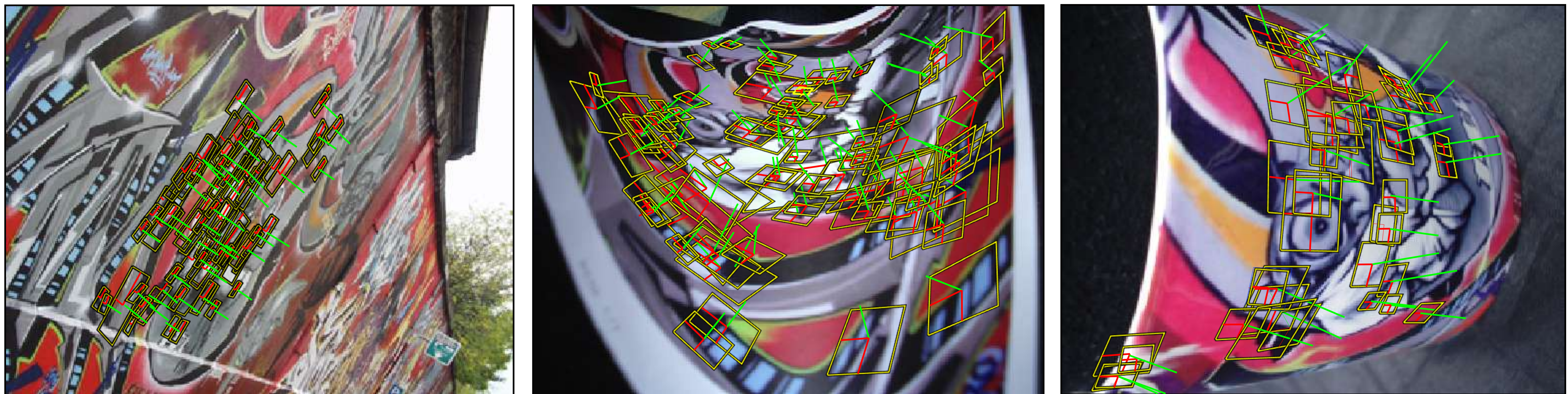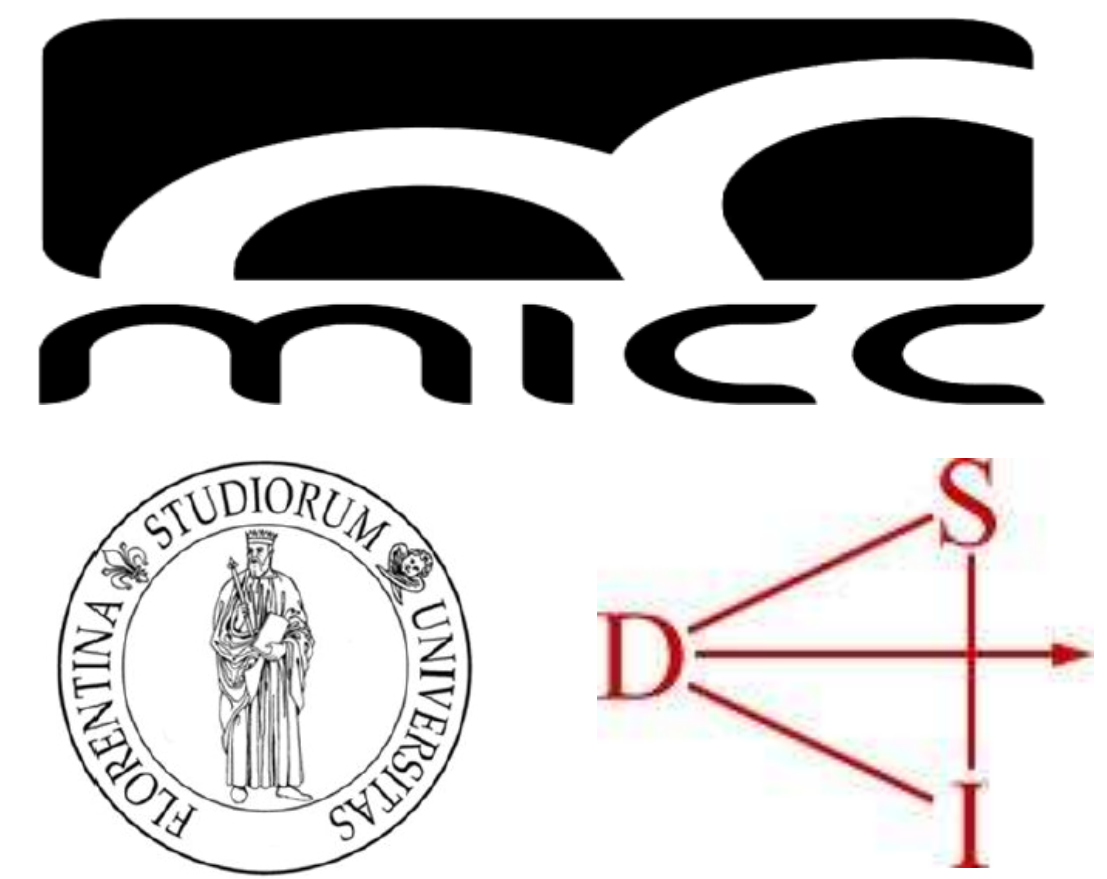# LOCAL 3D SURFACE POSE ESTIMATION BY NUISANCE RESIDUAL LEARNING

