

# Seeing the Objects Behind the Dots: Recognition in Videos from a Moving Camera

---

Bjorn Ommer, Theodor Mader, Joachim M. Buhmann  
Present by: Alisa Kunapinun

# Abstract

---

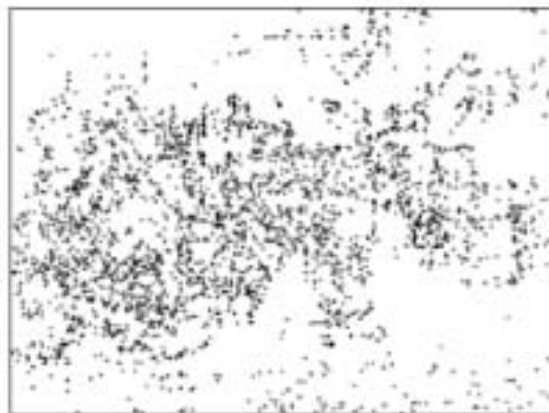
- This paper also significantly contributes to the systems design aspect—showing how all of these subtasks can be combined in a computer vision system so that they mutually benefit from another.

# Keywords

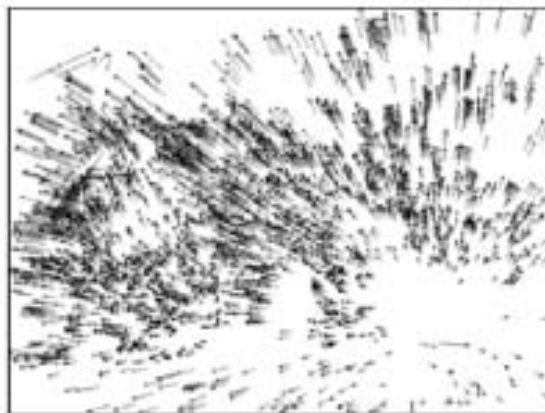
---

- Object Recognition
- Segmentation
- Tracking
- Videos Analysis
- Compositionality
- Visual Learning

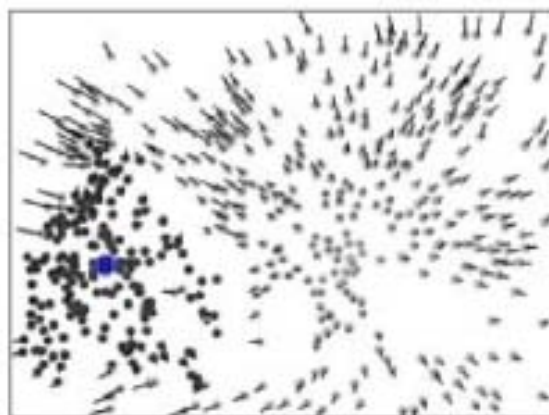
# Process on this paper



(a)



(b)



(c)

