

Tutorial on Datasets, Annotations and Metrics in Computer Vision

iV&LSS 2015





iV&L Net Training School 2015

Dr. Sergio Escalera Guerrero

ChaLearn, Computer Vision Center (UAB) and University of Barcelona Head of Human Pose Recovery and Behavior Analysis Group Vice-chair of International Association on Pattern Recognition IAPR-TC12: Visual and multimedia systems sergio@maia.ub.es

Leuven, Juve 1-4, 2015



INDEX

1.Computer Vision Datasets

2. Annotations and metrics

3.PASCAL Challenges

4.IMAGENET Challenges

5.ChaLearn Challenges

6.Other large scale CV DBs



How many CV DBs are available?

- Countless
 - Several research papers propose new datasets
 - Most of them are used for comparative among works in the field
 - Some of them are then selected for future challenges:
 - E.g. OpenCV CVPR 2015 Vision challenge:
 - Vision Challenge
 - OpenCV is launching a community-wide challenge to update and extend the OpenCV library. An award pool of \$50,000 will be provided to the best performing algorithms in the following 11 CV application areas:
 - image segmentation
 image registration
 human pose estimation
 SLAM (Simultaneous localization and mapping)
 multi-view stereo matching
 object recognition

face recognition
gesture recognition
action recognition
text recognition
tracking





Image/video datasets have increased in size in a dramatic way recently.



From: The Promise and Perils of Benchmark Datasets and Challenges. David Forsyth, Alyosha Efros, Fei-Fei Li, Antonio Torralba and Andrew Zisserman. Frontiers in Computer Vision Workshop, CVPR 2011.



Which is the recognition tasks in CV databases?

OpenCV challenge is very representative of the different topics:

image segmentation
image registration
human pose estimation
SLAM
multi-view stereo matching
object recognition

face recognition
gesture recognition
action recognition
text recognition
tracking

- Which is the input data:
 - RGB
 - Other visual representations: Depth, Thermal
 - Multi-modal multi-disciplinary in several cases: in combination with inertial sensors or audio, involving signal processing, pattern recognition, machine learning, computer vision, NLP, etc.



Examples of data for different computer vision topics

» image segmentation

Segmentation Competition

• Segmentation: Generating pixel-wise segmentations giving the class of the object visible at each pixel, or "background" otherwise.



http://pascallin.ecs.soton.ac.uk/challenges/VOC/



http://grand-challenge.org/All_Challenges/



Examples of data for different computer vision topics

» image registration and multi-view stereo matching





http://www.dir-lab.com/ReferenceData.html



Examples of data for different computer vision topics

Human Layout Pascal 2007-2010 competitions

Person Layout: Predicting the bounding box and label of each part of a person (head, hands, feet).



http://pascallin.ecs.soton.ac.uk/challenges/VOC/

Human Pose Estimation: New Benchmark and State of the Art Analysis CVPR 2014



Dataset	#training #test		img. type		
Full body pose datasets					
Parse [16]	100	205	diverse		
LSP [12]	1,000	1,000	sports (8 types)		
PASCAL Person Layout [6]	850	849	everyday		
Sport [21]	649	650	sports		
UIUC people [21]	346	247	sports (2 types)		
LSP extended [13]	10,000	-	sports (3 types)		
FashionPose [2]	6,530	775	fashion blogs		
J-HMDB [11]	31,838	-	diverse (21 act.)		
Upper body pose datasets					
Buffy Stickmen [8]	472	276	TV show (Buffy)		
ETHZ PASCAL Stickmen [3]	-	549	PASCAL VOC		
Human Obj. Int. (HOI) [23]	180	120	sports (6 types)		
We Are Family [5]	350 imgs.	175 imgs.	group photos		
Video Pose 2 [18]	766	519	TV show (Friends)		
FLIC [17]	6,543	1,016	feature movies		
Sync. Activities [4]	-	357 imgs.	dance / aerobics		
Armlets [9]	9,593	2,996	PASCAL VOC/Flickr		
MPII Human Pose (this paper)	28,821	11,701	diverse (491 act.)		



Examples of data for different computer vision topics

» object recognition

Visual Object Classes Challenge 2012 (VOC2012)



Classification/Detection Competitions

- 1. Classification: For each of the twenty classes, predicting presence/absence of an example of that class in the test image.
- 2. Detection: Predicting the bounding box and label of each object from the twenty target classes in the test image.

20 classes





Examples of data for different computer vision topics

 face recognition (identification, age estimation, gender recognition, facial expression analysis)

> ChaLearn Looking at People http://gesture.chalearn.org/

Age Recognition	PROFILE	GAME	ACHIEVEMENTS	GALLERY	RANKING				
UPLOAD IMAGE									
Upload some imag	es so other users can try to guess	s your age:							
	Upi	oad Images							
	.3 MiB 6.5 put the age: 25 Inpu emove image Ren	KIB It the age: 45 nove image	7.2 KIB Input the age: 48 Remove Image						
Privacy Policy	P	roject developed in col	aboration with:		900				



Examples of data for different computer vision topics

» gesture recognition





http://gesture.chalearn.org/



Examples of data for different computer vision topics

» text recognition





	C
	Frands white to the Corner
	formations may night as withhalf to mant as with
	and fing!
Successful and	I of a longer prover by at a land in mark of at these to fail
	in his pyleform
Co WARAN	. I how and good of states have made and a state of the states of the st
	A stay of the heat by stay and and growing and
Recontinenter	- the fire
	+ to I taken to be descriptioned to
	I Say and many source was and
	of therease to find in his polyticity .
V	6. Il my depays - flags of and of the same of
	in all friends around out prover has a set
	he find - her lightfrom
	Winners and The Uptal Constances
	[(c) a - and if the series as he has not by the series the her
	- the order is not proved and an and
	welly be fitting as adjudging of fails of the and dames



Examples of data for different computer vision topics

» Tracking







frameN: X1, Y1, X2, Y2, X3, Y3, X4, Y4

http://www.votchallenge.net/vot2014/



Examples of data for different computer vision topics

» Image retrieval



Microsoft Research

MSR-Bing Image Retrieval Challenge (IRC)

http://research.microsoft.com/en-us/projects/irc/



Human Pose Recovery and Behavior Analysis Group

INDEX

1.Computer Vision Datasets 2.Annotations and metrics 3.PASCAL Challenges 4.IMAGENET Challenges 5.ChaLearn Challenges 6.Other large scale CV DBs



image segmentation

Segmentation Competition

• Segmentation: Generating pixel-wise segmentations giving the class of the object visible at each pixel, or "background" otherwise.



