



Pacman

Probabilistic and Compositional Representations for Object Manipulation



RoMANS

ROBOTIC MANIPULATION FOR NUCLEAR SORT AND SEGREGATION



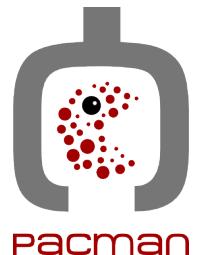
UoB Highly Occluded Object Challenge (UoB-HOOC)

1st International Workshop on Recovering 6D Object Pose
17 Dec 2015, In conjunction with ICCV 2015, Santiago, Chile

Challenge Organizers: Aleš Leonardis, Krzysztof Walas,
Hector Basevi, Jinpeng Wang



Outline



1. Dataset Description

- capturing devices and scene views
- scenes composition
- multi-camera system calibration
- annotations and scenes characterization

2. The Challenge

- training data
- expected results
- performance evaluation

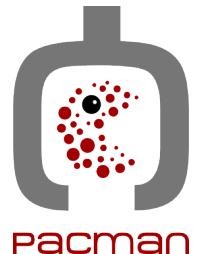
3. Results

- base line method
- submitted results

4. Summary and Future Work

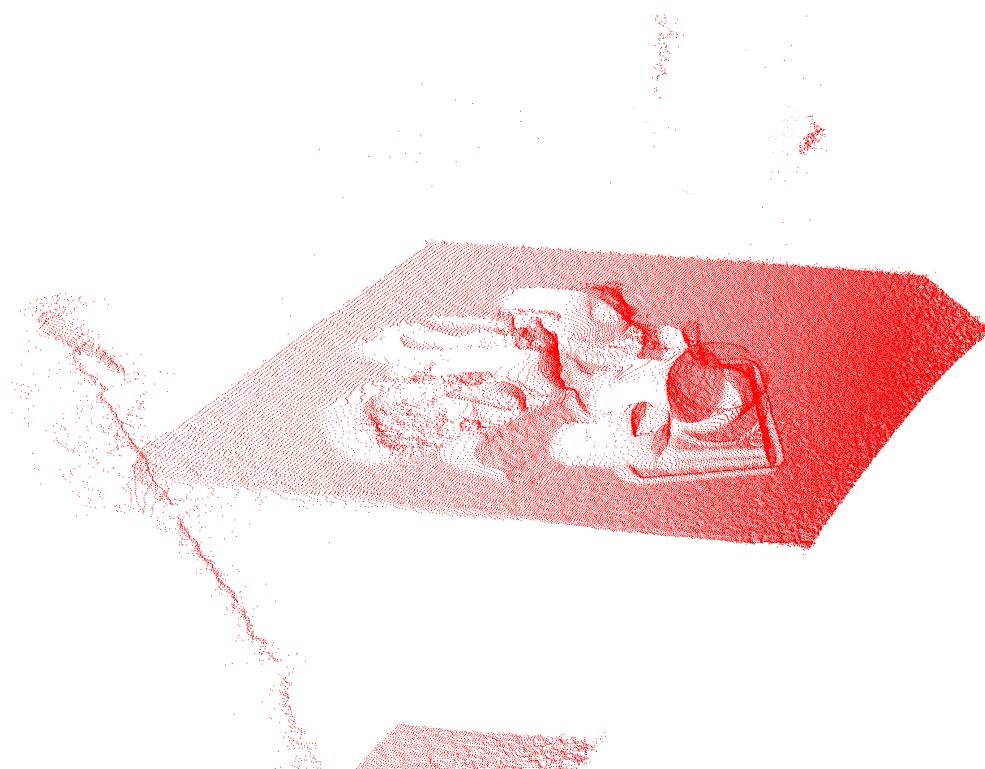


Dataset Description



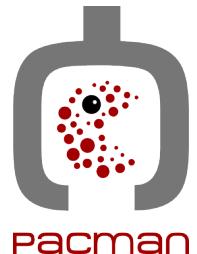
Scene views:

- data registered from three different Kinects
- all devices in fixed positions (mounted on tripods)



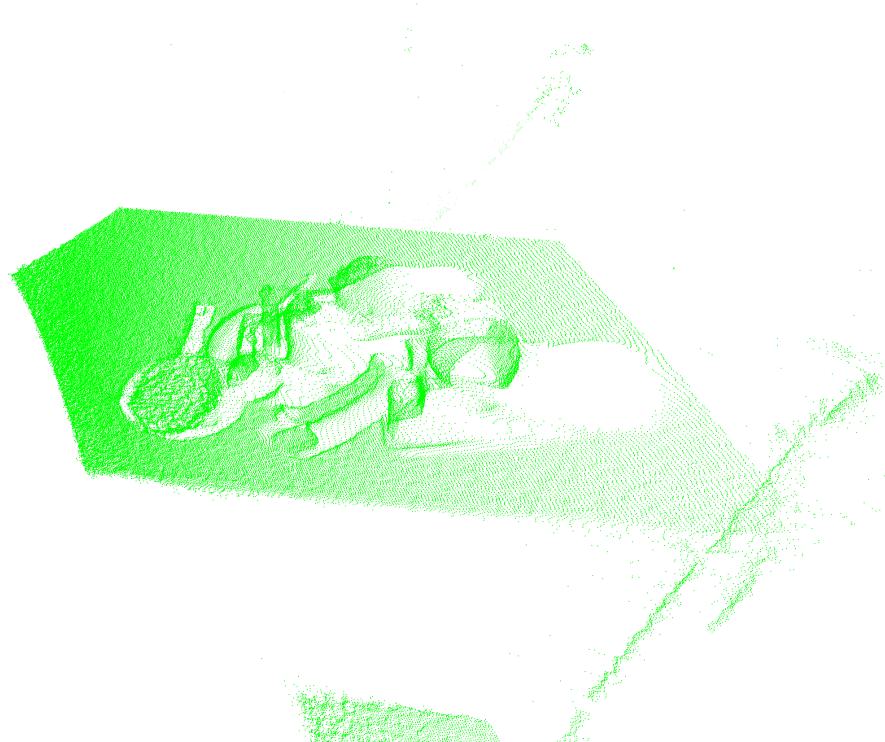


Dataset Description



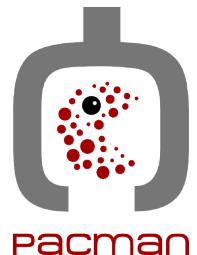
Scene views:

- data registered from three different Kinects
- all devices in fixed positions (mounted on tripods)



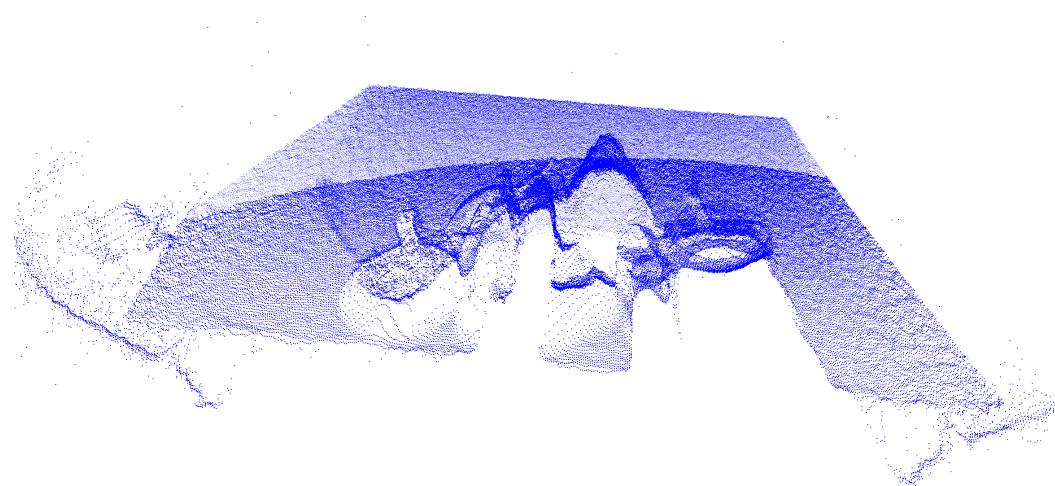


Dataset Description



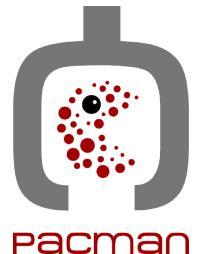
Scene views:

- data registered from three different Kinects
- all devices in fixed positions (mounted on tripods)



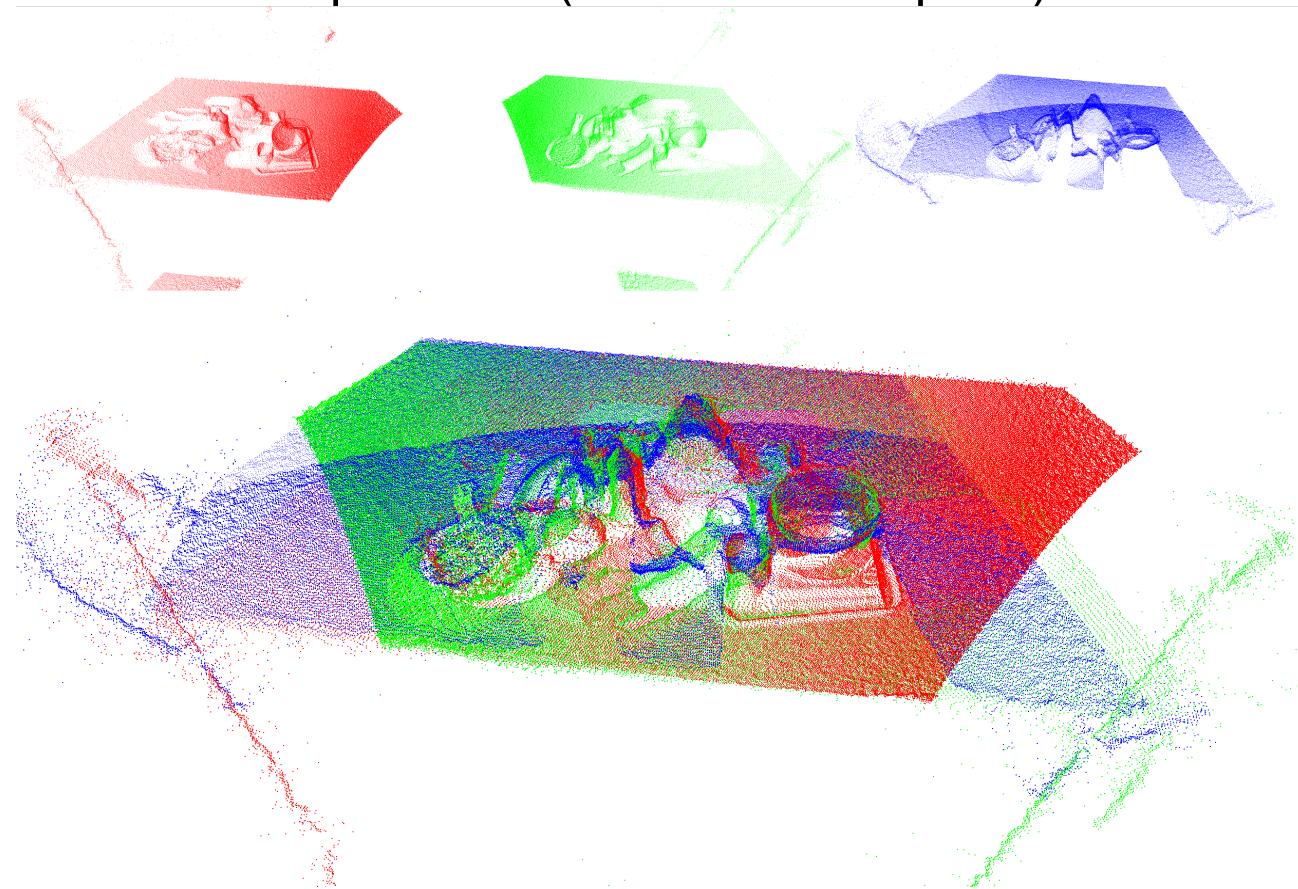


Dataset Description



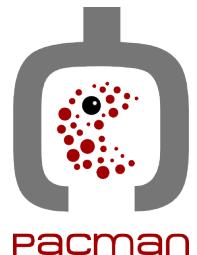
Scene views:

- data registered from three different Kinects
- all devices in fixed positions (mounted on tripods)





Dataset Description



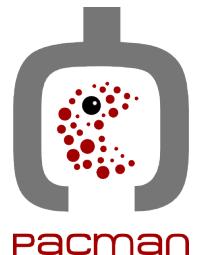
Scene compositions:

- up to 20 objects selected from 25 categories no repetitions
- high level of occlusions
- objects stacked one on top of another
- many symmetric objects





Dataset Description



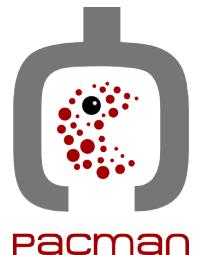
Calibration:

- infra-red image with long exposure (less noise)
- camera calibration toolbox for Matlab -- A2 chessboard
- intrinsic parameters, extrinsics computed for each scene separately
- all views transformed to the same reference frame associated with the chessboard lying on a table



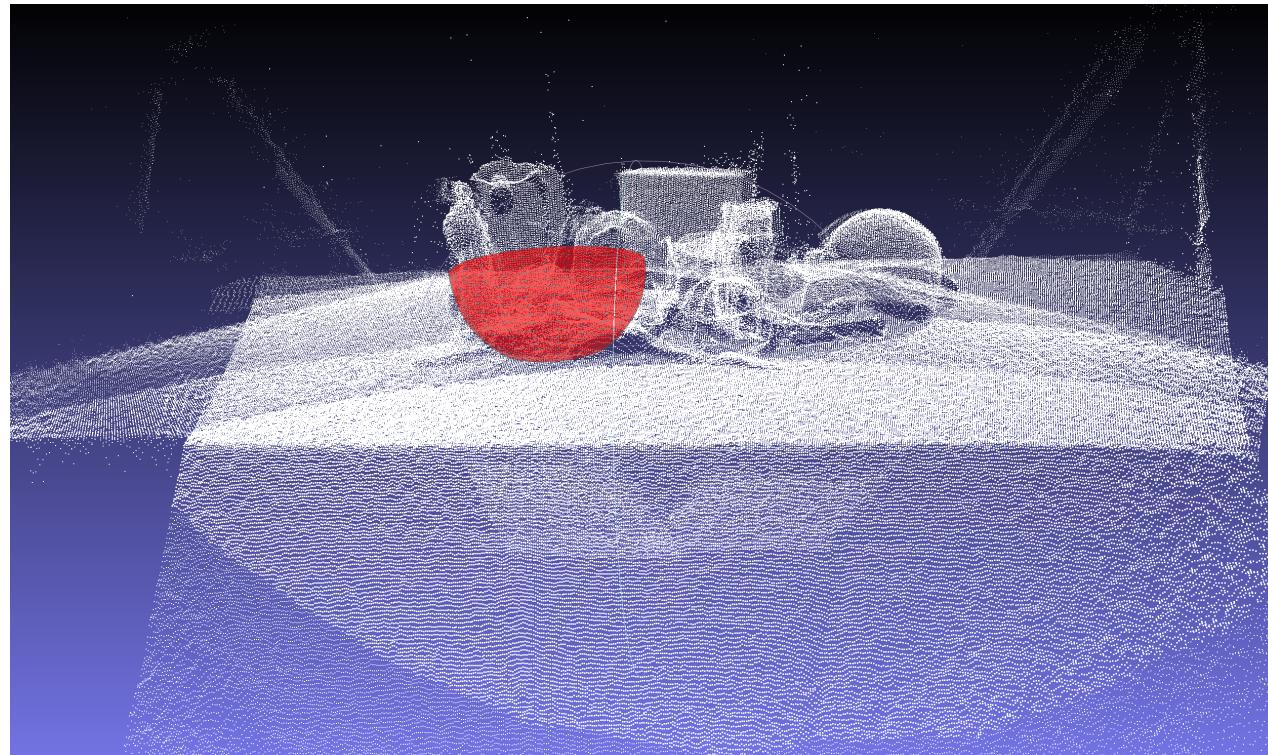


Dataset Description



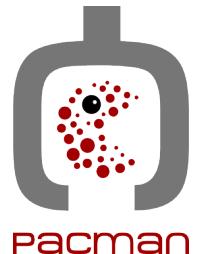
Annotations:

- using a set of 3D computer graphics models
- with 3 views of the scene loaded
- placing an object and establishing its scale



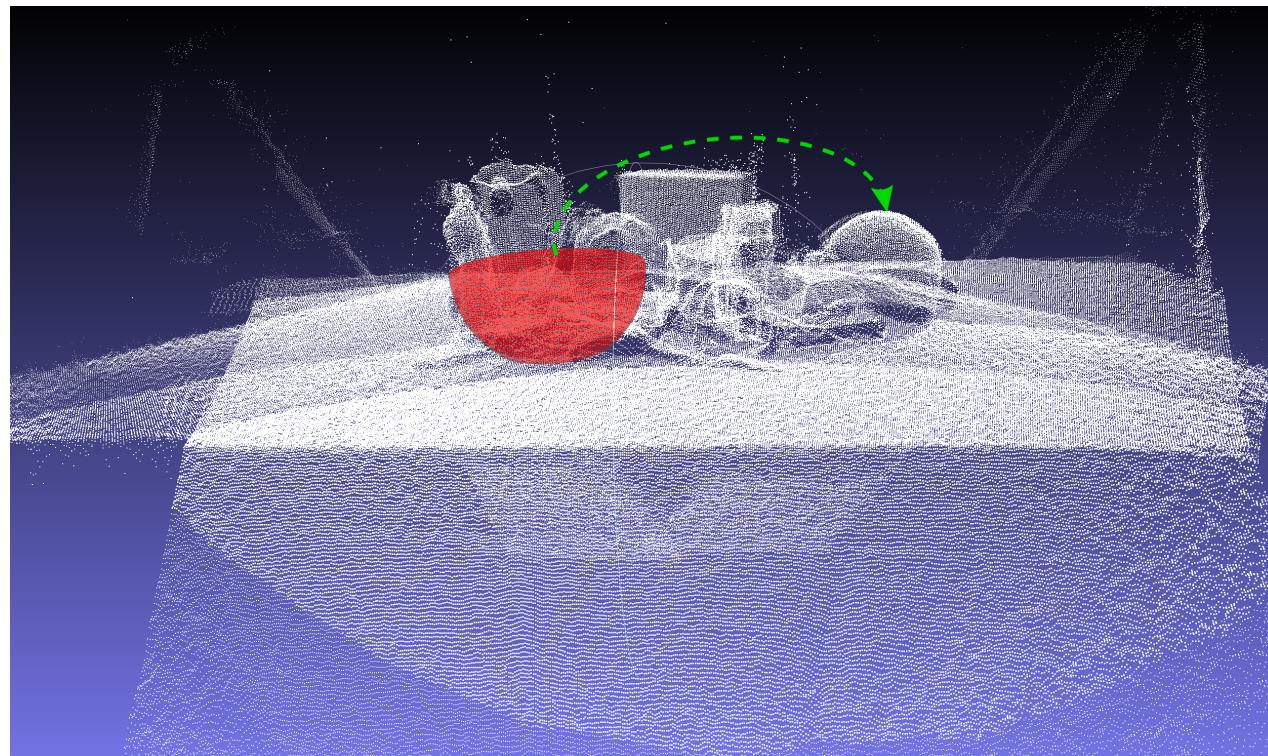


Dataset Description



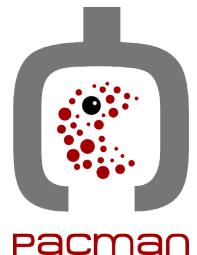
Annotations:

- using a set of 3D computer graphics models
- with 3 views of the scene loaded
- placing an object and establishing its scale



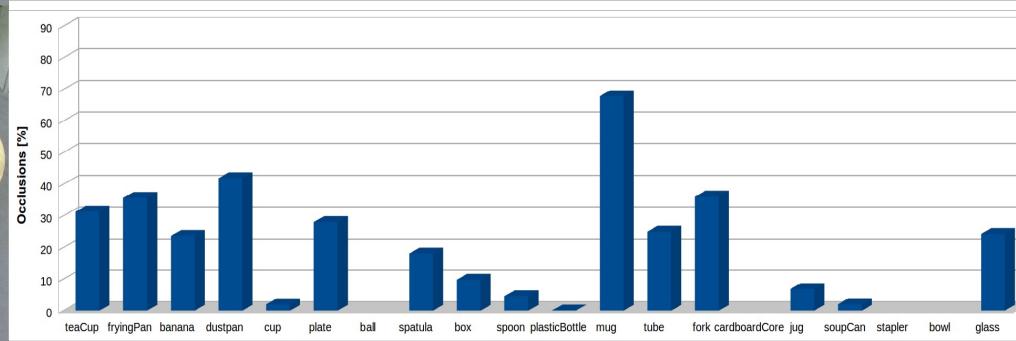
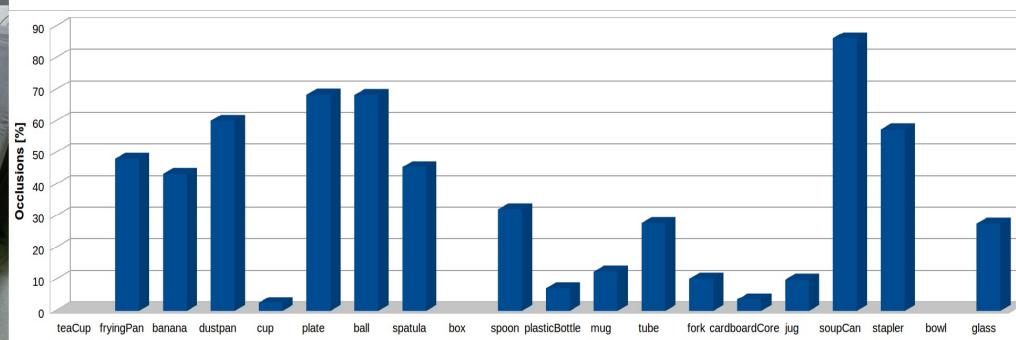
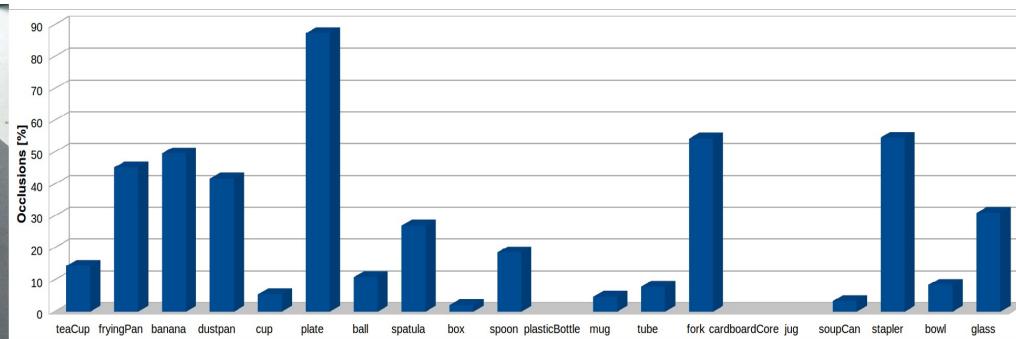


Dataset Description



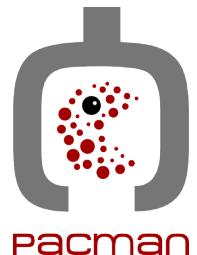
Characterisations – object occlusion caused by other objects (points):