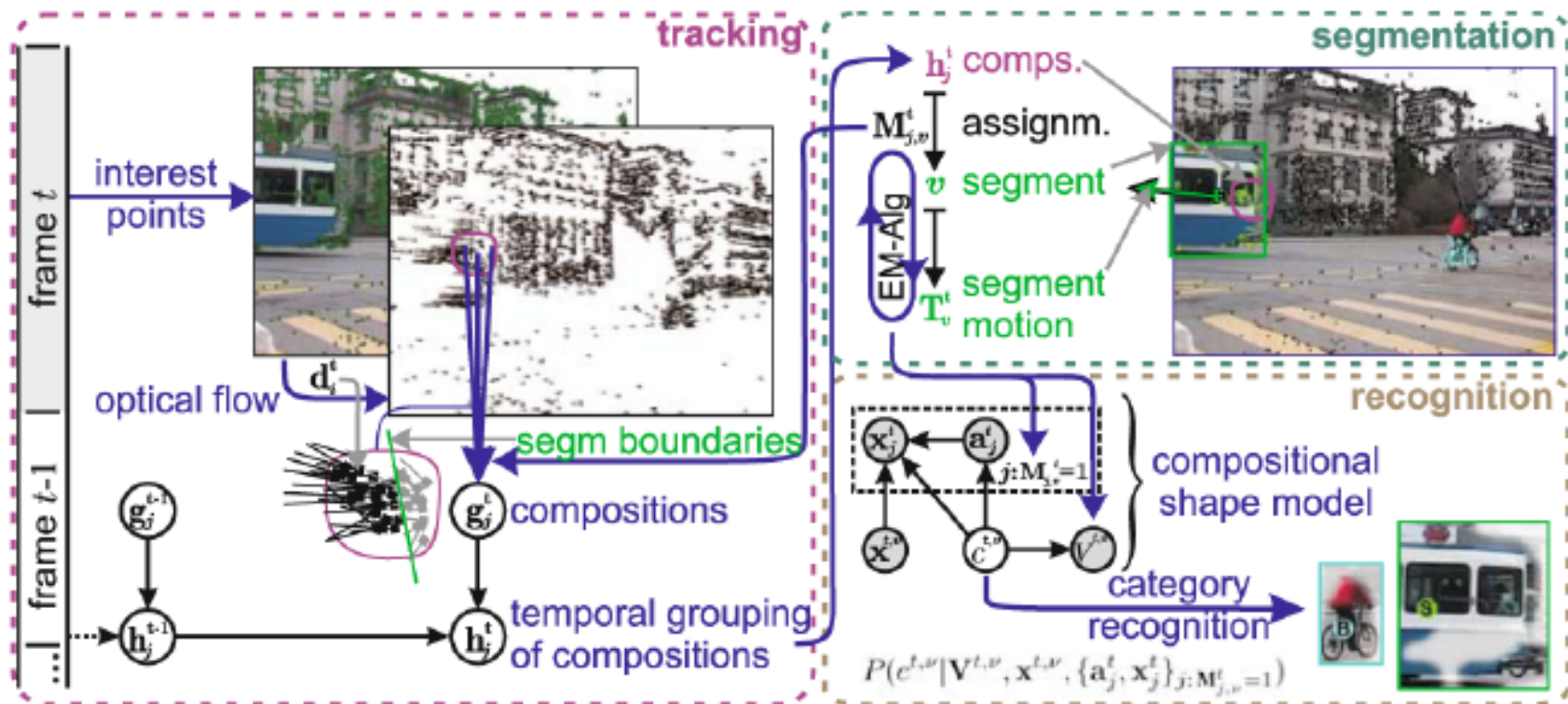


(d)

# Problem are investigated

- Category-level recognition
- Reducing supervision during learning
- Segmentation of videos from a moving camera
- Tracking without manual interaction
- Object models and shape representation

# Process Pipeline



# Region Tracking and Object Segmentation

- Tracking Object Regions
  - Compositions as Spatial Groupings of Parts
  - A composition represents all its constituent interest points
  - Tracking Compositions
  - Temporal Grouping of Composition

$$\Gamma^t(j) = \{i : \text{IP } i \text{ in neighborhood of } j\text{-th comp.}\}. \quad (1)$$

$$\mathbf{g}_j^t := \frac{1}{|\Gamma^t(j)|} \sum_{i \in \Gamma^t(j)} \mathbf{d}_t^i. \quad (2)$$

$$\mathbf{x}_j^{t+1} := \mathbf{x}_j^t + \frac{1}{|\Gamma^t(j)|} \sum_{i \in \Gamma^t(j)} \mathbf{d}_t^i. \quad (3)$$

$$\mathbf{h}_j^t = \eta \mathbf{g}_j^t + (1 - \eta) \mathbf{h}_j^{t-1}. \quad (4)$$

# Joint Tracking and Segmentation of Objects Based on Floating Image Regions

- Problem: assemble the object regions into the different objects and into background.
- Solving: Expectation-Maximization approach (EM)
- After use EM-Algorithm: Using Segmentation to Refine Object Region Tracking Algorithm



