

RWTH AACHEN UNIVERSITY

Computer Vision II - Lecture 14

Articulated Tracking I

08.07.2014

Bastian Leibe
RWTH Aachen
<http://www.vision.rwth-aachen.de>
leibe@vision.rwth-aachen.de

Computer Vision II, Summer'14

RWTH AACHEN UNIVERSITY

Outline of This Lecture

- Single-Object Tracking
- Bayesian Filtering
 - Kalman Filters, EKF
 - Particle Filters
- Multi-Object Tracking
 - Data association
 - MHT, (JPDAF, MCMCDA)
 - Network flow optimization
- Articulated Tracking
 - GP body pose estimation
 - (Model-based tracking, AAMs)
 - Pictorial Structures

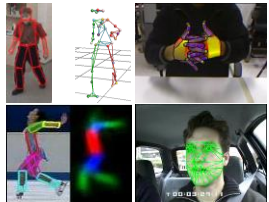


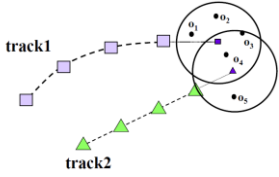
Image sources: Andreas Es, Deva Bhargava, Jan Matthies

Computer Vision II, Summer'14

RWTH AACHEN UNIVERSITY

Recap: Linear Assignment Formulation

- Form a matrix of pairwise similarity scores
- Example: Similarity based on motion prediction
 - Predict motion for each trajectory and assign scores for each measurement based on inverse (Mahalanobis) distance, such that closer measurements get higher scores.



	ai1	ai2
1	3.0	
2	5.0	
3	6.0	1.0
4	9.0	8.0
5		3.0

- Choose at most one match in each row and column to maximize sum of scores

Slide credit: Robert Collins. B. Leibe

3

RWTH AACHEN UNIVERSITY

Recap: Linear Assignment Problem

- Formal definition
 - Maximize $\sum_{i=1}^N \sum_{j=1}^M w_{ij} z_{ij}$

subject to $\sum_{j=1}^M z_{ij} = 1; i = 1, 2, \dots, N$
 $\sum_{i=1}^N z_{ij} = 1; j = 1, 2, \dots, M$
 $z_{ij} \in \{0, 1\}$

} These constraints ensure that Z is a permutation matrix

- The permutation matrix constraint ensures that we can only match up one object from each row and column.
- Note: Alternatively, we can minimize cost rather than maximizing weights.

$$\arg \min_{z_{ij}} \sum_{i=1}^N \sum_{j=1}^M c_{ij} z_{ij}$$

Slide adapted from Robert Collins. B. Leibe

4

RWTH AACHEN UNIVERSITY

Recap: Optimal Solution

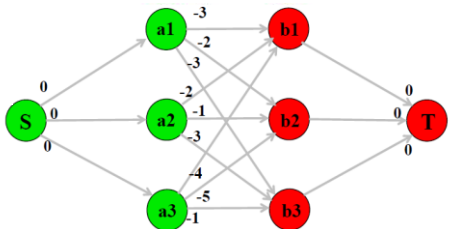
- Greedy Algorithm
 - Easy to program, quick to run, and yields "pretty good" solutions in practice.
 - But it often does not yield the optimal solution
- Hungarian Algorithm
 - There is an algorithm called Kuhn-Munkres or "Hungarian" algorithm specifically developed to efficiently solve the linear assignment problem.
 - Reduces assignment problem to bipartite graph matching.
 - When starting from an $N \times N$ matrix, it runs in $O(N^3)$.
 - ⇒ If you need LAP, you should use it.

Slide credit: Robert Collins. B. Leibe

5

RWTH AACHEN UNIVERSITY

Recap: Min-Cost Flow



- Conversion into flow graph
 - Transform weights into costs $c_{ij} \propto w_{ij}$
 - Add source/sink nodes with 0 cost.
 - Directed edges with a capacity of 1.

Slide credit: Robert Collins. B. Leibe

6

Computer Vision II, Summer'14 RWTH AACHEN UNIVERSITY

Recap: Min-Cost Flow

- Conversion into flow graph
 - Pump N units of flow from source to sink.
 - Internal nodes pass on flow ($\sum \text{flow in} = \sum \text{flow out}$).
- Find the optimal paths along which to ship the flow.

