

Consistency-based Self-supervised Learning for Temporal Anomaly Localization

European Conference on Computer Vision Workshops (ECCVW) – Tel Aviv, 2022

Aniello Panariello

Angelo Porrello

Simone Calderara

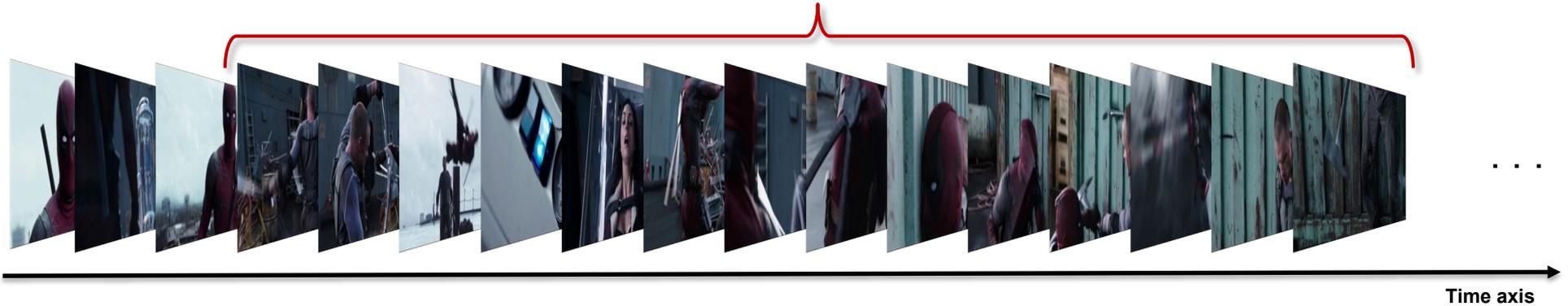
Rita Cucchiara



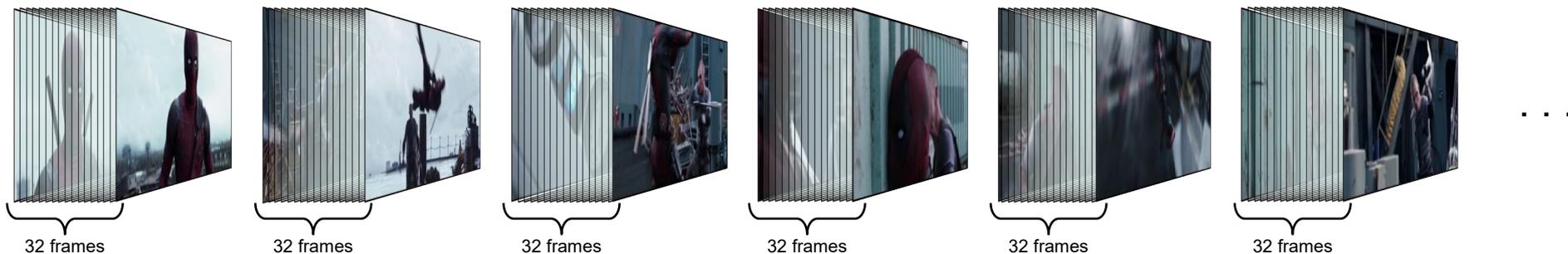