Annotation: label at pixel level Don't care regions: boundaries

• Per-pixel classification

$$performan@ = \frac{number_of_corrected_classified_pixels}{total_number_of_pixels}$$

• Per-class classification

$$performance = \frac{1}{C} \sum_{i=1}^{C} \frac{number_of_corrected_classified_pixels_class_i}{total_number_of_pixels_class_i}$$

Overlapping (jaccard index)

$$performance = \frac{1}{C} \sum_{i=1}^{C} \frac{TP}{TP + FP + FN}$$
(Intersection over the union)



image registration and multi-view stereo matching

Annotation: label at pixel level Don't care regions: boundaries

There is no clear standard

- Gary E. Christensen, Xiujuan Geng, Jon G. Kuhl, Joel Bruss, Thomas J. Grabowski, Imran A. Pirwani, Michael W. Vannier, John S. Allen, Hanna Damasio, Introduction to the Non-rigid Image Registration Evaluation Project (NIREP), Biomedical Image Registration, Lecture Notes in Computer Science Volume 4057, 2006, pp 128-135
- Some used measurements:
 - Overlapping (jaccard index)

performance = $\frac{1}{C} \sum_{i=1}^{C} \frac{TP_i}{TP_i + FP_i + FN_i}$

(Intersection over the union)

• Correlation based on intensity values



human pose estimation

Annotation: bounding boxes for limbs

- **"PCP"** metric: considers a body part to be localized correctly if the estimated body segment endpoints are within 50% of the ground-truth segment length from their true locations.
- **"PCPm"** metric: uses 50% of the mean ground-truth segment length over the entire test set as a matching threshold for applying "PCP".

Annotation: joint coordinates

- **"PCK"** metric: measures accuracy of the localization of the body joints. The threshold for matching of the joint position to the ground-truth is defined as a fraction of the person bounding box size.
- **"PCKh"** metric: as "PCK" but define the matching threshold as 50% of the head segment length.

http://www.robots.ox.ac.uk/~vgg/publications/papers/ferrari08.pdf http://ps.is.tuebingen.mpg.de/publications/168/get_file



object recognition

Annotation: bounding box

- Detection vs classification
- **Detection** $performance = \frac{1}{C} \sum_{i=1}^{C} \frac{TP_i}{TP_i + FP_i + FN_i}$

Where a hit for TP_i should satisfy $\frac{\Omega}{U} > \Theta$

• Multi-class classification

 $performan@ = \frac{number_of_corrected_classified_samples}{total_number_of_samples}$

~ sum(Diag (confusion matrix)) / (sum(sum(confusion matrix))

Face Analysis

- Face detection: detection
- Identification:
- Gender Recognition: binary classification
- Verification:
- Age estimation:

- multi-class classification

binary classification

Annotation: bounding box

Annotation: label per face/image

Annotation: label per face/image

Annotation: label per face/image

- **Regression deviation Annotation**: label/labels/levels per face/image
- Facial expression analysis:
 - multi-class classification
 - Regression in some cases
 - Multi-label



Multi-label Evaluation Metrics: Label-Based

Basic Strategy:

Calculate classic single-label metric on each label *independently*, and then combine metric values over all labels.

Label-based multi-label metrics are easy to compute, but ignore the relationships between different labels!

Exact Match: is the most strict metric, indicating the percentage of samples that have all their labels classified correctly.



Multi-label Evaluation Metrics: Instance-Based Basic Strategy:

Calculate metric value for each instance by addressing relationships among different class labels (especially the *ranking quality*), and then return the mean value over all instances.

Popular instance-based multi-label metrics [Schapire & Singer, MLJ00]:

Given the learned predictor $h(\cdot)$ or $f(\cdot, \cdot)$, and a test set $T = \{(x_i, Y_i | 1 \le i \le t)\}$

(1) Hamming loss

$$hloss_T(h) = \frac{1}{t} \sum_{i=1}^t \frac{1}{q} |h(\boldsymbol{x}_i \Delta Y_i)|$$

(3) *Coverage*

$$\operatorname{coverage}_{T}(f) = \frac{1}{t} \sum_{i=1}^{t} \max_{y \in Y_{i}} \operatorname{rank}_{f}(\boldsymbol{x}_{i}, y) - 1$$

 $\begin{bmatrix} rank_f(\boldsymbol{x}_i, y) \text{ returns the rank of } y \text{ derived from } f(\boldsymbol{x}_i, y) \end{bmatrix}$

Evaluates how many times an instance-label pair is misclassified (is the percentage of the wrong labels to the total number of labels. As a loss metric, 0 is better), being q a normalization value (e.g. the number of labels, cardinality of Y).

Evaluates how many steps are needed, on average, to go down the label list to cover all proper labels of the instance.

Useful metric for image retrieval (related also to taxonomy analysis)



gesture recognition (classification vs spotting)

- Multi-class classification Annotation: label per segmented sequence
- Spotting (can be multi-label depending on behavior taxonomy)





text recognition

Annotation: words locations and ground truth words



• WER – Word Error Rate (based on Levensthein distance)

$$WER = \frac{S+B+I}{N}$$

- S number of substitutions
- •B number of deletions
- I number of insertions
- •*N* number of sentence words

$$WER(i,j) = \min \begin{cases} WER(i-1,j) + 1\\ WER(i,j-1) + 1\\ WER(i-1,j-1) + \Delta(i,j) \end{cases}$$

text recognition

 $E(i,j) = \min\{E(i-1,j) + 1, E(i,j-1) + 1, E(i-1,j-1) + \texttt{diff}(i,j)\}$



		Р	0	L	Y	Ν	0	Μ	Ι	Α	L
	0	1	2	3	4	5	6	7	8	9	10
E	1	1	2	3	4	5	6	7	8	9	10
Х	2	2	2	3	4	5	6	7	8	9	10
Р	3	2	3	3	4	5	6	7	8	9	10
0	4	3	2	3	4	5	5	6	7	8	9
Ν	5	4	3	3	4	4	5	6	7	8	9
E	6	5	4	4	4	5	5	6	7	8	9
Ν	7	6	5	5	5	4	5	6	7	8	9
Т	8	7	6	6	6	5	5	6	7	8	9
Ι	9	8	7	7	7	6	6	6	6	7	8
Α	10	9	8	8	8	7	7	7	7	6	7
L	11	10	9	8	9	8	8	8	8	7	6

text recognition

 $E(i,j) = \min\{E(i-1,j) + 1, E(i,j-1) + 1, E(i-1,j-1) + \texttt{diff}(i,j)\}$





Tracking

Annotation: trajectory coordinates

- **Trajectories**: accumulated distance for detected localitions
- Bounding boxes:
 - Overlap
 - Hit per tracked bounding box based on overlapping threshold