Del Bimbo A., **Franco F.** and Pernici F.
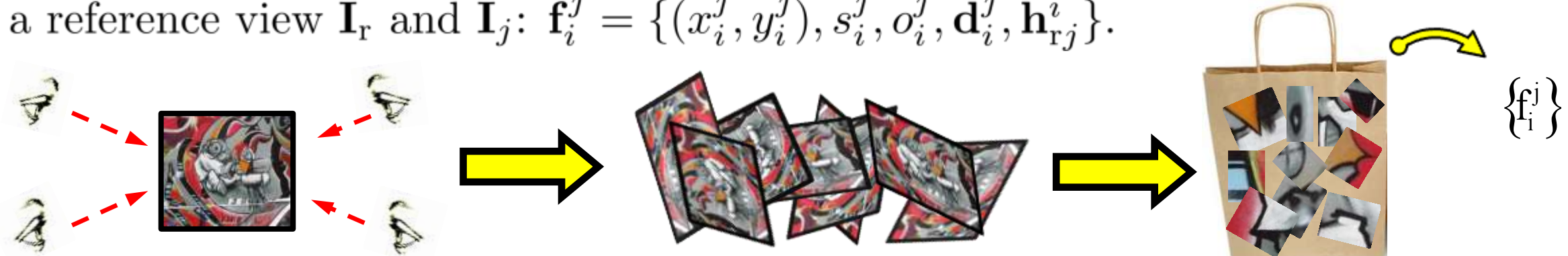
Media Integration and Communication Center (MICC), University of Florence, Italy
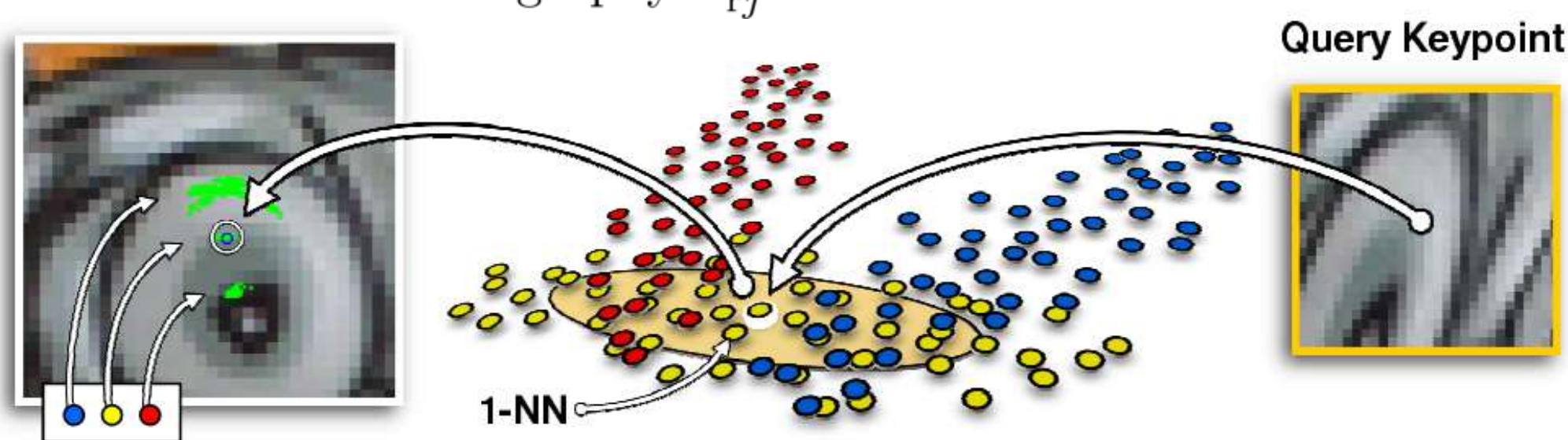
**Abstract:** We present a method of estimating the pose of an imaged scene surface element provided that it can be locally approximated by its tangent plane. The approach **simultaneously** learn the "nuisance residual" structure present in the detection and description steps of the SIFT algorithm allowing local perspective properties to be recovered through a homography. The estimated local poses can be applied to non rigid surfaces, with an accuracy representative of state-of-the-art for this challenging task.

