク分岐によって抽出される大域特徴とマージされることを可能にする。

40

【 0 0 4 5 】

ステップ S 2 0 1 では、カメラ 3、5、7、9 が、シーンの画像を取り込み、ここでシーンは、1 つまたは複数の物体を備える場合がある。一実施形態では、シーンは、街路シーンである。取り込まれた画像は、物体分類システムに入力され、ここで、シーン内の物

50

体は、その結果として、ニューラルネットワークによって識別され、ラベルを付与される。

#### 【0046】

一実施形態では、ステップ201の後に、「ダウンサンプリングを学習する」モジュールが設けられる。入力画像は、フル解像度で「ダウンサンプリングを学習する」モジュールに供給される。ニューラルネットワークは、ステップS203で、シーンの局所コンテキストを抽出し、入力画像のダウンサンプリングされた表現を生成する。

#### 【0047】

ステップS203の後に、一実施形態では、画像の経路は、2つの独立のサブネットワーク分岐に分かれる。低解像度の分岐である第1の分岐では、画像のダウンサンプリングされた表現が供給され、ステップS205で、シーンの大域コンテキストが、畳み込みニューラルネットワークによって抽出される。第2の分岐は、2つの分岐の特徴が結合されるステップS207で、空間詳細の回復を可能にするスキップ接続として実施される。

#### 【0048】

図2では、「ダウンサンプリングを学習する」モジュールは、ニューラルネットワークの多重解像度分岐の浅い低レベルの特徴を計算する。一実施形態に従ってこれがどのように達成されるのかは、図3のステップS303およびS305を参照して説明される。

#### 【0049】

入力画像は、3次元行列フォーマットで表現され得る。各画素は、3つのチャンネルを備え、各チャンネルは、3つの色すなわち赤、緑、および青(RGB)のうちの1つの強度に関連する数値を保持する。標準計算では、2D畳み込みが、 $h * w * c$ フィルタを使用して使用され得、 $h$ は高さ、 $w$ は幅、 $c$ はチャンネルである。 $h * w * 3$ (RGBを有する画像に対する畳み込み)は、 $h * w * 1$ (DW畳み込み)に類似するので、標準畳み込みが、ここでは申し分ない。通常は、特徴チャンネルの個数がCNNでは素早く増加し、したがって、 $h * w * 32 / 64 / 128$ が、めずらしくはないことに留意されたい。一実施形態では、深さ方向畳み込みが、増加する個数のチャンネルと共に使用される。

#### 【0050】

ステップS303では、畳み込みニューラルネットワークの第1の層は、前の段落で説明したように、より効率的な物体分類を達成するために、2次元プーリング畳み込み層である。プーリング層は、出力される特徴マップの次元が層のストライドに対する相対的な比率だけ減らされる、ダウンサンプリング動作を実行する。層のストライドは、畳み込みカーネルが画像にまたがってスキャンされる時に畳み込みカーネルがどれほど移動されるのかであると考えられ得る。たとえば、ストライド1(1画素)を有するプーリング畳み込み層は、特徴マップの空間次元に影響しないが、ストライド2を有するプーリング畳み込み層は、特徴マップの次元を2倍だけダウンサンプリングする。より小さい特徴マップは、より高速の推論を生じることができ、犠牲にされた予測正確さを犠牲にする。

#### 【0051】

この実施形態では、第1の畳み込み層は、入力ボリュームに複数のフィルタまたはカーネルを適用する。各フィルタはシーンと比較され、2次元特徴マップが生成される。特徴マップは、フィルタとシーン内の物体または特徴との間の一致の結果として引き起こされるアクティブ化の空間的配置を表す。すべての特徴マップが、深さ次元に沿って積み重ねられ、出力ボリュームを作る。たとえば、32個の出力チャンネルからなる畳み込みニューラルネットワークは、32個のフィルタを入力ボリュームに適用し、32の深さ次元を有する出力ボリュームをレンダリングするために積み重ねられる32個の特徴マップを生成する。

#### 【0052】

一実施形態では、物体分類は、複数の深さ方向の分離可能な畳み込み層を利用するニューラルネットワークを用いて実行される。深さ方向畳み込みでは、2次元畳み込みが、各入力チャンネルに対して別々に実行され、入力チャンネルごとの特徴マップが生成される。すべての特徴マップは、畳み込み手順の終りに一緒に積み重ねられる。たとえば、カメラ3

10

20

30

40

50

、5、7、9は、各画素が、1つは赤、1つは緑、1つは青の3つの値のセットを備えるカラーデジタル画像をニューラルネットワークに供給する。画像が深さ方向の分離可能な畳み込み層を用いて処理される時に、2次元畳み込みが、色ごとに別々に実行される。

【0053】

深さ方向畳み込みとそれに続く点方向畳み込みは、標準畳み込み層と比較して、物体検出の正確さを大幅には低下させない。情報は、まずチャンネル方向レベルで計算され、これに、チャンネル方向の情報を結合する、より安価な標準畳み込みが続く。畳み込みは、カーネル $1 \times 1$ のみを使用するので、より高速であり、より少数のパラメータを必要とする。

【0054】

一実施形態による「ダウンサンプリングを学習する」モジュール内の畳み込み層のパラメータが、表1に提示されている。

【表1】

表1

入力	演算子	出力チャンネル	繰返し	ストライド
1048 x 2048 x 3	Conv2d	32	1	2
512 x 1024 x 32	DWConv	48	1	2
256 x 512 x 48	Conv2d (1x1)	48	1	1
256 x 512 x 48	DWConv	64	1	2
128 x 256 x 64	Conv2d (1x1)	64	1	1

【0055】

一実施形態では、プーリング畳み込み層403は、フル解像度画像401に32個のフィルタを適用し、32個の特徴マップを生成する。

【0056】

ステップS305では、 $3 \times 3$ のカーネルサイズを有する2つの深さ方向の分離可能な畳み込みが使用される。各深さ方向の分離可能な畳み込みブロックは、深さ方向畳み込み層とそれに続く $1 \times 1$ 点方向畳み込み層とを備える。図4では、各深さ方向の分離可能な畳み込みの2つの層が、単一のブロックとして表されている。

【0057】

第1の深さ方向の分離可能な畳み込みブロック405は、32個の入力チャンネルと、48個の出力チャンネルと、ストライド2とを有する。

【0058】

第2の深さ方向の分離可能な畳み込みブロック407は、48個の入力チャンネルと、64個の出力チャンネルと、ストライド2とを有する。

【0059】

したがって、この実施形態では、「ダウンサンプリングを学習する」モジュールは、5つのニューラルネットワーク層を備える。

【0060】

図2では、第1のサブネットワーク分岐は、適当なラベルを用いてシーン内の物体にラベルを付けると説明された。これが一実施形態に従ってどのように達成されるのかが、図3のステップS307およびS309を参照して説明される。

