

特集 高速かつメモリ消費量の少ない局所特徴量*

CARD: Compact And Real-time Descriptors

安倍 満 Mitsuru AMBAI
 吉田 悠一 Yuichi YOSHIDA

This paper proposes a new algorithm, Compact And Real-time Descriptors (CARD), for an object recognition system in computer vision technology. Currently, the SIFT (Scale-Invariant Feature Transform) algorithm has been used in that technology field for a long time. However, the SIFT requires a long computation time and a large memory capacity. Therefore, we developed CARD which realizes a very short computation time, and operates perfectly with less memory capacity because of the compact expression of the local feature described by short binary codes. A new efficient algorithm based on lookup tables is presented for extracting histograms of the oriented gradients, and results in a computation time of approximately 16 times faster per descriptor than that of SIFT. Our lookup-table-based approach can handle arbitrary layouts of bins, such as the grid binning of SIFT and the log-polar binning of GLOH (Gradient Location and Orientation Histogram), thus yielding sufficient discrimination power. In addition, we introduced supervised sparse hashing to convert the extracted descriptors to short binary codes. This conversion is achieved very quickly by multiplying a very sparse integer weight matrix by the descriptors and aggregating the signs of their multiplications. The weight matrix is optimized in a training phase so as to make the Hamming distances between encoded training pairs reflect the visual dissimilarities between them. Experimental results demonstrate that CARD outperforms the current methods in terms of both computation time and memory usage.

Key words : Local Descriptors, SIFT, Matching, CARD

1. はじめに

画像間の対応付けはコンピュータビジョンにおける最も基本的なタスクであり、近年では局所特徴量が有効であると考えられている。Scale-Invariant Feature Transform (SIFT) ¹⁾ はその代表的な手法であり、スケールおよび回転不変な局所特徴量を抽出可能としている。SIFTの応用範囲は広く、3次元形状復元 ²⁾、特定物体認識 ³⁾、一般物体認識 ⁴⁾ など、様々な目的に利用されている。SIFTを始めとして、多くの局所特徴量が提案されているが、これらは計算速度が遅く、またメモリ消費量が大きいという課題を抱えている。そのため車載向けCPUや携帯端末などといった、PCよりも比較的性能の低い機器で動作させることが困難という課題が残されている。

以下、具体例としてSIFTを例に挙げて議論する。

Fig. 1はSIFTの計算速度を示したものであり、サイズが640×480の画像から3,762個のキーポイントを得たときの結果である。全体の処理に2秒近くを要しており、

そのうち約75%がキーポイントの方向推定と局所特徴量の抽出に費やされていることが分かる。また、SIFTは128次元という高次元の特徴ベクトルで表現されるため、大きくメモリを消費する。多くの場合、この特徴ベクトルはメモリ消費量を抑えるために符号無し8bit整数の配列に格納されるが、それでも1つのキーポイントあたり8×128=1024ビット必要である。

そこで本論文では、これら2つの課題に焦点を当てる。すなわち、高速に計算可能であり、128ビット程度のバイナリコードで表現可能な局所特徴量 Compact And Real-time Descriptors (CARD) を提案する。

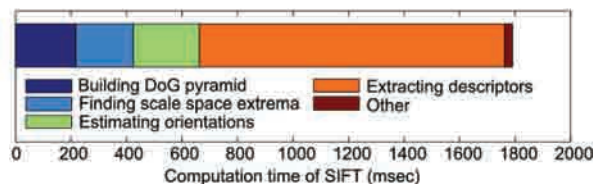


Fig. 1 Estimating orientations and extracting descriptors from each patch account for 75% of the SIFT computation time We used an Intel Core 2 Duo 2.66 GHz processor in this experiment.

* 画像の認識・理解シンポジウム (MIRU2011) の発表論文より、一部修正して転載。

また、この英訳論文が国際会議International Conference on Computer Vision (2011) に採択され、出版されている。

1. 安倍満, 吉田悠一, “高速かつメモリ消費量の少ない局所特徴量”, 画像の認識・理解シンポジウム (MIRU2011) 論文集, 2011, P: 1682-1689, 発行年: 2011-07-20

2. Mitsuru Ambai and Yuichi Yoshida, ‘CARD: Compact And Real-time Descriptors’, 2011 IEEE International Conference on Computer Vision, P97-104, 6-13 Nov. 2011

1.1 関連研究

1.1.1 局所特徴量の高速化

局所特徴量の高速化は、産業応用という観点から非常に意義が深く、幅広く研究が行われている。Bay⁵⁾らは、SIFTと同様に回転かつスケール不変な局所特徴量 Speeded-Up Robust Features (SURF) を提案し、SIFTに対して数倍の高速化を達成した。Takacs⁶⁾らは、携帯端末上での高速な特徴追跡とコンテンツ認識を目的としてRotation-Invariant, Fast Feature (RIFF) を提案した。同様に、Wagner⁷⁾らもまた、携帯端末上の動作を目的とした局所特徴量を提案し、拡張現実感における高速なカメラ姿勢推定を実現した。近年では、前記の研究事例よりもさらなる高速化に焦点を当てた研究^{8) 9)}が行われており、高速化という課題に対する注目度の高さを窺い知ることができる。しかしながら、これらの研究の多くは、SIFTよりも簡略化された局所特徴量を用いることで高速化を達成しているため、たとえば大規模な特定物体認識などに適用すると十分な認識性能が得られないという課題がある。

1.1.2 局所特徴量のコンパクト表現

近年、特徴ベクトルを短いバイナリコードに変換する Binary hashing^{10) -15)} が注目されている。数多くの手法が提案されているが、変換前の局所特徴量を $d \in R^D$ 、変換後のバイナリコードを $b \in \{0,1\}^B$ と記述するとき、Binary hashingは $b = \left\lfloor \frac{\text{sgn}(f(W^T d)) + 1}{2} \right\rfloor$ と一般的に定義できることが知られている¹³⁾。ここで、 D は局所特徴量の次元数、 B はバイナリコードのビット長、 $f(\cdot)$ は任意の関数、 $W \in R^{D \times B}$ は変換行列である。バイナリコード同士の類似度はハミング距離によって計算できる。これは変換前における局所特徴量同士の類似度の計算（コサイン距離、ベクトル間角度、ユークリッド距離など）よりもはるかに高速に計算可能である。従って、メモリ効率化だけでなく、最近傍探索の速度も向上できるという利点がある。