AlmageLab

University of Modena and Reggio Emilia, Italy

Anomalous Frames



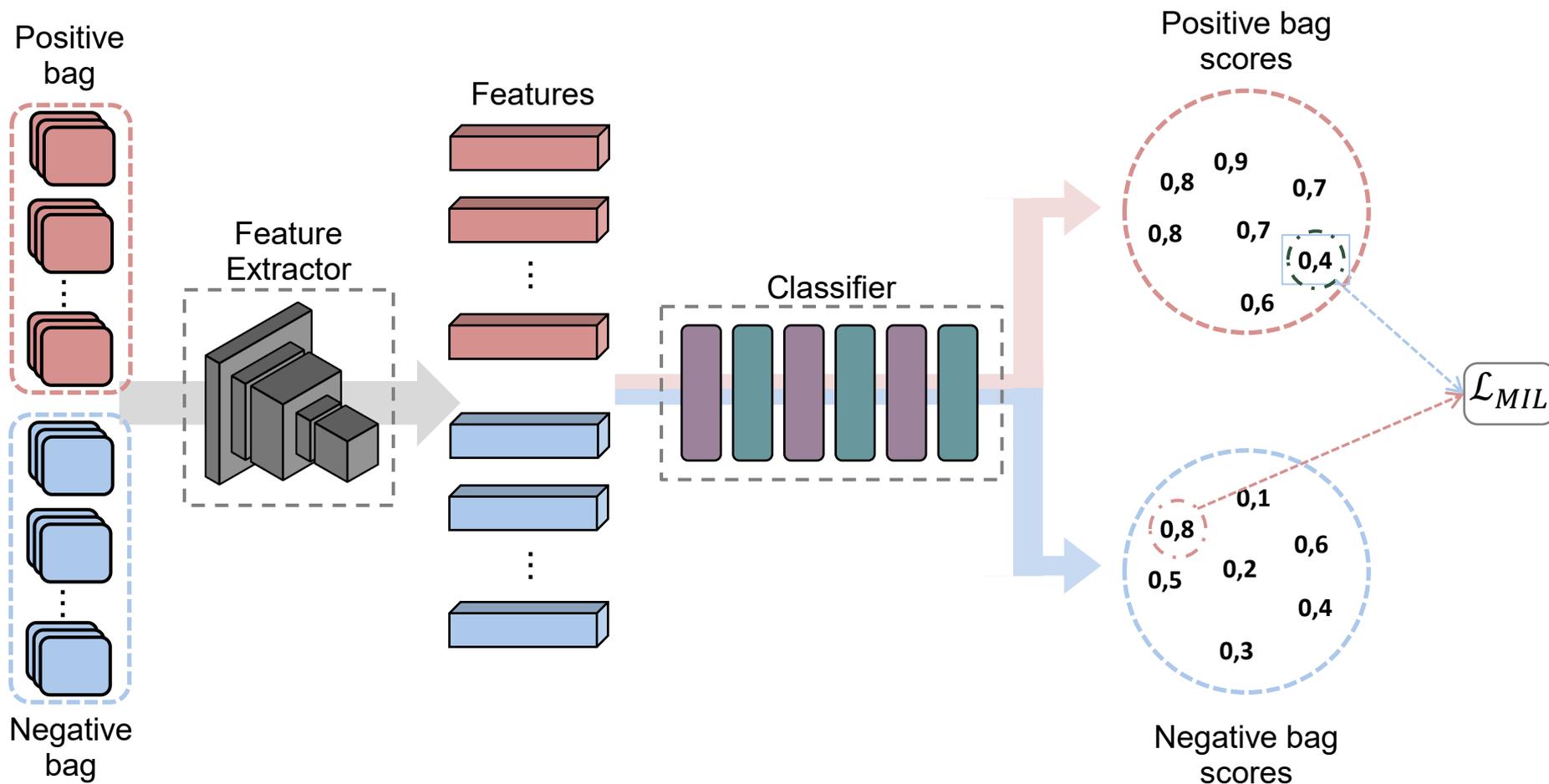
- Video Anomaly Detection (VAD) is the task of detecting anomalous (human) activities in videos.
- Recent popular approaches are weakly supervised, in which videos are labeled at video level.

