



Universitat Autònoma de Barcelona

MASTER IN COMPUTER VISION AND ARTIFICIAL INTELLIGENCE  
**REPORT OF THE RESEARCH PROJECT**  
OPTION: COMPUTER VISION

# **Pose and Face Recovery via Spatio-temporal GrabCut Human Segmentation**

**Author: Antonio Hernández Vela**  
**Date: 13/07/2010**  
**Advisor: Sergio Escalera Guerrero**

# Pose and Face recovery via Spatio-temporal GrabCut Human Segmentation

**Antonio Hernández**

*Computer Vision Center*

*Campus UAB, Edifici O*

*08193 Bellaterra, Barcelona, Spain*

ahernandez@cvc.uab.cat

**Supervisor:** Sergio Escalera and Petia Radeva

## Abstract

In this paper, we present a full-automatic Spatio-Temporal GrabCut human segmentation methodology which benefits from the combination of tracking and segmentation. GrabCut initialization is performed by a HOG-based subject detection, face detection, and skin color model for seed initialization. Spatial information is included by means of Mean Shift clustering whereas temporal coherence is considered by the historical of Gaussian Mixture Models. Moreover, full face and pose recovery is obtained by combining human segmentation with Active Appearance Models and Conditional Random Fields. Results over public data sets as well as in a new Human Limb data set show a robust segmentation and recovery of both face and pose using the presented methodology.

**Keywords:** Human segmentation, GrabCut, Pose recovery, Face modelling

## 1. Introduction

Human segmentation in uncontrolled environments is a hard task because of the constant changes produced in natural scenes: illumination changes, moving objects, changes in the point of view, or occlusions, just to mention a few. Because of the nature of the problem, a common way to proceed is to discard most part of the image so that the analysis can be performed on a reduced set of small candidate regions. In Dalal and Triggs (2005), the authors propose a full-body detector based on a cascade of classifiers (Viola and Jones (2004)) using HOG features. This methodology is currently being used in several works related to the pedestrian detection problem (Geronimo et al. (2009); Andriluka et al. (2008)). GrabCut (Rother et al. (2004)) has also shown high robustness in Computer Vision segmentation problems, defining the pixels of the image as nodes of a graph and extracting foreground pixels via iterated Graph Cut optimization. This methodology has been applied to the problem of human body segmentation with high success (Ferrari et al. (2008, 2009)). Many other works involving GrabCut -and graph-cuts in general- can be found (Chen et al. (2008); Lombaert et al. (2005); Kwolek (2009); Nagahashi et al. (2007)), what indicates that is a state-of-the-art methodology in image segmentation. In the case of working with sequences of images, this optimization problem can also be considered to have temporal coherence. In the work of Corrigan et al. (2008), the authors extended the Gaussian Mixture Model (GMM) of GrabCut algorithm so that the color space is complemented with the derivative in time of pixel intensities in order to include temporal information in the segmentation optimization process. However, the main problem of that method is that moving pixels corresponds to the boundaries between foreground and background regions, and thus, there is no clear discrimination.

Once a region of interest is determined, pose is often recovered by the determination of the body limbs together with their spatial coherence (also with temporal coherence in case of image sequences). Most of these approaches are probabilistic, and features are usually based on edges or 'appearance'. In Ramanan (2006), the author propose a probabilistic approach for limb detection based on edge learning complemented with color information. The image of probabilities is then formulated in a Conditional Random Field scheme and optimized using belief propagation. This work has obtained robust results and has been extended by other authors including local GrabCut segmentation and temporal refinement of the CRF model ( Ferrari et al. (2008, 2009)).

In this paper, we propose a full-automatic Spatio-Temporal GrabCut human segmentation methodology which benefits from the combination of tracking and segmentation. First, subjects are detected by means of a HOG-based cascade of classifiers. Face detection and skin color model are used to define a set of seeds used to initialize GrabCut algorithm. Spatial information is taken into account by means of Mean Shift clustering, whereas temporal information is considered taking into account the pixel probability membership to an historical of Gaussian Mixture Models. Moreover, the methodology is combined with Shape and Active Appearance Models to define three different meshes of the face, one near frontal view, and the other ones near lateral views. Temporal coherence and fitting cost are considered in conjunction with GrabCut segmentation to allow a smooth and robust face fitting in video sequences. Finally, the limb detection and a CRF model are applied on the obtained segmentation, showing high robustness capturing body limbs due to the accurate human segmentation. In order to test the proposed methodology, we use public data sets and present a new Human Limb data set useful for human segmentation, limb detection, and pose recovery purposes.

The rest of the paper is organized as follows: Section 2 describes the proposed methodology, presenting the spatio-temporal GrabCut segmentation, the Active Appearance models for face fitting, and the pose recovery methodology. Experimental results on public and novel data sets are performed in Section 3. Finally, Section 4 concludes the paper.

## 2. Full-body pose recovery

In this section, we present the Spatio-Temporal GrabCut methodology to deal with the problem of automatic human segmentation in video sequences. Then, we describe the Active Appearance Models used to recover the face, and the body pose recovery methodology based on the approach of Ramanan (2006). All methods presented in this section are combined to improve final segmentation and pose recovery. Figure 1 illustrates the different modules of the project.

### 2.1 GrabCut segmentation

In Rother et al. (2004), the authors proposed an approach to find a binary segmentation -Background, Foreground- of an image by formulating an energy minimization scheme as the one presented in Boykov and Jolly (2001); Boykov and Funka-Lea (2006); Kolmogorov and Zabih (2004), extended using color instead of just gray-scale information. Given a color image  $I$ , let us consider the array  $z = (z_1, \dots, z_n, \dots, z_N)$  of  $N$  pixels where  $z_i = (R_i, G_i, B_i)$ ,  $i \in [1, \dots, N]$  in RGB space. The segmentation is defined as array  $\alpha = (\alpha_1, \dots, \alpha_N)$ ,  $\alpha_i \in \{0, 1\}$ , assigning a label to each pixel of the image indicating if it belongs to background or foreground. A trimap  $T$  is defined by the user -in a semi-automatic way-, consisting on three regions:  $T_B$ ,  $T_F$  and  $T_U$ , each one containing initial background, foreground, and uncertain pixels, respectively. Pixels belonging to  $T_B$  and  $T_F$  are clamped as background and foreground respectively -that means GrabCut will not be able to modify these