【0061】

「ダウンサンプリングを学習する」モジュールの出力では、画像次元すなわち長さ $h$ および幅 $w$ が、 $n$ 分の1に縮小される。したがって、画像は、 $n^2$ 分の1でダウンサンプリングされる。このシステムの計算時間および物体分類の達成される正確さは、両方とも、因数 $n$ の値に反比例する。

【0062】

$n$ という因数は、2と32との間の範囲にわたる可能性があり、「ダウンサンプリング

10

20

30

40

50

を学習する」モジュールに備えられるプーリング層の個数によって決定される。

【0063】

物体分類のためのネットワークの上で説明された実施形態では、入力画像の空間次元は、 $n = 8$ 分の1に縮小される。したがって、画像は、「ダウンサンプリングを学習する」モジュールの出力で64分の1でダウンサンプリングされている。

【0064】

一実施形態では、深さ方向畳み込み層は、ボトルネックブロック内で使用され得る。ボトルネック残差ブロックは、表2に従って、入力を、 $c$ 個のチャネルから、高さ $h$ 、幅 $w$ 、拡張係数 $t$ 、畳み込みタイプのカーネルサイズ/ストライド $s$ 、および非線形関数 $f$ を有する $c'$ 個のチャネルに転送する。

10

【0065】

ボトルネックブロックでは、入力ボリュームは、拡張され、その後、深さ方向畳み込み層とそれに続く点方向畳み込みを用いてフィルタリングされる。

【0066】

一般に、まず、点方向畳み込みが適用される(行1)。その後、深さ方向畳み込みおよび点方向畳み込み(行2および行3が適用される)。

【0067】

一実施形態では、以下の処理が続く。

Conv2d 1 / 1

( $1 \times 1 \times c \times t * c$  個のパラメータ)

( $h \times w \times 1 \times 1 \times c \times t * c$  個の動作)

Conv2d 3 / s

( $3 \times 3 \times t * c \times c'$  個のパラメータ)

( $h / s \times w / s \times 3 \times 3 \times t * c \times c'$  個の動作)

20

【0068】

しかし、代替の実施形態では、以下が使用され得る。

Conv2d 1 / 1

( $1 \times 1 \times c \times t * c$  個のパラメータ)

( $h \times w \times 1 \times 1 \times c \times t * c$  個の動作)

DWConv 3 / s

( $3 \times 3 \times 1 \times t * c$  個のパラメータ)

( $h / s \times w / s \times 3 \times 3 \times 1 \times t * c$  個の動作)

Conv2d 1 / 1

( $1 \times 1 \times t * c \times c'$  個のパラメータ)

( $h / s \times w / s \times 1 \times 1 \times t * c \times c'$  個の動作)

30

【0069】

上の代替実施形態では、より少数の計算が要求され、パラメータは、より少数である。

【0070】

より高解像度での特徴の学習およびより低解像度へのそれらの射影は、特徴学習手順に利益を与える。さらに、ボトルネックブロックでの深さ方向畳み込みの利用は、計算効率をもたらし、メモリフットプリントを大幅に削減する。

40

【0071】

さらなる実施形態では、残差接続が、ボトルネックブロックに組み込まれ、ボトルネック残差ブロックを形成する。残差接続は、入力ボリュームおよび出力ボリュームが同一の空間次元および同一個数の特徴マップを有する場合に限って、ボトルネックブロック内で使用され得る。残差接続は、ボトルネックブロックの入力からその出力に接続された直線の層を表す。追加の接続層は、乗算層にまたがるより効率的な勾配伝搬を可能にし、ニューラルネットワークのトレーニングを改善する。

【0072】

一実施形態では、ステップS307で物体分類システムに使用されるボトルネック残差

50

ブロックは、表 2 に示された構造を有する。第 1 の層は、特徴マップの個数を  $t$  倍だけ増加させることによって入力ボリュームの次元を拡張する標準的な 2 次元畳み込み層である。

【表 2】

表 2

入力	演算子	出力
$h \times w \times c$	<i>Conv2d 1/1, f</i>	$h \times w \times tc$
$h \times w \times tc$	<i>DWConv 3/s, f</i>	$\frac{h}{s} \times \frac{w}{s} \times tc$
$\frac{h}{s} \times \frac{w}{s} \times tc$	<i>Conv2d 1/1, -</i>	$\frac{h}{s} \times \frac{w}{s} \times c'$

10

## 【0073】

ボトルネック残差ブロックの第 2 の層では、深さ方向畳み込みフィルタが、シーンから特徴を抽出するのに使用される。深さ方向畳み込みは、ストライド  $s$  を有し、したがって、特徴マップの出力次元は、空間的に  $s$  倍だけ縮小される。ストライドが 1 である畳み込み層では、特徴マップの空間サイズは影響を受けない。空間サイズの縮小は、後続のネットワーク層内での計算の回数を効果的に減少させる。畳み込みは、各深さで別々に計算されるので、*DWConv* は、標準的な畳み込みに対して大幅に改善する。空間サイズの縮小は、後続のネットワーク層内での計算（パラメータではない）の数を効果的に減少させる。

20

## 【0074】

点方向畳み込みがそれに続く深さ方向畳み込みは、標準的な畳み込み層と比較して、物体検出の正確さを大幅には低下させない。この情報は、まずチャンネル方向レベルで計算され、これに、チャンネル方向情報を結合する、より安価な標準的な畳み込みが続く。畳み込みはカーネル  $1 \times 1$  のみを使用するので、より高速であり、より少数のパラメータを必要とする。

## 【0075】

最後に、ボトルネック残差ブロックでは、深さ方向の分離可能な畳み込み層によって生成された出力ボリュームが、2 次元畳み込みの第 2 の層を使用して、その元々の低次元表現に戻って射影され得る。第 1 および第 2 の 2 次元畳み込みに使用される拡張係数  $t$  は、この 2 つの層に関して同一である。

30

## 【0076】

一実施形態では、ステップ S 307 のボトルネック残差ブロックには、ピラミッドプーリングモジュール S 309 が続く。ピラミッドプーリングモジュールは、画像をより微細なレベルからより粗なレベルへの区分に分割し、それらの中の局所特徴を集約する。空間ピラミッドプーリングは、テストのために任意のサイズの画像 / ウィンドウから表現を生成することと、トレーニング中に、変化するサイズまたはスケールの画像を供給することとを可能にする。

