

Level-Set Person Segmentation and Tracking with Multi-Region Appearance Models and Top-Down Shape Information *Appendix*

Esther Horbert, Konstantinos Rematas, Bastian Leibe

UMIC Research Centre, RWTH Aachen University

{horbert, leibe}@umic.rwth-aachen.de

This appendix contains a detailed derivation of our segmentation and tracking framework [27] and lists typical values for important parameters.

A. Derivation

Fig. 1 shows the results of our full model for the sequence WALKSTRAIGHT [26], illustrating the evolution of the three contours.



Figure 1. Results of our full model for the sequence WALKSTRAIGHT (125 frames, from [26]). This example nicely illustrates the evolution of the separating contours (dark red). The foreground contour Φ_f is initialized with a horizontal line at 50% of the object frame height, the background contour Φ_b with a line at 60%. However this is only the initialization and the lines can evolve to very different forms (as the person's contour does).

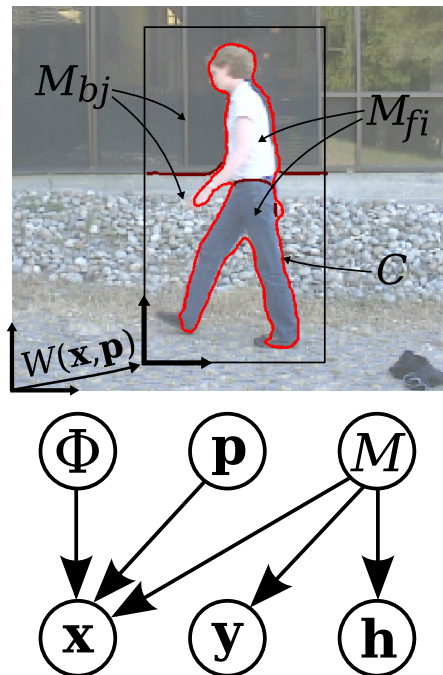


Figure 2. The generative model used in this approach, $M = \{M_{f1}, M_{f2}, M_{b1}, M_{b2}\}$.

\mathbf{x}	Pixel's coordinates inside reference frame
\mathbf{y}	Pixel's color
\mathbf{p}	Reference frame position
\mathbf{h}	Shape model
$W(\mathbf{x}, \mathbf{p})$	Warp with parameters \mathbf{p}
M_{f1}, M_{f2}	Foreground regions
M_{b1}, M_{b2}	Background regions
$P(\mathbf{y} M_k)$	Appearance models
Φ	Level set embedding function
$\{\Phi_c, \Phi_f, \Phi_b\}$	Embeddings for person and fore/background
C_k	Contour represented by the zero level set
$H_\epsilon(z)$	Smoothed Heaviside step function
$\delta_\epsilon(z)$	Smoothed Dirac delta function

Table 1. Notation used in this paper

The joint distribution for one pixel given by the model in Fig. 2 is:

$$P(\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i, \Phi, \mathbf{p}, M) = P(\mathbf{x}_i | \Phi, \mathbf{p}, M) P(\mathbf{y}_i | M) P(\mathbf{h}_i | M) P(M) P(\Phi) P(\mathbf{p}) \quad (1)$$

$$\text{(assumption: } \mathbf{y}, \mathbf{h} \text{ independent)} \quad (2)$$

$$P(\mathbf{x}_i, \Phi, \mathbf{p}, M | \mathbf{y}_i, \mathbf{h}_i) P(\mathbf{y}_i) P(\mathbf{h}_i) = P(\mathbf{x}_i | \Phi, \mathbf{p}, M) P(\mathbf{y}_i | M) P(\mathbf{h}_i | M) P(M) P(\Phi) P(\mathbf{p}) \quad (3)$$

$$P(\mathbf{x}_i, \Phi, \mathbf{p}, M | \mathbf{y}_i, \mathbf{h}_i) = P(\mathbf{x}_i | \Phi, \mathbf{p}, M) P(M | \mathbf{y}_i) P(M | \mathbf{h}_i) P(\Phi) P(\mathbf{p}) \quad (4)$$