Given a test image $I_t$ of one object taken from an arbitrary unknown point of view and a reference view $I_r$ of the object from which are taken multiple views with their homographies, we address the problem of finding the point of view from which $I_t$ is taken. In particular our goal is to estimate the homography $h_{rt}$ between $I_r$ and $I_t$, regardless of scale and perspective distortion that possibly affect the test image. To this end we assume that the surface of the 3D object has smooth curvature, so that the region around a keypoint can be considered locally planar.

**Training set generation.** Given a reference image $I_r$ of a planar surface, we take a full set of images taken from different viewpoints $I_j$ and extract SIFT keypoints in each view. Each keypoint $k_i^j$ in $I_j$ is associated with a representation feature that accounts for SIFT description, location, main orientation, scale at which it is detected, and the homographic transformation $h_{rj}^i$ between a reference view $I_r$ and $I_j$: $f_i^j = \{(x_i^j, y_i^j), s_i^j, o_i^j, d_i^j, h_{rj}^i\}$.



$\{f_i^j\}$

**Matching and geometry consistency check.** The set $C^t = \{f_i^j \in \mathcal{F} : i \in N_k(d^t)\}$ of corresponding features of keypoint $k^t$ is retrieved as the $k$-nearest neighbours $N_k$ in the space of SIFT appearance descriptors. Outliers are removed by back-projecting $f_i^j \in C^t$ in the coordinate system of $I_r$ using the inverse of its own homography $h_{rj}^{i-1}$.



Query Keypoint

1-NN

**Exploiting detection nuisance.** In order to obtain geometrically meaningful homographies from the dataset with respect to a test keypoint $k^t = \{(x^t, y^t), s^t, o^t, d^t\}$ and to **simultaneously** take into account the "residual" invariance of the SIFT detector and descriptor algorithm, the support region of retrieved features in $C^t$ are aligned. We calculate the shifts of: position $(u_i^j, v_i^j) = (x^t - \mu_x, y^t - \mu_y)$, and scale $(\sigma_i^j = \frac{s^t}{s_i^j})$, and orientation $(\theta_i^j = o^t - o_i^j)$ between features $f_i^j \in C^t$ and $k^t$, being $(\mu_x, \mu_y)$ the centroid of keypoints in $C^t$. The similarity transformation $s_{jt}^i = \begin{bmatrix} \sigma_i^j \cos \theta_i^j & -\sin \theta_i^j & u_i^j \\ \sin \theta_i^j & \sigma_i^j \cos \theta_i^j & v_i^j \\ 0 & 0 & 1 \end{bmatrix}$ accounts for scaling, translation and rotation and can be used to align the neighbourhoods of features $\{f_i^j\}$ to the region of input keypoint $k^t$, through the following homographic transformation: $h_{rt}^i = h_{rj}^i \cdot s_{jt}^i$

