

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7638398号
(P7638398)

(45)発行日 令和7年3月3日(2025.3.3)

(24)登録日 令和7年2月20日(2025.2.20)

(51)国際特許分類		F I			
G 0 6 T	7/20 (2017.01)	G 0 6 T	7/20	3 0 0 Z	
G 0 6 V	10/82 (2022.01)	G 0 6 V	10/82		
H 0 4 N	7/18 (2006.01)	H 0 4 N	7/18		K
G 0 6 T	7/00 (2017.01)	G 0 6 T	7/00	3 5 0 C	

請求項の数 15 (全22頁)

(21)出願番号	特願2023-565608(P2023-565608)	(73)特許権者	511151662 中興通訊股 ぶん 有限公司 ZTE CORPORATION 中華人民共和国広東省深 せん 市南山 区高新技术産業園科技南路中興通訊大厦 ZTE Plaza, Keji Road South, Hi-Tech Indu strial Park, Nanshan Shenzhen, Guangdong 518057 China
(86)(22)出願日	令和4年4月24日(2022.4.24)	(74)代理人	100112656 弁理士 宮田 英毅
(65)公表番号	特表2024-516642(P2024-516642 A)	(74)代理人	100089118 弁理士 酒井 宏明
(43)公表日	令和6年4月16日(2024.4.16)	(72)発明者	徐茜
(86)国際出願番号	PCT/CN2022/088692		
(87)国際公開番号	WO2022/228325		
(87)国際公開日	令和4年11月3日(2022.11.3)		
審査請求日	令和5年10月25日(2023.10.25)		
(31)優先権主張番号	202110459730.8		
(32)優先日	令和3年4月27日(2021.4.27)		
(33)優先権主張国・地域又は機関	中国(CN)		

最終頁に続く

(54)【発明の名称】 行動検出方法、電子機器およびコンピュータ読み取り可能な記憶媒体

(57)【特許請求の範囲】

【請求項1】

ビデオストリームから複数フレームのビデオ画像フレームデータを取得するステップと、
複数フレームの前記ビデオ画像フレームデータに基づき、前記ビデオストリームにおける
通行人の行動を検出するステップと、を含み、

複数フレームの前記ビデオ画像フレームデータに基づき、前記ビデオストリームにおける
通行人の行動を検出する前記ステップは、

複数フレームの前記ビデオ画像フレームデータを二次元畳み込みニューラルネットワーク
に入力し、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係と複数フ
レームの前記ビデオ画像フレームデータとに基づき、前記ビデオストリームにおける通行
人の行動を認識するステップを少なくとも含み、

複数フレームの前記ビデオ画像フレームデータに基づき、前記ビデオストリームにおける
通行人の行動を検出する前記ステップは、

複数フレームの前記ビデオ画像フレームデータを二次元畳み込みニューラルネットワー
クに入力し、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係と複数フ
レームの前記ビデオ画像フレームデータとに基づき、前記ビデオストリームにおける通行
人の行動を認識した後に、前記二次元畳み込みニューラルネットワークの出力データに基づ
き、前記ビデオストリームにおける通行人の行動の空間位置を検出するステップをさらに
含み、

前記二次元畳み込みニューラルネットワークは、少なくとも1つの畳み込み層を含み、

複数フレームの前記ビデオ画像フレームデータを二次元畳み込みニューラルネットワークに入力し、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係と複数フレームの前記ビデオ画像フレームデータとに基づき、前記ビデオストリームにおける通行人の行動を認識する前記ステップは、

前記少なくとも1つの畳み込み層により、複数フレームの前記ビデオ画像フレームデータについて特徴抽出を行い、複数フレームの前記ビデオ画像フレームデータに一つ一つに対応し、それぞれが複数の特徴チャンネルを含む複数の特徴マップを取得するステップと、
 複数の前記特徴マップの一部の特徴チャンネルを交換して、複数フレームの前記ビデオ画像フレームデータの時系列情報を融合する特徴データを取得するステップ、とを含み、前記時系列情報は、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係を特徴

10

付けるものであり、
 前記特徴データに基づき前記ビデオストリームにおける通行人の行動を認識するステップと、を含む、

行動検出方法。

【請求項2】

前記二次元畳み込みニューラルネットワークはさらに、少なくとも1つの全結合層とを含み、

前記特徴データに基づき前記ビデオストリームにおける通行人の行動を認識する前記ステップは、

前記少なくとも1つの全結合層により、前記特徴データに基づき前記ビデオストリームにおける通行人の行動を認識するステップを含む、

20

請求項1に記載の行動検出方法。

【請求項3】

前記二次元畳み込みニューラルネットワークは複数の直列接続された前記畳み込み層を含み、

前記少なくとも1つの畳み込み層により、複数フレームの前記ビデオ画像フレームデータについて特徴抽出を行う前記ステップは、

各前記畳み込み層について、前記畳み込み層の入力データを前記畳み込み層に入力して特徴抽出を行うステップを含む、

複数の前記特徴マップの一部の特徴チャンネルを交換して、複数フレームの前記ビデオ画像フレームデータの時系列情報を融合する特徴データを取得する前記ステップは、

30

複数の前記特徴マップの一部の特徴チャンネルを交換して第1データを取得するステップを含み、

前記畳み込み層が最初の畳み込み層である時に、前記畳み込み層の入力データは複数フレームの前記ビデオ画像フレームデータであり、

前記畳み込み層が最後の畳み込み層ではなく、かつ最初の畳み込み層でもない時に、前記第1データを次の畳み込み層の入力データとし、

前記畳み込み層が最後の畳み込み層である時に、前記第1データを前記特徴データとする、

請求項1に記載の行動検出方法。

40

【請求項4】

複数フレームの前記ビデオ画像フレームデータは順に並んだビデオ画像フレームデータをNフレーム含み、複数の前記特徴マップは順に並んだ特徴マップをN個含み、

複数の前記特徴マップの一部の特徴チャンネルを交換して第1データを取得する前記ステップは、

各前記特徴マップにおける複数の特徴チャンネルを、順に並んだN組の特徴チャンネルに分けるステップと、

N個の順に並んだ特徴マップにおけるi番目の特徴マップについて、i番目の特徴マップに対応するj番目の特徴マップを特定し、i番目の特徴マップはN個の前記順に並んだ特徴マップにおけるいずれか1つであり、j番目の特徴マップはN個の前記順に並んだ特

50

徴マップにおけるいずれか1つであるステップと、

i番目の特徴マップにおける第i組の特徴チャンネルをj番目の特徴マップにおけるいずれか1組の特徴チャンネルと交換し、前記第1データを取得するステップと、を含み、

N、i、jは正の整数である、

請求項3に記載の行動検出方法。

【請求項5】

前記少なくとも1つの全結合層により、前記特徴データに基づき前記ビデオストリームにおける通行人の行動を認識する前記ステップは、

前記少なくとも1つの全結合層により、前記特徴データに基づき分類特徴ベクトルを取得し、前記分類特徴ベクトルの各要素が一種類の行動類型に対応するステップと、

前記分類特徴ベクトルに基づき、各種行動類型の分類確率を特定するステップと、

前記各種行動類型の分類確率に基づき前記ビデオストリームにおける通行人の行動を認識するステップと、を含む、

請求項2に記載の行動検出方法。

【請求項6】

各種行動類型の分類確率に基づき前記ビデオストリームにおける通行人の行動を認識する前記ステップは、

各種行動類型の分類確率がフィルタ閾値を上回るかどうか判断するステップと、

少なくとも一種類の行動類型の分類確率が前記フィルタ閾値を上回った時に、対象行動を認識したと判断するステップと、

分類確率が前記フィルタ閾値を上回る行動類型を前記対象行動の類型であると特定するステップと、

各種行動類型の分類確率がいずれも前記フィルタ閾値を上回らない時に、対象行動を認識していないと判断するステップと、を含む、

請求項5に記載の行動検出方法。

【請求項7】

前記二次元畳み込みニューラルネットワークはさらに、少なくとも1つの全結合層とを含み、前記二次元畳み込みニューラルネットワークの出力データは、前記少なくとも1つの全結合層により前記特徴データに基づき取得した分類特徴ベクトルと、対象畳み込み層が出力した前記複数の特徴マップとを含み、前記対象畳み込み層は前記少なくとも1つの畳み込み層における1つであり、前記分類特徴ベクトルの各要素が一種類の行動類型に対応し、