40

## 【0077】

ピラミッドプーリングモジュールが、図 5 に、より詳細に示されている。ここでは、それは、4 つのカーネルからなり、各カーネルの解像度は、それぞれ  $32 \times 64$ 、 $16 \times 32$ 、 $8 \times 16$ 、および  $4 \times 8$  である。カーネルは、ステップ S 307 の特徴マップ出力に対して均等に分布される。結果の 4 つの特徴マップは、単一の特徴マップ出力を生成するために、双線形 (bilinearly) にアップサンプリングされ、一緒に加算される。特徴マップが加算されることに留意することが重要である。これは、メモリ内にすべての解像度を記憶する必要を回避する。

## 【0078】

表 3 は、一実施形態に従って大域コンテキストを取り込むのに使用され得る層の詳細を

50

示す。

【表 3】

表 3

入力	演算子	拡張係数	出力チャンネル	繰返し	ストライド
128 x 256 x 64	ボトルネック	6	64	3	2
64 x 128 x 64	ボトルネック	6	96	3	2
32 x 64 x 96	ボトルネック	6	128	3	1
32 x 64 x 128	PPM	-	128	-	-

10

【 0 0 7 9 】

これは、図 4 にも絵図的に示されている。「ダウンサンプリングを学習する」モジュールの後で、ネットワークは、2つの分岐に分岐する。第1のネットワーク分岐は、9つのボトルネック残差ブロック 4 0 9、4 1 1、4 1 3、4 1 5、4 1 7、4 1 9、4 2 1、4 2 3、4 2 5、および 4 2 7 と、これに続くピラミッドプーリングモジュール 4 2 7 とを備える。第2の分岐 4 4 7 は、スキップ接続を表す。

20

【 0 0 8 0 】

当業者によって了解されるように、ボトルネック残差演算子は、実際には、表 2 を参照して上で説明したように複数の層を介して実施される。しかし、図 4 では、簡易的に、ボトルネック残差演算子が単一のエンティティとして図示されている。

【表 4】

表 4

入力	ブロック	t	c	n	s
1024 x 2048 x 3	Conv2D	-	32	1	2
512 x 1024 x 32	DSCConv	-	48	1	2
256 x 512 x 48	DSCConv	-	64	1	2
128 x 256 x 64	ボトルネック	6	64	3	2
64 x 128 x 64	ボトルネック	6	96	3	2
32 x 64 x 96	ボトルネック	6	128	3	1
32 x 64 x 128	PPM	-	128	-	-
32 x 64 x 128	FFM	-	128	-	-
128 x 256 x 128	DSCConv	-	128	2	1
128 x 256 x 128	Conv2D	-	19	1	1

30

【 0 0 8 1 】

第1および第2のボトルネック残差層 4 0 9 および 4 1 1 は、64個の入力チャンネルと、64個の出力チャンネルと、ストライド1と、6の拡張係数(t)とを有する。

40

【 0 0 8 2 】

第3のボトルネック残差層 4 1 3 は、64個の入力チャンネルと、64個の出力チャンネルと、ストライド2と、6の拡張係数とを有する。

【 0 0 8 3 】

第4および第5のボトルネック残差層 4 1 5 および 4 1 7 は、64個の入力チャンネルと、64個の出力チャンネルと、ストライド1と、6の拡張係数とを有する。

【 0 0 8 4 】

第6のボトルネック残差層 4 1 9 は、64個の入力チャンネルと、96個の出力チャンネルと、ストライド2と、6の拡張係数とを有する。

50

【 0 0 8 5 】

第7および第8のボトルネック残差層4 2 1および4 2 3は、96個の入力チャンネルと、96個の出力チャンネルと、ストライド1と、6の拡張係数とを有する。

【 0 0 8 6 】

第9のボトルネック残差層4 2 5は、96個の入力チャンネルと、128個の出力チャンネルと、ストライド2と、6の拡張係数とを有する。

【 0 0 8 7 】

大域特徴抽出器の最後の層は、ピラミッドプーリング層4 2 7である。ピラミッドプーリング層4 2 7は、128個の入力チャンネルと128個の出力チャンネルとを有する。

【 0 0 8 8 】

－実施形態では、ダウンサンプリングを学習するモジュールは、主に、大域特徴抽出器によって抽出された大域コンテキストを洗練する処理を担う。より高い正確さおよびよりよい物体分類結果を達成するために、ダウンサンプリングを学習するモジュール内の深さ方向の分離可能なブロックは、直接に使用される。深さ方向の分離可能な畳み込みのボトルネック実施態様は、達成される分類正確さより動作の速度が重要である場合に、多数の層からなるネットワーク分岐を好むが、直接手法は、低下した動作速度を犠牲にして、物体分類のより高い精度をもたらす。しかし、ダウンサンプリングを学習することにおけるより少数のネットワーク層は、より多数の動作を補償する。

10

【 0 0 8 9 】

－実施形態では、図3のステップS 3 0 9で、単一の特徴融合モジュールが、ネットワークの2つの分岐の出力ボリュームをマージするのに使用される。特徴マップをマージするプロセスは、システムメモリ内に特徴を保持することを含む。単一の特徴融合ユニットは、低解像度デバイス要件に従う、より効率的な設計を考慮に入れたものである。

20

【 0 0 9 0 】

－実施形態では、2つの分岐の出力ボリュームは、特徴マップの数と空間次元との両方において異なる。より低解像度の分岐の特徴マップの空間次元は、フル解像度の分岐の特徴マップの空間次元より小さい。したがって、より低解像度のサブネットワーク分岐の出力は、アップサンプリング層4 2 9によって処理され、特徴マップは、4倍でアップスケールされる。特徴融合モジュールのアーキテクチャは、表5にも表されている。

【表5】

30

表5

分岐 -1	分岐 -4
-	アップサンプリング x 4
-	DWConv (膨張 4)
Conv2d 1/1 -	3/1, f
	Conv2d 1/1, -
Add, f	

40

【 0 0 9 1 】

さらなる実施形態では、アップサンプリング層に、1とは異なる膨張係数を有する深さ方向畳み込み層4 3 1が続く。膨張畳み込み層は、特徴マップ上の物体の間の空間を増大させる。膨張させる深さ方向畳み込みは、カーネルのサイズだけに影響し、具体的には、カーネルのサイズが、指数関数的に増大される。たとえば、1の膨張係数を有する深さ方向畳み込みは、3 x 3の元々のカーネルサイズを有するが、2の膨張係数を有する畳み込みは、7 x 7のカーネルサイズを有し、4の膨張係数を有する畳み込みは、15 x 15のカーネルサイズを有する。膨張が、カーネルサイズ7 x 7のものであるが、これが9つの計算だけを有することに留意されたい。