Binary hashingの性能は、 W と $f(\cdot)$ の定義によって大きく異なる。最も単純な例はRandom projections^{12) 16)} であり、 $f(\cdot)$ は恒等関数、 W の各要素は正規分布に基づく乱数で定義される。一方、学習を導入することで、乱数に基づく方法よりもより短いバイナリコードが得られることが知られている^{10) 13)}。例えばSpectral hashing¹⁰⁾ では、 $f(\cdot)$ に非線形関数、 W に学習データの主成分軸を用いている。以上、Binary hashingにつ

いて述べたが、この他にも情報エントロピーに基づく方法¹⁷⁾、次元削減に基づく方法¹⁸⁾ など、局所特徴量を圧縮するための様々な手法が提案されている。

以上述べた手法は、バイナリコードへの変換にかかる計算時間についての議論が欠けている。リアルタイムなアプリケーションにおいては、この変換も高速に行える必要がある。Sparse random projections^{14) 15)} はこの課題に焦点を当てている。Achlioptas¹⁴⁾ らは、 W の要素として3種類の整数 $\{-1,0,1\}$ をそれぞれ $\left\{ \frac{1}{6}, \frac{2}{3}, \frac{1}{6} \right\}$ の確率で選択することで、Binary Hashingにおける $W^T d$ の計算量を大幅に削減した。これに対し、Li¹⁵⁾ らは3種類の整数 $\{-1,0,1\}$ をそれぞれ $\left\{ \frac{1}{2\sqrt{D}}, \frac{1}{\sqrt{D}}, \frac{1}{2\sqrt{D}} \right\}$ の確率で選択し、 W をより疎にすることでさらなる高速化が可能であると指摘している。なお、 $f(\cdot)$ としては恒等関数を用いる。Sparse random projectionsは学習に基づいていないため、Random projectionsと同様に、バイナリコード化すると性能が低下するという欠点がある。

1.2 提案手法(CARD)の特徴

高速な局所特徴量抽出：2種類のルックアップテーブルを用いることで、1つのキーポイントあたりSIFTの約16倍の速度で局所特徴量を計算可能とした。CARDはSIFTおよびGradient Location and Orientation Histogram (GLOH)と同様に勾配ヒストグラム特徴量であり、十分な識別能力を有することを示した。

Learning-based sparse hashing：Sparse random projectionsに学習の概念を導入した。変換前における距離と、変換後のバイナリコード間のハミング距離が相関するように W を学習することで、Sparse random projectionsの性能を改善した。

2. 高速な局所特徴抽出

2.1 キーポイント検出

スケール不変性はキーポイント検出における重要な性質である。Lowe¹⁾ はDifference of Gaussians (DoG) ピラミッドの極値をスケール不変なキーポイントとして定義した。Dufournaud¹⁹⁾ らは、分散を徐々に大きくした複数のガウシアンフィルタを入力画像に適用することで、スケール不変なHarrisコーナー検出器を実現した。これらの欠点は、異なる分散のガウシアンフィルタを幾度も適用しなければならず、計算が重いという点である。

これに対して我々は、入力画像を一定の割合で縮小した多段のピラミッド画像を作成し、各層における画像から独立に特徴点²⁰⁾を抽出する方法を採用した。この方法でも、実用上十分なスケール不変性が得られることが報告されている^{6) 7)}。0段目を入力画像とするとき、 $n-1$ 段目の画像を $\frac{1}{\sqrt{2}}$ 倍に縮小することで n 段目の画像を生成した。実際の計算においては、 n 段目の画像は $n-2$ 段目の画像を $\frac{1}{2}$ にダウンサンプリングするだけで得られるため、高速にピラミッド画像を生成できる。

2.2 勾配ヒストグラムに基づく局所特徴量

SIFTやGLOHは勾配ヒストグラムに基づく局所特徴量であり、他に比べて高い識別能力を持つことが知られている²¹⁾。本節では、これを高速に算出するアルゴリズムについて述べる。提案するアルゴリズムは、(a) ピクセルの線形補間処理が不要、(b) ルックアップテーブルによる効率的な計算が可能という2つの特徴を持つ。

2.2.1 従来法(SIFT, GLOH)の概要

提案アルゴリズムを説明する前に、従来手法(SIFT, GLOH)による局所特徴量抽出の概要を述べる。まずは記号を整理する。入力画像を $I(x,y)$ 、キーポイントの座標を $(p_x, p_y)^T$ と記述する。ピクセルの勾配の強度 $m(x,y)$ と方向 $\theta(x,y)$ は次のように計算できる。

$$m(x,y) = \sqrt{I_x(x,y)^2 + I_y(x,y)^2} \tag{1}$$

$$\theta(x,y) = \text{angle}(I_x(x,y), I_y(x,y)) \tag{2}$$

ここで、 I_x および I_y は $I(x,y)$ の x,y 方向微分であり、 $\text{angle}(x,y)$ は方向をラジアンで返す関数である。角度は u 軸の正方向から v 軸正方向に向かって反時計回りに $(0, 2\pi)$ の範囲の値をとる。

推定された方向とスケールに基づき、キーポイント周辺の領域は変形(拡大/縮小, 回転)される。これにより、局所特徴量はスケールおよび回転に関する不変性を得る。変形された周辺領域は K 個の小領域に分割される。例えば、SIFTでは 4×4 のグリッド分割、GLOHでは放射状(log-polar)分割が行われる。各小領域から L 次元の方向ヒストグラム h_k を抽出する。

$$h_k = (h_{k,0}, h_{k,1}, \dots, h_{k,(L-1)})^T \in R^L \tag{3}$$

ここで、 k は小領域のインデックスである。 h_k を全て統合することで、 KL 次元の局所特徴量 $d \in R^{KL}$ を得る。

$$d = (h_0^T, h_1^T, \dots, h_{K-1}^T)^T \tag{4}$$

記述を明確にするために、次節以降ではキーポイント周辺領域の分割の仕方を「ブロック分割パターン」、各小領域を「ブロック」と記載することにする。

なお、勾配ヒストグラムを求めるときに必要な、ラジアン表現の角度を0から $N-1$ に量子化する関数 $Q_N(\theta)$ は次のように定義できる。

$$Q_N(\theta) = \left\lfloor \frac{N\theta}{2\pi} + \frac{1}{2} \right\rfloor \bmod N \tag{5}$$