Slide credit: Robert Collins B. Leibe 7

Computer Vision II, Summer'14 RWTH AACHEN UNIVERSITY

Recap: Min-Cost Flow

- Conversion into flow graph
 - Pump N units of flow from source to sink.
 - Internal nodes pass on flow ($\sum \text{flow in} = \sum \text{flow out}$).
- Find the optimal paths along which to ship the flow.

Slide credit: Robert Collins B. Leibe 8

Computer Vision II, Summer'14 RWTH AACHEN UNIVERSITY

Recap: Using Network Flow for Tracking

- Complication 1
 - Tracks can start later than frame1 (and end earlier than frame4)
- Connect the source and sink nodes to all intermediate nodes.

Slide credit: Robert Collins B. Leibe 9

Computer Vision II, Summer'14 RWTH AACHEN UNIVERSITY

Recap: Using Network Flow for Tracking

- Complication 2
 - Trivial solution: zero cost flow!

Slide credit: Robert Collins B. Leibe 10

Computer Vision II, Summer'14 RWTH AACHEN UNIVERSITY

Recap: Network Flow Approach

Solution: Divide each detection into 2 nodes

Zhang, Li, Nevatia, [Global Data Association for Multi-Object Tracking using Network Flows](#), CVPR'08.

image source: [Zhang, Li, Nevatia, CVPR'08]

11

Computer Vision II, Summer'14 RWTH AACHEN UNIVERSITY

Recap: Min-Cost Formulation

- Objective Function

$$\mathcal{T}^* = \operatorname{argmin}_{\mathcal{T}} \sum_i C_{in,i} f_{in,i} + \sum_i C_{i,out} f_{i,out} + \sum_{i,j} C_{i,j} f_{i,j} + \sum_i C_i f_i$$
- subject to
 - Flow conservation at all nodes

$$f_{in,i} + \sum_j f_{j,i} = f_i = f_{out,i} + \sum_j f_{i,j} \quad \forall i$$
 - Edge capacities

$$f_i \leq 1$$

Slide credit: Laura Leal B. Leibe 12

RWTH AACHEN UNIVERSITY

Topics of This Lecture

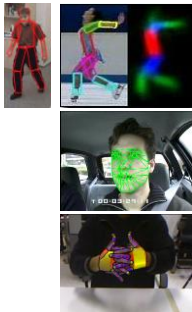
- **Articulated Tracking**
 - Motivation
 - Classes of Approaches
- **Body Pose Estimation as High-Dimensional Regression**
 - Representations
 - Training data generation
 - Latent variable space
 - Learning a mapping between pose and appearance
- **Review: Gaussian Processes**
 - Formulation
 - GP Prediction
 - Algorithm
- **Applications**
 - Articulated Tracking under Egomotion

13

RWTH AACHEN UNIVERSITY

Articulated Tracking

- **Examples**
 - Recover a person's body articulation
 - Track facial expressions
 - Track detailed hand motion
 - ...
- **Common properties**
 - Detailed parameterization in terms of joint locations or joint angles
 - Two steps
 - Pose estimation (in single frame)
 - Tracking (using dynamics model)
 - Challenging problem
 - High-dimensional
 - Hitting the limits of sensor data

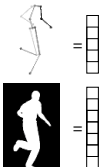
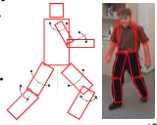


14
image sources: T. Svoboda, D. Ramanan, I. Matthews, J. Oikonomidis

RWTH AACHEN UNIVERSITY

Basic Classes of Approaches


- **Global methods**
 - Entire body configuration is treated as a point in some high-dimensional space.
 - Observations are also global feature vectors.
 - ⇒ View of pose estimation as a high-dimensional regression problem.
 - ⇒ Often in a subspace of "typical" motions...
- **Part-based methods**
 - Body configuration is modeled as an assembly of movable parts with kinematic constraints.
 - Local search for part configurations that provide a good explanation for the observed appearance under the kinematic constraints.
 - ⇒ View of pose estimation as probabilistic inference in a dynamic Graphical Model.

15
image sources: T. Jaegle, D. Ramanan, T. Svoboda

RWTH AACHEN UNIVERSITY

Why Is It Difficult?



- **Challenges**
 - Poor imaging, motion blur, occlusions, etc.
 - Difficult to extract sufficiently good figure-ground information
 - Mapping is generally multi-modal: an image observation can represent more than one pose!

