

---

# TRECVID-2008 High-Level Feature task: Overview

---

Wessel Kraaij  
TNO // Radboud University

George Awad  
NIST

# Outline

---

- Task summary
- Evaluation details
  - Inferred Average precision
  - Participants
- Evaluation results
  - Pool analysis
  - Results per category
  - Results per feature
  - Significance tests category A
  - comparison with TV2007
- Global Observations
- Issues

# High-level feature task (1)

---

- Goal: Build benchmark collection for visual concept detection methods
- Secondary goals:
  - encourage generic (scalable) methods for detector development
  - semantic annotation is important for search/browsing
- Participants submitted runs for 20 LSCOM features
- TRECVID 2008 video data
  - Netherlands Institute for Sound and Vision (~**200 hours** of news magazine, science news, news reports, documentaries, educational programming and archival video in MPEG-1).
  - 100 hours for development.
  - 100 hours for test.
  - TRECVID 2003, 2005 & TRECVID 2007 annotated data.

# High-level feature task (2)

---

- NIST evaluated 20 features using a 50% random sample of the submission pools (Inferred AP)
- Six training types were allowed
  - A : Systems trained on only common TRECVID development collection data
  - B : Systems trained on only common development collection data but not on (just) common annotation of it.
  - C : System is not of type A or B
  - a : same as A but no training data specific to any sound and vision data has been used
  - b : same as B but no training data specific to any sound and vision data has been used
  - c : same as C but no training data specific to any sound and vision data has been used

# TV2007 vs TV2008 dataset

---

	TV2007	TV2008
Dataset length (hours)	~100	~200
Number of shots	18,142	35,766
Number of unique program titles	47	77

# 20 LSCOM features evaluated

---

- 1 Classroom
- 2 Bridge
- 3 Emergency\_Vehicle
- 4 Dog
- 5 Kitchen
- 6 Airplane\_flying
- 7 Two people
- 8 Bus
- 9 Driver
- 10 Cityscape
- 11 Harbor
- 12 Telephone
- 13 Street
- 14 Demonstration\_Or\_Protest
- 15 Hand
- 16 Mountain
- 17 Nighttime
- 18 Boat\_ship
- 19 Flower
- 20 Singing

Features were selected to be better suited to sound and vision data

# Evaluation

---

- Each feature assumed to be binary: absent or present for each master reference shot
- Task: Find shots that contain a certain feature, rank them according to confidence measure, submit the top 2000
- NIST pooled and judged top results from all submissions
- Evaluated performance effectiveness by calculating the *inferred average precision* of each feature result
- Compared runs in terms of **mean** *inferred average precision* across the 20 feature results.

# Inferred average precision (infAP)

---

- Developed\* by Emine Yilmaz and Javed A. Aslam at Northeastern University
- Estimates average precision surprisingly well using a surprisingly small sample of judgments from the usual submission pools
- Experiments on TRECVID 2005 & 2006 & 2007 feature submissions confirmed quality of the estimate in terms of actual scores and system ranking

\* J.A. Aslam, V. Pavlu and E. Yilmaz, *Statistical Method for System Evaluation Using Incomplete Judgments* Proceedings of the 29th ACM SIGIR Conference, Seattle, 2006.

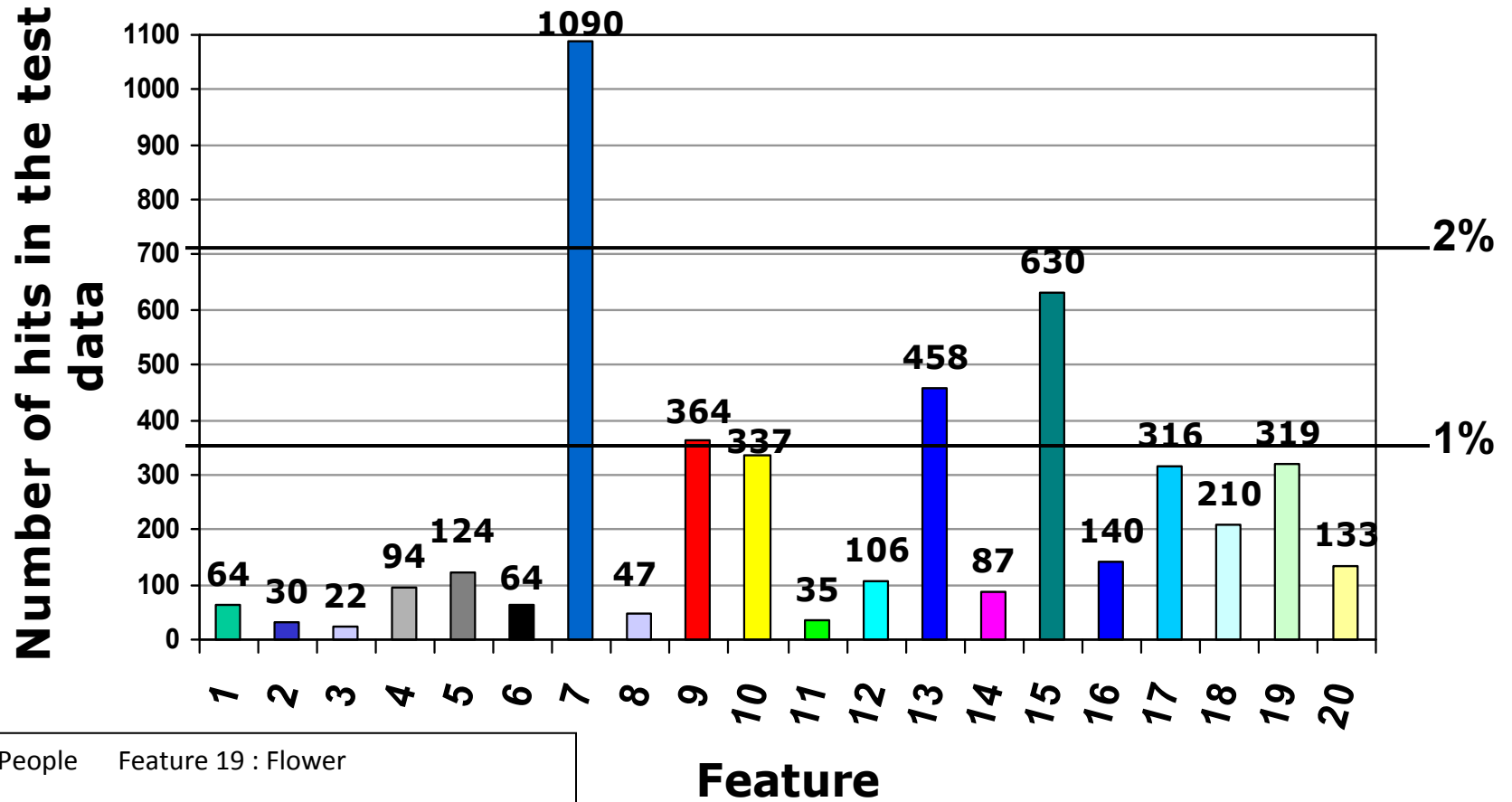


# 2008: Inferred average precision (infAP)

---

- Submissions for each of 20 features were pooled down to about average 130 items (so that each feature pool contained ~ 6777 shots)
  - varying pool depth per feature
- A 50% random sample of each pool was then judged:
- 67,774 total judgments (TV7: 66,293)
- Judgment process: one assessor per feature, watched complete shot while listening to the audio.
- infAP was calculated using the judged and unjudged pool by trec\_eval