2.2.2 ピクセル補間を用いない局所特徴量抽出

従来手法では、キーポイント周辺の領域を変形(拡大/縮小, 回転)させることにより、スケールおよび回転に関する不変性を獲得しているが、この処理には膨大なピクセル補間処理が伴う。

これに対し、我々はブロック分割パターン自体をキーポイント周辺領域の方向(以下、「主軸」と記述)に応じて回転させることで、計算負荷の高いピクセル補間処理を回避した。主軸の角度(ラジアン)を p とするとき、回転済みブロック分割パターンは3つのパラメータを持つ関数 $\psi(u,v,p)$ で記述できる。ここで $(u,v)^T$ はキーポイントの座標 $(p_x, p_y)^T$ を中心としたときの相対座標である。

$\psi(u,v,p)$ は0から $K-1$ のブロック番号を返す関数であり、ブロック分割パターンは uv -平面上に現れ、 p に応じて回転するように定義する。Fig. 2に $\psi(u,v,p)$ の例を示した。以上に基づく、キーポイント周辺における座標位置 $(x,y)^T$ のピクセルが属するブロックは次のように決定できる。

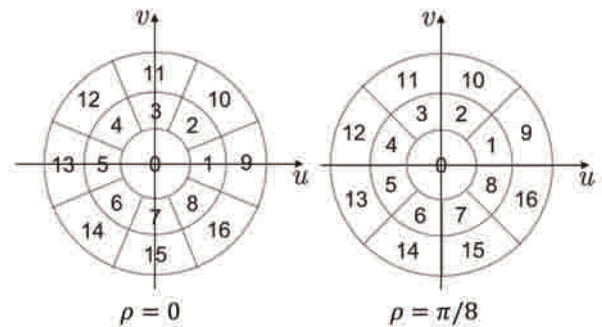


Fig. 2 Parametric representation of a log-polar binning pattern.

$$k = \psi(x - p_x, y - p_y, \rho) \quad (6)$$

座標位置 $(x, y)^T$ におけるピクセルの方向は L 段階に量子化され、投票処理を行うことで h_k が得られる。ただしブロック分割パターンが回転しているため、 p による回転の影響を打ち消すためにピクセルの勾配方向 $\theta(x, y)$ を補正する必要がある。勾配方向を補正し、量子化した値 l は次式で計算できる。

$$l = Q_L(\theta(x, y) - \rho) \quad (7)$$

結果として、キーポイント $(p_x, p_y)^T$ から局所特微量 d を抽出するアルゴリズムは次のように簡略化できる。

1. h_0, h_1, \dots, h_{K-1} をゼロで初期化。
2. キーポイント $(p_x, p_y)^T$ の周辺全てのピクセル $(x, y)^T$ について次の2つの処理(a)(b)を適用。
 (a) 式 (6) (7) を計算し、 k, l を取得。
 (b) $h_{k,l} \leftarrow h_{k,l} + m(x, y)G_{\sigma_1}(x - p_x, y - p_y)$
3. h_0, h_1, \dots, h_{K-1} を統合し、 d を取得。
4. d の長さを1に正規化

ブロック分割パターンとしては、GLOHで採用されているような放射状 (log-polar) 分割が最も良いと報告されている²¹⁾。そこで、本論文でもこれを適用する。このとき、関数 $\psi(u, v, \rho)$ は次のように定義できる。

$$\psi(u, v, \rho) = \begin{cases} 0 & , r_0 \leq r < r_1 \\ 1 + Q_8(\text{angle}(u, v) - \rho) & , r_1 \leq r < r_2 \\ 9 + Q_8(\text{angle}(u, v) - \rho) & , r_2 \leq r \leq r_3 \end{cases} \quad (8)$$

ここで、 $r = \sqrt{u^2 + v^2}$ であり、 $r_1 r_2 r_3$ は分割のためのしきい値である。Fig. 2のように、キーポイント周辺領域は17個のブロックに分割される。なお、提案手法では、放射状分割だけでなくSIFTで採用されているグリッド分割など、任意のパターンを扱うことができる。

2.2.3 ルックアップテーブルによる高速化

実際には、式 (6) (7) の計算は後述する2つのルックアップテーブル (ψ および χ) で置き換えられるため、大幅な高速化が可能である。

ルックアップテーブルに基づくアルゴリズムを説明する前に、キーポイントの主軸の求め方について整理しておく。主軸は、キーポイント周辺領域の勾配を M 段階に量子化し、勾配ヒストグラム \hat{h} の最大値を探すことで得られる。

1. $\theta(x, y)$ を M 段階に量子化： $i_\theta(x, y) = Q_M(\theta(x, y))$.

2. $\hat{h} = (\hat{h}_0, \dots, \hat{h}_{M-1})^T$ をゼロで初期化。
3. キーポイント $(p_x, p_y)^T$ の周辺全てのピクセル $(x, y)^T$ について次の2つの処理 (a) (b) を適用。
 (c) $i = i_\theta(x, y)$.
 (d) $\hat{h}_i \leftarrow \hat{h}_i + m(x, y)G_{\sigma_2}(x - p_x, y - p_y)$.
4. 分散 σ_s の正規分布で h を平滑化。
5. 勾配ヒストグラム h の最大値を主軸として取得。 $i_p = \text{argmax}(h_i)$