Usage of training – validation – test sets

- Partitions:
 - Training used for learning methods
 - Validation useful for tuning parameters, support generalization, and avoid or delay the appearance of overfitting
 - Test only used for final generalization performance
 - Use of many splits of the data:
 - N-fold cross-validation
 - Random vs stratified
 - Confidence interval is useful to analyze the stability of the results
 - Statistical significance analysis

Janez Demsar, Statistical Comparisons of Classifiers over Multiple Data Sets, Journal of Machine Learning Research, Volume 7, 12/1/2006, Pages 1-30



Which aspects of the data set are interesting

Looking for generalization: -Ground truth quality

- don't care regions
- Inter-labeller varability
- •etc.
- -Training set variability
- -Amount of data



INDEX

1.Computer Vision Datasets 2.Annotations and metrics **3.PASCAL Challenges**4.IMAGENET Challenges 5.ChaLearn Challenges 6.Other large scale CV DBs



PASCAL

The 2005-2012 Visual Object Challenges

A. Zisserman, C. Williams, M. Everingham, L. v.d. Gool







Classification: is there an X in this image?,

Detection: where are the X's ?

Segmentation: which pixels belong to X ?

PASCAL

The challenge organizer perspective

Selecting the data

real images from flickr (no selection / cleaning)

Annotation, how

Occluded Object is significantly occluded within BB

Truncated ____ Object extends beyond BB



Difficult Not scored in evaluation

Pose Facing left

Annotation, who

- Annotation parties
- The Amazon turk

Experimental setting

- At least 500 images per object
- Equally divided among training/validation and test
- Increased along years (enables to measure progress)

Software supplied

- Includes baseline classifier/detector/segmenter
- Generates precision-recall curve and computes accuracy scores
- On train/validation/test and other datasets

Means that results on VOC can be consistently compared in publications



PASCAL

Evaluation

Option 1

Release test data and annotation (most liberal) and participants can assess performance Cons: open to abuse

Option 2

Release test data, but test annotation withheld - participants submit results and organizers assess performance (use an evaluation server)

Option 3

No release of test data - participants have to submit software and organizers run this and assess performance

Human Pose Recovery and Behavior Analysis Group

PASCAL

Detection Challenge: Progress 2008-2010



 Results on 2008 data improve for best 2009 and 2010 methods for all categories, by over 100% for some categories

– Caveat: Better methods or more training data?

INDEX

1.Computer Vision Datasets
2.Annotations and metrics
3.PASCAL Challenges
4.IMAGENET Challenges
5.ChaLearn Challenges
6.Other large scale CV DBs


IMAGENET Large Scale Visual Recognition Challenge (ILSVRC)

http://image-net.org/challenges/LSVRC/

Thx to: Olga Russakovsky, Stanford University

Large-scale recognition







COLUMN STREET, STRE



PASCAL VOC 2005-2012 20 object classes 22,591 images

IMAGENET

Classification: person, motorcycle



Action: riding bicycle

Everingham, Van Gool, Williams, Winn and Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010. **PBA**

IMAGENET

IMAGENET ILSVRC 2010-2014

20 object classes

22,591 images

200 object classes DETECTION 1000 object classes CLASSIF.

517,840 images

1,431,167 images



http://image-net.org/challenges/LSVRC/

IMAGENET

ILSVRC types of image annotations

Image classification

- <u>one</u> object class per image
- <u>no</u> bounding boxes

Steel drum



Single-object localization

- <u>one</u> object class per image
- bounding boxes around <u>all</u> <u>instances</u> of this class

Steel drum



Object detection

- <u>all</u> target object classes
- bounding boxes around <u>all</u> <u>instances</u>



Statistics of ILSVRC2014 released annotated images:

1000 object classes 1,331,167 images 1000 object classes 573,966 images 657,231 bounding boxes 200 object classes 81,799 images 228,981 bounding boxes

IMAGENET

ILSVRC large-scale annotation

Artificial Artificial Intelligence



J. Deng, W. Dong, R. Socher, L.-J. Li, L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. CVPR, 2009







Challenge procedure every year

- 1. Training data released: images and annotations
- 2. Test data released: images only (annotations hidden)
- 3. Participants train their models on train data
- 4. Submit text file with predictions on test images
- 5. Evaluate and release results, and run a workshop

http://image-net.org/challenges/LSVRC/2014/eccv2014

IMAGENET

Participation in ILSVRC over the years





IMAGENET



IMAGENET

ILSVRC image classification task

Steel drum

PBA









Deep learning impact on ILSVRC classification accuracy

IMAGENET



Massive drop in error with a deep learning method

Deep learning is here to stay

Year 2012

SuperVision





Year 2014

VGG

image

conv-64 maxpool

conv-128 conv-128 maxpool conv-256

conv-256

maxpool

conv-512 maxpool

conv-512 conv-512 maxpool FC-4096

FC-4096 FC-1000

softmax



<u>ē</u> —

35/36 teams used deep learning

20/36 teams used open-source Caffe implementation

[Krizhevsky NIPS 2012]

[Szegedy arxiv 2014] [Simonyan arxiv 2014] [He arxiv 2014]

(Highest accuracy in percent of any method in ILSVRC 2012-2014)

Easiest and hardest categories for image classification

red fox (100) hen-of-the-woods (100) ibex (100) goldfinch (100) flat-coated retriever (100) Easiest tiger (100) hamster (100) porcupine (100) stingray (100) Blenheim spaniel (100) muzzle (71) hatchet (68) water bottle (68) velvet (68) loupe (66)

Hardest



















ladle (65)



... and 111 more categories with 100% accuracy!