labels-, whereas those belonging to  $T_U$  are actually the ones the algorithm will be able to label. Color information is introduced by GMMs. A full covariance GMM of  $K$  components is defined for background pixels ( $\alpha_i = 0$ ), and another one for foreground pixels ( $\alpha_j = 1$ ), parametrized as follows

$$\boldsymbol{\theta} = \{\pi(\alpha, k), \mu(\alpha, k), \Sigma(\alpha, k), \alpha \in \{0, 1\}, k = 1..K\}, \quad (1)$$

being  $\pi$  the weights,  $\mu$  the means and  $\Sigma$  the covariance matrices of the model. We also consider the array  $\mathbf{k} = \{k_1, \dots, k_i, \dots, k_N\}$ ,  $k_i \in \{1, \dots, K\}$ ,  $i \in [1, \dots, N]$  indicating the component of the background or foreground GMM (according to  $\alpha_i$ ) the pixel  $z_i$  belongs to. The energy function for segmentation is then

$$\mathbf{E}(\boldsymbol{\alpha}, \mathbf{k}, \boldsymbol{\theta}, \mathbf{z}) = \mathbf{U}(\boldsymbol{\alpha}, \mathbf{k}, \boldsymbol{\theta}, \mathbf{z}) + \mathbf{V}(\boldsymbol{\alpha}, \mathbf{z}), \quad (2)$$

where  $\mathbf{U}$  is the likelihood potential, based on the probability distributions  $p(\cdot)$  of the GMM:

$$\mathbf{U}(\boldsymbol{\alpha}, \mathbf{k}, \boldsymbol{\theta}, \mathbf{z}) = \sum_i -\log p(z_i | \alpha_i, k_i, \boldsymbol{\theta}) - \log \pi(\alpha_i, k_i) \quad (3)$$

and  $\mathbf{V}$  is a regularizing prior assuming that segmented regions should be coherent in terms of color, taking into account a neighborhood  $C$  around each pixel

$$\mathbf{V}(\boldsymbol{\alpha}, \mathbf{z}) = \gamma \sum_{\{m,n\} \in C} [\alpha_n \neq \alpha_m] \exp(-\beta \|z_m - z_n\|^2) \quad (4)$$

With this energy minimization scheme and given the initial trimap  $T$ , the final segmentation is performed using a minimum cut algorithm ( Boykov and Kolmogorov (2001); Boykov and Jolly

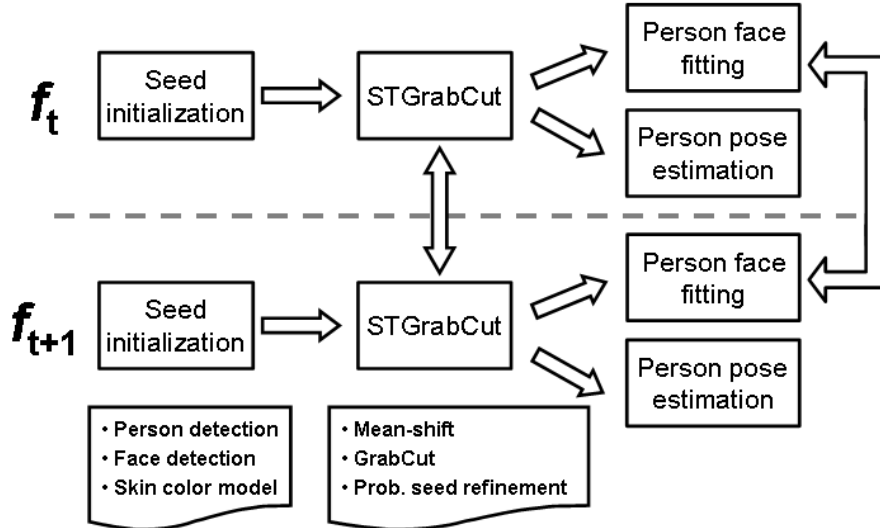


Figure 1: Overall block diagram of the methodology.

(2001); Boykov and Funka-Lea (2006)). The classical semi-automatic GrabCut algorithm is summarized in Algorithm 1.

---

**Algorithm 1 Original GrabCut algorithm.**

---

- 1: Trimap  $T$  initialization with manual annotation.
  - 2: Initialize  $a_i = 0$  for  $n \in T_B$  and  $a_i = 1$  for  $n \in T_U \cup T_F$ .
  - 3: Initialize Background and Foreground GMMs from sets  $a_n = 0$  and  $a_n = 1$  respectively, with  $k$ -means.
  - 4: Assign GMM components to pixels.
  - 5: Learn GMM parameters from data  $z$ .
  - 6: Estimate segmentation: Graph-cuts.
  - 7: Repeat from step 4, until convergence.
- 