まず、式 (6) をルックアップテーブルで置き換えることを考える。 i_p は量子化されているため、 $\psi(u, v, \rho)$ の計算はルックアップテーブル化が可能である。そのために、まず $Q_n(\theta)$ の逆関数を導入する。

$$Q_N^{-1}(i) = (2\pi(i \bmod N)) / N \quad (9)$$

$Q_N^{-1}(i)$ は量子化済の勾配方向を、当該量子化レベルにおける代表値 (ラジアン) に変換する関数である。量子化済の主軸 i_p は、 $Q_N^{-1}(i_p)$ を用いてラジアンに戻すことができる。そこで、式 (8) に $\rho = Q_M^{-1}(i_p)$ を代入することで、ブロック分割パターンテーブル $\bar{\psi}(u, v, i_p)$ が得られる。

$$\bar{\psi}(u, v, i_p) = \psi(u, v, Q_M^{-1}(i_p)) \quad (10)$$

u, v, i_p は有限の範囲内の整数しか取りえないため、 $\bar{\psi}(u, v, i_p)$ は事前に計算が可能である。Fig. 3に $\bar{\psi}(u, v, i_p)$ を図示した。 $\bar{\psi}(u, v, i_p)$ を用いると、式 (6) を次のように置き換えることができる。

$$k = \bar{\psi}(x - p_x, y - p_y, i_p) \quad (11)$$

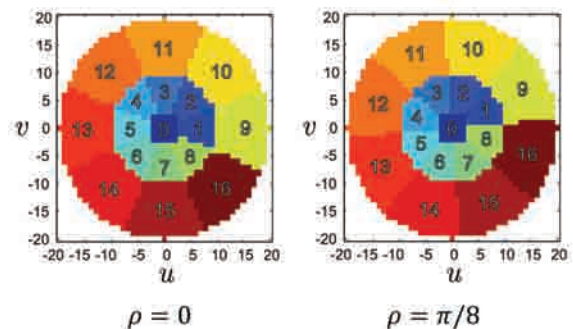


Fig. 3 Spatial binning table

次に、式 (7) をルックアップテーブルで置き換えることを考える。主軸を求める処理において、 M 段階に量子化された勾配方向 $i_\theta(x, y)$ が得られているが、これは式 (7) を計算するときに再利用が可能である。そし

て、この再利用は次の近似を導入することで、ルックアップテーブルにより実現できる。

$$l = Q_L(\theta(x,y) - \rho) \approx Q_L(Q_M^{-1}(i_\theta(x,y) - i_\rho)) = \chi(i_\theta(x,y), i_\rho) \quad (12)$$

式(12)では、 $i_\theta(x,y)$ から i_ρ を引くことにより回転を補正し、合成関数 $Q_L(Q_M^{-1}(\cdot))$ により量子化レベル数を M から L に変換している。この近似は、 $M \gg L$ であれば妥当である。 $i_\theta(x,y)$ は0から $M-1$ の範囲内の整数しか取りえないため、関数 χ は事前に計算可能であり、 $M \times M$ の大きさのテーブルとしてメモリ上に保持できる。従って、式(6)(7)を式(11)(12)で置き換えることで、大幅な高速化が達成される。(なお、正規分布 G_{01}, G_{02}, G_{03} もテーブル化が可能である。)

実験ではパラメータ $\{r_1, r_2, r_3, \sigma_1, \sigma_2, \sigma_3, M, L\}$ としてそれぞれ $\{3, 10, 20, 15, 10, 3, 40, 8\}$ を用いた。この設定の場合、 $8 \times 17 = 136$ 次元の局所特徴量が得られる。

3. Learning based sparse hashing

3.1 コスト関数

抽出された局所特徴量 $d \in R^D$ は、次式に示すBinary hashing関数によりバイナリコード $b \in \{0, 1\}^B$ に変換する。

$$b = (\text{sgn}(W^T d) + 1) / 2 \quad (13)$$

ここで D は局所特徴量の次元数、 B はビット長、 $W \in R^{D \times B}$ は変換行列である。バイナリコード間の距離は、ハミング距離をビット長で正規化した値で与えられる。

$$D_h(b_u, b_v) = \frac{1}{B} \|b_u - b_v\|_1 \quad (14)$$

Learning based sparse hashingでは、変換前の局所特徴量間の距離と、変換後のバイナリコード間の距離がなるべく一致するように W を学習することを考える。そこで、次のようにコスト関数を定義する。

$$C(W) = \sum_{(u,v) \in P} (D_\theta(d_u, d_v) - D_h(b_u, b_v))^2 \quad (15)$$

ここで P は学習用ペアの集合である。 $D_\theta(d_u, d_v)$ はベクトル間の角度を π で正規化したものであり、次のように定義される。

$$D_\theta(d_u, d_v) = \arccos\left(\frac{\langle d_u, d_v \rangle}{\|d_u\| \|d_v\|}\right) / \pi \quad (16)$$

実行速度が重視されるアプリケーションにおいては、式(13)を高速に計算しなければならない。バイナリコ

ード化のボトルネックは $W^T d$ の計算部分であり、 $D \times B$ 回の掛算と $(D-1) \times B$ 回の加算を要する。そこで Sparse random projections^{14) 15)} の考えに基づき、 W が疎で整数であるという制約条件を導入する。すなわち、 W の要素は $\{-1, 0, 1\}$ のいずれかの値をとり、非ゼロの要素数が S 個であるという制約の下で、 $C(W)$ を最小化する。これにより、 $W^T d$ は合計 S 回の加算と引算で計算できるようになる。 S を小さく設定することで、変換に要する計算時間を十分小さくできる。次章の実験で示すが、興味深いことに W の要素のうち 90% をゼロに設定した場合でも、性能を落とすことなくバイナリコード化ができることが分かった。これはバイナリコード変換に要する負荷を大幅に減らすことに繋がるため、応用上重要な性質であるといえる。

3.2 Greedy アルゴリズムによる最適化

コスト関数 $C(W)$ は非線型の不連続関数であるから、Levenberg-Marquardt法やGauss-Newton法などのような勾配に基づく手法が適用できない。そこで、Greedyアルゴリズムによる最適化アルゴリズムを適用した。我々の戦略では、ランダムに2つの要素を W から選択し、この2つの要素を最適な値に置き換えることを考える。

W は $w_{ij} \in \{-1, 0, 1\}$ および $\sum |w_{ij}| = S$ という2つの制約条件を満たす必要があるが、この制約のおかげで、選んだ2つの要素が取りえる値は少ない組み合わせに限られる。従って、次に示すGreedyアルゴリズムで最適化が行える。

