

Overview

Problem

Previous approaches (Leopar¹ - [Same authors, CVPR'08]) showed that we can efficiently estimate the 3D pose of a poorly textured object by learning the patch appearance. Leopar¹ is performed in 4 steps:

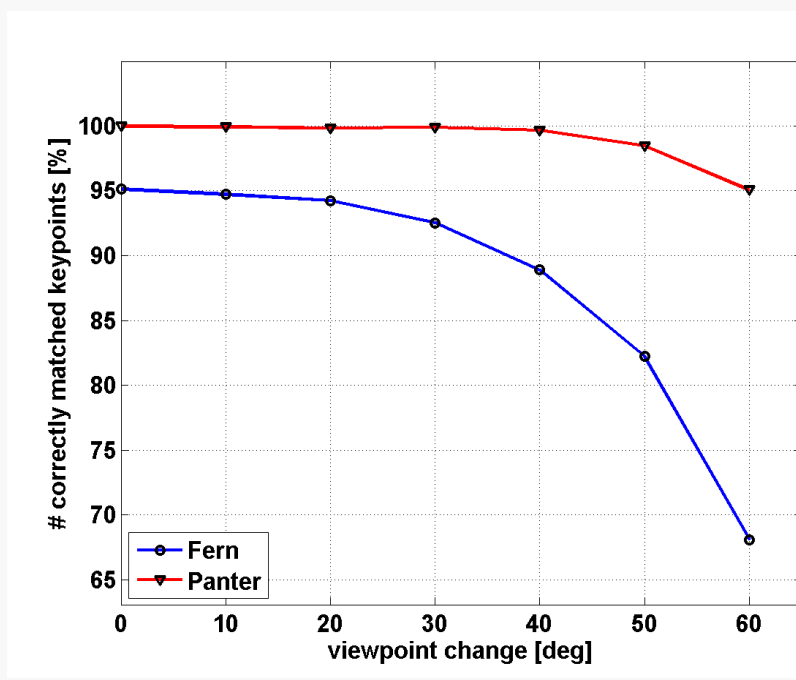
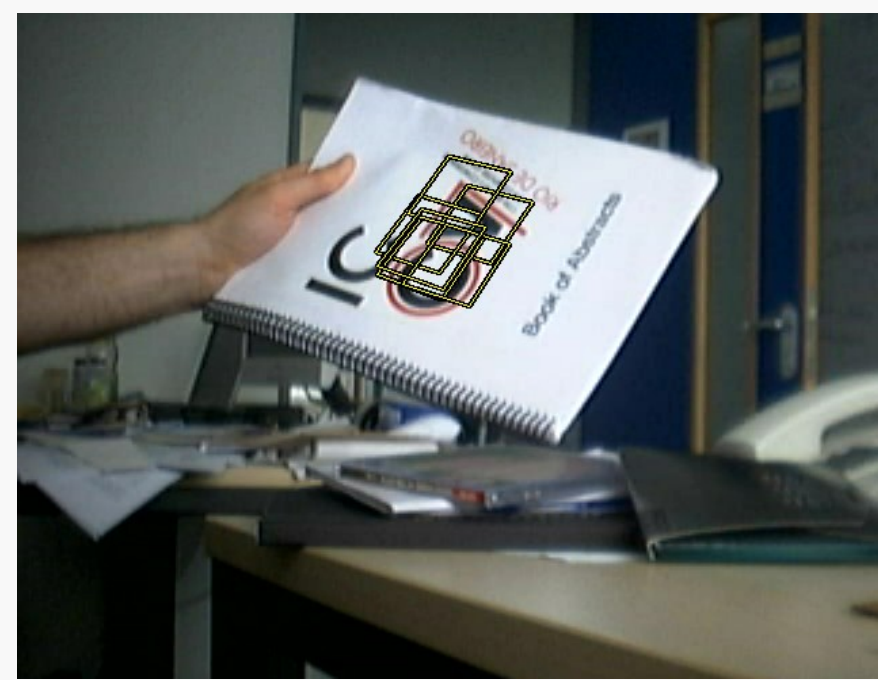
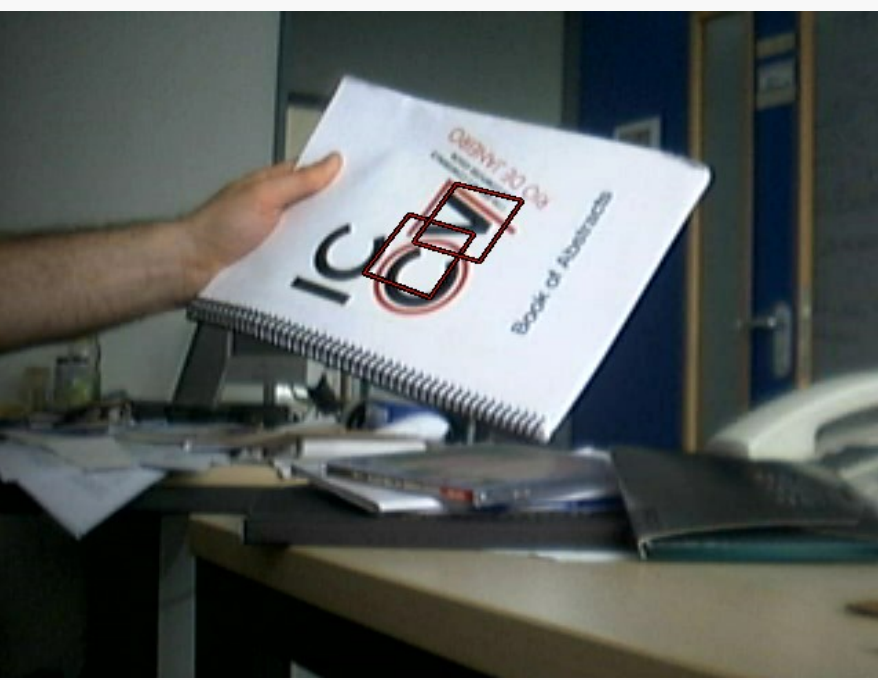
- Pre-classify feature points with e.g. Ferns²
- Rough orientation estimation of the patch with respect to the feature point identity
- Rectification refinement by applying a Template Matching algorithm³
- Outlier removal by simple correlation measure