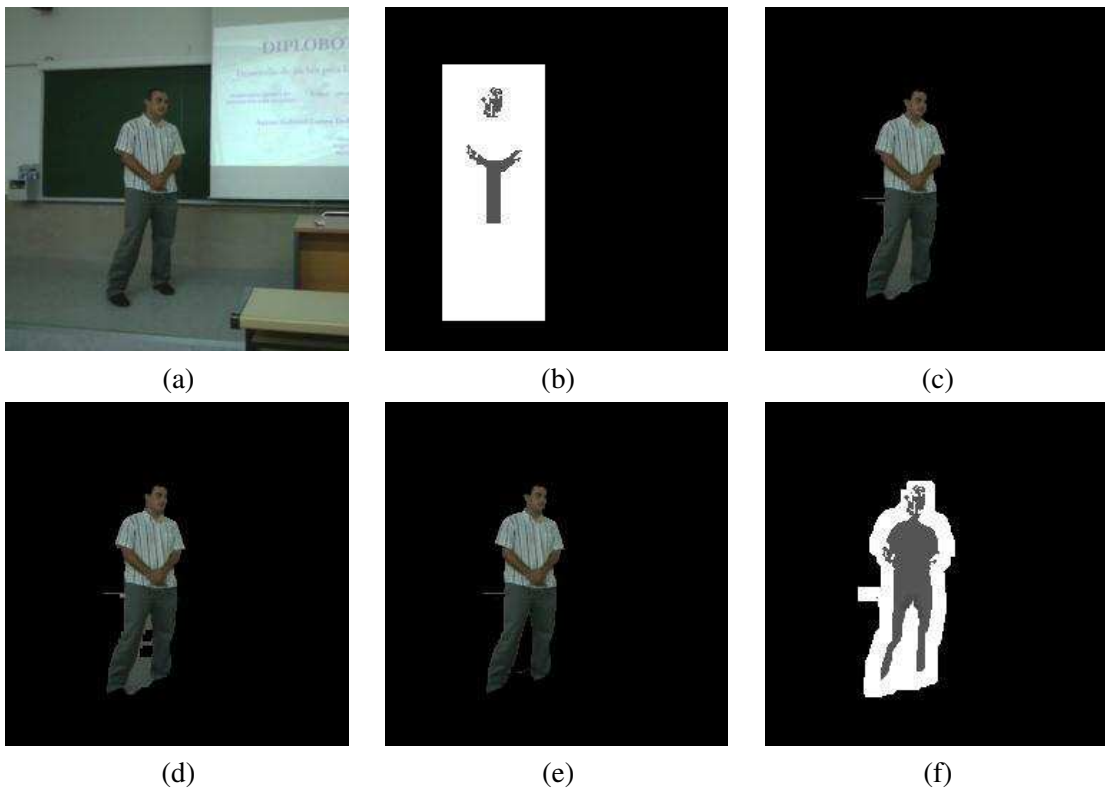


Figure 2: STGrabcut pipeline example: (a) Original frame, (b) Seed initialization, (c) GrabCut, (d) Probabilistic re-assignment, (e) Refinement and (f) Initialization mask for  $f_{t+1}$ .

## 2.2 Automatic initialization

Our proposal bases on the previous GrabCut framework, focusing on human body segmentation, being fully automatic, and extending it by taking into account temporal coherence. We refer to each frame of the video as  $f_t$ ,  $t \in \{1, \dots, M\}$  being  $M$  the length of the sequence. Given a frame  $f_t$ , we first apply a person detector based on a cascade of classifiers using HOG features (Dalal and Triggs (2005)). Then, we initialize the trimap  $T$  from the bounding box  $B$  returned by the detector:  $T_U = \{z_i \in B\}$ ,  $T_B = \{z_i \notin B\}$ . Furthermore, in order to increase the accuracy of the segmentation algorithm, we include Foreground seeds exploiting spatial and appearance prior information. On one hand, we define a small central region  $R$  inside  $B$  and set these pixels as Foreground. On the other, we apply a face detector based on a cascade of classifiers using Haar-like features Viola and Jones (2004) over  $B$ , and learn a skin color model  $h_{skin}$ . All pixels inside  $B$  fitting in  $h_{skin}$  are also set to foreground. Therefore, we initialize  $T_F = \{z_i \in R\} \cup \{z_i \in \delta(z_i, h_{skin})\}$ , where  $\delta$  returns the set of pixels belonging to the color model defined by  $h_{skin}$ . An example of seed initialization is shown in Figure 2(b).

## 2.3 Spatial extension

Once we have initialized the trimap, we can apply the iterative minimization algorithm shown in steps 4 to 7 of original GrabCut (Algorithm 1). However, instead of applying  $k$ -means for the initialization of the GMMs we propose to use Mean-Shift clustering, which also takes into account spatial coherence. Given an initial estimation of the distribution modes  $m_h(\mathbf{x}^0)$  and a kernel function  $g$ , Mean-shift iteratively updates the mean-shift vector with the following formula:

$$\mathbf{m}_h(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{x}_i g(\|\frac{\mathbf{x}-\mathbf{x}_i}{h}\|^2)}{\sum_{i=1}^n g(\|\frac{\mathbf{x}-\mathbf{x}_i}{h}\|^2)} \quad (5)$$

until it converges, where  $\mathbf{x}_i$  contains the value of pixel  $z_i$  in CIELuv space and its spatial coordinates, and returns the centers of the clusters (distribution modes) found. After convergence, we obtain a segmentation  $\alpha^t$  and the updated foreground and background GMMs  $\theta^t$  at frame  $f_t$ , which are used for further initialization at frame  $f_{t+1}$ . The result of this step is shown in Figure 2(c). Finally, we refine the segmentation of frame  $f_t$  eliminating false positive foreground pixels. By definition of the energy minimization scheme, GrabCut tends to find convex segmentation masks having a lower perimeter, given that each pixel on the boundary of the segmentation mask contributes on the global cost. Therefore, in order to eliminate these background pixels (commonly in concave regions) from the foreground segmentation, we re-initialize the trimap  $T$  as follows

$$\begin{aligned} T_B &= \{z_i | \alpha_i = 0\} \cup \\ &\quad \left\{ z_i \mid \sum_{k=t-j}^t \frac{p(z_i | \alpha_i = 0, k_i, \theta^k)}{j} > \sum_{k=t-j}^t \frac{p(z_i | \alpha_i = 1, k_i, \theta^k)}{j} \right\} \\ T_F &= \{z_i \in \delta(z_i, h_{skin})\} \\ T_U &= \{z_i | \alpha_i = 1\} \setminus T_B \setminus T_F \end{aligned} \quad (6)$$

where the pixel background probability membership is computed using the GMM models of previous segmentations. This formulation can also be extended to detect false negatives. However,

in our case we focus on false positives since they appear frequently in the case of human segmentation. The result of this step is shown in Figure 2(d). Once the trimap has been redefined, false positive foreground pixels still remain, so the new set of seeds is used to iterate again GrabCut algorithm, resulting in a more accurate segmentation, as we can see in Figure 2(e).

## 2.4 Temporal extension

Considering  $A$  as the binary image representing  $\alpha$  at  $f_t$  -the one obtained before the refinement-, we initialize the trimap for  $f_{t+1}$  as follows

$$\begin{aligned} T_F &= \{z_i \in I \mid z_i \in A \ominus ST_e, \alpha(z_i) = 1\} \\ T_U &= \{z_i \in I \mid z_i \in A \oplus ST_d, \alpha(z_i) = 1\} \setminus T_F \\ T_B &= \{z_i, z_i \in I\} \setminus (T_F \cup T_U) \end{aligned} \quad (7)$$

where  $\ominus$  and  $\oplus$  are erosion and dilation operations with their respective structuring elements  $ST_e$  and  $ST_d$ , and  $\alpha_i := \alpha(z_i)$ . The structuring elements are simple squares of a given size depending on the size of the person and the degree of movement we allow from  $f_t$  to  $f_{t+1}$ , assuming smoothness in the movement of the person. An example of a morphological mask is shown in Figure 2(f). Spatial information could be also included in the mean-shift algorithm in conjunction with color and spatial information. However, we included this information explicitly to be anisotropic. The whole segmentation methodology is detailed in the ST-GrabCut Algorithm 2.