【 0 0 9 2 】

50

出力ボリュームの合計は、ボリュームが同一の空間次元と同一個数の特徴マップとを有する場合に限って実施され得る。したがって、2次元畳み込みの層が、2つの分岐の出力、それぞれ433および435が一緒に加算される前に、これらを一般化するのに使用される。この2つの畳み込み層は、2つの分岐の特徴マップが同一の次元を有することを保証する。

【0093】

上で説明された実施形態では、畳み込み層433は、128個の入力チャンネルおよび128個の出力チャンネルを有するが、畳み込み層435は、48個の入力チャンネルおよび128個の出力チャンネルを有する。

【0094】

最終ステップでは、特徴マップは、437で単に直接に一緒に加算される。したがって、メモリ内に記憶される必要があるパラメータの個数は、増加しない。

【0095】

一実施形態では、特徴融合ユニットには、分類器モジュール、図3のステップS311が続く。分類器のアーキテクチャは、表4と図4のブロック439、441、443、445、および447に関して説明され得る。分類器は、2つの深さ方向の分離可能な畳み込み演算439および441と点方向畳み込み443とを使用する。点方向畳み込み層443は、使用される都市景観セグメンテーションデータセット内の19個のクラスのうち1つごとに1つの出力チャンネルの、19個の出力チャンネルのみを備える。分類器の最後の2つの層は、画像の当初の空間次元を復元するアップサンプリング層445と、クラスラ

【0096】

一実施形態による分類器モジュール内の畳み込み層のパラメータが、表5に提示されている。

【表6】

表5

入力	演算子	出力チャンネル	繰返し	ストライド
128 x 256 x 128	DWConv	128	1	1
128 x 256 x 128	Conv2d (1x1)	128	1	1
128 x 256 x 128	DWConv	128	1	1
128 x 256 x 128	Conv2d(1x1)	128	1	1
128 x 256 x 128	Conv2d (1x1)	19	1	1
128 x 256 x 19	アップサンプリング	19	1	1
1048 x 2048 x 19	Soft-max	19	1	1

【0097】

図6は、入力画像内に示された物体にラベルを付けることのできる、物体分類ニューラルネットワークをトレーニングする例のプロセスのフロー図を示す。ニューラルネットワークは、トレーニングの多数のサンプルを処理することと、すべてのサンプルについて、ニューラルネットワークによって生成された出力とトレーニングサンプル内で指定されるターゲット出力との間の誤差に従って各パラメータの重みを調整することとによってトレーニングされ得る。トレーニングされた後に、ニューラルネットワークは、システム、たとえば図2のニューラルネットワークシステム内で展開され得る。トレーニング手順は、1つまたは複数のコンピュータによって実行され得る。

【0098】

ステップS601では、トレーニングシステムが、トレーニングデータのセットを入手する。各データセットは、トレーニング画像とターゲット出力とを備える。入力画像は、1つまたは複数の物体を示す画像の表現である。たとえば、画像は、自律車両または半自律車両の付近に配置された物体を含むことができる。トレーニングデータセットによって表されるトレーニング画像は、互いとは異なり、同様の物体を含んでも含まなくてもよい。トレーニングデータセットは、ニューラルネットワークをトレーニングするのに使用され得る、有限個数のシーンを備える。一実施形態では、標準的なデータ増補技法が、数サンプル画像を拡張するのに使用され得る。増補技法は、0.5から2までの範囲内のランダムなスケール係数と、水平フリップと、変更された色相と、変更された飽和度と、変更された輝度と、変更されたコントラストとを含むが、これに限定はされない。

10

**【0099】**

トレーニングデータセットのトレーニングターゲット出力は、ニューラルネットワークによって生成されるべき物体分類ネットワークの所望の出力を表す。ターゲット出力は、ニューラルネットワークの実際の出力と比較され、重み付けパラメータは、ターゲット出力と生成された出力との間の誤差が縮小されるようにするために調整される。ステップS603では、ニューラルネットワークが、内部パラメータの現在値を使用してサンプル入力画像を処理し、出力画像を生成する。

**【0100】**

ステップS605では、ニューラルネットワークの予測された出力が、トレーニングデータセットのターゲット出力と比較され、予測の誤差が推定される。

20

**【0101】**

その結果、ステップS607では、各内部パラメータの重みが、予測された出力とターゲット出力との間の誤差が最小値まで減らされるようにするために調整される。

**【0102】**

ステップ609では、ニューラルネットワークが、トレーニングデータの異なるセットを与えられ、トレーニングは、トレーニング手順を繰り返し、予測された出力とターゲット出力とのより小さい誤差が達成されるようにするために、ニューラルネットワークの内部パラメータを調整するために、ステップS603に戻る。

**【0103】**

一実施形態では、モデルは、交差エントロピー損失を使用してトレーニングされ、トレーニング中に、重み付き補助損失 (weighted auxiliary loss) が、ダウンサンプリングを学習するモジュールおよび大域特徴抽出モジュールの終りで使用される。この形での損失の重み付けは、セマンティックセグメンテーションの意味のある特徴がダウンサンプリングを学習するモジュールおよび大域特徴抽出モジュールによって抽出され、ネットワークの他の副部分とは独立に学習されることを保証する。一実施形態では、補助損失の重みに0.4がセットされた。一実施形態では、これが、407および427の後に追加の出力を生成することによって達成される (すなわち、新しいsoftmax層が、このステージに導入され、ネットワークの分岐からフォークする)。softmax層の出力が評価される (タスクはセグメンテーションでもある)。3つの層の出力が、重みを更新するのに使用される。勾配降下法が使用され、これが0と1との間の確率値を与えるので、softmaxが、トレーニング中に使用される。推論中に、softmaxとargmaxとの両方の関数が単調に増加するので、高コストのsoftmax計算が、argmaxに置換される。argmaxは、1または0、すなわち、物体が存在するまたは存在しない、を用いてデータにラベルを付ける。

30

40

**【0104】**

ニューラルネットワークのトレーニングプロセスの一実施形態では、ネットワーク枝刈りが実施される。トレーニングの最初のステージでは、特徴マップの個数が2倍にされ、トレーニングは、上で説明された手順を使用して行われる。パラメータの個数は、元々のサイズの1.75倍、1.5倍、1.25倍、および1倍に徐々に減らされ、ここで、トレーニング手順は、パラメータのそれぞれの減少の後に繰り返される。ネットワーク枝刈

