

令和 3 年 6 月 16 日現在

機関番号：17102

研究種目：挑戦的研究（萌芽）

研究期間：2017～2020

課題番号：17K20008

研究課題名（和文）異なる大規模データベースで学習された深層畳込ネットを融合した特徴抽出

研究課題名（英文）Extracting Fusion Features of Deep Convolutional Neural Networks Trained on Different Large-scale Datasets

研究代表者

松川 徹（Tetsu, Matsukawa）

九州大学・システム情報科学研究所・助教

研究者番号：80747212

交付決定額（研究期間全体）：（直接経費） 4,800,000円

研究成果の概要（和文）：複数の学習済みCNNから得られる特徴を融合し、異なる画像認識タスクに転用する手法を開発した。まず、全結合層の重みの正規直交化により、CNNを学習したデータベースとは異なるタスクに特徴を転用した場合に性能が向上することを示した。この正規直交化した全結合層を画素へ局在化した特徴マップを2つのCNNから取得し、双線形プーリングにより融合を行った。この融合特徴は、最上位畳込層の特徴チャンネルを融合するよりも少ない計算量で得られるが、それと同等以上の認識性能を示した。また、対象タスクの学習データを用いて、畳み込み層の特徴を判別的に集積する方法を開発した。

研究成果の学術的意義や社会的意義  
クラス代表ベクトルの正規直交化は、対象タスクの学習データを利用せず学習済みのモデルのみから実行可能であるが、全結合層の特徴をそのまま利用した場合よりも対象タスクでの認識性能が高く、最上位畳み込み層よりも出力される次元数が少ない特徴抽出を可能とする。これを特徴融合に利用することで、計算量が削減された高性能な融合特徴を抽出できる。また、判別の特徴集積は、小サンプルデータにおいて高速に学習でき、CNNをファインチューニングするよりも有効であった。このように本研究成果は、学習済みの巨大な深層学習のモデルを、限られた計算資源で小サンプルデータ問題に適用する場合に有用である。

研究成果の概要（英文）：We have developed feature extraction methods from multiple pre-trained CNNs for transferring the learned representation to various image recognition tasks. First, we have shown that the orthonormalization of the weights of a fully connected layer improves the performance when transferring them to various tasks. We have extracted the localized feature maps of the orthonormalized fully connect layer from two CNNs, and fused them with bilinear pooling. The proposed fusion features achieved comparable or better performance compared with the fusion features of the top convolutional layers with a smaller computational cost. We have also developed a method, which conducts discriminative pooling of the convolutional features with training data of a target task.

研究分野：パターン認識

キーワード：畳込ネット 学習済みモデル 特徴転移 正規直交化 融合 双線形プーリング 画像識別 カメラ間人物照合

## 1. 研究開始当初の背景

近年、深層畳込ニューラルネットワーク(CNN)が画像認識分野で注目を集めている。大量の学習データにより学習された CNN は、画像中に写る物体のカテゴリ名称の認識を極めて正確に行え、画像から物体や背景などの関係性の説明文を自動生成するなど、広範囲な問題に対して適用され始めている。しかしながら、CNN が高い性能を発揮するためには、時に百万枚を超える膨大な量のラベル付き学習データを必要とし、その入手には相当に高い費用と労力を要する。

この弱点を克服するため、ImageNet などの大規模なデータベース(DB)で学習された CNN の中間層の出力を他のタスクへの特徴量として転用する CNN 特徴転移というアプローチが提案された。このアプローチでは、CNN の学習には膨大な量のラベル付データを要するが、転用先のタスクでは、CNN より抽出する特徴に比較的少量の学習サンプルを用いて学習した識別器を適用するのみで、画像認識を行う。学習タスクと転用先のタスクの性質が類似する場合には特に高い性能を発揮するが、大きく異なる場合には期待通りの性能が得られない。