前記二次元畳み込みニューラルネットワークの出力データに基づき、前記ビデオストリームにおける通行人の行動の空間位置を検出する前記ステップは、

前記対象畳み込み層が出力する複数の特徴マップと前記分類特徴ベクトルとに基づき対象行動の空間位置を特定するステップを含む、

請求項1に記載の行動検出方法。

【請求項8】

前記対象畳み込み層が出力する複数の特徴マップと前記分類特徴ベクトルとに基づき対象行動の空間位置を特定する前記ステップは、

前記対象畳み込み層が出力する複数の特徴マップと前記分類特徴ベクトルとに基づき、前記対象行動のエッジ輪郭を特定するステップと、

前記対象行動のエッジ輪郭に基づき前記対象行動の空間位置を特定するステップと、を含む、

請求項7に記載の行動検出方法。

【請求項9】

複数フレームの前記ビデオ画像フレームデータは所定時間長さの複数フレームのビデオ画像フレームから収集し、

前記対象畳み込み層が出力する複数の特徴マップと前記分類特徴ベクトルとに基づき前記対象行動のエッジ輪郭を特定する前記ステップは、

10

20

30

40

50

前記対象畳み込み層が出力する複数の特徴マップに対する前記分類特徴ベクトルの微分係数を計算し、重みマップを取得することと、

前記重みマップと前記対象畳み込み層が出力する複数の特徴マップを乗算し、複数種類の行動類型に対応する第1空間予測マップを取得することと、

前記第1空間予測マップに基づき、分類信用度が最も高い行動類型に対応する第1空間予測マップを抽出し、第2空間予測マップとすることと、

前記第2空間予測マップに基づき第3空間予測マップを生成し、前記第3空間予測マップのサイズは前記ビデオ画像フレームのサイズと同一であることと、

前記第3空間予測マップのエッジを抽出し、前記対象行動のエッジ輪郭を特定することと、を含む、

請求項8に記載の行動検出方法。

【請求項10】

前記対象行動のエッジ輪郭に基づき前記対象行動の空間位置を特定する前記ステップは、前記対象行動のエッジ輪郭を複数フレームの前記ビデオ画像フレームにおいて描くことを含む、

請求項9に記載の行動検出方法。

【請求項11】

前記対象行動のエッジ輪郭に基づき前記対象行動の空間位置を特定する前記ステップは、前記対象行動のエッジ輪郭を複数フレームの前記ビデオ画像フレームにおいて描く前記ステップの後に、前記対象行動のエッジ輪郭が描かれた複数フレームの前記ビデオ画像フレームをビデオ生成キャッシュ領域に取り込むことをさらに含む、

請求項10に記載の行動検出方法。

【請求項12】

複数フレームの前記ビデオ画像フレームデータを二次元畳み込みニューラルネットワークに入力し、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係と複数フレームの前記ビデオ画像フレームデータとに基づき、前記ビデオストリームにおける通行人の行動を認識した結果は、対象行動を認識したということ、または対象行動を認識していないということを含み、対象行動を認識していない時に、

エッジ輪郭が描かれたビデオ画像フレームがビデオ生成キャッシュ領域に記憶されているかどうか判断するステップと、

エッジ輪郭が描かれたビデオ画像フレームが前記ビデオ生成キャッシュ領域に記憶されている時に、前記ビデオ生成キャッシュ領域に記憶されている、エッジ輪郭が描かれたビデオ画像フレームに基づきビデオセグメントを生成するステップと、

前記ビデオ生成キャッシュ領域から前記ビデオセグメントを取り出すステップと、を含む、

請求項11に記載の行動検出方法。

【請求項13】

前記ビデオストリームから複数フレームのビデオ画像フレームデータを取得する前記ステップは、

前記ビデオストリーム内の目下のビデオ画像フレームにおける前景画像領域の面積を特定することと、

前記前景画像領域の面積が面積閾値を上回る時に、隣接する2つのビデオ画像フレームの運動量を特定することと、

隣接する2つのビデオ画像フレームの運動量が運動量閾値を上回る時に、目下のビデオ画像フレームをサンプリング開始点と決定することと、

所定時間長さの連続する複数フレームのビデオ画像フレームから、所定数の前記ビデオ画像フレームを均一サンプリングして前処理し、複数フレームの前記ビデオ画像フレームデータを取得することと、を含む、

請求項12に記載の行動検出方法。

【請求項14】

10

20

30

40

50

少なくとも1つのプロセッサと、

少なくとも1つのコンピュータプログラムが記憶され、前記少なくとも1つのコンピュータプログラムが前記少なくとも1つのプロセッサにより実行される時に、前記少なくとも1つのプロセッサに請求項1から1.3のいずれか一項に記載の行動検出方法を実現させるメモリと、

前記プロセッサと前記メモリとの間に接続され、前記プロセッサと前記メモリが情報のやり取りを実現するように配置された少なくとも1つのI/Oインターフェースと、を含む、

電子機器。

【請求項15】

コンピュータプログラムが記憶され、前記コンピュータプログラムがプロセッサにより実行される時に請求項1から1.3のいずれか一項に記載の行動検出方法を実現する、

コンピュータ読み取り可能な記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本願は2021年4月27日に提出された中国特許出願第202110459730.8号の優先権を主張し、当該中国特許出願の内容を参照によりここに援用する。

【0002】

本願は画像認識分野に関するものであり、特に行動検出方法、電子機器およびコンピュータ読み取り可能な記憶媒体に関するものである。

【背景技術】

【0003】

インテリジェントビデオ監視はコンピュータビジョン技術に基づき、ビデオデータをインテリジェント分析することができ、目下のところ、警備およびインテリジェント交通などの分野に広く適用されており、大幅に警備の反応速度向上、人的資源の節減を図っている。通行人はインテリジェントビデオ監視の重点フォロー対象であり、通行人の各種行動（例えば、異常行動など）の検出および認識は警備分野の重要なニーズの1つである。

【発明の概要】

【発明が解決しようとする課題】

【0004】

いくつかの関連技術において、インテリジェントビデオ監視はインテリジェントビデオ分析技術を使用して、膨大な量の監視ビデオから通行人の各種行動を検出して認識することで、公共安全应急管理に重要な参考を与え、公共安全突発事件の危害の低減に利することができる。しかし、通行人の行動を検出および認識する関連の技術は、いくつかの実際の適用シーンにおける配置のニーズを満たすことができない。

【課題を解決するための手段】

【0005】

第1態様において、本願実施例は、

ビデオストリームから複数フレームのビデオ画像フレームデータを取得するステップと、複数フレームの前記ビデオ画像フレームデータに基づき、前記ビデオストリームにおける通行人の行動を検出するステップと、を含み、

複数フレームの前記ビデオ画像フレームデータに基づき、前記ビデオストリームにおける通行人の行動を検出するステップは、

複数フレームの前記ビデオ画像フレームデータを二次元畳み込みニューラルネットワークに入力し、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係と複数フレームの前記ビデオ画像フレームデータとに基づき、前記ビデオストリームにおける通行人の行動を認識するステップを少なくとも含む、行動検出方法を提供する。

【0006】

第2態様において、本願実施例は、

10

20

30

40

50

1つまたは複数のプロセッサと、

1つまたは複数のコンピュータプログラムが記憶され、前記1つまたは複数のコンピュータプログラムが前記1つまたは複数のプロセッサにより実行される時に、前記1つまたは複数のプロセッサに第1態様において本願実施例が提供する前記行動検出方法を実現させるメモリと、

前記プロセッサと前記メモリとの間に接続され、前記プロセッサと前記メモリが情報のやり取りを実現するように配置された1つまたは複数のI/Oインターフェースと、を含む、電子機器を提供する。

【0007】

第3態様において、本願実施例は、

コンピュータプログラムが記憶され、前記コンピュータプログラムがプロセッサにより実行される時に、第1態様において本願実施例が提供する前記行動検出方法を実現する、コンピュータ読み取り可能な記憶媒体を提供する。

