## Detection and representation for image understanding

#### Matthieu Cord

Laboratoire d'Informatique de Paris 6 (LIP6) Université Pierre et Marie Curie (UPMC)

October, 2011



## Outline

- 1
- Introduction
- Multimedia Group
- Projets
- 2 Text Detection in Natural Images
  - Motivation
  - Related Works
  - Hypotheses Generation
  - SnooperText
  - Hypothesis Validation
  - Experiments
- Biologically inspired image classification
  - Image classification: state-of-the-art
  - Proposed model
  - Results
  - Conclusion

< 3 > < 3 >

# LIP6/DAPA/MALIRE

#### People

- LIP6 lab in Paris
  - $\bullet~\sim$  200 permanent researchers,  $\sim$  250 Phd students
- ② DAPA department: Databases and Machine learning
  - $\sim$  35 permanent researchers,  $\sim$  50 Phd students
- **③** MALIRE team: MAchine Learning and Information REtrieval
  - $\sim$  20 permanent researchers,  $\sim$  35 Phd students,  $\sim$  10 Post-docs
- Our research group
  - $\bullet~2$  permanent researchers (M. Cord, N.Thome),  $\sim~10$  Phd/Post-docs

イロト 不得 トイヨト イヨト

# LIP6/DAPA/MALIRE

#### **Research Topics: image representations and similarities**

- Image Representations:
  - Deep learning for image understanding
  - Computer-vision-based image representations (beyond BoW)
- Kernel methods
  - Similarity design (image search, actor video retrieval)
  - Similarity learning (feature combination)
- Content-based image and video retrieval systems in huge multimedia databases
- Active/Interactive strategies for image and video retrieval
- Object (Text) Detection with context
- Hybrid strategies for web archiving

イロト 不得下 イヨト イヨト

# Outline



- Multimedia Group
- Projets
- 2 Text Detection in Natural Images
- Biologically inspired image classification

- 4 週 ト - 4 三 ト - 4 三 ト

## Main Projects

#### Retrieval, object detection

#### • ANR ITowns: street digitalization, visu., understand., search ...



## Main Projects

#### Retrieval, object detection

• ANR ITowns: street digitalization, visu., understand., search ...



## Main Projects

#### Retrieval, object detection

#### • ANR ITowns: street digitalization, visu., understand., search ...



Matthieu.cord@lip6.fr

Detection and representation for image understanding

## Main Projects

#### Retrieval, object detection

#### • ANR ITowns: street digitalization, visu., understand., search ...



- 4 回 ト - 4 回 ト

## Main Projects

#### Retrieval, object detection

• ANR ITowns: street digitalization, visu., understand., search ...



#### **Deep Learning**

- ANR ASAP: french consortium for learning deep representations
- Bilateral project MERLION: partnership with IPAL, Singapore

(人間) トイヨト イヨト

## Main Projects

#### Retrieval, object detection

• ANR ITowns: street digitalization, visu., understand., search ...



#### **Deep Learning**

- ANR ASAP: french consortium for learning deep representations
- Bilateral project MERLION: partnership with IPAL, Singapore

#### Image and video analysis

• French-Bresilian CAPES-COFECUB Project (S. Avila and R. Minetto's Phd)

Matthieu.cord@lip6.fr

# Main Projects

#### **Object detection**

- ANR Geopeuple: old Maps Analysis (EHESS, COGIT/IGN)
- ullet Interpreting Population Evolution: from 18th century  $\rightarrow$  nowdays

10rme la Baraque 1 + uena de ta Touraudrie

#### Contextual object detection

## Main Projects

#### Scalable data migration preservation

- European Project (IP) SCAPE: SCAlable Preservation Environments
- LIP6 (BD/MALIRE): automation on the web page versioning
  - Combining structural and visual feature



## Main Projects

#### **SCAPE: SCAlable Preservation Environments**



significant change  $\rightarrow$  generate a new versionning

- LIP6 (BD/MALIRE): Combining structural and visual feature
  - Metric/kernel learning (MS student + PhD)

Matthieu.cord@lip6.fr

## Outline

#### Introduction

#### Text Detection in Natural Images

- Motivation
- Related Works
- Hypotheses Generation
- SnooperText
- Hypothesis Validation
- Experiments



