

# A Benchmark Dataset and Evaluation Methodology for Video Object Segmentation

## *Supplemental Material*

F. Perazzi<sup>1,2</sup> J. Pont-Tuset<sup>1</sup> B. McWilliams<sup>2</sup>  
<sup>1</sup>ETH Zurich

L. Van Gool<sup>1</sup> M. Gross<sup>1,2</sup> A. Sorkine-Hornung<sup>2</sup>  
<sup>2</sup>Disney Research

### 1. Contents

In this document we include additional material related to the dataset and the evaluation.

- **Dataset.** In Figure 1 we visualize a representative frame of each sequence with ground-truth masks overlaid in red. We refer the reader to the accompanying video for a complete visualization of the dataset. A summary of the requirements detailed in Section 3 (*manuscript*) can be found in Table 1.
- **Attributes** In Table 3 we summarize the attributes assigned to each specific sequence. We refer the reader to Table 1 (*manuscript*) for a comprehensive description of the attributes. In Section 2 we provide an in-depth description of the statistical procedure we adopted to determine the attribute dependencies, visualized in Figure 2 (*right, manuscript*).
- **Evaluation and Running Times.** In Figure 2 we visualize the mean performance of all evaluated approaches, based the per-sequence region-similarity  $\mathcal{J}$  and contour accuracy  $\mathcal{F}$ . The results are an estimator of the expected segmentation *difficulty* for a specific sequence. Sequences are sorted with respect to the estimated difficulty. In Tables 4, 5, 6 we summarize the per-sequence performance of each method, respectively in terms of region similarity  $\mathcal{J}$ , contour accuracy  $\mathcal{F}$  and temporal stability  $\mathcal{T}$ . In Figures 3, 4, 5, 6. Estimated running times are given in Table 2.

### 2. Attributes Dependency - In Depth

We consider the presence or absence of each attribute in a video sequence to be represented as a binary random variable  $X = (X_1, X_2, \dots, X_d)$ . The dependencies between the attributes can be modelled by a pairwise Markov random field (MRF) defined on the undirected graph  $G = (V, E)$  where  $V = \{1, \dots, d\}$  is the set of vertices and  $E$  is the (unknown) set of edges. Each variable  $X_s$  is associated

Dataset	SIZE	HD-Q	VARY	DENSE-GT	OBJ
OURS	✓	✓	✓	✓	✓
MoSeg [1]	✓				✓
BVSD [5]	✓	✓	✓		
SegTrack [6]				✓	✓
SegTrack v2 [2]				✓	✓

Table 1: Summary of requirements fulfilled by datasets most relevant to video object segmentation. *From left:* large overall size of the dataset (SIZE), high-resolution videos (HD-Q), variety of content and challenges (VARY), pixel-accurate, per-frame ground-truth (DENSE-GT) and object presence (OBJ). A detailed overview of the requirements is described in Section 3. Our dataset is the only one meeting all requirements.

with a vertex  $s \in V$ . The pairwise MRF associated with  $G$  is the family of distributions which factorise as

$$\mathbb{P}_\theta(x) \propto \exp \left\{ \sum_{(s,t) \in E} \theta_{s,t} x_s x_t \right\}.$$

The absence of an edge between  $s$  and  $t$  means that  $X_s$  and  $X_t$  are *independent* conditioned on their respective Markov blankets<sup>1</sup>. In other words, given the state of the neighbours of  $s$  and  $t$ , knowing  $t$  gives us no information about  $s$  and *vice-versa*.

Equivalently,  $\theta$  can be viewed as a  $\binom{d}{2}$ -dimensional vector which indexes all distinct pairs of vertices but is non-zero only when the vertex pair  $(s, t)$  belongs to the edge set  $E$  of the graph. Recovering  $E$  is equivalent to recovering the neighbourhood set  $\mathcal{N}(r) := \{t \in V | (r, t) \in E\}$  for each  $r \in V$ . Estimating the neighbourhood set  $\mathcal{N}(r)$  is equivalent to estimating the support (i.e. location of non-zero entries) of the  $(d - 1)$  dimensional sub-vector  $\theta_{\setminus r} := \{\theta_{ru}, u \in V \setminus r\}$ .

<sup>1</sup>For a MRF this consists of the neighbours of  $s$  and  $t$ , respectively

	Preprocessing			Unsupervised						Semi-Supervised					
	MCG	SF-LAB	SF-MOT	NLC	CVOS	TRC	MSG	KEY	SAL	FST	TSP	SEA	HVS	JMP	FCP
Time	1498.8s	85.3s	84.9s	880.3s	6190.4s	16612.3	2614.9s	55191.1s	149.4s	2319.2s	2870.3s	299.1s	634.1s	576.6s	2261.7s

Table 2: *Running times.* Estimated running times (in seconds) for each of the evaluated approaches. Due to the substantial processing power required to carry out this large scale evaluation, we used multiple machines and a cluster with thousands nodes and different CPU, therefore while computing times have been normalized to comparable processing power, they should be considered an approximate estimate.

Following [4], given a collection of  $n$  observations  $\mathcal{X}^n = \{x^{(1)}, \dots, x^{(n)}\}$  of  $d$ -dimensional binary vectors  $x^{(i)}$ , the support of each  $\theta_{\setminus r}$  can be estimated by solving the following minimization problem

$$\min_{\theta_{\setminus r} \in \mathbb{R}^{d-1}} -\frac{1}{n} \sum_{i=1}^n \log \mathbb{P}_{\theta}(x_r^{(i)} | x_{\setminus r}^{(i)}) + \lambda \|\theta_{\setminus r}\|. \quad (1)$$

Since the random variables are binary, minimizing the penalised negative log likelihood above corresponds to solving  $\ell_1$  penalised logistic regression treating  $x_{\setminus r}^{(i)} \in \mathbb{R}^{n \times (d-1)}$  as covariates and  $x_r^{(i)}$  as the response.

The solution to (1) can be highly sensitive to the regularization strength  $\lambda$  which controls the sparsity of the solution. In order to determine the correct degree of sparsity we employ *stability selection* [3]. Briefly, this amounts to performing the above procedure on multiple  $n/2$ -sized subsamples of the data and computing the proportion of times each edge is selected. Setting an appropriate threshold on this selection probability allows us to control the number of wrongly estimated edges according to Theorem 1 in [3]. For example, for a threshold value of 0.6 and choosing a value of  $\lambda$  which on average selects neighbourhoods of size 4, the number of wrongly selected edges is at most 4 (out of  $16^2 = 256$  possible edges).

## References

- [1] T. Brox and J. Malik. Object segmentation by long term analysis of point trajectories. In *ECCV*, 2010. [1](#)
- [2] F. Li, T. Kim, A. Humayun, D. Tsai, and J. M. Rehg. Video segmentation by tracking many figure-ground segments. In *ICCV*, 2013. [1](#)
- [3] N. Meinshausen and P. Bühlmann. Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4):417–473, 2010. [2](#)
- [4] P. Ravikumar, M. J. Wainwright, J. D. Lafferty, et al. High-dimensional ising model selection using  $\ell_1$ -regularized logistic regression. *The Annals of Statistics*, 38(3):1287–1319, 2010. [2](#)
- [5] P. Sundberg, T. Brox, M. Maire, P. Arbelaez, and J. Malik. Occlusion boundary detection and figure/ground assignment from optical flow. In *CVPR*, 2011. [1](#)
- [6] D. Tsai, M. Flagg, and J. M. Rehg. Motion coherent tracking with multi-label MRF optimization. In *BMVC*, 2010. [1](#)



Figure 1: Sample sequences from our dataset, with ground truth segmentation masks overlaid. Please refer to the accompanying video for a complete visualization of the dataset.