**Learning with kernel regression.** Let $D \in \mathbb{R}^{128}$ denotes a real valued random vector descriptor, and $H \in \mathbb{R}^8$ a real valued random output homographic transformation, with joint distribution $p_{D,H}(d, h)$. Given a training set consisting of N input-output pairs: $\{(d_i^j, h_{rt}^i)_1, ..., (d_i^j, h_{rt}^i)_N\}$ for which the probability of taking a specific value obeys $p_{D,H}(d, h)$, we wish to learn the function $\Phi : D \mapsto H$ defined as $h = \Phi(d) = [\Phi_1(d), ..., \Phi_8(d)]$ that maps from the space of input vector descriptors to the space of homographies. According to this, we can define a local regression estimate of $\Phi(d^t)$ as $\Phi_{\hat{\theta}}(d^t)$, which minimizes the cost: $L(h_{rt}^i, \Phi_\theta(d^t)) = \sum_{d_i^j \in C^t} K_\Sigma(d^t, d_i^j)[h_{rt}^i - \Phi_\theta(d_i^j)]^2$ and $\Phi_\theta$ is some parameterized function. We used the constant function model: $\Phi_\theta(d) = \theta_0$ in which case the final form of the kernel estimation is the Nadaraya-Watson weighted average: $h_{rt} = \hat{\Phi}(d^t) = \hat{\theta}_0 = \frac{\sum_i^N K_\Sigma(d^t, d_i^j)h_{rt}^i}{\sum_i^N K_\Sigma(d^t, d_i^j)}$. The kernel $K_\Sigma$ embeds feature descriptors of a cluster data into a vector space. Local neighborhood metric is specified by a Gaussian kernel function centered on the query descriptor: $K_\Sigma(d^t, d_i^j) = ((2\pi)^n |\Sigma|)^{\frac{1}{2}} \exp^{-\frac{1}{2}(d_i^j - d^t)^T \Sigma^{-1}(d_i^j - d^t)}$ where $\Sigma = N^{-1} \sum_{i=1}^N (d_i^j - d^t) \cdot (d_i^j - d^t)^T$.

**Experimental Results:** Several experiments were performed in order to assess the effectiveness of the method and compare it against "Leopar" and "Caspar" [1] [2]. In the first set of experiment we compare the average overlap between the quadrangles obtained with our method and those obtained using the ground truth homography (Fig. 1). Similarly to "Leopar" this overlap is very close to 100% for our method. The second set shows the comparison of the mean reprojection error for the quadrangle corners (Fig. 2). The error of the patch corner is less than one pixel in average and outperforms other methods. In the last set we prove the stability of our method with respect to the number of multiview descriptor-homography pairs used in the Kernel based regression (Fig. 3). Error bars show that we obtain a mean reprojection error of less than 1 pixel over almost the range from 30 to 250 descriptors. It also shows that the minimum error is obtained in the range from 140 to 165 descriptor-homography pairs. Our current implementation runs at about 15 frame per second using 400.000 features in the database and extracting about 200 SIFT keypoints in the input images, on a standard notebook with an Intel Centrino Core Duo with 2.4GHz and 3Gb RAM. The average times for the most expensive steps are shown in the table and compares favorably with state-of-the-art methods.

| step | average time |
|---|---|
| SIFT Point Extraction | 0.025 sec |
| Geometry Consistency Check | 0.021 sec |
| Kernel Based Regression | 0.036 sec |

