

h p 0 A q` h °` ú . Y Y w h Š w ì í

CoHOG › ã ” w U |

æ y œ — y Z ± y G 5 y ô @ y è y y a y ° à y l y 8 y

y Ê y 0 G ¶ G ¶ ã Ø C J ¶ Z ∈ J β 464-8601 j œ j Ê y 0 ¢ • à Æ Y Ê  
yy 0 Þ \_ ™ ¶ G ¶ & A Ø C ¶ æ β 500-8288 0 Þ j 0 Þ ¢ ¢ Û 1 , è 38 j •

E-mail: ysnakamura@murase.m.is.nagoya-u.ac.jp, f ddeguchi, ide, murase@is.nagoya-u.ac.jp,  
yyttakahashi@gifu.shotoku.ac.jp

K' ± ` Web í w G " w ^ h p , @ p ' X U g b " h Š w O A s A É U [ w ° m q ` o | ^ h p t é l o M " ú .  
› Y Y b " U [ U • [ ' • } H R | i - h p , 0 A q ` h ° ` ú . Y Y x ] œ t Z ∈ U æ ~ • o V h U | ^ h p ›  
0 A q ` h ° ` ú . Y Y w Z ∈ x , , q œ r æ ~ • o M s M } ^ h p , 0 A q ` h ° ` ú . Y Y p x | ^ h p ¢ w 7 ' ,  
s Ñ è " Ü T ' ~ ' • " Y › ã q ^ V › ã w Ò M › @ L \$ t b ; b " \ q U O A q s " } Š C p x | 2 æ U  
Z s r p ô M Q ó , Ê m CoHOG ¢ Co-occurrence Histograms of Oriented Gradients , ì M º t | Á ` h ì í  
CoHOG › ã › Š b " } ì í CoHOG › ã x | ^ h p ¢ w Á t - - p w ì í - M º w ž l Î µ ã - ã Ü p K  
" } Y Y t x | BoF ¢ Bag of Features £ ¯ q q \$ " É ¢ SVM › ; M " } î g p x | Web Í T ' ) B ` h 10 § Á ° æ |  
- 1,000 Š w ^ h p › ; M | í Ó Á Y § ¢ Ñ é " › ã S ' | SIFT › ã q z ± b " \ q t ' " | ì í CoHOG › ã  
w ® Q › - Y ` h }  
© " ë " Á ° ` ú . Y Y | ^ h p | ì í CoHOG › ã | BoF ¯ q

## A Study on Spatio-Temporal CoHOG Features for Recognition of Generic Objects in Video

Shogo NAKAMURA<sup>y</sup>, Daisuke DEGUCHI<sup>y</sup>, Tomokazu TAKAHASHI<sup>yy</sup>,

Ichiro IDE<sup>y</sup>, and Hiroshi MURASE<sup>y</sup>

<sup>y</sup> Graduate School of Information Science, Nagoya University

Furo-cho, Chikusa-ku, Nagoya-shi, Aichi, 464-8601 Japan

<sup>yy</sup> Faculty of Economics and Information, Gifu Shotoku Gakuen University

Nakauzura 1{38, Gifu-shi, Gifu, 500-8288 Japan

E-mail: ysnakamura@murase.m.is.nagoya-u.ac.jp, f ddeguchi, ide, murase@is.nagoya-u.ac.jp,  
yyttakahashi@gifu.shotoku.ac.jp

**Abstract** Recognizing objects in videos is one of the important technologies to search a large amount of videos efficiently on the Web. Recently, generic object recognition has been actively studied for still images, but almost not for videos. As for the generic object recognition in a video, it is important to use both the shape features and the motion features obtained from multiple frames in the video efficiently. In this paper, we propose spatio-temporal CoHOG (Co-occurrence Histograms of Oriented Gradients) features. This is an extension of CoHOG features that provide a high performance for pedestrian detection and others. The spatio-temporal CoHOG features are co-occurrence histograms of oriented spatio-temporal gradients in local regions in a video. In the recognition, a BoF (Bag of Features) representation and a kernel SVM are employed. We conducted an experiment on 1,000 videos including 10 categories collected from the Web. The experimental results showed the effectiveness of the spatio-temporal CoHOG features compared with conventional optical flow features and SIFT features.