Camera 1 Camera 2 Camera 3

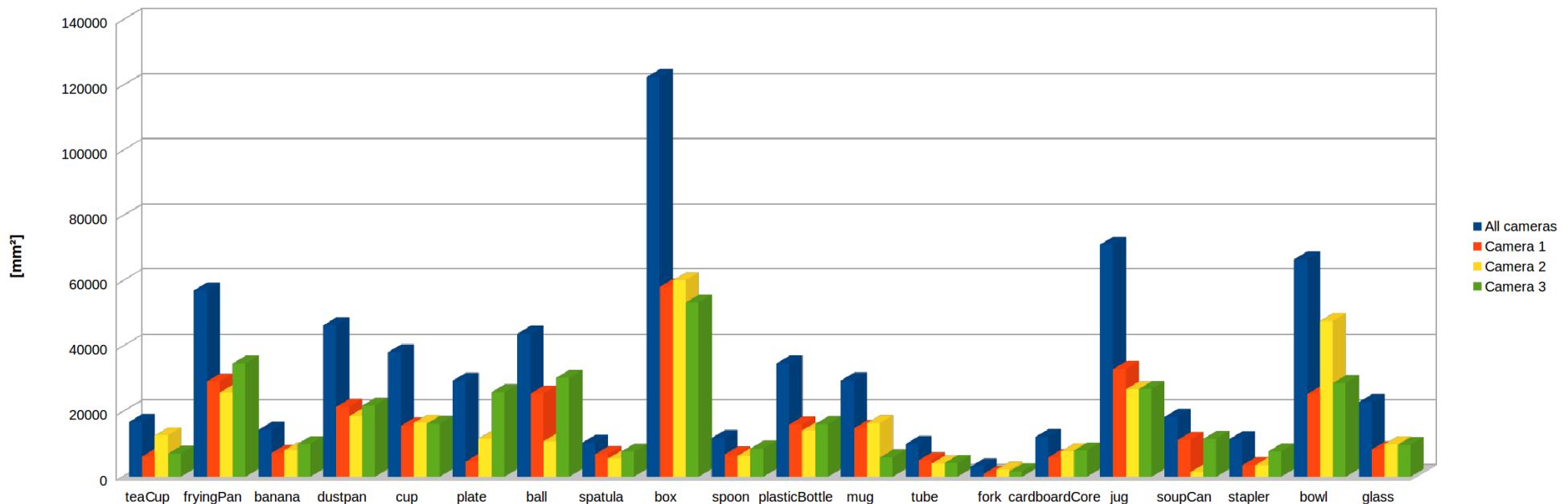


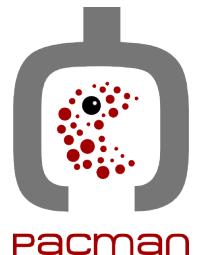


Dataset Description



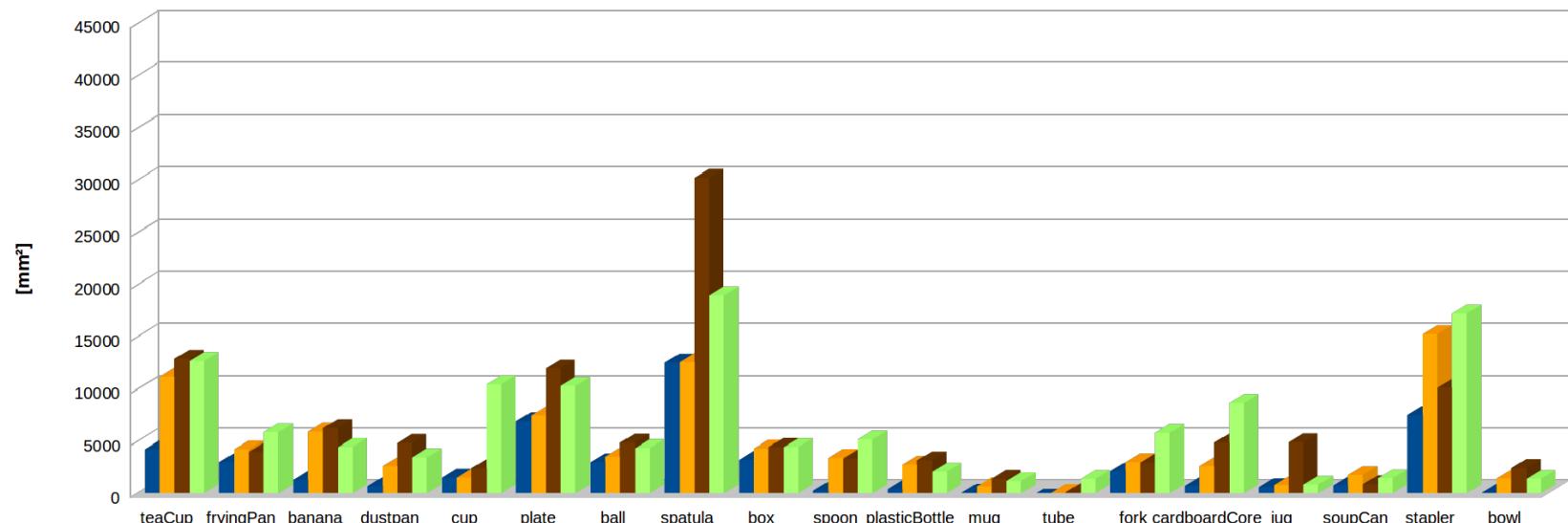
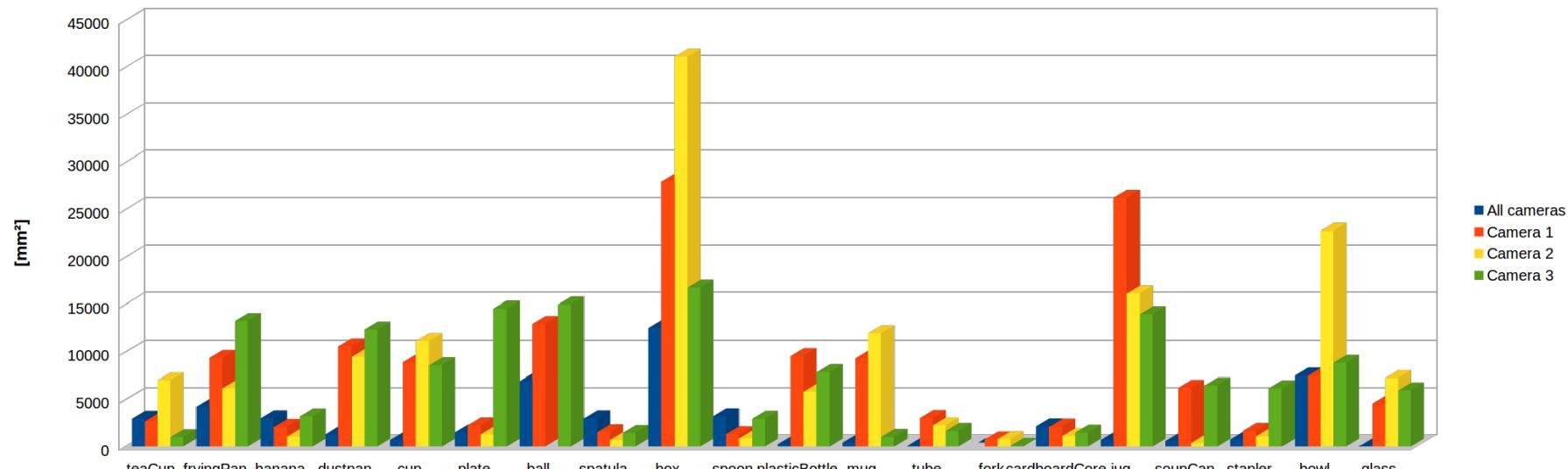
Characterisations – object visibilities (area):





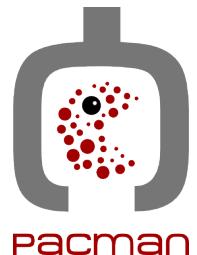
Dataset Description

Characterisations – object visibilities overlap (area):

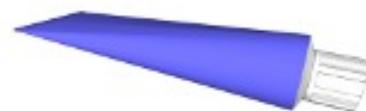
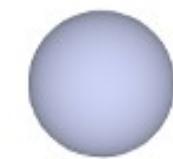




Training Data

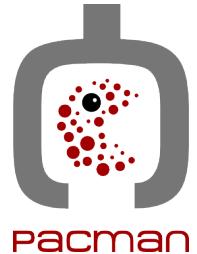


Scene may contain up to 20 objects from 25 different categories:



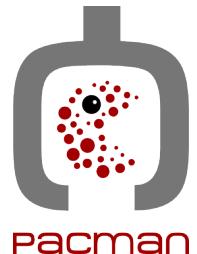


Dataset Description



Content of the repository:

- pointclouds in camera reference frame, ordered pointclouds
- pointclouds in a common reference frame placed on a table top
- RGB images of the scenes (calibration data relating RGB and D image is not provided, RGB image should be used for illustration purposes only)
- depth data together with intrinsic and extrinsic parameters of the camera
- set of generic objects which has to be fit into test sceneData registered from three different Kinects
- scripts for downloading training data

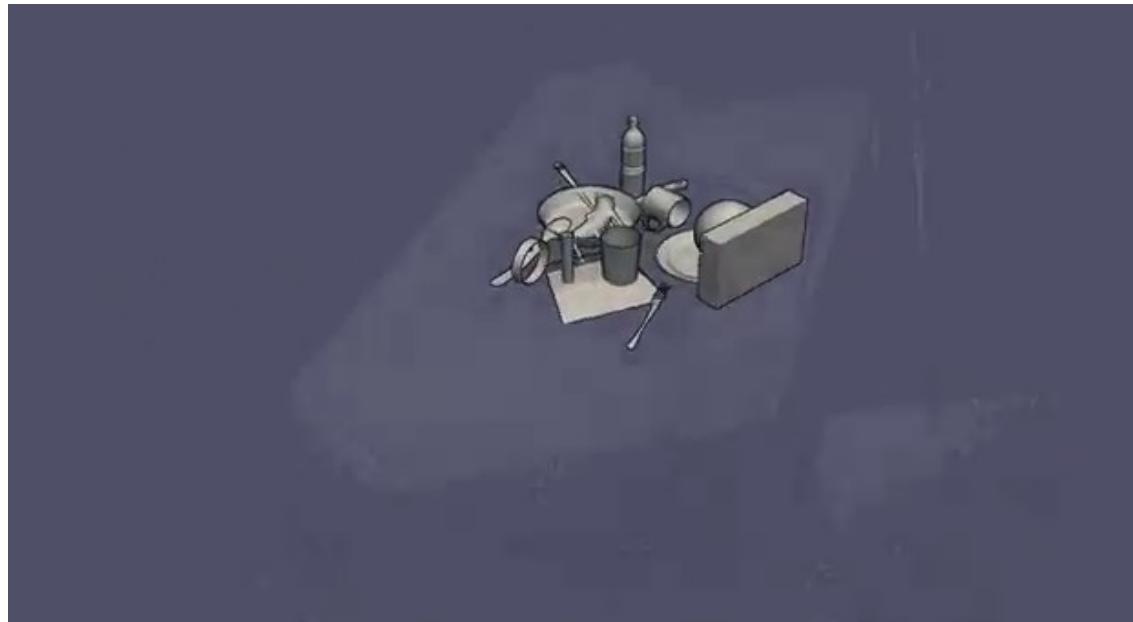


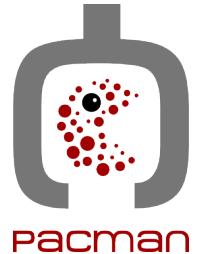
Expected result

Given a set of computer graphics models of generic representatives:

- establish category of each object in a registered pointcloud
- find its 6D pose and appropriate scale

In other words, the generic exemplar has to be fitted as accurately as possible into the object in the registered scene .





