#### The Bag of Words Torn Open: Instance Retrieval goes Deep

Al Ukraine 2016 Kharkiv, Ukraine

James Pritts

Center for Machine Perception Czech Technical University in Prague







## Who are we?





Ondřej Chum Associate Professor



**Giorgos Tolias** Post-Doctoral candidate



Jiří Matas Professor



Filip Radenović PhD candidate



James Pritts PhD candidate

#### Goals

Introduce the Instance Retrieval Problem

- Compare two ways to learn an image encoding Bag-of-words (BoW) descriptor: ~1,000,000D vector
  - *Convolutional Neural Network (CNN)* descriptor 512D vector
- Demonstrate state-of-the-art retrieval performance

#### Part 1: The Instance Retrieval Task

Significant viewpoint and/or scale change Significant illumination change Severe occlusions Visually similar but different objects











Significant viewpoint and/or scale change
Significant illumination change
Severe occlusions
Visually similar but different objects



Significant viewpoint and/or scale change Significant illumination change

Severe occlusions

Visually similar but different objects











Significant viewpoint and/or scale change Significant illumination change Severe occlusions











#### **Instance Retrieval Demo**

<u>Click Here</u>

## Notional Instance Retrieval System



## Notional Instance Retrieval System



# Part 2: The Bag of Words (BoW) representation

## Bag of Words: Off-line stage

**Keypoint Detection** Local Appearance SIFT Description [Lowe'04] Keypoint descriptor mage gradient **Visual Vocabulary** graffiti Geom. Vocabulary  $x_1, y_1, B_1$ Local Geometry **Visual Words**  $x_2, y_2, B_5$  $x_{3}, y_{3}, B_{3}$  $word_1, word_2, word_8, \dots$ word<sub>948534</sub>, word<sub>998125</sub>  $x_N, y_N, B_N$ graffiti graffit

#### Quantization by K-Means



centres

Re-compute cluster centres as centroids

### Quantization by Approximate K-Means



- + fast O(N log k)
- + reasonable quantization
- Can be inconsistent when ANN fails

Philbin, Chum, Isard, Sivic, and Zisserman – CVPR 2007 Object retrieval with large vocabularies and fast spatial matching

## Quantization by Hierarchical K-means



Nistér & Stewénius: Scalable recognition with a vocabulary tree. CVPR 2006

## Bag-of-Words Image Representation



#### Bag of Words : On-line Stage





#### BoW and Inverted File



# BoW and Inverted File score = $\frac{\mathbf{q}^{\mathsf{T}}\mathbf{x}}{||\mathbf{x}||}$



#### BoW and Inverted File

Efficient (fast) Linear complexity (in # documents) Can be interpreted as voting



#### **Efficient Scoring**



### Word Weighting

Words (in text) common to many documents are less informative - 'the', 'and', 'or', 'in', ...

 $idf_{\chi} = \log \frac{\# \text{ documents}}{\# \text{ docs containing } (\mathbf{x})}$ 

Images are represented by weighted histograms  $tf_{\chi} idf_{\chi}$  (rather than just a histogram of  $tf_{\chi}$ )



Words that are too frequent (virtually in every document) can be put on a stop list (ignored as if they were not in the document)

Baeza-Yates, Ribeiro-Neto. Modern Information Retrieval. ACM Press, 1999.

#### Bag of Words : On-line Stage





#### **Query Expansion**



#### Results



New query

Chum, Philbin, Sivic, Isard, Zisserman: Total Recall..., ICCV 2007

#### Query Expansion: Step by Step







Query Image

Retrieved image

Originally not retrieved

#### Query Expansion: Step by Step







#### Query Expansion: Step by Step







#### The Bag of Words solution

Significant viewpoint scale change Significant illumination change Severe occlusions Visually similar but different objects

covariant local features, invariant descriptors color-normalized feature descriptors locality of the features, geometric verification Feature discriminability & geometric verification

#### **\*\*** Encoding is learned, but representation has many assumptions



Filip Radenović Giorgos Tolias Ondřej Chum

Center for Machine Perception, CTU in Prague

ECCV 2016



Filip Radenović Giorgos Tolias Ondřej Chum

Center for Machine Perception, CTU in Prague

ECCV 2016

CNN Image Retrieval

compact image descriptors Nearest Neighbor search



#### **CNN** Image Retrieval

compact image descriptors Nearest Neighbor search



#### CNN Learning (Fine-Tuning)

start with CNN trained for different but similar task (reasonable parameters) re-train with data relevant to your task