通 ト イヨ ト イヨト

## Introduction

- Text detection is a challenging task in computer vision;
- Existing approaches are dedicated to specific contexts;
- Text detection in urban scenes is hard:
  - $\rightarrow$  Font variations;
  - $\rightarrow$  Strong background clutter;
  - $\rightarrow$  Natural noise;
  - $\rightarrow$  Perspective distortion, blurring, illumination changes, etc;
- State of the art OCR's fail in urban scenes images;

Co-supervision PhD student R. Minetto with Prof. J. Stolfi, University of Campinas (UNICAMP), Brazil Based on a strategy developed with CMM in ANR itowns

Matthieu.cord@lip6.fr

Detection and representation for image understanding

#### Motivation

## Motivation



Matthieu.cord@lip6.fr

#### Detection and representation for image understanding

3

(日) (周) (三) (三)

## Motivation



#### **Related Works**

- Hinnerk Becker [1] (Bottom-up approach)
- Alex Chen et al. [1] (Top-down approach)
- Epshtein et al. [2] (Bottom-up approach)
- Chen et al. [3] (Bottom-up approach)

[1] S.M. Lucas. Text Locating Competition Results. ICDAR 2005.

[2] Boris Epshtein, Eyal Ofek and Yonatan Wexler. **Detecting Text in Natural Scenes** with Stroke Width Transform. CVPR 2010.

[3] Huizhong Chen, Sam S. Tsai, Georg Schroth, David M. Chen, Radek Grzeszczuk and Bernd Girod. Robust Text Detection in Natural Images with Edge-enhanced Maximally Stable Extremal Regions. ICIP 2011.

イロト イポト イヨト イヨト

- Image segmentation:
  - $\rightarrow$  Toggle mapping
- Character classification:
  - $\rightarrow$  Rotation invariant image descriptors
- Character grouping:
  - ightarrow Geometric criteria
- Multi-resolution



#### Mono-resolution v.s. Multi-resolution segmentation

- Coarser levels:
  - $\rightarrow$  detects large text areas
  - $\rightarrow$  ignores texture details
- Finer levels:
  - $\rightarrow$  detects small regions
  - $\rightarrow$  analyses more accurately the local image content













d) Multi=resolution

Detection and representation for image understanding

 $\exists \rightarrow$ 

#### Problem

- Analyzes the image content locally
  - $\rightarrow$  Prone to false positives



## Generation/validation process: SnooperText



- [Minetto2010] extends previous work [Fabrizio2009]
- Hybrid scheme: hypothesis generation/validation paradigm
  - Hypothesis generation: multiresolution bottom-up approach
     → improves segmentation robustness over [*Fabrizio*2009]
  - Hypothesis validation: top-down strategy
    - $\rightarrow$  To remove false positives by analyzing globally the window content

[Minetto2010] R. Minetto, N. Thome, M. Cord, J. Fabrizio, B. Marcotegui, SnooperText: A Multiresolution System for Text Detection in Complex Visual Scenes, ICIP 2010. [Fabrizzio2009] J. Fabrizio, B. Marcotegui, and M. Cord, text segmentation in natural scenes using togglemapping. ICIP 09.

# Hypothesis Validation

#### **Fuzzy HOG**

- Idea: analyze each candidate text region globally
- Fuzzy HOG: a global HOG descriptor with different weight masks
- Eliminate the regions with non "text-like" periodical patterns



# Hypothesis Validation

## Fuzzy HOG

- Idea: analyze each candidate text region globally
- Fuzzy HOG: a global HOG descriptor with different weight masks
- Eliminate the regions with non "text-like" periodical patterns



#### HOG idea

- Images of complex objects typically have different HOG's in different parts;
- Humans:
  - $\rightarrow$  different gradient orientation distributions in the head, torso, legs, etc;



Figure: Image from: Histograms of Oriented Gradients for Human Detection. Navneet Dalal and Bill Triggs. CVPR 2004

## HOG of some isolated letters



◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□ ◆ ○ ◆

#### Text HOG idea

- Text-lines of Roman letters: ≠ HOG's in the top, middle and bottom parts:
   → The image is divided into an array of cells with one HOG to each cell;
- Top and bottom parts: Large proportion of horizontal strokes
   → gradients pointing mostly in the vertical direction;
- Middle part: Large proportion of vertical strokes
   → gradients pointing mostly in the horizontal direction;
- All parts: Amall amount of diagonal strokes
- The concanetation of the 3 HOG's is the descriptor of the full region.