Performance measures

Common statistical measures for multi-label classification¹:

$$Accuracy(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{|Y_i \cap Z_i|}{|Y_i \cup Z_i|} \quad (1)$$

$$Precision(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{|Y_i \cap Z_i|}{|Z_i|} \quad (2)$$

$$Recall(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{|Y_i \cap Z_i|}{|Y_i|} \quad (3)$$

$$F_1(H, D) = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (4)$$

where:

H – multi-label classifier

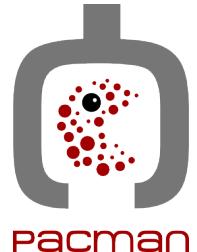
D – multi-label dataset consisting of $|D|$ multi-label examples $(x_i, Y_i); i = 1 \dots |D|$

$Z_i = H(x_i)$ – a set of labels predicted by H for example x_i ;

¹ Grigoris Tsoumakas and Ioannis Katakis. Multi-label classification: An overview. Int J Data Warehousing and Mining, 2007:1–13, 2007.



Performance measures



Dataset characterization:

$$LD(D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{|Y_i|}{|L|} \quad (5)$$

where:

L – a set of disjoint labels;

$|L|$ – cardinality of the set of labels;

$Y \subseteq L$ – set of labels for current example (scene);

D – multi-label dataset consisting of $|D|$ multi-label examples $(x_i, Y_i); i = 1 \dots |D|$

Hamming Loss:

$$HammingLoss(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{|Y_i \Delta Z_i|}{|L|} \quad (6)$$

where:

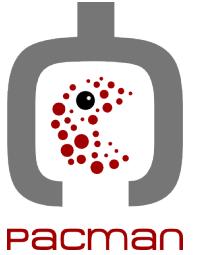
Δ – symmetric difference of two sets;

$Z_i = H(x_i)$ – a set of labels predicted by H for example x_i ;

H – multi-label classifier.



Performance measures



Detection measure:

$$DetectionError(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{1}{3} \cdot \left(\frac{trans_{dist}}{0.2} + \frac{rot_{dist}}{\pi} + scale_{dist} \right) \cdot \left(1 - \frac{|Y_i|}{|L|} \right) \quad (7)$$

$$T = \begin{bmatrix} s_x r_{11} & s_y r_{12} & s_z r_{13} & t_x \\ s_x r_{21} & s_y r_{22} & s_z r_{23} & t_y \\ s_x r_{31} & s_y r_{32} & s_z r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (8)$$

$$T = Trans \cdot Rot \cdot Scale \quad (9)$$

$$trans_{dist} = \sqrt{(t_{x1} - t_{x2})^2 + (t_{y1} - t_{y2})^2 + (t_{z1} - t_{z2})^2}. \quad (10)$$

$$rot_{dist} = \Phi(Q_1, Q_2) = \arccos(\|Q_1 \cdot Q_2\|), \quad (11)$$

where:

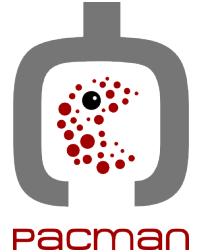
$$Q_1 \cdot Q_2 = w_1 w_2 + x_1 x_2 + y_1 y_2 + z_1 z_2$$

$$Q = \begin{bmatrix} \frac{1}{2}\sqrt{1 + r_{11} + r_{22} + r_{33}} = w \\ \frac{r_{32} - r_{23}}{4w} = x \\ \frac{r_{13} - r_{31}}{4w} = y \\ \frac{r_{21} - r_{12}}{4w} = z \end{bmatrix} \quad (12)$$

$$scale_{dist} = \left| \frac{s_1^3 - s_2^3}{s_1^3} \right| \quad (13)$$



Performance measures



Average Distance (AD) measure²:

$$m = \operatorname{avg}_{\mathbf{x} \in \mathcal{M}} \|(\mathbf{R}\mathbf{x} + \mathbf{T}) - (\tilde{\mathbf{R}}\mathbf{x} + \tilde{\mathbf{T}})\| \quad (14)$$

where:

m – a matching score,

\mathcal{M} – a 3D model,

\mathbf{x} – vertex of a 3D model \mathcal{M} ,

\mathbf{R} – a ground truth rotation matrix,

\mathbf{T} – a ground truth translation vector,

$\tilde{\mathbf{R}}$ – an estimated rotation matrix,

$\tilde{\mathbf{T}}$ – an estimated translation vector,

$$m = \operatorname{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \|(\mathbf{R}\mathbf{x}_1 + \mathbf{T}) - (\tilde{\mathbf{R}}\mathbf{x}_2 + \tilde{\mathbf{T}})\| \quad (15)$$

where:

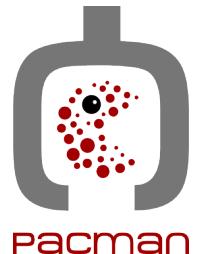
\mathbf{x}_1 – vertex of a 3D model \mathcal{M} in a ground truth pose,

\mathbf{x}_2 – vertex of a 3D model \mathcal{M} in an estimated pose,

² Stefan Hinterstoisser, Vincent Lepetit, Slobodan Ilic, Stefan Holzer, Gary R. Bradski, Kurt Konolige, Nassir Navab, Model Based Training, Detection and Pose Estimation of Texture-Less 3D Objects in Heavily Cluttered Scenes. ACCV 2012



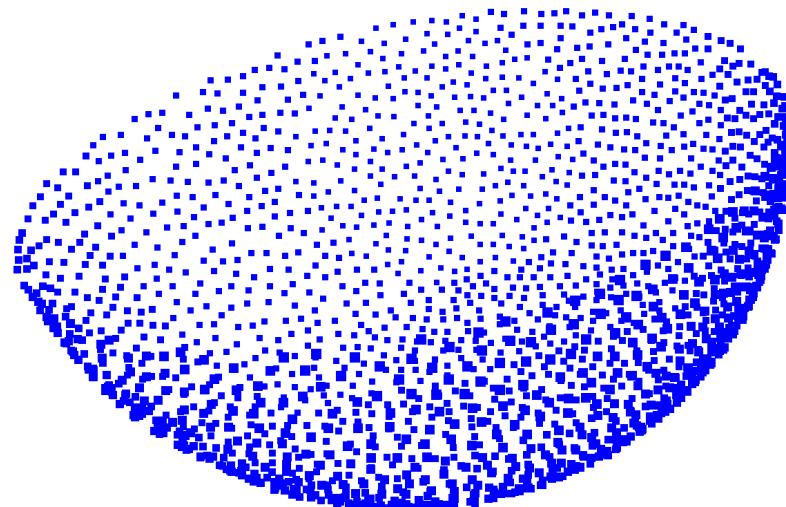
Performance measures



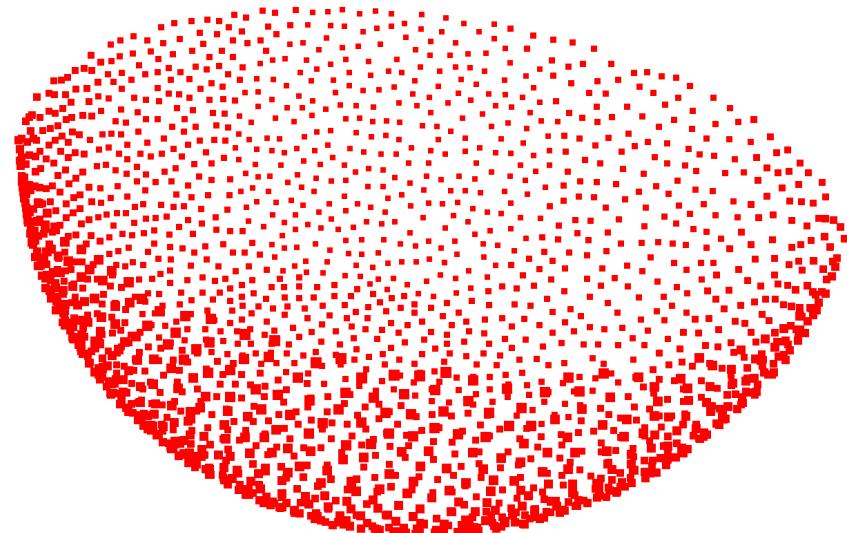
Improved Average Distance (IAD) measure:

$$m = \text{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \|(\mathbf{R}\mathbf{x}_1) - (\tilde{\mathbf{R}}\mathbf{x}_2)\| + \|\mathbf{T} - \tilde{\mathbf{T}}\| \quad (16)$$

Actual object pose

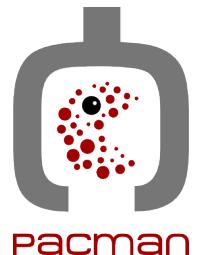


Predicted object pose





Performance measures

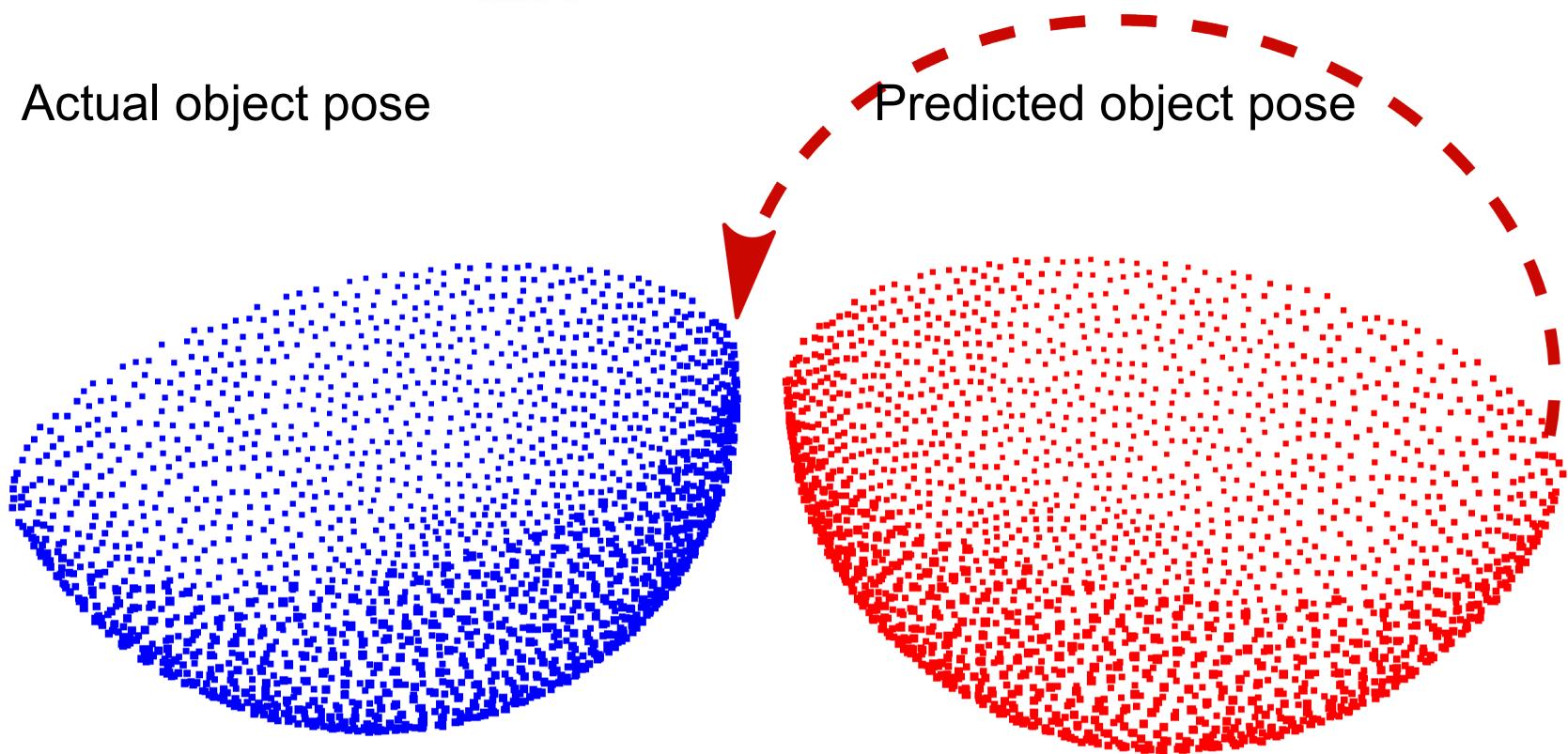


Improved Average Distance (IAD) measure:

$$m = \text{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \|(\mathbf{R}\mathbf{x}_1) - (\tilde{\mathbf{R}}\mathbf{x}_2)\| + \|\mathbf{T} - \tilde{\mathbf{T}}\| \quad (16)$$