#### **CNN Image Retrieval**

compact image descriptors Nearest Neighbor search



#### CNN Learning (Fine-Tuning)

start with CNN trained for different but similar task (reasonable parameters) re-train with data relevant to your task

#### Bag of Words

state-of-the-art retrieval performance couples well with SfM



#### **CNN Image Retrieval**

compact image descriptors Nearest Neighbor search



#### **CNN Learning (Fine-Tuning)**

start with CNN trained for different but similar task (reasonable parameters) re-train with data relevant to your task

#### Bag of Words

state-of-the-art retrieval performance couples well with SfM

#### Unsupervised training data generation

no human interaction

#### **CNN Image Retrieval**

compact image descriptors Nearest Neighbor search



#### CNN Learning (Fine-Tuning)

start with CNN trained for different but similar task (reasonable parameters) re-train with data relevant to your task

#### Bag of Words

state-of-the-art retrieval performance couples well with SfM

#### Unsupervised training data generation

no human interaction

Hard Examples







hard **positives** 

hard **negatives** 

## "Lots of Training Examples"



Large Internet photo collection





Convolutional Neural Network (CNN)



## Off-the-shelf CNN

- Target application: classification
- Training dataset: ImageNet
- Architecture: AlexNet & VGG



Images from ImageNet.org

- Directly applicable to other tasks
  - Fine-grain classification



Images from ImageNet.org

#### **Object detection**





Images from PASCAL VOC 2012

#### Image retrieval







## Annotations for CNN Image Retrieval

CNN pre-trained for classification task used for retrieval

[Gong et al. ECCV'14, Babenko et al. ICCV'15, Kalantidis et al. arXiv'15, Tolias et al. ICLR'16]



Fine-tuned CNN using a dataset with landmark classes



• NetVLAD: Weakly supervised fine-tuned CNN using GPS tags [Arandjelovic et al. CVPR'16]



We propose: automatic annotations for CNN training





## **Retrieval and SfM**





























[Schonberger et al. CVPR'15] [Radenovic et al. CVPR'16]

## CNN learns from BoW – Training Data

#### Camera Orientation Known Number of Inliers Known



[Schonberger et al. CVPR'15] 7.4M images  $\rightarrow$  713 training 3D models

## Hard Negative Examples

Negative examples: images from different 3D models than the query Hard negatives: closest negative examples to the query Only hard negatives: as good as using all negatives, but faster

increasing CNN descriptor distance to the query

query



the most similar







naive hard negatives

top k by CNN











diverse hard negatives

top k: one per 3D model





## Hard Positive Examples

**Positive examples:** images from the same 3D model as the query **Hard positives:** positive examples not close enough to the query



## **CNN Siamese Learning**



## **CNN Siamese Learning**



**Contrastive vs. Triplet loss: Contrastive better with our data** 

Contrastive loss more strict, requires accurate training data Triplet loss less sensitive to inaccurate annotation

#### Whitening and dimensionality reduction



- 1. PCA<sub>w</sub> PCA of an independent set of descriptors [Babenko et al. ICCV'15, Tolias et al. ICLR'16]
- L<sub>w</sub> We propose to learn whitening using labeled training data and linear discriminant projections [Mikolajczyk & Matas ICCV'07]

#### Whitening and dimensionality reduction



- 1. PCA<sub>w</sub> PCA of an independent set of descriptors [Babenko et al. ICCV'15, Tolias et al. ICLR'16]
- L<sub>w</sub> We propose to learn whitening using labeled training data and linear discriminant projections [Mikolajczyk & Matas ICCV'07]
- 3. End-to-end Learning Performs comparable or worse than L<sub>w</sub>, while slowing down the convergence

#### Whitening and dimensionality reduction



- 1. PCA<sub>w</sub> PCA of an independent set of descriptors [Babenko et al. ICCV'15, Tolias et al. ICLR'16]
- L<sub>w</sub> We propose to learn whitening using labeled training data and linear discriminant projections [Mikolajczyk & Matas ICCV'07]
- 3. End-to-end Learning Performs comparable or worse than L<sub>w</sub>, while slowing down the convergence



(512D)

Nearest neighbors used on CNN descriptors Can use any fast NN search, like ANN

## Experiments – datasets

- Oxford 5k dataset [Philbin et al. CVPR'07]
- Paris 6k dataset [Philbin et al. CVPR'08]
- Holidays dataset [Jegou et al. ECCV'10]







