Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Human Action Recognition via Keypoint Tracking

Yunus Emre Kara

Department of Computer Engineering Boğaziçi University

MS Thesis Presentation, 17.01.2011

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00	00			
000		00	000	
0000	00		000	



1 Introduction

Challenges Approaches Datasets Our Approach

The Generic Keypoint Tracker (2)

Feature Extraction (3)





Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Action Recognition vs. Action Detection

- The two concepts are usually confused.
- Detection answers the question:
 "Is there an action in this sequence?"
- Recognition answers the question:
 "What is the action, given that there exists an action in this sequence?"
- In this work, we deal with the recognition of actions.

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
• 0 00 000				
			000	

Challenges

Action variations

- Actions can have large variations in performance
 - An action can differ in speed
 - Stride length can differ
- Antropometric differences
- Avoiding obstacles
- Temporal variations
 - Segmenting the actions in time is difficult
 - The length of an action can differ for different people

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00	0			
000		000	000	

Challenges

Environment and recording settings

- Person localization: Harder in cluttered or dynamic environments
 - Dynamic backgrounds
 - Moving camera
- Parts of the person can be occluded
- Lighting conditions can affect appearance
- An action has different observations for different viewpoints

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
•0				
000		00	000	
0000	00		000	

Approaches

Global Representations

- Person is localized first
- The region of interest is encoded as a whole
- Bobick and Davis¹ use MEI and MHI for recognizing actions
- Efros et al.² calculate optical flow measurements
- Gorelick et al.³ stack silhouettes of the consecutive frames to form spacetime volumes.

¹Bobick, A. and J. Davis, "The recognition of human movement using temporal templates", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 23, No. 3, pp. 257-267, Mar. 2001.

²Efros, A. A., A. C. Berg, G. Mori and J. Malik, "Recognizing action at a distance", Proceedings Ninth IEEE International Conference on Computer Vision, Vol. 2, No. October, pp. 726-733, 2003.

³Gorelick, L., M. Blank, E. Shechtman, M. Irani and R. Basri, "Actions as spacetime shapes.", IEEE transactions on pattern analysis and machine intelligence, Vol. 29, No. 12, pp. 2247-2253, Dec. 2007.

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
000			000	
0000			000	

Approaches

Local Representations

- Some patches in the images are encoded separately
- Patches are extracted from the neighborhoods of keypoints
- The representation is found by combining the information of the patches
- Laptev ⁴ proposed space-time interest points by extending the Harris interest point detector to 3D
- Willems et al.⁵ use integral videos to find salient points.
- Sun et al.⁶ find SIFT descriptors and track interest points using these descriptors.

⁴Laptev, I., "On Space-Time Interest Points", International Journal of Computer Vision, Vol. 64, No. 2-3, pp. 107-123, Sep. 2005.

⁵Willems, G., T. Tuytelaars and L. Van Gool, "An efficient dense and scale-invariant spatio-temporal interest point detector", Computer Vision ECCV 2008, pp. 650-663, 2008.

⁶Sun J, Wu X, Yan S, et al. "Hierarchical spatio-temporal context modeling for action recognition", Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. 2009:2004-2011.

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
0000 0000			000 000	

Datasets

KTH

- 2391 sequences, 25 people, 4 scenarios, 3-4 repeats, low resolution
- 6 actions: walking, running, jogging, boxing, hand waving, hand clapping

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
0000 0000			000 000	

Datasets

URADL

• 150 sequences, 5 people, 3 repeats, high resolution

 10 actions: Answering a phone, Chopping a banana, Dialing a phone, Drinking water, Eating banana, Eating snack chips, Looking up a phone number in a phone book, Peeling a banana, Using silverware, Write a phone number on a white board

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00 00 000 0000			0 0 000 000	

Datasets

WeCare

720 sequences, 6 people, 10 repeats, moderate resolution

 8 actions: Walking, Jumping, Sitting on the armchair, Standing up from the armchair, Lying on the gym mat, Standing up from the gym mat, Falling onto the armchair, Falling onto the gym mat

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00 00 000 ●000			0 0 000 000	

Outline

- A local approach
- Tracks keypoints
- Uses their trajectories



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
000 0 00 0			000 000	

What is a keypoint?