【図面の簡単な説明】

【0008】

【図1】図1は本願実施例における行動検出方法のフローチャートである。

【図2】図2は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図3】図3は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図4】図4は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図5】図5は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図6】図6は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図7】図7は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図8】図8は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図9】図9は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図10】図10は本願実施例の行動検出方法における一部のステップのフローチャートである。

【図11】図11は本願実施例における電子機器の組成概略ブロック図である。

【図12】図12は本願実施例におけるコンピュータ読み取り可能な媒体の組成概略ブロック図である。

【図13】図13は本願実施例における事例の行動検出装置およびシステム構成の概略図である。

【発明を実施するための形態】

【0009】

当業者が本願の技術案をよりよく理解できるように、以下に図面を組み合わせる本願が提供する行動検出方法、電子機器およびコンピュータ読み取り可能な記憶媒体について詳細に説明する。

【0010】

以下では図面を参照して例示的な実施例についてより十分に説明するが、前記例示的な実施例は異なる形式で具体化されてよく、本明細書にて説明された実施例に限定されると解釈すべきではない。むしろ、これらの実施例を提供するのは、本願を詳らかに十全なものにするのと同時に、当業者に本願の範囲を十分に理解させることを目的とする。

【0011】

10

20

30

40

50

矛盾することがなければ、本願の各実施例および実施例における各特徴は互いに組み合わせることができる。

【0012】

本明細書にて使用する、「および/または」という用語は1つまたは複数の関連列挙アイテムの任意のおよび全ての組み合わせを含む。

【0013】

本明細書にて使用する用語は特定の実施例を説明するためにのみ用いられ、本願を限定する意図はない。本明細書にて使用する単数形の「1つ」および「当該」という用語は、文脈が別途明らかに示さない限り複数形も含むことを意図する。さらに、本明細書にて「含む」および/または「...からなる」という用語を使用する場合は、特定の特徴、全体、ステップ、操作、素子および/またはコンポーネントが存在することを指すが、1つまたは複数のその他の特徴、全体、ステップ、操作、素子、コンポーネントおよび/またはそのグループが存在すること、または追加することを排除しない。

10

【0014】

別途限定しない限り、本明細書にて使用する全ての用語（技術および科学用語を含む）の意味は当業者が一般的に理解する意味と同じである。さらに、一般的な辞書において限定されたこれらの用語は、関連技術および本願の背景での意味と一致する意味を有すると解釈されるべきであり、本明細書にて明確に限定しない限り、理想的された、または過度の形式上の意味を有すると解釈されない。

【0015】

いくつかの関連技術では、主にフレーム間差分法を使用して通行人の行動を検出する。つまり、ビデオにおける連続する画像フレームのフレーム間差分の変化を分析することで異常行動領域を大まかに位置決めしてから、大まかな位置決め領域について行動認識を行い、異常行動が存在するかどうか、または異常行動の種類を特定する。フレーム間差分法に基づく通行人の異常行動検出および認識技術は光線の変化に敏感であり、かつ連続する画像フレームのフレーム間差分変化は全てが通行人の異常行動によって生じるものではなく、異常ではない多くの行動も画像フレームに大きな変化をもたらす可能性があり、ある種の異常行動は画像フレームの大きな変化をもたらすことはない。このほか、フレーム間差分法は連続した画像フレームを使用して異常行動について位置決めするため、大まかな位置決め領域について行動認識を行う際にも連続する画像フレームに基づいており、行動認識を行う際に使用する連続する画像フレームのフレーム数が少なければタイムドメイン情報の無駄を招き、行動認識を行う際に使用する連続する画像フレームのフレーム数が多ければ異常行動検出および認識の時間コストとリソース消費の増加を招くことになる。よって、当該フレーム間差分法は、通行人が少なく背景が単一のシーンにより適している。

20

30

【0016】

別のいくつかの関連技術では、各通行人の行動を分析することで異常行動が存在するかどうか特定し、例えば、通行人の検出または骨格点の分析により、各通行人の空間位置を特定して追跡し、各通行人の時間緯度での軌跡を取得してから、通行人個々の運動画像シーケンスまたは骨格点運動シーケンスをまとめて、行動認識を行い、異常行動の種類を特定する。通行人の検出または骨格点の分析をする際には、監視装置の画角に厳しい要求があり、監視装置がハイアングルであれば通行人の骨格点の特定において問題があり、監視装置がアイアングルであれば通行人同士で互いに遮られることによって通行人の空間位置を特定できないため、誤検出および検出漏れが生じる可能性がある。このほか、通行人の検出と骨格点検出はともに大量の計算リソースの消費を必要とし、処理速度も遅いため、異常行動検出および認識のリアルタイム分析要件を満たすことができない。

40

【0017】

いくつかの関連技術では三次元畳み込みまたはデュアルネットワークを使用して、異常行動の時系列情報を学習した上で異常行動の検出と認識を実現するが、三次元畳み込みおよびデュアルネットワークは短時間の時系列情報しか学習できず、例えば3×5×5の三

50

次元畳み込みの一回の演算は3フレームの画像しか関連付けることができず、デュアルネットワークにおけるオプティカルフロー計算は隣接するフレームの計算によって得られ、三次元畳み込みおよびデュアルネットワークはいずれも実行時に大量のリソースを消費する。いくつかの関連技術では二次元畳み込みを使用して特徴を重ねてから、三次元畳み込みを使用して時系列情報を融合して異常行動検出と認識を実現するが、三次元畳み込みを使用しているため、実行速度の向上には限りがある。いくつかの関連技術では、単一フレームのビデオ画像フレームを直接使用して、または複数フレームのビデオ画像フレームの特徴を重ねてから異常行動について分類を行っているが、フレームとフレームの間の関連関係を無視し、時系列情報を無駄にしており、検出と認識の精度は低下することになる。

【0018】

以上を踏まえ、関連の異常行動検出と認識技術は、いくつかの実際のシーンにおける配置のニーズを満たすのが困難である。

【0019】

この点に鑑み、第1態様において、図1を参照すれば分かるように、本願実施例は、ステップS100およびステップS200を含む行動検出方法を提供する。

【0020】

ステップS100では、ビデオストリームから複数フレームのビデオ画像フレームデータを取得する。

【0021】

ステップS200では、複数フレームの前記ビデオ画像フレームデータに基づき、前記ビデオストリームにおける通行人の行動を検出する。

【0022】

上記ステップS200は少なくともステップS300を含む。

【0023】

ステップS300では、複数フレームの前記ビデオ画像フレームデータを二次元畳み込みニューラルネットワークに入力し、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係と複数フレームの前記ビデオ画像フレームデータとに基づき、前記ビデオストリームにおける通行人の行動を認識する。

【0024】

いくつかの実施の形態において、ステップS100におけるビデオストリームは監視装置により取得したものである。本願実施例において、ビデオストリームは監視装置によりリアルタイムで取得してよく、監視装置により取得してからコンピュータ読み取り可能な媒体に記憶してもよい。本願実施例ではこれについて特に限定しない。なお、本願実施例において、毎回行われる行動検出はいずれも一定時間長さのビデオストリームに対応し、つまり、複数フレームのビデオ画像フレームデータは一定時間長さのビデオストリームから取得する。

【0025】

いくつかの実施の形態において、ステップS100における複数フレームのビデオ画像フレームデータは、本願実施例における二次元畳み込みニューラルネットワークで処理できるデータである。本願実施例では、ビデオストリームのデコード後に複数フレームのビデオ画像を取得することができ、各フレームのビデオ画像フレームデータはいずれも1フレームのビデオ画像フレームに対応する。

【0026】

通行人の行動には時間的な連続性があり、これに応じて、ビデオストリームのデコードで得られる複数フレームのビデオ画像フレームの間にも時系列上関連関係があり、複数フレームのビデオ画像フレームにそれぞれ対応する複数フレームのビデオ画像フレームデータの間にも時系列上関連関係がある。本願実施例にて使用する二次元畳み込みニューラルネットワークは、各フレームのビデオ画像フレームデータの特徴を学習できるだけでなく、複数フレームのビデオ画像フレームデータ間の時系列関連関係を学習することもできることから、複数フレームのビデオ画像フレームデータ間の時系列関連関係と複数フレーム