1. 次の二つの処理 (a) (b) により、 W を初期化。
 - (a) 全ての要素を1,-1でランダムに初期化。
 - (b) ランダムに $(DB-S)$ 個の要素を選び、0に設定。
2. ランダムに2つの要素 w_{uv}, w_{pq} を選択。
3. w_{uv}, w_{pq} を更新。
 - (a) (w_{uv}, w_{pq}) が両方とも非ゼロの場合
 $(1,-1), (-1,1), (1,1), (-1,-1)$ をチェックし、 $C(W)$ を最小化する組み合わせを採用。
 - (b) w_{uv}, w_{pq} が両方ともゼロの場合
 ステップ2に戻って繰り返し。
 - (c) w_{uv}, w_{pq} のいずれか一方がゼロの場合
 $(0,1), (0,-1), (1,0), (-1,0)$ をチェックし、 $C(W)$ を最小化する組み合わせを採用。
4. 収束するまでステップ2に戻って繰り返し。

Fig. 4に $C(W)$ が減少する様子を示した。3つの異なる

初期値から始めた結果を図示しているが、どの場合においても $C(W)$ が正しく減少していることが分かる。以上により、局所特徴量はバイナリコードへ変換できる。なお、多くの Binary hashing 手法と同様に、実際に適用するときには d の平均を求め、事前に重心を原点にずらしておく。

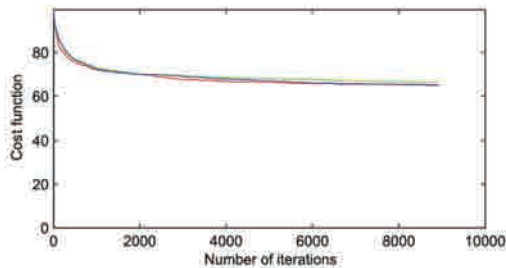


Fig. 4 The cost value $C(W)$ is minimized by the algorithm proposed in Section 3.2. In this example, 25,000 training pairs were used. The bit length B was set to 64. 90% of the elements in W were set to zero.

4. 実験結果

本節では、提案する局所特徴量をバイナリコード化した場合としなかった場合の両方について評価を行う。

以下、 W におけるゼロ要素が占める割合 $\left(1 - \frac{S}{DB}\right)$ を「ゼロ要素率」と記述することにし、ゼロ要素率が性能にどのように影響するかを議論する。

4.1 弁別性

まず、バイナリコード化する前の、局所特徴量の弁別性能を確かめる。そこで Fig. 5 に、局所特徴量間の角度の確率密度関数を示した。左が CARD、右が SIFT の結果であり、SIFT では DoG フィルタをキーボ

イント検出器として用いた。2 画像間の正確な Homography 行列が既知である Mikolajczyk らによるデータベース^{21) 22)}を用いてマッチング結果として正しいペアと正しくないペアを別々に収集し、それぞれから得た確率密度関数を独立に示した。

CARD と SIFT から得られた確率密度関数は非常に似ていることが分かる。正しくないマッチングペアの確率密度分布は分散の小さい正規分布とみなせる。一方、正しいマッチングペアの確率密度分布は広い分布を持つようである。これは、Mikolajczyk らによるデータベースにはアフィン歪みや射影歪みなどのような複雑な変形が含まれており、正しいマッチング位置であっても、局所特徴量が十分な弁別性能を発揮できなかったことによると考えられる。これは、SIFT、CARD ともに同様である。

正しくないマッチングペアの確率密度分布は正規分布で近似できるため、適切なしきい値 t によって誤マッチングを棄却できると考えられる。CARD の場合、この正規分布の平均と標準偏差はそれぞれ 0.49 および 0.05 であるから、3シグマの法則に従うと、 $t = 0.49 - 3 \times 0.05 = 0.34$ と設定することで 99.7% の誤マッチングが棄却できると考えられる。

4.2 バイナリコードの性能

局所特徴量の隣接関係は、バイナリコード変換前と変換後で保存されねばならない。たとえば、バイナリコード変換前の空間において最近傍に位置する局所特徴量は、変換後のハミング空間においても最近傍であることが望まれる。そこで、Mean Average Precision (MAP) によって Binary hashing による近傍探索の性

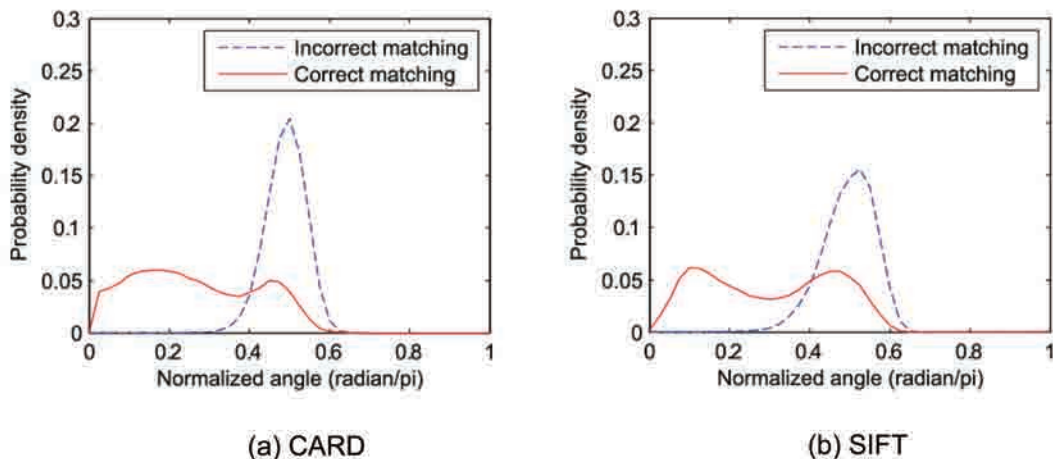


Fig. 5 Probability density of normalized angles. The x-axis is the radian normalized by π .