- Keypoint: Important points; such as corners, blobs, t-junctions
- Usage: object detection, image registration, camera calibration, image mosaicing, ...
- <u>Methods</u>: Harris, FAST, SIFT, SURF, ...
- Keypoint descriptor: Gives information about the neighborhood of a keypoint

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
000 0000			000 000	

Why do we use trajectories?



• Keypoints are robust for tracking objects

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
000 0000			000 000	

Why do we use trajectories?



- Keypoints are robust for tracking objects
- Localization and background subtraction are not required

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
000 0000			000 000	

Why do we use trajectories?



- Keypoints are robust for tracking objects
- Localization and background subtraction are not required
- Trajectories of tracked keypoints give information about the motion

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00				
000 0000			000 000	

Why do we use trajectories?



- Keypoints are robust for tracking objects
- Localization and background subtraction are not required
- Trajectories of tracked keypoints give information about the motion
- The motion information is invariant to the appearance

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Some advantages of our approach

Invariant to some viewpoint changes

Invariant to partial occlusions

• Robust to the temporal length variations.

• Invariant to the location of the action.

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Introduction



3 Feature Extraction



5 Conclusions

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Outline of the Generic Keypoint Tracker



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00	00			
000		00	000	

Keypoint detection

- Finding point correspondences between two images:
 - Keypoint detection
 - Descriptor extraction
 - Keypoint matching
- Our algorithm allows the use of external keypoint detector and descriptor extraction modules
- Built-in "Keypoint Matching" module
- We use <u>SURF method</u>⁷ both for <u>keypoint detection</u> and descriptor extraction steps.
- SURF keypoint descriptor: 64 dimensional vector extracted from the keypoint's local patch

¹Bay, H., A. Ess, T. Tuytelaars and L. V. Gool, "Speeded-Up Robust Features (SURF)", Computer Vision and Image Understanding, Vol. 110, No. 3, pp. 346-359, 2008.

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	○ 0			
000		00	000	

- Gaps in time are allowed \Rightarrow robust to occlusions
- Missing keypoints are interpolated linearly
- Trajectory descriptor: Calculated using the descriptors of member keypoints
- We compare the descriptors of keypoints with the descriptors of trajectories
- The descriptor distance of <u>the best match</u> is <u>divided by</u> the descriptor distance of <u>the second best match</u>.
 Small ratio ⇒ Mark for matching
- For matching, we must mark in both ways. The matched trajectory of a keypoint must also be matched with that keypoint.

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	00			
000		00	000	
0000	00		000	



• Current trajectories (blue curves)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	00			
000		00	000	
0000	00		000	



- Current trajectories (blue curves)
- Newly detected keypoints (red points)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	00			



- Current trajectories (blue curves)
- Newly detected keypoints (red points)
- Spatial neighborhood of a trajectory

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	00			



- Current trajectories (blue curves)
- Newly detected keypoints (red points)
- Spatial neighborhood of a trajectory
- Best matching keypoint (green point)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	00			



- Current trajectories (blue curves)
- Newly detected keypoints (red points)
- Spatial neighborhood of a trajectory
- Best matching keypoint (green point)
- Spatial neighborhood of the best matching keypoint

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	00			



- Current trajectories (blue curves)
- Newly detected keypoints (red points)
- Spatial neighborhood of a trajectory
- Best matching keypoint (green point)
- Spatial neighborhood of the best matching keypoint
- Best matching trajectory (green curve)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000	•	00	000	



• Matching is done for each trajectory (green curves)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000	•	00	000	



• Matching is done for each trajectory (green curves)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000	•	00	000	



• Matching is done for each trajectory (green curves)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000	•	00	000	



- Matching is done for each trajectory (green curves)
- Some trajectories may not be matched with a keypoint

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000	•	00	000	



- Matching is done for each trajectory (green curves)
- Some trajectories may not be matched with a keypoint

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	•			



- Matching is done for each trajectory (green curves)
- Some trajectories may not be matched with a keypoint
- Unmatched keypoints initiate new trajectories (green points)

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	● 0		000	

Elimination and Storing

Stationary points

The stationary points are found by looking at the diagonal length of the trajectory's bound box.