**Key words** Generic object recognition, Video, Spatio-temporal CoHOG features, BoF representation

## 1. x a Š t

Ùâ| Web ÍtxG" w^hþU O`|f•'`®p'X  
Ugb"U[U{Š'•oM"}Ugb"hŠw«æw°m  
q`o|^hþtq•"ú.U•['•"}^hþæwú.›  
ÝÝb"\\qUpV•y|â"²xxüw\_hM^hþ›0›  
ts`Zb\\qUpV"}«Qy|x^›Ug«æqb"  
\\qp|\$ 1w'Os^hþ›UgpV"}°`\$t|\\•'  
wUgx^hþtÇ5`h»¬çÄ©µÄ£›;Moæ~•"}  
`T`sU'\\w'Os»¬xâ"²U \$tÇZ"‹w  
pK"hŠ|`G••srUj¼pY`XUgpVsMÔûU  
K"\\|f‹f‹»¬UÇ5`sM‹w‹M}fwhŠ|^h  
þæwú.›-‰;pÝÝb"U[UžAqs"}  
Web Íw^hþt ‡•"ú.x7'pK"hŠ|› wú  
.t'`sMÝÝ OUžApK"}îH„t ‡•"ú.  
›-‰;U°`\$sÊ¶pÝÝb" Oxo`ú.ÝÝqzy  
•|Ùâ|ætz€Uæ~•oM" [1]}°`ú.ÝÝx\$Â°  
æ°p\_hèwìææ"³ãĩUGVMú.›{OhŠ|ôM  
ÝÝp›~"\\qUE`X|r>b,V]JU M}  
HR|i-hþ›0Äq`h°`ú.ÝÝ Ot b"Z  
€U|ætae~•oM"}E`\$s Oq`o| SIFT çScale-  
Invariant Feature Transform ¶2]›;Mh‹wU•['•"}  
SIFT x|hþws8•°!|=t\$HsÁt›ÄpK"\\q  
UÇ'•oS"\\`Z`h›Ä:w\*%¬¬wæ¿'•Ä«µ  
½ßsrw Ý›Ä› 128íiw›ÄÖ«ÄçpG\\b"‹w  
pK"}\\t0`o|^hþ›0Äq`h°`ú.ÝÝtS  
Mox Ý›ÄtCQo|`V›Ä›;M"\\qU ®pK"}  
Klaser 'x| HOG çHistograms of Oriented Gradients ¶3]  
w¬›ì M²t!Á'h›Ä› Š'oM" [4]}\\w›Ä  
x|^hþw ÝØCiZpsXj€Ñè"Ü w^VØC‹  
~"\\qUpV|^hÝÝt ®q^•oM"}°M|ó:w  
›Ä" wzl›q"\\qp0ÄwG\\ó—›²í^d"Z€U  
æ~•oM"}f•'wZ€w°mq`o| Watanabe 'x-  
wžl›b;`h CoHOG çCo-occurrence Histograms of  
Oriented Gradients £› Š'oM" [5]}hþí æw7's  
•" p¬ wzl›q"\\qpóvs ÝU¬qDóqs  
"\\2æ UZt ®q^•oM"}  
ŠC p x|\\•'w›ÄwQí›Ê^ù~dhìí Co-  
HOG ›Ä› Šb"}é.\$tx|HRw CoHOG tSZ"  
¬ qžl›ì M²t!Áb"}fOb"\\qp| Ýq^  
Vt0`oôMG\\¬›~"}  
Žñ| 2...pìí CoHOG ›Äqf•›;MhÝÝ Ot  
mMo\\,"}f`o| 3...pîgwMOqAL\\,\"} 4  
...p x| 3...p\\,hîgwßotCQo|ìí CoHOG ›  
Äq w›Ä›wù`hÔùw@LtmMoÐ\*`hAL\\  
,"}7™t| 5...pŠC ›‰b•}