50

りは、同一の性能を保ちながらネットワーク内のパラメータの個数を効果的に削減する。これは、ニューラルネットワークアーキテクチャ内で使用されていない特徴を除去することによって達成される。さらに、ネットワーク枝刈りは、必要な特徴だけがネットワークによって学習されるので、より効率的な学習を可能にする。

#### 【0105】

したがって、枝刈りは、通常、パラメータを減らすのに使用される。しかし、本明細書で説明される実施形態では、パラメータの個数が許容できるものなので、枝刈りは、性能を高めるために実行される。したがって、ネットワークの表現力は、2倍にされる（パラメータの個数を2倍にする）。今や、ネットワークは、より多数のフィルタが存在する（個数を2倍にする）ので、はるかにより低速である。しかし、枝刈りは、フィルタの元の個数にもう一度達するのに使用される。ここで、枝刈りは、フィルタの個数を拡大し（より多数を可能にし）、その後、良好である1回を選択する。

10

#### 【0106】

したがって、ニューラルネットワークのトレーニングプロセスのこの実施形態では、フィルタの個数は、トレーニング手順の初めに2倍にされる。これがどのように行われるのかは、図7を参照して説明される。

#### 【0107】

ステップS701では、層ごとのフィルタの個数が2倍にされる。ステップS703では、フィルタのチャンネル数が、前の層の出力のチャンネル数と一致させられる。たとえば、各層が、入力サイズ $h \times w$ と出力サイズ $h' \times w'$ を有するものとする。すべてのこれらの層が、使用されるフィルタの個数に関する深さ/チャンネルを有する。すなわち、 $h \times w \times c$ を与えられて、サイズ $3 \times 3$ （たとえば）の $c'$ 個のフィルタが、 $h' \times w' \times c'$ を作るのに使用され得る。ステップS701と同様に、前の層のフィルタの個数が2倍にされる場合に、 $h \times w \times 2c$ 個の入力があり、サイズ $3 \times 3 \times 2c$ の $2c'$ 個のフィルタが使用される。

20

#### 【0108】

ステップS705では、ネットワークが、図6に関して説明されたようにトレーニングされる。ステップS707では、容量に達したかどうかすなわち、フィルタのターゲット数に達したかどうか決定される。

#### 【0109】

ステップS709では、第1の層内の最弱のフィルタ（画像入力を有する）が、識別され、除去される。これらのフィルタは、多数の異なる行列を識別する可能性があり、一実施形態では、11和が使用される。たとえば、上の次元を使用すると、 $h \times w \times 1.5c$ が、今は入力である。次の層のフィルタは、サイズ $3 \times 3 \times 2c$ であり、したがって、ステップS711では、前の層の除去されるフィルタに関する重みが、サイズ $3 \times 3 \times 1.5c$ のフィルタを得るために除去される。その後、ステップS813では、現在の層内の最弱のフィルタが、決定され、このプロセスが継続される。

30

#### 【0110】

ステップS715では、ステップS711およびS713でのフィルタの除去によって影響を受ける層がまだあるかどうかを調べるためにチェックされる。ある場合には、このプロセスはステップS711にループバックする。最弱のフィルタおよび重みのすべてが除去された後に、このプロセスは、トレーニングステップS707に戻る。ネットワークがトレーニングされた後に、フィルタの個数がさらに減らされる（たとえば、ここでは、フィルタの個数が、 $2 \times$ から $1.75 \times$ に、 $1.5 \times$ に、 $1.25 \times$ に、 $1 \times$ に減らされる）べきかについてチェックされる。さらなる削減が要求される場合には、このプロセスは、さらなるフィルタを除去するためにステップS709に移る。そうではない場合には、このプロセスは終了し、ネットワークは、トレーニングされ、枝刈りされる。

40

#### 【0111】

一実施形態では、バッチ正規化が、トレーニング中にダウンサンプリングを学習するモジュールのすべての層の後で使用される。バッチ正規化は、ニューラルネットワークの各

50

層が他の層とはより独立に学習することを可能にする。バッチ正規化は、前の層の特徴マップ内のアクティブ化をスケールする。ニューラルネットワーク全体のすべてのアクティブ化が、所与の範囲内なので、大きすぎる値または小さすぎる値に関連するアクティブ化はない。これは、より高い学習速度と改善された特徴学習とを可能にする。

#### 【0112】

一実施形態では、ニューラルネットワークノードのセットの脱落が、soft-max層の前に実施される。脱落は、ネットワーク内のニューロン間の相互に依存する学習を減らすトレーニング手法である。トレーニング中に、ノードのランダムなセットが、ネットワークから脱落され、その結果、ネットワークの縮小された版が作成されるようになる。ネットワークの縮小された版は、ニューラルネットワークの他のセクションとは独立に学習し、したがって、ニューロンが互いの間の共依存関係を展開するのを防ぐ。

10

#### 【0113】

上を実証するために、上の切除研究が、都市景観データセットを使用して行われ [ M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, および B. Schiele, The Cityscapes dataset for semantic urban scene understanding. CVPR, 2016年 ]、都市景観テストセット、すなわち都市景観ベンチマークサーバに関する性能を報告する。

#### 【0114】

実験は、CUDA 9.0およびcuDNN V7を用い、Nvidia Titan X (Maxwell, 3072個のCUDAコア)またはNvidia Titan Xp (Pascal, 3840個のCUDAコア)を有するワークステーション上で実行された。ReLUが、ReLU6と比較して、達成される、より高速のトレーニングおよびよりよい正確さに起因して、非線形関数として使用された。トレーニング中に、バッチ正規化が、すべての層で使用され、脱落が、soft-max層の前のみで使用される。推論中に、バッチ正規化のパラメータは、親層の重みおよびバイアスとマージされる。深さ方向畳み込み層では、我々は、L2正則化が不要であることを見出した。一実施形態では、他の層に関して、L2正則化は0.00004である。

20

#### 【0115】

ラベル付けされたトレーニングデータが制限されたので、ランダムスケール0.5から2までと、水平フリップと、変更された色相と、変更された飽和度と、変更された輝度と、変更されたコントラストという標準的なデータ増補技法が、すべての実験で適用された。

30

#### 【0116】

Fast-SCNNのモデルは、Pythonを使用するTensorFlow機械学習プラットフォームを用いてトレーニングされる。0.9の運動量およびバッチサイズ12を有する確率的勾配降下法(SGD)が使用される。さらに、ポリ学習レート(poly learning rate)が、ベースレート0.045および電力0.98を用いて適用された。エポックの最大個数は、事前トレーニングが使用されないので1000にセットされる。