Sequence	AC	BC	CS	DB	DEF	EA	FM	HO	IO	LR	MB	OCC	OV	SC	SV
bear						✓									
blackswan															
bmx-bumps			✓			✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
bmx-trees	✓	✓			✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
boat		✓		✓		✓									✓
breakdance	✓				✓	✓	✓	✓	✓		✓		✓		
breakdance-flare			✓		✓		✓	✓			✓				
bus						✓		✓				✓		✓	
camel			✓		✓				✓						
car-roundabout				✓											
car-shadow	✓	✓				✓				✓					
car-turn		✓													✓
cows			✓		✓			✓	✓			✓			
dance-jump					✓	✓	✓	✓			✓	✓			✓
dance-twirl			✓		✓			✓	✓		✓		✓	✓	
dog			✓		✓	✓	✓	✓			✓				
dog-agility	✓				✓	✓	✓	✓	✓		✓	✓	✓	✓	
drift-chicane	✓			✓		✓	✓	✓	✓		✓				✓
drift-straight	✓		✓			✓	✓	✓	✓		✓	✓	✓	✓	
drift-turn	✓			✓			✓	✓	✓				✓		✓
elephant			✓	✓	✓	✓									
flamingo			✓	✓				✓	✓					✓	
goat	✓	✓			✓	✓									
hike					✓			✓		✓					
hockey					✓			✓	✓					✓	
horsejump-high					✓			✓	✓			✓		✓	
horsejump-low						✓	✓	✓	✓			✓		✓	
kite-surf			✓			✓		✓	✓		✓	✓		✓	✓
kite-walk			✓	✓				✓	✓			✓		✓	
libby			✓	✓				✓			✓	✓			✓
lucia			✓					✓				✓			
mallard-fly	✓			✓	✓	✓	✓			✓	✓		✓		✓
mallard-water				✓		✓			✓	✓					
motocross-bumps	✓	✓					✓	✓	✓				✓		✓
motocross-jump	✓					✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
motorbike						✓	✓	✓	✓	✓	✓		✓	✓	
paragliding								✓	✓	✓					✓
paragliding-launch						✓	✓	✓	✓						✓
parkour	✓					✓		✓	✓	✓	✓	✓			✓
rhino		✓				✓									✓
rollerblade			✓		✓		✓	✓	✓	✓	✓				
scooter-black							✓	✓	✓						✓
scooter-gray	✓						✓	✓	✓	✓		✓		✓	
soapbox	✓					✓		✓	✓		✓				✓
soccerball								✓	✓	✓	✓	✓			
stroller		✓		✓			✓	✓	✓						✓
surf		✓	✓				✓	✓	✓				✓		✓
swing				✓			✓	✓	✓			✓			✓
tennis				✓			✓	✓	✓		✓				✓
train						✓		✓							✓

Table 3: *List of attributes for each video in the dataset.*.. Left to right: appearance changes (AC), background clutter (BC), camera shake (CS), dynamic background (DB), non-linear deformation (DEF), edge ambiguity (EA), fast-motion (FM), heterogeneous object (HO), interacting objects (IO), low resolution (LR), motion blur (MB), occlusions (OCC), out-of-view (OV), shape complexity (SC), scale variation (SV). See Table 1 in the paper for the description of each attribute.