16
Slide credit: Raquel Urtasun, B. Leibe

RWTH AACHEN UNIVERSITY

Topics of This Lecture

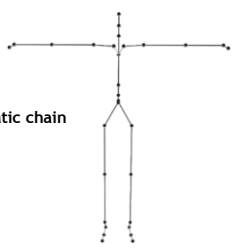
- **Articulated Tracking**
 - Motivation
 - Classes of Approaches
- **Body Pose Estimation as High-Dimensional Regression**
 - Representations
 - Training data generation
 - Latent variable space
 - Learning a mapping between pose and appearance
- **Review: Gaussian Processes**
 - Formulation
 - GP Prediction
 - Algorithm
- **Applications**
 - Articulated Tracking under Egomotion

17

RWTH AACHEN UNIVERSITY

Body Representation

- **The body can be approximated as kinematic tree**
- **Parametrization via**
 - Joint locations
 - Joint angles
 - Relative joint angles along kinematic chain
 - ...
- **Example using in the following**
 - 3D joint locations of 20 joints
 - ⇒ 60-dimensional space




18
Slide adapted from Raquel Urtasun, B. Leibe, image source: R. Urtasun

Computer Vision II, Summer'14

Image Representation

- Many possibilities...
- Popular choice: Silhouettes
 - Easy to extract using background modeling techniques.
 - Capture important information about body shape.

⇒ We will use them as an example for today's lecture...



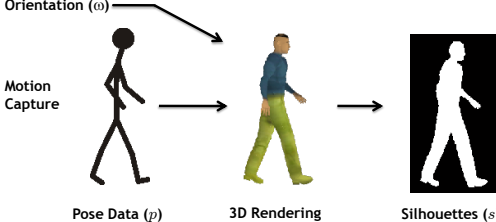
B. Leibe Video source: Hedvig Sidhenblad

19

Computer Vision II, Summer'14

Another Advantage of Silhouette Data

- Synthetic training data generation possible!
 - Create sequences of „Pose + Silhouette“ pairs
 - Poses recorded with Mocap, used to animate 3D model
 - Silhouette via 3D rendering pipeline



Orientation (θ)

Motion Capture

Pose Data (p)

3D Rendering

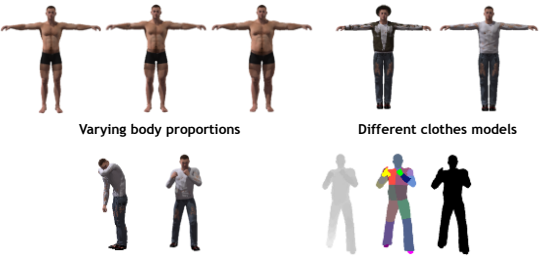
Silhouettes (s)

Slide adapted from Stefan Gammeter B. Leibe

20

Computer Vision II, Summer'14

Synthetic Training Data Generation



Varying body proportions

Different clothes models

Animate with MoCap data

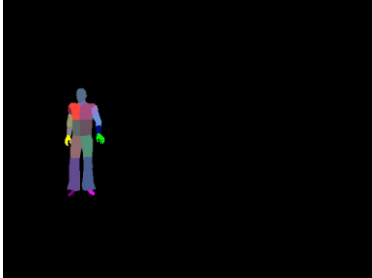
Resulting synthetic training data (depth, body part labels, silhouette)

Image source: Umer Raf

21

Computer Vision II, Summer'14

Synthetic Training Data Generation



Example training sequence

Video source: Umer Raf

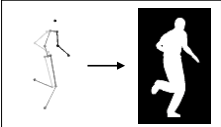
22

Computer Vision II, Summer'14

Learning a Mapping b/w Pose and Appearance


- Appearance prediction
 - Regression problem
 - High-dimensional data on both sides

⇒ Low-dim. representation needed for learning!



- 3D joint locations • segm. image
- ~60-dim. • ~2500-dim.