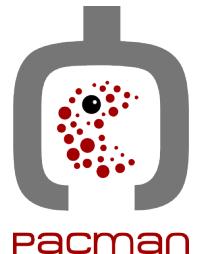
Actual object pose

Predicted object pose





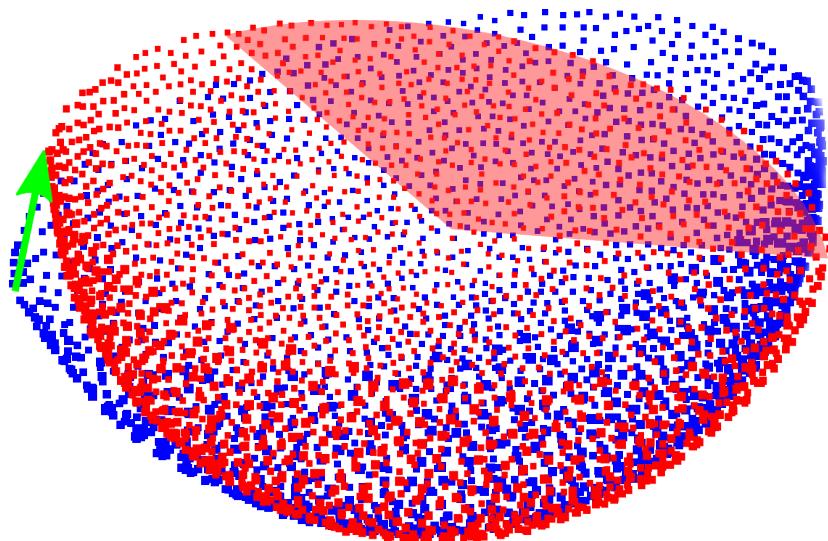
Performance measures

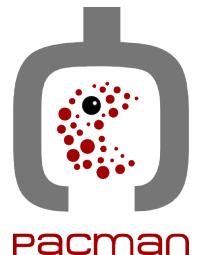


Improved Average Distance (IAD) measure:

$$m = \text{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \|(\mathbf{R}\mathbf{x}_1) - (\tilde{\mathbf{R}}\mathbf{x}_2)\| + \|\mathbf{T} - \tilde{\mathbf{T}}\| \quad (16)$$

Actual object rotation/Predicted object rotation



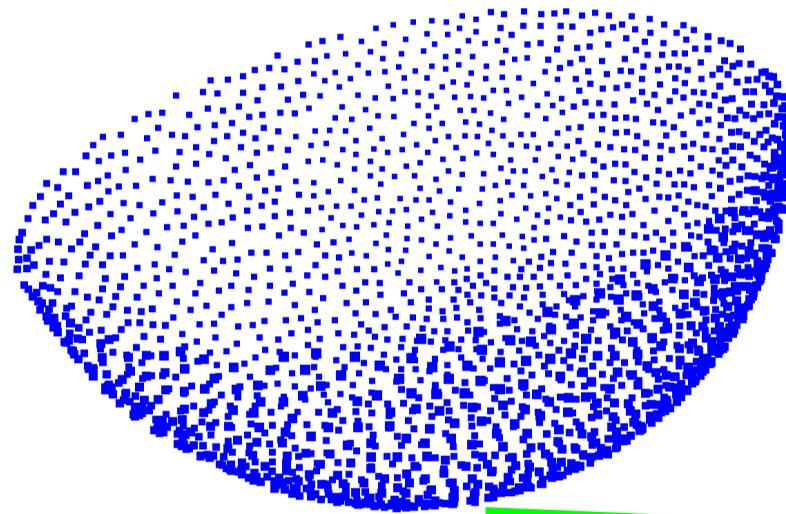


Performance measures

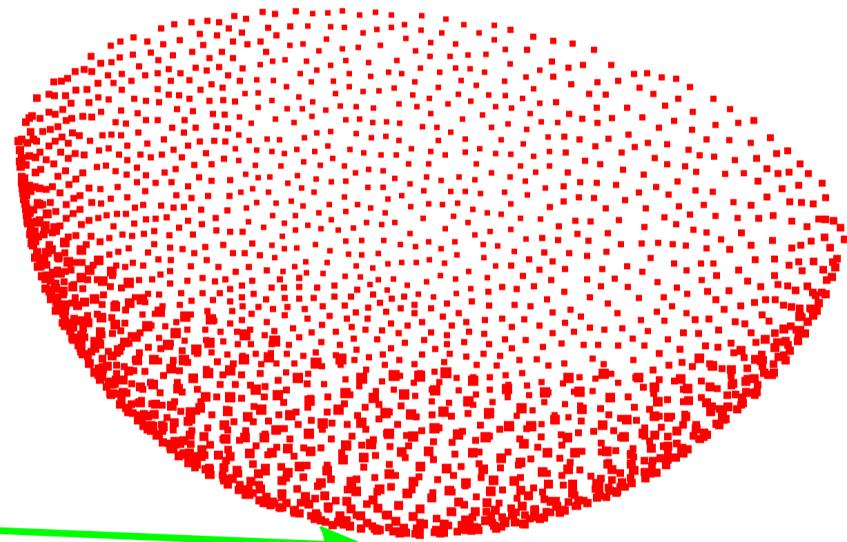
Improved Average Distance (IAD) measure:

$$m = \text{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \|(\mathbf{R}\mathbf{x}_1) - (\tilde{\mathbf{R}}\mathbf{x}_2)\| + \|\mathbf{T} - \tilde{\mathbf{T}}\| \quad (16)$$

Actual object translation

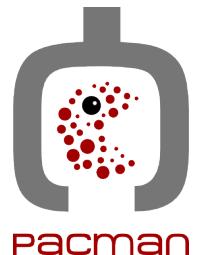


Predicted object translation

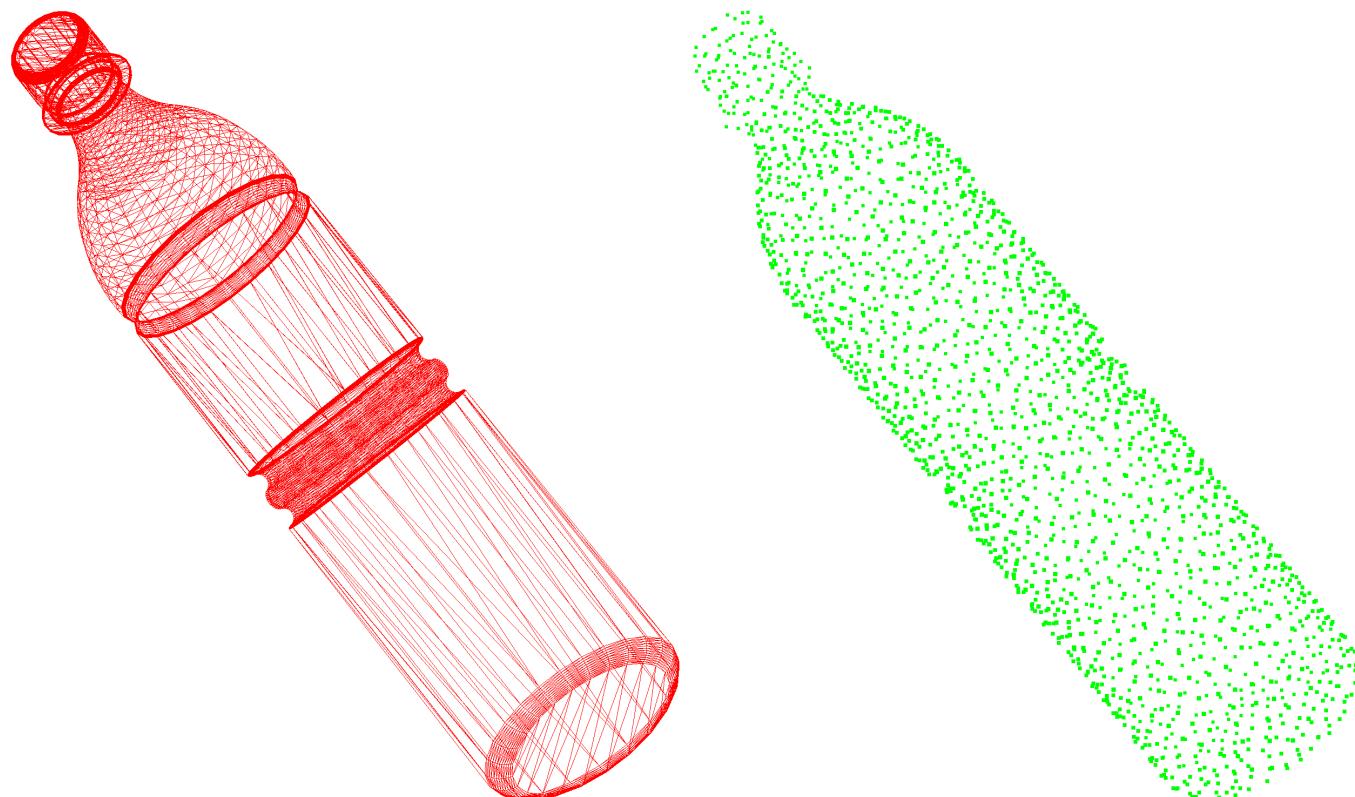




Performance measures



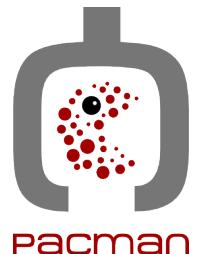
Non-uniform discretisation of object for meshes – resampling³:



³M.; Cignoni, P.; Scopigno, R., "Efficient and Flexible Sampling with Blue Noise Properties of Triangular Meshes," in Visualization and Computer Graphics, IEEE Transactions on , vol.18, no.6, pp.914-924, June 2012



Results

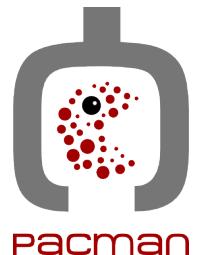


Baseline method based on Brachmann et al., Learning 6D Object Pose Estimation using 3D Object Coordinates, ECCV 2014 :

- decision forests for calculating class probabilities and estimated object coordinates for each pixel in a RGB-D image
- decisions via binary linear classifiers on RGB or D local gradients
- forests trained to separate pixels by object class and discretised object coordinates
- estimated object coordinates used to perform rough pose estimation
- RANSAC-based algorithm and energy function for pose refinement and final pose selection

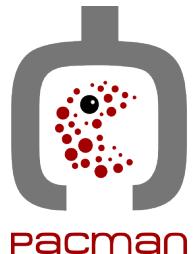


Results



Baseline method – modifications to Brachmann et al. method:

- Data pre-processing:
 - UoB-HOOC data transformed using calibration data to correct for lens spherical aberration
- Training:
 - on synthetic images of generic objects in randomised poses
 - no RGB data, decisions on depth features alone
- Decision forests:
 - different forests trained for each camera
 - data integrated to perform pose estimation using all views
- Pose estimation:
 - rough pose estimation and final pose selection process modified for compatibility with symmetric objects



Results

Baseline – workflow on selected scene:

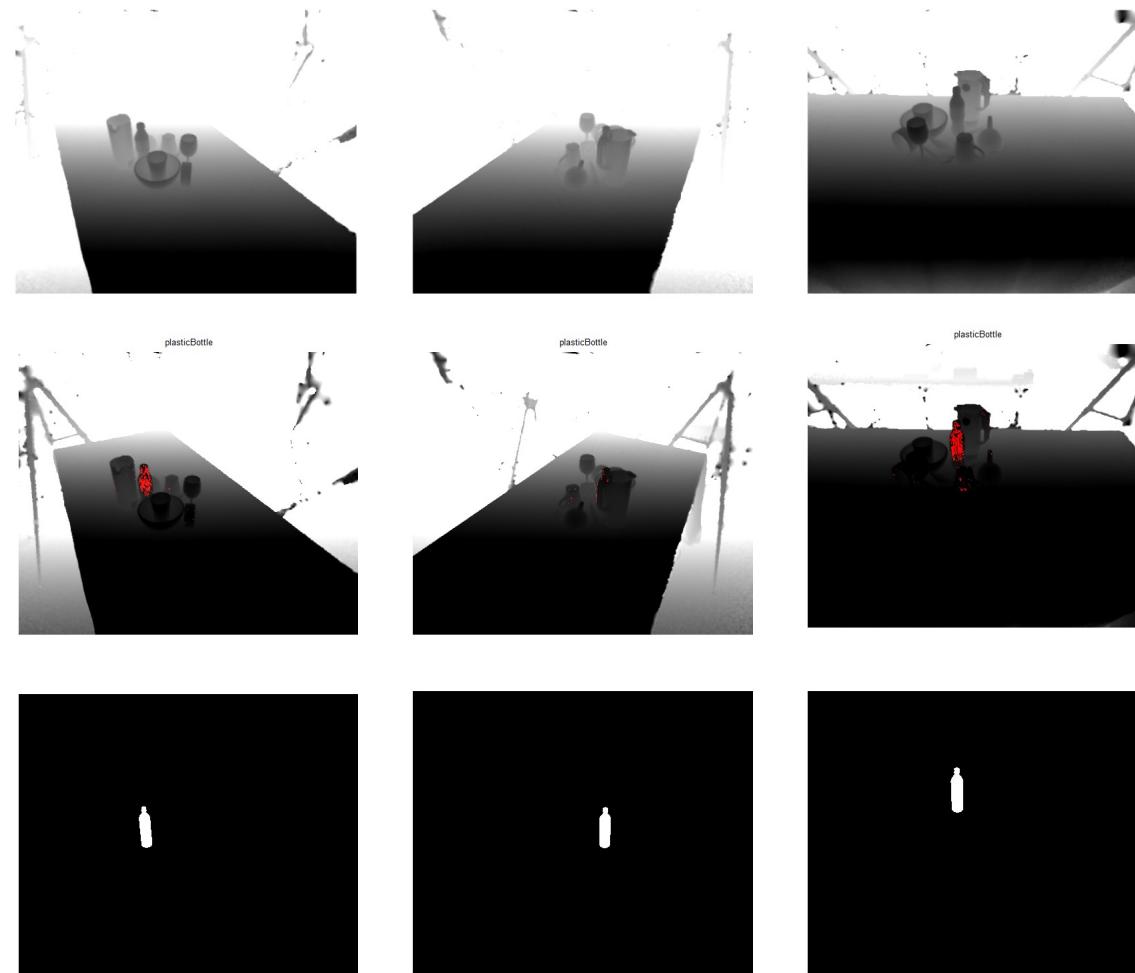
corrected
depth
maps



probability
maps from
decision
forests

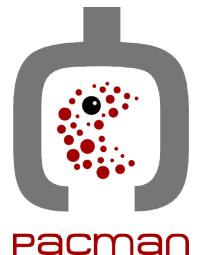


estimated
depth from
probable
depth data

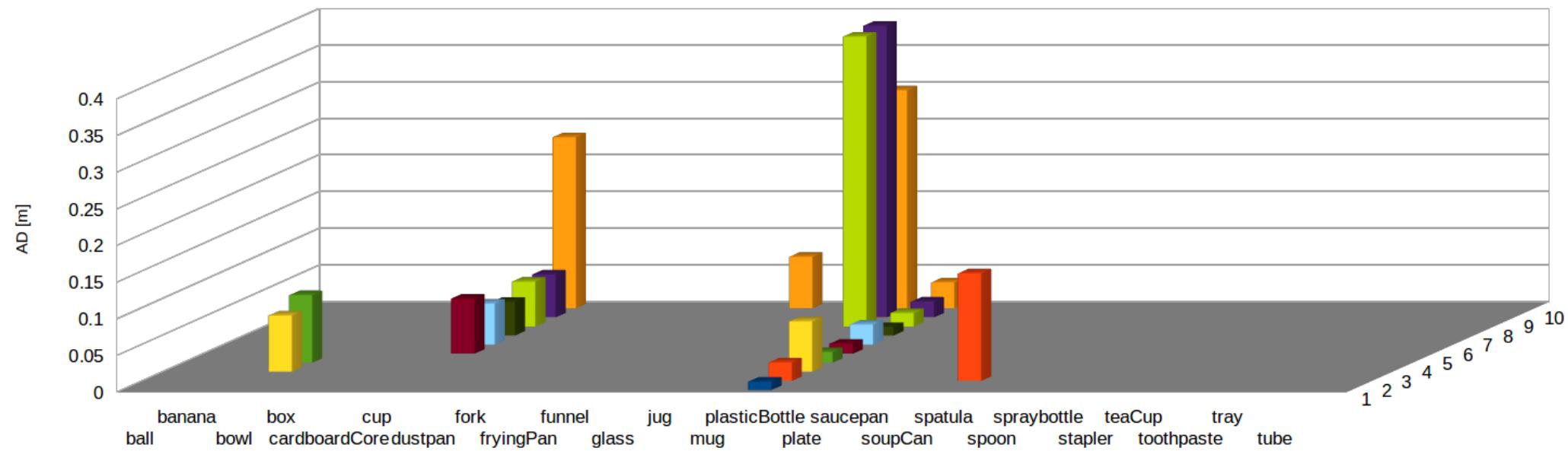




Preliminary Results

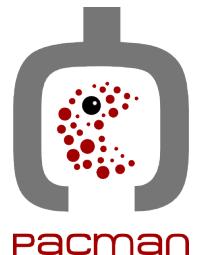


Up to 10 objects – Birmingham:

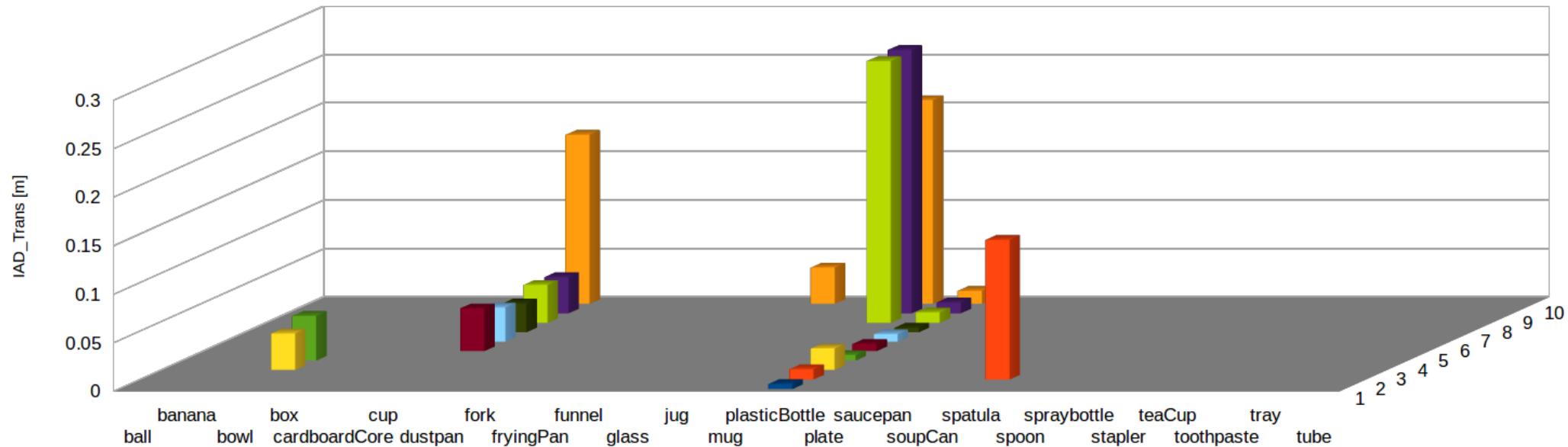




Preliminary Results

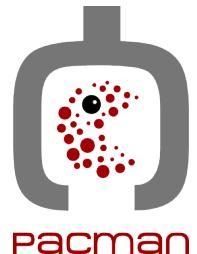


Up to 10 objects – Birmingham:

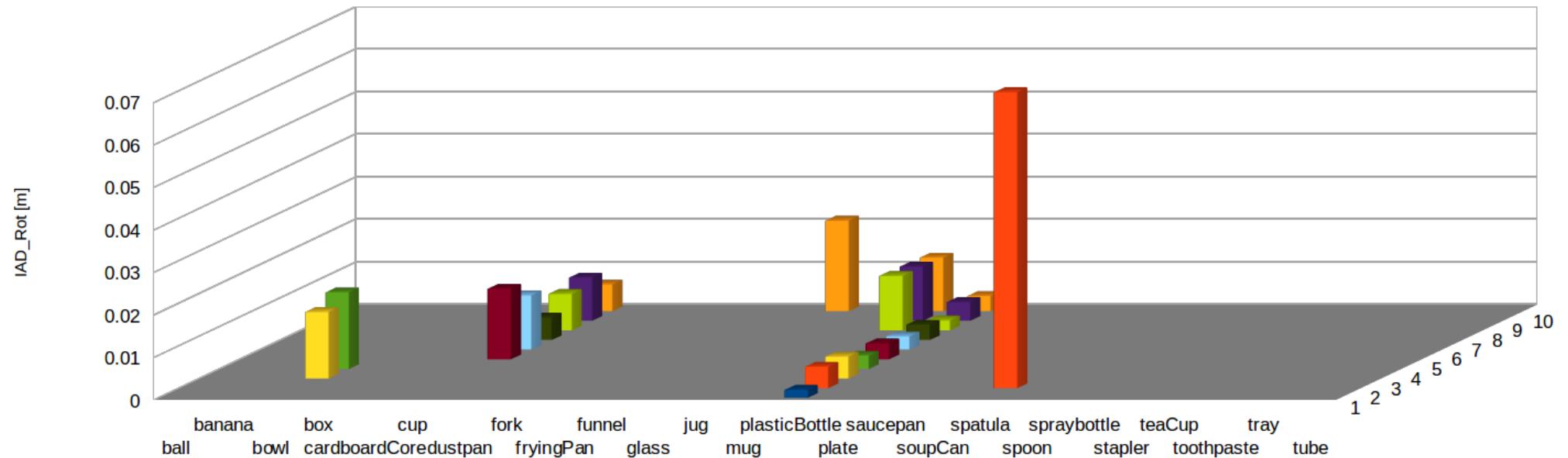




Preliminary Results

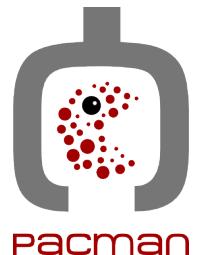


Up to 10 objects – Birmingham:

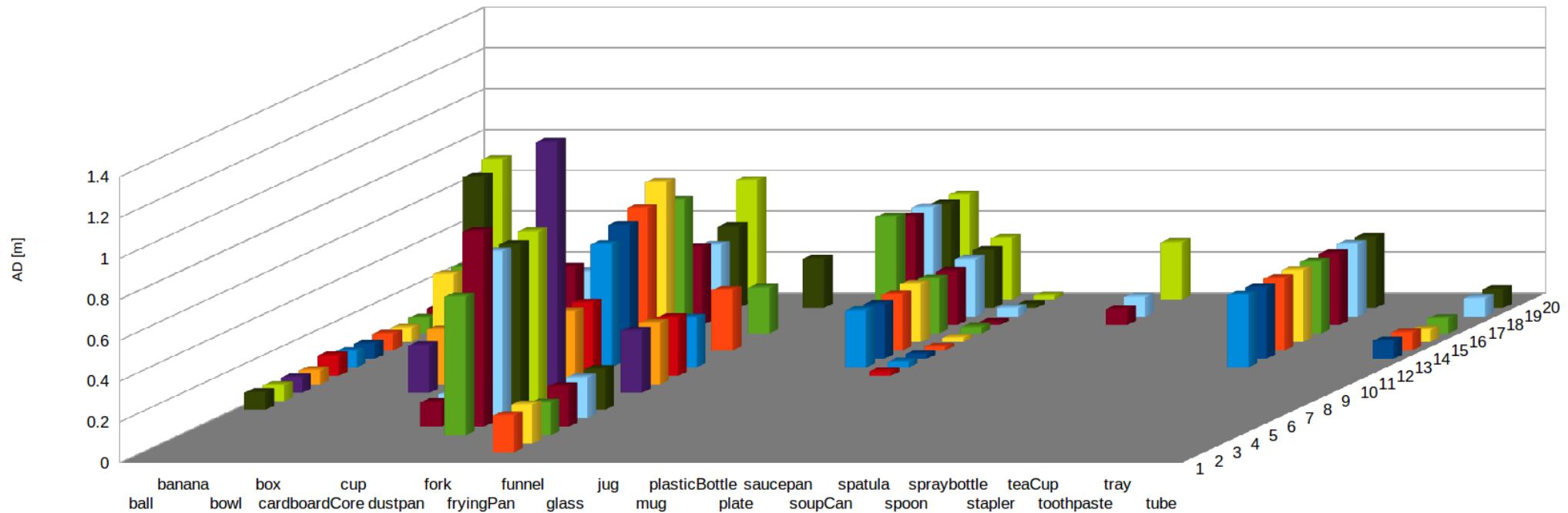




Preliminary Results

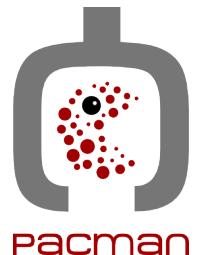


Up to 20 objects – Birmingham:

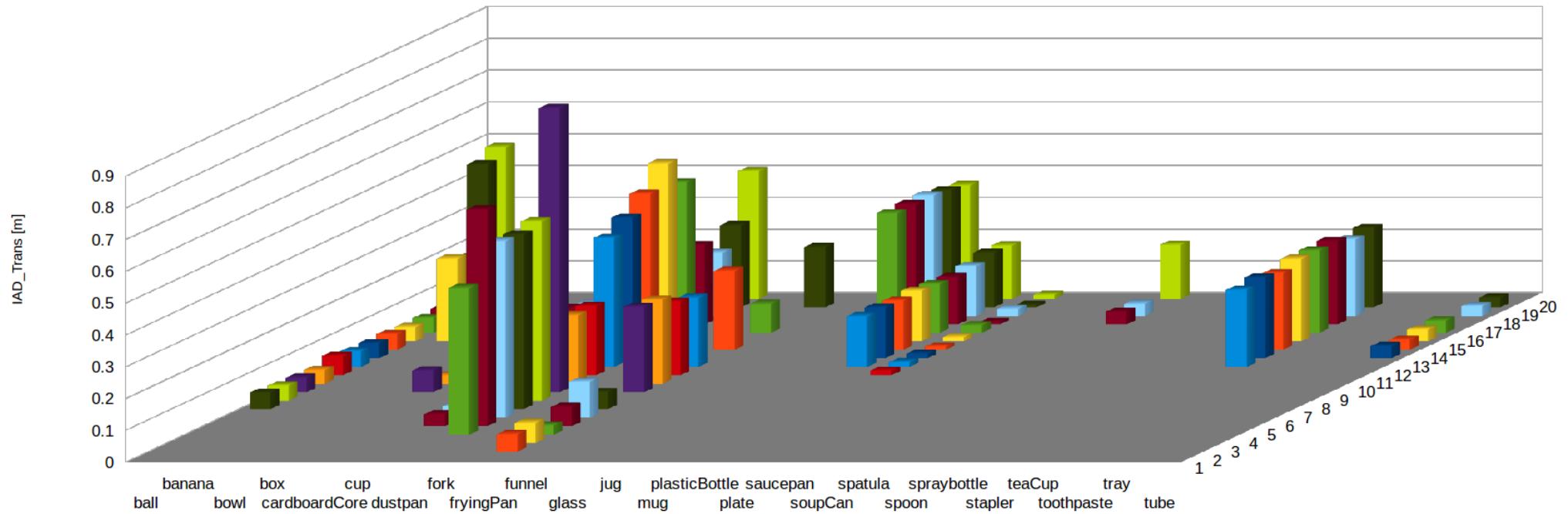




Preliminary Results

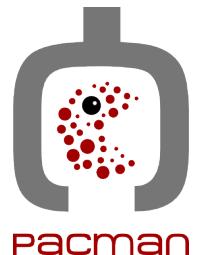


Up to 20 objects – Birmingham:

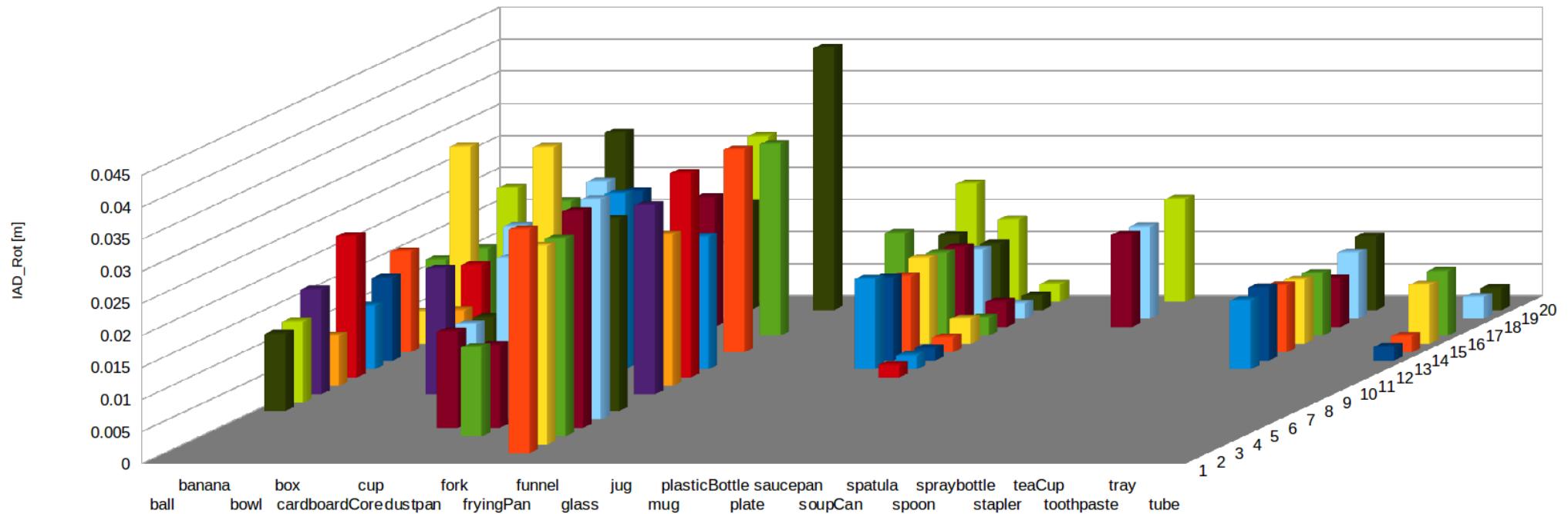




Preliminary Results

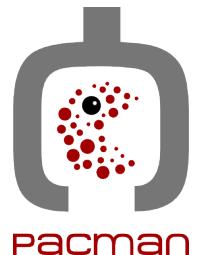


Up to 20 objects – Birmingham:

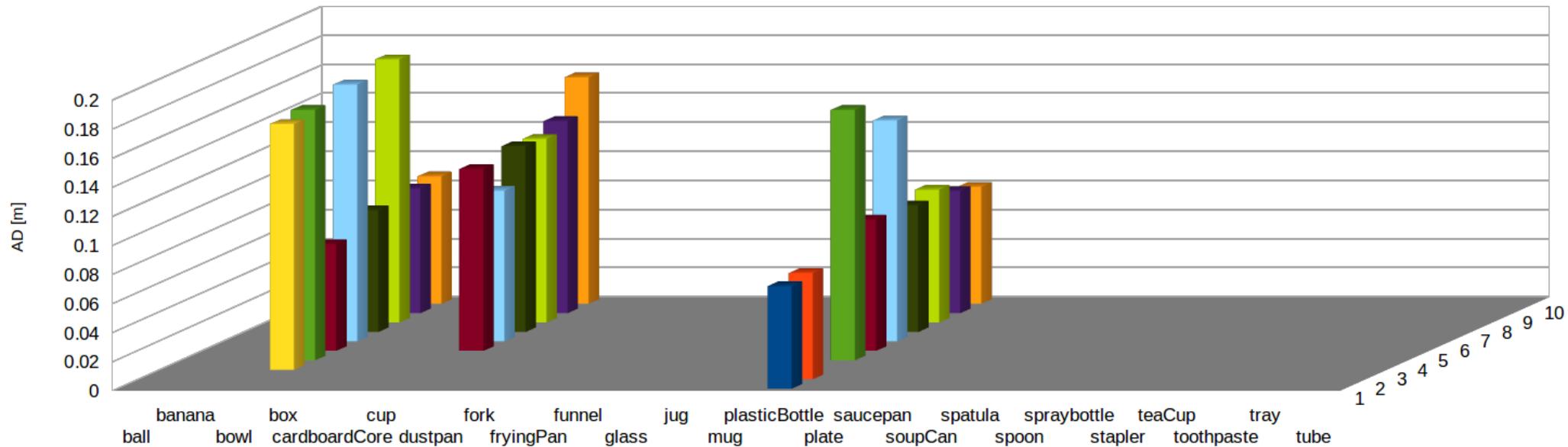




Preliminary Results

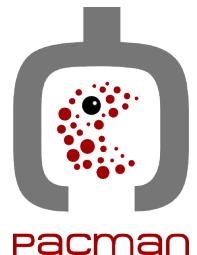


Up to 10 objects – Innsbruck:

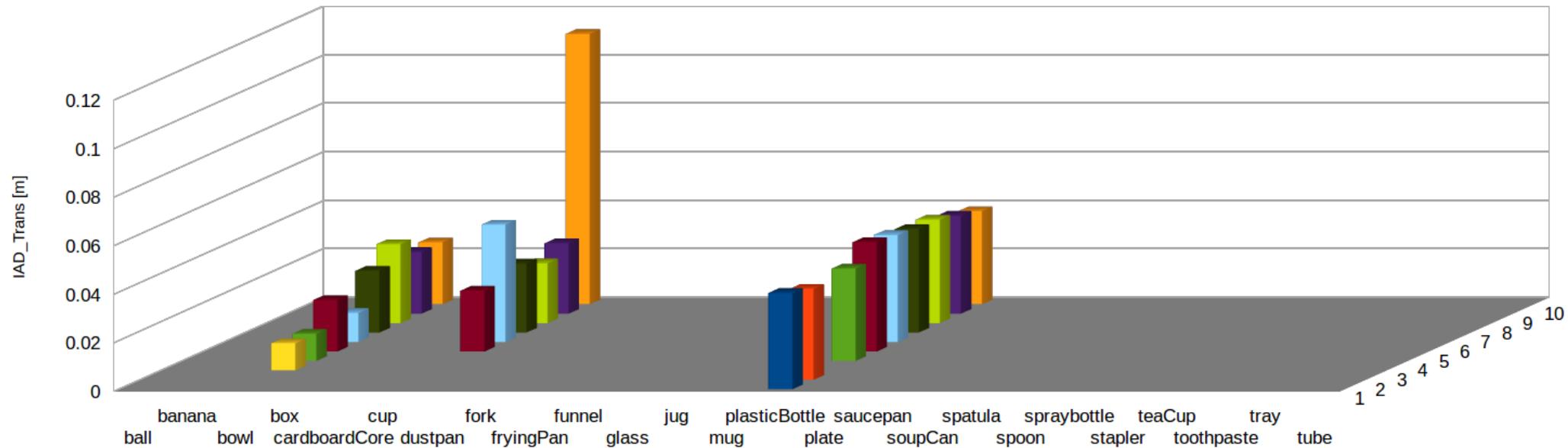




Preliminary Results

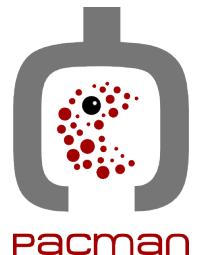


Up to 10 objects – Innsbruck:

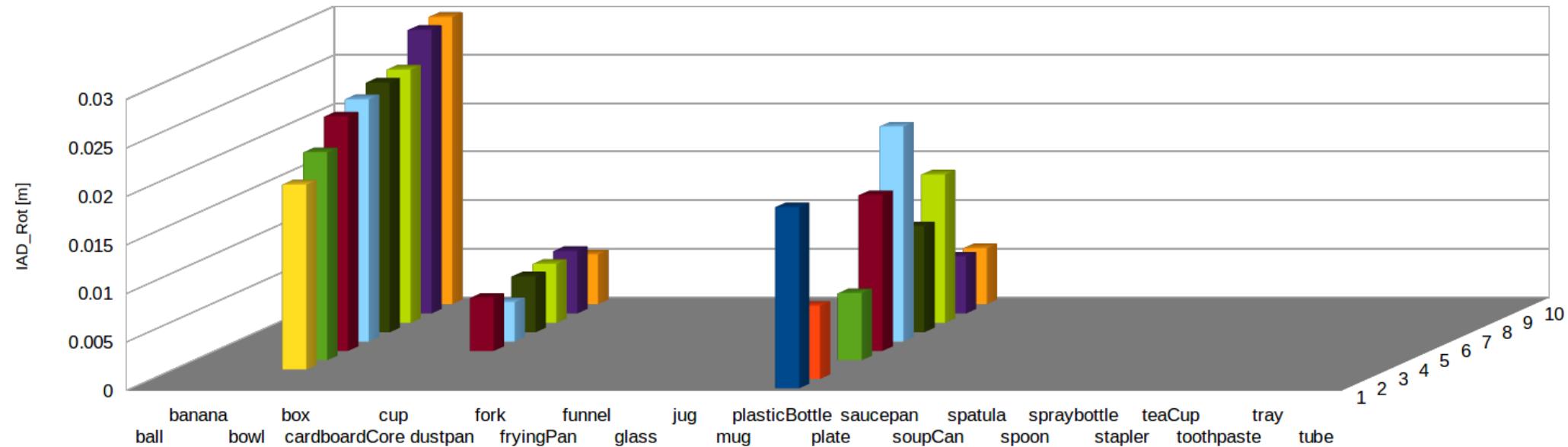




Preliminary Results

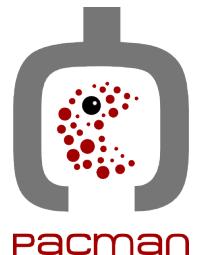


Up to 10 objects – Innsbruck:

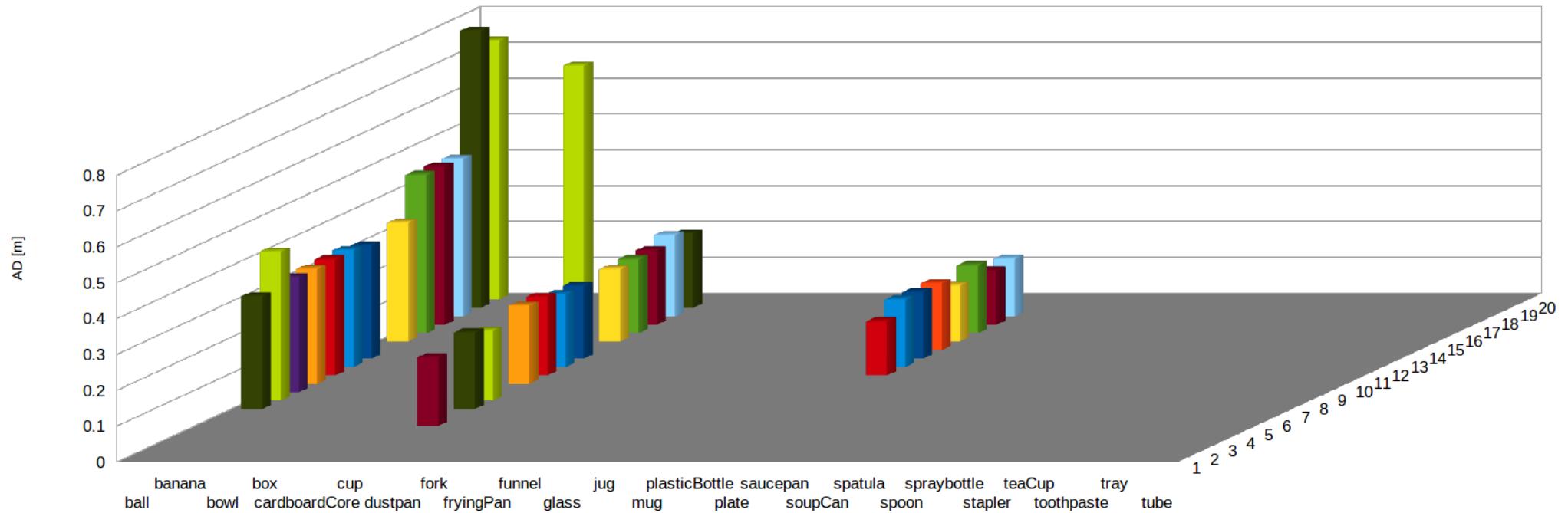




Preliminary Results

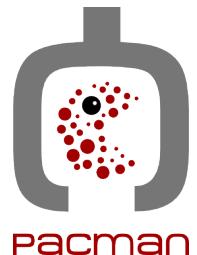


Up to 20 objects – Innsbruck:

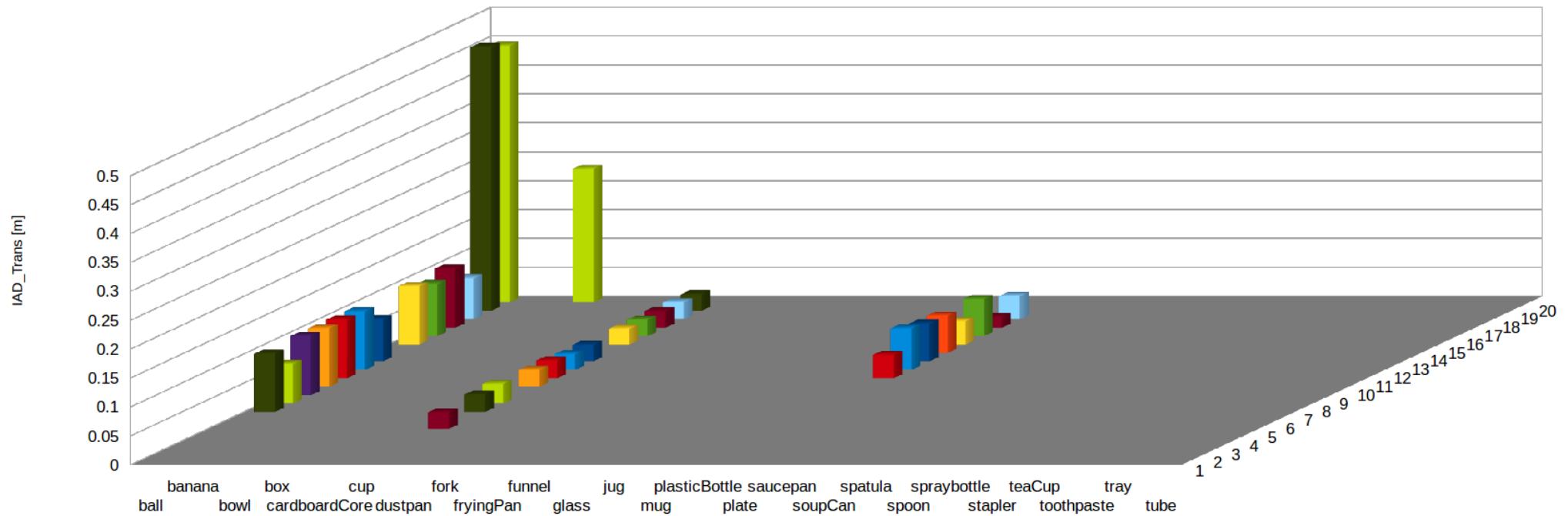




Preliminary Results

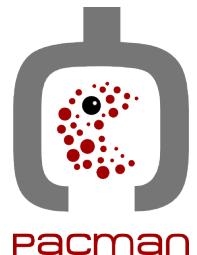


Up to 20 objects – Innsbruck:

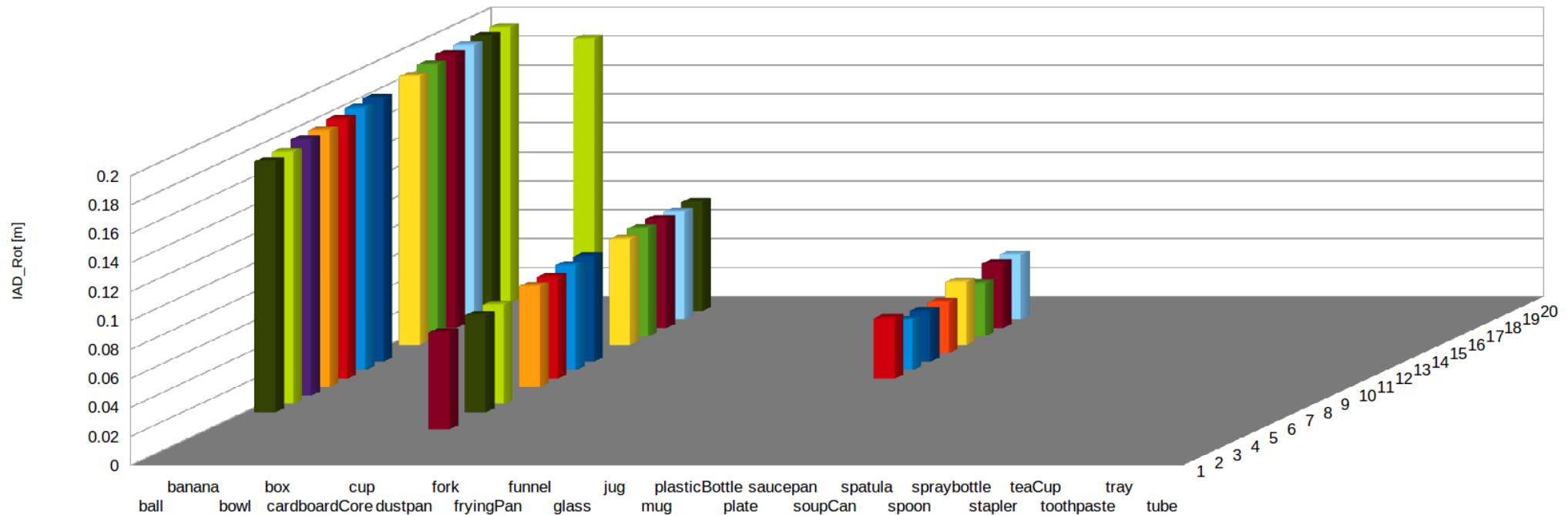




Preliminary Results

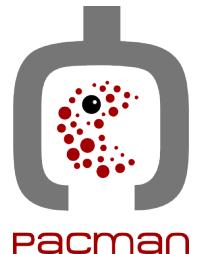


Up to 20 objects – Innsbruck:





Summary

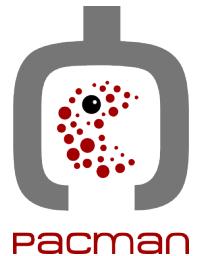


UoB HOOC :

- high level of occlusions
- symmetric objects
- learning on synthetic data
- testing on real data
- different evaluation measures



Future work



UoB HOOC :

- more attributes
- different object types
- standardized measures