Sequence	Preprocessing			Unsupervised							Semi-Supervised				
	MCG	SF-LAB	SF-MOT	NLC	CVOS	TRC	MSG	KEY	SAL	FST	TSP	SEA	HVS	JMP	FCP
bear	<b>0.937</b>	0.126	0.556	<b>0.906</b>	0.864	0.873	0.851	0.891	0.657	0.898	0.778	0.912	<b>0.938</b>	0.929	0.906
blackswan	<b>0.871</b>	0.524	0.547	<b>0.874</b>	0.422	0.569	0.526	0.842	0.222	0.732	0.872	<b>0.933</b>	0.917	0.930	0.908
bmx-bumps	<b>0.490</b>	0.030	0.281	<b>0.635</b>	0.368	0.350	0.353	0.309	0.188	0.241	0.290	0.198	<b>0.428</b>	0.336	0.300
bmx-trees	<b>0.473</b>	0.021	0.468	<b>0.212</b>	0.121	0.162	0.188	0.193	0.194	0.180	0.095	0.113	0.178	0.229	<b>0.248</b>
boat	<b>0.619</b>	0.068	0.171	0.007	0.056	0.130	0.144	0.065	0.271	<b>0.361</b>	0.656	<b>0.793</b>	0.782	0.705	0.613
breakdance	<b>0.713</b>	0.204	0.649	<b>0.673</b>	0.183	0.114	0.237	0.549	0.422	0.467	0.056	0.329	0.550	0.478	<b>0.567</b>
breakdance-flare	<b>0.733</b>	0.206	0.629	<b>0.804</b>	0.317	0.245	0.157	0.559	0.476	0.616	0.040	0.131	0.499	0.430	<b>0.723</b>
bus	0.749	0.290	<b>0.797</b>	0.629	0.664	0.684	<b>0.885</b>	0.785	0.739	0.825	0.515	0.752	0.809	0.668	<b>0.832</b>
camel	<b>0.795</b>	0.015	0.620	0.768	<b>0.850</b>	0.778	0.756	0.579	0.320	0.562	0.654	0.649	<b>0.876</b>	0.640	0.734
car-roundabout	<b>0.786</b>	0.306	0.708	0.509	<b>0.871</b>	0.552	0.630	0.640	0.500	0.808	0.614	0.708	<b>0.777</b>	0.726	0.717
car-shadow	0.701	0.135	<b>0.765</b>	0.645	0.759	0.449	<b>0.880</b>	0.589	0.538	0.698	0.636	<b>0.775</b>	0.699	0.645	0.723
car-turn	<b>0.865</b>	0.071	0.593	0.833	0.820	0.805	0.621	0.806	0.611	<b>0.851</b>	0.323	<b>0.909</b>	0.810	0.834	0.724
cows	<b>0.811</b>	0.139	0.684	<b>0.883</b>	0.562	0.833	0.799	0.337	0.623	0.791	0.595	0.707	0.779	0.756	<b>0.812</b>
dance-jump	0.470	0.251	<b>0.582</b>	0.718	0.341	0.303	0.065	<b>0.748</b>	0.291	0.598	0.132	0.662	<b>0.680</b>	0.490	0.522
dance-twirl	<b>0.644</b>	0.093	0.624	0.347	0.452	0.366	0.366	0.380	0.372	<b>0.453</b>	0.099	0.117	0.318	0.444	<b>0.471</b>
dog	<b>0.621</b>	0.195	0.532	<b>0.809</b>	0.753	0.786	0.331	0.692	0.566	0.708	0.313	0.581	0.722	0.673	<b>0.774</b>
dog-agility	<b>0.663</b>	0.060	0.354	<b>0.652</b>	0.193	0.138	0.110	0.132	0.055	0.280	0.079	0.354	0.457	<b>0.699</b>	0.453
drift-chicane	<b>0.806</b>	0.046	0.396	0.324	0.313	0.722	<b>0.758</b>	0.188	0.244	0.667	0.018	0.119	0.331	0.243	<b>0.457</b>
drift-straight	<b>0.753</b>	0.171	0.427	0.473	0.344	0.431	0.575	0.194	0.268	<b>0.683</b>	0.197	0.513	0.295	0.618	<b>0.668</b>
drift-turn	<b>0.856</b>	0.161	0.359	0.154	0.615	0.412	<b>0.638</b>	0.255	0.349	0.533	0.162	0.667	0.276	<b>0.717</b>	0.606
elephant	<b>0.686</b>	0.099	0.640	0.518	0.494	0.760	0.689	0.675	0.510	<b>0.824</b>	0.666	0.553	0.742	<b>0.750</b>	0.655
flamingo	<b>0.850</b>	0.247	0.517	0.539	0.783	0.731	0.794	0.692	0.570	<b>0.817</b>	0.666	0.583	<b>0.811</b>	0.530	0.717
goat	<b>0.641</b>	0.057	0.138	0.010	0.074	<b>0.793</b>	0.736	0.705	0.257	0.554	0.444	0.535	0.580	<b>0.731</b>	0.677
hike	<b>0.900</b>	0.335	0.657	<b>0.918</b>	0.878	0.756	0.603	0.895	0.683	0.889	0.679	0.776	<b>0.877</b>	0.664	0.874
hockey	<b>0.775</b>	0.080	0.538	0.810	<b>0.817</b>	0.674	0.713	0.515	0.566	0.468	0.413	<b>0.714</b>	0.698	0.677	0.647
horsejump-high	<b>0.649</b>	0.263	0.596	<b>0.834</b>	0.830	0.364	0.734	0.370	0.568	0.578	0.236	0.637	<b>0.765</b>	0.586	0.676
horsejump-low	0.545	0.125	<b>0.618</b>	0.651	<b>0.743</b>	0.705	0.682	0.630	0.388	0.526	0.291	0.498	0.551	<b>0.663</b>	0.607
kite-surf	<b>0.654</b>	0.059	0.208	0.453	0.357	0.501	0.419	<b>0.585</b>	0.193	0.272	0.366	0.487	0.405	0.500	<b>0.577</b>
kite-walk	<b>0.736</b>	0.668	0.420	<b>0.813</b>	0.447	0.052	0.597	0.197	0.725	0.649	0.447	0.498	<b>0.765</b>	0.509	0.682
libby	<b>0.655</b>	0.097	0.443	<b>0.635</b>	0.169	0.073	0.050	0.611	0.470	0.507	0.070	0.226	<b>0.553</b>	0.295	0.316
lucia	<b>0.820</b>	0.119	0.760	<b>0.876</b>	0.840	0.669	0.417	0.847	0.706	0.644	0.377	0.626	0.776	<b>0.836</b>	0.801
mallard-fly	<b>0.799</b>	0.022	0.310	<b>0.617</b>	0.380	0.293	0.033	0.585	0.227	0.601	0.200	<b>0.557</b>	0.436	0.536	0.541
mallard-water	<b>0.755</b>	0.035	0.017	0.761	0.245	0.190	0.045	<b>0.785</b>	0.085	0.087	0.623	<b>0.865</b>	0.704	0.751	0.687
motocross-bumps	<b>0.827</b>	0.236	0.460	0.614	0.603	0.502	0.466	<b>0.689</b>	0.351	0.617	0.133	0.470	0.534	<b>0.761</b>	0.306
motocross-jump	<b>0.760</b>	0.204	0.428	0.251	0.245	0.338	<b>0.618</b>	0.288	0.491	0.602	0.123	0.386	0.099	<b>0.583</b>	0.511
motorbike	<b>0.688</b>	0.116	0.572	0.714	0.387	0.723	<b>0.737</b>	0.572	0.335	0.558	0.340	0.451	0.687	0.506	<b>0.713</b>
paragliding	<b>0.877</b>	0.140	0.743	0.880	0.890	0.816	<b>0.933</b>	0.861	0.568	0.725	0.735	0.863	0.907	<b>0.951</b>	0.866
paragliding-launch	<b>0.599</b>	0.255	0.501	<b>0.628</b>	0.591	0.555	0.513	0.559	0.539	0.506	0.301	0.577	0.537	<b>0.589</b>	0.571
parkour	<b>0.815</b>	0.288	0.491	<b>0.901</b>	0.146	0.345	0.295	0.410	0.392	0.458	0.070	0.121	0.240	<b>0.342</b>	0.322
rhino	<b>0.864</b>	0.253	0.615	0.682	0.520	<b>0.846</b>	<b>0.902</b>	0.675	0.685	0.776	0.694	0.736	<b>0.812</b>	0.716	0.794
rollerblade	0.554	0.004	<b>0.564</b>	<b>0.814</b>	0.406	0.566	0.801	0.510	0.141	0.318	0.098	0.138	0.461	<b>0.726</b>	0.449
scooter-black	<b>0.704</b>	0.125	0.613	0.162	<b>0.759</b>	0.435	0.579	0.502	0.348	0.522	0.378	<b>0.793</b>	0.624	0.626	0.504
scooter-gray	0.653	0.039	<b>0.703</b>	<b>0.586</b>	0.327	0.357	0.345	0.363	0.421	0.325	0.133	0.241	0.433	0.123	<b>0.483</b>
soapbox	<b>0.680</b>	0.163	0.579	0.634	<b>0.832</b>	0.294	0.672	0.757	0.332	0.410	0.247	<b>0.783</b>	0.684	0.759	0.449
soccerball	<b>0.856</b>	0.052	0.732	0.829	0.242	0.350	0.370	<b>0.878</b>	0.378	0.843	0.029	0.653	0.065	0.096	<b>0.820</b>
stroller	<b>0.600</b>	0.313	0.573	<b>0.850</b>	0.619	0.720	0.678	0.759	0.466	0.580	0.369	0.464	<b>0.662</b>	0.656	0.597
surf	<b>0.944</b>	0.521	0.653	0.775	0.273	0.464	0.770	<b>0.893</b>	0.312	0.475	0.814	0.821	0.759	<b>0.941</b>	0.843
swing	<b>0.709</b>	0.387	0.677	<b>0.851</b>	0.533	0.413	0.622	0.710	0.569	0.431	0.098	0.511	0.104	0.115	<b>0.648</b>
tennis	<b>0.714</b>	0.032	0.562	<b>0.871</b>	0.494	0.196	0.590	0.762	0.480	0.388	0.074	0.482	0.576	<b>0.765</b>	0.623
train	0.535	0.205	<b>0.548</b>	0.729	<b>0.903</b>	0.876	0.887	0.450	0.620	0.831	0.648	0.854	0.846	<b>0.873</b>	0.841
Mean	<b>0.724</b>	0.173	0.532	<b>0.641</b>	0.514	0.501	0.543	0.569	0.426	0.575	0.358	0.556	0.596	0.607	<b>0.631</b>

Table 4: Results of region similarity ( $\mathcal{J}$ ) for each video sequence in the dataset. The best performing method of each category is highlighted in bold.