Marginalization over the models M yields the pixel-wise posterior probability of shape Φ and location \mathbf{p} given a pixel $\{\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i\}$:

$$\sum_{k \in \{f1, f2, b1, b2\}} P(\mathbf{x}_i, \Phi, \mathbf{p}, M_k | \mathbf{y}_i, \mathbf{h}_i) = P(\mathbf{x}_i, \Phi, \mathbf{p} | \mathbf{y}_i, \mathbf{h}_i) = P(\Phi, \mathbf{p} | \mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i) P(\mathbf{x}_i) \quad (5)$$

$$= \sum_{k \in \{f1, f2, b1, b2\}} \left\{ P(\mathbf{x}_i | \Phi, \mathbf{p}, M_k) \frac{P(\mathbf{y}_i | M_k) P(M_k)}{\sum_l P(\mathbf{y}_i | M_l) P(M_l)} P(M_k | \mathbf{h}_i) \right\} P(\Phi) P(\mathbf{p}) \quad (6)$$

It follows

$$P(\Phi, \mathbf{p} | \mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i) = \frac{1}{P(\mathbf{x}_i)} \sum_{k \in \{f1, f2, b1, b2\}} \left\{ P(\mathbf{x}_i | \Phi, \mathbf{p}, M_k) \frac{P(\mathbf{y}_i | M_k) P(M_k)}{\sum_l P(\mathbf{y}_i | M_l) P(M_l)} P(M_k | \mathbf{h}_i) \right\} P(\Phi) P(\mathbf{p}) \quad (7)$$

We use a smoothed Heaviside step function H_ϵ to select the respective regions and a smoothed Dirac delta function δ_ϵ to select the contours:

$$H_c = H_\epsilon(\Phi_c(\mathbf{x}_i)), \quad \tilde{H}_c = 1 - H_\epsilon(\Phi_c(\mathbf{x}_i)) \quad (8)$$

$$H_f = H_\epsilon(\Phi_f(\mathbf{x}_i)), \quad \tilde{H}_f = 1 - H_\epsilon(\Phi_f(\mathbf{x}_i)) \quad (9)$$

$$H_b = H_\epsilon(\Phi_b(\mathbf{x}_i)), \quad \tilde{H}_b = 1 - H_\epsilon(\Phi_b(\mathbf{x}_i)) \quad (10)$$

$$P(M_k) = \frac{\eta_k}{\eta}, k \in \{f1, f2, b1, b2\}, \quad M_f = M_{f1} \cup M_{f2}, M_b = M_{b1} \cup M_{b2} \quad (11)$$

The number of pixels in the four respective regions can be obtained as follows:

$$N = \sum \eta_k, \quad \eta_{f1} = \sum_{i=1}^N H_c H_f, \quad \eta_{f2} = \sum_{i=1}^N H_c \tilde{H}_f, \quad (12)$$

$$\eta_{b1} = \sum_{i=1}^N \tilde{H}_c H_b, \quad \eta_{b2} = \sum_{i=1}^N \tilde{H}_c \tilde{H}_b \quad (13)$$

and thus the probability of pixel x_i for each region:

$$P(\mathbf{x}_i|\Phi, \mathbf{p}, M_{f1}) = \frac{H_c H_f}{\eta_{f1}}, \quad P(\mathbf{x}_i|\Phi, \mathbf{p}, M_{f2}) = \frac{H_c \tilde{H}_f}{\eta_{f2}}, \quad (14)$$

$$P(\mathbf{x}_i|\Phi, \mathbf{p}, M_{b1}) = \frac{\tilde{H}_c H_b}{\eta_{b1}}, \quad P(\mathbf{x}_i|\Phi, \mathbf{p}, M_{b2}) = \frac{\tilde{H}_c \tilde{H}_b}{\eta_{b2}}. \quad (15)$$

