# TRECVID-2005 Low-level (camera motion) feature task

Wessel Kraaij

TNO

&

Tzveta Ianeva

NIST

# Task definition

ØTRECVID 2005 pilot task

ØAbility to detect camera movement features:

    q Pan (left or right ) or track

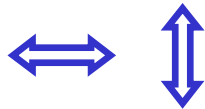    q Tilt (up or down) or boom

    q Zoom (in or out) or dolly

# Task definition ...

Ø Camera movement features are usually combined

    q Pan & Tilt

    q Pan & Zoom

    q Tilt & Zoom
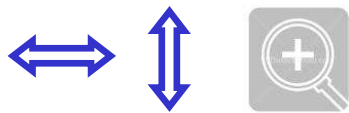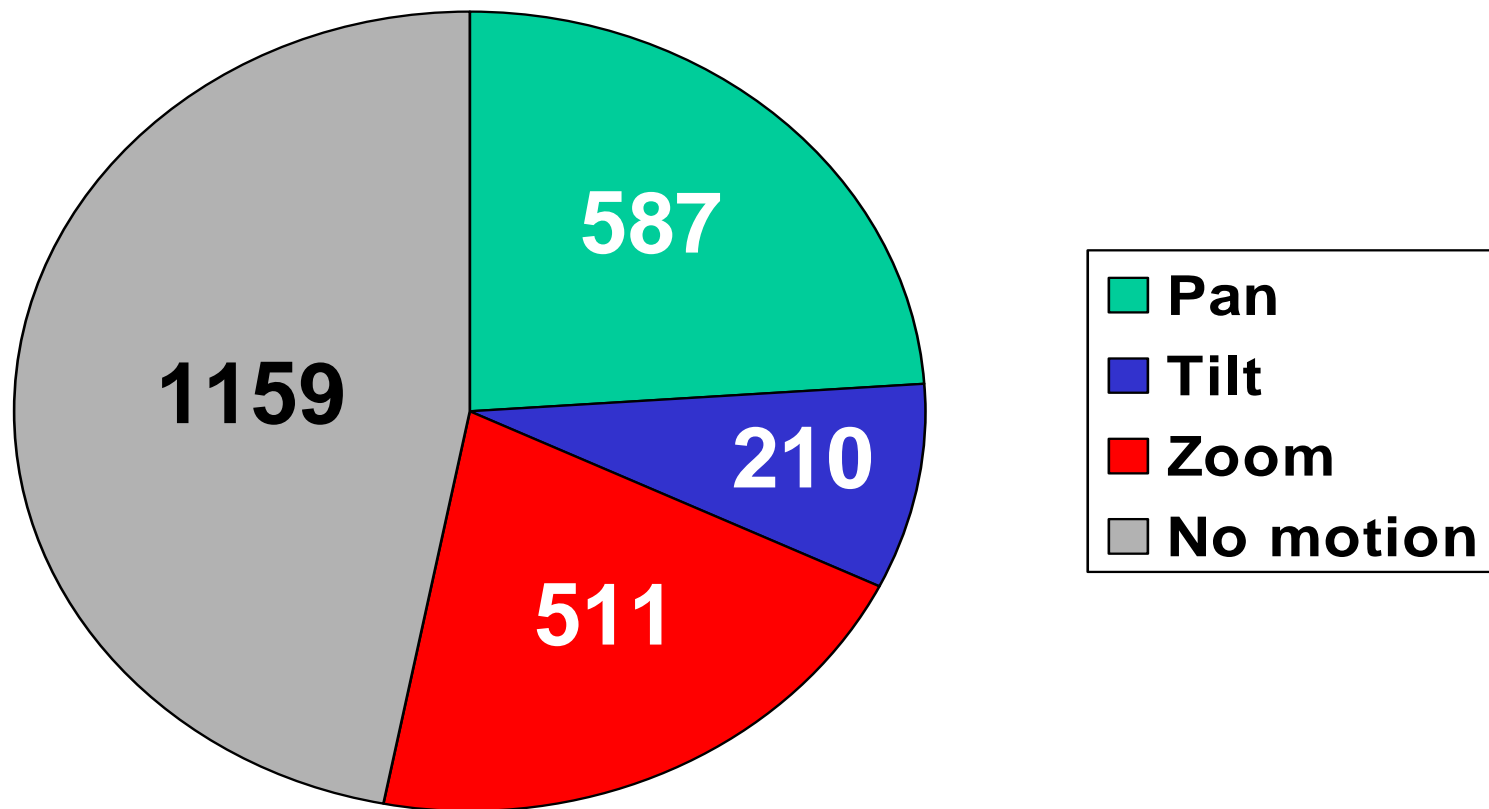
# Task definition ...

q  Pan & Tilt & Zoom



Ø  Submissions provide complete judgments for test set by specifying all shots identified as positive by the system

Ø  No Training data provided by NIST

Ø  Tool to create development data developed by Werner Bailer at Joanneum Researh

# Ground truth creation at NIST

Ø  Watch randomly chosen subset of test data (~5000 shots)

Ø  Keep only shots with "clear" examples of (no) motion (~2226)

Ø  No-motion shots seem to more clearly exhibit no motion
than shots with motion features exhibit motion⌐ *#FP will tend to be small, #FN will tend to be high*

Ø  Define test subset for each feature by combining

  Ø  shots exhibiting the feature

  Ø  shots exhibiting no motion (same for all features)

Ø  No adjustments to subset sizes or true:false ratios

  Ø  Pan       587:1159

  Ø  Tilt        210:1159

  Ø  Zoom    511:1159

# Truth data distribution (number of shots)

# Truth and evaluation issues

Ø**Why feature groups?**

ØPerceptual limits in truth creation

ØCost of creating truth data

ØMany shots with lots of small camera movement – not what's wanted when user asks for a "pan", etc.