# Frequency of hits varies by feature



Feature 7 : 2 People    Feature 19 : Flower  
Feature 15 : Hand    Feature 17 : Nighttime  
Feature 13 : Street  
Feature 9 : Driver  
Feature 10 : Cityscape

# 2008 : 43/115 Participants

---

Asahikasei Co.	-- ** FE RU --
Bilkent University	CD -- FE -- **
Brno University of Technology	CD ED FE ** SE
Beijing University of Posts and Telecommunications	CD ** FE -- --
Carnegie Mellon University	-- ED FE RU **
Columbia University	CD -- FE -- SE
COST292 Team (Delft Univ.)	CD ** FE RU SE
Florida International Univ.	-- ** FE -- --
Fudan University	CD ED FE -- SE
IBM T. J. Watson Research Center	CD ** FE ** SE
INRIA-LEAR	CD -- FE -- --
MMIS (Open Univ.)	** -- FE -- SE
Microsoft Research Asia	** ** FE ** SE
NHKSTRL	** ED FE RU **
National Institute of Informatics	CD ** FE RU SE
IRIM	** ** FE RU **
ISM (The Institute of Statistical Mathematics)	-- -- FE -- --
IUPR-DFKI	** -- FE -- --

\*\* : group didn't submit any runs

-- : group didn't participate

# 2008: 43 Participants (continued)

---

JOANNEUM RESEARCH Forschungsgesellschaft mbH	** ** FE RU --
LIG (Laboratoire d'Informatique de Grenoble)	** -- FE -- **
Laboratoire LIRIS (LYON)	-- ** FE ** **
University of Twente and CWI	-- -- FE -- SE
LSIS_GLOT (CNRS LSIS)	-- -- FE -- --
Marburg	** ** FE ** **
MCG-ICT-CAS	CD ED FE -- SE
Mediamill (Univ. of Amsterdam)	-- ** FE -- SE
MESH	-- -- FE - SE
National Taiwan University	** ** FE -- SE
Oxford Univ.	** -- FE -- SE
PKU-ICST (Peking Univ.)	** ** FE ** SE
PicSom(Helsinki University of Technology)	CD -- FE RU SE
Queensland University of Technology	-- -- FE RU --
REGIM	-- ** FE RU SE
SJTU	-- ED FE -- SE
SP-UC3M (Universidad Carlos III de Madrid)	-- -- FE -- SE
Thu-intel	CD ** FE RU SE
Tokyo Institute of Technology	-- ED FE RU --
University of Electro-Communications	** ** FE RU **
University of Karlsruhe (TH)	-- -- FE -- --



# Number of runs of each training type

Tr-Type	2008	2007
A	152 (76%)	146 (89.5%)
B	15 (7.5%)	7 (4.3%)
C	22(11%)	6 (3.7%)
a	9(4.5%)	4 (2.5%)
b	0	0
c	2(1%)	0
Total runs	200	163

- 1- More interest in special training data (B, C).
- 2- More interest in un-related Data (a, c).
- 3- The common data (A) still is the most popular.

System training type:

**A** - Only on common dev. collection and the common annotation.

**B** - Only on common dev. collection but not on (just) the common annotation.

**C** - not of type A or B.

**a , b, c** – Same as A, B, & C respectively but without using any specific training data from Sound and Vision dataset.

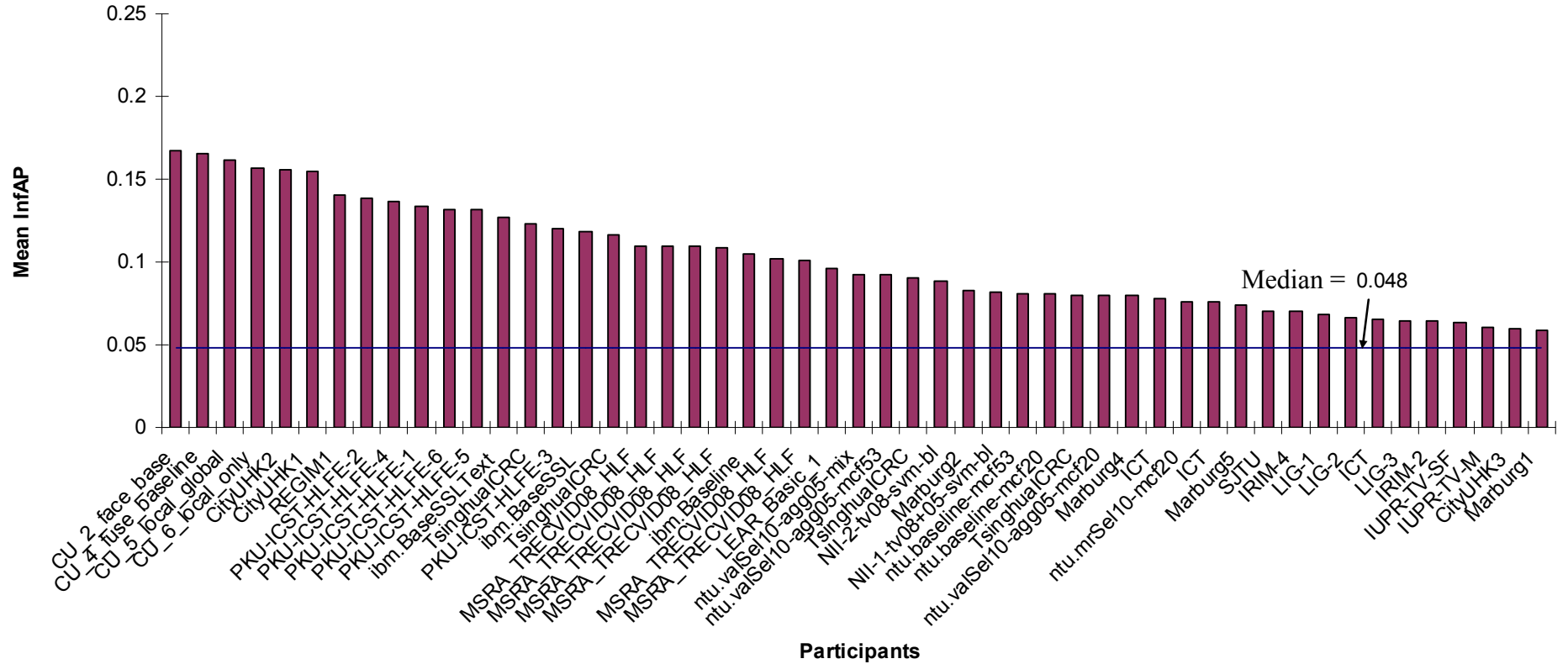
# True shots contributed uniquely **by team** for each feature