---

### Algorithm 2 Spatio-Temporal GrabCut algorithm.

---

- 1: Person detection on  $f_1$ .
  - 2: Face detection and skin color model learning.
  - 3: Trimap  $T$  initialization with detected bounding box and learnt skin color model.
  - 4: Initialize  $a_i = 0$  for  $n \in T_B$  and  $a_i = 1$  for  $n \in T_U \cup T_F$ .
  - 5: Initialize Background and Foreground GMMs from sets  $a_n = 0$  and  $a_n = 1$  respectively, with Mean-shift.
  - 6: **for**  $t = 1 \dots M$
  - 7:   Person detection on  $f_t$ .
  - 8:   Assign GMM components to pixels of  $f_t$ .
  - 9:   Learn GMM parameters from data  $z$ .
  - 10:   Estimate segmentation: Graph-cuts.
  - 11:   Repeat from step 8, until convergence.
  - 12:   Re-initialize trimap  $T$  (equation 6).
  - 13:   Assign GMM components to pixels.
  - 14:   Learn GMM parameters from data  $z$ .
  - 15:   Estimate segmentation: Graph-cuts.
  - 16:   Repeat from step 12, until convergence.
  - 17:   Initialize trimap  $T$  using segmentation obtained in step 11 after convergence (equation 7) for  $f_{t+1}$ .
  - 18: **end for**
-

## 2.5 Face fitting

Once we have properly segmented the body region, next step consists of fitting the face and the body limbs. For the case of face recovery, we base our procedure on mesh fitting using Active Appearance Models (AAM), that benefits from Active Shape Models and color and texture information Cootes et al. (a).

Active Appearance Model is generated by combining a model of shape and texture variation. First, a set of points are marked on the face of the training images that are aligned, and a statistical shape model is build Cootes et al. (b). Each training image is warped so the points match those of the mean shape. This is raster scanned into a texture vector,  $\mathbf{g}$ , which is normalized by applying a linear transformation,  $\mathbf{g} \mapsto (\mathbf{g} - \mu_g \mathbf{1})/\sigma_g$ , where  $\mathbf{1}$  is a vector of ones, and  $\mu_g$  and  $\sigma_g^2$  are the mean and variance of elements of  $\mathbf{g}$ . After normalization,  $\mathbf{g}^T \mathbf{1} = 0$  and  $|\mathbf{g}| = 1$ . Then, principal component analysis is applied to build a texture model. Finally, the correlations between shape and texture are learnt to generate a combined appearance model. The appearance model has parameter  $\mathbf{c}$  controlling the shape and texture according to

$$x = \bar{x} + \mathbf{Q}_s \mathbf{c} \quad (8)$$

$$g = \bar{g} + \mathbf{Q}_g \mathbf{c} \quad (9)$$

where  $\bar{x}$  is the mean shape,  $\bar{g}$  the mean texture in a mean shaped patch, and  $\mathbf{Q}_s, \mathbf{Q}_g$  are matrices designing the modes of variation derived from the training set. A shape  $\mathbf{X}$  in the image frame can be generated by applying a suitable transformation to the points,  $\mathbf{x} : \mathbf{X} = S_t(\mathbf{x})$ . Typically,  $S_t$  will be a similarity transformation described by a scaling  $s$ , an in-plane rotation,  $\theta$ , and a translation  $(t_x, t_y)$ .

Once constructed the AAM, it is deformed on the image to detect and segment the face appearance as follows. During matching, we sample the pixels in the region of interest  $\mathbf{g}_{im} = T_u(\mathbf{g}) = (u_1 + 1)\mathbf{g}_{im} + u_2 \mathbf{1}$ , where  $\mathbf{u}$  is the vector of transformation parameters, and project into the texture model frame,  $\mathbf{g}_s = T_u^{-1}(\mathbf{g}_{im})$ . The current model texture is given by  $\mathbf{g}_m = \bar{g} + \mathbf{Q}_g \mathbf{c}$ , and the difference between model and image (measured in the normalized texture frame) is as follows

$$\mathbf{r}(\mathbf{p}) = \mathbf{g}_s - \mathbf{g}_m \quad (10)$$

Given the error  $E = |\mathbf{r}|^2$ , we compute the predicted displacements  $\delta \mathbf{p} = -\mathbf{R}\mathbf{r}(\mathbf{p})$ , where  $\mathbf{R} = \left( \frac{\partial \mathbf{r}^T}{\partial \mathbf{p}} \frac{\partial \mathbf{r}}{\partial \mathbf{p}} \right)^{-1} \frac{\partial \mathbf{r}^T}{\partial \mathbf{p}}$ . The model parameters are updated  $\mathbf{p} \mapsto \mathbf{p} + k\delta \mathbf{p}$ , where initially  $k = 1$ . The new points  $\mathbf{X}'$  and model frame texture  $\mathbf{g}'_m$  are estimated, and the image is sampled at the new points to obtain  $\mathbf{g}'_{mi}$  and the new error vector  $\mathbf{r}' = T_u^{-1}(\mathbf{g}'_{mi}) - \mathbf{g}'_m$ . A final condition guides the end of each iteration: if  $|\mathbf{r}'|^2 < E$ , then we accept the new estimate, otherwise, we set to  $k = 0.5$ ,  $k = 0.25$ , and so on. The procedure is repeated until no improvement is made to the error.

Taking into account the discontinuity that appears when a face moves from frontal to profile view, we use three different AAM corresponding to three meshes of 21 points: frontal view  $\mathfrak{S}_F$ , right lateral view  $\mathfrak{S}_R$ , and left lateral view  $\mathfrak{S}_L$ . In order to include temporal and spatial coherence, meshes at frame  $f_{t+1}$  are initialized by the fitted mesh points at frame  $f_t$ . Additionally, we include a temporal change-mesh control procedure, as follows

$$\mathfrak{S}^{t+1} = \min_{\mathfrak{S}^{t+1}} \{E_{\mathfrak{S}_F}, E_{\mathfrak{S}_R}, E_{\mathfrak{S}_L}\}, \mathfrak{S}^{t+1} \in \nu(\mathfrak{S}^t) \quad (11)$$

where  $\nu(\mathfrak{S}^t)$  corresponds to the meshes contiguous to the mesh  $\mathfrak{S}^t$  fitted at time  $t$  (including the same mesh), and  $E_{\mathfrak{S}_i}$  is the fitting error cost of mesh  $\mathfrak{S}_i$ . This constraint avoids false jumps and



imposes smoothness in the temporal face behavior (e.g. a jump from right to left profile view is not allowed).