- Training with Motion-capture stimuli
 - Real dynamics from human actors
 - Synthesized silhouettes for training
 - Background subtraction for test

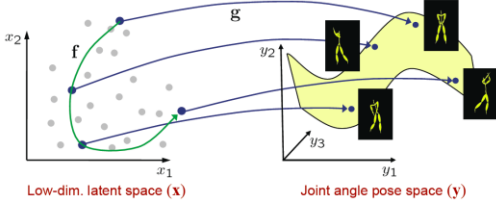


T. Jaeggli, E. Koller-Meier, L. Van Gool, "Learning Generative Models for Monocular Body Pose Estimation", ACCV 2007. image source: T. Jaeggli

23

Computer Vision II, Summer'14

Latent Variable Models



Low-dim. latent space (x)

Joint angle pose space (y)

- Joint angle pose space is huge!
 - Only a small portion contains valid body poses.

⇒ Restrict estimation to the subspace of valid poses for the task

- Latent variable models: PCA, FA, GPLVM, etc.

B. Leibe image source: R. Urtasun

24

Computer Vision II, Summer '14

Example: Subspace of Walking Motion

- Pose modeling in a subspace
 - Pose model has 60 (highly dependent) DoF
 - But gait is cyclic, can be represented by a 2D latent space
 - Capture the dependency by dimensionality reduction (PCA, FA, CCA, LLE, GPLVM, ...)

B. Leibe image sources: S. Gammeter, T. Jaeggli

Computer Vision II, Summer '14

Articulated Motion in the Latent Space

- Regression from latent space to
 - Pose $p(\text{pose} | \mathbf{z})$
 - Silhouette $p(\text{silhouette} | \mathbf{z})$
- Regressors need to be learned from training data.

Slide adapted from Stefan Gammeter B. Leibe

Computer Vision II, Summer '14

Learning a Generative Mapping

T. Jaeggli, E. Koller-Meier, L. Van Gool, "Learning Generative Models for Monocular Body Pose Estimation", ACCV 2007.

Slide credit: Tobias Jaeggli

Computer Vision II, Summer '14

Example Results

- Difficulties
 - Changing viewpoints
 - Low resolution (50 px)
 - Compression artifacts
 - Disturbing objects (umbrella, bag)

B. Leibe Video sources: Hedvig Sidenblad, Tobias Jaeggli

Computer Vision II, Summer '14

Representing Multiple Activities

- Learn multiple models
 - One model per activity
 - Separate LLE embedding
 - Separate dynamics
- Learn transition function
 - Link the LLE spaces
 - Find similar pose pairs
 - Learn smooth transition

$$p(x_t, a_t | x_{t-1}, a_{t-1}) \propto \begin{cases} p_{\text{noswitch}}^{a_t} & \text{if } a_t = a_{t-1} \\ p_{\text{switch}}^{a_t} & \text{else} \end{cases}$$

Slide credit: Tobias Jaeggli B. Leibe

Computer Vision II, Summer '14

Switching b/w Multiple Activities

- Activity switching
 - Low-res. traffic scene
 - Transition from Walking to Running

B. Leibe Videos by Tobias Jaeggli

Computer Vision II, Summer'14

Topics of This Lecture

- Articulated Tracking
 - Motivation
 - Classes of Approaches
- Body Pose Estimation as High-Dimensional Regression
 - Representations
 - Training data generation
 - Latent variable space
 - Learning a mapping between pose and appearance
- Review: Gaussian Processes
 - Formulation
 - GP Prediction
 - Algorithm
- Applications
 - Articulated Tracking under Egomotion

31

Computer Vision II, Summer'14

Classification vs. Regression

In classification: $y \in \{-1, 1\}$ In regression: $y \in \mathbb{R}$

Slide credit: Raquel Urtasun B. Leibe

32

Computer Vision II, Summer'14

Gaussian Process Regression

- “Regular” regression: $y = f(x)$

- GP regression: $p(y|x) \sim \mathcal{N}(\mu(x), \sigma(x))$

Slide credit: Stefan Gammeter B. Leibe

33

Computer Vision II, Summer'14

Gaussian Process Regression

- GP Regression
 - Very easy to apply
 - Automatic confidence estimate of the result
 - Well-suited for pose regression tasks
- In the following, I will give a quick intro to GPs
 - Focus on main concepts and results
 - A far more detailed discussion will be given in the Advanced Machine Learning lecture (next semester).

B. Leibe

34

Computer Vision II, Summer'14

Gaussian Process

- Gaussian distribution
 - Probability distribution over scalars / vectors.
- Gaussian process (generalization of Gaussian distrib.)
 - Describes properties of functions.
 - Function: Think of a function as a long vector where each entry specifies the function value $f(x_i)$ at a particular point x_i .
 - Issue: How to deal with infinite number of points?
 - If you ask only for properties of the function at a finite number of points...
 - Then inference in Gaussian Process gives you the same answer if you ignore the infinitely many other points.
- Definition
 - A Gaussian process (GP) is a collection of random variables any finite number of which has a joint Gaussian distribution.