49



restaurant (64) letter opener (59)











PBA

IMAGENET

ILSVRC single-object localization task

Steel drum





ILSVRC single-object localization task



Steel drum

Output





ILSVRC single-object localization task

Steel drum



Output (bad localization)



Output



Output (bad classification)





ILSVRC single-object localization task

Steel drum



Output





IMAGENET

ILSVRC over the years

Classification





IMAGENET

(Highest *accuracy* in percent of any method in ILSVRC 2012-2014)

55

Easiest and hardest categories for single-object localization



Russakovsky* and Deng* et al., ImageNet Large Scale Visual Recognition Challenge, http://arxiv.org/abs/1409.0575

IMAGENET

New in 2013

ILSVRC object detection task

Fully annotated 200 object classes across 120,000 images



Allows evaluation of generic object detection in cluttered scenes at scale

IMAGENET

ILSVRC object detection task

<u>All</u> instances of <u>all</u> target object classes expected to be localized on <u>all</u> test images



Evaluation modeled after PASCAL VOC:

- Algorithm outputs a list of bounding box detections with confidences
- A detection is considered correct if overlap with ground truth is big enough
- Evaluated by average precision per object class
- Winners of challenge is the team that wins the most object categories

IMAGENET

ILSVRC object detection data



PBA









ILSVRC detection since 2013

IMAGENET



1.9x increase in object detection average precision in one year

(Highest average precision in percent of any method in ILSVRC 2012-2014)

Easiest and hardest categories for object detection

butterfly (93) volleyball (83) dog (84) rabbit (83) frog (82) Easiest basketball (80) snowplow (80) bird (78) tiger (77) zebra (77) horizontal bar (14) spatula (13) lamp (15) flute (15) nail (13)

backpack (8)



ladle (9)







ski (12)



60 Russakovsky* and Deng* et al., ImageNet Large Scale Visual Recognition Challenge, http://arxiv.org/abs/1409.0575



Q: So how well do current methods work on large-scale object recognition?

• Well, they work much better than last year!

IMAGENET

- Work very well on
 - Classifying and detecting animals
 - Recognizing objects with distinctive patterns
- Don't work as well on
 - Thin objects
 - Untextured objects



What is human accuracy on ILSVRC2014 classification? Human vs computer accuracy on ILSVRC2014 classification

	Annotator 1
Total number of images	1500
GoogLeNet classification error	6.8%
Human classification error	?



Human vs computer accuracy on ILSVRC2014 classification

	Annotator 1
Total number of images	1500
GoogLeNet classification error	6.8%
Human classification error	5.1%



Human vs computer accuracy on ILSVRC2014 classification

	Annotator 1
Total number of images	1500
GoogLeNet classification error	6.8%
Human classification error	5.1%

- Annotator 1 achieved better accuracy than GoogLeNet by 1.7%
- Task required *significant* amount of training for humans



Human vs computer accuracy on ILSVRC2014 classification

	Annotator 1	Annotator 2
Total number of images	1500	258
GoogLeNet classification error	6.8%	5.8%
Human classification error	5.1%	12.0%

- Annotator 1 achieved better accuracy than GoogLeNet by 1.7%
- Task required *significant* amount of training for humans

Q: Are current methods close to humanlevel classification accuracy?

Current methods are not as good as humans yet,

<u>but</u>

current methods are better than non-domain-expert humans on fine-grained classification!



Future: MORE OBJECTS MORE CONTEXT INFORMATION

ILSVRC object detection: *all instances of the 200 target objects*





Future: MORE OBJECTS MORE CONTEXT INFORMATION





1.Computer Vision Datasets
2.Annotations and metrics
3.PASCAL Challenges
4.IMAGENET Challenges **5.ChaLearn Challenges**6.Other large scale CV DBs



ChaLearn Looking at People

ChaLearn http://www.chalearn.org/





Mission:

Machine Learning is the science of building hardware or software that can achieve tasks by learning from examples. The examples often come as {input, output} pairs. Given new inputs a trained machine can make predictions of the unknown output.

Examples of machine learning tasks include:

- automatic reading of handwriting
- assisted medical diagnosis
- automatic text classification (classification of web pages; spam filtering)
- financial predictions

We organize challenges to stimulate research in this field. **The web sites of past challenges remain open** for post-challenge submission as ever-going benchmarks.

ChaLearn is a tax-exempt organization under section 501(c)(3) of the US IRS code. DLN: 17053090370022.



ChaLearn Looking at people (multimedia datasets, http://gesture.chalearn.org/)



ChaLearn Looking at People Challenges and Workshops

CVPR 2011	Workshop and Challenge on Multi-modal Sign Language Recognition
CVPR 2012	Workshop and Challenge on Multi-modal Sign Language Recognition

- **ICPR 2012** Workshop and Challenge on Multi-modal Sign Language Recognition
- ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition
- **ECCV 2014** Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting
- CVPR 2015Workshop and Challenge on human pose recovery, action/interactionspotting, cultural event recognition
- ICCV 2015 Workshop and Challenge on age estimation, action spotting and cultural event recognition
- **~ 2016** Workshop and Challenge on Multi-modal speed interviews analysis

And so on! Let us know about your opinion! <a>sergio@maia.ub.es



ChaLearn Looking at People

ChaLearn Looking at People Challenges and Workshops

CVPR 2011Workshop and Challenge on Multi-modal Sign Language RecognitionCVPR 2012Workshop and Challenge on Multi-modal Sign Language RecognitionICPR 2012Workshop and Challenge on Multi-modal Sign Language Recognition



Microsoft[®] Research




ChaLearn Looking at People Challenges and Workshops

CVPR 2011Workshop and Challenge on Multi-modal Sign Language RecognitionCVPR 2012Workshop and Challenge on Multi-modal Sign Language RecognitionICPR 2012Workshop and Challenge on Multi-modal Sign Language Recognition

- **1.** Body language gestures (like scratching your head, crossing your arms).
- 2. Gesticulations performed to accompany speech.
- 3. Illustrators (like Italian gestures).
- 4. Emblems (like Indian Mudras).
- 5. Signs (from sign languages for the deaf).
- 6. Signals (like referee signals, diving signals, or Marshalling signals to guide machinery or vehicle).
- 7. Actions (like drinking or writing).
- 8. Pantomimes (gestures made to mimic actions).
- 9. Dance postures.

Evaluation metric: levenstein edition distance

Microsoft[®] Research





ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition



Multi-modal ChaLearn Gesture Recognition Challenge and Workshop

http://gesture.chalearn.org/ sunai.uoc.edu/chalearn Web of the competition Data

The challenge features a **quantitative evaluation** of automatic gesture recognition from a multi-modal dataset recorded with **Kinect** (providing RGB images of face and body, depth images of face and body, skeleton information, joint orientation and audio sources), **including 13,858 Italian gestures from near 30 users.**

The emphasis of this edition of the competition will be on multi-modal automatic learning of a vocabulary of 20 types of Italian anthropological/cultural gestures performed by different users, with the aim of performing user independent continuous gesture recognition combined with audio information.