能を評価した。MAPはPrecision-Recall曲線によって囲まれる領域の面積を近似したものであり、次の手順によって計算した。まずCaltech101データセット²³⁾の各カテゴリから1枚ずつ画像を抽出し、136次元の局所特徴量を計算した。ここから50,000個の学習用サンプルと10,000個のテスト用サンプルを、重複の無いようにランダムに抽出した。それぞれのテストサンプルにおいて、学習用サンプルとの距離(式(16))を計算し、距離が t 以下の候補を「真の近傍」として定義した。次に、学習用サンプルを用いて W を学習し、全ての学習用サンプルとテスト用サンプルをバイナリコードへ変換した。それぞれのテストサンプルについて、学習用サンプルをハミング距離(式(14))の近いものから順に並べたランキングリストを作成した。このランキングリストの上位に真の近傍が正しく現れるか否かをMAPにより評価した。

Fig. 6(a)は縦軸にMAP、横軸に学習ペアの数を示したものである。Learning based sparse hashingで用いる学習ペアとしては、学習用サンプルからランダムに選択する方法と、ペア間の距離が近いものと遠いものを優先的に選択する方法の2種類を試した。また、学習を用いないSparse random projectionsの結果も比較対象として示した。いずれの場合においても、ゼロ要素率は0.3、ビット長は64である。Fig. 6(a)に示されている通り、学習ペアをランダムに選択したときに限り、学習ペアを増やすとMAPが向上した。すなわち、学習によってSparse random projectionsの最近傍探索性能が改善できることが分かった。ただし、学習ペアの選び方としてペア間の距離を活用した方法ではMAPの改善は見られなかった。学習ペアの選択方法は

Learning based sparse hashingの性能を決める重要な検討項目であるといえる。ここでは単純にランダムに選択することを推奨する。

Fig. 6(b)はLearning based sparse hashing, Spectral hashing¹⁰⁾, Random projections¹²⁾¹⁶⁾, Sparse random projections¹⁴⁾¹⁵⁾の比較結果を示したものである。Li¹⁵⁾らが W の要素を $1-1/\sqrt{D}$ の確率でゼロに設定することを推奨していることに基づき、Sparse random projectionsではゼロ要素率を $1-\frac{1}{\sqrt{136}} \approx 0.9$ とした。また、Learning based sparse hashingでは2種類のゼロ要素率0.3, 0.9を用いた。Fig. 6(b)から分かる通り、Supervised sparse hashingは高いゼロ要素率であっても、広い範囲のビット長において従来手法よりも高いMAPを達成できていることが分かる。

Fig. 6(c)は縦軸にMAP、横軸にゼロ要素率を示したものである。興味深いことに、Learning based sparse hashingのゼロ要素率を0.9付近まで高めても、MAPが低下しないことが分かった。学習を用いないSparse random projectionsの場合、Li¹⁵⁾らによって同様の知見が既に報告されているが、学習を導入した場合でも同じ効果が認められるようである。従ってゼロ要素率を0.9に設定することで、近傍探索性能を損なうことなく、バイナリコード化を高速に処理できるようになる。これは応用上、非常に重要な利点であるといえる。

4.3 二視点のマッチング性能

Mikolajczyk²¹⁾²²⁾らの方法に従って2視点間のマッチング性能を評価した。Mikolajczykらはマッチング性能評価用のデータベースとして8種類のカテゴリを提供しており、それぞれカテゴリは同一の対象を異なる

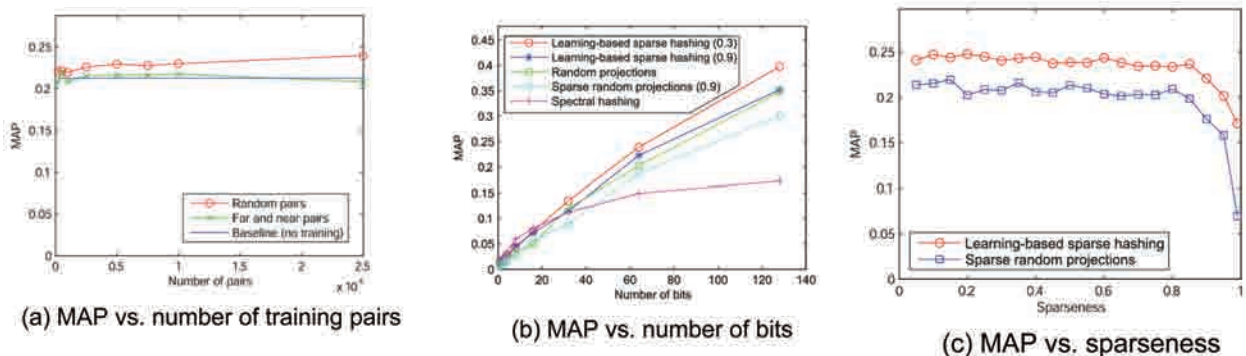


Fig. 6 (a) MAP was improved by increasing the number of randomly chosen training pairs. (b) Different binary hashing functions are compared in this figure. Learning-based sparse hashing was able to improve the performance of sparse random projections. (c) Interestingly, MAPs provided by learning-based sparse hashing were not dropped until sparseness was set to greater than about 0.9.

る条件で撮影した6枚の画像No.1~No.6で構成されている。そして、No.1と残りの画像No.2~No.6との間のHomography行列がそれぞれ与えられている。

評価実験では、No.1とNo.2の画像間の対応点を、最も近い局所特徴量と2番目に近い局所特徴量の距離の比に対してしきい値処理を適用することで求め、Homography行列によってその対応点が正しいか否かを検証した。検証結果に従って1-precision vs. recall曲線をプロットすることにより、局所特徴量の性能を評価した。Wを学習するために、8種類のカテゴリの画像No.4~No.6から局所特徴量を抽出し、ランダムに50,000個の学習サンプルを選定した。ここからランダムに25,000ペアを抽出し、Wを学習した。ゼロ要素率

は0.9とし、異なる3種類のビット長 $B=32,64,128$ について評価した。

Fig. 7に4つカテゴリ *bark*, *boat*, *wall*, *graffiti* についての結果を示した。カテゴリ *bark*, *boat* は回転およびスケール変化の評価用であり、カテゴリ *wall*, *graffiti* は射影変換を含む視点位置変化の評価用として用意されている。実験では、GLOH, SIFT, PCA-SIFT¹⁸⁾, モーメント不変量²¹⁾, CARDを比較した。提案手法であるCARDについては、バイナリコード化を適用した場合としなかった場合の両方について評価した。バイナリコード化の手法としては、学習の効果を確かめるためにLearning based sparse hashingとSparse random projectionsの両方を適用した。いずれの方法において