Fusing the pixel-wise posteriors with a logarithmic opinion pool yields

$$P(\Phi, \mathbf{p}|\mathbf{x}, \mathbf{y}, \mathbf{h}) \quad (16)$$

$$= \prod_{i=1}^N \sum_{k \in \{f1, f2, b1, b2\}} \left[P(\mathbf{x}_i|\Phi, \mathbf{p}, M_k) \frac{P(\mathbf{y}_i|M_k)P(M_k)}{\sum_l P(\mathbf{y}_i|M_l)P(M_l)} P(M_k|\mathbf{h}_i) \right] P(\Phi)P(\mathbf{p}) \quad (17)$$

$$= \prod_{i=1}^N \left[\frac{H_c H_f}{\eta_{f1}} \frac{P(\mathbf{y}_i|M_{f1})P(M_{f1})}{\sum_l P(\mathbf{y}_i|M_l)P(M_l)} P(M_{f1}|\mathbf{h}_i) + \frac{H_c \tilde{H}_f}{\eta_{f2}} \frac{P(\mathbf{y}_i|M_{f2})P(M_{f2})}{\sum_l P(\mathbf{y}_i|M_l)P(M_l)} P(M_{f2}|\mathbf{h}_i) \right. \quad (18)$$

$$\left. + \frac{\tilde{H}_c H_b}{\eta_{b1}} \frac{P(\mathbf{y}_i|M_{b1})P(M_{b1})}{\sum_l P(\mathbf{y}_i|M_l)P(M_l)} P(M_{b1}|\mathbf{h}_i) + \frac{\tilde{H}_c \tilde{H}_b}{\eta_{b2}} \frac{P(\mathbf{y}_i|M_{b2})P(M_{b2})}{\sum_l P(\mathbf{y}_i|M_l)P(M_l)} P(M_{b2}|\mathbf{h}_i) \right] P(\Phi)P(\mathbf{p}) \quad (19)$$

$$= \prod_{i=1}^N \left[H_c H_f P_{f1} + H_c \tilde{H}_f P_{f2} + \tilde{H}_c H_b P_{b1} + \tilde{H}_c \tilde{H}_b P_{b2} \right] P(\Phi)P(\mathbf{p}) \quad (20)$$

$$= \prod_{i=1}^N P(\mathbf{x}_i|\Phi, \mathbf{p}, \mathbf{y}_i, \mathbf{h}_i) P(\Phi)P(\mathbf{p}) \quad (21)$$

where

$$P_k = \frac{P(\mathbf{y}_i|M_k)P(M_k)P(M_k|\mathbf{h}_i)}{\eta_k \sum_l P(\mathbf{y}_i|M_l)P(M_l)} = \frac{P(\mathbf{y}_i|M_k)P(M_k|\mathbf{h}_i)}{\sum_l \eta_l P(\mathbf{y}_i|M_l)}, \quad k \in \{f1, f2, b1, b2\} \quad (22)$$

$$P(\mathbf{x}_i|\Phi, \mathbf{p}, \mathbf{y}_i, \mathbf{h}_i) = H_c H_f P_{f1} + H_c \tilde{H}_f P_{f2} + \tilde{H}_c H_b P_{b1} + \tilde{H}_c \tilde{H}_b P_{b2}. \quad (23)$$

$P(\mathbf{y}_i|M_k)$ is computed from the appearance models, *i.e.* color histograms.

Eq. (21) contains $P(M_k|h)$ for four regions, but the detector only provides probabilities for two regions: foreground and background. However, (21) selects a model for each region by use of H , H_f and H_b respectively. We set

$$P(M_f|h) = P(M_{f1}|h) + P(M_{f2}|h) = H_f P(M_f|h) + \tilde{H}_f P(M_f|h) \quad (24)$$

$$P(M_b|h) = P(M_{b1}|h) + P(M_{b2}|h) = H_b P(M_b|h) + \tilde{H}_b P(M_b|h). \quad (25)$$

Thus it holds:

$$\text{Either } P(M_f|h) = P(M_{f1}|h) \text{ or } P(M_f|h) = P(M_{f2}|h) \text{ and} \quad (26)$$

$$\text{either } P(M_b|h) = P(M_{b1}|h) \text{ or } P(M_b|h) = P(M_{b2}|h). \quad (27)$$

This means in practice $P(M_{f1}|h)$ and $P(M_{f2}|h)$ are both set to $P(M_f|h)$ (background accordingly), which is possible because for each pixel one of the subregions is selected.