#### 【0117】

都市景観は、ドイツ内の50の異なる都市からの街路シーン内の画像の多様なセットを含む、セマンティックセグメンテーションの大規模データセットである。合計して、それは、25000枚の注釈付き1024x2048画素画像であり、そのうちの5000枚は、高い画素正確さでのラベルを有し、20000枚は、弱く注釈付けされている。本明細書で提示される実験では、5000枚の画像だけすなわち、都市景観評価サーバ上で評価され得る2975枚の画像のトレーニングセットと、500枚の画像の妥当性検査セットと、1525枚の試験画像とが、高いラベル品質を伴って使用された。

40

#### 【0118】

一実施形態では、ImageNetを用いる事前トレーニングが使用された。

#### 【0119】

都市景観は、30個のクラスラベルをも付与するが、19個のクラスだけが、評価に使

50

用される。結果は、平均インターセクションオーバーユニオン (mean intersection-over-union) (mIoU) として報告され、ランタイム評価は、転送推論時間 (forward inference time) を測定するためにシングルスレッド式CPUおよびGPU内で実行される。バースインのために、100フレームが使用され、フレーム毎秒 (fps) 測定のために、100フレームの平均値が報告される。

【0120】

fast-SCNNの全体的な性能は、都市景観の差し控えられたテストセットに対して評価される。表6には、Nvidia Titan X (Maxwell、3072個のCUDA) と、「\*」によって表されるNvidia Titan Xp (Pascal、3840個のCUDAコア) との両方に関して、異なる解像度でのfps単位で比較されたランタイムがある。fast-SCNNの2つの版すなわち、soft-max出力 (我々のprob) および物体ラベル出力 (我々のcls) が示されている。

10

【表7】

表6

	1024 x 2048	512 x 1024	256 x 512
我々の prob	47.5	149.4	310.7
我々の cls	57.8	171.6	379.7
我々の prob*	78.5	244.3	475.0
我々の cls*	91.0	273.5	531.4

20

【0121】

都市景観テストセットを使用するFast-SCNNのクラスおよびカテゴリmIoUが、表7に提示されている。Fast-SCNNは、68.0% mIoUを達成する。このモデルは、少メモリ組み込みデバイス用に設計され、1.1百万パラメータだけを使用する。Fast-SCNNの結果は、定量分析のために図9に表示されている。第1の列には、入力RGB画像があり、第2の列は、グラウンドトゥールスラベルであり、第3の列は、Fast-SCNNの出力である。

30

【表8】

表7

モデル	クラス	カテゴリ	パラメータ
Fast-SCNN	68.0	84.7	01.11

【表9】

表8

入力サイズ	クラス	FPS
1024 x 2048	68.0	91.0
512 x 1024	62.8	273.5
256 x 512	51.9	531.4

40

【0122】

Fast-SCNNが、少ない容量を有するように特に設計されているので、少メモリであることの理由は、組み込みデバイス上での実行を可能にし、よりよい一般化が期待される。提案されるネットワークの性能は、事前トレーニングの有無を伴い、追加の弱くら

50

ベル付けされたデータの有無に関連して評価された。結果は、表9に提示されている。事前トレーニングに関して、ImageNetデータベースが使用され、特徴融合モデルは、平均プーリングによって置換され、分類モジュールは、softmax層だけを備える。ImageNetに対する事前トレーニングは、しばしば、正確さと一般性を押し上げる。ImageNet事前トレーニングを用いるFast-SCNNの正確さは、都市景観の妥当性検査セットに関して69.15% mIoUであるが、Fast-SCNNは、事前トレーニングなしで68.62% mIoUを達成する。

【0123】

さらに、都市景観の都市道路とImageNetの分類タスクとの間のオーバーラップが制限されるので、Fast-SCNNが、両方の領域の制限された能力に起因して利益を得ない可能性があることと仮定することは、穏当である。したがって、都市景観によって付与される追加の20000個の粗にラベル付けされた画像が、類似する領域からのものなので、組み込まれた。それでも、粗なトレーニングデータ(ImageNetありまたはImageNetなし)を用いてトレーニングされたFast-SCNNは、互いと同様に、事前トレーニングなしの元々のFast-SCNNに対するわずかな改善のみを伴って実行する。

10

【0124】

低容量Fast-SCNNが、Imagenetを用いる事前トレーニングから大きくは利益を得ないと結論することができる。同様の結果が、積極的なデータ増補およびより多数のエポックを使用することによって達成され得る。

20

【表10】

表9

モデル	クラス
Fast-SCNN	68.62
Fast-SCNN + ImageNet	69.15
Fast-SCNN + 粗	69.22
Fast-SCNN + 粗 + ImageNET	69.19

30

【0125】

図9に、トレーニング曲線を示す。粗データを用いるFast-SCNNは、弱いラベル品質のゆえに反復に関して低速でトレーニングする。ImageNet事前トレーニング版の両方は、早期エポック(トレーニングセットのみに関して400エポックまで、追加の粗にラベル付けされたデータを用いてトレーニングされる時に100エポックまで)に関してよりよく実行する。これは、Fast-SCNNが一からトレーニングされる時に、同様の正確さに達するためにより長くトレーニングされる必要があることを意味する。

【0126】

上の実施形態は、高解像度画像(1024x2048画素)に関するリアルタイムより高速の物体分類(91.0fps)のためのネットワークに関する。多重分岐ネットワークの計算コストの共有は、ランタイム効率をもたらす。上のアーキテクチャでは、スキップ接続が、空間詳細の回復に関して有益であることを示す。スキップ接続は、小さいサイズの境界および物体の周囲で特に有益である、図10。

40

【0127】

さらに、上の研究は、十分に長くトレーニングされた場合に、追加の補助タスクに対するモデルの大規模事前トレーニングが、低容量ディープ畳み込みニューラルネットワークに関して必要ではないことを示す。

【0128】

図11は、実施形態に従って方法を実施するのに使用され得るハードウェアの概略図で

50

ある。これが、一例にすぎず、他の配置が使用され得ることに留意されたい。

【0129】

ハードウェアは、計算セクション900を備える。この特定の例では、このセクションの構成要素は、一緒に説明される。しかし、これらが、必ずしも同一位置に配置されないことを了解されたい。

【0130】

コンピューティングシステム900の構成要素は、処理ユニット913（中央処理装置、CPUなど）と、システムメモリ901と、システムメモリ901を含む様々なシステム構成要素を処理ユニット913に結合するシステムバス911とを含むがこれに限定されない。システムバス911は、メモリバスもしくはメモリコントローラと、周辺バスと、様々なバスアーキテクチャのいずれかを使用するローカルバスなどを含む複数のタイプのバス構造のいずれとすることもできる。計算セクション900は、バス911に接続された外部メモリ915をも含む。