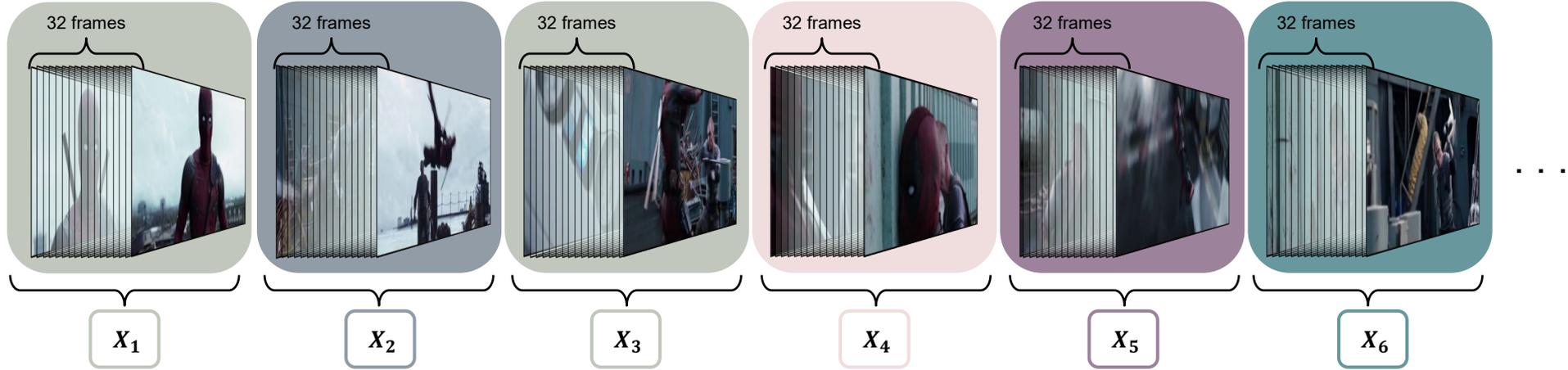


Videos get split in 32 frames clips.

Standard video split: Every clip gets processed by the model.