We now specify $P(\Phi)$ as the internal energy of the level set embedding function(s). It contains a geometric prior that rewards a signed distance function: the gradient of the level set function is a normal distribution with mean 1. It thus makes the level set embedding function numerically stable without the need for periodic re-initializations [20]. The second term (as in [10]) describes the length of the contour and rewards a smoother contour. This is very useful for cluttered scenes, where pixels with foreground or background appearance can form very small regions, which can easily result in a very uneven contour.

$$P(\Phi) = \prod_{i=1}^N \left[\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(|\nabla\Phi(\mathbf{x}_i)| - 1)^2}{2\sigma^2}\right) \exp\left(-\lambda|\nabla H_\epsilon(\Phi)|\right) \right] \quad (28)$$

where σ and λ are the weights of the priors.

Maximizing the posterior is equivalent to minimizing its negative logarithm:

$$\begin{aligned} \mathcal{E}(\Phi) &= -\log(P(\Phi, \mathbf{p} | \mathbf{x}, \mathbf{y}, \mathbf{h})) \\ &\propto -\left(\sum_{i=1}^N \left\{ \log(P(\mathbf{x}_i | \Phi, \mathbf{p}, \mathbf{y}_i, \mathbf{h}_i)) - \frac{(|\nabla \Phi| - 1)^2}{2\sigma^2} - \lambda |\nabla H_\epsilon(\Phi)| \right\} + N \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \log(P(\mathbf{p})) \right) \end{aligned} \quad (29)$$

A.1. Derivation of the Segmentation Framework

For segmentation we optimize (29) w.r.t Φ , so the last two terms can be dropped, the rest is then differentiated by calculus of variation:

$$\begin{aligned} \frac{\partial \Phi_c}{\partial t} &= -\frac{\partial \mathcal{E}(\Phi_c)}{\partial \Phi_c} = \frac{\delta H_f P_{f1} + \delta \tilde{H}_f P_{f2} - \delta H_b P_{b1} - \delta \tilde{H}_b P_{b2}}{P(\mathbf{x} | \Phi, \mathbf{p}, \mathbf{y}, \mathbf{h})} - \frac{1}{\sigma^2} \left[\nabla^2(\Phi_c) - \text{div} \left(\frac{\nabla \Phi_c}{|\nabla \Phi_c|} \right) \right] \\ &\quad + \lambda \delta_\epsilon(\Phi_c) \text{div} \left(\frac{\nabla \Phi_c}{|\nabla \Phi_c|} \right) \end{aligned} \quad (30)$$

$$\begin{aligned} \frac{\partial \Phi_f}{\partial t} &= -\frac{\partial \mathcal{E}(\Phi_f)}{\partial \Phi_f} = \frac{H_c \delta_f P_{f1} - H_c \delta_f P_{f2} + \tilde{H}_c H_b P_{b1} + \tilde{H}_c \tilde{H}_b P_{b2}}{P(\mathbf{x} | \Phi, \mathbf{p}, \mathbf{y}, \mathbf{h})} - \frac{1}{\sigma^2} \left[\nabla^2(\Phi_f) - \text{div} \left(\frac{\nabla \Phi_f}{|\nabla \Phi_f|} \right) \right] \\ &\quad + \lambda \delta_\epsilon(\Phi_f) \text{div} \left(\frac{\nabla \Phi_f}{|\nabla \Phi_f|} \right) \end{aligned} \quad (31)$$