ChaLearn Looking at People Challenges and Workshops

ICMI 2013

Workshop and Challenge on Multi-modal gesture recognition Gesture categories (1/2)



(1) Vattene



(2) Viene qui





(3) Perfetto





(4) E un furbo





(5) Che due palle



(10) Nonme me friega 75 niente

(6) Che vuoi



(8) Sei pazzo

(9) Cos hai combinato



ChaLearn Looking at People Challenges and Workshops

ICMI 2013

Workshop and Challenge on Multi-modal gesture recognition Gesture categories (2/2)



(11) Ok



(12) Cosa ti farei



(17) Tanto tempo fa



(13) Basta



(18) Buonissimo



(14) Le vuoi prendere



(19) Si sono messi d'accordo



(15) Non ce ne piu



(20) Sono stufo



ChaLearn Looking at People Challenges and Workshops

ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition Data and modalities



- Framerate 20FPS
- RGB: 640x480
- Depth: 640x480
- Audio: Kinect 20 michropone array
- Users: 27
- Italians: 81%

- Total number of sequences: 956 \in [1,2] min.
- Total number of gestures: 13,858
- Total number of frames: 1.720.800
- Noisy gestures

Data structure information: *S. Escalera, J. Gonzàlez, X. Baró, M. Reyes, O. Lopes, I. Guyon, V. Athistos, H.J. Escalante, "Multi-modal Gesture Recognition Challenge 2013: Dataset and Results", ICMI 2013.*



ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition

Easy and challenging aspects of the data.

	Easy
	Fixed camera
	Near frontal view acquisition
	Within a sequence the same user
	Gestures performed mostly by arms and hands
	Camera framing upper body
	Several available modalities: audio, skeletal model, user mask,
	depth, and RGB
	Several instances of each gesture for training
	Single person present in the visual field
,	
	Challenging
1	Within each sequence:
	Continuous gestures without a resting pose
	Many gesture instances are present
	Distracter gestures out of the vocabulary may be present in terms
	of both gesture and audio
	Between sequences:
	High inter and intra-class variabilities of gestures in terms of both
	gesture and audio
	Variations in background, clothing, skin color, lighting, tempera-
	ture, resolution
	Some parts of the body may be accluded
1	some parts of the body may be occluded

Evaluation metric: levenstein edition distance



ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition







ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition

• Participation

• The challenge attracted high level of participation, with a total of **54 teams** and near **300 total number of entries**.

• Finally, 17 teams successfully submitted their prediction in final test set, while providing also their code for verification and summarizing their method by means of a fact sheet questionnaire.

• After verifying the codes and results of the participants, the final scores of the top rank participants on both validation and test sets were made public.

• In the end, the final error rate on the test data set was around 12%.

Top rank results on validation and test sets.

\mathbf{TEAM}	Validation score	Test score
IVA MM	0.20137	0.12756
WWEIGHT	0.46163	0.15387
${ m ET}$	0.33611	0.16813
MmM	0.25996	0.17215
PPTK	0.15199	0.17325
LRS	0.18114	0.17727
MMDL	0.43992	0.24452
TELEPOINTS	0.48543	0.25841
CSI MM	0.32124	0.28911
SUMO	0.49137	0.31652
GURU	0.51844	0.37281
AURINKO	0.31529	0.63304
STEVENWUDI	1.43427	0.74415
JACKSPARROW	0.86050	0.79313
JOEWAN	0.13653	0.83772
MILAN KOVAC	0.87835	0.87463
IAMKHADER	0.93397	0.92069



ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition



Validation and test scores histograms.



ChaLearn Looking at People Challenges and Workshops

ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition

Team methods and results. Early and late refer to early and late fusion of features/classifier outputs. HMM: Hidden Markov Models. KNN: Nearest Neighbor. RF: Random Forest. Tree: Decision Trees. ADA: Adaboost variants. SVM: Support Vector Machines. Fisher: Fisher Linear Discriminant Analysis. GMM: Gaussian Mixture Models. NN: Neural Networks. DGM: Deep Boltzmann Machines. LR: Logistic Regression. DP: Dynamic Programming. ELM: Extreme Learning Machines.

TEAM	Test score	Rank position	Modalities	Segmentation	Fusion	Classifier
IVA MM	0.12756	1	Audio,Skeleton	Audio	None	HMM,DP,KNN
WWEIGHT	0.15387	2	Audio,Skeleton	Audio	Late	RF,KNN
ET	0.16813	3	Audio,Skeleton	Audio	Late	Tree,RF,ADA
MmM	0.17215	4	Audio,RGB+Depth	Audio	Late	SVM,Fisher,GMM,KNN
PPTK	0.17325	5	Skeleton,RGB,Depth	Sliding windows	Late	GMM,HMM
LRS	0.17727	6	Audio,Skeleton,Depth	Sliding windows	Early	NN
MMDL	0.24452	7	Audio,Skeleton	Sliding windows	Late	DGM+LR
TELEPOINTS	0.25841	8	Audio,Skeleton,RGB	Audio,Skeleton	Late	HMM,SVM
CSI MM	0.28911	9	Audio,Skeleton	Audio	Early	HMM
SUMO	0.31652	10	Skeleton	Sliding windows	None	RF
GURU	0.37281	11	Audio,Skeleton,Depth	DP	Late	DP, RF, HMM
AURINKO	0.63304	12	Skeleton,RGB	Skeleton	Late	ELM
STEVENWUDI	0.74415	13	Audio,Skeleton	Sliding windows	Early	DNN,HMM
JACKSPARROW	0.79313	14	Skeleton	Sliding windows	None	NN
JOEWAN	0.83772	15	Skeleton	Sliding windows	None	KNN
MILAN KOVAC	0.87463	16	Skeleton	Sliding windows	None	NN
IAMKHADER	0.92069	17	Depth	Sliding windows	None	RF

Data structure information: *S. Escalera, J. Gonzàlez, X. Baró, M. Reyes, O. Lopes, I. Guyon, V. Athistos, H.J. Escalante, "Multi-modal Gesture Recognition Challenge 2013: Dataset and Results", ICMI 2013.*



ChaLearn Looking at People Challenges and Workshops

ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition





ICMI 2013 Workshop and Challenge on Multi-modal gesture recognition





ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

•<u>Track 1: Human Pose Recovery</u>: More than 8,000 frames of continuous RGB sequences are recorded and labeled with the objective of performing human pose recovery by means of <u>recognizing more than 120,000 human limbs</u> of different people.

•<u>Track 2: Action/Interaction Recognition</u>: 235 performances of 11 action/interaction categories are recorded and manually labeled in continuous RGB sequences of different people performing natural isolated and collaborative behaviors.

•<u>Track 3: Gesture Recognition</u>: The gestures are drawn from a vocabulary of Italian sign gesture categories. The emphasis of this third track is on multi-modal automatic learning of a set of 20 gestures performed by several different users, with the aim of performing <u>user independent continuous gesture spotting</u>.



