

Social behavior recognition in continuous video

Xavier P. Burgos-Artizzu¹

¹ Electrical Engineering California Institute of Technology {xpburgos,perona}@caltech.edu

Summary • Largest behavior dataset to date. • Social behavior. Segmentation and classification of actions in continuous video. Novel trajectory features. • Temporal context in behavior analysis. **Caltech Resident-Intruder** Mouse Dataset (CRIM13) http://www.vision.caltech.edu/Video Datasets/CRIM13 Walk away from intruder Circle around Chase intruder $p=3.4\%, \mu = 1.2, \sigma = 0.5$ $p=3.3\%, \mu = 1.6, \sigma = 0.6$ $p=1.0\%, \mu=0.8, \sigma=0.3$ $p=0.3\%, \mu = 2.4, \sigma = 2.5$ Copulation, courts intruder $(p=3.4\%, \mu = 2.3, \sigma = 17.4 p=4.2\%, \mu = 3.2, \sigma = 40.5 p=0.3\%, \mu = 4.0, \sigma = 30.2 p=1.6\%, \mu = 9.5, \sigma = 104$ Sniff any body part of intruder Up Stands in its back legs Human intervenes

 $p=7.6\%, \mu = 2.6, \sigma = 25.4$ $p=1.2\%, \mu = 3.5, \sigma = 10.5$ $p=14.4\%, \mu = 2.7, \sigma = 27.7$ $p=3.8\%, \mu = 2.1, \sigma = 12.0$

Largest behavior dataset to date: 88 hours and 8 million frames.

Continuous, unsegmented videos: ~10 min, 140 behavior instances/video.

Real-world scenario: high variance in behavior frequencies and durations, both inter/intra-class.

Frame-by-frame expert annotation: experts spent 350 hours annotating videos. **13 Meaningful behaviors:** basis for a state-of-the-art neuroscience study of behavior, published in Nature.

Focused in social interaction: dataset collected to study neurophysiological mechanisms involved in aggression and courtship in mice.

Synchronized top and side views: allow to study behavior from different views.

Spatio-temporal features benchmark

Detector+ Descriptor	Performance	Codebook size	fps
Harris3D+ Pca-Sift	20.9%	250	2.7
Harris3D+ Hog3D	18.7%	250	4.0
Harris3D+ Hog/Hof	15.5%	500	1.1
Cuboids+ Pca-Sift	24.6%	250	4.5
Cuboids+ Hog3D	18.2%	250	8.7
Cuboids+ Hog/Hof	19.8%	500	1.6
Cuboids+ Pca-Sift multi-scale	16.4%	1000	0.8
LTP	22.2%	-	15

Differences are small.

• We use Cuboids+Pca-Sift.

Piotr Dollár²

² Interactive Visual Media Microsoft Research, Redmond pdollar@microsoft.com

Dayu Lin³

³ Dept. of Phsychiatry and Neuroscience NYU Medical Center dayu.lin@nyumc.org

Proposed method



Dir change Dir difference

Acceleration

 $\mathbf{CDir}_{\mathbf{i}}(\mathbf{t}) = Dir_{i}(t) - Dir_{i}(t-1)$ $\mathbf{DDir}(\mathbf{t}) = |Dir_1(t) - Dir_2(t)|$ **Velocity** $\begin{bmatrix} \mathbf{V}\mathbf{x}_{\mathbf{m}_{i}}(\mathbf{t}) \\ \mathbf{V}\mathbf{y}_{\mathbf{m}_{i}}(\mathbf{t}) \end{bmatrix} = \frac{1}{\Delta t} \begin{bmatrix} x_{m_{i}}(t-1) - x_{m_{i}}(t-2) & x_{m_{i}}(t) - x_{m_{i}}(t-1) & x_{m_{i}}(t+1) - x_{m_{i}}(t) \\ y_{m_{i}}(t-1) - y_{m_{i}}(t-2) & y_{m_{i}}(t) - y_{m_{i}}(t-1) & y_{m_{i}}(t+1) - y_{m_{i}}(t) \end{bmatrix} * \begin{bmatrix} 0.25 \\ 0.5 \\ 0.25 \end{bmatrix}$ $\begin{bmatrix} \mathbf{A}\mathbf{x}_{\mathbf{m}_{\mathbf{i}}}(\mathbf{t}) \\ \mathbf{A}\mathbf{y}_{\mathbf{m}_{\mathbf{i}}}(\mathbf{t}) \end{bmatrix} = \begin{bmatrix} \frac{Vx_{m_{i}}(t+1) - Vx_{m_{i}}(t-1)}{2\Delta t} \\ \frac{Vy_{m_{i}}(t+1) - Vy_{m_{i}}(t-1)}{2\Delta t} \end{bmatrix}$

David J. Anderson⁴

Pietro Perona¹

⁴ Division of Biology California Institute of Technology wuwei@caltech.edu

confidence (all, past, post)

 $\begin{bmatrix} \mu(\mathbf{i} - \frac{\mathbf{s}\mathbf{z}}{2}, \mathbf{i} + \frac{\mathbf{s}\mathbf{z}}{2}) & \mu(\mathbf{i} - \frac{\mathbf{s}\mathbf{z}}{2}, \mathbf{i}) & \mu(\mathbf{i}, \mathbf{i} + \frac{\mathbf{s}\mathbf{z}}{2}) \end{bmatrix}$ where, $\mu(\text{start}, \text{end}) = \frac{\sum_{start}^{end} h_{t-1}^k}{(end - start)}$

lst order $\left[\begin{array}{c} \frac{\partial \mathbf{h}_{t-1}^{\mathbf{k}}}{\partial t}\big(i-\frac{\mathbf{s}\mathbf{z}}{2},i+\frac{\mathbf{s}\mathbf{z}}{2}\big) & \frac{\partial \mathbf{h}_{t-1}^{\mathbf{k}}}{\partial t}\big(i-\frac{\mathbf{s}\mathbf{z}}{2},i\big) & \frac{\partial \mathbf{h}_{t-1}^{\mathbf{k}}}{\partial t}\big(i,i+\frac{\mathbf{s}\mathbf{z}}{2}\big)\end{array}\right]$ derivative where, $\frac{\partial \mathbf{h}_{t-1}^{k}}{\partial \mathbf{t}}(\mathbf{start}, \mathbf{end}) = h_{t-1}^{k}(end) - h_{t-1}^{k}(start)$ (all, past, post)

Results

CJUICS				
	Features used	No Temporal Context	With Temporal Context	
	TF	52.3%	58.3%	
	STF (both)	29.3%	43.0%	
	STF (top)	26.6%	39.3%	
	STF (side)	28.2%	39.1%	
	(Full method) TF + STF	53.1%	61.2%	

Human agreement **69.7%**

Proposed approach 5 62.6%

> **Ground truth Proposed method**

Conclusions

- CRIM13 is a realistic and challenging dataset.
- Novel trajectory features are more discriminative than spatio-temporal features.

- Temporal context is crucial in continuous videos. • Our method's performance is not far from that of trained human annotators.



Office of Naval Researc

• Spatio-temporal features are not sufficient.

 Novel weak trajectory features (TF) perform 20% better than STF.

 Combination of both features achieves best results.

 Temporal context is crucial, improving final classification 8.1% (14% in average).

Comparison with human annotations