# EM-Algorithm

COMPSEGMENTATION( $\{\mathbf{h}_j^t\}_j, \{\mathbf{T}_v^{t-1}\}_{v=1,\dots,K}$ )

1 Initialization:  $\forall v : \mathbf{T}_v^t \leftarrow \mathbf{T}_v^{t-1}$

2 **repeat**

3 E-Step:  $\triangleright$  update assignments:

4  $\mathbf{M}_{j,v}^t \leftarrow \mathbf{1} \left\{ v = \operatorname{argmin}_{\hat{v}} \mathcal{R} \left( \mathbf{T}_{\hat{v}}^t, \mathbf{h}_j^t, \mathbf{x}_j^t \right) \right\}$

5 M-Step:  $\triangleright$  update segments:

6 **for**  $v = 1, \dots, K$

7 **do** Solve with Levenberg-Marquardt

(start with  $\hat{\mathbf{T}}_v^t \leftarrow \mathbf{T}_v^t$ ):

$\mathbf{T}_v^t \leftarrow \operatorname{argmin}_{\hat{\alpha}, \hat{s}, \hat{\delta}_x, \hat{\delta}_y} \sum_j \mathbf{M}_{j,v}^t \mathcal{R} \left( \hat{\mathbf{T}}_v^t, \mathbf{h}_j^t, \mathbf{x}_j^t \right)$

8 **until** convergence of  $\mathbf{M}_{j,v}^t$

9 **return**  $\mathbf{M}^t, \{\mathbf{T}_v^t\}_{v=1,\dots,K}$

# Using Segmentation to Refine Object Region Tracking Algorithm

```
COMPOSITIONTRACKING( $\{\mathbf{h}_j^{t-1}, \mathbf{x}_j^t\}_j, \{\mathbf{T}_v^{t-1}\}_{v=1,\dots,K}$ )  
1  Detect interest points  $i$  in frame  $t$   
2  for all compositions  $j \triangleright$  update comps with IP flow:  
3  do  $\Gamma^t(j) \leftarrow \{i : \|\mathbf{x}_j^t - \bar{\mathbf{x}}_i^t\| \leq w\}$   
4      $\mathbf{g}_j^t \leftarrow \frac{1}{|\Gamma^t(j)|} \sum_{i \in \Gamma^t(j)} \mathbf{d}_t^i$   
5      $\mathbf{h}_j^t \leftarrow \eta \mathbf{g}_j^t + (1 - \eta) \mathbf{h}_j^{t-1}$   
6   $\mathbf{M}^t, \{\mathbf{T}_v^t\}_v \leftarrow \text{COMPSEGMENTATION}(\{\mathbf{h}_j^t\}_j, \{\mathbf{T}_v^{t-1}\}_v)$   
7  for all compositions  $j \triangleright$  update comps with segmentat.:  
8  do  $\Gamma^t(j) \leftarrow \{i : i \in \Gamma^t(j) \wedge$   
            $1 = \mathbf{M}_{j, \text{argmin}_v^{\hat{\varphi}} \mathcal{R}(\mathbf{T}_v^t, \mathbf{d}_t^i, \bar{\mathbf{x}}_i^t)}^t$   
9      $\mathbf{g}_j^t \leftarrow \frac{1}{|\Gamma^t(j)|} \sum_{i \in \Gamma^t(j)} \mathbf{d}_t^i$   
10     $\mathbf{h}_j^t \leftarrow \eta \mathbf{g}_j^t + (1 - \eta) \mathbf{h}_j^{t-1}$   
11     $\mathbf{x}_j^{t+1} \leftarrow \mathbf{x}_j^t + \frac{1}{|\Gamma^t(j)|} \sum_{i \in \Gamma^t(j)} \mathbf{d}_t^i$   
12  return  $\{\mathbf{h}_j^t, \mathbf{x}_j^{t+1}\}_j, \{\mathbf{T}_v^t\}_v$ 
```