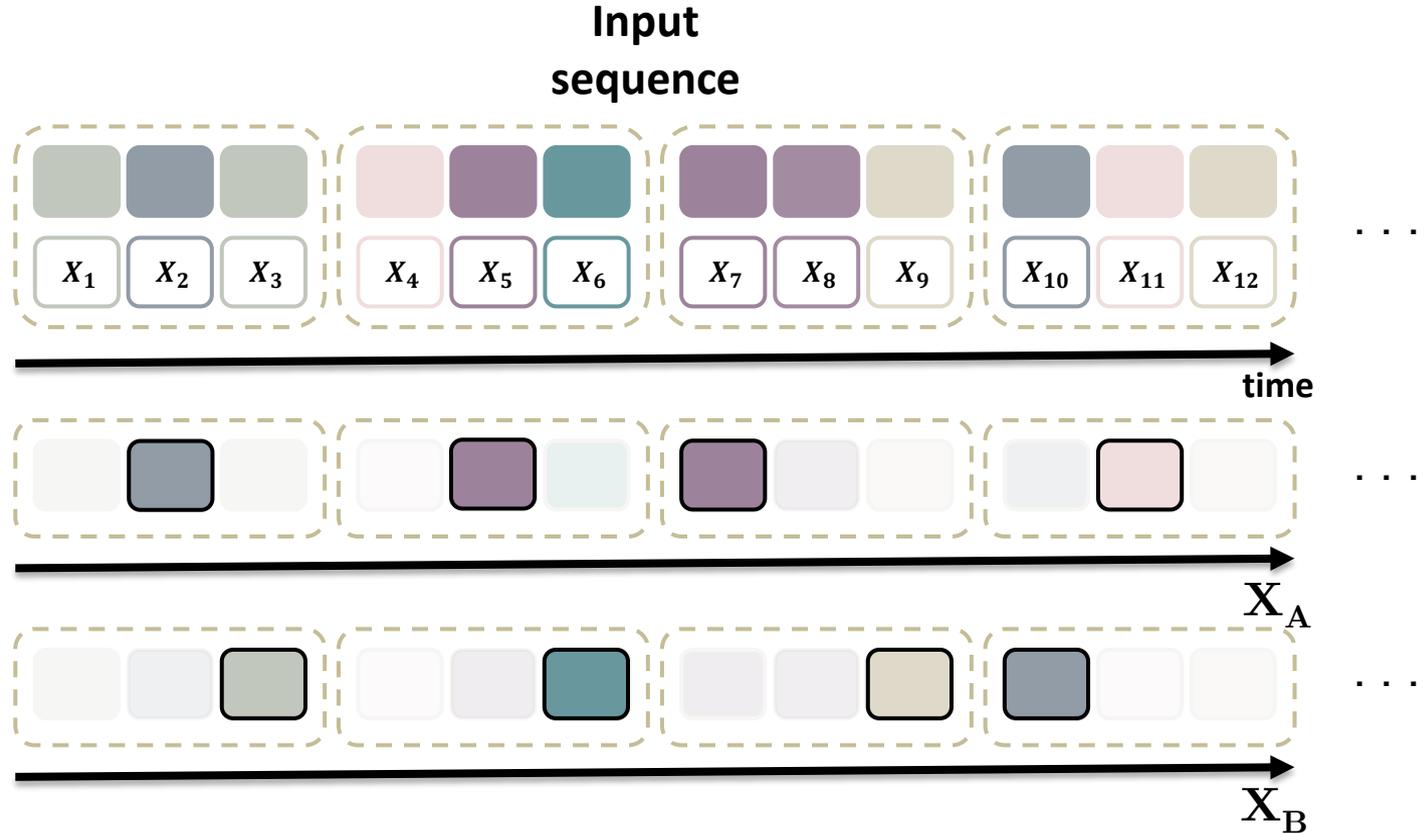
In weakly-supervised methods, clips from anomalous videos are grouped into **positive bags**, while normal ones into **negative bags**.