一般に画像認識の認識対象は、複数の意味概念の集まりであることが多い。例えば、画像からその内容を説明する文章を自動生成するタスクであれば、複数の物体間や背景などとその関係性が認識対象である。また、動物の種別認識などの単一物体の認識においても、物体は顔や胴体などのパーツやその模様の特徴などの複数の意味概念から構成される。このように多様な意味概念から構成される対象の高精度な認識には、単一の DB から学習された CNN 特徴のみでは不十分であると考えられる。

この問題に対処するため、可能な限り豊富な概念を含む巨大な学習画像 DB を用いることで汎用性を向上させることが考えられる。しかし、CNN は 1000 種類程度のカテゴリであっても 100 万枚単位のラベル付画像群と数日の学習時間が必要であり、1 つの CNN で多数の概念を学習することは、計算時間の面で現実的ではない。また、類似した認識カテゴリの識別に特化した CNN は認識精度を向上させるという報告もある。従って、認識タスク毎の複数の大規模 DB で個別に CNN を学習し、それらを対象タスクに転用することが、学習効率、精度面で望ましい。

しかしながら従来の研究では、一つの大規模 DB で学習された CNN を他の認識タスクに転用することに留り、CNN 特徴の連結や学習 DB の結合といった単純な統合方法を除き、複数の DB で学習された CNN の特徴転移はほとんど検討されていない。

## 2. 研究の目的

本研究では、複数の大規模 DB で学習された CNN を融合し、異なるタスクの特徴に転用する手法を開発する。画像カテゴリ認識、同一人物検索などの広範なタスクへ転用可能な汎用的な融合特徴抽出手法を目指す。

## 3. 研究の方法

本研究期間では、複数 CNN の融合に必要な以下の 3 点を検討する。

### (1) 各 CNN から取得する特徴の評価

#### ① 有効な層の組み合わせ

CNN の中間層は、一般に低層ほど低次のパターン (エッジや模様など)、上位層ほど高次のパターン (顔や目などの意味概念) を抽出する。複数の CNN から特定の中間層の特徴を抽出する場合、個々の CNN で抽出する特徴の層はそれぞれ可変となる。転用タスクの認識精度を向上させる組み合わせを調査する。

#### ② 局在化した全結合層の有効性

CNN のモデルでは、最終の数層に全結合層を採用するものが多い。全結合層は認識対象カテゴリの意味概念の情報を反映しやすい。この全結合層の出力を画像中に局在化して対応づける技術により、局在化された認識対象カテゴリの意味概念を特徴に用いた場合の性能を評価する。

### (2) 融合方式の開発

2 つの CNN から抽出する、特徴チャンネルの相関関係を記述する融合方式を開発し、対象タスクでの認識性能の向上を行う。以下の A, B の 2 種類の相関関係を評価する。

#### A. 画像中の同一位置での相関

2 つの特徴マップからこれを記述可能な双線形プーリング(BP) [1]を適用する。

#### B. 画面全体での相関

画面全体で特徴を集積した後、BP へ入力したものを評価する。

### (3) 転用先のタスクの学習データに基づく特徴抽出の最適化

転用先タスクでの学習データを用いることで、対象タスクに対してより有効な特徴抽出方法を開発する。

## 4. 研究成果

### (1) CNN から取得する特徴の評価

#### ① 有効な層の組み合わせ

以下の A. B. のタスクにおいて有効な層の組み合わせの評価を行った。

#### A. 画像識別タスク

人物の活動に関する複合的なクラスを含む DB として, 99 種類の文化的な行事と 1 種類の背景のクラスを含む Cultural, 61 種類の日常的な出来事のクラスを含む WIDER, 40 種類の人物行動のクラスを含む Action, 17 種類の人物や動物の視覚的関係のクラスを含む Phrase を用いた. 詳細物体の DB として, 200 種類の鳥のクラスを含む CUB を用いた. 場所を対象とした DB として, 67 種類の屋内クラスを含む Indoor を用いた. 2つの大規模 DB ImageNet 及び Places で各々学習された CNN から特徴を抽出した. 抽出した特徴を線形 SVM により識別した. 数種類の CNN 構造を評価した内, 代表的な構造の ResNet50 の結果を示す.