In order to obtain a more accurate pose estimation, after fitting the mesh, we take advantage of its variability to differentiate among a set of head poses. We have defined five different head poses: right, middle-right, frontal, middle-left, and left. In order to define this set, the fitted frontal meshes in the training set are classified in three different poses: middle-right, frontal, and middle left, whereas the training samples of the left and right meshes are directly classified in full-left and full-right poses, respectively. In order to learn the five different head poses, training images are aligned, and PCA is applied to save the 20 most representative eigenvectors. Then, a new mesh is projected to that new space and classified to one of the five different head poses according to a 3-Nearest Neighbor rule.

Figure 3 shows examples of the AAM model fitting in images (obtained from Huang et al. (2007)) for the three different meshes.



Figure 3: From left to right: left, frontal, and right mesh fitting using AAM.

## 2.6 Pose recovery

Considering the refined segmented body region obtained using the proposed ST-GrabCut algorithm, we construct a pictorial structure model (Felzenszwalb and Huttenlocher; Ronfard et al. (2002)). We use the method of Ramanan (Ramanan (2006); Ferrari et al. (2009)), which captures the appearance and spatial configuration of body parts. A person's body parts are tied together in a tree-structured conditional random field. Parts,  $l_i$ , are oriented patches of fixed size, and their position is parameterized by location  $(x, y)$  and orientation  $\phi$ . The posterior of a configuration of parts  $L = l_i$  given a frame  $f_t$  is

$$P(L|f_t) \propto \exp \left( \sum_{(i,j) \in E} \Psi(l_i, l_j) + \sum_i \Phi(l_i|f_t) \right) \quad (12)$$

The pairwise potential  $\Psi(l_i, l_j)$  corresponds to a spatial prior on the relative position of parts and embeds the kinematic constraints. The unary potential  $\Phi(l_i|I)$  corresponds to the local image evidence for a part in a particular position. Inference is performed over tree-structured conditional random field.

Since the appearance of the parts is initially unknown, a first inference uses only edge features in  $\Phi$ . This delivers soft estimates of body part positions, which are used to build appearance models of the parts and background (color histograms). Inference is then repeated with  $\Phi$  using both edges and appearance. This parsing technique simultaneously estimates pose and appearance of parts. For each body part, parsing delivers a posterior marginal distribution over location and orientation  $(x, y, \phi)$  (Ramanan (2006); Ferrari et al. (2009)).

### 3. Results

Before the presentation of the results, we discuss the data, methods and parameters of the comparative, and validation measurements.

- *Data*: We use the public image sequences of the Chroma Video Segmentation Ground Truth (cVSG, Tiburzi et al. (2008)), a corpus of video sequences and segmentation masks of people. Chroma based techniques have been used to record Foregrounds and Backgrounds separately, being later combined to achieve final video sequences and accurate segmentation masks almost automatically. Some samples of the sequence we have used for testing are shown in Figure 4(a). The sequence has a total of 307 frames. This image sequence includes several critical factors that make segmentation difficult: object textural complexity, object structure, uncovered extent, object size, Foreground and Background velocity, shadows, background textural complexity, Background multimodality, and small camera motion.

As a second database we have also used a set of 30 videos corresponding to the defense of undergraduate thesis at the University of Barcelona to test the methodology in a different environment (UBDataset). Some samples of this data set are shown in Figure 4(b).

Moreover, we present the Human Limb data set, a new data set composed by 227 images from 25 different people. At each image, 14 different limbs are labeled (see Figure 4(c)), including the "don't care" label between adjacent limbs, as described in Figure 5. Backgrounds are from different real environments with different visual complexity. This data set is useful for human segmentation, limb detection, and pose recovery purposes<sup>1</sup>.

- *Methods*: We test the classical semi-automatic GrabCut algorithm for human segmentation comparing with the proposed ST-GrabCut algorithm. We also test the mesh fitting and body pose recovery methodologies on the obtained segmentations.

- *Validation measurements*: In order to evaluate the robustness of the methodology for human body segmentation, face and pose fitting, we use the ground truth masks of the images to compute the overlapping factor  $O$  as follows

$$O = \frac{\sum M_{GC} \cap M_{GT}}{\sum M_{GC} \cup M_{GT}} \quad (13)$$

where  $M_{GC}$  and  $M_{GT}$  are the binary masks obtained for spatio-temporal GrabCut segmentation and the ground truth mask, respectively.

#### 3.1 Spatio-temporal GrabCut Segmentation

First, we test the proposed ST-GrabCut segmentation on the sequence from the public cVSG corpus. The results for the different experiments are shown in Table 1. In order to avoid the manual

---

1. The data set is public at <http://www.maia.ub.es/~sergio/Code.html>



(a)



(b)



(c)

Figure 4: (a) Samples of the cVSG corpus and (b) UBDataset image sequences, and (c) Human-Limb data set.

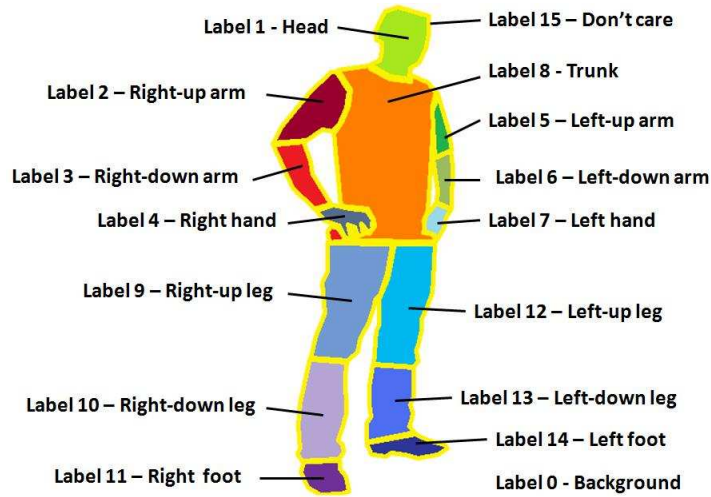


Figure 5: Human Limb data set labels description.