Ø**Implications of test set construction on measures**

ØLack of randomness makes generalization hard

ØVarying true:false ratios make precision harder for tilt than pan and zoom

ØGreater clarity of no-motion shots would make false positive less likely then false negatives and higher precision easier to achieve than higher recall

# No motion shots

# Truth data costly to create – lot's of shaky shots



Hard to judge



Not what a user wants

# 12 Participating Groups

Carnegie Mellon University ( CMU )  - USA

City University of Hong Kong ( CUHK ) - China

Fudan University ( FUDAN ) - China

Institute for Infocomm Research ( IIR ) - Singapore

JOANNEUM RESEARCH ( Joanneum ) - Austria

KDDI & R&D Laboratories, Inc. ( KDDI ) - Japan

LaBRI ( LaBRI ) - France

Tsinghua University ( Tsinghua ) - China

University of Central Florida / University of Modena ( UCF ) – USA/Italy

University of Iowa ( Uiowa ) - USA

University of Marburg ( MARBURG ) - Germany

Univ. of Amsterdam & TNO ( UvA ) - Netherlands

# NIST baseline runs

ØAll features true for all shots (TrueForAllShots)

ØRandom run with true distribution of Pan, Tilt, Zoom as in truth data (TruthDataDistrib)
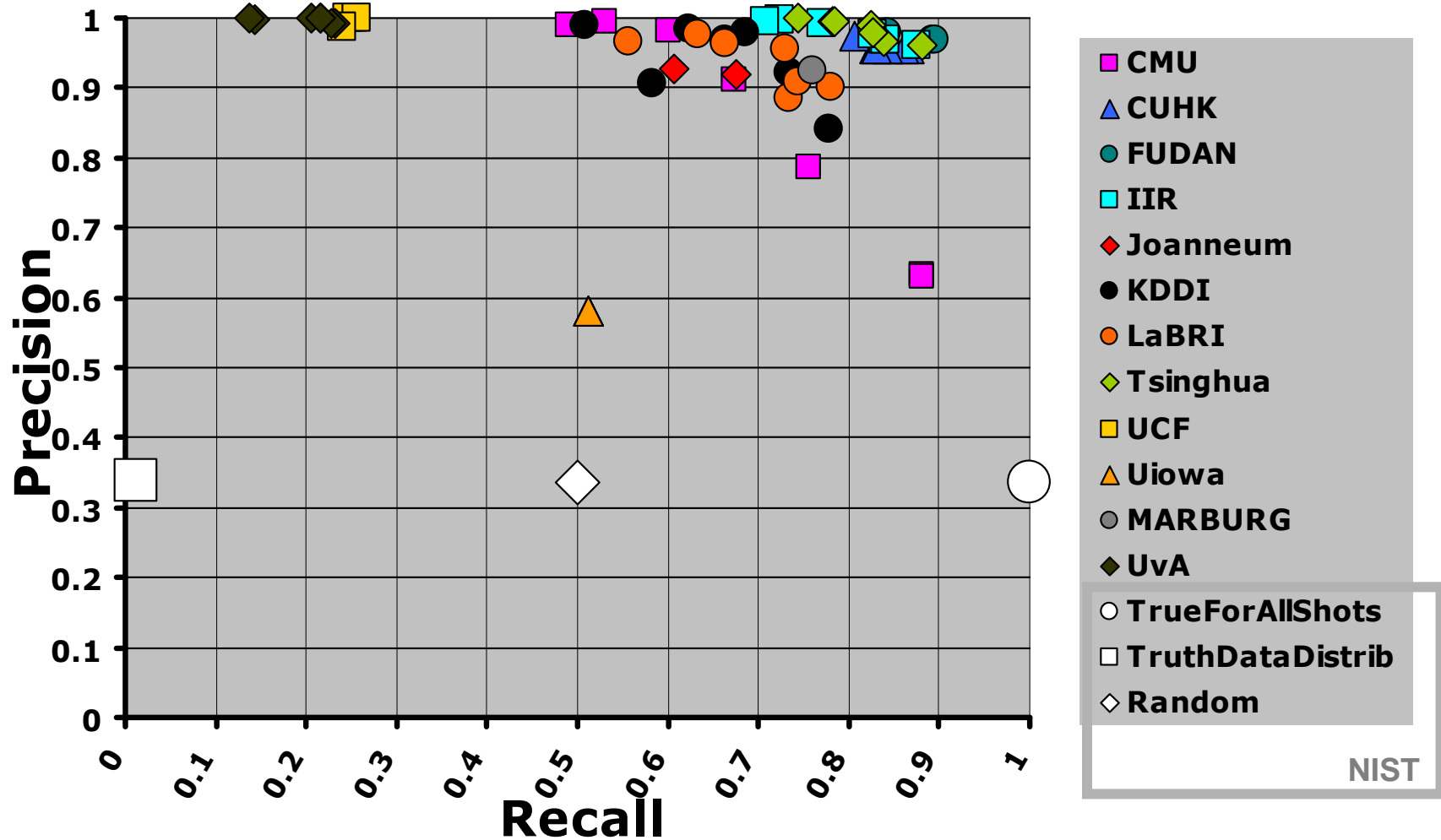
ØFeatures randomly true/false for each shot (Random)