Workshop and Challenge on multi-modal gesture spotting, human pose **ECCV 2014** recovery, action/interaction spotting

•Track 1: Human Pose Recovery: More than 8,000 frames of continuous RGB sequences are recorded and labeled with the objective of performing human pose recovery by means of recognizing more than 120,000 human limbs of different people.

Training frames	Validation frames	Test frames	Sequence duration	FPS
4,000	2,000	2,236	1-2 min	15
Modalities	Num. of users	Limbs per body	Labeled frames	Labeled limbs
RGB	14	14	8,234	124,761

Human pose recovery data characteristics.

•9 videos (RGB sequences) and a total of 14 different actors. Stationary camera with the same static background.

•15 fps rate, resolution 480x360 in BMP file format.

• For each actor **14 limbs** (if not occluded) were manually tagged: Head, Torso, R-L Upper-arm, R-L Lower-arm, R-L Hand, R-L Upper-leg, R-L Lower-leg, and R-L Foot.

• Limbs are manually labeled using binary masks and the minimum bounding box containing each subject is defined.





• The actors appear in a wide range of different poses and performing different actions/gestures which vary the visual appearance of human limbs. So there is a large variability of human poses, self-occlusions and many variations in clothing and skin color. 86



ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

•<u>Track 1: Human Pose Recovery</u>: More than 8,000 frames of continuous RGB sequences are recorded and labeled with the objective of performing human pose recovery by means of <u>recognizing more than 120,000 human limbs</u> of different people.

$$H_{i,n} = \frac{A_{i,n} \bigcap B_{i,n}}{A_{i,n} \bigcup B_{i,n}},$$
 Overlap evaluation
$$H_{i,n} = \begin{cases} 1 & \text{if } \frac{A_n \bigcap B_n}{A_n \bigcup B_n} \ge 0.5\\ 0 & \text{otherwise} \end{cases}$$

 $B_{i,head}$

 $A_{i,head}$

J



$$J_{i,head} = \frac{A_{i,head} \cap B_{i,head}}{A_{i,head} \cup B_{i,head}} = 0.82$$
$$J_{i,head} > 0.5 \longrightarrow HR_{i,head} = 1$$





ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

•**Track 2: Action/Interaction Recognition:** 235 performances of 11 action/interaction categories are recorded and manually labeled in continuous RGB sequences of different people performing natural isolated and collaborative behaviors.

Training actions	Validation actions	Test actions	Sequence duration	FPS
150	90	95	$9 \times 1-2 \min$	15
Modalities	Num. of users	Action categories	interaction categories	Labeled sequences
RGB	14	7	4	235

Action and interaction data characteristics.

- 235 action/interaction samples performed by 14 actors.
- Large **difference in length** about the performed actions and interactions.
- Several distracter actions out of the 11 categories are also present.

• **11** action categories, containing isolated and collaborative actions: Wave, Point, Clap, Crouch, Jump, Walk, Run, Shake Hands, Hug, Kiss, Fight. There is a high intra-class variability among action samples.

Overlap evaluation







ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

Track 2: Action/Interaction Recognition

Action categories



Wave





Clap







Jump



Walk



Run



Point



Hug



Interaction categories



Shake Hands

Kiss

Fight



ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

<u>Track 3: Gesture Recognition</u>: The gestures are drawn from a vocabulary of Italian sign gesture categories. The emphasis of this third track is on multi-modal automatic learning of a set of 20 gestures performed by several different users, with the aim of performing <u>user independent continuous gesture spotting</u>.

Training seq.	Validation seq.	Test seq.	Sequence duration	FPS	
393 (7,754 gestures)	287 (3,362 gestures)	276 (2,742 gestures)	1-2 min	20	
Modalities	Num. of users	Gesture categories	Labeled sequences	Labeled frames	
RGB, Depth, User mask, Skeleton	27	20	13,858	1,720,800	
Main abaractoristics of the Montalhane mature dataset					

Main characteristics of the *Montalbano* gesture dataset.





Depth

User mask

Skeletal model



ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

•Largest dataset in the literature with a large duration of each individual performance showing no resting poses and self-occlusions.

• There is **no information about the number of gestures to spot** within each sequence, and **several distracter gestures** (out of the vocabulary) are present.

• High intra-class variability of gesture samples and low inter-class variability for some gesture categories.





ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

• State of the art comparison

	Labeling at pixel precision	Number of limbs	Number of labeled limbs	Number of frames	Full body	Limb annotation	Gesture- action annotation	Number of gestures- actions	Number of gest-act. samples
Montalbano[8]	No	16	27532800	$1\ 720\ 800$	Yes	Yes	Yes	20	13858
HuPBA $8K + [7]$	Yes	14	124761	8 2 3 4	Yes	Yes	Yes	11	235
LEEDS SPORTS[4]	No	14	28000	2000	Yes	Yes	No	-	-
UIUC people[10]	No	14	18186	1299	Yes	Yes	No	-	-
Pascal VOC[2]	Yes	5	8500	1218	Yes	Yes	No	-	-
BUFFY[3]	No	6	4488	748	No	Yes	No	-	-
PARSE[11]	No	10	3050	305	Yes	Yes	No	-	-
MPII Pose[12]	Yes	14	-	40522	Yes	Yes	Yes	20	491
FLIC[13]	No	29	-	5003	No	No	No	-	-
H3D[14]	No	19	-	2000	No	No	No	-	-
Actions[15]	No	-	-	-	Yes	No	Yes	6	600
HW[5]	-	-	-	-	-	No	Yes	8	430

Comparison of public dataset characteristics.

ChaLearn LAP data sets, public available at: <u>http://sunai.uoc.edu/chalearnLAP/</u> ChaLearn LAP challenges and news: <u>http://gesture.chalearn.org/</u>



93

ChaLearn Looking at People Challenges and Workshops

ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

• **Connectivity:** During the Challenge period, the download page had a total of 2.895 visits from 920 different users of 59 countries.





ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

• Track1 results



Team	Accuracy	Rank position	Features	Pose model
ZJU	0.194144	1	HOG	tree structure
Seawolf Vision	0.182097	2	HOG	tree structure

Track 1 Pose Recovery results.

Both winner participants applied a similar approach based on [*].

 Mixture of templates for each part. This method incorporates the co-occurrence relations, appearance and deformation into a model represented by an objective function of pose configurations. Model is tree-structured, and optimization is conducted via dynamic programming.