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
	00			

Elimination and Storing

- We eliminate trajectories
 - for narrowing down the search space
 - and removing noise
- We store trajectories if
 - they are ready for feature extraction
 - and not modified for some time
- We do these operations after each frame



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	

Introduction

2 The Generic Keypoint Tracker

Feature Extraction
 Normalizing Against Time
 Normalizing Against Spatial Position
 Trajectory Feature Extraction
 Extracting Sub-trajectories
 Bag-of-trajectories

4 Results



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Feature extraction outline



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
		•		
000		00	000	
0000	00		000	

Normalizing Against Time



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
		•		
000		00	000	
0000	00		000	

Normalizing Against Time





Normalized Against Time



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
		•		
000		00	000	
0000	00		000	

Normalizing Against Spatial Position



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
		•		
000		00	000	
0000	00		000	

Normalizing Against Spatial Position



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
		○ ○		

Trajectory Feature Extraction

- Start time: t₁
- End time: t_n
- Mean x coordinate of keypoints (x_m)
- Mean y coordinate of keypoints (y_m)
- Length of the path until halftime (red path)
- Length of the path since halftime (blue path)



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
		00		

Trajectory Feature Extraction

- x and y components of the vector passing through the first and the last keypoints (L_{N,x} and L_{N,y})
- Maximum $|E_i|$, with sign • $\frac{L_i \cdot L_N}{||L_N||^2}$ _____

The vector passing through the first and i^{th} keypoints: $L_i = (x_i - x_1, y_i - y_1)$



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
		•		

Extracting Sub-trajectories

- Trajectories differ in length
- We divide them into smaller subparts
- These subparts are used in describing videos



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Bag-of-trajectories

- Trajectory features are normalized (t-statistic or min-max)
- Normalized features are clustered using k-means
- Image sequence descriptor: Normalized histogram of the clusters of trajectories



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	



(2) The Generic Keypoint Tracker

3 Feature Extraction

4 Results

Parameters Experiment Setup Validation Results Test Results



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00			0	
			000 000	

Parameters

Parameters

- σ the minimum required sample count
- u the maximum allowed time since last observation
- ho the minimum age to be eliminated
- ω $\;$ the maximum number of keypoints for a sub-trajectory
- n normalization method before bag-of-trajectories
- K cluster count for K-means of bag-of-trajectories
- C SVM cost parameter
- k k-NN k value depending on the method

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00			0	
000 0000			000 000	

Experiment Setup

Experiment Setup

- A total of 23994 experiments
 - 4300 experiments on the KTH dataset
 - 15394 experiments on the URADL dataset
 - 4300 experiments on the WeCare dataset
- 10-fold cross validation on KTH and WeCare datasets
- 5-fold cross validation on the URADL datasets
- SVM with χ^2 kernel is tried with 11 different cost parameter values
- k-NN is tried with all k's in the range 1-32

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
			000	

Validation Results

Validation Results (1)



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Validation Results

Validation Results (2)



Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
			000	
0000	00		000	

Validation Results

Validation Results (3)



36/43

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000			000	

Test Results

KTH dataset

Accuracy: 87.25%

 <u>Classes:</u> C₁: walking, C₂: jogging, C₃: running, C₄: boxing, C₅: hand clapping, C₆: hand waving

				Actua	l Class		
		C_1	<i>C</i> ₂	<i>C</i> ₃	<i>C</i> ₄	C_5	C_6
	C_1	130	8	0	0	0	0
	C_2	14	126	26	0	0	0
uc	C_3	0	10	118	0	0	0
ctić	C_4	0	0	0	139	15	5
šdie	C_5	0	0	0	4	129	28
Pre	C_6	0	0	0	0	0	111

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00	00			
000		00	000	
			000	

Test Results

URADL dataset

- Accuracy: 88%
- <u>Classes</u>: C₁: answer phone, C₂: chop banana, C₃: dial phone, C₄: drink water, C₅: eat banana, C₆: eat snack, C₇: lookup in phonebook, C₈: peel banana, C₉: use silverware, C₁₀: write on whiteboard