Approach	Mean overlapping
K-means	0.5356
Mean-shift	0.5424
Morphology	0.6229
ST-GrabCut	0.8747

Table 1: GrabCut and ST-GrabCut Segmentation results on cVSG corpus.

initialization of classical GrabCut algorithm, for all the experiments, seed initialization is performed applying the commented person HOG detection, face detection, and skin color model. First row of Table 1 shows the overlapping performance of eq.(13) applying GrabCut segmentation with  $k$ -means clustering to design the GMM models. Second row shows the overlapping performance considering Mean Shift clustering to design the GMM models. One can see a slight improvement when using the second strategy. This is mainly due to the fact that Mean Shift clustering takes into account spatial information of pixels in clustering time, which better defines contiguous pixels of image to belong to GMM models of foreground and background. Third performance in Table 1 shows the overlapping results considering the morphology refinement based on previous segmentation. In this case, we obtain near 10% of performance improvement respect the previous result. Finally, last result of Table 1 shows the full-automatic ST-GrabCut segmentation overlapping performance. One can see that it achieves about 25% of performance improvement in relation with the previous best performance. Some segmentation results obtained by the GrabCut algorithm for the cVSG corpus are shown in Figure 6. Note that the ST-GrabCut segmentation is able to robustly segment convex regions. We have also applied the ST-GrabCut segmentation methodology on the image sequences of UBdataset. Some segmentations are shown in Figure 6.

### 3.2 Face fitting

In order to measure the robustness of the spatio-temporal AAM mesh fitting methodology, we performed the overlapping analysis of meshes in both un-segmented and segmented image sequence of



Figure 6: Segmentation examples of (a) UBDataset sequence 1, (b) UBDataset sequence 2 and (c) cVSG sequence.

the public cVSG corpus. Overlapping results are shown in Table 3. One can see that the mesh fitting works fine in unsegmented images, obtaining a final mean overlapping of 89.60%. However, note that combining the temporal information of previous fitting and the ST-GrabCut segmentation, the face mesh fitting considerably improves, obtaining a final of 96.36% of overlapping performance. Some example of face fitting using the AAM meshes for different face poses of the cVSG corpus are shown in Figure 7.

Finally, we have tested the classification of the five face poses on the cVSG corpus, obtaining the percentage of frames of the subject at each pose. The obtained percentages are shown in Table 3.

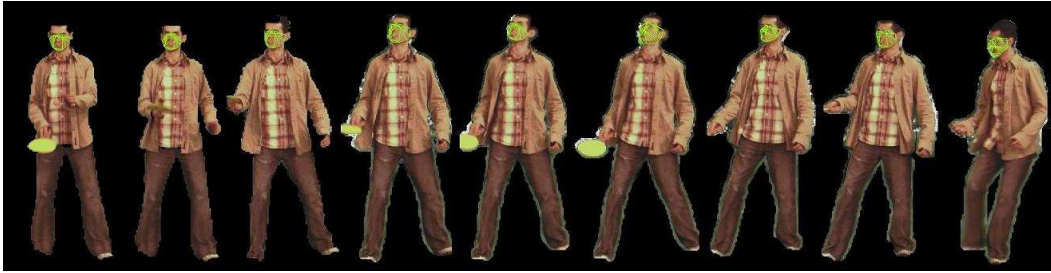


Figure 7: Samples of the segmented cVSG corpus image sequences fitting the different AAM meshes.

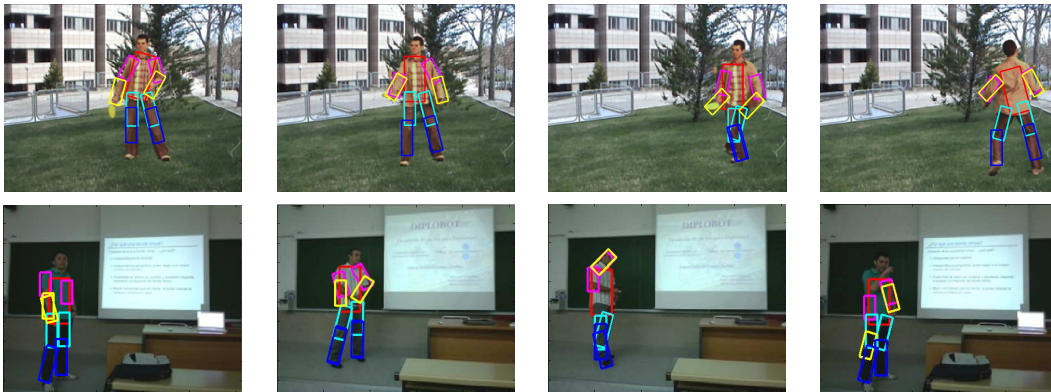


Figure 8: Pose recovery results in cVSG sequence.

### 3.3 Body limbs recovery

Finally, we combine the previous segmentation and face fitting with a full body pose recovery (Ramanan (2006)). In order to show the benefit of applying previous ST-GrabCut segmentation, we perform the overlapping performance of full pose recovery with and without human segmentation, always within the bounding box obtained from HOG person detection. Results are shown in Table 4. One can see that pose recovery considerably increases its performance when reducing the region of search based on ST-GrabCut segmentation. Some examples of pose recovery within the human segmentation regions for cVSG corpus and UBdataset are shown in Figure 8. One can see that in most of the cases body limbs are correctly detected. Only in some situations, occlusions or changes in body appearance can produce a wrong limb fitting.

In Figure 9 we show the application of the whole framework to perform temporal tracking, segmentation and full face and pose recovery. The colors correspond to the body limbs. The colors

Approach	Mean overlapping
Mesh fitting without segmentation	0.8960
ST-Grabcut & Temporal mesh fitting	0.9636

Table 2: AAM mesh fitting on original images and segmented images of the cVSG corpus.

Face view	Percentage of frames
Left view	0.1300
Near Left view	0.1470
Frontal view	0.2940
Near Right view	0.1650
Right view	0.2340

Table 3: Face pose percentages on the cVSG corpus.