Leopar¹ gives much better results concerning the accuracy and repeatability of the pose than affine region detectors but still suffers from a decreasing robustness towards large viewpoint changes. Reasons for that:

- Error-prone pre-classification of initial feature points in (a) (by Ferns²) without taking into account the pose of the patch
- Rough estimation of the pose of the patch in (b) is limited to orientations

Solution (Panter)

- No pre-classification but simultaneous estimation of keypoint identity and pose
- Estimation of the initial rectification not only from a small set of orientations but from a much larger set of real homography transformations



Patches extracted with Leopar¹

Patches extracted with Panter

Leopar¹ vs. Panter

Results

- Fast (~8-17fps) and very accurate tracking by detection
- More robust to large perspective distortions and scale changes than Leopar¹
- Only little texture and few feature points necessary to estimate the pose

More robust tracking by detection with better patch

retrieval than Leopar [CVPR'08]

!!! Pose estimation possible with only ONE feature point!!!

Simultaneous Estimation of Patch Identities and Pose

I. Extraction of Rough Rectification Information

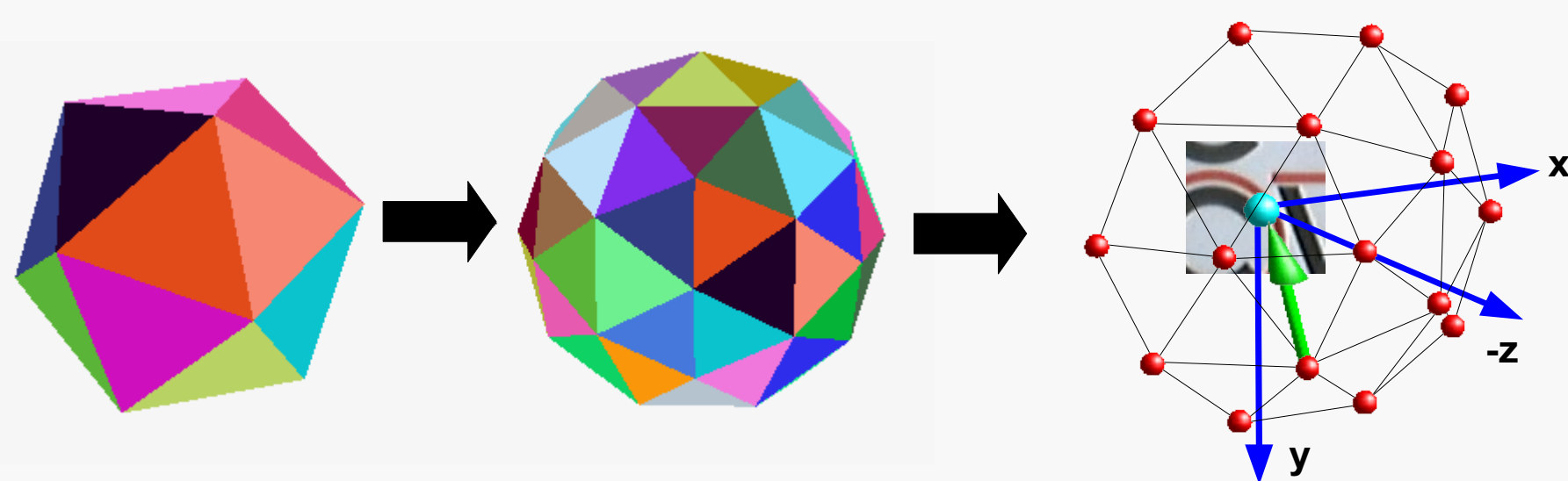