Sequence	Preprocessing			Unsupervised							Semi-Supervised				
	MCG	SF-LAB	SF-MOT	NLC	CVOS	TRC	MSG	KEY	SAL	FST	TSP	SEA	HVS	JMP	FCP
bear	<b>0.934</b>	0.180	0.451	0.850	0.845	0.832	0.781	0.775	0.495	<b>0.860</b>	0.635	0.899	<b>0.905</b>	0.904	0.845
blackswan	<b>0.873</b>	0.491	0.500	<b>0.820</b>	0.695	0.654	0.700	0.787	0.430	0.736	0.857	<b>0.957</b>	0.910	0.945	0.905
bmx-bumps	<b>0.597</b>	0.106	0.283	<b>0.734</b>	0.409	0.325	0.410	0.453	0.313	0.349	0.338	0.254	<b>0.525</b>	0.397	0.340
bmx-trees	<b>0.605</b>	0.145	0.557	0.330	0.118	0.189	0.263	<b>0.366</b>	0.206	0.348	0.138	0.125	0.282	0.309	<b>0.324</b>
boat	<b>0.536</b>	0.289	0.135	0.036	0.108	0.403	<b>0.485</b>	0.000	0.264	0.197	0.682	0.764	<b>0.807</b>	0.607	0.460
breakdance	<b>0.670</b>	0.265	0.635	<b>0.661</b>	0.191	0.121	0.231	0.463	0.300	0.411	0.070	0.389	0.473	<b>0.511</b>	0.473
breakdance-flare	<b>0.781</b>	0.215	0.626	<b>0.808</b>	0.335	0.301	0.230	0.585	0.512	0.694	0.116	0.167	0.625	0.523	<b>0.738</b>
bus	0.497	0.296	<b>0.505</b>	0.406	0.535	0.542	<b>0.657</b>	0.635	0.570	0.584	0.477	<b>0.724</b>	0.682	0.604	0.539
camel	<b>0.728</b>	0.085	0.498	0.719	<b>0.873</b>	0.698	0.629	0.437	0.432	0.590	0.529	0.614	<b>0.871</b>	0.711	0.617
car-roundabout	<b>0.512</b>	0.246	0.491	0.250	<b>0.678</b>	0.451	0.602	0.362	0.301	0.625	0.435	<b>0.710</b>	0.551	0.619	0.478
car-shadow	0.588	0.233	<b>0.628</b>	0.546	0.617	0.474	<b>0.858</b>	0.459	0.441	0.540	0.513	<b>0.755</b>	0.594	0.625	0.642
car-turn	<b>0.760</b>	0.113	0.431	0.634	0.703	<b>0.741</b>	0.677	0.632	0.485	0.731	0.379	<b>0.883</b>	0.605	0.742	0.614
cows	<b>0.736</b>	0.168	0.554	<b>0.807</b>	0.499	0.721	0.621	0.293	0.499	0.681	0.544	0.677	0.632	<b>0.700</b>	0.667
dance-jump	0.372	0.204	<b>0.497</b>	0.567	0.282	0.272	0.038	<b>0.569</b>	0.262	0.462	0.186	0.567	<b>0.571</b>	0.526	0.418
dance-twirl	<b>0.604</b>	0.121	0.589	0.365	0.444	0.376	0.325	0.317	0.301	<b>0.471</b>	0.128	0.213	0.516	<b>0.520</b>	0.427
dog	<b>0.582</b>	0.140	0.495	0.707	<b>0.761</b>	0.695	0.304	0.633	0.418	0.659	0.295	0.543	0.635	0.596	<b>0.672</b>
dog-agility	<b>0.517</b>	0.153	0.288	<b>0.551</b>	0.262	0.122	0.076	0.095	0.102	0.265	0.083	0.410	0.446	<b>0.654</b>	0.315
drift-chicane	<b>0.906</b>	0.154	0.364	0.312	0.397	0.823	<b>0.886</b>	0.192	0.206	0.731	0.033	0.159	<b>0.547</b>	0.338	0.477
drift-straight	<b>0.593</b>	0.152	0.312	0.385	0.330	0.408	<b>0.509</b>	0.053	0.167	0.470	0.213	<b>0.500</b>	0.266	0.473	0.479
drift-turn	<b>0.700</b>	0.207	0.199	0.185	<b>0.480</b>	0.310	0.459	0.018	0.231	0.442	0.217	0.512	0.216	<b>0.631</b>	0.488
elephant	<b>0.506</b>	0.080	0.446	0.251	0.359	0.546	0.505	0.324	0.231	<b>0.569</b>	0.523	0.399	<b>0.579</b>	0.542	0.430
flamingo	<b>0.876</b>	0.253	0.622	0.610	<b>0.806</b>	0.663	0.776	0.589	0.621	0.763	0.544	0.563	<b>0.790</b>	0.650	0.641
goat	<b>0.526</b>	0.171	0.225	0.133	0.241	<b>0.724</b>	0.657	0.552	0.187	0.400	0.404	0.470	0.546	<b>0.617</b>	0.576
hike	<b>0.926</b>	0.312	0.540	<b>0.943</b>	0.922	0.804	0.702	0.925	0.691	0.918	0.675	0.796	0.878	0.744	<b>0.912</b>
hockey	<b>0.742</b>	0.198	0.430	<b>0.808</b>	0.789	0.651	0.761	0.560	0.559	0.584	0.579	0.721	<b>0.778</b>	0.726	0.612
horsejump-high	<b>0.703</b>	0.326	0.561	<b>0.881</b>	0.841	0.405	0.748	0.392	0.613	0.621	0.343	0.655	<b>0.807</b>	0.653	0.699
horsejump-low	<b>0.548</b>	0.152	0.516	0.659	<b>0.709</b>	0.672	0.637	0.533	0.419	0.490	0.356	0.548	0.572	<b>0.696</b>	0.533
kite-surf	<b>0.447</b>	0.221	0.286	0.448	0.241	0.422	<b>0.521</b>	0.504	0.368	0.346	0.268	0.285	<b>0.375</b>	0.309	0.362
kite-walk	0.486	<b>0.524</b>	0.285	<b>0.662</b>	0.438	0.014	0.577	0.128	0.526	0.561	0.435	0.355	<b>0.624</b>	0.359	0.411
libby	<b>0.733</b>	0.244	0.581	<b>0.748</b>	0.185	0.086	0.118	0.730	0.529	0.718	0.091	0.209	<b>0.641</b>	0.365	0.389
lucia	<b>0.743</b>	0.233	0.736	<b>0.872</b>	0.801	0.663	0.491	0.819	0.691	0.568	0.453	0.542	0.782	<b>0.870</b>	0.708
mallard-fly	<b>0.824</b>	0.071	0.338	<b>0.660</b>	0.391	0.332	0.019	0.631	0.293	0.633	0.235	<b>0.607</b>	0.441	0.579	0.539
mallard-water	<b>0.701</b>	0.115	0.034	0.692	0.254	0.225	0.000	<b>0.733</b>	0.115	0.079	0.585	<b>0.886</b>	0.646	0.755	0.557
motocross-bumps	<b>0.710</b>	0.242	0.338	0.560	0.567	0.497	0.466	<b>0.674</b>	0.300	0.610	0.184	0.520	0.548	<b>0.743</b>	0.302
motocross-jump	<b>0.568</b>	0.274	0.290	0.303	0.186	0.307	0.393	0.237	0.388	<b>0.453</b>	0.116	0.404	0.137	<b>0.539</b>	0.386
motorbike	<b>0.642</b>	0.247	0.264	0.571	0.380	0.541	0.594	<b>0.726</b>	0.391	0.585	0.406	0.481	<b>0.823</b>	0.578	0.632
paragliding	0.742	0.133	<b>0.798</b>	0.744	0.744	0.724	<b>0.909</b>	0.681	0.541	0.675	0.634	0.744	0.857	<b>0.907</b>	0.727
paragliding-launch	<b>0.200</b>	0.188	0.171	0.243	0.182	0.157	0.196	<b>0.253</b>	0.169	0.185	0.122	0.180	<b>0.206</b>	0.176	0.183
parkour	<b>0.811</b>	0.304	0.553	<b>0.916</b>	0.158	0.421	0.401	0.374	0.359	0.478	0.094	0.278	0.323	<b>0.418</b>	0.292
rhino	<b>0.767</b>	0.193	0.472	0.431	0.469	0.739	<b>0.826</b>	0.429	0.487	0.634	0.499	0.658	<b>0.658</b>	0.653	0.647
rollerblade	0.577	0.114	<b>0.586</b>	<b>0.868</b>	0.475	0.687	0.822	0.351	0.211	0.411	0.143	0.155	0.552	<b>0.759</b>	0.576
scooter-black	<b>0.548</b>	0.177	0.451	0.228	0.557	0.304	<b>0.565</b>	0.420	0.257	0.395	0.411	<b>0.722</b>	0.574	0.529	0.363
scooter-gray	<b>0.494</b>	0.126	0.434	<b>0.467</b>	0.212	0.266	0.272	0.367	0.333	0.321	0.122	0.275	<b>0.545</b>	0.123	0.437
soapbox	<b>0.647</b>	0.190	0.453	0.658	<b>0.754</b>	0.389	0.633	0.719	0.307	0.355	0.336	<b>0.750</b>	0.690	0.677	0.423
soccerball	<b>0.898</b>	0.043	0.754	0.855	0.262	0.377	0.401	<b>0.924</b>	0.355	0.900	0.059	0.697	0.074	0.127	<b>0.836</b>
stroller	<b>0.635</b>	0.426	0.488	<b>0.874</b>	0.606	0.691	0.662	0.751	0.417	0.558	0.404	0.525	0.708	<b>0.718</b>	0.581
surf	<b>0.881</b>	0.524	0.372	0.673	0.515	0.637	0.804	<b>0.820</b>	0.395	0.445	0.641	0.732	0.652	<b>0.872</b>	0.713
swing	<b>0.599</b>	0.373	0.571	<b>0.778</b>	0.493	0.417	0.611	0.614	0.502	0.491	0.087	0.409	0.091	0.109	<b>0.538</b>
tennis	<b>0.751</b>	0.151	0.450	<b>0.927</b>	0.547	0.301	0.670	0.818	0.530	0.567	0.114	0.537	0.579	<b>0.818</b>	0.652
train	0.403	0.347	<b>0.504</b>	0.521	<b>0.831</b>	0.766	0.770	0.464	0.440	0.660	0.589	0.713	0.688	<b>0.770</b>	0.736
Mean	<b>0.654</b>	0.218	0.452	<b>0.593</b>	0.490	0.478	0.525	0.503	0.383	0.536	0.346	0.533	0.576	<b>0.586</b>	0.546