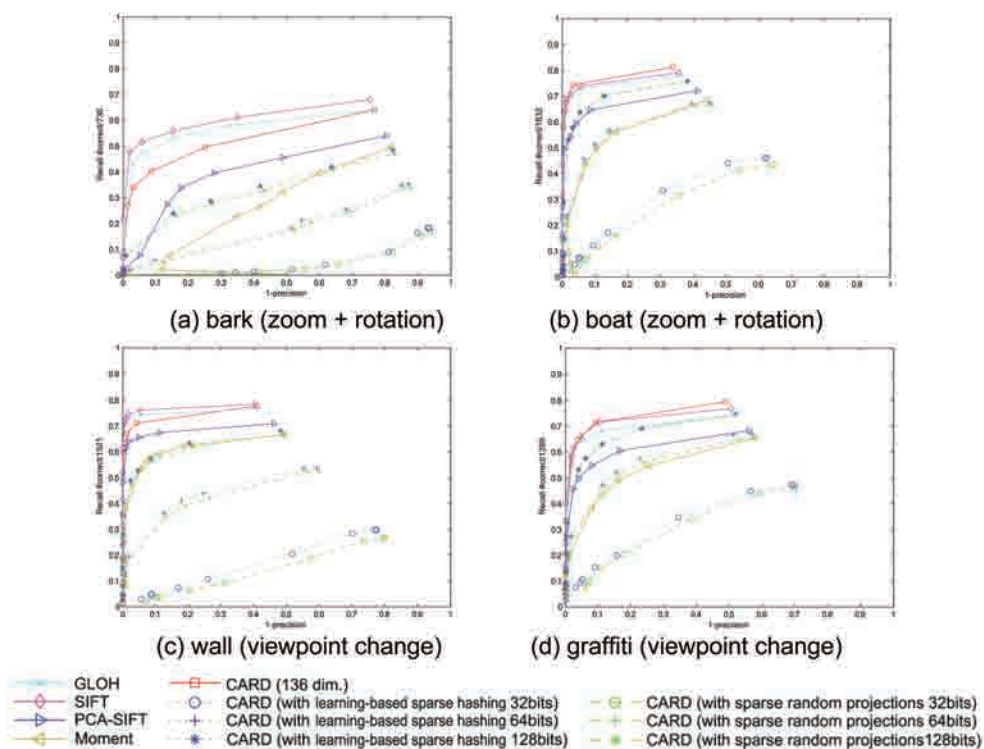


Fig. 7 Different descriptors were compared according to the Mikolajczyk's method.

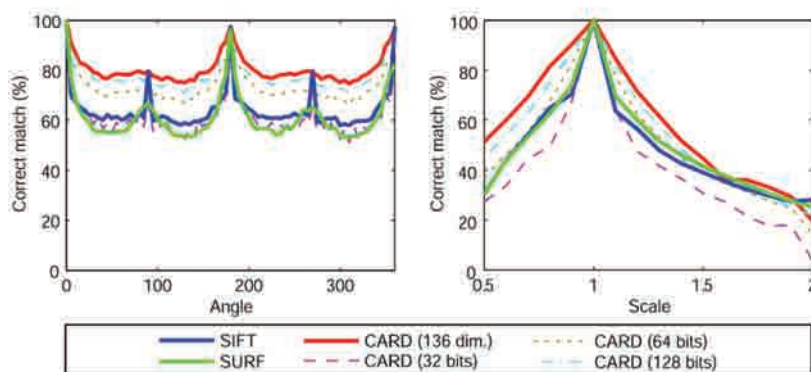


Fig. 8 Scale and orientation invariance of SIFT, SURF, and CARD are compared in this figure.

も、ゼロ要素率として0.9を設定した。これらのバイナリコード化手法では乱数が用いられているため、同じ実験を10回繰り返して得た1-PrecisionとRecallの平均をプロットした。本実験では、いずれの場合においても2.1節で記載したキーポイント検出器を用いた。

バイナリコード化前のCARDの結果を見ると、全てのカテゴリにおいてGLOH, SIFTに匹敵していることが分かる。バイナリコード化を適用した後も、128bitあればモーメント不変量よりも性能が高く、またPCA-SIFTに匹敵する性能が得られていることが分かる。また、同じビット長におけるLearning based sparse hashingとSparse random projectionsの結果を比較すると、 W を学習することによりマッチング性能が改善できていることが読み取れる。

Fig. 8にSIFT, SURF, CARDのスケールおよび回転不変性の性能評価結果を示した。SIFTではDoG, SURFではHarr-like特徴量をキーポイント検出器として適用した。Mikolajczykによるデータベースのカテゴリ *graffiti* の画像No.1を変形させ、回転およびスケール変換させた画像セットを生成した。オリジナルの画像と変形後の画像間において、最近傍となる局所特徴量を探索し、RANSACによるHomography行列推定処理を適用することでアウトライアを除去した。最近傍探索で得られたマッチング数に対するインライア数の割合を求め、これをFig. 8にプロットした。図に示されている通り、提案手法では短いビット長においてもSIFT, SURFに匹敵する性能が得られていることが分かる。

4.4 計算時間

本節では、Intel Core 2 Duo 2.66GHzを用いてCARDの計算時間を評価した結果について述べる。1つの局所特徴量あたりの計算時間をFig. 9に示した。ここでは、キーポイント検出における計算時間は排除して評価している。x軸はCARDの計算時間であり、SIFTの場合を100%として正規化して示してある。なお、SIFTを計算するソフトウェアとしてHessによる実装²⁴⁾を用いた。図中上段は、バイナリコード化手法としてRandom projectionsを用いた結果であり、下段はゼロ要素率を0.9としてLearning based sparse hashingを用いた結果である。ここではビット長 $B=128$ とした。疎行列によって大幅な高速化が達成されていることが分かる。局所特徴量抽出の部分もまた、

ルックアップテーブル化によって大幅に高速化されているため、全体でSIFTの約6%の計算量(約16倍)の高速化を達成できた。

最後に、Table 1にキーポイント検出器を含めた総合的な評価結果をまとめておく。SURFの実装としてはOpenCVを用いた。実験では、640×480の大きさの画像を用いた。5つの異なる画像に対してSIFT, SURF, CARDを適用し、計算時間とキーポイント数の平均値をまとめた。CARDは従来手法よりも高速であり、SIFTの約8倍、SURFの約1.6倍という結果を得た。

Table 1 Computation time of SIFT, SURF, and CARD (640×480 size image).