10

20

30

40

50

のビデオ画像フレームデータとに基づき通行人の行動検出を行うことができる。

【0027】

いくつかの実施の形態において、複数フレームのビデオ画像フレームに基づき前記ビデオストリームにおける通行人の行動を検出する前記ステップでは、ビデオストリームにおける通行人の行動についてのみ認識してよく、前記の通行人の行動についての認識は、通行人の行動が存在するかどうか特定すること、通行人の行動の種類を特定することなどを含むがこれらに限定されない。いくつかの実施の形態において、複数フレームのビデオ画像フレームに基づき前記ビデオストリームにおける通行人の行動を検出する前記ステップでは、ビデオストリームにおける通行人の行動を認識してから、通行人の行動の空間位置を検出してよい。本願実施例では通行人の行動の種類について特に限定しない。いくつかの実施の形態において、通行人の行動は例えば、転倒、殴り合いなどの異常行動を含んでよく、駆け足、飛び跳ねなどの正常行動を含んでもよい。

10

【0028】

本願実施例が提供する行動検出方法では、二次元畳み込みニューラルネットワークに基づき通行人の行動を検出し、使用する二次元畳み込みニューラルネットワークは、各フレームのビデオ画像フレームデータの特徴を学習できるだけでなく、複数フレームのビデオ画像フレームデータ間の時系列関連関係を学習することもできることから、複数フレームのビデオ画像フレームデータ間の時系列関連関係と複数フレームのビデオ画像フレームデータとに基づき、ビデオストリームにおける通行人の行動を認識することができる。三次元畳み込みまたはデュアルネットワークを使用して行動検出を行うものと比べて、本願実施例における、二次元畳み込みニューラルネットワークを使用した通行人の行動検出は、計算量が少なく、実行速度が速く、リソース消費が少なく、実際の配置におけるリアルタイム性を満たすことができる。単一フレームのビデオ画像フレームを直接使用するもの、または複数フレームのビデオ画像フレームの特徴を重ねてから行動検出を行うものと比べて、本願実施例における二次元畳み込みニューラルネットワークは通行人の行動の時系列特徴を学習していることから、誤検出、検出漏れを効果的に回避することができる。このほか、本願実施例では、通行人の行動の種類を直接認識すること、または通行人の行動の種類を認識してから、通行人の行動の空間位置を特定することができ、通行人の行動の領域を大まかに位置決めしてから、行動の種類を特定することによる適用シーンの制限を回避することができ、行動検出のシーン適応性を大幅に向上させている。

20

30

【0029】

本願実施例では二次元畳み込みニューラルネットワークの構造について特に限定しない。いくつかの実施の形態において、二次元畳み込みニューラルネットワークは、少なくとも1つの畳み込み層と少なくとも1つの全結合層とを含む。いくつかの実施の形態において、複数フレームのビデオ画像フレームデータはバッチ処理方式で二次元畳み込みニューラルネットワークに入力され、二次元畳み込みニューラルネットワークにおける畳み込み層は入力データについて特徴抽出を行うことができ、全結合層は畳み込み層により取得した特徴データに基づき今回の通行人の行動検出に対応するビデオストリームにおける通行人の行動を特定する。

40

【0030】

関連して、いくつかの実施の形態において、前記二次元畳み込みニューラルネットワークは、少なくとも1つの畳み込み層と少なくとも1つの全結合層とを含み、図2を参照すれば分かるように、ステップS300はステップS310およびS320を含む。

【0031】

ステップS310では、前記少なくとも1つの畳み込み層により、複数フレームの前記ビデオ画像フレームデータについて特徴抽出を行い、特徴データを取得し、前記特徴データは複数フレームの前記ビデオ画像フレームデータの時系列情報を融合し、前記時系列情報は複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係を特徴付けるものである。

50

【0032】

ステップS320では、前記少なくとも1つの全結合層により、前記特徴データに基づき前記ビデオストリームにおける通行人の行動を認識する。

【0033】

なお、いくつかの実施の形態において、前記特徴データは複数フレームのビデオ画像フレームデータの時系列情報を融合するというは、特徴データは各フレームのビデオ画像フレームデータの特徴を特徴付けることができ、複数フレームのビデオ画像フレームデータ間の時系列関連関係を特徴づけることもできることをいう。

【0034】

いくつかの実施の形態において、前記二次元畳み込みニューラルネットワークはプーリ 10
ング層をさらに含む。

【0035】

いくつかの実施の形態において、前記二次元畳み込みニューラルネットワークにおける 20
複数の畳み込み層は直列接続されており、各畳み込み層が入力データについて特徴抽出を行い、入力データに対応する特徴マップを取得する。複数フレームのビデオ画像フレームデータがバッチ処理方式で二次元畳み込みニューラルネットワークに入力される時に、各畳み込み層が入力データについて特徴抽出を行えば複数の特徴マップを取得できる。本願実施例において、バッチ処理方式で二次元畳み込みニューラルネットワークに入力される複数フレームのビデオ画像フレームデータは時系列順に並べられ、各畳み込み層が取得する複数の特徴マップも時系列順に並べられる。いくつかの実施の形態では、各畳み込み層後に、複数の特徴マップの一部の特徴チャンネルを交換することで、異なる時系列情報の相互融合を実現し、最終的に特徴データにおいて複数フレームの前記ビデオ画像フレームデータの時系列データの融合を実現する。

【0036】

関連して、いくつかの実施の形態において、前記二次元畳み込みニューラルネットワークは複数の直列接続された前記畳み込み層を含み、図3を参照すれば分かるように、ステップS310はステップS311およびステップS312を含む。

【0037】

ステップS311では、各前記畳み込み層について、前記畳み込み層の入力データを前記畳み込み層に入力して特徴抽出を行い、複数の特徴マップを取得し、複数の前記特徴マップは複数フレームの前記ビデオ画像フレームデータに一对一に対応し、各前記特徴マップは複数の特徴チャンネルを含む。 30

【0038】

ステップS312では、複数の前記特徴マップの一部の特徴チャンネルを交換して第1データを取得する。

【0039】

いくつかの実施の形態では、前記畳み込み層が最初の畳み込み層である時に、前記畳み込み層の入力データは複数フレームの前記ビデオ画像フレームデータである。

【0040】

いくつかの実施の形態では、前記畳み込み層が最後の畳み込み層ではなく、かつ最初の畳み込み層でもない時に、前記第1データを次の畳み込み層の入力データとする。 40

【0041】

いくつかの実施の形態では、前記畳み込み層が最後の畳み込み層である時に、前記第1データを前記特徴データとする。

【0042】

なお、いくつかの実施の形態において、複数の特徴マップの一部の特徴チャンネルの交換はデータの移動を行うだけで、加算乗算操作はないため、時系列情報のやり取りを行う際に計算量は増加せず、データ移動の速度が速く、行動検出の実行効率に影響しない。

【0043】

関連して、いくつかの実施の形態において、複数フレームの前記ビデオ画像フレームデ 50

ータは順に並んだビデオ画像フレームデータをNフレーム含み、複数の前記特徴マップは順に並んだ特徴マップをN個含み、図4を参照すれば分かるように、ステップS312はステップS3121からステップS3123を含む。

【0044】

ステップS3121では、各前記特徴マップにおける複数の特徴チャンネルを、順に並んだN組の特徴チャンネルに分ける。

【0045】

ステップS3122では、N個の順に並んだ特徴マップにおけるi番目の特徴マップについて、i番目の特徴マップに対応するj番目の特徴マップを特定し、i番目の特徴マップはN個の前記順に並んだ特徴マップにおけるいずれか1つであり、j番目の特徴マップはN個の前記順に並んだ特徴マップにおけるいずれか1つである。

10

【0046】

ステップS3123では、i番目の特徴マップにおける第i組の特徴チャンネルをj番目の特徴マップにおけるいずれか1組の特徴チャンネルと交換し、前記第1データを取得し、N、i、jは正の整数である。