Table 5: Results of boundary precision ( $\mathcal{F}$ ) for each video sequence in the dataset. The best performing method of each category is highlighted in bold.

Sequence	Preprocessing				Unsupervised						Semi-Supervised				
	MCG	SF-LAB	SF-MOT	NLC	CVOS	TRC	MSG	KEY	SAL	FST	TSP	SEA	HVS	JMP	FCP
bear	<b>0.222</b>	0.633	0.840	0.151	<b>0.059</b>	0.272	0.156	0.068	0.448	0.227	0.077	<b>0.047</b>	0.086	0.051	0.114
blackswan	<b>0.309</b>	0.332	0.322	0.110	0.058	0.219	0.145	<b>0.049</b>	0.660	0.225	0.049	0.032	0.060	<b>0.029</b>	0.064
boat	0.813	<b>0.437</b>	0.971	0.559	1.159	0.350	0.163	<b>0.015</b>	0.382	0.177	0.067	<b>0.055</b>	0.125	0.062	0.136
bus	0.741	0.757	<b>0.397</b>	0.178	0.146	0.194	0.154	<b>0.143</b>	0.369	0.270	0.293	<b>0.109</b>	0.306	0.193	0.156
camel	0.593	0.438	<b>0.341</b>	0.232	<b>0.123</b>	0.172	0.129	0.138	0.380	0.161	0.084	<b>0.055</b>	0.117	0.062	0.212
car-roundabout	0.545	1.178	<b>0.428</b>	0.352	<b>0.064</b>	0.382	0.291	0.161	0.536	0.242	0.158	<b>0.071</b>	0.255	0.078	0.283
car-shadow	1.162	1.270	<b>0.227</b>	0.361	<b>0.180</b>	0.452	0.206	0.313	0.793	0.353	<b>0.206</b>	0.230	0.351	0.274	0.339
car-turn	<b>0.274</b>	0.714	0.477	0.235	0.118	0.202	0.204	<b>0.108</b>	0.566	0.214	0.796	0.117	0.135	<b>0.065</b>	0.256
cows	0.494	1.022	<b>0.457</b>	0.147	<b>0.133</b>	0.148	0.196	0.412	0.511	0.282	0.179	<b>0.044</b>	0.164	0.055	0.163
dance-jump	1.083	<b>0.630</b>	0.870	0.316	0.459	0.576	<b>0.110</b>	0.214	0.586	0.242	0.272	0.286	0.324	<b>0.173</b>	0.506
drift-straight	0.969	1.170	<b>0.775</b>	0.599	0.900	0.638	0.543	<b>0.292</b>	0.950	0.482	0.826	0.396	0.823	<b>0.317</b>	0.597
drift-turn	<b>0.391</b>	0.987	0.984	0.850	0.334	0.475	0.402	<b>0.150</b>	1.002	0.258	0.633	0.168	0.703	<b>0.128</b>	0.328
elephant	1.103	0.960	<b>0.557</b>	0.315	0.118	0.236	0.236	<b>0.085</b>	0.426	0.139	0.097	0.076	0.213	<b>0.075</b>	0.404
flamingo	0.523	<b>0.319</b>	0.769	0.138	0.173	0.215	0.382	<b>0.113</b>	0.486	0.175	0.118	<b>0.069</b>	0.133	0.089	0.182
hike	0.331	0.873	<b>0.282</b>	0.158	0.125	0.230	0.251	<b>0.117</b>	0.412	0.247	0.141	0.122	0.120	<b>0.092</b>	0.164
hockey	0.519	1.052	<b>0.378</b>	0.227	<b>0.159</b>	0.228	0.211	0.162	0.377	0.276	0.403	0.103	0.258	<b>0.102</b>	0.228
kite-surf	<b>0.420</b>	0.464	0.888	0.944	0.248	0.432	0.507	<b>0.233</b>	0.568	0.404	0.278	0.125	0.497	<b>0.117</b>	0.212
kite-walk	0.409	<b>0.322</b>	0.829	0.221	0.127	<b>0.002</b>	0.328	0.366	0.356	0.301	0.241	0.173	0.185	<b>0.154</b>	0.166
mallard-water	0.815	<b>0.579</b>	1.579	0.242	0.394	0.641	<b>0.000</b>	0.184	1.070	0.230	0.287	<b>0.123</b>	0.295	0.219	0.317
motocross-bumps	<b>0.606</b>	0.756	0.919	0.541	<b>0.327</b>	0.566	0.481	0.344	0.903	0.329	0.628	0.289	0.767	<b>0.211</b>	0.486
paragliding-launch	0.498	0.598	<b>0.324</b>	0.259	0.273	0.347	0.331	<b>0.213</b>	0.602	0.703	0.660	0.208	0.316	<b>0.180</b>	0.329
rhino	<b>0.308</b>	0.396	0.700	0.188	0.064	0.153	0.093	<b>0.056</b>	0.390	0.138	0.066	<b>0.037</b>	0.093	0.037	0.151
scooter-black	1.020	0.910	<b>0.419</b>	0.760	<b>0.320</b>	0.577	0.364	0.558	0.790	0.475	0.960	<b>0.216</b>	0.599	0.282	0.423
soapbox	0.983	0.613	<b>0.504</b>	0.390	<b>0.154</b>	0.413	0.214	0.160	0.613	0.158	0.763	<b>0.112</b>	0.314	0.126	0.379
stroller	1.195	0.968	<b>0.864</b>	0.205	<b>0.116</b>	0.235	0.366	0.127	0.546	0.184	0.134	<b>0.130</b>	0.363	0.189	0.376
surf	<b>0.234</b>	1.338	0.728	0.364	0.169	0.375	0.223	<b>0.086</b>	1.093	0.398	0.207	0.238	0.291	<b>0.127</b>	0.276
train	1.038	0.739	<b>0.377</b>	0.576	<b>0.056</b>	0.110	0.070	0.270	0.396	0.159	0.249	0.069	0.106	<b>0.047</b>	0.448
Mean	0.652	0.758	<b>0.637</b>	0.356	0.243	0.327	0.250	<b>0.190</b>	0.600	0.276	0.329	0.137	0.296	<b>0.131</b>	0.285

Table 6: Results of temporal stability ( $\mathcal{T}$ ) for each video sequence in the dataset. The best performing method of each category is highlighted in bold. Please note that this measure is only computed on those sequences without occlusions and strong deformations.