[*] Yang, Y., Ramanan, D.: Articulated human detection with flexible mixtures of parts. IEEE TPAMI (2013)



ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

Team name	Accuracy	Rank	Features
CUHK-SWJTU	0.507173	1	Improved trajectories [\star]
ADSC	0.501164	2	Improved trajectories [\star]
SBUVIS	0.441405	3	Improved trajectories [\star]
DonkeyBurger	0.342192	4	MHI, STIP
UC-T2	0.121565	5	Improved trajectories [\star]
MindLAB	0.008383	6	MBF

Track2 results

]	V	2
]		
-		
ļ	Ø₿	HH
	Since the Oas	

Team name	Dimension reduction	Clustering	Classifier	Temporal coherence	Gesture representation
CUHK-SWJTU	PCA	-	SVM	Sliding windows	Fisher Vector
ADSC	-	-	SVM	Sliding windows	-
SBUVIS	-	-	SVM	Sliding windows	-
DonkeyBurger	-	Kmeans	Sparse code	Sliding windows	-
UC-T2	PCA	-	Kmeans	Sliding windows	Fisher Vector
MindLAB	-	\mathbf{Kmeans}	\mathbf{RF}	Sliding windows	BoW

* Wang, H., Schmid, C.: Action recognition with improved trajectories. ICCV (2013)



ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

• Track3 results

Percentage of methods using each independent modality

Team	Accuracy	Rank	Modalities
LIRIS	0.849987	1	SK, Depth, RGB
CraSPN	0.833904	2	SK, Depth, RGB
JY	0.826799	3	SK, RGB
CUHK-SWJTU	0.791933	4	RGB
Lpigou	0.788804	5	Depth, RGB
stevenwudi	0.787310	6	SK, depth
Ismar	0.746632	7	SK
Quads	0.745449	8	SK
Telepoints	0.688778	9	SK, Depth, RGB
TUM-fortiss	0.648979	10	SK, Depth, RGB
CSU-SCM	0.597177	11	Skeleton, Depth, mask
iva.mm	0.556251	12	Skeleton, RGB, depth
Terrier	0.539025	13	Skeleton
Team Netherlands	0.430709	14	Skeleton, Depth, RGB
VecsRel	0.408012	15	Skeleton, Depth, RGB
Samgest	0.391613	16	Skeleton, Depth, RGB, mask
YNL	0.270600	17	Skeleton





ECCV 2014Workshop and Challenge on multi-modal gesture spotting, human pose
recovery, action/interaction spotting

• Track3 results Percentage of methods using each gesture classification strategy

Team	Gesture representation	Classifier
LIRIS	-	DNN
CraSPN	BoW	Adaboost
JY	-	MRF, KNN
CUHK-SWJTU	Fisher Vector, VLAD	SVM
Lpigou	-	CNN
stevenwudi	-	HMM, DNN
Ismar	-	\mathbf{RF}
Quads	Fisher Vector	SVM
Telepoints	-	SVM
TUM-fortiss	-	RF, SVM
CSU-SCM	2DMTM	SVM, HMM
iva.mm	BoW	SVM, HMM
Terrier	-	RF
Team Netherlands	-	SVM, RT
VecsRel	-	DNN
Samgest	-	HMM
YNL	Fisher Vector	HMM, SVM





ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting

Team	Features	Fusion	Temp. segmentation	Dimension reduction
LIRIS	RAW, SK joints	Early	Joints motion	-
CraSPN	HOG, SK	Early	Sliding windows	-
JY	SK, HOG	\mathbf{Late}	MRF	PCA
CUHK-SWJTU	Improved trajectories	-	Joints motion	PCA
Lpigou	RAW, SK joints	Early	Sliding windows	Max-pooling CNN
stevenwudi	RAW	Late	Sliding windows	-
Ismar	SK	-	Sliding windows	-
Quads	SK quads	-	Sliding windows	-
Telepoints	STIPS, SK	Late	Joints motion	-
TUM-fortiss	STIPS	Late	Joints motion	-
CSU-SCM	HOG, Skeleton	Late	Sliding windows	-
iva.mm	Skeleton, HOG	Late	Sliding windows	-
Terrier	Skeleton	-	Sliding windows	-
Team Netherlands	MHI	Early	DTW	Preserving projections
VecsRel	RAW, skeleton joints	Late	DTW	-
Samgest	Skeleton, blobs, moments	Late	Sliding windows	-
YNL	Skeleton	-	Sliding windows	-

For more details of the challenge and the results: Sergio Escalera, Xavier Baró, Jordi Gonzàlez, Miguel Ángel Bautista, Meysam Madadi, Miguel Reyes, Víctor Ponce-López, Hugo J. Escalante, Jamie Shotton, Isabelle Guyon, ChaLearn Looking at People Challenge 2014: Dataset and Results, ChaLearn Looking at People, European Conference on Computer Vision, 2014.



ECCV 2014 Workshop and Challenge on multi-modal gesture spotting, human pose recovery, action/interaction spotting









CVPR 2015Workshop and Challenge on action recognition and cultural event
recognition

Action/Interaction Recognition (second round)

Action categories





Point

Wave











Walk



Run





Hug

Shake Hands

Kiss

Fight



ChaLearn Looking at People Challenges and Workshops

CVPR 2015Workshop and Challenge on action recognition and cultural eventrecognition



Dataset	#Images	#Classes	Year
Action Classification Dataset [8]	5,023	10	2010
Social Event Dataset [11]	160,000	149	2012
Event Identification Dataset [1]	594,000	24,900	2010
Cultural Event Dataset	11,776	50	2015

First state of the art data set for cultural event recognition in still images



CVPR 2015 Workshop and Challenge on action recognition and cultural event recognition

Team name	Accuracy	Rank	Features
CUHK-SWJTU	0.507173	1	Improved trajectories [\star]
ADSC	0.501164	2	Improved trajectories [\star]
SBUVIS	0.441405	3	Improved trajectories [\star]
DonkeyBurger	0.342192	4	MHI, STIP
UC-T2	0.121565	5	Improved trajectories [\star]
MindLAB	0.008383	6	MBF

2014

Action/Interaction Track									
Rank	Team name	Score	Features Dimension reduction Clustering			Classification	Temporal coherence	Action representation	
1	MMLAB	0.5385	IDT [19]	PCA	-	SVM	-	Fisher Vector	
2	FIKIE	0.5239	IDT	PCA	-	HMM	Appearance+Kalman filter	-	
Cultural Event Trac						ck	•		
Rank	Team name	Score	Features	Features			Classification		
1	MMLAB	0.855	Multiple C	Multiple CNN			Late weighted fusion of CNNs predictions.		
2	UPC-ST	0.767	Multiple C	Multiple CNN			SVM and late weighted fusion.		
3	MIPAL_SNU	0.735	Discrimina	Discriminant regions [18] + CNNs			Entropy + Mean Probabilities of all patches		
4	SBU_CS	0.610	CNN-M [2]			SPM [III] based on LSSVM [II6]			
5	MasterBlaster	0.58	CNN			SVM, KNN, LR and One Vs Rest			
6	Nyx	0.319	Selective-s	Selective-search approach [III] + CNN			Late fusion AdaBoost		



ICCV 2015 Workshop and Challenge on action recognition, cultural event recognition, and apparent age recognition

BULLCHARGECAPE



ACTION RECOGNITION

•Large video collection (25 hours for testing) in which actions to detect are rare.