• 100k distractor dataset [Philbin et al. CVPR'07] Training 3D models do not contain any landmark from these datasets

• **Protocol:** mean Average Precision (mAP)

## Experiments – Learning (AlexNet)

 Careful choice of positive and negative training images makes a difference



#### Experiments – Over-fitting and Generalization

 We added Oxford and Paris landmarks as 3D models and repeated fine-tuning



# Only +0.3 mAP on average over all testing datasets

	Mathad		Б	Ox	f5k	Oxf	105k	Pa	r6k	Par	106k	Hol	Hol
	Method		D	$\mathtt{Crop}_\mathcal{I}$	${\tt Crop}_{\mathcal{X}}$	$\mathtt{Crop}_\mathcal{I}$	$Crop_{\mathcal{X}}$	$\mathtt{Crop}_\mathcal{I}$	$Crop_{\mathcal{X}}$	$\mathtt{Crop}_\mathcal{I}$	${\tt Crop}_{\mathcal{X}}$		101k
State-ot-the-art	Compact representations												
	mVoc/BoW [11]		128	48.8	_	41.4	_	_	_	_	_	65.6	_
	Neural codes <sup>†</sup> [14]	$(\mathbf{fA})$	128	—	55.7	—	52.3	—	-	—	_	78.9	—
	$MAC^{\ddagger}$	$(\mathbf{V})$	128	53.5	55.7	43.8	45.6	69.5	70.6	53.4	<b>55.4</b>	72.6	56.7
	CroW [24]	$(\mathbf{V})$	128	<b>59.2</b>	_	51.6	_	74.6	-	63.2	_	-	—
	$\star$ MAC	$(\mathbf{fV})$	128	75.8	76.8	68.6	70.8	77.6	78.8	68.0	69.0	73.2	58.8
	$\star$ R-MAC	$(\mathbf{fV})$	128	72.5	76.7	64.3	69.7	78.5	80.3	<b>69.3</b>	71.2	79.3	65.2
	MAC <sup>‡</sup>	$(\mathbf{V})$	256	54.7	56.9	45.6	47.8	71.5	72.4	55.7	57.3	76.5	61.3
	SPoC [23]	$(\mathbf{V})$	256	—	53.1	—	50.1	—	-	—	_	80.2	—
	R-MAC [25]	$(\mathbf{A})$	256	56.1		47.0	_	72.9	_	60.1		-	
	CroW [24]	$(\mathbf{V})$	256	65.4	—	59.3	_	77.9	_	67.8	_	83.1	—
	NetVlad [35]	$(\mathbf{V})$	256		· ~		_	—	67.7	—	_	86.0	—
	NetVlad [35]	$(\mathbf{fV})$	255		<b>5</b> .	5	-	_	73.5		-	84.3	-
	* MAC	$(\mathbf{f}\mathbf{A})$	256		00.0	×0.0	58.0	68.9	72.2	54.7	58.5	76.2	63.8
NetVLAD 256D	* R-MAC	$(\mathbf{f}\mathbf{A})$	256	62.5	68.9	53.2	61.2	74.4	76.6	61.8	64.8	81.5	70.8
	* MAC	$(\mathbf{I} \mathbf{V})$	256	77.4	78.2	70.7	72.6	80.8	81.9	72.2	73.4	77.3	62.9
	* R-MAC	$(\mathbf{IV})$	256	74.9	78.2	67.5	72.1	82.3	83.5	74.1	75.6	81.4	69.4
	MAC <sup>*</sup>	$(\mathbf{V})$	512	56.4	58.3	47.8	49.2	72.3	72.6	58.0	59.1	76.7	62.7
VS.	$\begin{bmatrix} R-MAC & [25] \\ CnoW & [24] \end{bmatrix}$	$(\mathbf{V})$	512 519	00.9 68 9		01.0	_	83.0 70.6	_	75.7	_	×1 0	_
	1  MAC	$(\mathbf{V})$	512	70.7	80.0	03.4 79.0	75 1	19.0	- -	71.0	75.9	04.9 70 5	67.0
	$\star$ MAC	$(\mathbf{I} \mathbf{V})$ $(\mathbf{f} \mathbf{V})$	512	77.0	80.0	60.2	74.1	82.8	02.9 85 0	74.0 76 4	77.0	79.5 82.5	71 5
	$[\star n-mAC \qquad (IV)] \frac{12}{12} (1.0 \ 00.1] \frac{09.2}{14.1} \frac{14.1}{33.8} \frac{35.0}{50.0} \frac{10.4}{17.9} \frac{17.9}{82.5} \frac{82.5}{11.5}$												
		(2.4.)	10	Exti	reme	short	codes	3					
	Neural codes' [14]	$(\mathbf{f}\mathbf{A})$	16	-	41.8	_	35.4	-	-	-	-	<b>60.9</b>	-
	* MAC	$(\mathbf{f} \mathbf{V})$	16	56.2	57.4	45.5	47.6	57.3	62.9	43.4	48.5	51.3	25.6
	* R-MAC	$(\mathbf{f} \mathbf{V})$	16	46.9	52.1	37.9	41.6	58.8	63.2	45.6	49.6	54.4	31.7
	Neural codes' [14]	$(\mathbf{C}\mathbf{I}\mathbf{I})$	32				46.7	-	-	- F1 C	-	<b>72.9</b>	41.0
	$\star$ MAC	$(\mathbf{I} \mathbf{V})$	<b>პ</b>	6	9.4	2 P	59.5	63.9	69.5 67.4	51.0	56.3	62.4	41.8
	* R-MAC	$(\mathbf{I} \mathbf{V})$	32				55.1	63.9	07.4	52.7	00.0	08.0	49.0
		Re-rai	nkin	g(R)	and	query	expa	ansion	ı (QE	)			
Concurrent work:	BoW(1M) + QE[6]		—	82.7	—	76.7	-	80.5	-	71.0		-	—
	BoW(16M) + QE [50]		—	84.9		79.5	_	82.4	_	77.3	_	-	_
Gordo et al. ECCV'16	$ \mathrm{HQE}(65\mathrm{k}) [8] $	1 (***	_	88.0	—	84.0	_	82.8	_		_	-	—
	R-MAC+R+QE [25]	$ (\mathbf{V}) $	512	77.3	_	73.2	-	86.5	-	79.8		-	_
	CroW + QE [24]	$(\mathbf{V})$	512	72.2	-	67.8	-	85.5	-	79.7	-	-	_
	* MAC+R+QE	$(\mathbf{I} \mathbf{V})$	512 512	85.0	85.4	81.8	82.3	86.5	87.0	18.8	79.6	-	_
	★ K-MAC+K+QE	$(1 \mathbf{V})$	512	82.9	84.5	77.9	80.4	85.6	86.4	78.3	79.7	-	_