# Evaluation Measures

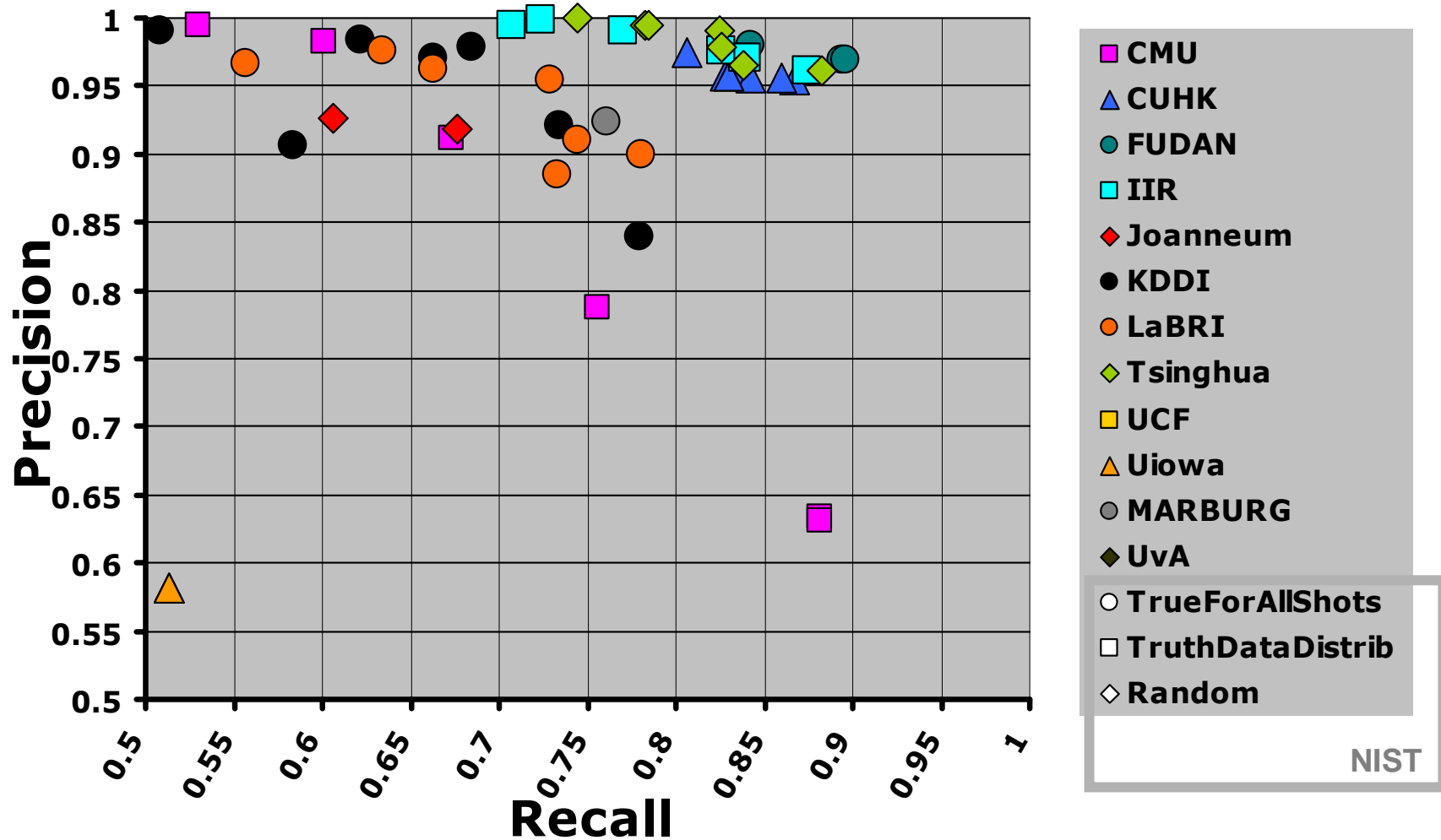$$\text{Precision} = \frac{\text{\# True positives}}{\text{\# True positives} + \text{\# False positives}}$$

$$\text{Recall} = \frac{\text{\# True positives}}{\text{\# True positives} + \text{\# False negatives}}$$