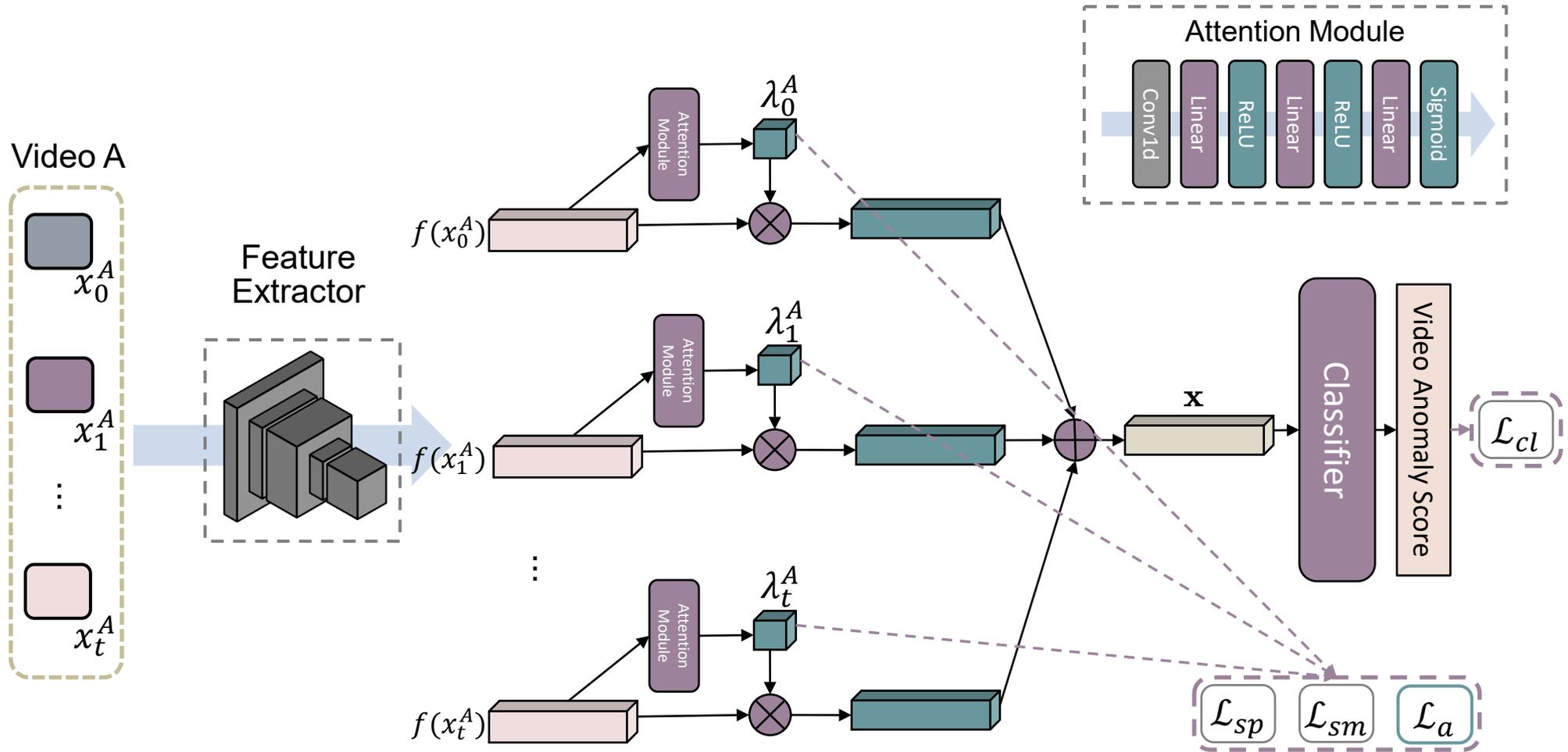




Videos get split in 32 frames clips.

Custom sampling: We sample only one clip from every window of three clips.





$$\mathcal{L} = \mathcal{L}_{cl} + \mathcal{L}_{sm} + \mathcal{L}_{sp} + \mathcal{L}_a$$

$$\mathcal{L} = \boxed{\mathcal{L}_{cl}} + \mathcal{L}_{sm} + \mathcal{L}_{sp} + \mathcal{L}_a$$

$$\mathcal{L}_{cl} = \text{BCE}(\boxed{f(x_i)}, \boxed{y_i})$$

Classifier output

True video label

The **classification loss** is a binary cross-entropy at video-level.

This is used to ensure the model can recognize if a video contains anomalies or not.

$$\mathcal{L} = \mathcal{L}_{cl} + \boxed{\mathcal{L}_{sm}} + \mathcal{L}_{sp} + \mathcal{L}_a$$

$$\mathcal{L}_{cl} = \text{BCE}(f(x_i), y_i)$$

$$\mathcal{L}_{sm} = \sum_{t=1}^{T-1} (\lambda_t - \lambda_{t+1})^2$$

The **smooth loss** impose adjacent attention coefficients to vary as little as possible.

$$\mathcal{L} = \mathcal{L}_{cl} + \mathcal{L}_{sm} + \boxed{\mathcal{L}_{sp}} + \mathcal{L}_a$$

$$\mathcal{L}_{cl} = \text{BCE}(f(x_i), y_i)$$

$$\mathcal{L}_{sm} = \sum_{t=1}^{T-1} (\lambda_t - \lambda_{t+1})^2$$

$$\mathcal{L}_{sp} = \|\lambda_t\|_1$$

The **sparsity loss** penalizes the ℓ_1 norm of the attention weights.
This reflects the rarity of the anomalies in videos.