【0047】

本願実施例ではステップS3122においてi番目の特徴マップに対応するj番目の特徴マップを特定することを如何に実行するかについて特に限定しない。いくつかの実施の形態では、iとjの代数関係に基づき、i番目の特徴マップに対応するj番目の特徴マップを特定し、いくつかの実施の形態では、iとjの隣接関係に基づき、i番目の特徴マップに対応するj番目の特徴マップを特定し、いくつかの実施の形態では、N個の前記順に並んだ特徴マップから特徴マップを1つランダムに指定し、i番目の特徴マップに対応するj番目の特徴マップとする。

20

【0048】

いくつかの実施の形態では、二次元畳み込みニューラルネットワークの全結合層により分類特徴ベクトルを取得し、分類特徴ベクトルの各要素が一種類の行動類型の分類確率を特徴付ける。各種行動類型の分類確率に基づき、今回の歩行者の行動検出に対応するビデオストリームから歩行者の行動を特定することができる。今回の歩行者の行動検出に対応するビデオストリームから歩行者の行動を特定することは、今回の歩行者の行動検出に対応するビデオストリームに、検出しようとする対象行動があるかどうか判断することと、存在する対象行動の類型を判断することと、を含むが、これらに限定されない。

30

【0049】

関連して、いくつかの実施の形態では、図5を参照すれば分かるように、ステップS320はステップS321からステップS323を含む。

【0050】

ステップS321では、前記少なくとも1つの全結合層により、前記特徴データに基づき分類特徴ベクトルを取得し、前記分類特徴ベクトルの各要素が一種類の行動類型に対応する。

【0051】

ステップS322では、前記分類特徴ベクトルに基づき、各種行動類型の分類確率を特定する。

40

【0052】

ステップS323では、前記各種行動類型の分類確率に基づき前記ビデオストリームにおける歩行者の行動を認識する。

【0053】

いくつかの実施の形態では、ステップS322において、ステップS321で取得した分類特徴ベクトルを分類器に入力し、各種行動類型の分類確率を取得する。

【0054】

関連して、いくつかの実施の形態では、図6を参照すれば分かるように、ステップS323はステップS3231からステップS3234を含む。

50

【 0 0 5 5 】

ステップ S 3 2 3 1 では、各種行動類型の分類確率がフィルタ閾値を上回るかどうか判断する。

【 0 0 5 6 】

ステップ S 3 2 3 2 では、少なくとも一種類の行動類型の分類確率が前記フィルタ閾値を上回った時に、対象行動を認識したと判断する。

【 0 0 5 7 】

ステップ S 3 2 3 3 では、分類確率が前記フィルタ閾値を上回る行動類型を前記対象行動の類型であると特定する。

【 0 0 5 8 】

ステップ S 3 2 3 4 では、各種行動類型の分類確率がいずれも前記フィルタ閾値を上回らない時に、対象行動を認識していないと判断する。

【 0 0 5 9 】

二次元畳み込みニューラルネットワークの畳み込み層は空間不変性を有し、つまり、畳み込み層を介して取得した特徴マップと初期画像との間には空間対応関係があり、トレーニングを経て取得した二次元畳み込みニューラルネットワークの畳み込み層は、特徴マップにおける分類に関連する領域特徴値を大きくし、分類に関係のない領域特徴値を小さくすることができる。いくつかの実施の形態では、対象行動を認識した時に、二次元畳み込みニューラルネットワークの畳み込み層が出力する特徴マップについて輪郭分析を行うことで、対象行動に関連する領域のエッジ輪郭を特定して、認識した対象行動の空間位置を特定することができる。いくつかの実施の形態では、二次元畳み込みニューラルネットワークの少なくとも1つの畳み込み層の中から畳み込み層を1つ指定して対象畳み込み層とし、対象畳み込み層が出力する特徴マップと、二次元畳み込みニューラルネットワークの全結合層が取得した分類特徴ベクトルとに基づき輪郭分析を行う。

【 0 0 6 0 】

いくつかの実施の形態では、図7を参照すれば分かるように、ステップ S 3 0 0 の後に、ステップ S 2 0 0 はステップ S 4 0 0 をさらに含む。

【 0 0 6 1 】

ステップ S 4 0 0 では、前記二次元畳み込みニューラルネットワークの出力データに基づき、前記ビデオストリームにおける通行人の行動の空間位置を検出する。

【 0 0 6 2 】

いくつかの実施の形態において、ステップ S 4 0 0 は、ステップ S 3 0 0 により認識した通行人の行動の状況下で実行してよい。いくつかの実施の形態では、毎回ステップ S 3 0 0 を実行した後にいずれもステップ S 4 0 0 を実行し、つまり、通行人の行動を認識したかどうかに関わらずステップ S 4 0 0 を実行する。いくつかの実施の形態では、異なる適用シーンに対して、ステップ S 4 0 0 の実行またはスキップを選択してよい。

【 0 0 6 3 】

関連して、いくつかの実施の形態では、前記二次元畳み込みニューラルネットワークは、少なくとも1つの畳み込み層と少なくとも1つの全結合層とを含み、前記二次元畳み込みニューラルネットワークの出力データは、前記少なくとも1つの全結合層により取得した分類特徴ベクトルと、対象畳み込み層が出力した複数の特徴マップとを含み、前記対象畳み込み層は前記少なくとも1つの畳み込み層における1つであり、前記分類特徴ベクトルの各要素が一種類の行動類型に対応し、図8を参照すれば分かるように、ステップ S 4 0 0 はステップ S 4 1 0 を含む。

【 0 0 6 4 】

ステップ S 4 1 0 では、前記対象畳み込み層が出力する複数の特徴マップと前記分類特徴ベクトルとに基づき対象行動の空間位置を特定する。

【 0 0 6 5 】

関連して、いくつかの実施の形態では、図9を参照すれば分かるように、ステップ S 4 1 0 はステップ S 4 1 1 およびステップ S 4 1 2 を含む。

10

20

30

40

50

ステップS 4 1 1では、前記対象畳み込み層が出力する複数の特徴マップと前記分類特徴ベクトルとに基づき、前記対象行動のエッジ輪郭を特定する。

【0066】

ステップS 4 1 2では、前記対象行動のエッジ輪郭に基づき前記対象行動の空間位置を特定する。

【0067】

関連して、いくつかの実施の形態において、上記ステップS 4 1 2は、前記対象畳み込み層が出力する複数の特徴マップに対する前記分類特徴ベクトルの微分係数を計算し、重みマップを取得することと、前記重みマップと前記対象畳み込み層が出力する複数の特徴マップを乗算し、複数種類の行動類型に対応する第1空間予測マップを取得することと、前記第1空間予測マップに基づき、分類信用度が最も高い行動類型に対応する第1空間予測マップを抽出し、第2空間予測マップとすることと、前記第2空間予測マップに基づき第3空間予測マップを生成し、前記第3空間予測マップのサイズは前記ビデオ画像フレームのサイズと同一であることと、前記第3空間予測マップのエッジを抽出し、前記対象行動のエッジ輪郭を特定することと、を含んでよい。

10

【0068】

関連して、いくつかの実施の形態において、複数フレームの前記ビデオ画像フレームデータは所定時間長さの複数フレームのビデオ画像フレームから収集し、上記ステップS 4 1 2は、前記対象畳み込み層が出力する複数の特徴マップに対する前記分類特徴ベクトルの微分係数を計算して、ゼロ未満の微分係数値をゼロとし、重みマップを取得することと、前記重みマップと前記対象畳み込み層が出力する複数の特徴マップを乗算して、ゼロ未満の積をゼロとし、複数種類の行動類型に対応する第1空間予測マップを取得することと、前記第1空間予測マップに基づき、分類信用度が最も高い行動類型に対応する第1空間予測マップを抽出し、第2空間予測マップとすることと、前記第2空間予測マップを正規化処理することと、正規化処理後の第2空間予測マップのサイズを前記ビデオ画像フレームのサイズにスケールリングして二値化処理を行い、第3空間予測マップを取得することと、前記第3空間予測マップについてエッジ抽出を行い、前記対象行動のエッジ輪郭を特定することと、を含んでよい。