10

【0131】

システムメモリ901は、読取専用メモリなど、揮発性メモリ/または不揮発性メモリの形のコンピュータ記憶媒体を含む。スタートアップ中などにコンピュータ内の要素の間での情報の転送を助けるルーチンを含む基本入出力システム（BIOS）903が、通常はシステムメモリ901内に記憶される。さらに、システムメモリは、CPU913によって使用中のオペレーティングシステム905と、アプリケーションプログラム907と、プログラムデータ909とを含む。

20

【0132】

また、インターフェース925が、バス911に接続される。インターフェースは、コンピュータシステムがさらなるデバイスから情報を受信するためのネットワークインターフェースとすることができる。インターフェースは、ユーザがある種のコマンドなどに応答することを可能にするユーザインターフェースとすることもできる。

【0133】

この例では、ビデオインターフェース917が設けられる。ビデオインターフェース917は、グラフィックス処理メモリ921に接続されたグラフィックス処理ユニット919を備える。

【0134】

グラフィックス処理ユニット（GPU）919は、この多重並列呼出しの動作に起因して、上で説明される方法に特によく適する。したがって、一実施形態では、処理は、CPU913とGPU919との間で分割され得る。

30

【0135】

一実施形態では、GPUは、低電力GPUチップであるNVIDIA Jetson TX2である。

【0136】

一実施形態では、専用コンピューティングデバイス900は、各カメラ（図1参照）へのリンクを設けられる。上で図2から図4までに関して説明されるアーキテクチャは、浮動小数点計算の必要を回避し、したがって、コンピューティングデバイスは、車両上のカメラと結合されるなど、低電力位置によく適する。

40

【0137】

上で説明されるアーキテクチャは、GPUを使用する携帯電話機にも特に役立つ。

【0138】

ある種の実施形態が説明されたが、これらの実施形態は、例としてのみ提示され、本発明の範囲を限定することは意図されていない。実際に、本明細書で説明される新規のデバイスおよび方法は、様々な他の形で具現化され得、さらに、本明細書で説明されるデバイス、方法、および製品の形態における様々な省略、置換、および変更が、本発明の趣旨から逸脱せずに行われ得る。添付の特許請求の範囲およびその均等物は、本発明の範囲および趣旨に含まれ得るものとしてそのような形態または修正を包含することが意図されてい