$$\approx_{(\Phi_f: H > 0)} \frac{\delta_f (P_{f1} - P_{f2})}{P(\mathbf{x} | \Phi, \mathbf{p}, \mathbf{y}, \mathbf{h})} - \frac{1}{\sigma^2} \left[\nabla^2(\Phi_f) - \text{div} \left(\frac{\nabla \Phi_f}{|\nabla \Phi_f|} \right) \right] + \lambda \delta_\epsilon(\Phi_f) \text{div} \left(\frac{\nabla \Phi_f}{|\nabla \Phi_f|} \right) \quad (32)$$

$$\frac{\partial \Phi_b}{\partial t} = -\frac{\partial \mathcal{E}(\Phi_b)}{\partial \Phi_b} \quad \text{accordingly.}$$

We evolve the two additional level set functions interleaved with the original level set function Φ_c . In this way, the four appearance models are optimized at the same time, which leads to more robust and accurate segmentation results.

In our implementation we use $\sigma^2 = 50$, $\lambda = 2$, $\tau = 2$, $\epsilon = 6$ and

$$H_\epsilon(x) = \begin{cases} 0 & \text{if } x < -\epsilon \\ \frac{x}{2\epsilon} + \frac{1}{2\pi} \sin\left(\frac{\pi x}{\epsilon}\right) + \frac{1}{2} & \text{if } |x| < \epsilon \\ 1 & \text{if } x > \epsilon \end{cases} \quad (33)$$

$$\delta_\epsilon(x) = \begin{cases} \frac{1}{2\epsilon} (1 + \cos\left(\frac{\pi x}{\epsilon}\right)) & \text{if } |x| < \epsilon \\ 0 & \text{else} \end{cases} \quad (34)$$

A.2. Derivation of the Tracking Framework

In preparation for differentiation w.r.t \mathbf{p} some terms in (29) can be dropped:

$$\mathcal{E}(\Phi) \propto -\left(\sum_{i=1}^N \{ \log(P(\mathbf{x}_i | \Phi, \mathbf{p}, \mathbf{y}_i, \mathbf{h}_i)) \} + \log(P(\mathbf{p})) + \text{const.} \right) \quad (35)$$

Now the warp $\mathbf{W}(\mathbf{x}_i, \Delta \mathbf{p})$ is introduced into (35), *i.e.* pixels \mathbf{x}_i are warped with parameters \mathbf{p} . $P(\mathbf{p})$ is dropped for the moment, this is handled with drift correction, as in [3]:

$$\mathcal{E}(\Phi) \propto -\sum_{i=1}^N \log \left\{ P(\mathbf{W}(\mathbf{x}_i, \Delta \mathbf{p}) | \Phi, \mathbf{p}, \mathbf{y}_i, \mathbf{h}_i) \right\} \quad (36)$$

We maximize w.r.t \mathbf{p} :

$$\mathbf{p} = \arg \max_{\mathbf{p}} \left\{ \sum_{i=1}^N \log P(\mathbf{W}(\mathbf{x}_i, \Delta \mathbf{p}) | \Phi, \mathbf{p}, \mathbf{y}_i, \mathbf{h}_i) \right\} \quad (37)$$

We use a second order Newton optimization scheme as in [4]: With the short-hand notation $P(\dots) = P(\mathbf{W}(\mathbf{x}_i, \Delta\mathbf{p})|\Phi, \mathbf{p}, \mathbf{y}_i, \mathbf{h}_i)$:

$$\Delta\mathbf{p} = \left[\sum_{i=1}^N \frac{\left(\frac{\partial P(\dots)}{\partial \mathbf{p}} \right)^2}{P(\dots)} \right]^{-1} \sum_{i=1}^N \frac{\partial P(\dots)}{\partial \mathbf{p}} \quad (38)$$