Given an image patch $p_{i,j}$ of Keypoint i under Pose j we want to find i' and j' (using linear classifiers for efficiency). Therefore, in the first step we assume that $i' = i$:

$$\forall j': a_{i,j'}^T p_{i,j'} = \begin{cases} +1 & j = j' \\ -1 & \text{otherwise} \end{cases}$$

$a_{i,j}$ can be learned in the training phase by a simple linear least squares method:

$$a_{i,j} = (PWW^T P^T)^{-1} PWW^T y$$

where P is the matrix made of patches column vectors samples $p_{i,j}$, y is the row vector made of +1 and -1 values and W is a diagonal matrix containing the equation weights $w_i = N-1$ for positive and $w_i = 1$ for negative examples under pose j (N is the number of poses). To ensure an equal distribution of the training samples of all important viewpoints we approximate a regular polyhedron with an initial icosahedron. The vertices of the approximated regular polyhedron serve as discretized viewpoints.



Training samples are only taken from a local neighborhood of the vertices.

At run-time we get for each patch $p_{i,j}$ and for each keypoint i in the database a list Γ_i of possible pose indices. We select for each keypoint i the best pose j that maximizes:

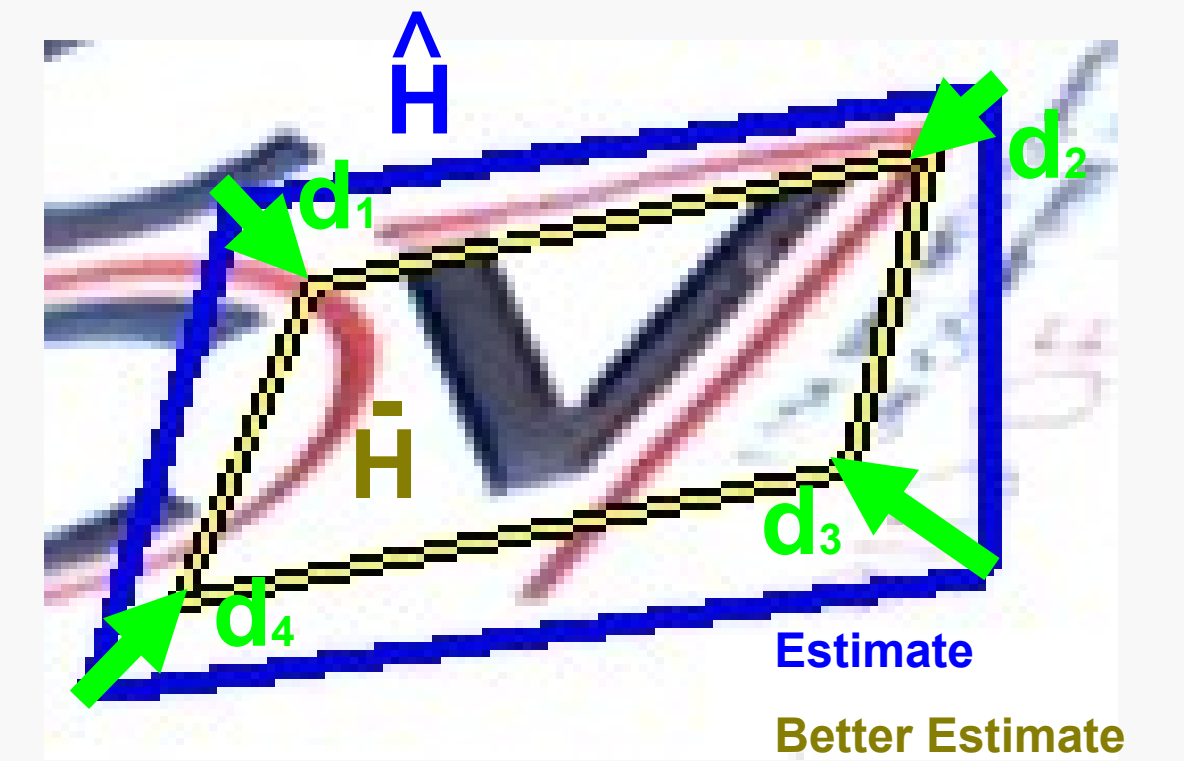
$$\forall i : \arg \max_{j \in \Gamma_i} p^T \bar{p}_{i,j}$$