Approach	Mean overlapping
Limb recovery without segmentation	0.7919
ST-Grabcut & Limb recovery	0.8760

Table 4: Overlapping of body limbs based on ground truth masks.

increase in intensity based on the instant of time of its detection. One can see the robust detection and temporal coherence based on the smooth displacement of face and limb detections.



Figure 9: Application of the whole framework (pose and face recovery) on an image sequence.

### 3.4 Human Limb data set

In this last experiment, we test our methodology on the presented Human Limb data set. From the 14 total limb annotations, we grouped them into six categories: trunk, up-arms, up-legs, low-arms, low-legs, and head, and we tested the full pose recovery framework. In this case, we tested the body limb recovery with and without applying the ST-GrabCut segmentation, and computed three different overlapping measures: %, that corresponds to the overlapping percentage defined in

		Trunk	Up-arms	Up-legs	Low-arms	Low-legs	Head	Mean
%	No segmentation	0.58	0.53	0.59	0.50	0.48	0.67	0.56
	STGrabCut*	0.58	0.53	0.58	0.50	0.56	0.67	<b>0.57</b>
Wins	No segmentation	106	104	108	109	68	120	102.5
	STGrabCut*	121	123	119	118	159	107	<b>124.5</b>
Match	No segmentation	133	127	130	121	108	155	129
	STGrabCut*	125	125	128	117	126	157	<b>129.66</b>

Table 5: Overlapping percentages between body parts -Intersection over Union-, wins -comparing the highest overlappings with and without segmentation-, and matchings -considering only overlappings greater than 0.6-.

\* STGrabCut was used without taking into account temporal information.

eq.(13), Wins, that corresponds to the number of Limb regions with higher overlapping comparing both strategies, and Match, that corresponds to the number of limb recoveries with overlapping superior to 0.6. The results are shown in Table 5. One can see that because of the reduced region where the subjects appear, in most cases there is no significant difference applying the limb recovery procedure with or without previous segmentation. Moreover, the segmentation algorithm is not working at maximum performance due to the same reason, since very small background regions are present in the images, and thus the background color model is quite poor. On the other hand, looking at the mean average overlapping in the last column of the table, one can see that ST-GrabCut improves for all overlapping measures the final limb overlapping. In particular, in the case of the Low-legs recovery is when a more clear improvement appears using ST-GrabCut segmentation. The part of the image corresponding to Low-legs is where more background influence exists, and thus the limb recovery has the highest confusion. However, as ST-GrabCut is able to properly segment the concave regions of the Low-legs regions, a significant improvement is obtained when applying the limb recovery methodology. Some results are illustrated on the images of Figure 10, where the images on the bottom correspond to the improvements obtained using the ST-GrabCut algorithm. Finally, Figure 11 show examples of the face fitting methodology applied on the human body limb data set.

#### 4. Conclusion

In this paper, we presented an evolution of the semi-automatic GrabCut algorithm for dealing with the problem of human segmentation in image sequences. The new full-automatic ST-GrabCut algorithm uses a HOG-based person detector, face detection, and skin color model to initialize GrabCut seeds. Spatial coherence is introduced via Mean Shift clustering, and temporal coherence is considered based on the historical of Gaussian Mixture Models. The segmentation procedure is combined with Shape and Active Appearance models to perform full face and pose recovery.

One of the problems we tackled with higher influence in the final results was the GrabCut bias towards convex segmentations. Since human silhouettes have many concavities, the segmentations obtained with GrabCut contain many errors. With our probabilistic reassignment of pixels based on color information from a set of frames we are able to correct large false positive regions in the segmentation, which results in a great improvement as showed by quantitative results.

Furthermore, we have demonstrated that pose recovery is enhanced by a prior segmentation of the subject, due to the large amount of background noise removed by the segmentation algorithm.



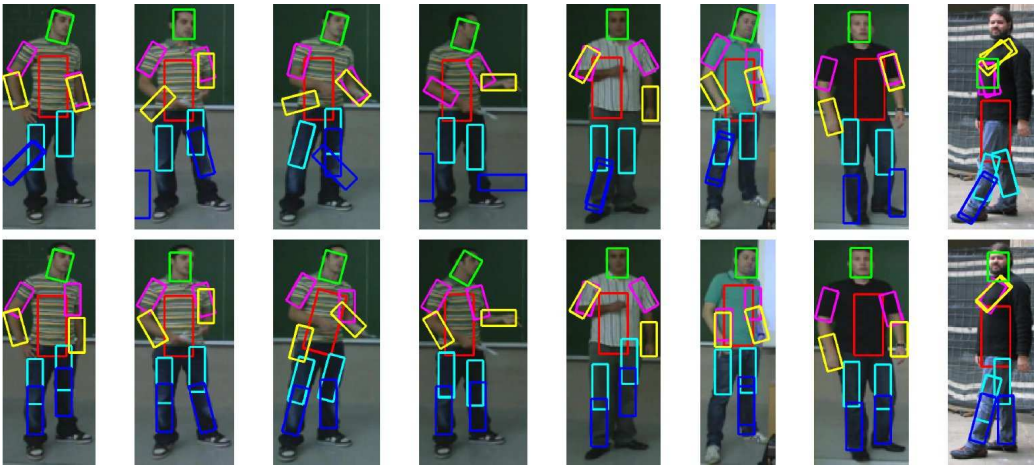


Figure 10: Human Limb data set results. Up row: limb recovery without ST-GrabCut segmentation. Down row: limb recovery with ST-GrabCut segmentation.



Figure 11: Application of face recovery on human body limb data set.

This general and full-automatic human segmentation, pose recovery, and tracking methodology showed higher performance than classical approaches in public image sequences and a novel Human Limb data set from uncontrolled environments, which makes it useful for general human face and gesture analysis applications.

One of the limitations of the method is that it depends on the initialization of the ST-GrabCut algorithm. Moreover, due to its sequential application, false seed labeling can accumulate segmentation errors along the video sequence.

As future work, we want to include the spatio-temporal coherence of the segmentations inside the graph-cuts framework. More specifically, we plan to modify not only the potentials of the energy function -and adding new ones-, but also the graph topology. Instead of just segmenting one frame at a time, graph-cuts can be applied on  $N$ -dimensional volumes, allowing us to segment multiple frames at a time by building a 3-dimensional graph. Additionally, the pixel connections neighbourhood can also be changed in order to increase or decrease the influence area of one pixel label to the surrounding ones. Furthermore, we plan to add shape constraints on the segmentation method as in the recent published works Varun Gulshan and Zisserman (2010); Brian L. Price (2010), trying to obtain more accurate segmentations.