Slide credit: Bernt Schiele B. Leibe

35

Computer Vision II, Summer'14

Gaussian Process

- Example prior over functions $p(f)$
 - Represents our prior belief about functions before seeing any data.
 - Although specific functions don't have mean of zero, the mean of $f(x)$ values for any fixed x is zero (here).
 - Favors smooth functions
 - I.e. functions cannot vary too rapidly
 - Smoothness is induced by the covariance function of the Gaussian Process.
 - Learning in Gaussian processes
 - Is mainly defined by finding suitable properties of the covariance function.

Slide credit: Bernt Schiele B. Leibe Image source: Rasmussen & Williams, 2006

36

Gaussian Process

- A Gaussian process is completely defined by

- Mean function $m(\mathbf{x})$ and

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$$

- Covariance function $k(\mathbf{x}, \mathbf{x}')$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

- We write the Gaussian process (GP)

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

Gaussian Process: Squared Exponential

- Typical covariance function

- Squared exponential (SE)

- Covariance function specifies the covariance between pairs of random variables

$$\text{cov}[f(\mathbf{x}_p), f(\mathbf{x}_q)] = k(\mathbf{x}_p, \mathbf{x}_q) = \exp\left\{-\frac{1}{2}|\mathbf{x}_p - \mathbf{x}_q|^2\right\}$$

- Remarks

- Covariance between the **outputs** is written as a function between the **inputs**.
- The squared exponential covariance function corresponds to a Bayesian linear regression model with an **infinite** number of basis functions.
- For any positive definite covariance function $k(\cdot, \cdot)$, there exists a (possibly infinite) expansion in terms of basis functions.

Gaussian Process: Prior over Functions

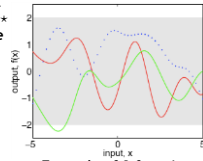
- Distribution over functions:

- Specification of covariance function implies distribution over functions.

- I.e. we can draw samples from the distribution of functions evaluated at a (finite) number of points.

- Procedure

- We choose a number of input points X_*
- We write the corresponding covariance matrix (e.g. using SE) element-wise:
 $K(X_*, X_*)$
- Then we generate a random Gaussian vector with this covariance matrix:
 $f_* \sim \mathcal{N}(0, K(X_*, X_*))$



Example of 3 functions sampled
39
Image source: Rasmussen & Williams, 2006

GP Prediction with Noisy Observations

- Assume we have a set of observations:

$$\{(\mathbf{x}_n, f_n) \mid n = 1, \dots, N\} \text{ with noise } \sigma_n$$

- Joint distribution of the training outputs \mathbf{f} and test outputs \mathbf{f}_* according to the prior:

$$\begin{bmatrix} \mathbf{t} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X, X_*) & K(X_*, X_*) \end{bmatrix}\right)$$

- $K(X, X_*)$ contains covariances for all pairs of training and test points.

- To get the **posterior** (after including the observations)

- We need to restrict the above prior to contain only those functions which agree with the observed values.
- Think of generating functions from the prior and rejecting those that disagree with the observations (obviously prohibitive).

Result: Prediction with Noisy Observations

- Calculation of posterior:

- Corresponds to **conditioning the joint Gaussian prior distribution on the observations**:

$$f_* | X_*, X, \mathbf{t} \sim \mathcal{N}(\bar{f}_*, \text{cov}[f_*]) \quad \bar{f}_* = \mathbb{E}[f_* | X, X_*, \mathbf{t}]$$

- with:

$$\bar{f}_* = K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} \mathbf{t}$$

$$\text{cov}[f_*] = K(X_*, X_*) - K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} K(X, X_*)$$

- This is the key result that defines Gaussian process regression!**

- The predictive distribution is a Gaussian whose mean and variance depend on the test points X_* and on the kernel $k(\mathbf{x}, \mathbf{x}')$, evaluated on the training data X .