$$\mathcal{L} = \mathcal{L}_{cl} + \mathcal{L}_{sm} + \mathcal{L}_{sp} + \boxed{\mathcal{L}_a}$$

$$\mathcal{L}_{cl} = \text{BCE}(f(x_i), y_i)$$

$$\mathcal{L}_{sm} = \sum_{t=1}^{T-1} (\lambda_t - \lambda_{t+1})^2$$

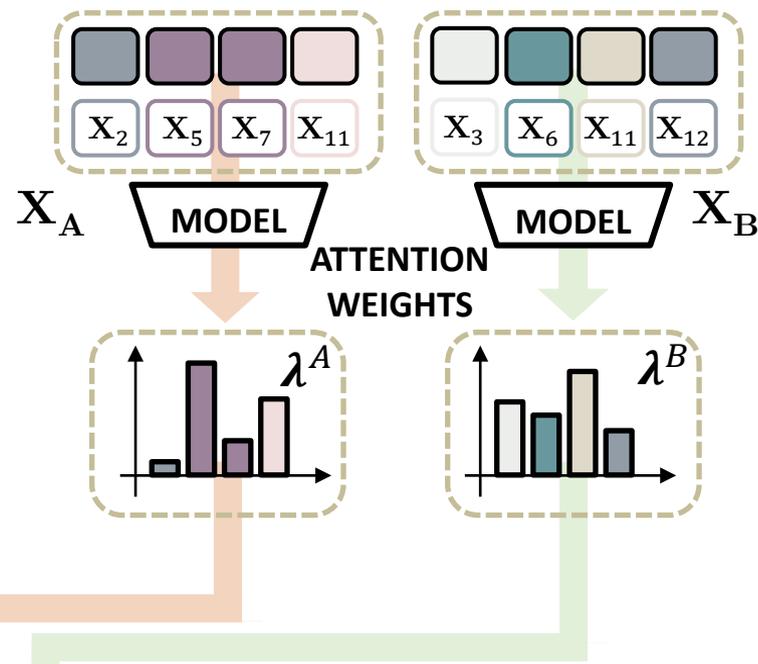
$$\mathcal{L}_{sp} = \|\lambda_t\|_1$$

$$\mathcal{L}_a = \sum_{t=1}^T (\lambda_t^A - \lambda_t^B)^2$$

The alignment loss is a consistency-based regularization term.

1. Sample two slightly different version of the video.
2. Classify them and extract the attention weights of the clips.
3. The attention weights of the two sampling should be close together.