## 2. ìí CoHOG ›Ä›;MhÝÝ O

2.1 ìí CoHOG ›Ä›;MhÝÝ Ow"A  
\$ 2 t| Owv•›Ôb}‡c|Ö—^hþT'o•à

\$ 1 Web Íw^hþ

\$ 2 ìí CoHOG ›Ä›;MhÝÝ Owv•

›`Z`|f\\T'ìí CoHOG ›Ä›`Zb"}o•  
à x|È`hó:Ñè"ÜT's"\\ÝÝ›æOo•p  
K"}f`o|o•à T'`Z^•h›Ä› BoF çBag of  
Features¶6]¬q`|§Â°æ¬t;™`h 2«âµÝ +›  
;Mo0Ä§Â°ætb"ôTS›‰Zb"}7™t|¶o  
•à wôTST'|Ö—^hþw0Ä§Â°ætb"ôT  
S››b"}Žñp|µ`ŠtmMoÄ`X\\,"}  
2.2 ìí CoHOG ›Äw¬Z  
o•à T'°`pìí Át¬¬çÒé¿«£›±  
ĩÓæĩ¬`|Òé¿«¬tìí CoHOG ›Ä›`Zb"}  
\$ 3 t|Òé¿«w±ĩÓæĩ¬w7›Ôb}Ž<pìí  
CoHOG ›Ä"w-‰M OtmMo\\,"}  
‡c|Òé¿«°w¶hÉpìí ¬ g(x;y;t) ›-‰b  
"}l(x;y;t)›hÉwKSq`o| g(x;y;t) wæAE'Íw'  
Ot{Š"}  
8  
≥ g<sub>x</sub>(x;y;t)=l(x+1;y;t) l(x-1;y;t)  
≥ g<sub>y</sub>(x;y;t)=l(x;y+1;t) l(x;y-1;t) (1)  
≥ g<sub>t</sub>(x;y;t)=l(x;y;t+1) l(x;y;t-1)

Ít|Òé¿«›ó:w¬¬çç£tũÄ`|µ·çp É¬  
g(x;y;t) ›-‰b"}f`o|µ É¬ g(x;y;t) ›Y 20

\$ 3 o•à T'wÖé¿ «w±ĩÓæĩ¬

\$ 4 Y 20 Ø.¢\$w:Èx" =M²>¬b£

Ø.wÖ:pK"p<sub>12</sub>M²t" =b"}Y 20 Ø.wÖ:w2  
<sup>a</sup>x| ' =  $\frac{1+}{2}$  <sup>5</sup> q`o ( 1; ' ; 0);(0; 1; ' );( ' ; 0; 1)  
 pK"}\$ 4 t|Y 20 Ø.›Ôb}

Òé¿ «°w±.çp-%`h¬ M²>b;`o|¬ M  
 ²wžlîµÄ¬âÜ›^Rb"}\•Uîí CoHOG ›Ä  
 pK"}Ž<p|\$ 5 t lo^Rwv•›†ìb"}  
 ¢a£ ‡c|,j.çqw•" t"|" 63 "wÖž›  
 Rb"¢hi`|,j.çx ‹ %£sS| ² tÖž  
 › R`sMwx|0¶Qt"ÑÖQ› ‡b"hŠpK"}  
 ¢b£ ¢Öžw•" tK"·çwÊ^ù~d›Òé¿ «  
 °T'»•`|fw¬ M²wÖž›îµÄ¬âÜtd®°o  
 MX}

¢c£ ¢ÖžwîµÄ¬âÜ›ÈAb"\qp74\$ŝîµ  
 Ä¬âÜqb"}hi`|,j.ç%œwÖžxfw·çx  
 w¬ M²w^pîµÄ¬âÜ›^Rb"}‡h|¬ \$SU  
 0ts`M·çxb;`sM}

îí CoHOG ›Ä"x (12 12 62)+12=8,940 îí  
 qs"|‡xtîiUGVXs"}fwhŠ| PCA ¢Principal  
 Component Analysis £t" d íityVb"}  
 2.3 îí CoHOG ›Äw BoF ¬q