#### B. カメラ間人物照合タスク

125-500 人以下の人物の画像から同一人物を同定する 4つのカメラ間人物照合 DB を評価用 DB に用いた. 現在の最大規模の人物照合 DB である MSMT17 で学習した ResNet50 に基づくパーツベースモデル (PCB) を用いた. PCB から抽出した特徴量に対して, 距離計量学習手法 XQDA を適用して評価した. この距離計量学習の正当性を確認するための理論解析を行い, 国際会議 MVA2019 で成果発表を行った.

図 1 にタスク A の結果を示す. CNN の各層の特徴を大域的な平均プーリング (GAP) により集積した後, ImageNet と Places DB で学習された特徴を連結している. Action DB では, 最上位畳み込み層 (Conv5\_x) の連結が最良の結果である. 鳥類の詳細カテゴリを認識する CUB DB では, ImageNet における Conv5\_x と Places における低層 (Conv2\_x-Conv4\_x) の組み合わせが最良である. 他の 4つの DB では, Action DB と同様の結果となった. これらの結果より, 人物の活動に関するクラスを識別するには, 場所の意味情報が手がかりとなるが, 同一概念の詳細なカテゴリの認識では, 場所の高レベル情報よりも詳細な見た目の差異を分類するための低レベル特徴が有効と思われる. このようなタスク間での場所の高レベル情報の有効性の差異により, 層の組み合わせに有効性の違いが生じたと考えられる.

タスク B では詳細な見た目の差異への着目が有効と考えられ, 一つの CNN から抽出した低層と高層の特徴の組み合わせの有効性を確認した. 国際会議 ICPR2020 の成果発表では, 抽出した全ての層の連結が最も有効な層の性能を上回ることを示した.

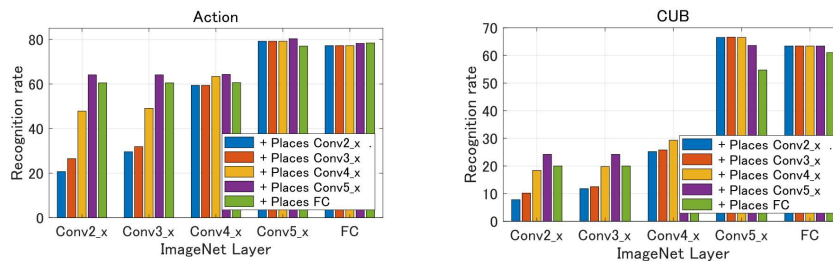


図 1 層の組み合わせの評価. 横軸に示す ImageNet DB で学習された CNN の各層の特徴に対して各棒グラフで示す Places DB で学習された各層の特徴を連結した特徴量の認識率.

#### A. 局在化した全結合層の有効性

全結合層により予測されるクラスの意味概念の情報を画素に局在化させるクラス活性化マップ (CAM) [2] を特徴として利用した場合の性能をタスク A で評価した. CAM の特徴チャンネルの次元数は CNN の学習クラス数に対応し, 多くの CNN 構造で最上位畳み込み層の特徴チャンネルの次元数よりも少ない. 例えば, ResNet50 の Conv5\_x 層の特徴チャンネル数 2048 に対して, ImageNet のクラス数は 1000 である.

最上位全結合層の特徴は学習タスクのクラス代表ベクトルに対応するが, 異なるクラスは共通の概念を含み得るため, これらは互いに相関を持つ. この傾向により, CAM で抽出される特徴は, クラス代表ベクトルの少数の次元のみが重視して得られ, 異なるタスクで特徴抽出を行う際の汎用性が損なわれている. この問題を解決するため, 事前学習クラスの代表ベクトルを特異値分解 (SVD) により正規直交化を施した特徴を開発した. 画像の認識と理解技術に関する会議 MIRU2021 では, この成果を報告する.