where $\bar{p}_{i,j}$ is the mean of all positive examples of the patch of keypoint i under the pose j .

II. Iterative Refinement with Linear Predictors³

For each keypoint i in the database and its corresponding best rough pose estimate j we try to refine the patch rectification:

$$\begin{pmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{pmatrix} = B_i (p(H_j) - p_i^*)$$



where H_j is the initial homography estimate with respect to keypoint i , $p(H_j)$ the normalized intensity value vector of the patch under matrix H_j and p_i^* the normalized intensity vector of the reference patch. This equation has to be applied iteratively to converge to the right solution.

Matrix B_i can simply be learned by warping the patch of interest by small random pose changes δR and computing the normalized intensity changes δp .

$$B_i = XD^T (DD^T)^{-1} \quad \begin{array}{l} X: \text{the training matrix with the } \delta R \text{ vectors} \\ D: \text{the training matrix with the } \delta p \text{ vectors} \end{array}$$

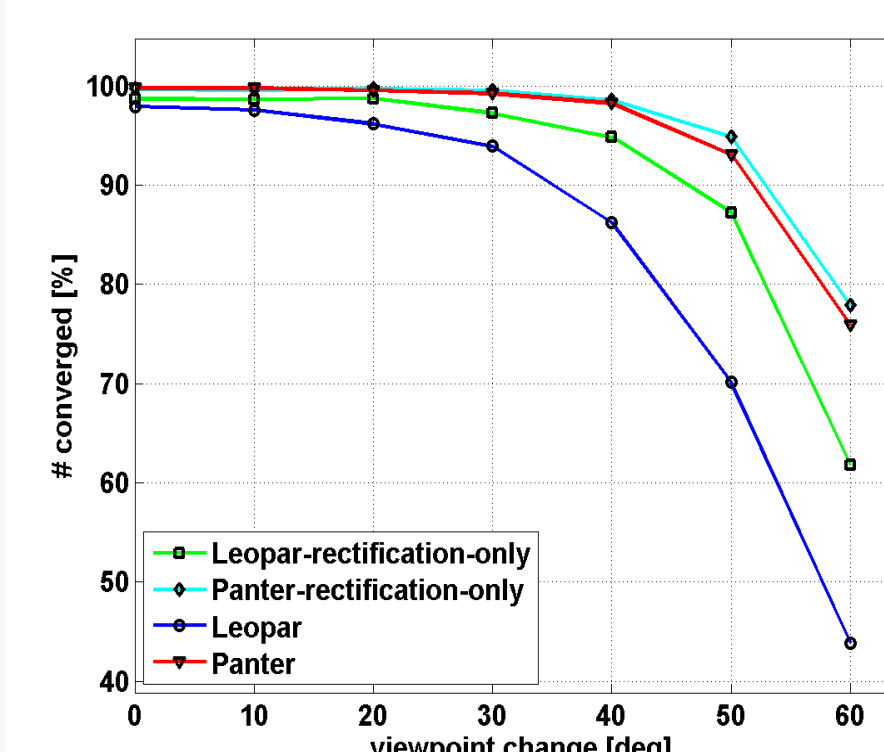
III. Establishing Patch Identity & Pose

Thanks to the accuracy of the Linear predictors, we can remove the outliers by simple cross-correlation between the warped patch and the reference patch:

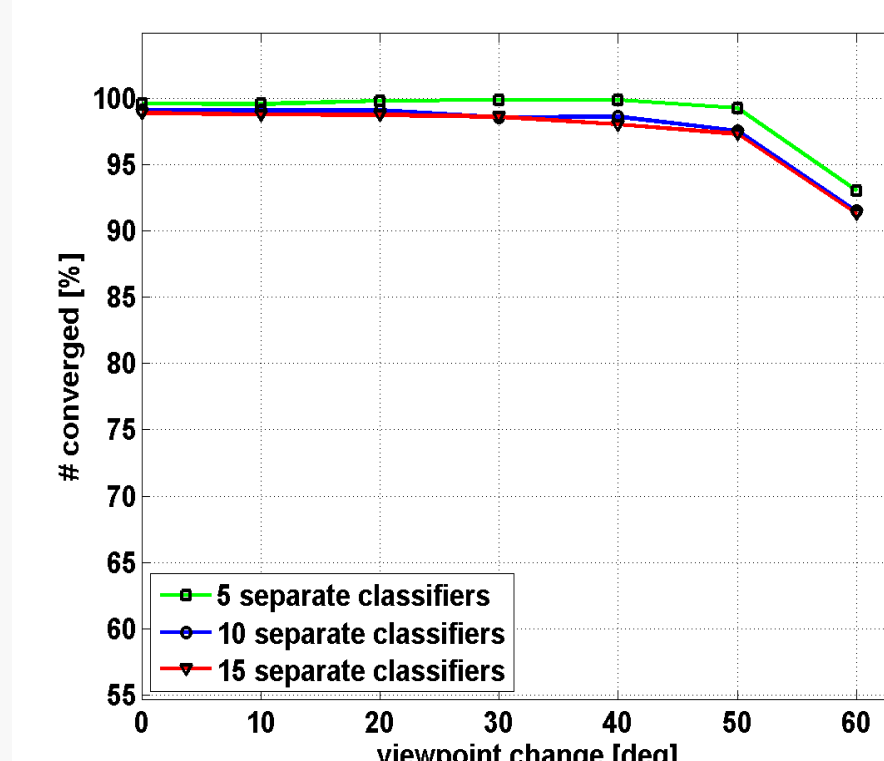
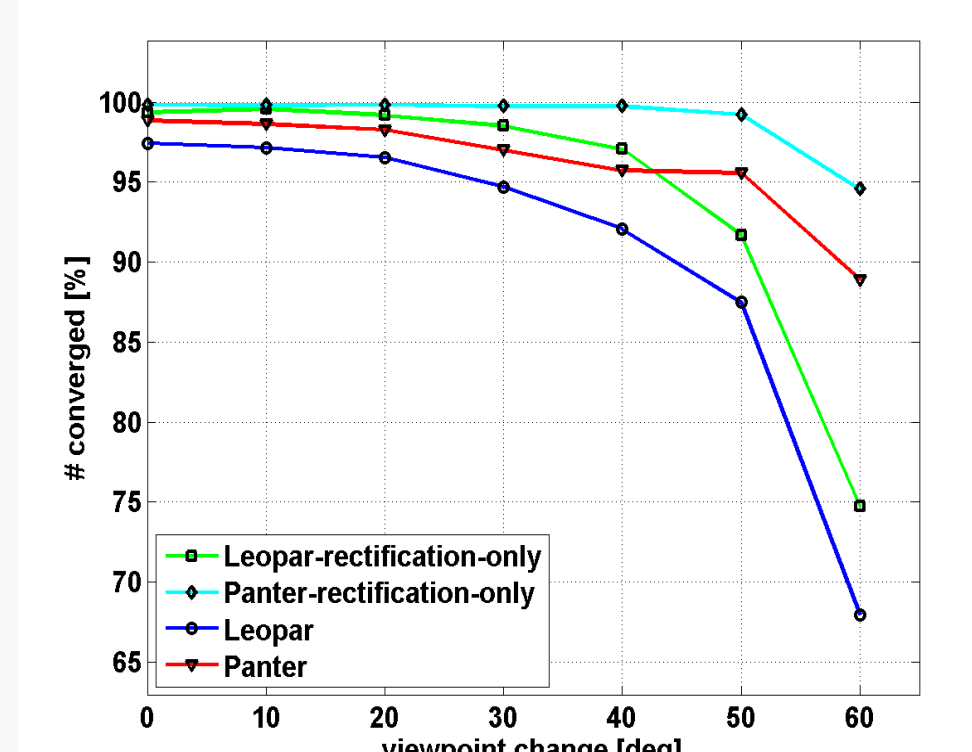
$$p(H_{final,i})^T \cdot p_i^* \geq \tau_{ncc}$$