---

- BUPT - 1 shot (feature 7 : 2 people)
- IUPR-DFKI – 1 shot (feature 15 : hand)
- NHKSTRL – 1 shot (feature 17 : Nighttime)
- Queensland University – 1 shot (feature 18 : Boat\_Ship)
- Institute of Image Communication and Information Processing – 1 shot (feature 8 : Bus)
- Asahikasei – 1 shot (feature 7 : 2 people)
- Delft University of Technology – 2 shots (feature 12 : Telephone , feature 13 : Street)
- CNRS LSIS – 2 shots (feature 7 : 2 people)
- MESH – 2 shots (feature 15 : hand , feature 16 : Mountain)
- UTC – 2 shots (feature 7 : 2 people)
- École Nationale d'Ingénieurs de Sfax ENIS – 6 shots (feature 17 : nighttime , feature 4 : dog , feature 12 : Telephone)

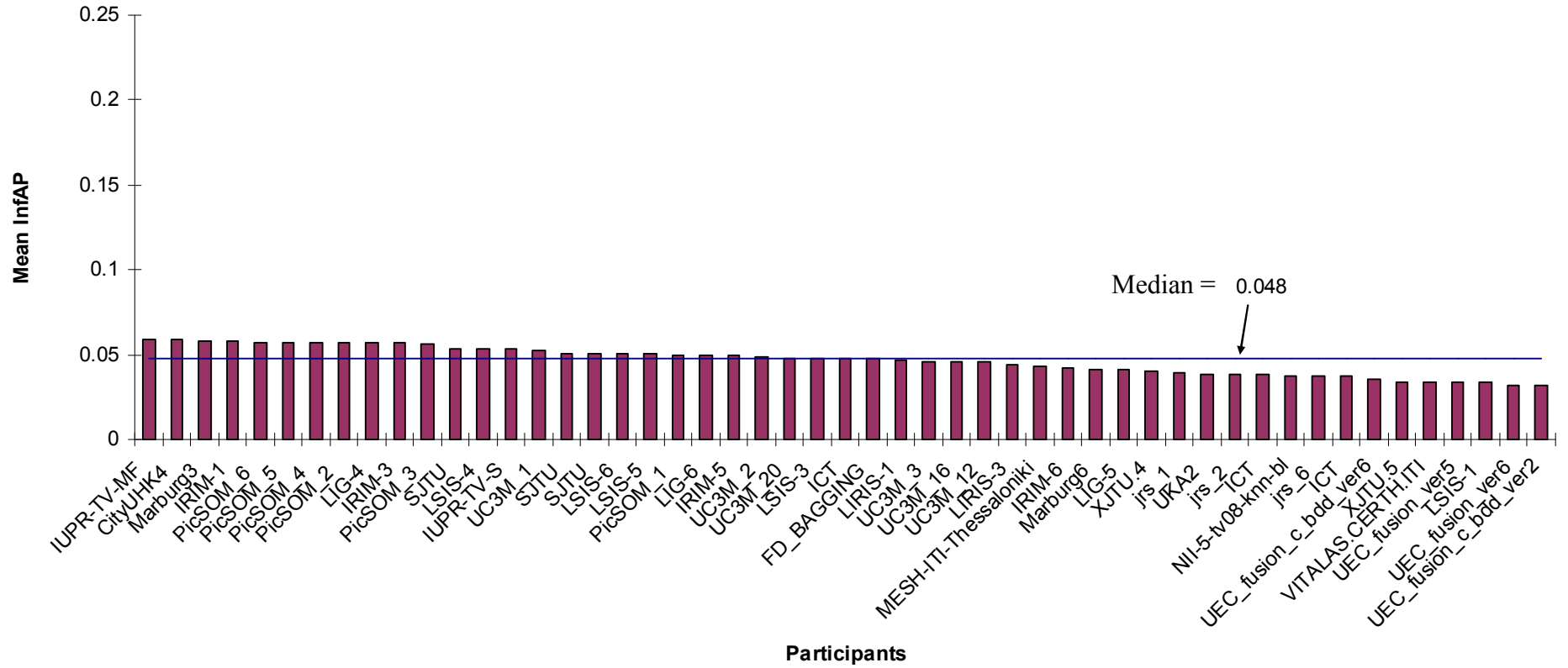
-Unlike TRECVID 2007 where only two groups found different unique true shots.

Category A results (Top 50)