GP Regression Algorithm

- Very simple algorithm

input: X (inputs), \mathbf{y} (targets), k (covariance function), σ_n^2 (noise level), \mathbf{x}_* (test input)

$$2: L := \text{cholesky}(K + \sigma_n^2 I)$$

$$\alpha := L^{-T} (L \mathbf{y})$$

$$4: \bar{f}_* := \mathbf{k}_*^T \alpha \quad \left. \vphantom{\bar{f}_*} \right\} \text{predictive mean eq. (2.25)}$$

$$\mathbf{v} := L \mathbf{k}_* \quad \left. \vphantom{\mathbf{v}} \right\} \text{predictive variance eq. (2.26)}$$

$$6: \mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v}$$

$$\log p(\mathbf{y}|X) := -\frac{1}{2} \mathbf{y}^T \alpha - \sum_i \log L_{ii} - \frac{n}{2} \log 2\pi \quad \text{eq. (2.30)}$$

$$8: \text{return: } \bar{f}_* \text{ (mean), } \mathbb{V}[f_*] \text{ (variance), } \log p(\mathbf{y}|X) \text{ (log marginal likelihood)}$$

- Based on the following equations (Matrix inv. \leftrightarrow Cholesky fact.)

$$\bar{f}_* = \mathbf{k}_*^T (K + \sigma_n^2 I)^{-1} \mathbf{t}$$

$$\text{cov}[f_*] = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K + \sigma_n^2 I)^{-1} \mathbf{k}_*$$

$$\log p(\mathbf{t}|X) = -\frac{1}{2} \mathbf{t}^T (K + \sigma_n^2 I)^{-1} \mathbf{t} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{N}{2} \log 2\pi$$

RWTH AACHEN UNIVERSITY

Computational Complexity

- Complexity of GP model
 - Training effort: $\mathcal{O}(N^3)$ through matrix inversion
 - Test effort: $\mathcal{O}(N^2)$ through vector-matrix multiplication
- Complexity of basis function model
 - Training effort: $\mathcal{O}(M^3)$
 - Test effort: $\mathcal{O}(M^2)$
- Discussion
 - Exact GP methods become infeasible for large training sets.
 - ⇒ Need to use approximate techniques whenever #training examples > 2500-3000.

43

Computer Vision II, Summer'14

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Articulated Tracking
 - Motivation
 - Classes of Approaches
- Body Pose Estimation as High-Dimensional Regression
 - Representations
 - Training data generation
 - Latent variable space
 - Learning a mapping between pose and appearance
- Review: Gaussian Processes
 - Formulation
 - GP Prediction
 - Algorithm
- Applications
 - Articulated Tracking under Egomotion

44

Computer Vision II, Summer'14

RWTH AACHEN UNIVERSITY

Articulated Multi-Person Tracking using GP

- Idea: Only perform articulated tracking where it's easy!
- Multi-person tracking
 - Solves hard data association problem
- Articulated tracking
 - Only on individual "tracklets" between occlusions
 - GP regression on full-body pose

45

Computer Vision II, Summer'14

[Gammeter, Ess, Jaegeli, Schindler, Leibe, Van Gool, ECCV'08]

RWTH AACHEN UNIVERSITY

Articulated Multi-Person Tracking

- Multi-Person tracking
 - Recovers trajectories and solves data association
- Articulated Tracking
 - Estimates detailed body pose for each tracked person

46

Computer Vision II, Summer'14

[Gammeter, Ess, Jaegeli, Schindler, Leibe, Van Gool, ECCV'08]

RWTH AACHEN UNIVERSITY

Articulated Tracking under Egomotion

- Guided segmentation for each frame
 - No reliance on background modeling
 - Approach applicable to scenarios with moving camera
 - Feedback from body pose estimate to improve segmentation

47

Computer Vision II, Summer'14

[Gammeter, Ess, Jaegeli, Schindler, Leibe, Van Gool, ECCV'08]

RWTH AACHEN UNIVERSITY

Summary: Articulated Tracking with Global Models

- Pros:
 - View as regression problem (pose ↔ appearance)
 - Lots of machine learning techniques available
 - Research focus on handling the ambiguities
 - Training on MoCap data possible
 - Accurate models for human dynamics
- Cons:
 - High-dimensional problem
 - Global model
 - Can handle only those articulations it has previously seen
 - Not robust against partial occlusion
 - Difficult to get good appearance representation
 - MoCap data ⇒ Can synthesize silhouettes, but not appearance
 - Restricted to background subtraction

48

Computer Vision II, Summer'14

B. Leibe