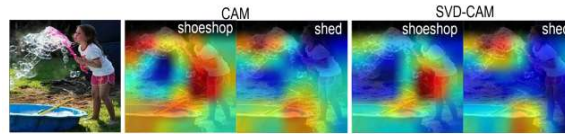


図2 同一クラス (shoeshop/shed) に対する CAM

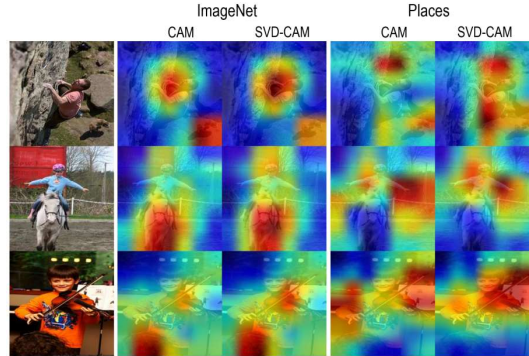


図3 CAMの全クラスの絶対値の平均

図2にPlaces DBで学習されたCAMの例を示す. この例では, SVD-CAMは多くの場所クラスで共通すると思われる背景部分での中程度の大きさを持つ値を減少させ, 人とそれが持つ物体の部位などへは大きな値を残している. 全クラスを通した傾向を解析するため, CAMの絶対値を全クラスで平均化したものを図3へ示す. SVD-CAMは, 画像内で高い値を持つ領域がCAMよりも広い. 特に人物内やその周辺で高い値を持つ領域が増加し, また, 背景などの影響も抑えられている. この傾向は, Placesで学習されたモデルで顕著であり, 場所を認識するタスクで学習したモデルであってもクラス代表ベクトルの正規直交化が, 人物や物体をより重視して特徴を抽出する変更を行ったことが示唆される.

表1に各層より得られる特徴マップのチャンネル情報を大域的な平均プーリング(GAP)を用いて集積した特徴を用いたタスクAの認識結果を示す. GAPの線形性により, CAMをGAPで集積したものが全結合層(FC)層に対応し, 同様にSVD-CAMにGAPを適用して集積したものが, 正規直交化した全結合層(SVD-FC)層に対応する. FC層と比較してSVD-FC層の特徴は全てのDBで高い認識性能を示している. FC層の重みベクトルの特異値を解析したところ, いずれのモデルも50%の次元で80%以上の累積寄与率に到達しており, 少ない次元に特異値が集中していた. これより正規直交化を行っていない場合, 少数の基底を重視した特徴であったため, これを改善したSVD-FCはFCよりも高い性能を示したと考えられる. 最上位畳み込み層(Conv5\_x)の特徴と比較して, SVD-FCはI+Pの場合, 次元数が1/3以下である. SVD-FCはConv5\_xに対して, やや低い性能であるが, その差は, 最大の性能を達成しているモデルに対して, 全てのDBで0.7%以下であった.

表1 ResNet50の各層の特徴をGAPにより集積した特徴の認識率. I, PはそれぞれImageNet, Places DBで事前学習されたモデルから抽出した特徴, P+Iはそれらを連結したものを示す. 赤字, 青字はそれぞれ1番目, 2番目に最良の層, 太字\*は各層で最良のモデルを示す.