20

【0069】

いくつかの実施の形態において、上記ステップS 4 1 2は、前記対象行動のエッジ輪郭を複数フレームの前記ビデオ画像フレームにおいて描くことを含む。

30

【0070】

いくつかの実施の形態では、対象行動を認識して対象行動の空間位置を特定した後に、対象行動の輪郭が描かれたビデオ画像フレームをビデオ生成キャッシュ領域に書き込み、さらにビデオファイルを生成してファイルシステムに記憶する。

【0071】

関連して、いくつかの実施の形態において、上記ステップS 4 1 2は、前記対象行動のエッジ輪郭を複数フレームの前記ビデオ画像フレームにおいて描いた後に、前記対象行動のエッジ輪郭が描かれた複数フレームの前記ビデオ画像フレームをビデオ生成キャッシュ領域に取り込むことをさらに含む。

40

【0072】

いくつかの実施の形態において、毎回の通行人の行動検出は一定時間長さのビデオストリームに対応し、ステップS 100により一定時間長さのビデオストリームから複数フレームのビデオ画像フレームデータを取得する。本願実施例では当該一定時間長さについて特に限定しない。いくつかの実施の形態では、検出を要する通行人の行動が一般的に発生する時間長さに基づき当該一定時間長さを決定する。例えば、殴り合いは2秒、転倒は1秒とするなどである。

【0073】

いくつかの実施の形態において、通行人の行動が持続する時間長さは上述の一定時間長さよりも長いことがあり得る。いくつかの実施の形態では、毎回の検出で対象行動を認識

50

すると、対象行動の輪郭が描かれたビデオ画像フレームをビデオ生成キャッシュ領域に書き込み、ある回の検出で対象行動が認識されなくなり、対象行動が終了したことを表すまで続け、このような場合、ビデオ生成キャッシュ領域におけるビデオ画像フレームはビデオセグメントに変換され、対象行動の開始から終了までの全過程が記録されたビデオセグメントを取得することができ、当該ビデオセグメントに基づき、対象行動の開始時間、終了時間、持続時間などの情報を特定することもできる。関連して、図10を参照すれば分かるように、いくつかの実施の形態では、複数フレームの前記ビデオ画像フレームデータを二次元畳み込みニューラルネットワークに入力し、複数フレームの前記ビデオ画像フレームデータ間の時系列関連関係と複数フレームの前記ビデオ画像フレームデータとに基づき、前記ビデオストリームにおける通行人の行動を認識した結果は、対象行動を認識した

10

【0074】

ステップS501では、エッジ輪郭が描かれたビデオ画像フレームがビデオ生成キャッシュ領域に記憶されているかどうか判断する。

【0075】

ステップS502では、エッジ輪郭が描かれたビデオ画像フレームが前記ビデオキャッシュ領域に記憶されている時に、前記ビデオキャッシュ領域に記憶されている、エッジ輪郭が描かれたビデオ画像フレームに基づきビデオセグメントを生成する。

【0076】

ステップS503では、前記ビデオ生成キャッシュ領域から前記ビデオセグメントを取り出す。

20

【0077】

いくつかの実施の形態において、ビデオ生成キャッシュ領域から前記ビデオセグメントを取り出す前記ステップは、ビデオ生成キャッシュ領域におけるビデオセグメントをファイルシステムに記憶することと、ビデオ生成キャッシュ領域を空にすることと、を含む。

【0078】

いくつかの実施の形態において、上記ステップS100は、前記ビデオストリームを取得することと、前記ビデオストリームをデコードして、連続する複数のビデオ画像フレームを取得することと、連続する複数のビデオ画像フレームについてサンプリングを行い、複数の検出対象ビデオ画像フレームを取得することと、複数の前記検出対象ビデオ画像フレームについて前処理を行い、複数フレームの前記ビデオ画像フレームデータを取得することと、を含む。

30

【0079】

なお、本願実施例ではビデオストリームを如何にデコードするかについて特に限定しない。いくつかの実施の形態では、グラフィックスプロセッサ(GPU、Graphics Processing Unit)を使用してビデオストリームをデコードする。

【0080】

なお、本願実施例では連続する複数のビデオ画像フレームを如何にサンプリングするかについて特に限定しない。いくつかの実施の形態では、複数のビデオ画像フレームにおいてランダムにサンプリングを行う。いくつかの実施の形態では、所定の間隔で複数のビデオ画像フレームにおいてサンプリングを行う。いくつかの実施の形態では、複数のビデオ画像フレームにおいて連続してサンプリングを行う。

40

【0081】

なお、前記の所定の間隔で複数のビデオ画像フレームにおいてサンプリングを行うことは、前記の複数のビデオ画像フレームにおいて連続してサンプリングを行うことと比べて、より多くの時系列情報を取得することができるため、検出精度が向上する。

【0082】

本願実施例では複数の検出対象のビデオ画像フレームを如何に前処理するかについて特に限定しない。いくつかの実施の形態において、複数の検出対象のビデオ画像フレームに

50

対して前処理を行うことは、各検出対象のビデオ画像フレームのサイズを所定のサイズに調整することと、所定のサイズに調整された検出対象ビデオ画像フレームについて色空間変換処理、画素値正規化処理、平均値を引いて標準偏差で除す処理を行い、複数フレームのビデオ画像フレームデータを取得することと、を含む。

【0083】

いくつかの実施の形態において、毎回の通行人の行動検出は一定時間長さのビデオストリームに対応し、均一サンプリング方式で一定時間長さのビデオ画像フレームにおいて所定のフレーム数のビデオ画像フレームデータを取得する。

【0084】

関連して、いくつかの実施の形態において、上記ステップS100は、前記ビデオストリーム内の目下のビデオ画像フレームにおける前景画像領域の面積を特定することと、前記前景画像領域の面積が面積閾値を上回る時に、隣接する2つのビデオ画像フレームの運動量を特定することと、隣接する2つのビデオ画像フレームの運動量が運動量閾値を上回る時に、目下のビデオ画像フレームをサンプリング開始点と決定することと、所定時間長さの連続する複数フレームのビデオ画像フレームから、所定数の前記ビデオ画像フレームを均一サンプリングして前処理し、複数フレームの前記ビデオ画像フレームデータを取得することと、を含む。

10

【0085】

本願実施例では目下のビデオ画像フレームにおける前景画像領域の面積を如何に特定するかについて特に限定しない。いくつかの実施の形態では、フレーム間差分法を使用して目下のビデオ画像フレームの前景画像を取得する。

20

【0086】

本願実施例では隣接する2つのビデオ画像フレームの運動量を如何に特定するかについて特に限定しない。いくつかの実施の形態では、まばらなオプティカルフローを使用して隣接する2つのビデオ画像フレームの運動量を計算する。

【0087】

いくつかの実施の形態では、ステップS100からステップS200により行動検出を行う前に、二次元畳み込みニューラルネットワークをトレーニングするステップをさらに含み、ビデオストリームを取得し、ビデオストリームをデコードしてビデオ画像フレームを生成し、データクレンジングを行ってサンプルビデオセグメントを取得し、サンプルビデオセグメントにおける通行人の行動の種類をマークし、検出を要する通行人の行動がないサンプルビデオセグメントを背景としてマークし、マークされたサンプルビデオセグメントを使用して二次元畳み込みニューラルネットワークをトレーニングし、トレーニングされた二次元畳み込みニューラルネットワークについて定量化操作を行い、フォーマット変換する。

30

【0088】

第2態様において、図11を参照すれば分かるように、本願実施例は、

1つまたは複数のプロセッサ101と、

1つまたは複数のコンピュータプログラムが記憶され、前記1つまたは複数のコンピュータプログラムが1つまたは複数のプロセッサ101により実行される時に、1つまたは複数のプロセッサ101に第1態様において本願実施例が提供する前記行動検出方法を実現させるメモリ102と、

40

プロセッサ101とメモリ102との間に接続され、プロセッサ101とメモリ102が情報のやり取りを実現するように配置された1つまたは複数のI/Oインターフェース103と、を含む、電子機器を提供する。