50

る。

【図1】

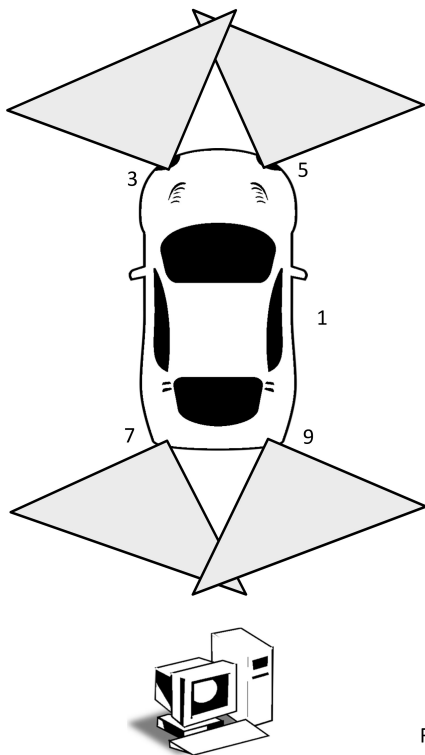


Figure 1

【図2】

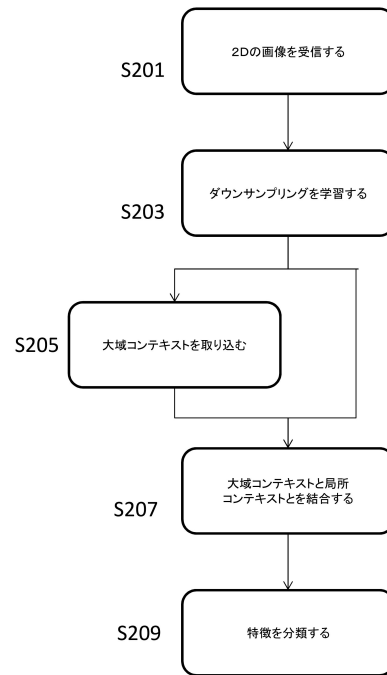


Figure 2

【 図 3 】

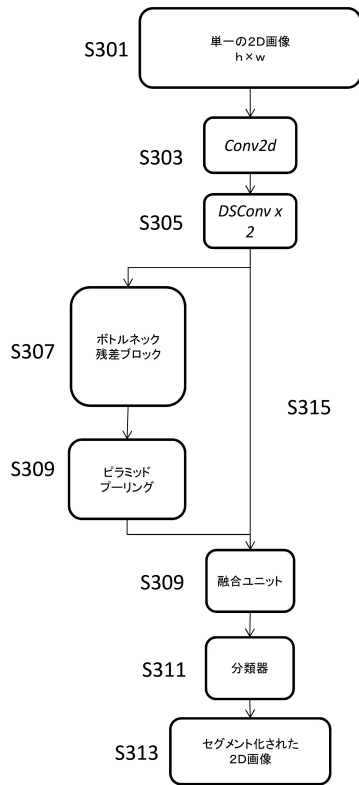


Figure 3

【 図 4 】

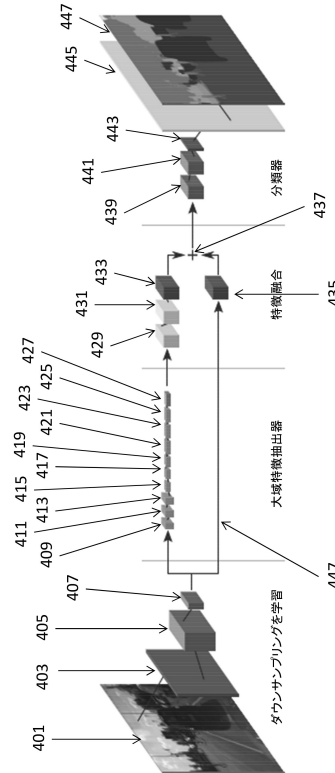


Figure 4

【 図 5 】

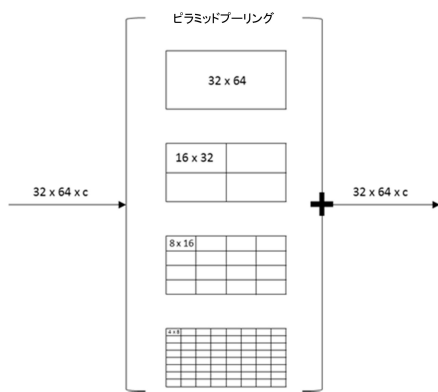


Figure 5

【 図 6 】

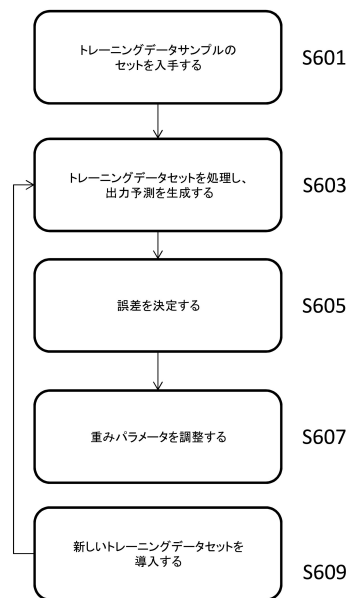


Figure 6

【 図 7 】

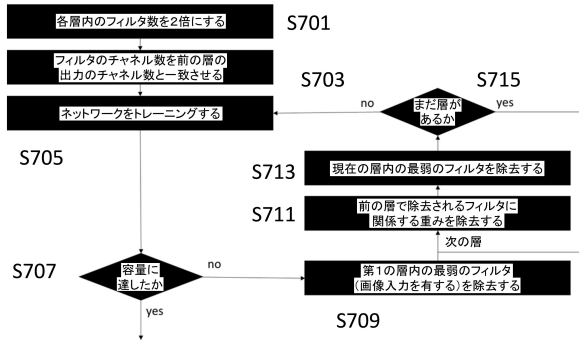


Figure 7

【 図 8 】

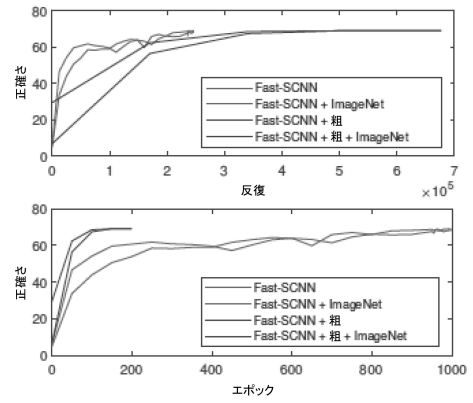


Figure 8

【 図 9 】

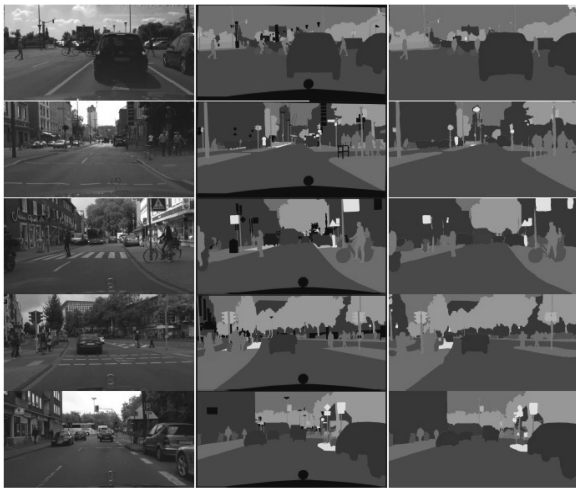


Figure 9

【 図 10 】

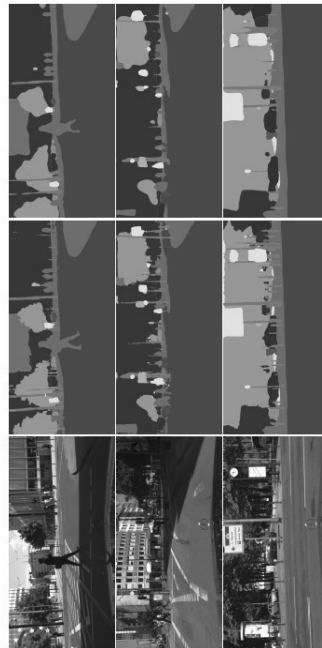


Figure 10

【 図 11 】

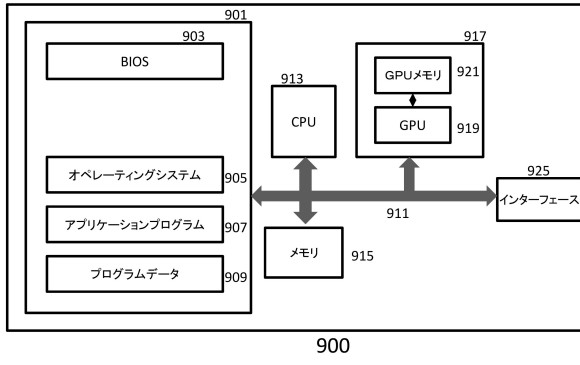


Figure 11

## フロントページの続き

- (72)発明者 ステファン リウィッキ  
イギリス国、 シービー４・０ジーゼット、 ケンブリッジシャー、 ケンブリッジ、 ミルトン  
・ロード、 ケンブリッジ・サイエンス・パーク ２０８、 トーシバ・リサーチ・ヨーロッパ・  
リミテッド、 ケンブリッジ・リサーチ・ラボラトリー内
- (72)発明者 ロベルト シボラ  
イギリス国、 シービー４・０ジーゼット、 ケンブリッジシャー、 ケンブリッジ、 ミルトン  
・ロード、 ケンブリッジ・サイエンス・パーク ２０８、 トーシバ・リサーチ・ヨーロッパ・  
リミテッド、 ケンブリッジ・リサーチ・ラボラトリー内
- (72)発明者 ルドラ プラサド パウデル カルマタ  
イギリス国、 シービー４・０ジーゼット、 ケンブリッジシャー、 ケンブリッジ、 ミルトン  
・ロード、 ケンブリッジ・サイエンス・パーク ２０８、 トーシバ・リサーチ・ヨーロッパ・  
リミテッド、 ケンブリッジ・リサーチ・ラボラトリー内

審査官 岡本 俊威

- (56)参考文献 国際公開第２０１８／１４５０２８(WO, A1)  
特開２０２０-０７１８６２(JP, A)

## (58)調査した分野(Int.Cl., DB名)

G 0 6 T      7 / 0 0 - 7 / 9 0  
G 0 6 N      3 / 0 4  
G 0 6 N      3 / 0 8