•Many examples for training (900 clips for 'Bull Charge Cape' and 500 for 'Horse Riding').

•High intra-class variability: different points of view, zoom level, action direction, color, occlusions.

•Actions are not related only to the human pose but with scene understanding.

•Videos produced during a 60 year period (1945 to 2012).



ICCV 2015 Workshop and Challenge on action recognition, cultural event recognition, and apparent age recognition

Number of	Number of	Number of	Number of test	Number of	Number of
countries for	images	categories	images	validation images	training images
all the events	per				
	category				
45	>200	90	5000	5000	15000
	Numberofcountriesforall the events45	NumberofNumber ofcountriesforimagesall the eventspercategory200	NumberofNumber ofNumber ofcountriesimagescategoriesall the eventspercategory45>20090	NumberofNumber ofNumber ofNumber ofcountriesimagescategoriesimagesall the eventspercategorycategory45>200905000	Number of countries for all the eventsNumber of imagesNumber of imagesNumber of imagesall the events categoryper

CULTURAL EVENT RECOGNITION

•First database on cultural events.

•More than 25,000 images representing 90 different categories.

•High intra- and inter-class variability.

•For this type of images, different cues can be exploited like garments, human poses, crowds analysis, objects and background scene.

•The evaluation metric will be the recognition accuracy.



ICCV 2015 Workshop and Challenge on action recognition, cultural event recognition, and **apparent age recognition**

Range of labeled ages	Information from the labelers	Contains	Contains	Number of	Number of	Number of
		real age	estimated	labelers	actors	images
			age by the			
			labelers			
0-85	Nationality, age, and gender of	YES	YES	> 3600	>2000	5000
	the labelers					

AGE ESTIMATION

•More than 5,000 faces from more than 2000 different people.

•Images with background.

•Non-controlled environments.

•Non-labeled faces neither landmarks, making the estimation problem even harder.

•One of the first datasets in the literature including estimated age labeled by many users to define the ground truth with the objective of estimating the age.

•The evaluation metric will be pondered by the mean and the variance of the labeling by the participants.

•The dataset also provides for each image the real age although not used for recognition (just for analysis purposes). In the same way for all the labelers we have their nationality, age, and gender, which will allow analyzing demographic and other interesting studies among the correlation of labelers.

INDEX

1.Computer Vision Datasets
2.Annotations and metrics
3.PASCAL Challenges
4.IMAGENET Challenges
5.ChaLearn Challenges
6.Other large scale CV DBs



Other CV datasets

Microsoft COCO Common Objects in Context

Home People

cocodataset@outlook.com

Explore Dataset

News

MS COCO Captions Challenge

Participate to the MS COCO Captioning Challenge organized with the LSUN Challenge at CVPR 2015

Click here for info



- a woman is playing a frisbee with a dog.
- a woman is playing frisbee with her large dog.
- a girl holding a frisbee with a dog coming at her.
- a woman kneeling down holding a frisbee in front of a white dog.
- a young lady is playing frisbee with her dog.

What is Microsoft COCO?



Microsoft COCO is a new image recognition, segmentation, and captioning dataset. Microsoft COCO has several features:

- Object segmentation
 - Object segmentation
- Recognition in Context
 - Multiple objects per image
- More than 300,000 images
- More than 2 Million instances
- 80 object categories
- 5 captions per image

Collaborators

Tsung-Yi Lin Cornell Tech Michael Maire TTI Chicago Serge Belongie Cornell Tech Lubomir Bourdev Facebook AI Ross Girshick Microsoft Research James Hays Brown University Pietro Perona Caltech Deva Ramanan UC Irvine Larry Zitnick Microsoft Research





facebook

Brown University

UCIrvine University of California, Irvine

Microsoft Research



Other CV datasets







http://groups.csail.mit.edu/vision/SUN/


Human Pose Recovery and Behavior Analysis Group

Other CV datasets



Welcome to the MS COCO Captioning Challenge 2015!



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

1. Introduction

The MS COCO Captioning Challenge is designed to spur the development of algorithms producing image captions that are informative and accurate. Teams will be competing against each other by training their algorithms on the MS COCO 2014 dataset and their results will be scored by **human judges**.



Other CV datasets

• CHALEARN SPEED INTERVIEWS WORLD CHALLENGE (2016?)





http://gesture.chalearn.org/speed-interviews



amazon



Other CV datasets

Platforms for CV competitions: KAGGLE

kaggle

Sign up Login

The Home of Data Science

COMPETITIONS • CUSTOMER SOLUTIONS • JOBS BOARD

Get started »





Other CV datasets

Platforms for CV competitions: CODALAB

Accelerating reproducible computational research.

CodaLab is an ecosystem for conducting computational research in a more efficient, reproducible, and collaborative manner.



Worksheets allow you to capture complex research pipelines in a reproducible way and create "executable papers". Use any data format or programming language — great for the power user!

» Explore worksheets



Competitions bring together the entire community to tackle the most challenging data and computational problems today. You can win prizes and also create your own competition.

ALPHA

» Explore competitions

Contribute your skills to help develop the CodaLab platform!



ChaLearn Looking at people news

ICCV 2015 COMPETITIONS AND WORKSHOP!! ALREADY STARTED!

WANT TO COLLABORATE IN CV CHALLENGE ORGANIZATIONS? MAIL US. There are many tasks to do.





Final remarks



ChaLearn LAP challenges and news: http://gesture.chalearn.org/

Organization of ChaLean Looking at People requires:

- -Good ideas to solve real problems focused on humans
- -Collecting data
- -Labeling tools
- -Dissemination and repositories
- -Baseline designs based on state of the art approaches
- -Online platform for the competition
- -Sponsoring
- -Presentation of the results in a relevant events

-Organization of special issues and challenge report documents, making competition data public for the scientific community

For each competition many organizers contribute. Our plan is to perform yearly challenges. Feel free to contact us if you want to be included in our ChaLearn LAP mailing list or collaborate in some aspect propose ideas related to ChaLearn Looking at People competitions:

sergio@maia.ub.es





Thank you!



XX CENTURY

XXI CENTURY ?