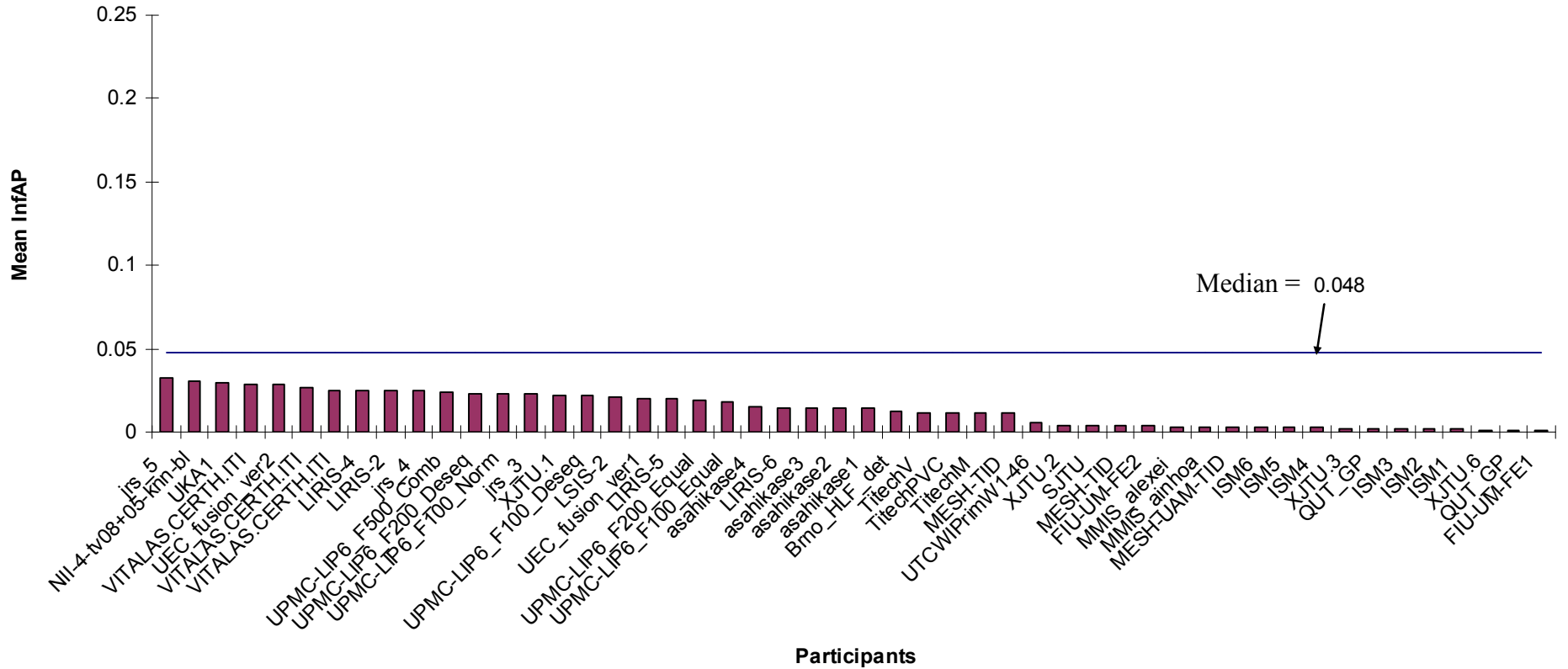




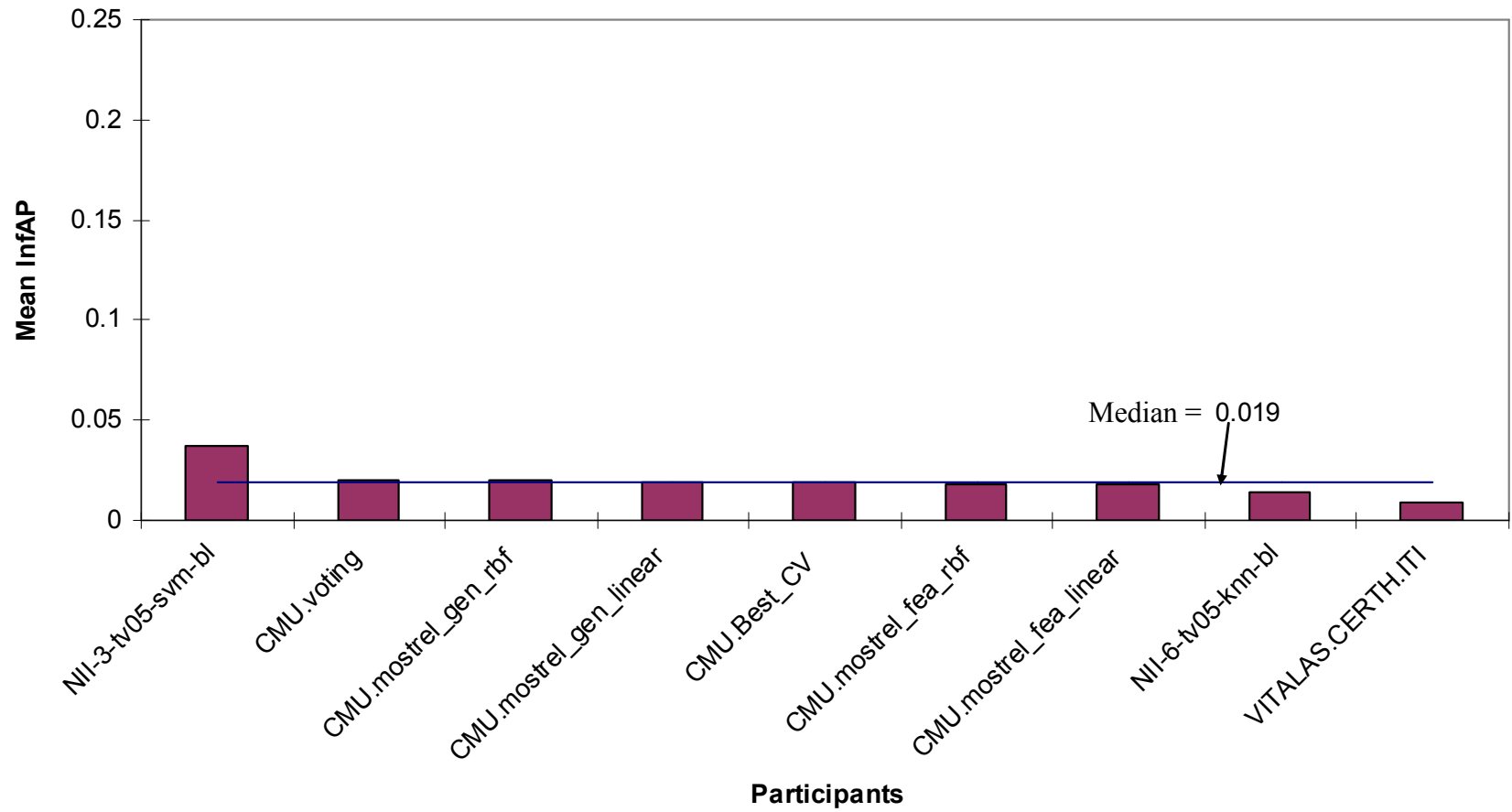
### Category A results (Middle 50)