### Teacher vs. Student

Method	Oxf5k	Oxf105k	Par6k	Par106k		
BoW(16M)+R+QE	84.9	79.5	82.4	77.3		
CNN(512D)	79.7	73.9	82.4	74.6		
CNN(512D)+R+QE	85.0	81.8	86.5	78.8		

Our CNN with re-ranking (R) and query expansion(QE) surpasses its teacher on all datasets!!!

### Teacher vs. Student

#### top 10 (correct | incorrect)





BoW



#### first incorrect at rank 127



### Teacher vs. Student



#### top 10 (correct | incorrect)



### CNN descriptors

Significant viewpoint scale change	lots of training data
Significant illumination change	lots of training data
Severe occlusions	lots of training data
Visually similar but different objects	lots of training data

#### **CNN** descriptors

Significant viewpoint scale change Significant illumination change Severe occlusions Visually similar but different objects lots of training data lots of training data lots of training data lots of training data

#### versus

## Bag of Words

Significant viewpoint scale change Significant illumination change Severe occlusions Visually similar but different objects covariant local features, invariant descriptors color-normalized feature descriptors locality of the features, geometric verification Feature discriminability & geometric verification

## **CNN descriptor learning**

- Proposed a method to generate the necessary "lots of training examples" without any human interaction
- Strong supervision for hard negative, hard positive mining, and supervised whitening
- Data and trained networks available at: <u>cmp.felk.cvut.cz/~radenfil/projects/siamac.html</u>
- For more details about the paper visit **Poster O-1A-01**

# So Is the Bag-of-Words REALLY torn?

<u>Click Here</u>

# So is the Bag-of-Words REALLY torn?



## Not yet, but don't mess with tape ;)

## Questions?

- Thanks for your attention
- Interested students should ask about our PhD program

Center for Machine Perception Czech Technical University in Prague <u>http://cmp.felk.cvut.cz</u>

Contact Jiri Matas or Ondrej Chum