ÝÝ›æO²rgq`o|o•à w¶Öé¿ «p^R`h  
 îí CoHOG ›Ä› BoF ¬qb"} BoF ¬qxhþ›Át  
 ›Ä"wîµÄ¬âÜp¬qb" OpK"|°`ú.ÝÝw  
 üúp¿X;M'•oM"}\w Ox¶6^ŠqÝÝ^Št  
 üT•oS"|°`tŽ<w qpæO}

¶6^Špx|‡c¶¶6;hþT'ó:wÁt›Ä"›  
 Zb"}f`o|¶¶6;hþw¶Át›Ä"› k-means «  
 âµ»æĩ¬b"\qt"|" visual word ›\Rb"} visual  
 word x|Át›Ä"›¬b›ÄÖ «Äç›Ö «Äç" =`  
 h‹wpK"}±Át›Ä"›7‹"Äb" visual word q`

\$ 5 îí CoHOG ›Äw¬Z

o¬qb"\qt"|"¶¶6;hþx visual word wZqÄ  
 SwîµÄ¬âÜpG\^•"}ÝÝ^Špx|‡c¶6^Š  
 q%7t|¶Ö—;hþT'Át›Ä"›¬Zb"}f`o|  
 ¶6^Šp\R^•h visual word ›;Mo|¶Ö—;hþ›  
 visual word wZqÄSwîµÄ¬âÜp¬qb"}  
 2.4 Ý +w¶6

¶6^hþwo•à p^R`hîí CoHOG ›Äw BoF  
 ¬q›;Mo|Ý +›¶6b"}sS|Ý +x\$Â°æ¬  
 t;™°oSX}Ý +tx|§"Éç SVM ¢Support Vector  
 Machine [7] ›;M"}§"Éç SVM xôMÝÝQó›ÈI  
 oM"\qUÆ'•|7'shþÝÝðJt ;^•oM"}  
 ŠZ€px|§"Éç :q`o ² §"Éç›;M"} ²  
 §"ÉçxŽ<wÜp¬^•"}  
 k(x;y)=exp  $\sum_i \frac{(x_i - y_i)^2}{x_i + y_i}$  ! (2)

² §"Éçxhþü¬tSMo7‹QóU'M§"ÉçpK  
 "qC ^•oM" [8]} tx|¶¶6Ö «ÄçwÊ^ù~d  
 tSZ" É ² 'mwo:›f b"}  
 2.5 o•à wÝÝ

Ä²t¶6°oSMhÝ +›;Mo|o•à wÝÝ›æ  
 O}Ý +wZ—ALq`ox|ümo ØT'wÖøÇV'  
 mU¬'•"}\w‹0Ä\$Â°ætb"ôTSqb"}  
 2.6 w ü

o•à t wôTS ç›;Mo|Ö—^hþwÝÝ›æO}  
 Ö—^hþwo•à :› T qb"q|74\$ŝôTS ¢xí  
 Üt'lo>b"}  
 ¢ =  $\frac{1}{T} \sum_{t=1}^T \alpha$  (3)

^hþtéloM"ú.›ÝÝ`hMÔù| ¢>0qs"\$  
 Â°æ›ÝÝALqb"}hi`|éloM"ú.U 1 "q  
 MOÚEs'y| ¢U7Gw\$Â°æ›ÝÝALqb"}‡h|  
 «±æ›)Qo^hþ›Ug`hMÔùx| ¢›fw‡‡b;  
 b"}  
 | 3 |