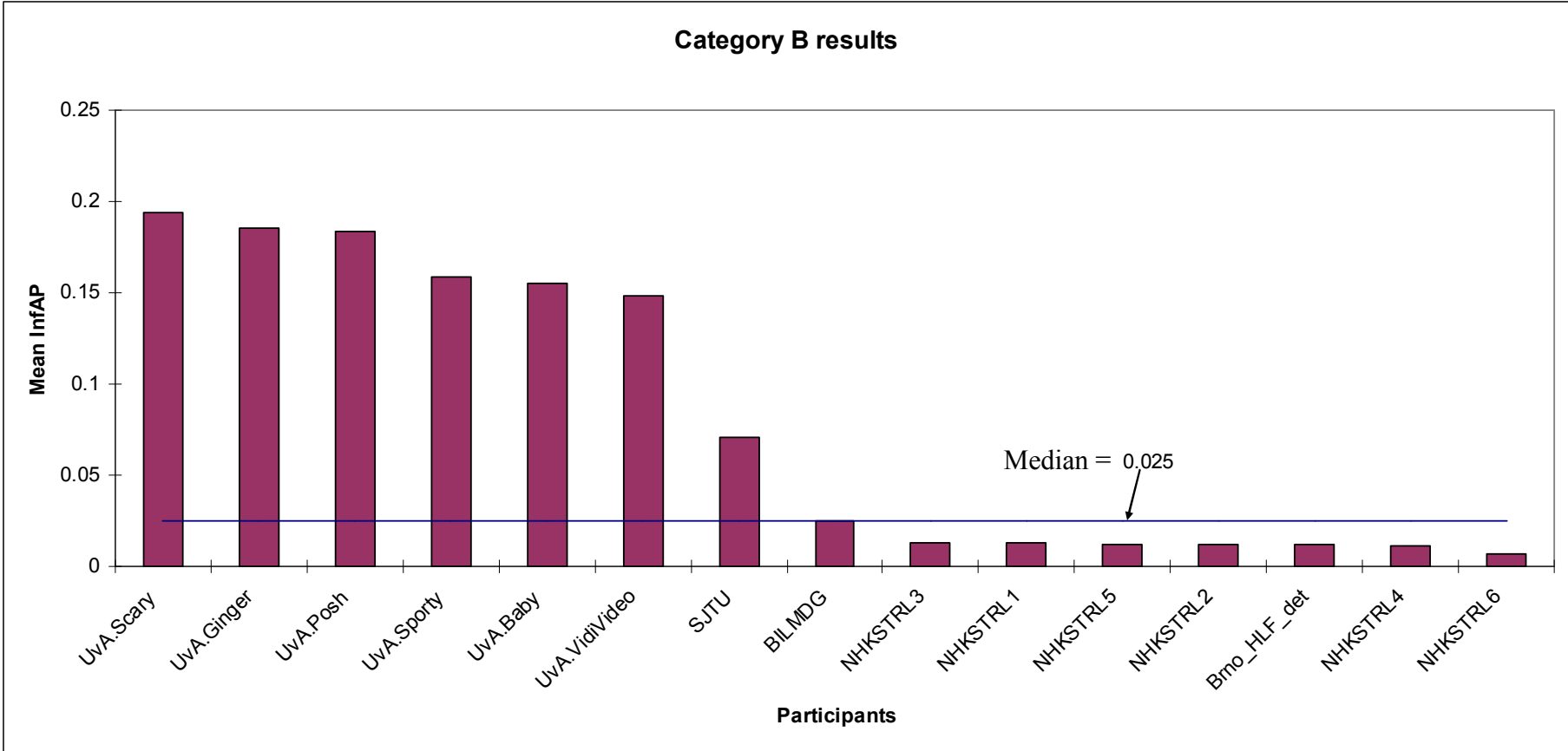


### Category A results (Bottom 50)

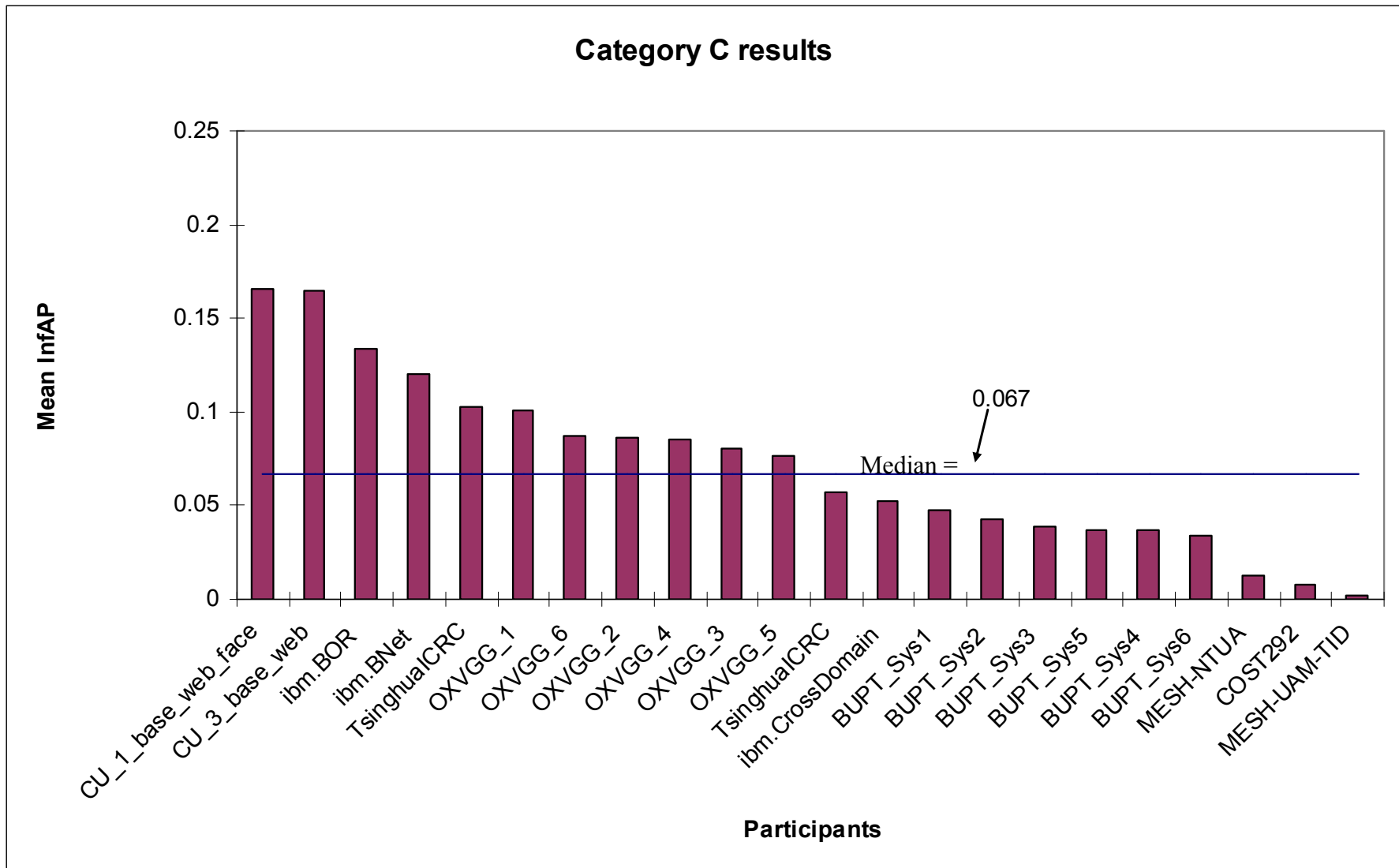


### Category a results





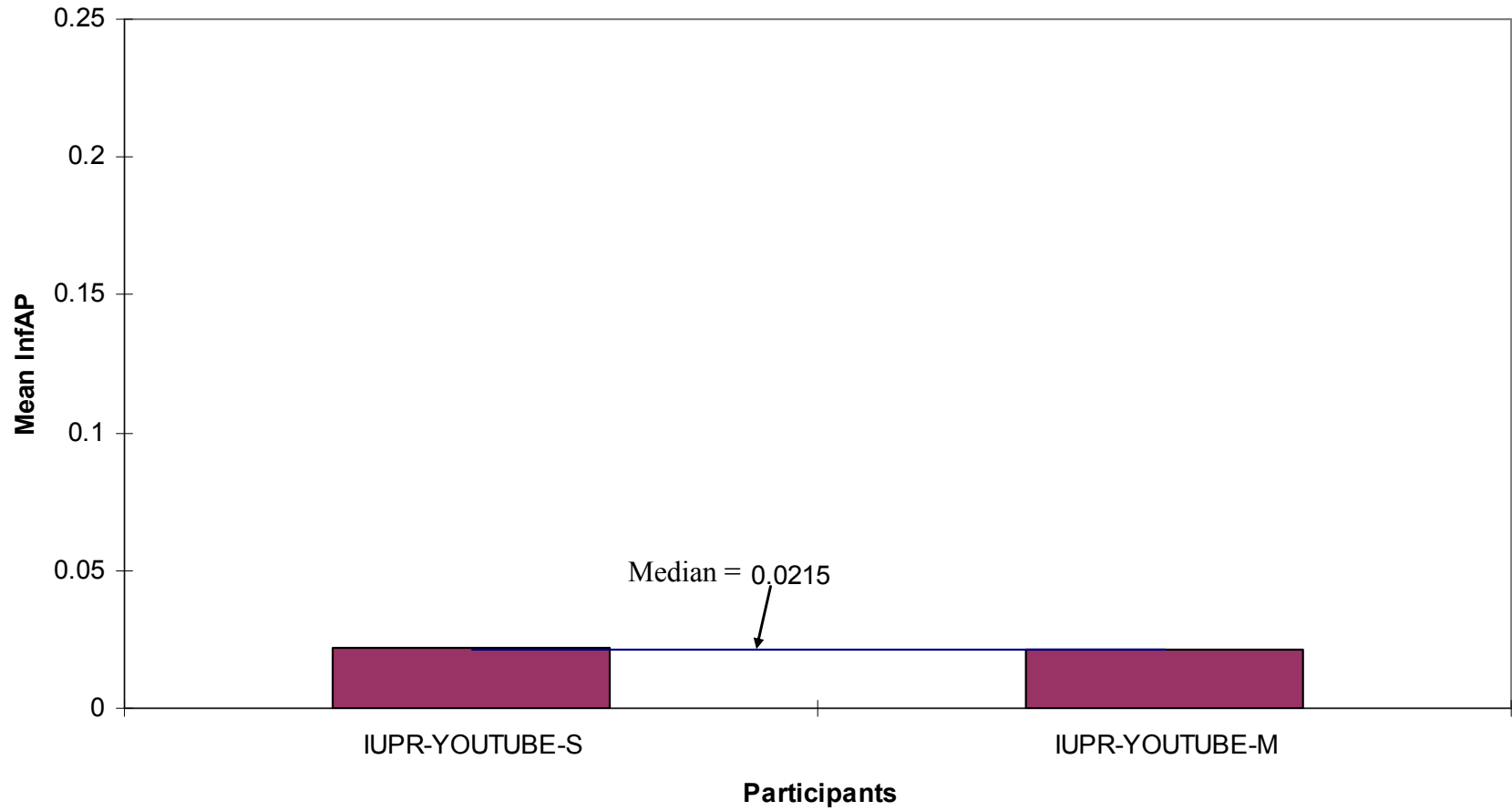
First year where B runs are better than A runs



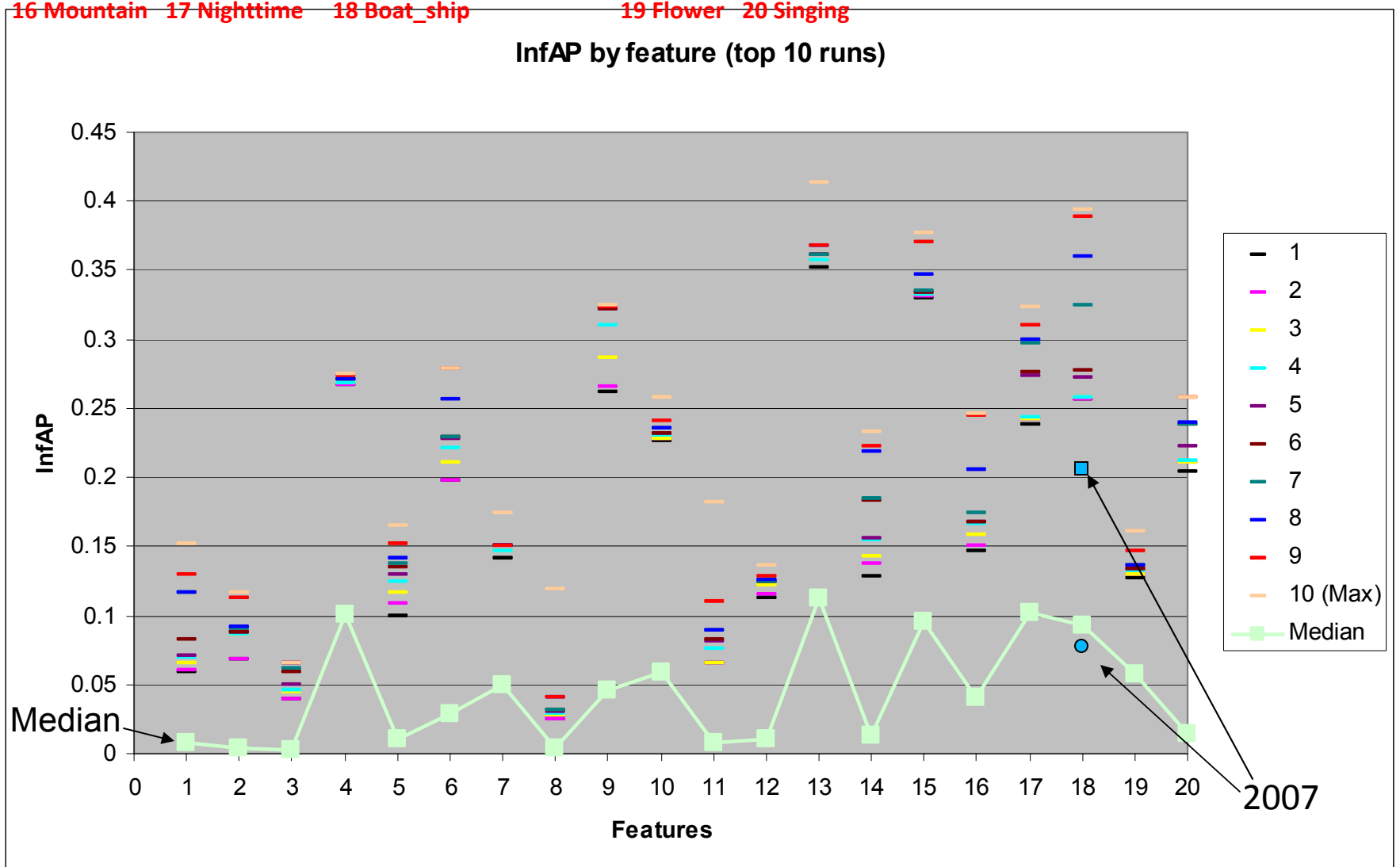
First year where C runs are on par with A runs

Data from e.g. Flickr, Youtube, Peekaboom

### Category c results



- 1 Classroom    2 Bridge    3 Emergency\_Vehicle    4 Dog    5 Kitchen    6 Airplane\_flying    7 Two people    8 Bus
- 9 Driver    10 Cityscape    11 Harbor    12 Telephone    13 Street    14 Demonstration\_Or\_Protest    15 Hand
- 16 Mountain    17 Nighttime    18 Boat\_ship    19 Flower    20 Singing



Which, if any, differences are significant, i.e. not due to chance?