Another line of research would be to feedback the segmentation methodology with the pose information of the subject. It seems probable that this iterative bidirectional information flow would improve both segmentation and pose recovery results. Moreover, this feedback from the pose information could be used to perform a multi-label segmentation of the image. This way, instead of just obtaining a binary segmentation -background, foreground-, we could be able to segment each body part separately.

Finally, we also plan to extend the limb recovery approach so that more complex poses and gestures can be recognized, and feed a gesture recognition system (Alon et al. (2009)) with the temporal aggregation of the recovered poses along the sequence in order to look for motion patterns of the limbs. With this information, we would be able to perform action recognition as well as human behaviour analysis, which can be very useful for psychological studies.

## Acknowledgments

I would like to thank Miguel Reyes for his work in the creation of the face meshes, and specially for his extremely nice and hard-working collaboration in the conference paper we got accepted in the 23rd IEEE Conference on Computer Vision and Pattern Recognition, as well as in the journal version of the paper.

I would also like to thank Victor Ponce for the creation of the Human Body Limb Dataset which resulted really useful for our work.

I would like to express my gratitude to all my colleagues at the Computer Vision Center, technical and administration staff, my supervisors S. Escalera and P. Radeva who guided and advised my work intensively, and in general to all the people who helped me. Special thanks to E. Artola, Ll. P. de las Heras, J. Almazan, D. Fernandez, M. Piñol and C. Davesa for their encouragements towards me and the nice moments we shared along the master.

My most special thanks to my family and friends, specially to my parents who always did their best for me.

This work has been supported in part by projects TIN2009-14404-C02 and CONSOLIDER-INGENIO CSD 2007-00018.

## References

- Jonathan Alon, Vassilis Athitsos, Quan Yuan, and Stan Sclaroff. A unified framework for gesture recognition and spatiotemporal gesture segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(9):1685–1699, 2009. ISSN 0162-8828. doi: <http://dx.doi.org/10.1109/TPAMI.2008.203>.
- Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. People-tracking-by-detection and people-detection-by-tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR'08)*, 2008. URL <http://www.mis.tu-darmstadt.de/node/382>.
- Yuri Boykov and Gareth Funka-Lea. Graph cuts and efficient n-d image segmentation. *Int. J. Comput. Vision*, 70(2):109–131, 2006. ISSN 0920-5691. doi: <http://dx.doi.org/10.1007/s11263-006-7934-5>.
- Yuri Boykov and Vladimir Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26:359–374, 2001.
- Yuri Y. Boykov and Marie-Pierre Jolly. Interactive graph cuts for optimal boundary region segmentation of objects in n-d images, 2001.
- Scott Cohen Brian L. Price, Bryan Morse. Geodesic graph cut for interactive image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- Daniel Chen, Brenden Chen, George Mamic, Clinton Fookes, and Sridha Sridharan. Improved grabcut segmentation via gmm optimisation. In *DICTA '08: Proceedings of the 2008 Digital Image Computing: Techniques and Applications*, pages 39–45, Washington, DC, USA, 2008. IEEE Computer Society. ISBN 978-0-7695-3456-5. doi: <http://dx.doi.org/10.1109/DICTA.2008.68>.
- T.F. Cootes, J. Edwards, and C. Taylor. Active appearance models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6):681–685, a.
- T.F. Cootes, C.J. Taylor, D. Cooper, and J. Graham. Active shape models - their training and application. *Computer Vision and Image Understanding*, 61(1):38–59, b.
- D. Corrigan, S. Robinson, and A. Kokaram. Video matting using motion extended grabcut. *IET Conference Publications*, 2008. doi: 10.1049/cp:20081076.
- N. Dalal and B. Triggs. Histogram of oriented gradients for human detection. volume 2, pages 886–893, 2005.
- P. Felzenszwalb and D. Huttenlocher. Pictorial structures for object recognition. *International Journal of Computer Vision*, 61(1).
- V Ferrari, M. Marin-Jimenez, and A. Zisserman. Progressive search space reduction for human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2008.

- V Ferrari, M. Marin-Jimenez, and A. Zisserman. Pose search: retrieving people using their pose. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- D. Geronimo, A. Lopez, and A. Sappa. Survey of pedestrian detection for advanced driver assistance systems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2009.
- Gary B. Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. (07-49), October 2007.
- Vladimir Kolmogorov and Ramin Zabih. What energy functions can be minimized via graph cuts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26:65–81, 2004.
- Bogdan Kwalek. Object segmentation in video via graph cut built on superpixels. *Fundam. Inf.*, 90(4):379–393, 2009. ISSN 0169-2968.
- Herve Lombaert, Yiyong Sun, Leo Grady, and Chenyang Xu. A multilevel banded graph cuts method for fast image segmentation. In *In ICCV05*, pages 259–265, 2005.
- Tomoyuki Nagahashi, Hironobu Fujiyoshi, and Takeo Kanade. Image segmentation using iterated graph cuts based on multi-scale smoothing. In *ACCV'07: Proceedings of the 8th Asian conference on Computer vision*, pages 806–816, Berlin, Heidelberg, 2007. Springer-Verlag. ISBN 3-540-76389-9, 978-3-540-76389-5.
- Deva Ramanan. Learning to parse images of articulated bodies. In *NIPS*, 2006.
- Remi Ronfard, Cordelia Schmid, and Bill Triggs. Learning to parse pictures of people. In *In European Conference on Computer Vision*, pages 700–714, 2002.
- C. Rother, V. Kolmogorov, and A. Blake. Grabcut: interactive foreground extraction using iterated graph cuts. 23(3):309–314, 2004.
- Fabrizio Tiburzi, Marcos Escudero, Jesus Bescos, and J.M. Martinez. A ground-truth for motion-based video-object segmentation: <http://www-gti.ii.uam.es/cvsg>. *IEEE International Conference on Image Processing (Workshop on Multimedia Information Retrieval)*, 2008.
- Antonio Criminisi Andrew Blake Varun Gulshan, Carsten Rother and Andrew Zisserman. Geodesic star convexity for interactive image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- Paul Viola and Michael J. Jones. Robust real-time face detection. volume 57, pages 137–154, 2004.