# Object Representations for Category-Level Recognition

- Compositional, Appearance-Based Model: Use multi-class SVM
- Recognition Using the Motion of Dot Patterns: Use SVM
- Global Shape and Local Appearance Combined
- Processing Pipeline for Training: Use SVM

# Output



(a)



(b)



(c)



(d)



(e)



(f)

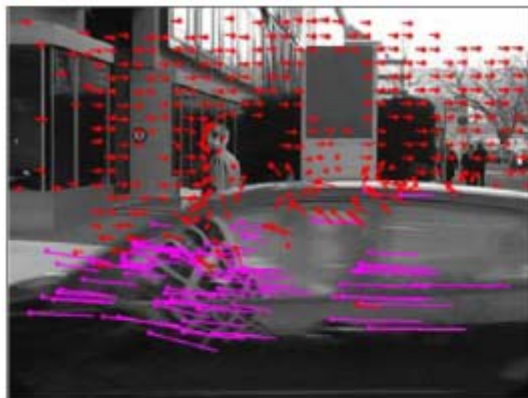


(g)



(h)

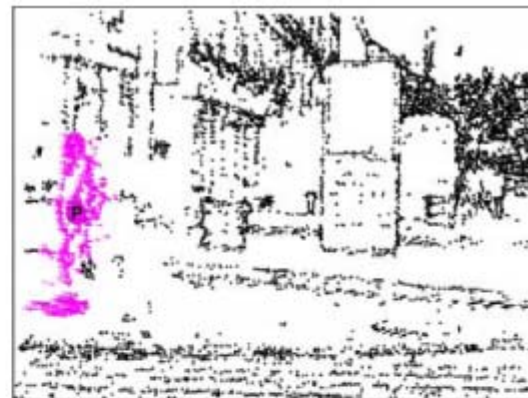
# Process Output



(a)



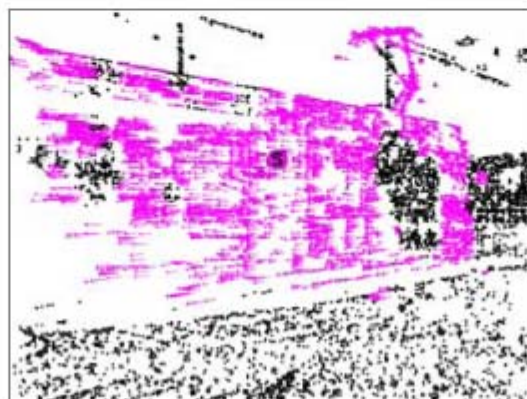
(b)



(c)



(d)



(e)

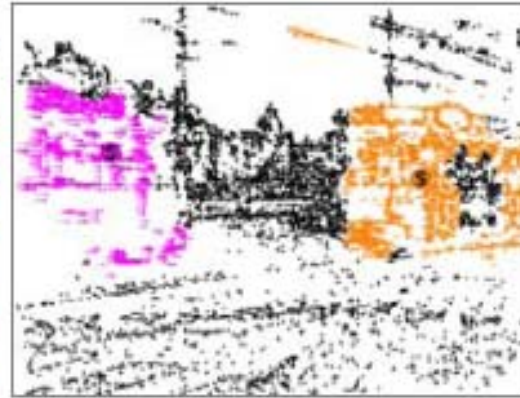


(f)

# Process Output(2)



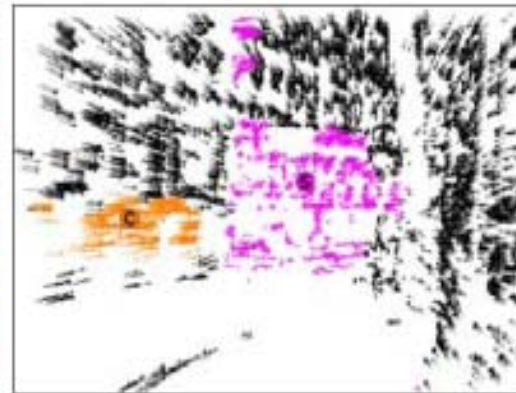
(a)