where

$$\begin{aligned} \frac{\partial P(\dots)}{\partial \mathbf{p}} &= (\mathbf{J}_c H_f + H_c \mathbf{J}_f) P_{f1} + (\mathbf{J}_c \tilde{H}_f - H_c \mathbf{J}_f) P_{f2} - \mathbf{J}_c H_b P_{b1} - \mathbf{J}_c \tilde{H}_b P_{b2} \\ &= \mathbf{J}_c (H_f P_{f1} + \tilde{H}_f P_{f2} - H_b P_{b1} - \tilde{H}_b P_{b2}) + \mathbf{J}_f (H_c P_{f1} - H_c P_{f2}) \end{aligned} \quad (39)$$

with

$$\mathbf{J}_c = \frac{\partial H_c}{\partial \Phi_c} \frac{\partial \Phi_c}{\partial \mathbf{x}} \frac{\partial \mathbf{W}}{\partial \Delta \mathbf{p}} = \delta_\epsilon(\Phi_c(\mathbf{x}_i)) \nabla \Phi_c(\mathbf{x}_i) \frac{\partial \mathbf{W}}{\partial \Delta \mathbf{p}}, \quad (40)$$

$$\mathbf{J}_f = \frac{\partial H_f}{\partial \Phi_f} \frac{\partial \Phi_f}{\partial \mathbf{x}} \frac{\partial \mathbf{W}}{\partial \Delta \mathbf{p}} = \delta_\epsilon(\Phi_f(\mathbf{x}_i)) \nabla \Phi_f(\mathbf{x}_i) \frac{\partial \mathbf{W}}{\partial \Delta \mathbf{p}} \quad (41)$$

We assume that the background does not move with the foreground, thus $H_\epsilon(\Phi_b(\mathbf{x}_i))$ is constant w.r.t the derivation $\frac{\partial}{\partial \mathbf{p}}$. Eq. (39) illustrates that both the person's contour and the division line of the foreground contribute to the warp, whereas the division line of the background does not contribute: \mathbf{J}_c and \mathbf{J}_f contain the factor $\delta_\epsilon(\Phi_k)$, the Dirac delta function of the respective level set function, which is only greater than zero in a narrow band around the contour (with width ϵ). For parameters $\mathbf{p} = (s, t_x, t_y)^T$ with scale and translation the warp is:

$$\mathbf{W}(\mathbf{x}; \mathbf{p}) = \begin{pmatrix} 1+s & 0 & t_x \\ 0 & 1+s & t_y \end{pmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{pmatrix} (1+s) \cdot x + t_x \\ (1+s) \cdot y + t_y \\ 1 \end{pmatrix} \quad (42)$$

The term $\frac{\partial \mathbf{W}}{\partial \mathbf{p}}$ is the Jacobian of the warp and with $\mathbf{W}(\mathbf{x}; \mathbf{p}) = (W_x(\mathbf{x}; \mathbf{p}), W_y(\mathbf{x}; \mathbf{p}))^T$:

$$\frac{\partial \mathbf{W}}{\partial \mathbf{p}} = \begin{pmatrix} x & 1 & 0 \\ y & 0 & 1 \end{pmatrix} \quad (43)$$

In our implementation the rigid registration typically converges after 10 to 40 iterations.

References

- [3] C. Bibby and I. Reid. Robust Real-Time Visual Tracking using Pixel-Wise Posteriors. *ECCV*, 2008.
- [4] C. Bibby and I. Reid. Real-time Tracking of Multiple Occluding Objects using Level Sets. *CVPR*, 2010.
- [10] D. Cremers, M. Rousson, and R. Deriche. A Review of Statistical Approaches to Level Set Segmentation Integrating Color, Texture, Motion and Shape. *IJCV*, 72:195-215, 2007.
- [20] C. Li, C. Xu, C. Gui, and M. Fox. Level Set Evolution without Re-initialization: A New Variational Formulation. *CVPR*, 2005.
- [26] H. Sidenbladh and M. J. Black. Learning the statistics of people in images and video. *IJCV*, 54(1-3):183-209, 2003.
- [27] E. Horbert, K. Rematas and B. Leibe. Level-Set Person Segmentation and Tracking with Multi-Region Appearance Models and Top-Down Shape Information. *ICCV*, 2011.