(a) bicycle+person

(b) bus

(c) car

(d) cat

(e) cow

(f) dog

(g) horse

(h) motorbike+person

(i) person

(j) sheep

\$ 6 í g p - ; ` h ^ h p w «

### 3. í g

#### 3.1 í g Ä " » . ç Ä

Ä " » . ç Ä x | YouTube [9] T ' ^ h p » ) B ` o i T M ` h } ^ h p w ° ` ú . w § Ä ° æ x | PASCAL Visual Object Classes Challenge 2006 [10] p - ; ^ • h Ä " » . ç Ä q % ° w 10 § Ä ° æ q ` h } é . \$ t x | bicycle, bus, car, cat, cow, dog, horse, motorbike, person, shepp K " } \$ 6 t Ä " » . ç Ä w ° æ } Ö b } Ä " » . ç Ä w ^ h p x 0 Ä ú . U é " ' O t ~ " Z ` h } ` T ` | ! « ç " ' ä i • h Ø Ž • w x ^ Z ` p ú . w ° æ U é l o M s M Ñ è " Ü U O ` h } ‡ h | Ÿ s " § Ä ° æ w ú . U % ò t é l o M " ^ h p ( ) B ` h } « Q y | bicycle, motorbike U ‡ • " ^ h p t m M o x | b , o person « % ò t é l o M " « w p K l h } f w h Š | \ • ' w ^ h p » person w Ý + t ¶ 6 ^ d " M x | b , o Ü ' Ä Ÿ Ö Ä " » q ` o { l h } Ä " » . ç Ä w ^ h p : x 1,000 Š p K " | ¶ § Ä ° æ 100 Š ¢ person x bicycle, motorbike q % ò t é l o M " ^ h p « Š " q 300 Š £ q ` h } ‡ h | ^ h p w Ñ è " Ü ± ¶ x 1,280 720 T M , 920 1,080 [pixels] | Ñ è " Ü : x 45 T M 270 p K l h }

#### 3.2 í g M O

Š b " í CoHOG › Ä ; M h O q | ! Ó Ä Ÿ § ç

Ñ é " › Ä | SIFT › Ä ; M " Ž < w 2 m w O › z ± ` h } y z ± O 1 • Ý Ý O x | 2.2 ... › ! Ó Ä Ÿ § ç Ñ é " › Ä t " V ö Q | 2.3 ... w B o F - q › æ ~ s M « w q ` h } ! Ó Ä Ÿ § ç Ñ é " › Ä q ` o x | i € Ñ è " Ü w ! Ó Ä Ÿ § ç Ñ é " › ; M h } Ñ é " w " Z t x Ö é ç « Ü ç ½ i - O › b ; ` h } ¶ Ñ è " Ü w ¶ Ö é ç « t 0 ` o Ñ é " › - % ° ` | o • à p Ñ é " w M ² î µ Ä - ä Ü › ^ R ` h } y z ± O 2 • Ý Ý O x | 2.2 ... › SIFT › Ä t " V ö Q h « w q ` h } SIFT › Ä : w U Z t x | DoG ¢ Di erence of Gaussian £ ; M h sparse s O q - æ ç Ä ± i Ó æ i - • ä i ¼ Ü ± i Ó æ i - ; M h dense s O U K " U | Š Z € p x ° ` ú . Ý Ý p " " @ q ^ • " T M w dense s O ; M h } µ - " ç ! = t \$ H s › Ä q b " h Š | › Ä : U Z x ó : µ - " ç w - æ ç Ä ± i Ó æ i - t " " æ l h } ¶ O ; M " M w í ä Ý " » x | í w " O t f ` h } í CoHOG › Ä t m M o x | o • à w Ñ è " Ü : › f = 5 | Ö é ç « ± ¶ › w B = h B = 80 | • ç ± ¶ › w C = h C = 5, f C = 1 | PCA & ; T M w í i : › d = 500 q ` h } ! Ó Ä Ÿ § ç Ñ é " › Ä t m M o x | Ö é ç « ± ¶ › 20 20 | M ² î µ Ä - ä Ü w i i : › 36 q ` h } SIFT › Ä t m M o x | ± i Ó æ i - ' › 80 | µ - " ç › 8, 16, 24, 32 w 4 " q ` h } ! í CoHOG › Ä q SIFT › Ä w B o F t s Z " visual word w : x 500 q ` h }