*Given the imbalance in class properties, it's easier to achieve a high precision than a high recall. The use of $F_{\beta=1}$ seems not appropriate*
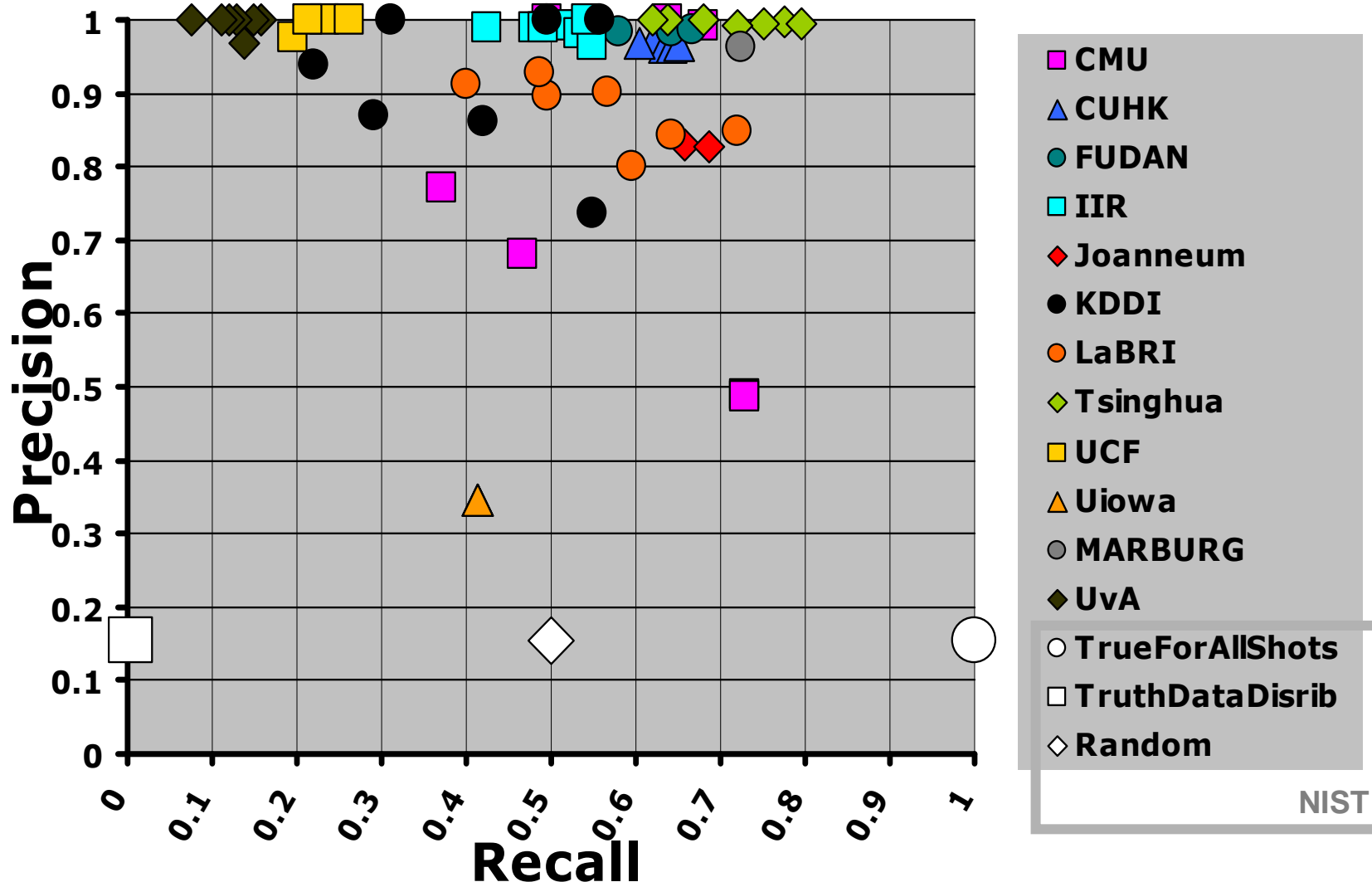
# **Pan**: recall and precision by system

# **Pan**: recall and precision by system (zoomed)

# **Tilt**: recall and precision by system



**Legend:**
- ■ CMU
- ▲ CUHK
- ● FUDAN
- ■ IIR
- ◆ Joanneum
- ● KDDI
- ● LaBRI
- ◆ Tsinghua
- ■ UCF
- ▲ Uiowa
- ● MARBURG
- ◆ UvA
- ○ TrueForAllShots
- □ TruthDataDisrib
- ◇ Random

NIST

# Tilt: recall and precision by system (zoomed)

# Zoom: recall and precision by system



Legend:
- ■ CMU
- ▲ CUHK
- ● FUDAN
- ■ IIR
- ◆ Joanneum
- ● KDDI
- ● LaBRI
- ◆ Tsinghua
- ■ UCF
- ▲ Uiowa
- ● MARBURG
- ◆ UvA
- ○ TrueForAllShots
- □ TruthDataDistrib
- ◇ Random

NIST

# **Zoom**: recall and precision by system (zoomed)



Legend:
- ■ CMU
- △ CUHK
- ● FUDAN
- □ IIR
- ◆ Joanneum
- ● KDDI
- ● LaBRI
- ◇ Tsinghua
- ■ UCF
- △ Uiowa
- ● MARBURG
- ◆ UvA
- ○ TrueForAllShots
- □ TruthDataDistrib
- ◇ Random

NIST

# **Mean** recall and precision over all 3 features by system



**Legend:**
- ■ CMU
- ▲ CUHK
- ● FUDAN
- ■ IIR
- ◆ Joanneum
- ● KDDI
- ● LABRI
- ◆ Tsingua
- ■ UCF
- ▲ Uiowa
- ● MARBURG
- ◆ UvA
- ○ TrueForAllShots
- □ TruthDataDistrib
- ◇ Random

NIST

# **Mean** recall and precision over all 3 features by system (zoomed)



Legend:
- ■ CMU
- ▲ CUHK
- ● FUDAN
- ■ IIR
- ◆ Joanneum
- ● KDDI
- ● LABRI
- ◆ Tsingua
- □ UCF
- △ Uiowa
- ● MARBURG
- ◆ UvA
- ○ TrueForAllShots
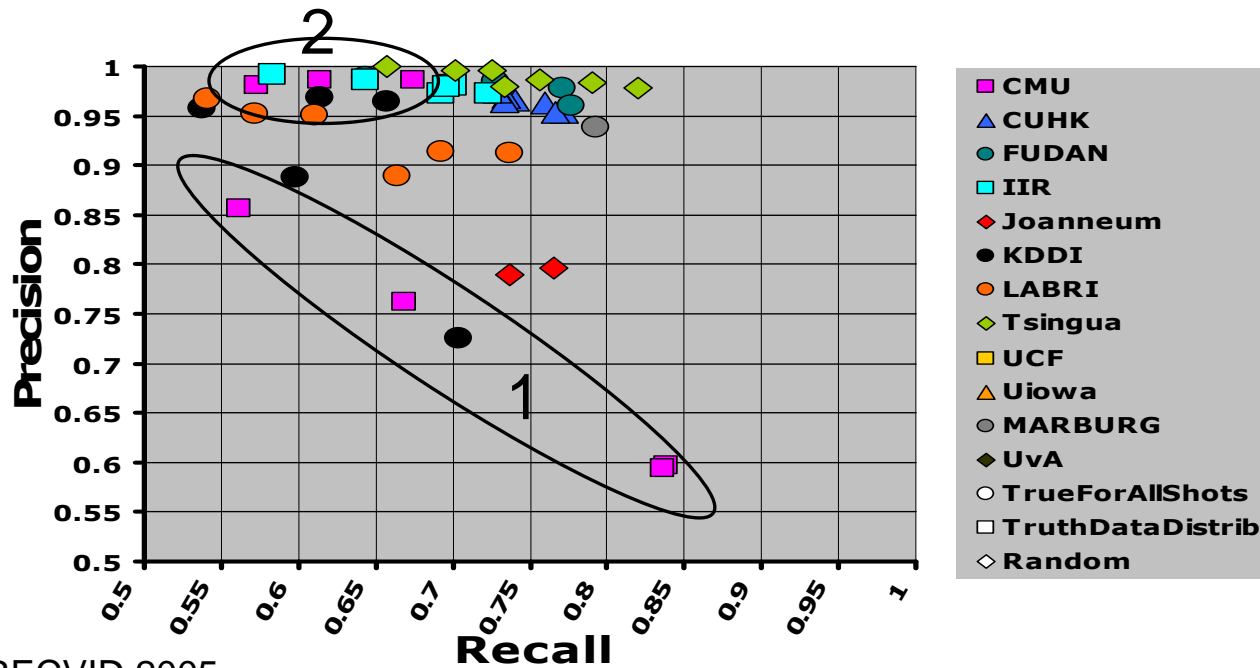- □ TruthDataDistrib
- ◇ Random

