

Pose Estimation from Local Features

David Demirdjian¹ Raquel Urtas

Raquel Urtasun² Christopher Wilkens²



¹ Toyota Technical Center, Cambridge, USA

² Massachusetts Institute of Technology, Cambridge, USA

Introduction

Our goal is to perform pose estimation from single images. Our main contribution is a generalization of the Implicit Shape Model – introduced earlier for object recognition by Leibe & Schiele – to learn the relationship between articulated pose and the appearance of reference patches.

Recovering accurate pose from single images requires a more sophisticated model than the Implicit Shape Model. In our approach, flexibility in the model is obtained by introducing hyper-parameters that are learned from training data and optimally weigh the spatial and pose contribution of patches. In our framework, the appearance of a patch is modeled using multiple descriptors. A task-specific metric in the patch space is learned for best performance on pose estimation.

Pose Estimation from Local Features

Individual local features provide rich and useful information on the global pose. e.g.

Location of 'landmarks' in the neighborhood of local features
 Scale, configuration, ...



Approach

We employ a technique similar to the Implicit Shape Model introduced

- **Database construction**. Learn a set of local patches (and associated descriptors) and pose *x* relative to their location in the image.
- Vocabulary tree. Build a vocabulary tree using *k*-means clustering using the *Pyramid Match Toolkit* by J. Lee et al.

 Pose estimation. For each local feature of a query image, retrieve its kNNs in the database and compute the density by kernel-based estimation. Pose is found by maximizing the density (voting)

Benefits of our approach:

Robustness to occlusions Does not require segmentation or normalization Fast

Relative/absolute pose

In order to incorporate geometric relationships, the relative pose *x* associated with a local feature *I* is given with respect to the feature location I_i and scale σ_i . At testing, the absolute pose is *reconstructed* from the feature location 1 and scale σ .





Density estimation

Given a query image *I*, descriptors *d* are estimated at multiple locations *l* and scales σ and matched against the database. The probability of observing pose x given image Z is:

$$\phi(x) = p(x | Z) = \sum_{j=1}^{n} p(x | e_j, l_j, \sigma_j)$$

Density associated with a single local feature

$$p(x \mid e, l, \sigma) = \sum_{i} p(x \mid I_i, l, \sigma) p(I_i \mid e)$$

Implicit Shape Model

War

The probability of observing pose *x* given a patch *I* at location *l* and scale σ is modeled using a Gaussian model. The model gives more confidence to pose 'parts' m_k closer to patch locations *l*

$$p(x \mid I, l, \sigma) = N(x \mid F(x_{I}, l, \sigma), \Sigma_{I})$$
ping function
$$F(x_{I}, l, \sigma) = \frac{\sigma}{\sigma_{I}} (x_{I} - L_{I}) + L \qquad L = (l, ..., l)^{T}$$

$$L_{I} = (l_{I}, ..., l_{I})^{T}$$

$$\Sigma_{I} = \frac{1}{(s^{2})^{2}} diag(\Lambda_{1}, ..., \Lambda_{N}) \qquad \Lambda_{k} = (m_{k} - l_{k})(m_{k} - l_{i})$$

 \boldsymbol{s} is the covariance scale. It corresponds to the patch reliability and is learned from training data

Task-Oriented Patch Metric

The patch likelihood $p(I_i|e)$ is defined as a linear combination of kernels

$$p(I_i | e) = d^T K(I_i | e) \quad d = (d_1, \dots, d_N) \quad \text{with} \quad |d| = 1$$
$$K = (K_1, \dots, K_N)$$
$$\forall e, \sum K_j(I_i | e)$$

Experiments

Articulated Body Pose

Evaluation on the HumanEva database [L. Sigal and M. J. Black.]. We used a subset of the database containing sequences of *Walking, Jogging* and *Boxing*. Sequences were evenly divided into training and testing sets



						ISM-pose
	SIFT	SC	CC	Multi-D	Walking	4.16 (2.71)
Leibe et al.	11.32 (7.58)	12.36 (5.10)	13.52 (6.29)		Jogging	4.29 (1.68)
ISM-pose	5.01 (3.58)	5.13 (3.62)	5.47 (4.21)	4.73 (2.58)	Boxing	7.22 (5.03)

Mean and variance of the errors (in pixels) on the HumanEva testing set.





Ternary diagram showing the average error vs. descriptor weights d

Error Mean (top) and variance (bottom) vs. covariance scale *s*

Fiducial

Estimation of the 2D location of 20 fiducial points (e.g. nose, mouth, eye corners) BiolD database – 1521 images – 23 users

- Database generated from 1000 training images ~ 90000 patches.
- Testing set consists of 100 images (subjects different from training images)



Occlusions Occlusions were synthetically generated (black rectangles) in real images



Example of results for the occlusion experiment