Layer	Cultural			WIDER			Action			Phrase			CUB			Indoor		
	I	P	I+P	I	P	I+P	I	P	I+P	I	P	I+P	I	P	I+P	I	P	I+P
Conv2_x	23.0	21.4	<b>22.3*</b>	22.5	21.8	<b>25.2*</b>	20.6	20.0	<b>22.6*</b>	18.6	18.3	<b>20.7*</b>	7.7	7.2	<b>18.8*</b>	30.6	28.7	<b>31.0*</b>
Conv3_x	40.2	35.9	<b>39.6*</b>	30.7	28.7	<b>42.6*</b>	33.3	29.0	<b>30.7*</b>	<b>33.6*</b>	30.4	32.0	12.9	10.8	<b>33.2*</b>	48.0	42.8	<b>46.8*</b>
Conv4_x	66.3	62.6	<b>67.6*</b>	45.0	42.4	<b>45.9*</b>	65.9	53.7	<b>63.4*</b>	<b>67.0*</b>	56.2	64.5	<b>34.4*</b>	21.1	29.3	69.7	67.3	<b>71.3*</b>
Conv5_x	<b>67.5</b>	<b>64.5</b>	<b>71.8*</b>	<b>49.0</b>	<b>48.3</b>	<b>52.3*</b>	<b>79.3</b>	<b>64.1</b>	<b>80.1*</b>	<b>77.5</b>	<b>67.8</b>	<b>78.8*</b>	<b>66.5*</b>	<b>24.1</b>	<b>63.6</b>	<b>76.0</b>	<b>84.9</b>	<b>85.4*</b>
FC	63.9	57.2	<b>68.4*</b>	46.7	44.9	<b>50.3*</b>	77.5	60.7	<b>78.5*</b>	76.5	63.6	<b>77.3*</b>	<b>63.4*</b>	19.9	60.9	73.8	<b>84.3*</b>	83.5
SVD-FC	<b>67.0</b>	<b>61.7</b>	<b>71.1*</b>	<b>48.5</b>	<b>46.5</b>	<b>51.7*</b>	<b>79.0</b>	<b>63.1</b>	<b>80.0*</b>	<b>77.4</b>	<b>67.0</b>	<b>78.6*</b>	<b>66.1*</b>	22.2	<b>63.4</b>	<b>75.7</b>	<b>85.1*</b>	<b>85.1*</b>

## (2) 融合方式の開発

双線形プーリング(BP) [1]に対して行列パワー正規化(MPN) [3]を適用したMPN-BPはCNN特徴の有効な特徴チャンネル融合法であることが知られる. MPNの有効性をタスクBにおける手動設計特徴で評価し, IEEE Trans. on PAMIで発表した論文のAppendixに示した.

BPの出力次元数は入力特徴チャンネルの次元数の2乗のオーダーで増加する. そこで, 部分行列平方根[4]を用い, BPの精度を保ち出力次元数の削減が行えるコンパクト型BP[5](CBP)に対してMPNを適用可能としたMPN-CBPを特徴融合に用いる. MPN-CBPの出力次元数はユーザーが設定でき, これを全ての条件で5000と設定するが, 特徴抽出の計算量は, MPN-CBPに入力する特徴マップの特徴チャンネル次元数に依存する.



そこで、前節のSVD-CAMをMPN-CBPの入力特徴マップとして用い、入力特徴チャンネル次元数を低く抑える。各層から得られた特徴をMPN-CBPで融合した特徴を用いたタスクAの認識結果を表2に示す。表1の結果と比較することで、MPN-CBPによるI+Pの融合が、GAPにより得られる特徴の単純連結よりも有効な特徴融合方法であることが確認できる。MPN-CBPに入力する層としてSVD-CAMは、Conv5\_xよりも高い性能が同等である。なおConv5\_xの特徴に対して1/3程度の次元数であるI+Pの条件では、SVD-CAMを用いたMPN-CBPは、Conv5\_xを用いた場合に対して20%程度特徴抽出時間が短かった。

SVD-CAMをMPN-CBPの入力に直接用いる場合、3.(2)A.画像中の同一の位置における相関に対応し、SVD-CAMにGAPを適用したSVD-FCに対してMPN-CBPを適用する場合、3.(2)B.画面全体での相関に対応する。SVD-FCは、CulturalとWIDERで有効であった。これらは、青森ねぶた祭や交通事故などの比較的、高レベルな概念を認識するタスクであるため、画面内で存在している意味情報の共起関係が有効であるためと考えられる。その他のデータでは、SVD-CAMが有効であり、これは特にCUBで有効であった。画面内に小さく映る鳥の詳細な種類を識別するタスクであるCUBは、同一の位置に着目した詳細な見た目情報が、有効であるためと考えられる。