NIST

# General points

o NIST did not provide training data: some training data was available from other sources and some training data was produced by participants

o Input:

  n MPEG motion vectors: optimal for compression, not optimal for modeling real motion

  n Frame to frame motion analysis
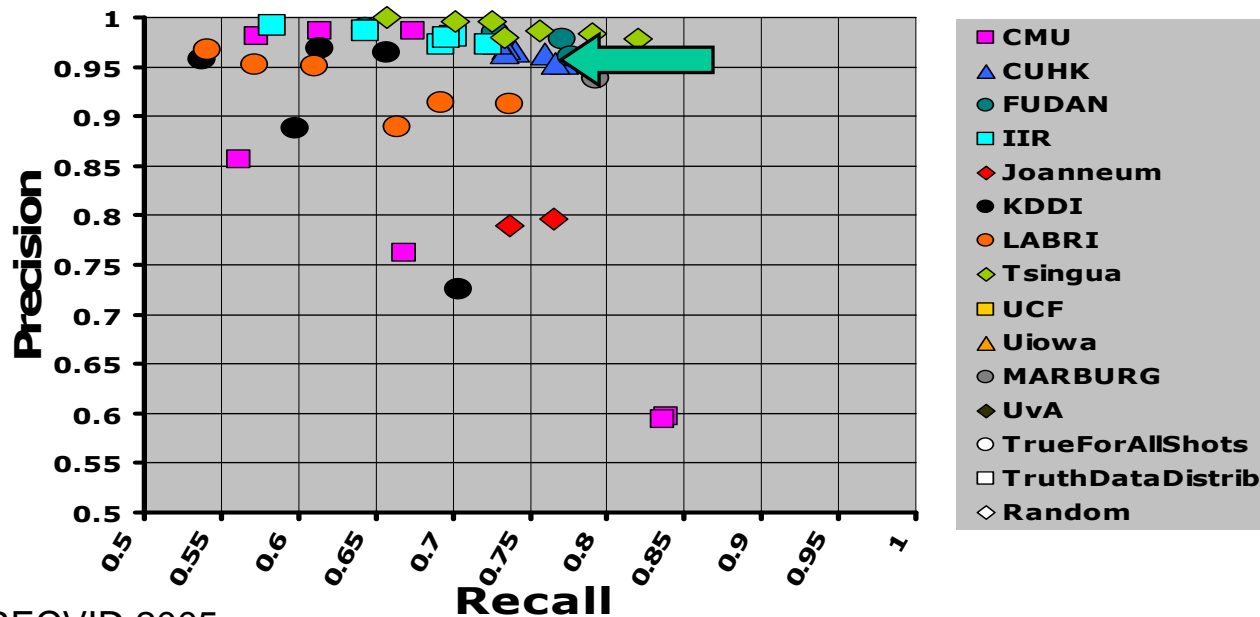
o Distinguish "jitter" from intended motion

# CMU

○ Approach

1. Probabilistic model (fitted using EM) based on MPEG motion vectors

2. Optical Flow model: extract the most consistent motion from the optical flows (frame to frame)
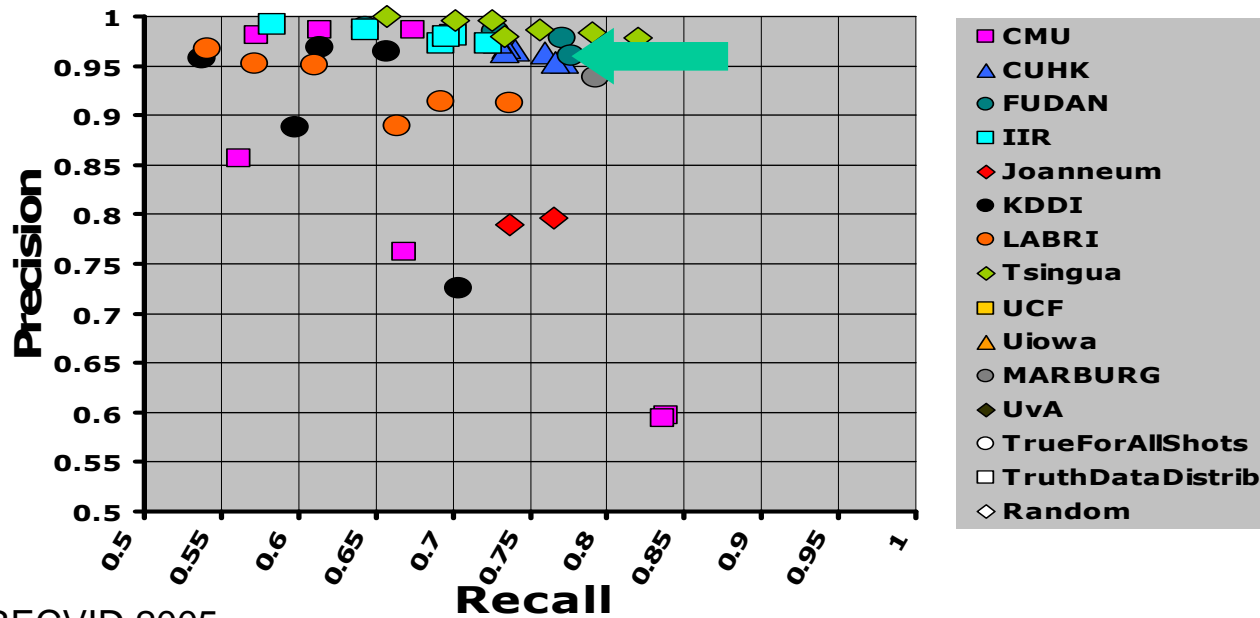
# CUHK