Figure: Top, middle and bottom HOGs for the text "RECOGNITION". The arrows show the contribution of specific letters strokes to the final descriptor.

200

#### Sharp cells

- Cells defined by sharp boundaries:
  - $\rightarrow$  HOG may change with small vertical displacements



Wo



*w*<sub>2</sub>



#### **Fuzzy cells**

• To avoid this problem, we used "fuzzy" cells :



w<sub>0</sub>







#### Dalal et al. masks to human recognition

- Gaussian weight functions:
  - $\rightarrow$  Problem: Sharp boundaries.



Figure: Weight functions for a single block of  $1 \times 3$  cells ( $\sigma_x = W/2$ ,  $\sigma_y = H/2$ ).



Figure: Weight functions for a single block of  $1 \times 3$  cells ( $\sigma_x = W/4$ ,  $\sigma_y = H/4$ ).



Figure: Weight functions for  $1 \times 3$  single-cell blocks. Each with height H/2 and overlapped with stride H/4 ( $\sigma_x = W/4$ ,  $\sigma_y = H/8$ ).

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

## Text HOG descriptor scheme



## F-HOG of text and non-text regions



#### Dataset

- 499 color images (training/testing)
- Captured with different digital cameras and resolutions
- Images from indoor and outdoor scenes
- Groundtruth available (XML)













# MetricsPrecisionRecallRanking $p = \frac{\sum_{r_e \in E} m(r_e, T)}{|E|}$ $r = \frac{\sum_{r_t \in T} m(r_t, E)}{|T|}$ $f = \frac{1}{\alpha/p + (1 - \alpha)/r}$

• m(r, R): best match for a rectangle r in a set of rectangles R.

< □ > < @ > < 注 > < 注 > ... 注

- T: set of manually identified text regions (groundtruth);
- E: set of text regions reported by the detector;
- *f*: harmonic mean of precision and recall ( $\alpha = 0.5$ )

#### Performances results

System	Precision (p)	Recall (r)	f
Our System	0.73	0.61	0.67
Epshtein et al. (CVPR 2010)	0.73	0.60	0.66
Chen et al. (ICIP 2011)	0.73	0.60	0.66
SnooperText (ICIP 2010)	0.63	0.61	0.61
Hinnerk Becker <sup>1</sup>	0.62	0.67	0.62
Alex Chen	0.60	0.60	0.58
Ashida	0.55	0.46	0.50
HWDavid	0.44	0.46	0.45
Wolf	0.30	0.44	0.35
Qiang Zhu	0.33	0.40	0.33
Jisoo Kim	0.22	0.28	0.22
Nobuo Ezaki	0.18	0.36	0.22
Todoran	0.19	0.18	0.18
Full	0.01	0.06	0.08

- イロン スピン スピン スピン

#### Successfull detections



p = 0.96, r = 0.64, f = 0.77



p = 0.90, r = 0.90, f = 0.90



p = 0.93, r = 0.93, f = 0.93



p = 0.68, r = 0.56, f = 0.61

#### **Failures**



#### p = 0.00, r = 0.00, f = 0.00





## p = 0.17, r = 0.17, f = 0.17



p = 0.00, r = 0.00, f = 0.00 p = 0.71, r = 0.95, f = 0.81

## iTowns

#### Performances

- ICDAR metrics;
- Text Detection + F-HOG: precision improvement of 23%

System	Precision (p)	Recall (r)	f
Our System	0.69	0.49	0.55
SnooperText (ICIP 2010)	0.46	0.49	0.48

◆□▶ ◆□▶ ◆注▶ ◆注▶ 注 のへで

## iTowns - Detection results



## iTowns - Detection results



◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ● ●

## iTowns - Detection results





# itowns KeyWord Search

Jonathan Guyomard, Frederic Precioso, Nicolas Thome, Matthieu Cord Text detection + OCR (Tesseract)

- Textual query to image retrieval;
- World matching by Edit distance.