(b)



(c)



(d)



# Experiments

---

- Recognition Performance on Videos with Substantial Camera Motion
  - use 10-fold cross-validation and train on 16 randomly
  - Object models are learn on a randomly drawn subset of 15 frames per train video

# Experiment Output

- Baseline Performance of Appearance w/o Compositions and Shape—Bag-of-Parts
  - 53.0 ± 5.6% of all frames correctly.
- Compositional Segmentation and Recognition w/o Shape Model
  - 64.9 ± 5.4% per frame.



# Experiment Output(2)

- Comparing the Different Algorithm

Object model	Per frame	Per video
<i>Dataset of Ommer and Buhmann (2007)</i> <i>(car, bicycle, pedestrian, streetcar):</i>		
Approach of Ommer and Buhmann (2007)	$74.3 \pm 4.3$	$87.4 \pm 5.8$
Compositional motion (13)	$52.6 \pm 1.1$	$68.2 \pm 3.4$
Appearance-only: bag-of-parts	$53.0 \pm 5.6$	$58.9 \pm 6.5$
Segment. w/o shape: bag-of-comps	$64.9 \pm 5.4$	$78.9 \pm 5.8$
Shape: $P(c^{t,v}   V^{t,v})$ (11)	$74.4 \pm 5.3$	$88.4 \pm 5.2$
Compositional appear + location (12)	$79.6 \pm 5.5$	$90.7 \pm 5.3$
Combined shape + appear (14)	$81.4 \pm 2.9$	$94.5 \pm 4.9$
<i>Dataset (Ommer and Buhmann 2007) plus additional category</i> <i>"cow" from (Magee and Boyle 2002):</i>		
Compositional appearance (12)	$76.5 \pm 2.4$	$88.4 \pm 2.3$

# Experiment Output(3)

- Comparing Different Object Models

---

True classes →	Bicycle	Car	Pedest	Streetcar
Bicycle	<b>74.3</b>	3.2	13.7	2.9
Car	7.8	<b>84.1</b>	4.2	5.9
Pedestrian	13.3	2.5	<b>80.0</b>	3.9
Streetcar	4.7	10.2	2.2	<b>87.3</b>

---

# Experiment Output(4)

- Computational Demands
  - recognizes objects in videos of 768x576 pixel
  - using the combined shape and appearance model at the order of 1 fps on a 3 GHz Pentium 4 desktop PC.

Processing step	Comp. demand
Tracking and segmentation, Algorithm 2:	
IPs $i$ , flow $d_i^j$ (Algorithm 2, line 1)	27.7%
Updating comps (Algorithm 2, line 2–5)	5.2%
EM estimation Algorithm 1, i.e. (Algorithm 1, line 6)	4.9%
Updating comps with segm. (Algorithm 2, line 7–11)	0.3%
Feature extraction and recognition:	
Computing loc feat histos to represent $a_j^f$ (Sect. 3.2)	36.5%
Computing all individual probs in (14)	12.3%
Eval. GM of Fig. 5, i.e. calc. product in (14)	0.09%
Video stream ops, writing of results, etc.	12.9%

# Experiment Output(5)

- Action Recognition using KTH

True classes →	Box	Help	Hwav	Jog	Run	Walk
Boxing	<b>84.5</b>	0.0	5.5	0.0	0.0	0.0
Hand clapping	1.0	<b>87.0</b>	16.5	0.0	0.0	0.0
Hand waving	12.5	13.0	<b>75.5</b>	0.0	0.0	0.0
Jogging	0.0	0.0	0.5	<b>93.0</b>	0.0	0.0
Running	2.0	0.0	0.0	3.0	<b>92.3</b>	5.0
Walking	0.0	0.0	2.0	4.0	7.7	<b>95.0</b>