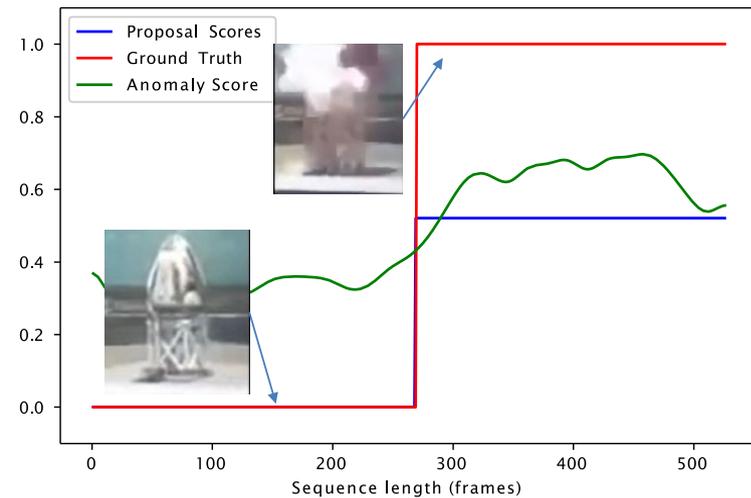
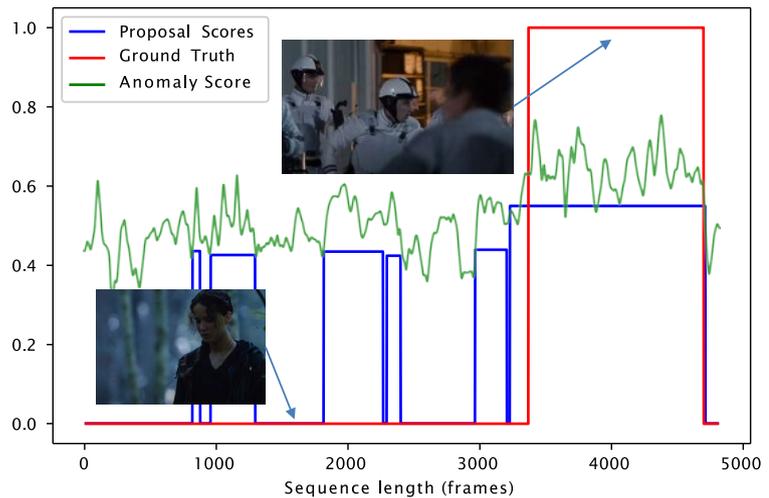
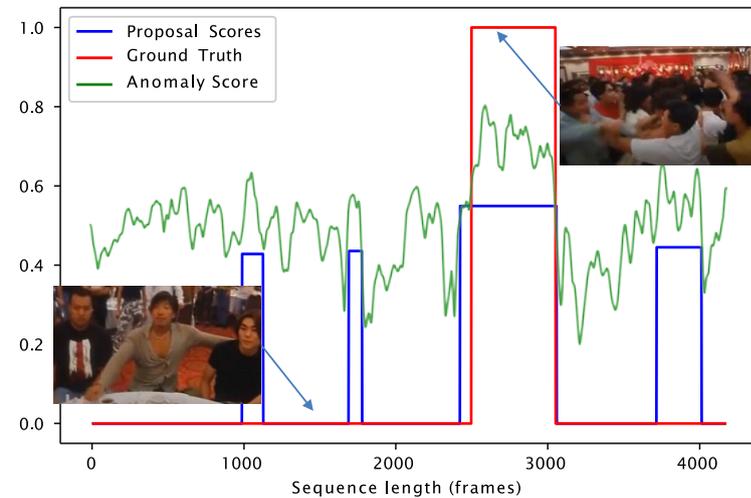
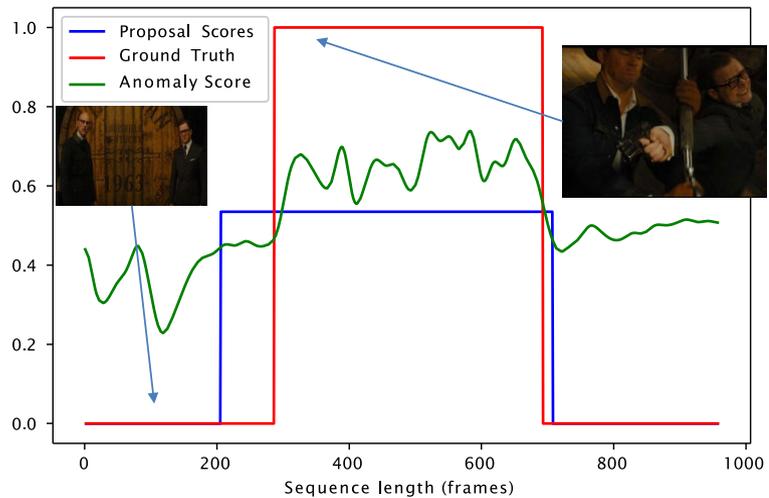
With this additional regularization term we enforce smoothness over a wider temporal horizon.



$$\mathcal{L}_a = \sum_{t=1}^T (\lambda_t^A - \lambda_t^B)^2$$

	Video Level		Segment Level		Frame Level Proposal		Frame Level	
	AUC%	AP%	AUC%	AP%	AUC%	AP%	AUC%	AP%
Align Loss								
-	97.91	98.36	84.39	66.75	85.14	68.01	84.57	65.96
✓	97.79	98.28	85.49	66.87	90.23	71.68	85.65	66.05

Qualitative Results



- Alignment loss allows to learn effective frame-level scores in weakly supervised settings.
- A base network equipped also with other common regularization techniques brings even more improvements.

Code: github.com