## itowns KeyWord Search



▲□▶ ▲圖▶ ▲匡▶ ▲匡▶ 三臣 - のへで

## itowns KeyWord Search



# SnooperTrack: Extension to videos

## SnooperText: Conclusion

- Combines bottom-up & top-down mechanisms
- Efficient in various contexts: urban images, standard databases
- Computational time may make approach difficult to scale up: 640  $\times$  480 pixel images  $\sim$  1 minute

## SnooperTrack: Motivations

- Combining detection & tracking:
  - Speedup text detection in image sequences
  - Discard false positives
  - Improves detection accuracy
- Detection: SnooperText
- Tracking: Particle Filtering (HoG)
- Merging detection & tracking with a combination of position, size and appearance features

Experiments

## SnooperTrack: Results





#### $Loading \ ./images/text/trackingtext.avi$

[*MinettolCIP*11] Rodrigo Minetto, Nicolas Thome, Matthieu Cord, Neucimar Leite, Jorge Stolfi, SnooperTrack: Text Detection and Tracking for Outdoor Videos, ICIP 2011

Matthieu.cord@lip6.fr

Detection and representation for image understanding

## Outline

#### Introduction

Text Detection in Natural Images

#### Biologically inspired image classification

- Image classification: state-of-the-art
- Proposed model
- Results
- Conclusion

通 ト イヨ ト イヨト

## Image classification: state-of-the-art

#### State of the art Model: Bag-of-words (BOW)



#### Credit: Prof. Shih-Fu Chang

Some recent improvements:

- Spatial Information [Lazebnik06]
- Finer coding of local descriptors: soft-assignement, sparse coding [*Wang*10]
- sum v.s. max pooling [Boureau10]
- gives state of the art classification performances

[Lazebnik06] P.Lazebnik.S, Schmid.C. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories CVPR2006.

[Bourgeau10]Y.-L. Boureau, F. Bach, Y. LeCun, and J. Ponce. Learning mid-level features for recognition CVPR2010. [Wang10]J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong.Locality-constrained linear coding for image classification CVPR2010.

Matthieu.cord@lip6.fr

Detection and representation for image understanding

## Image classification: state-of-the-art

#### **Deep Networks**

• Convolutional networks: [LeCun98], improvements [Jarrett09, Lee09]



- $\ominus$  Learning: training lower layers hard

[Jarrett09]K. Jarrett, K. Kavukcuoglu, M. Ranzato, and Y. LeCun.What is the best multi-stage architecture for object recognition? In Proc ICCV2009.

[Lee09]H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. ICML2009

## Image classification: state-of-the-art

## Biologically-inspired Methods [Fidler08, Serre07, Mutch08]



- Mimics feedforward properties of primate visual cortex V1 simple cells
- Based on the HMAX model [Serre07, Mutch08]
  - $\bullet \ \oplus \ \mathsf{Deep} \ \mathsf{models}$
  - $\oplus$  Trainable with real images

[Fidler08]S. Fidler, B. Boben, and A. Leonardis. Similarity-based cross-layered hierarchical representation for object categorization CVPR2008.

[Serre07] T. Serre, et.al, Robust object recognition with cortex-like mechanisms, PAMI, 2007.

[*Mutch0*8] Mutch.J and Lowe.D.G, Object class recognition and localization using sparse features with limited receptive fields, IJCV, 2008

Matthieu.cord@lip6.fr

Detection and representation for image understanding

## **Bio framework**

## HMAX-like architecture [Serre07, Mutch08]



 [Serre07] T. Serre, et.al, Robust object recognition with cortex-like mechanisms, PAMI, 2007.

 [Mutch08] Mutch.J and Lowe.D.G, Object class recognition and localization using sparse features with limited receptive fields,

 Matthieu.cord@lip6.fr
 Detection and representation for image understanding

 47 / 60

## Contribution: Deep multiscale high level filters

• High level filters (yellow) should fit complex visual structures with multiple local scales (red)



Joint work with Christian Theriault (Post Doc) and Nicolas Thome LIP6

Matthieu.cord@lip6.fr

Detection and representation for image understanding

## Proposed model

#### Learning of deep multiscale high level filters



Learning deep scale filters

#### Matthieu.cord@lip6.fr

#### Detection and representation for image understanding

## Proposed model



<u>व</u> • २००

# Proposed model

#### Main parametrization

The scale depth S of high level filter

- $S \uparrow \Longrightarrow \Delta \downarrow$  (More fitted, less scale invariance)
- $S\downarrow \Longrightarrow \Delta\uparrow$  (Less fitted, more scale invariance)