Attr	Unsupervised							Semi-Supervised				
	NLC	CVOS	TRC	MSG	KEY	SAL	FST	TSP	SEA	HVS	JMP	FCOP
LR	<b>0.65 +0.03</b>	0.40 +0.14	0.45 +0.05	0.48 +0.06	0.53 +0.05	0.32 +0.13	0.52 +0.06	0.29 +0.08	0.46 +0.11	0.46 +0.16	0.50 +0.13	<b>0.58 +0.06</b>
SV	<b>0.52 +0.16</b>	0.42 +0.12	0.42 +0.11	0.50 +0.05	0.49 +0.10	0.33 +0.14	0.49 +0.11	0.23 +0.18	0.48 +0.09	0.44 +0.21	<b>0.57 +0.04</b>	0.51 +0.16
SC	<b>0.59 +0.06</b>	0.49 +0.02	0.46 +0.06	0.52 +0.01	0.50 +0.11	0.45 -0.05	0.51 +0.09	0.32 +0.05	0.50 +0.09	0.56 +0.05	0.51 +0.14	<b>0.58 +0.07</b>
FM	<b>0.62 +0.01</b>	0.36 +0.25	0.40 +0.16	0.44 +0.15	0.49 +0.12	0.34 +0.13	0.49 +0.13	0.17 +0.31	0.39 +0.28	0.40 +0.31	0.49 +0.18	<b>0.54 +0.13</b>
CS	<b>0.59 +0.06</b>	0.41 +0.12	0.54 -0.06	0.53 -0.00	0.51 +0.06	0.36 +0.09	0.53 +0.05	0.34 +0.01	0.42 +0.17	0.54 +0.06	0.60 +0.00	<b>0.60 +0.03</b>
IO	<b>0.61 +0.03</b>	0.54 -0.06	0.46 +0.05	0.57 -0.07	0.53 +0.05	0.41 +0.02	0.48 +0.16	0.34 +0.03	0.53 +0.03	0.55 +0.06	0.58 +0.04	<b>0.59 +0.07</b>
DB	<b>0.52 +0.14</b>	0.37 +0.18	0.38 +0.14	0.43 +0.14	0.51 +0.07	0.34 +0.10	0.52 +0.06	0.39 -0.06	0.57 -0.03	0.59 -0.01	0.59 +0.01	<b>0.61 +0.01</b>
MB	<b>0.59 +0.05</b>	0.35 +0.23	0.31 +0.28	0.33 +0.30	0.50 +0.08	0.32 +0.15	0.47 +0.15	0.14 +0.32	0.39 +0.24	0.42 +0.25	0.50 +0.15	<b>0.52 +0.15</b>
DEF	<b>0.67 -0.11</b>	0.51 -0.01	0.48 +0.01	0.51 +0.05	0.57 -0.01	0.45 -0.08	0.56 +0.01	0.31 +0.10	0.49 +0.14	0.58 +0.00	0.58 +0.03	<b>0.60 +0.04</b>
OCC	<b>0.68 -0.08</b>	0.42 +0.13	0.42 +0.11	0.46 +0.11	0.51 +0.08	0.43 -0.02	0.52 +0.08	0.26 +0.14	0.46 +0.13	0.51 +0.12	0.46 +0.21	<b>0.58 +0.07</b>
HO	<b>0.64 -0.04</b>	0.49 +0.05	0.45 +0.14	0.53 +0.01	0.53 +0.12	0.42 -0.01	0.54 +0.10	0.27 +0.29	0.49 +0.23	0.53 +0.21	0.55 +0.17	<b>0.59 +0.13</b>
EA	0.50 +0.24	0.40 +0.19	0.45 +0.07	0.45 +0.16	0.48 +0.15	0.35 +0.12	<b>0.51 +0.10</b>	0.31 +0.06	0.51 +0.07	0.53 +0.10	0.54 +0.09	<b>0.57 +0.10</b>
OV	<b>0.49 +0.17</b>	0.34 +0.20	0.32 +0.21	0.39 +0.17	0.42 +0.18	0.30 +0.15	0.48 +0.10	0.20 +0.19	0.43 +0.14	0.39 +0.24	<b>0.59 +0.00</b>	0.51 +0.13
AC	0.53 +0.13	0.41 +0.12	0.36 +0.17	0.47 +0.08	0.41 +0.19	0.33 +0.12	<b>0.54 +0.04</b>	0.17 +0.23	0.45 +0.12	0.41 +0.22	<b>0.57 +0.04</b>	0.50 +0.16
BC	0.45 +0.22	0.45 +0.06	0.50 -0.02	0.54 -0.01	0.52 +0.05	0.42 +0.00	<b>0.57 -0.00</b>	0.41 -0.07	0.58 -0.04	<b>0.61 -0.04</b>	0.60 -0.00	0.58 +0.05

Table 7: Attribute-based aggregate performance. For each method, the respective left column corresponds to the average region similarity  $\mathcal{J}$  over all sequences with that specific attribute (e.g., AC), while the right column indicates the performance gain (or loss) for that method for the remaining sequences without that respective attribute. The best performing method of each category is highlighted in bold.

