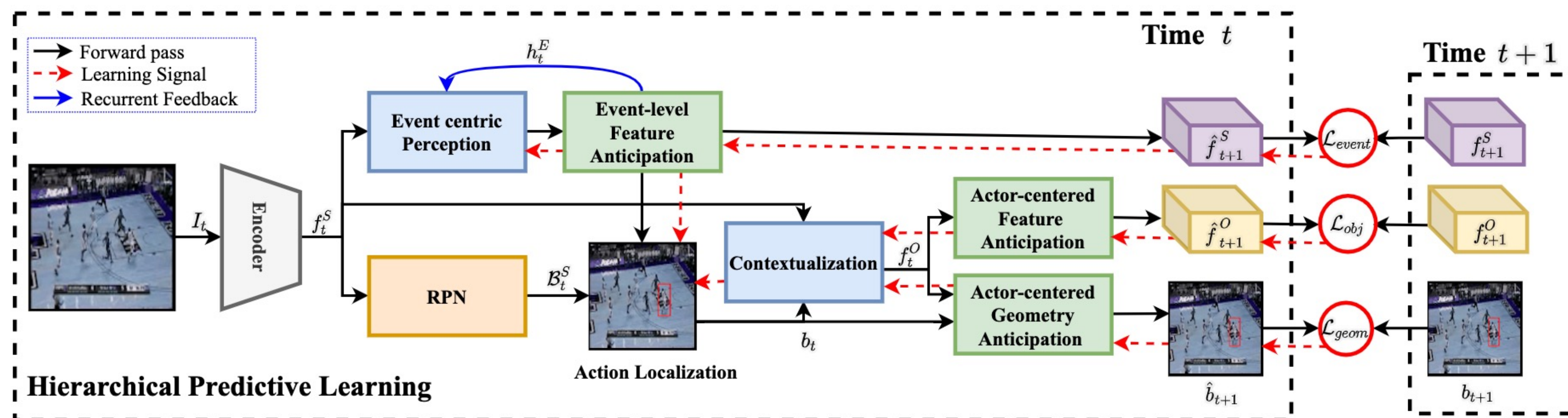


# Actor-centered Representations for Action Localization in Streaming Videos

## Motivation

- Event perception tasks such as action recognition and localization are important for visual understanding
- Progress has largely been driven by the use of large-scale, annotated training data in a supervised manner.
- **Goal:** Can we learn robust representations for video sequences for localizing the action in streaming videos?

## Overall Framework



## Quantitative Evaluation

Data: UCF Sports, THUMOS'13, JHMDB Metrics: mAP@IoU, Recall

**Comparison to state-of-the-art approaches on action localization across different supervision needs, datasets and IOU thresholds**

Approach	Supervision		UCF Sports		JHMDB		THUMOS'13	
	Spatial	Label	$\sigma=0.2$	$\sigma=0.5$	$\sigma=0.2$	$\sigma=0.5$	$\sigma=0.2$	$\sigma=0.5$
Tube CNN [12]	✓	✓	0.47	-	-	0.77	0.47	0.41
Action Tubelets [14]	✓	✓	0.53	0.27	-	-	0.48	-
Action Tubes [9]	✓	✓	0.56	0.49	0.55	0.45	-	-
MRSTL [48]	✓	✓	-	-	-	0.37	-	0.68
MENET [24]	✓	✓	-	-	-	<b>0.82</b>	-	<b>0.84</b>
HISAN [27]	✓	✓	-	-	-	0.77	-	0.73
ACAR-Net [26]	✓	✓	-	-	-	-	-	0.84
ALSTM [32]	✗	✓	-	-	-	-	0.06	-
VideoLSTM [21]	✗	✓	-	-	-	-	0.37	-
Actor Supervision [6]	✗	✓	-	<b>0.48</b>	-	<b>0.36</b>	<b>0.46</b>	-
Soomro <i>et al</i> [37]	✗	✗	0.46*	0.30*	<b>0.43*</b>	<b>0.22*</b>	0.21*	0.06*
PredLearn [3] ( $k=k_{gt}$ )	✗	✗	0.55	0.32	0.30	0.10	0.31	0.10
AC-HPL (Ours, $k=k_{gt}$ )	✗	✗	<b>0.70</b>	<b>0.59</b>	<b>0.43</b>	0.15	<b>0.38</b>	<b>0.20</b>

**Generalization to outside domain training data without finetuning**

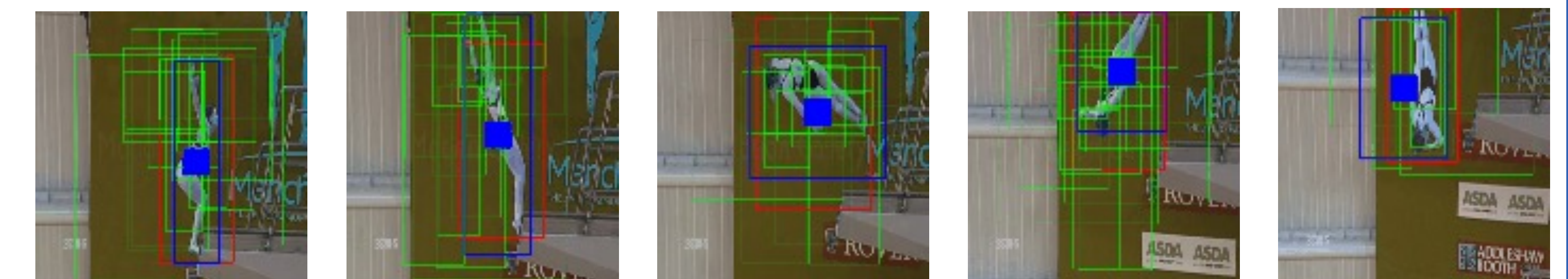
Test Data →	UCF Sports		JHMDB		THUMOS'13	
	AC-HPL	PredLearn	AC-HPL	PredLearn	AC-HPL	PredLearn
Train Data ↓	$\sigma=0.5$		$\sigma=0.2$		$\sigma=0.2$	
UCF Sports	<b>0.59</b>	0.32	<b>0.39</b>	0.19	<b>0.38</b>	0.20
JHMDB	<b>0.48</b>	0.23	<b>0.43</b>	0.30	<b>0.35</b>	0.26
THUMOS'13	<b>0.50</b>	0.27	<b>0.40</b>	0.24	<b>0.38</b>	0.31

## Contributions

- Introduce *hierarchical predictive learning* for unsupervised action localization in *streaming* videos
- Introduce a novel, attention-driven formulation for learning robust, actor-centered event features
- Demonstrate extension to multi-actor group activity recognition and localization and generalization to data outside the training domain *without finetuning*

## Qualitative Results

**Generic Single Actor localization**



**Multi-Actor localization**