⇒ Learn filters with different values of  $S \in \{1, 2, 3, 4, 5, 6, 7\}$  to get both invariance and discrimination

## Experiments

#### Dataset: Caltech101

- 9,144 images
- 102 categories
- 30 to 800 images per categorie



- Standard evaluation protocol for baseline comparison:
  - Train with 15-30 images / class
  - Test on the reamaining images
  - Metric: Multi-class Accuracy

#### Results

## Experiments

## Caltech 101: Deep multiscale high level filters

#### Deep multiscale high level filters better fit visual content



#### Detection and representation for image understanding

## Classification results

#### Average accuracy results

Setup:

- 8 scales, 12 orientations
- 4080 S2 filters
- One-Against-All gaussian SVM

Model	S	15 images	30 images
[Mutch08]	1	48	54
Our model	3	56.4%	62.5%
Our model	7	55.6 %	62.1%
Our model	1-7	59.1 $\pm$ 0.2 %	$\textbf{66.9} \pm \textbf{0.8\%}$

A B A A B A

#### Average accuracy results

Model	15 images	30 images
Our model	59.1 $\pm$ 0.2 %	$\textbf{66.9} \pm \textbf{0.8\%}$

Deep biologically inspired architectures			
[Mutch08]	48	54	
[Ranzato&alCVPR07]	-	54	
[Kavukcuoglu&alNIPS10]	-	66.3	
[Mutch09]	57.7	64.4	
[ <i>Jarrett</i> 09]	-	65.6	
[Zeiler&alCVPR2010]	58.6	66.9	
[Fidler08]	60.5	66.5	
Shallow architectures			
[Lazebnik06]	56.4	64.6	
[Zhang&alCVPR06]	59.1	62.2	
[Wang10]	64.43	73.44	

## Integrate sparse code learning into the architecture

- We proposed an extension of the HMAX model by learning high level filters with deep scale range
- Our classification results indicates that these more fitted filters, combined with more invariant shallow filters, increase classification scores by nearly 12%
- Our next step is to include filter learning with sparse constraints: Fergus (2010), Olshausen (1996,2009), Ranzato(2007), Kavukcuoglu&Lecun (2010)

#### Thank you for your attention !

## **QUESTIONS** ?

#### People

Matthieu Cord, Nicolas Thome LIP6, Univ. UPMC-PARIS VI matthieu.cord@lip6.fr

- PhD students: Sandra Avila, Rodrigo Minetto, Hanlin Goh, Mar Law, Denis Pitzalis
- Post-Docs: Christian Theriault
- Research Inge. J. Guyomard, C. Sureda

 $http://webia.lip6.fr/{\sim}cord/$ 



## Text Detection: SnooperText

#### Hypothesis generation: Segmentation

- Toggle Mapping [Serra89]: morphological operator
- Efficient approach for segmenting characters



[Serra89] Jean Serra, Toggle mappings, From pixels to features, pp. 61-72, 1989, J.C. Simon (ed.), North-Holland, Elsevier.

## SnooperText: Multi-Resolution Segmentation

- Difficult context for segmentation:
  - $\rightarrow$  very small relevant regions
  - $\rightarrow$  large textured text regions
- Multi-resolution goal:
  - Coarser levels: detect large regions and ignoring texture details (high frequencies)
  - Finer levels: detect smaller regions (analyzing accurately the local image content)

Region sizes managed at different resolution levels:



at level *I*, detecting text regions with size  $s_l \in [m_l; m_l + \delta_l]$  ( $c_l$ : overlap)

#### ANNEXES

# SnooperText: Character extraction & grouping

#### Character extraction

- Each region described by 3 shape descriptors:
  - Pzeudo-Zernike Moments (PZM)
  - Fourier Descriptors
  - Polar Descriptor [Fabrizzio2010]
- Late fusion: Hierachical SVM classifier

#### **Character grouping**

- Each character is merged with neighboring characters
  - Constraints related to distance, relative size, etc (see [Retornaz07])

[Fabrizzio2010] J. Fabrizio, M. Cord, B. Marcotegui, Text extraction from street level images, ISPRS Workshop CMRT, 2009 [Retornaz07] Thomas Retornaz and Beatriz Marcotegui, Scene text localization based on the ultimate opening, ISMM, vol.1, pp. 177-188, 2007.

Matthieu.cord@lip6.fr

Detection and representation for image understanding