【0089】

プロセッサ101はデータ処理機能を有するデバイスであり、中央プロセッサ(CPU)などを含むがこれらに限定されず、メモリ102はデータ記憶機能を有するデバイスであり、ランダムアクセスメモリ(RAM。より具体的にはSDRAM、DDRなど)、読み取り専用メモリ(ROM)、電氣的に消去、プログラムが可能な読み取り専用メモリ(

50

E E P R O M)、フラッシュメモリ (F L A S H) を含むがこれらに限定されず、I / O インターフェース (読み書きインターフェース) 1 0 3 はプロセッサ 1 0 1 とメモリ 1 0 2 との間に接続され、プロセッサ 1 0 1 とメモリ 1 0 2 の情報のやり取りを実現することができ、データバス (B u s) などを含むがこれらに限定されない。

【 0 0 9 0 】

いくつかの実施の形態において、プロセッサ 1 0 1、メモリ 1 0 2、I / O インターフェース 1 0 3 はバス 1 0 4 により互いに接続されて、さらにはコンピュータ装置のその他のコンポーネントに接続される。

【 0 0 9 1 】

第 3 態様において、図 1 2 を参照すれば分かるように、本願実施例は、コンピュータプログラムが記憶され、前記コンピュータプログラムがプロセッサにより実行される時に、第 1 態様において本願実施例が提供する前記行動検出方法を実現する、コンピュータ読み取り可能な記憶媒体を提供する。

10

【 0 0 9 2 】

本願実施例が提供する技術案を当業者がより明確に理解できるように、以下では具体的な実例を通じて本願実施例が提供する技術案について詳細に説明する。

【 0 0 9 3 】

実例 1

図 1 3 は本願実施例における実例の行動検出装置およびシステム構成の概略図である。

【 0 0 9 4 】

20

図 1 3 に示すように、前記行動検出装置は行動認識モジュールと、行動位置検出モジュールと、ビデオ自動記憶モジュールと、を含む。本実例 1 において、前記行動検出装置はサーバに配置され、サーバは G P U と、C P U と、ネットワークインターフェースと、ビデオメモリと、内部メモリと、をさらに含み、内部バスにより互いに接続される。

【 0 0 9 5 】

実例 2

本実例 2 において行動検出を行う際のフローは以下の通りである。

【 0 0 9 6 】

二次元畳み込みニューラルネットワークをビデオメモリまたは内部メモリにロードして初期化し、入力画像サイズ制限、バッチ処理の大きさ、フィルタ閾値、面積閾値、行動の一般的な発生時間長さ、フレームレートなどのパラメータを配置する。

30

【 0 0 9 7 】

カメラからビデオストリームを取得して、取得したビデオストリームをシステムサーバの G P U に送りハードウェアデコードし、複数フレームのビデオ画像フレームを生成する。

【 0 0 9 8 】

複数フレームのビデオ画像フレームから行動の一般的な発生時間長さに対応するビデオ画像フレームセットを抽出し、そこから N フレームを均一サンプリングし、N は二次元畳み込みニューラルネットワークのトレーニング時に決定される。

【 0 0 9 9 】

均一サンプリングで取得した N フレームのビデオ画像フレームを前処理し、その長辺を二次元畳み込みニューラルネットワーク指定のサイズに固定し、等比例スケールを行い、短辺塗りつぶし画素値を 0 としてから R G B 色空間に変換し、最後に画素値を 0 から 1 の間に正規化して、平均値を引いて標準偏差で除し、前処理後の N フレームのビデオ画像フレームデータを取得する。

40

【 0 1 0 0 】

前処理を経て取得した N フレームのビデオ画像フレームを、N の二次元データセットとしてバッチ処理され二次元畳み込みニューラルネットワークに送られるものとし、二次元畳み込みニューラルネットワークの各畳み込み層は全てサイズが N C H W であるとの特徴を得られ、N フレームのサイズを C H W の特徴マップとして特徴付け、C は特徴チャンネル数を表し、H および W は特徴マップの幅と高さをそれぞれ表す。

50

【 0 1 0 1 】

各畳み込み層が出力するNフレームの特徴マップの特徴チャンネルをN組に分け、i番目の特徴マップについて、i番目の特徴マップ以外のその他の特徴マップから特徴マップを1つランダムに抽出し、j番目の特徴マップと記し、j番目の特徴マップのN組の特徴チャンネルにおいて1組の特徴チャンネルをランダムに選択し、これをi番目の特徴マップのi組の特徴チャンネルと交換することで、計算量を別途増やさずに時系列情報のやり取りを行う。

【 0 1 0 2 】

特徴チャンネルの交換を経た特徴マップを次の畳み込み層に送り計算し、最後の畳み込み層後に特徴データを取得し、全結合層により分類特徴ベクトルを生成する。分類特徴ベクトルを分類器に送り、各種行動タイプの分類確率を取得する。

10

【 0 1 0 3 】

少なくとも一種類の行動タイプの分類確率がフィルタ閾値を上回った時に、対象行動を認識したと判断する。

【 0 1 0 4 】

分類確率がフィルタ閾値を上回る行動タイプを対象行動の種類であると特定する。

【 0 1 0 5 】

各種行動タイプの分類確率がいずれもフィルタ閾値を上回らない時に、対象行動を認識していないと判断する。

【 0 1 0 6 】

対象行動を認識した時に、対象行動の空間位置を特定するステップは、

20

二次元畳み込みニューラルネットワークにおける対象畳み込み層の特徴マップと、全結合層が出力する分類特徴ベクトルを抽出し、対象畳み込み層の特徴マップに対する分類特徴ベクトルの微分係数を計算し、微分係数の値が0未満の微分係数の値を0とし、対象畳み込み特徴マップと空間の大きさが一致する重みマップを取得することと、

重みマップと対象畳み込み層の特徴マップを乗算して、積における0未満の値を0として第1空間予測マップを取得することと、

第1空間予測マップをN*Class*H*Wとの記憶形式に変換し、Classは行動類型数を表し、その後、類型緯度において第1空間予測マップについてsoftmax操作を行い、分類結果の信用度が最も高い類型に対応する緯度を抽出し、N*H*Wの空間予測マップを取得し、特徴マップにおける全ての要素の、バッチ処理数Nの所在緯度における最大値を計算し、H*W第2空間予測マップを取得することと、

30

第2空間予測マップにおける全ての要素から全ての要素の最小値を引き、さらに全ての要素の最大値で除して、第2空間予測マップを0から1の間に正規化することと、

正規化後の第2空間予測マップをビデオ画像フレームのサイズにスケールングして二値化処理を行い(要素値が0.5を上回れば1、そうでなければ0とする)、第3空間予測マップを取得することと、

第3空間予測マップについてエッジ輪郭抽出を行い、取得したエッジ輪郭は対象行動が存在する空間位置であることと、

輪郭境界を検出結果としてビデオ画像フレームにおいて描いて、ビデオ画像フレームをビデオ生成キャッシュ領域に取り込み、次の検出を開始することと、を含む。

40

ビデオセグメントを記憶するステップは、

対象行動を認識していない時に、ビデオ生成キャッシュ領域にビデオ画像フレームがあるかどうか判断する(つまり、前回の検出が対象行動を認識したかどうか判断する)ことと、

ビデオキャッシュ領域にビデオ画像フレームがある(つまり、前回の検出が対象行動を認識した)時に、ビデオキャッシュ領域におけるビデオ画像フレームに基づきビデオセグメントを生成することと、