# Significant differences among top 10 A-category runs (using randomization test, $p < 0.05$ )

---

## Run name (mean infAP) >

- CU\_2\_face\_base\_2 (0.167)
  - CU\_4\_fuse\_baseline\_4 (0.165)
  - CU\_5\_local\_global\_5 (0.162)
  - CU\_6\_local\_only\_6 (0.157)
  - CityUHK2\_2 (0.156)
  - CityUHK1\_1 (0.155)
  - REGIM1\_1 (0.140)
  - PKU-ICST-HLFE-2\_2 (0.138)
  - PKU-ICST-HLFE-4\_4 (0.137)
  - PKU-ICST-HLFE-1\_1 (0.134)
- > CU\_2\_face\_base\_2
    - > CU\_6\_local\_only\_6
      - > PKU-ICST-HLFE-1\_1
      - > PKU-ICST-HLFE-2\_2
      - > PKU-ICST-HLFE-4\_4
  - > CU\_4\_fuse\_baseline\_4
    - > CU\_6\_local\_only\_6
      - > PKU-ICST-HLFE-1\_1
      - > PKU-ICST-HLFE-2\_2
      - > PKU-ICST-HLFE-4\_4
  - > CU\_5\_local\_global\_5
    - > CU\_6\_local\_only\_6
      - > PKU-ICST-HLFE-1\_1
      - > PKU-ICST-HLFE-2\_2
      - > PKU-ICST-HLFE-4\_4
  - > CityUHK2\_2
    - > PKU-ICST-HLFE-1\_1
    - > PKU-ICST-HLFE-4\_4
  - > CityUHK1\_1
    - > PKU-ICST-HLFE-1\_1



# Significant differences among top 10 a-category runs (using randomization test, $p < 0.05$ )

---

- Run name (mean infAP)
  - NII-3-tv05-svm-bl\_3 (0.037)
  - CMU.voting\_6 (0.020)
  - CMU.mostrel\_gen\_rbf\_2 (0.020)
  - CMU.mostrel\_gen\_linear\_1 (0.019)
  - CMU.Best\_CV\_5 (0.019)
  - CMU.mostrel\_fea\_linear\_3 (0.018)
  - CMU.mostrel\_fea\_rbf\_4 (0.018)
  - NII-6-tv05-knn-bl\_6 (0.014)
  - VITALAS.CERTH.ITI\_5 (0.009)
- NII-3-tv05-svm-bl\_3
    - CMU.voting\_6
      - cMU.mostrel\_fea\_linear\_3
      - cMU.mostrel\_fea\_rbf\_4
      - VITALAS.CERTH.ITI\_5
    - CMU.Best\_CV\_5
      - VITALAS.CERTH.ITI\_5
    - cMU.mostrel\_gen\_rbf\_2
      - cMU.mostrel\_fea\_linear\_3
      - VITALAS.CERTH.ITI\_5
    - cMU.mostrel\_gen\_linear\_1
    - NII-6-tv05-knn-bl\_6

# Significant differences among top 10 B-category runs (using randomization test, $p < 0.05$ )

---

Run name (mean infAP)	UvA.Scary_1
□ UvA.Scary_1 (0.194)	➤ UvA.Ginger_4 , UvA.Posh_3
□ UvA.Ginger_4 (0.185)	➤ UvA.Sporty_2
□ UvA.Posh_3 (0.184)	➤ SJTU_5
□ UvA.Sporty_2 (0.159)	➤ BILMDG_1
□ UvA.Baby_5 (0.155)	➤ NHKSTRL3_3
□ UvA.VidiVideo_6 (0.148)	➤ NHKSTRL1_1
□ SJTU_5 (0.071)	➤ UvA.Baby_5
□ BILMDG_1 (0.025)	➤ SJTU_5
□ NHKSTRL3_3 (0.013)	➤ BILMDG_1
□ NHKSTRL1_1 (0.013)	➤ NHKSTRL3_3
	➤ NHKSTRL1_1
	➤ UvA.VidiVideo_6
	➤ SJTU_5
	➤ BILMDG_1
	➤ NHKSTRL3_3
	➤ NHKSTRL1_1

# Significant differences among top 10 C-category runs (using randomization test, $p < 0.05$ )

---

## Run name (mean infAP)

- CU\_1\_base\_web\_face\_1 (0.166) ■ CU\_1\_base\_web\_face\_1 ,
- CU\_3\_base\_web\_3 (0.165) ■ CU\_3\_base\_web\_3
- ibm.BOR\_1 (0.134) > ibm.BOR\_1
- lbm.BNet\_2 (0.120) > OXVGG\_1\_1
- TsinghuaCRC\_4 (0.103) > lbm.BNet\_2
- OXVGG\_1\_1 (0.101) > OXVGG\_6\_6
- OXVGG\_6\_6 (0.087) > OXVGG\_3\_3
- OXVGG\_2\_2 (0.086) > OXVGG\_4\_4
- OXVGG\_4\_4 (0.085) > OXVGG\_2\_2
- OXVGG\_3\_3 (0.080) > OXVGG\_3\_3
- OXVGG\_3\_3 (0.080) > TsinghuaCRC\_4
- OXVGG\_3\_3 (0.080) > OXVGG\_3\_3

# Significant differences among A/a category runs by group (using randomization test, $p < 0.05$ )

---

## Run name (mean infAP)

- A\_NII-2-tv08-svm-bl\_2 (0.088) ➤ A\_NII-1-tv08+05-svm-bl\_1, A\_NII-2-tv08-svm-bl\_2
- A\_NII-1-tv08+05-svm-bl\_1 (0.082) ➤ A\_NII-4-tv08+05-knn-bl\_4
- A\_NII-5-tv08-knn-bl\_5 (0.037) ➤ a\_NII-6-tv05-knn-bl\_6
- a\_NII-3-tv05-svm-bl\_3 (0.037) ➤ A\_NII-5-tv08-knn-bl\_5
- A\_NII-4-tv08+05-knn-bl\_4 (0.031) ➤ a\_NII-6-tv05-knn-bl\_6
- a\_NII-6-tv05-knn-bl\_6 (0.014) ➤ a\_NII-3-tv05-svm-bl\_3
- a\_NII-6-tv05-knn-bl\_6