- o Approach
    - n Motion features extracted from tracking image features in consecutive frames
    - n Estimation of 6 parameter affine model, transformation in p,t,z vector for each set of adjacent frames
    - n Rule based motion classification using empirical thresholds
    - n Interesting failure analysis

# Fudan

- Approach
  - Motion vectors from MPEG,SVM, motion accumulation method to filter out imperceptable movements
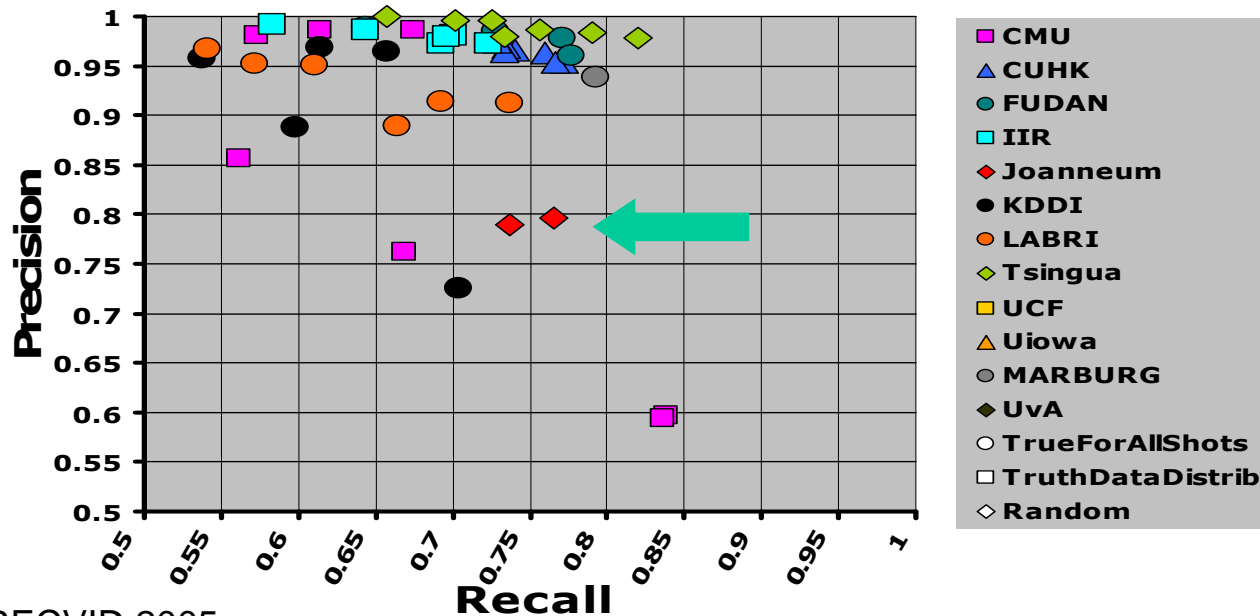  - Filter method seems to decrease precision though…

# Joanneum
## - presentation follows -

o  Approach
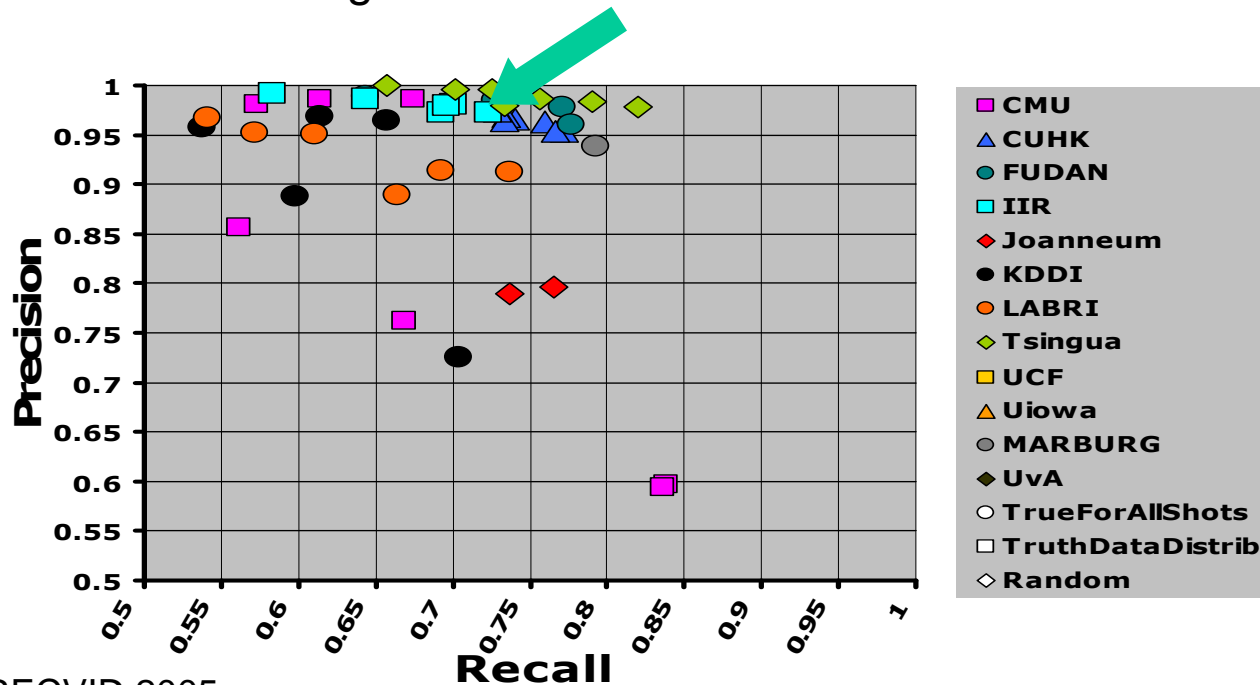
   n  Developed a training set , problems with annotation..