ビデオセグメントを記憶することと、

ビデオ生成キャッシュ領域を空にして次の検出を開始すること、または、

ビデオキャッシュ領域にビデオ画像フレームがない(つまり、前回の検出で対象行動を

50

認識していない)時に、次の検出を開始することと、を含む。

【0107】

本明細書にて開示した方法のうち全てまたはいくつかのステップ、装置における機能モジュール/ユニットはソフトウェア、ファームウェア、ハードウェアおよびその適切な組み合わせとして実施されてよいと当業者は理解できる。ハードウェアの実施の形態において、以上の説明で言及した機能モジュール/ユニットの間の区分は必ずしも物理的なコンポーネントの区分に対応せず、例えば、1つの物理的コンポーネントは複数の機能を有してよく、あるいは1つの機能またはステップは複数の物理的コンポーネントが連携して実行することができる。ある物理的コンポーネントまたは全ての物理的コンポーネントは、(中央処理装置、デジタル信号プロセッサまたはマイクロプロセッサのような)プロセッサにより実行されるソフトウェアとして実施されてよく、またはハードウェアとして、あるいは専用集積回路のような集積回路として実施されてもよい。このようなソフトウェアはコンピュータ読み取り可能な媒体に配置することができ、コンピュータ読み取り可能な媒体はコンピュータ記憶媒体(または非一時的な媒体)および通信媒体(または一時的な媒体)を含んでよい。当業者に知られているように、コンピュータ記憶媒体という技術用語は(コンピュータ読み取り可能な命令、データ構造、プログラムモジュールまたはその他のデータのような)情報を記憶するための任意の方法または技術において実施される揮発性および不揮発性、リムーバブルなおよび非リムーバブルな媒体を含む。コンピュータ記憶媒体はRAM、ROM、EEPROM、フラッシュメモリまたはその他のメモリ技術、CD-ROM、デジタル多機能ディスク(DVD)またはその他の光ディスクメモリ、磁気ボックス、磁気テープ、磁気ディスクメモリまたはその他の磁気記憶装置、または所望の情報を記憶しかつコンピュータによりアクセス可能な任意のその他の媒体を含むがこれらに限定されない。このほか、当業者であれば、通信媒体は一般的にコンピュータ読み取り可能な命令、データ構造、プログラムモジュールまたは搬送波あるいはその他の伝送機構のような変調データ信号におけるその他のデータを含み、また任意の情報伝送媒体を含むことができるということは公知の事項である。

10

20

【0108】

本明細書では例示的な実施例を開示し、かつ具体的な用語を使用しているが、これらは単に一般的な例示的な意味としてのみ使用され、またそのように解釈されるべきであり、制限の目的に用いられない。いくつかの実例において、別途明示しない限り、特定の実施例と組み合わせて説明された特徴、特性および/または要素を単独で使用することができ、またはその他の実施例と組み合わせて説明された特徴、特性および/または要素と組み合わせて使用することができるということは当業者にとって自明である。したがって、添付の請求項により説明された本願の範囲から逸脱しなければ、様々な形式および詳細において変更を加えることができると当業者は理解できる。

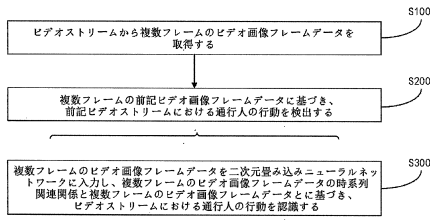
30

40

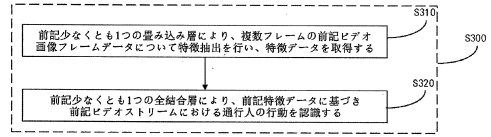
50

【図面】

【図1】

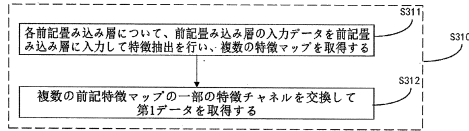


【図2】

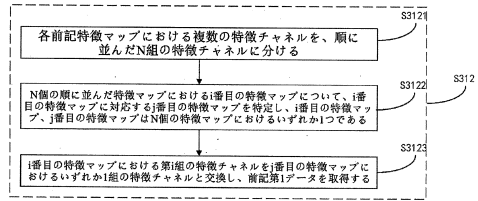


10

【図3】

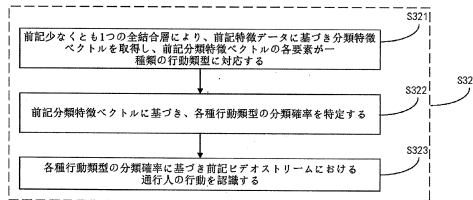


【図4】

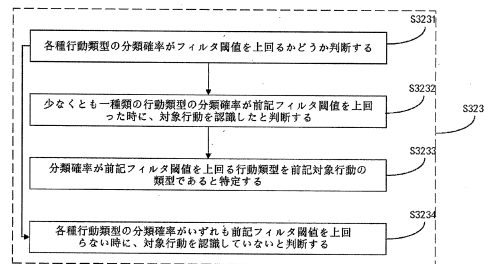


20

【図5】



【図6】

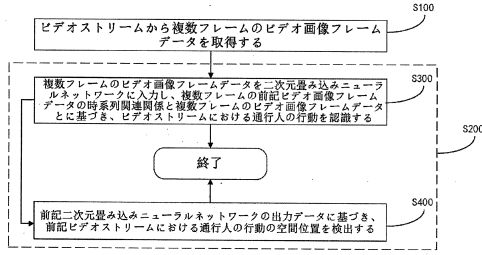


30

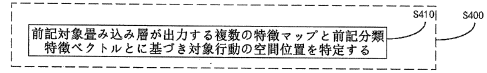
40

50

【図 7】

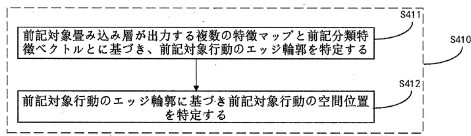


【図 8】

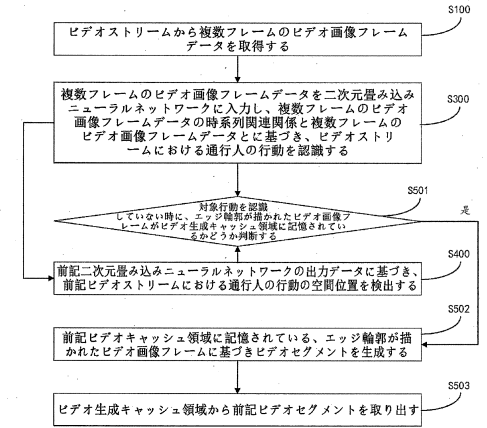


10

【図 9】

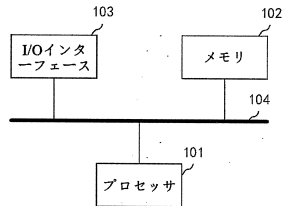


【図 10】

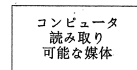


20

【図 11】



【図 12】

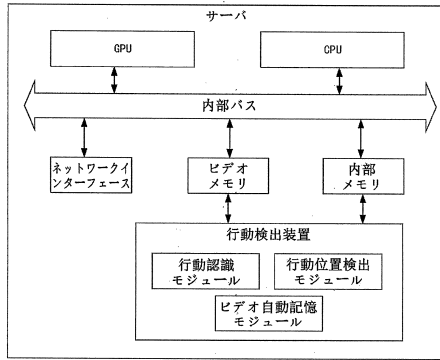


30

40

50

【図 13】



10

20

30

40

50

フロントページの続き

- 中国広東省深 せん 市南山区高新技术産業園科技南路中興通訊大厦
(72)発明者 賈霞
- 中国広東省深 せん 市南山区高新技术産業園科技南路中興通訊大厦
(72)発明者 劉明
- 中国広東省深 せん 市南山区高新技术産業園科技南路中興通訊大厦
(72)発明者 張羽豐
- 中国広東省深 せん 市南山区高新技术産業園科技南路中興通訊大厦
(72)発明者 林巍 ヨ
- 中国広東省深 せん 市南山区高新技术産業園科技南路中興通訊大厦
審査官 村山 絢子
- (56)参考文献 Ji Lin, Chuang Gan, Song Han , TSM: Temporal Shift Module for Efficient Video Understand
ing , 2019 IEEE/CVF International Conference on Computer Vision (ICCV) , 米国 , IEEE , 2
019年10月27日 , pp. 7082-7092 , <https://ieeexplore.ieee.org/document/9008827/>
- (58)調査した分野 (Int.Cl. , D B名)
- G 0 6 T 7 / 0 0 - 7 / 9 0
G 0 6 V 1 0 / 0 0 - 2 0 / 9 0
H 0 4 N 7 / 1 8