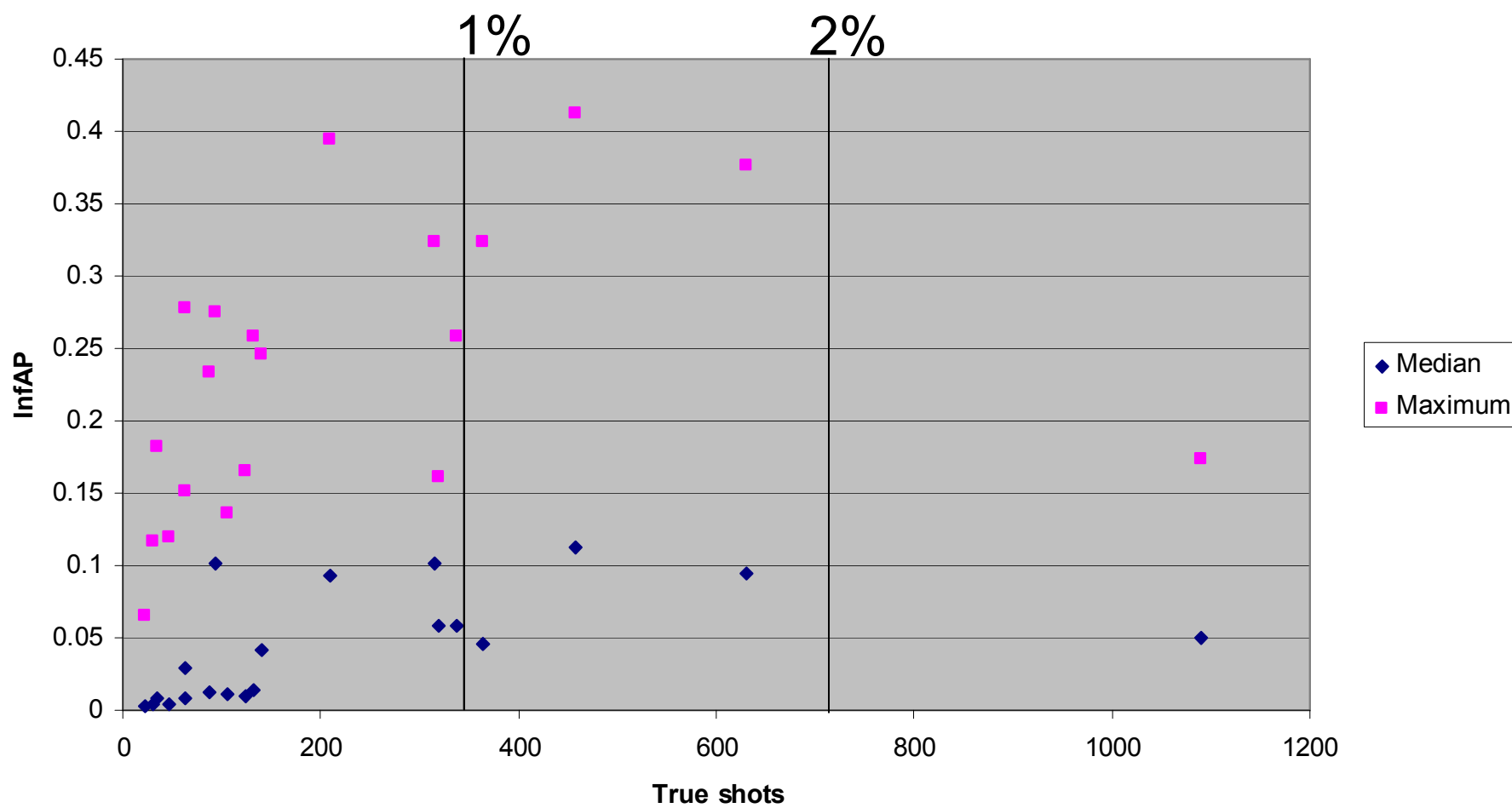
# Significant differences among A/a category runs by group (using randomization test, $p < 0.05$ )

---

## Run name (mean infAP)

- A\_VITALAS.CERTH.ITI\_2 (0.034)
  - A\_VITALAS.CERTH.ITI\_1 (0.029)
  - A\_VITALAS.CERTH.ITI\_4 (0.027)
  - A\_VITALAS.CERTH.ITI\_3 (0.025)
  - a\_VITALAS.CERTH.ITI\_5 (0.009)
- A\_VITALAS.CERTH.ITI\_2
    - A\_VITALAS.CERTH.ITI\_3
      - a\_VITALAS.CERTH.ITI\_5
    - A\_VITALAS.CERTH.ITI\_1
      - a\_VITALAS.CERTH.ITI\_5
    - A\_VITALAS.CERTH.ITI\_4
      - a\_VITALAS.CERTH.ITI\_5

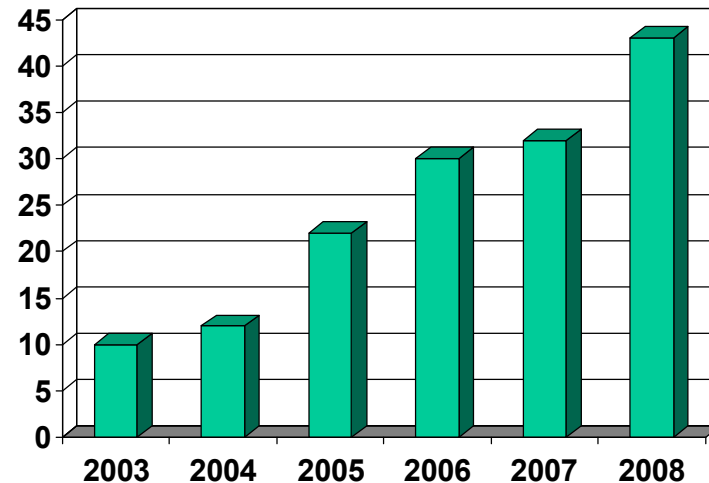
InfAP vs true shots in test data (across 20 features)



# General observations (1)

---

- Participation is still increasing
- Accepted as an important building block for search
- More interest in cat B and C submissions (using e.g. web data)
- Submissions in cat B achieve best performance
- Submissions in cat C are on par with cat A



# General observations (2)

---

- Hardly any feature specific approaches
- Large variety in classifier architectures and choices of feature representations
- Hardware: usually a single, cpu, but several medium and larger clusters
- Nr of classifiers used for fusion ranges between 1 and >1160
- Testing times vary between 10m and 150h per feature.
- Approx. 30% of the runs do some form of temporal analysis
- Approx 50% of the runs use salient/sift points
- Compiled metadata will be made available to participants