表2 ResNet50の各層の特徴をMPN-CBPにより集積した特徴の認識率。P+Iは、P,Iの各モデルから抽出した特徴チャンネルを連結したマップをMPN-CBPの入力としている。

Layer	Cultural			WIDER			Action			Phrase			CUB			Indoor		
	I	P	I+P	I	P	I+P	I	P	I+P	I	P	I+P	I	P	I+P	I	P	I+P
Conv2_x	49.7	46.8	<b>50.0*</b>	33.3	32.6	<b>34.1*</b>	34.7	33.4	<b>35.4*</b>	33.7	32.5	<b>34.5*</b>	21.6	19.9	<b>22.1*</b>	48.8	46.1	<b>49.4*</b>
Conv3_x	61.5	58.4	<b>60.2*</b>	39.4	38.6	<b>38.8*</b>	48.6	44.7	<b>47.0*</b>	<b>48.8*</b>	45.1	48.2	<b>34.9*</b>	29.2	33.4	60.8	59.3	<b>62.1*</b>
Conv4_x	68.4	65.2	<b>67.4*</b>	43.5	41.8	<b>42.4*</b>	72.6	59.7	<b>68.4*</b>	<b>73.8*</b>	62.7	69.3	<b>61.8*</b>	38.3	52.4	73.5	71.6	<b>73.8*</b>
Conv5_x	69.9	<b>67.6</b>	<b>72.6*</b>	48.2	48.2	<b>50.8*</b>	<b>80.5</b>	<b>67.2</b>	<b>81.3*</b>	78.4	<b>70.8</b>	<b>79.3*</b>	<b>72.8*</b>	<b>33.1</b>	<b>70.4</b>	75.8	<b>85.1*</b>	84.5
CAM	68.7	64.1	<b>71.9*</b>	48.4	47.9	<b>50.9*</b>	79.7	64.7	<b>80.1*</b>	<b>79.2*</b>	68.6	79.1	<b>71.2*</b>	29.2	69.4	74.5	<b>84.7*</b>	84.0
SVD-CAM	<b>70.2</b>	<b>67.3</b>	<b>73.1*</b>	48.4	48.8	<b>51.4*</b>	80.4	67.0	<b>81.3*</b>	<b>78.6</b>	<b>71.0</b>	<b>80.0*</b>	<b>72.4*</b>	<b>32.9</b>	<b>71.4</b>	75.4	<b>85.3</b>	<b>85.7*</b>
FC	67.8	62.2	<b>71.8*</b>	50.2	48.6	<b>52.9*</b>	77.3	59.7	<b>77.6*</b>	<b>79.1*</b>	64.5	78.9	<b>64.5*</b>	20.5	60.7	<b>75.9</b>	<b>84.7*</b>	83.6
SVD-FC	<b>69.8</b>	66.0	<b>73.5*</b>	<b>51.0</b>	<b>49.8</b>	<b>53.6*</b>	79.3	63.3	<b>79.6*</b>	79.2	68.5	<b>79.6*</b>	<b>67.0*</b>	23.7	63.9	<b>76.8</b>	<b>85.4</b>	<b>85.2*</b>

### (3) 転用先のタスクの学習データを利用した特徴抽出の最適化

転用先のタスクの学習データを用いて、学習済みモデルから畳み込み層の特徴を判別的に集積する方法を開発し、タスクBへ適用した。判別的な重みマップ学習法にランダム射影を導入し、学習の30倍以上の高速化を実現した。これに加え、距離計量学習XQDAと複数の層の特徴の利用により、学習済みモデルを小サンプルの学習データを用いてファインチューニングを行うよりも高精度な人物照合率を達成可能な特徴量を高速に学習できた。この研究成果を国際会議ICPR2020で発表した。