	SIFT	SURF	CARD (128 bits)
Computation time (msec)	1030.5	206.1	129.3
Number of keypoints	1553.8	1267.8	1658.0

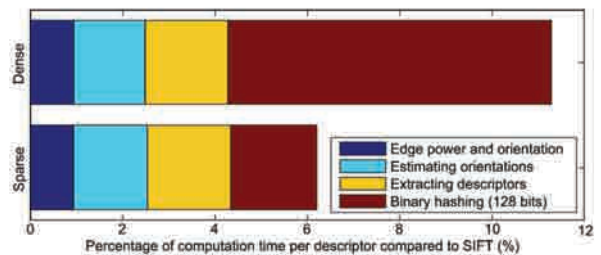


Fig. 9 Percentage of computation time per descriptor compared to SIFT. Computation time of SIFT is set as 100% on the x-axis.

5. おわりに

本論文では、高速に計算可能であり、128ビット程度のバイナリコードで表現可能な局所特徴量CARDを提案した。提案手法は少ない計算量・メモリ量で動作できるため、車載向けCPUや携帯端末などといった、低スペック機器に向いていると考えられる。残る課題としては、Learning based sparse hashingの最適化手法の改善や、アフィン変形への対応などが挙げられる。これらの課題を解決することで、応用範囲を拡大できると期待される。

<参考文献>

- 1) D. G. Lowe: "Distinctive image features from scale invariant keypoints", IJCV, 60, pp. 91-110 (2004).
- 2) N. Snavely, S. M. Seitz and R. Szeliski: "Photo tourism: exploring photo collections in 3D", SIG-GRAPH, pp. 835-846 (2006).
- 3) D. Nister and H. Stewenius: "Scalable recognition

- with a vocabulary tree”, CVPR, pp. 2161-2168 (2006).
- 4) S. Lazebnik, C. Schmid and J. Ponce: “Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories”, CVPR, pp. 2169-2178 (2006).
 - 5) H. Bay, A. Ess, T. Tuytelaars and L. Van Gool: “Speeded-up robust features (SURF)”, *Computer Vision and Image Understanding*, 110, pp. 346-359 (2008).
 - 6) G. Takacs, V. Chandrasekhar, S. Tsai, D. Chen, R. Grzeszczuk and B. Girod: “Unified real-time tracking and recognition with rotation-invariant fast features”, CVPR, pp. 934-941 (2010).
 - 7) D. Wagner, G. Reitmayr, A. Mulloni, T. Drummond and D. Schmalstieg: “Pose tracking from natural features on mobile phones”, *ISMAR*, pp. 125-134 (2008).
 - 8) M. Calonder, V. Lepetit, P. Fua, K. Konolige, J. Bowman and P. Mihelich: “Compact signatures for high-speed interest point description and matching”, *ICCV*, pp. 357-64 (2009).
 - 9) S. Taylor, E. Rosten and T. Drummond: “Robust feature matching in $2.3 \mu s$ ”, *CVPR Workshop*, pp. 15-22 (2009).
 - 10) Y. Weiss, A. Torralba and R. Fergus: “Spectral hashing”, *NIPS*, pp. 1753-1760 (2008).
 - 11) J. Wang, S. Kumar and S-F. Chang: “Semisupervised hashing for scalable image retrieval”, CVPR, pp.3424-3431 (2010).
 - 12) K. Min, L. Yang, J. Wright, L. Wu, X.-S. Hua and Y. Ma: “Compact projection: Simple and efficient near neighbor search with practical memory requirements”, CVPR, pp. 3477-3484 (2010).
 - 13) J. Wang, S. Kumar and S-F. Chang: “Sequential projection learning for hashing with compact codes”, *Proceedings of the international conference on Machine Learning* (2010).
 - 14) D. Achlioptas: “Database-friendly random projections”, *Proceedings of the twentieth ACM SIGMODSIGACT-SIGART symposium on Principles of database systems*, pp. 274-281 (2001).
 - 15) P. Li, T. J. Hastie and K. W. Church: “Very sparse random projections”, *Proceedings of the international conference on Knowledge Discovery and Data mining*, pp. 287-296 (2006).
 - 16) M. X. Goemans and D. P. Williamson: “Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming”, *Journal of the ACM*, 42, pp. 1115-1145 (1995).
 - 17) V. Chandrasekhar, G. Takacs, D. Chen, S. Tsai, R. Grzeszczuk and B. Girod: “CHoG: Compressed histogram of gradients a low bit-rate feature descriptor”, CVPR, pp. 2504-2511 (2009).
 - 18) Y. Ke and R. Sukthankar: “PCA-SIFT: A more distinctive representation for local image descriptors”, CVPR, pp. 506-513 (2004).
 - 19) Y. Dufournaud, C. Schmid and R. Horaud: “Matching images with different resolutions”, CVPR, pp.612-618 (2000).
 - 20) J. Shi and C. Tomasi: “Good features to track”, CVPR, pp. 593-600 (1994).
 - 21) K. Mikolajczyk and C. Schmid: “A performance evaluation of local descriptors”, *PAMI*, 27, pp. 1615-1630 (2005).
 - 22) “Affine covariant regions datasets”, <http://www.robots.ox.ac.uk/vgg/data/data-aff.html>.
 - 23) R. F. L. Fei-Fei and P. Perona: “Learning generative visual models from few training examples: an incremental bayesian approach tested on 101 object categories”, *CVPR Workshop*, p. 178 (2004).
 - 24) R. Hess: “An open-source SIFTLibrary”, *Proceedings of the international conference on Multimedia*, pp. 1493-1496 (2010).

<著 者>



安倍 満
(あんばい みつる)
株式会社デンソーアイティエーラ
ボラトリ 研究開発グループ
博士 (工学)
画像処理を応用した安全支援・
ユーザーインタフェースの研究
開発に従事



吉田 悠一
(よしだ ゆういち)
株式会社デンソーアイティエーラボ
ラトリ 研究開発グループ
画像処理を応用したユーザーイン
タフェースの研究開発に従事