**References:** [1] S. Hinterstoisser, S. Benhimane, N. Navab, P. Fua, and V. Lepetit, "Online learning of patch perspective rectification for efficient object detection," in CVPR, 2008. [2] A. Pagani and D. Stricker, "Learning local patch orientation with a cascade of sparse regressors" in BMVC, 2009.
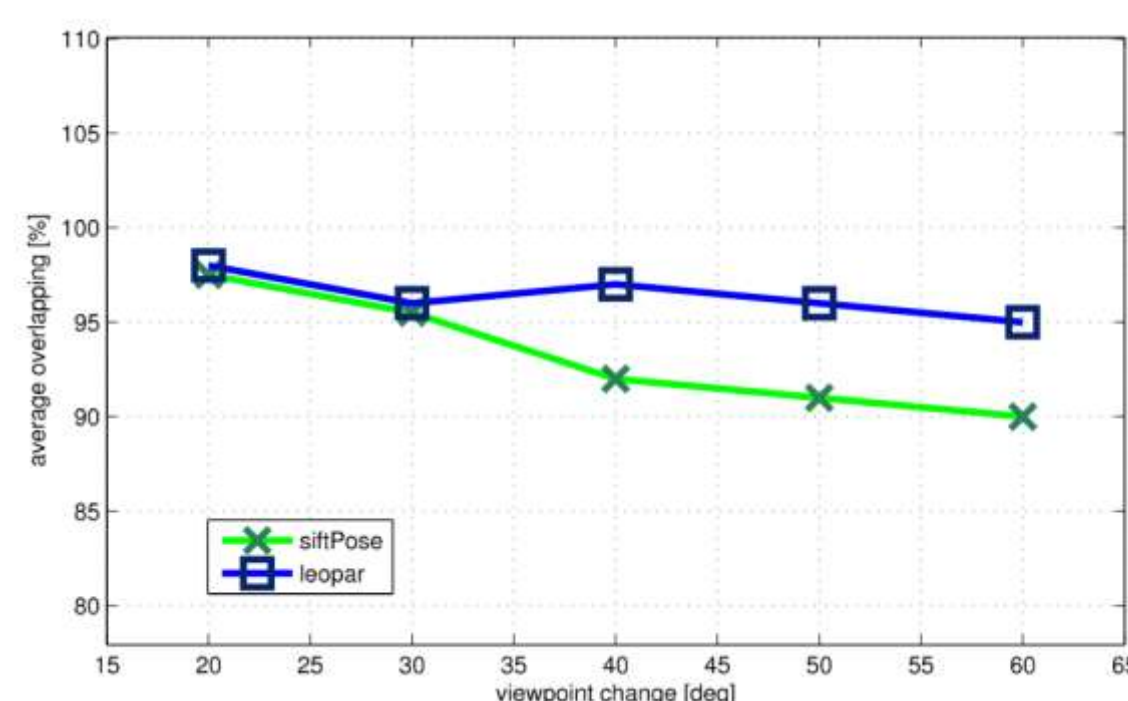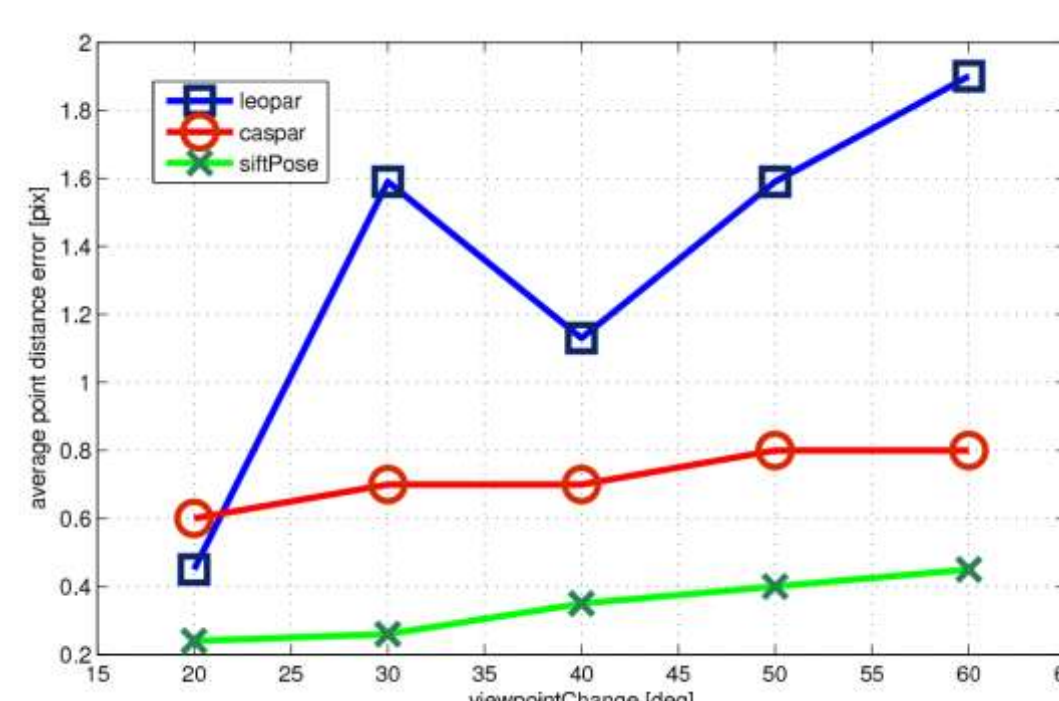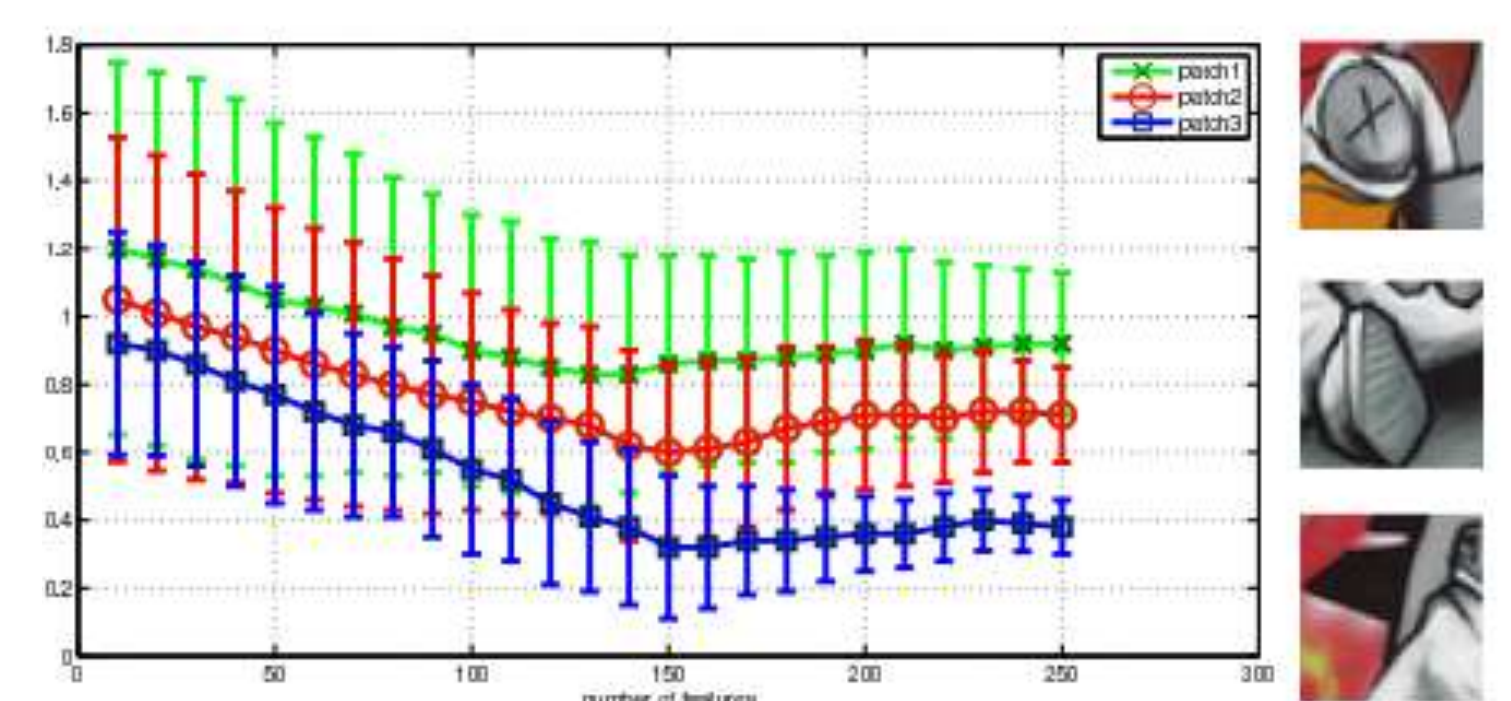
Figure 1



Figure 2



Figure 3