   n  Feature tracking, clustering trajectories, dominant cluster
       selection, camera motion detection, thresholding



Legend:
- ■ CMU
- ▲ CUHK
- ● FUDAN
- ■ IIR
- ◆ Joanneum
- ● KDDI
- ● LABRI
- ◆ Tsingua
- ■ UCF
- ▲ Uiowa
- ● MARBURG
- ◆ UvA
- ○ TrueForAllShots
- □ TruthDataDistrib
- ◇ Random

# IIR

o Approach

- n Annotated 24 video files
- n Estimated affine camera model based on MPEG motion vectors
- n Transformation of model parametersŁ series of p,t,z values for each shot
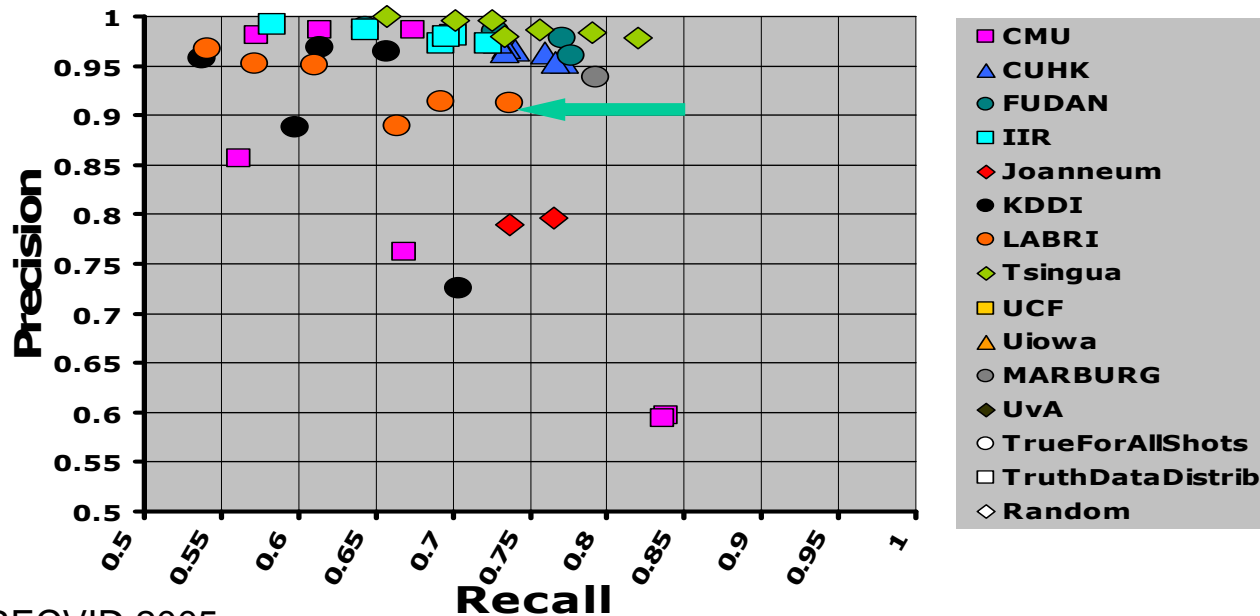- n Rule based classification of series using accumulation and thresholding



Legend:
- ■ CMU
- ▲ CUHK
- ● FUDAN
- ■ IIR
- ◆ Joanneum
- ● KDDI
- ● LABRI
- ◆ Tsingua
- □ UCF
- ▲ Uiowa
- ● MARBURG
- ◆ UvA
- ○ TrueForAllShots
- □ TruthDataDistrib
- ◇ Random