						Actua	I Class	s			
		C_1	C_2	<i>C</i> ₃	C_4	C_5	C_6	<i>C</i> ₇	C_8	C_9	C_{10}
	C_1	15	0	1	0	0	0	0	0	0	0
	C_2	0	14	0	0	0	0	0	0	0	0
	C_3	0	1	14	0	2	0	0	1	0	0
	<i>C</i> ₄	0	0	0	14	0	0	0	1	0	0
uo	C_5	0	0	0	0	12	0	0	2	0	0
cti.	C_6	0	0	0	0	0	12	0	2	0	0
edi	C ₇	0	0	0	0	0	2	14	0	1	0
Pr	<i>C</i> ₈	0	0	0	1	1	1	0	8	0	0
	C_9	0	0	0	0	0	0	1	1	14	0
	C_{10}	0	0	0	0	0	0	0	0	0	15

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00	00			
000			000	
0000			000	

Test Results

WeCare dataset

- Accuracy: 98.75%
- <u>Classes:</u> C₁: walking, C₂: jumping, C₃: sitting on the armchair, C₄: standing up from the armchair, C₅: lying on the gym mat, C₆: standing up from the gym mat, C₇: falling onto the armchair, C₈: falling onto the gym mat

		Actual Class							
		C_1	<i>C</i> ₂	<i>C</i> ₃	<i>C</i> ₄	C_5	C_6	C ₇	<i>C</i> ₈
	C_1	80	0	0	0	1	0	0	0
	C_2	0	20	0	0	0	0	0	0
	C_3	0	0	20	0	0	0	1	0
uo	C_4	0	0	0	40	0	0	0	0
ct i	C_5	0	0	0	0	18	0	0	0
edi	C_6	0	0	0	0	0	20	0	0
Pr	C_7	0	0	0	0	0	0	19	0
	C_8	0	0	0	0	1	0	0	20

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
000		00	000	
0000	00		000	

Conclusions (1)

- A novel method for tracking keypoints is proposed.
- A feature set for describing trajectories is proposed.
- Performance is comparable to the methods in the literature.
- The KTH dataset
 - Our approach: 87.25%
 - Schuldt et al.⁸: 71.72%
 - Niebles et al.⁹: 81.5%
 - Laptev et al.¹⁰: 91.8%
 - Messing et al.¹¹: 74%

⁸Schuldt, C., I. Laptev and B. Caputo, "Recognizing Human Actions: A Local SVM Approach", Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004., Vol. 3, pp. 32-36 Vol.3, IEEE, 2004.

⁹Niebles, J. C., H. Wang and L. Fei-Fei, "Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words", International Journal of Computer Vision, Vol. 79, No. 3, pp. 299-318, Mar. 2008.

¹⁰Laptev, I., M. Marszalek, C. Schmid and B. Rozenfeld, "Learning realistic human actions from movies", Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, pp. 1-8, IEEE, 2008.

¹¹Messing, R., C. Pal and H. Kautz, "Activity recognition using the velocity histories of tracked keypoints", IEEE 12th International Conference on Computer Vision, pp. 104-111, Sep. 2009. Human Action Recognition via Keypoint Tracking, Yunus Emre Kara

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00	00			
000		00	000	
0000			000	

Conclusions (2)

The URADL dataset

- Our approach: 88%
- Messing et al. (Velocity Histories): 63%
- Messing et al. (Latent Velocity Histories): 67%
- Messing et al. (Augmented Velocity Histories): 89%
- However, augmented velocity histories method of Messing et al. is highly dependent on the position of the action on the frame and it requires the position of the face.
- A new dataset is introduced: WeCare
- Our performance of the WeCare dataset is 98.75%
- More challenging scenarios should be added to the WeCare dataset for testing our work.

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions

Future directions

- More robust elimination module
- New trajectory features
- Improved recognition performance
- Human action <u>detection</u> system
- Trajectories in 3D space
- Real time recognition support

Introduction	The Generic Keypoint Tracker	Feature Extraction	Results	Conclusions
00	00			
000		00	000	
0000			000	

THANK YOU