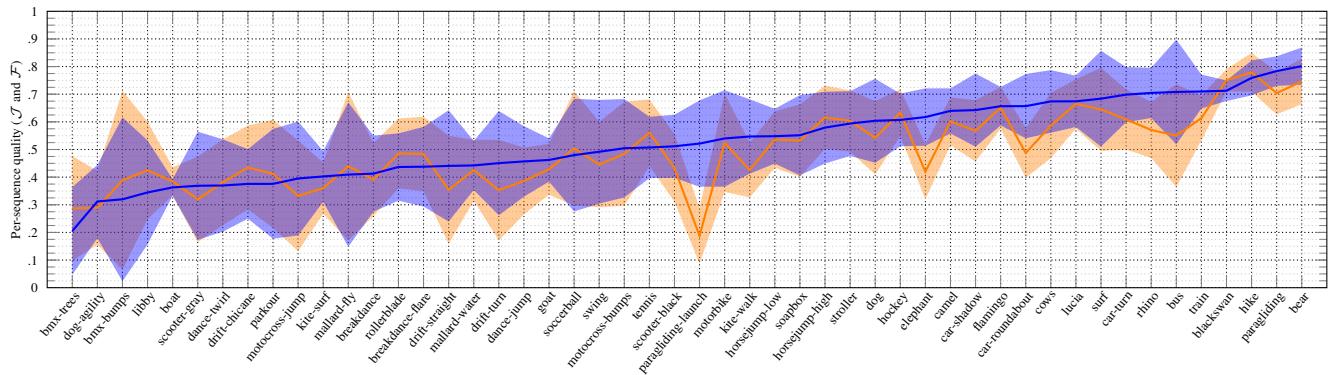
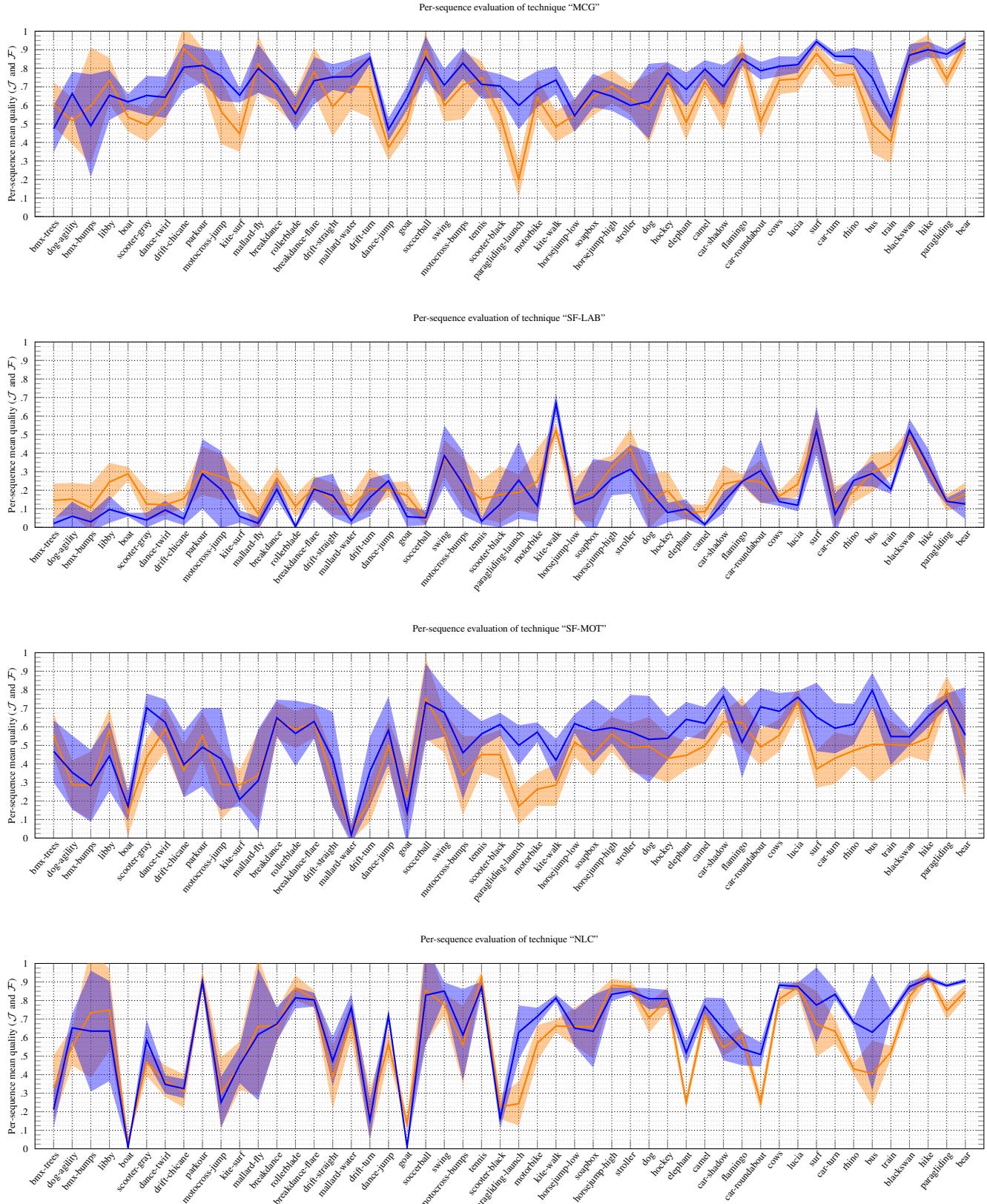
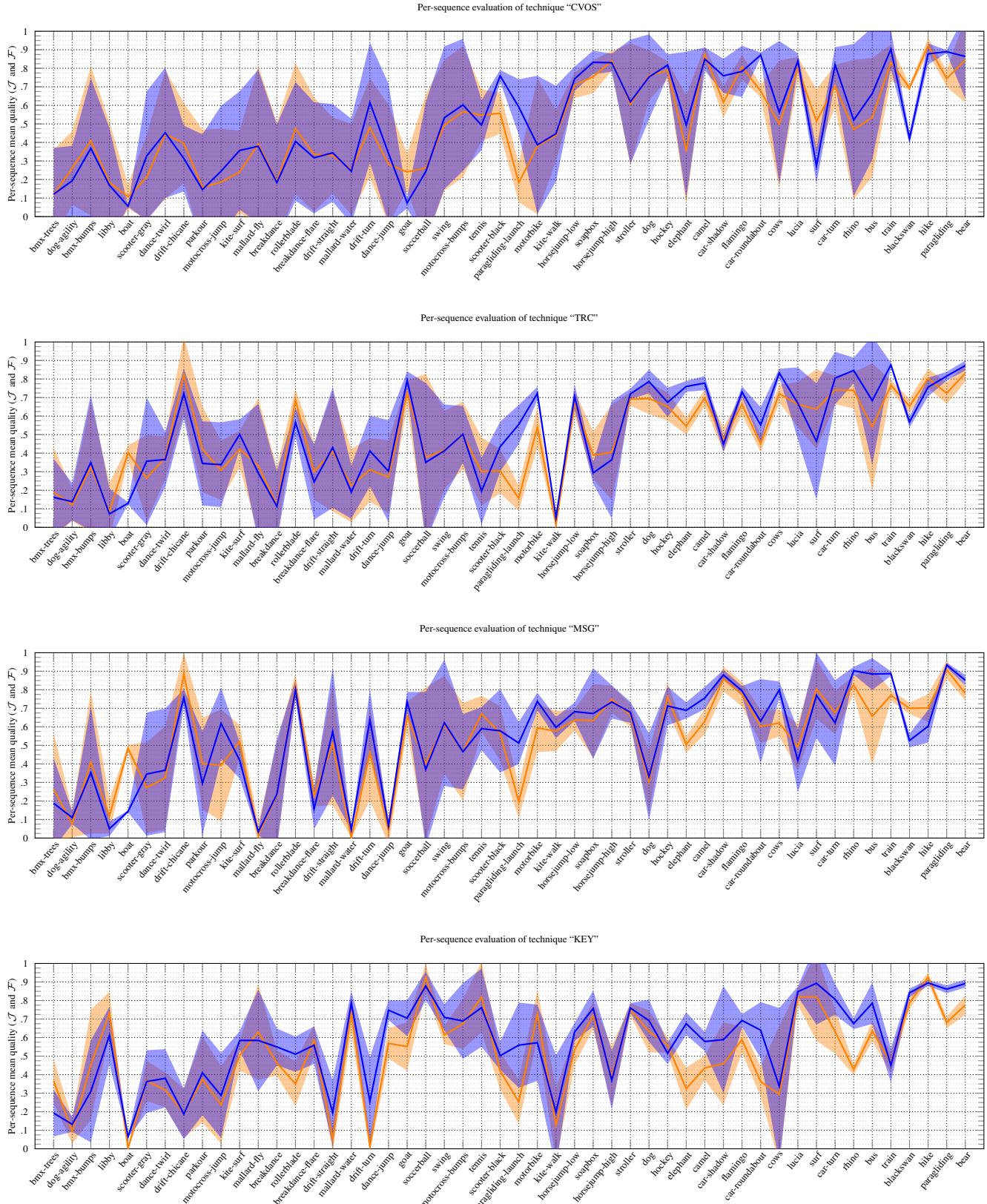