Axis labels: Precision (y-axis), Recall (x-axis)
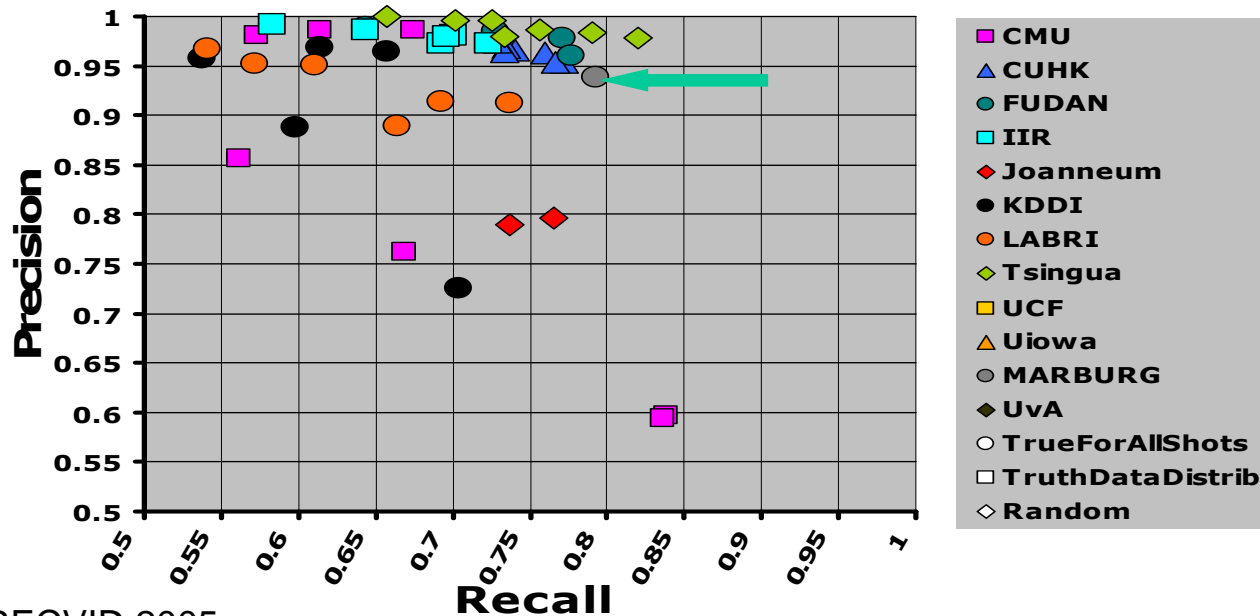
# LaBRI
## - presentation follows -

o Approach

   n Mpeg motion vector inputŁ  6 parameter affine model

   n Jitter suppression (statistical significance test)

   n Subshot segmentation (homogeneous motion)

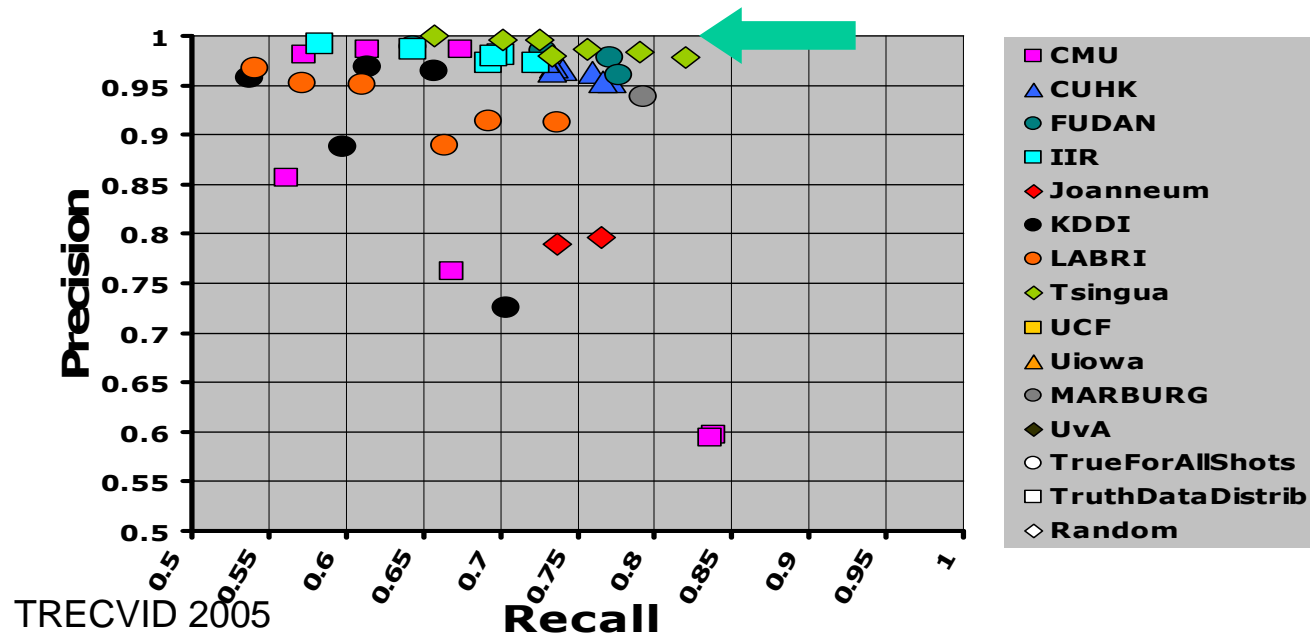   n Motion classification (using "a few annotated videos")

# Marburg

o Approach

  n 3D camera model estimated from MPEG motion vectors

  n Cleaning necessary, + exlusion of center, frame border

  n Optimal thresholds estimated on tv2005 training set

# Tsinghua

- ○ Approach
  - n Motion vector selection based spatial features, separating camera motion from object motion and accidental motion
  - n 4 parameter camera model (Iterative Least Squares) parameter estimation
  - n Rule based classification (FSA), using a range of thresholds for: 1.Continuous (speed) and noticable, 2,Minumum duration 3.Uninterrupted 4.Noticable in case in combination with other camera movement

# Observations

Ø This is clearly an easier task than the HLF task, though a high recall is hard to achieve.

Ø Truth data costly to create – lot's of shaky shots

  Ø Many hard to judge

  Ø Many not really what a user wants when s/he asks for a "pan" etc.

Ø Hard to generalize from small, constructed test subset to larger, more realistic test set

Ø Given the definition of our task and test set characteristics, F measure not appropriate

Ø Concentrate on within-feature system comparisons