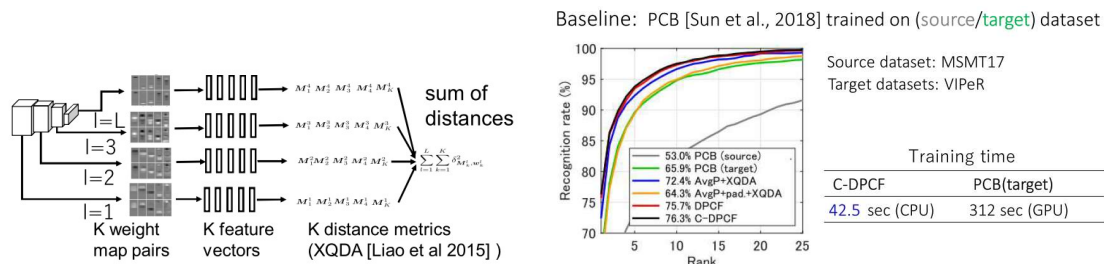


図4 左: 提案モデル(C-DPCF). 右: ファインチューニング(PCB(target))との比較.

### <引用文献>

- [1] T.-Y. Lin, A. RoyChowdhury, S. Maji, Bilinear convolutional neural networks for fine-grained visual recognition, IEEE Trans. on PAMI, vol.40, no. 6, pp.1309-1322, 2018.
- [2] B. Zhong, A. Khosla, A. Lapedriza, A. Oliva, A. Torralba, Learning deep features for discriminative localization, In CVPR2016.
- [3] P. Li, J. Xie, Q. Wang, W. Zhuo, Is second-order information helpful for large-scale visual recognition, In ICCV2017.
- [4] M. Gou, F. Xiong, O. Camps, M. Szanier, MoNet: Moments embedding network, CVPR2018.
- [5] Y. Gao, O. Beijbom, N. Zhang, T. Darrell, Compact bilinear pooling, In CVPR2016.

5. 主な発表論文等

〔雑誌論文〕 計1件（うち査読付論文 1件 / うち国際共著 0件 / うちオープンアクセス 0件）

1. 著者名 Tetsu Matsukawa, Takahiro Okabe, Einoshin Suzuki, Yoichi Sato	4. 巻 42
2. 論文標題 Hierarchical Gaussian Descriptors with Application to Person Re-Identification	5. 発行年 2020年
3. 雑誌名 IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)	6. 最初と最後の頁 2179 ~ 2194
掲載論文のDOI（デジタルオブジェクト識別子） 10.1109/TPAMI.2019.2914686	査読の有無 有
オープンアクセス オープンアクセスではない、又はオープンアクセスが困難	国際共著 -

〔学会発表〕 計3件（うち招待講演 0件 / うち国際学会 2件）

1. 発表者名 Tetsu Matsukawa, Einoshin Suzuki
2. 発表標題 Kernelized Cross-View Quadratic Discriminant Analysis for Person Re-Identification
3. 学会等名 Sixteenth International Conference on Machine Vision Applications (MVA 2019) (国際学会)
4. 発表年 2019年

1. 発表者名 Tetsu Matsukawa, Einoshin Suzuki
2. 発表標題 Convolutional Feature Transfer via Camera-specific Discriminative Pooling for Person Re-Identification
3. 学会等名 in 25th International Conference on Pattern Recognition (ICPR2020) (国際学会)
4. 発表年 2021年

1. 発表者名 松川徹, 鈴木英之進
2. 発表標題 正規直交クラス代表ベクトルとの内積表現によるConvNet特徴転移
3. 学会等名 画像の認識・理解シンポジウム (MIRU2021)
4. 発表年 2021年

〔図書〕 計0件

〔産業財産権〕

〔その他〕

-

6. 研究組織

	氏名 (ローマ字氏名) (研究者番号)	所属研究機関・部局・職 (機関番号)	備考
--	---------------------------	-----------------------	----

7. 科研費を使用して開催した国際研究集会

〔国際研究集会〕 計0件

8. 本研究に関連して実施した国際共同研究の実施状況

共同研究相手国	相手方研究機関
---------	---------