where $H_{final,i}$ is the final transformation obtained with the linear predictor and i is the corresponding Keypoint identity. In practice we use a threshold $\tau_{ncc} = 0.9$.

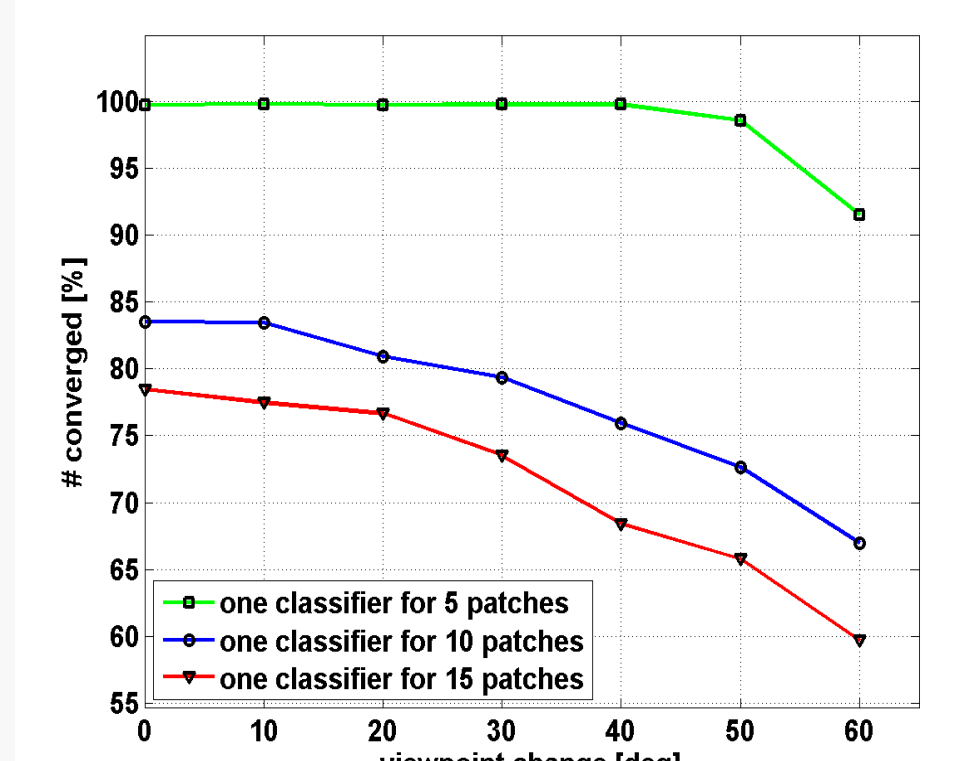
Experiments & Robust Real-Time Tracking by Detection



Panter outperforms Leopar¹ with respect to robustness against increasing viewpoint changes.



The performance of building one classifier for many patches is much lower than building one classifier for each patch.



The pose estimation with Panter is robust to high perspective distortions, to scale changes, to occlusion and even to some deformations.

References

- Hinterstößer, S. et al.: Online Learning of Patch Perspective Rectification for Efficient Object Detection, CVPR 2008, Anchorage, Alaska, USA.
- Ozuysal, M. et al.: Fast Keypoint Recognition in Ten Lines of Code, CVPR 2007, Minneapolis, Minnesota, USA.
- Jurie, F. et al.: Hyperplane approximation for template matching, PAMI, 2002