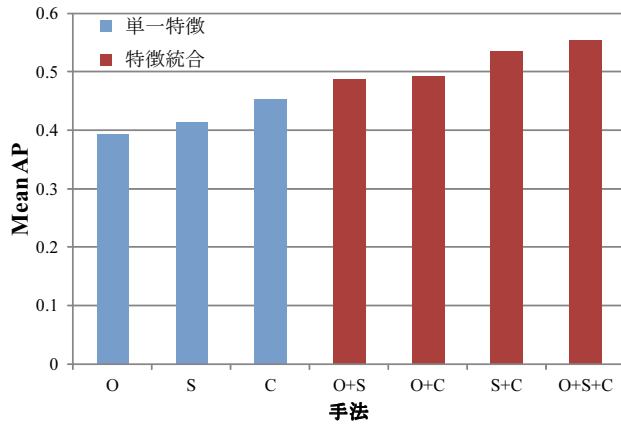


図 8 各特徴を用いた手法および特徴統合による手法の Mean AP (O はオプティカルフロー特徴, S は SIFT 特徴, C は時空間 CoHOG 特徴を表す)

で統合カーネルを作成する手法である．本来 MKL はカーネル選択のために作られたが, Varma らは MKL における各カーネルを特徴に対応付けて統合することで画像認識を行う手法を提案した [12]．そこで, 本節でもカーネル選択ではなく特徴統合のために MKL を用いる．統合カーネル  $K(\mathbf{x}, \mathbf{y})$  は特徴  $f$  のカーネル  $k_f(\mathbf{x}_f, \mathbf{y}_f)$  から次式のように作成できる．

$$K(\mathbf{x}, \mathbf{y}) = \sum_{f=1}^F \beta_f k_f(\mathbf{x}_f, \mathbf{y}_f) \quad (5)$$

$F$  は特徴数,  $\beta_f$  は特徴  $f$  の重みを表す ( $\beta_f \geq 0, \sum_{f=1}^F \beta_f = 1$ )．また, 各特徴のカーネル  $k_f(\mathbf{x}_f, \mathbf{y}_f)$  には, 2 節同様  $\chi^2$  カーネルを用いる．

#### 4.2.2 特徴の統合実験

特徴統合の有効性を調査するため, 各特徴のみを用いる手法と特徴統合を行う手法を比較した．実験条件は基本的に前節と同じであり, 各特徴  $f$  のカーネル重み  $\beta_f$  は均一重み ( $\beta_f = \frac{1}{F}$ ) とした．

各特徴を用いた手法および特徴統合による手法の Mean AP を図 8 に示す．この図から, 単一特徴を用いるよりも特徴統合を行う方が認識精度が高くなっていることがわかる．このことから, 特徴統合の効果を確認した．また, 特徴統合の中でも 3 種類すべての特徴を統合する手法が最も認識精度が良かった．すなわち, 時空間 CoHOG 特徴は形状と動きの双方を記述するものの, SIFT 特徴とオプティカルフロー特徴には認識に有効な情報が他にあることがわかる．特に, オプティカルフロー特徴よりも SIFT 特徴と統合した方が認識精度が向上することから, 時空間 CoHOG 特徴に形状情報が不足していると考えられる．そのため, 空間方向をより重視した勾配の量子化が必要である．

## 5. む す び

本報告では, 動画を対象とした一般物体認識のための時空間 CoHOG 特徴量を提案した．時空間 CoHOG 特徴は CoHOG における勾配と共起を時間方向に拡張したものである．単位区

間から一定間隔で時空間 CoHOG 特徴を抽出し, PCA により次元を圧縮し, BoF 表現することで単位区間を記述した．そして, カーネル SVM による認識を単位区間毎に行い, それらの結果を平均することにより, 入力動画の認識を行った．実験では, YouTube において収集した 1,000 本の動画像を使用し, 時空間 CoHOG 特徴とオプティカルフロー特徴, SIFT 特徴を比較した．オプティカルフロー特徴は, 各隣接フレーム間のオプティカルフローのヒストグラムを作成することによって抽出した．SIFT 特徴は, 各フレームからグリッドサンプリングにより特徴点を抽出し, BoF 表現することによって抽出した．実験の結果, 時空間 CoHOG 特徴を用いた手法が最も高い認識精度となることを確認した．

今後の課題としては, より離れた時間間隔 徳 久 藤 驪 産 き ち