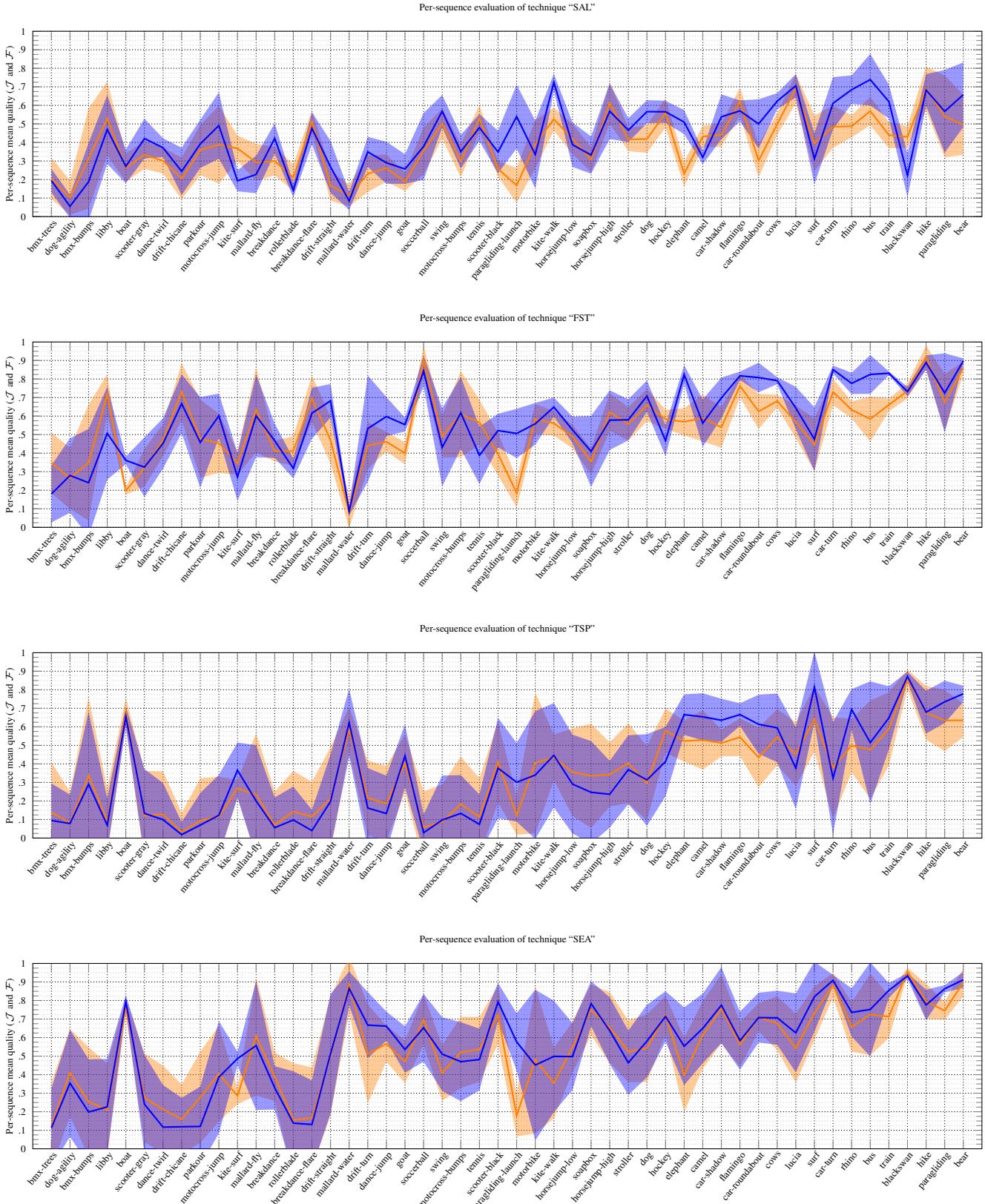
Figure 2: **Per-sequence mean performance:** Mean and variance of region Jaccard  $\mathcal{J}$  (blue) and boundary F measure  $\mathcal{F}$  (orange). Sequences are sorted by *difficulty*, *i.e.* mean performance of  $\mathcal{J}$  over all techniques.



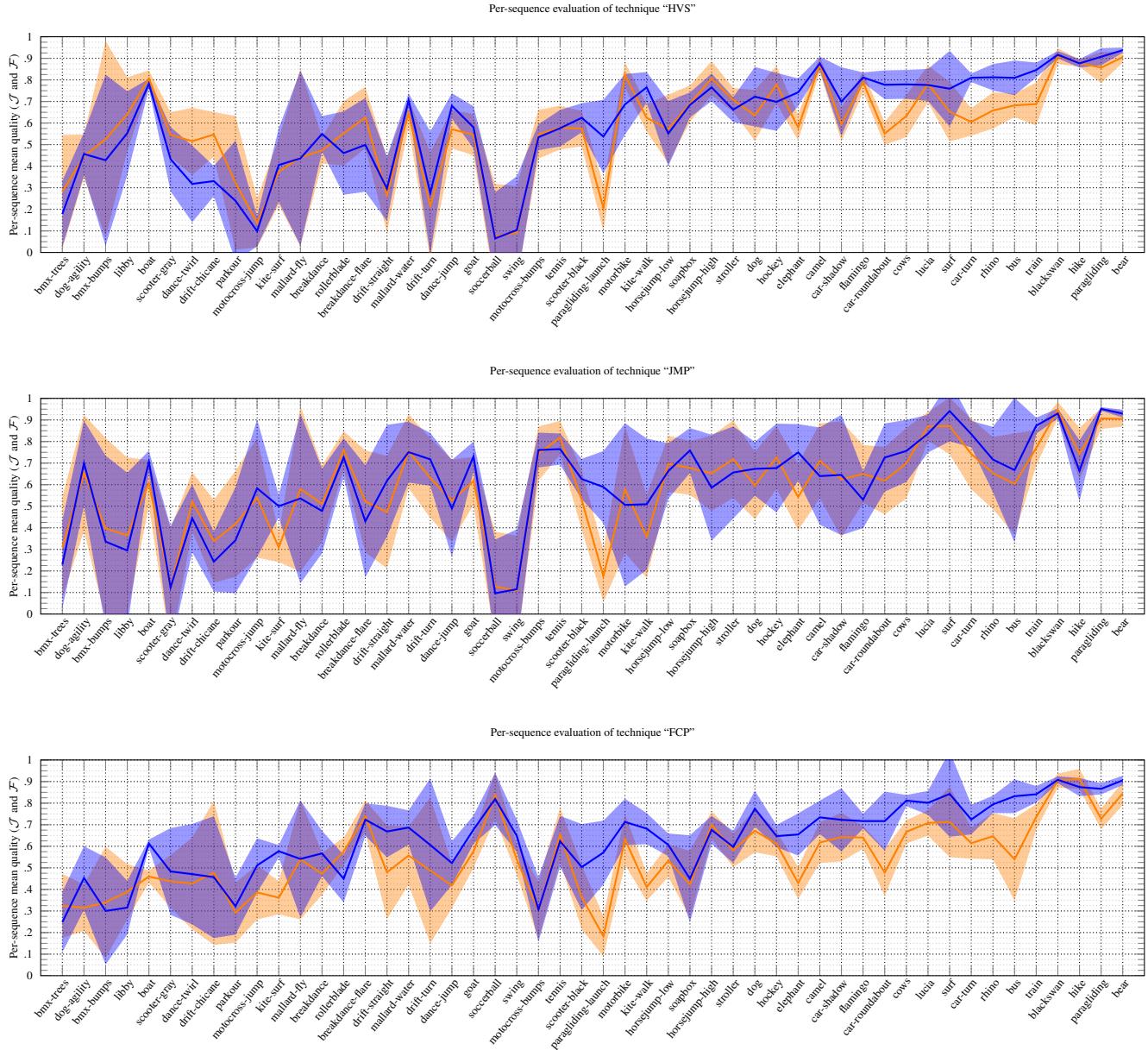
**Figure 3: Per-sequence per-method evaluation:** Mean and variance of region Jaccard  $\mathcal{J}$  (blue) and boundary F measure  $\mathcal{F}$  (orange)



**Figure 4: Per-sequence per-method evaluation:** Mean and variance of region Jaccard  $\mathcal{J}$  (blue) and boundary F measure  $\mathcal{F}$  (orange)



**Figure 5: Per-sequence per-method evaluation:** Mean and variance of region Jaccard  $\mathcal{J}$  (blue) and boundary F measure  $\mathcal{F}$  (orange)



**Figure 6: Per-sequence per-method evaluation:** Mean and variance of region Jaccard  $\mathcal{J}$  (blue) and boundary F measure  $\mathcal{F}$  (orange)