

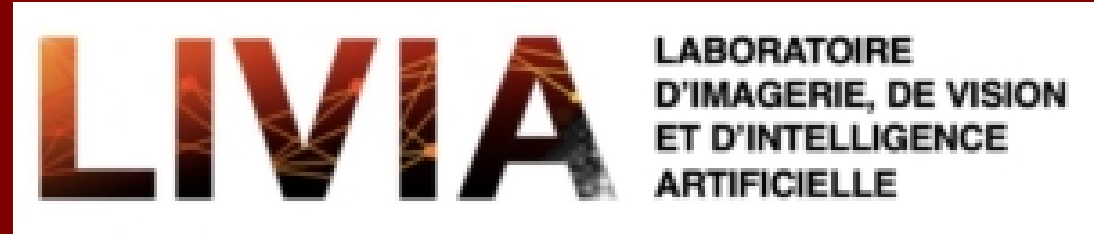
# TCAM: Temporal Class Activation Maps for Object Localization in Weakly-Labeled Unconstrained Videos

(#1703)

Soufiane Belharbi<sup>1\*</sup> & Ismail Ben Ayed<sup>1</sup> & Luke McCaffrey<sup>2</sup> & Eric Granger<sup>1</sup>

<sup>1</sup>LIVIA, Dept. of Systems Engineering, ETS Montreal, Canada

<sup>2</sup>Rosalind and Morris Goodman Cancer Research Centre, Department of Oncology, McGill University



## Context

Localization in unconstrained videos is challenging due to: (1) moving objects, (2) camera motion, (3) viewpoint changes, (4) decoding artifacts, (5) editing effects, and (6) costly annotation

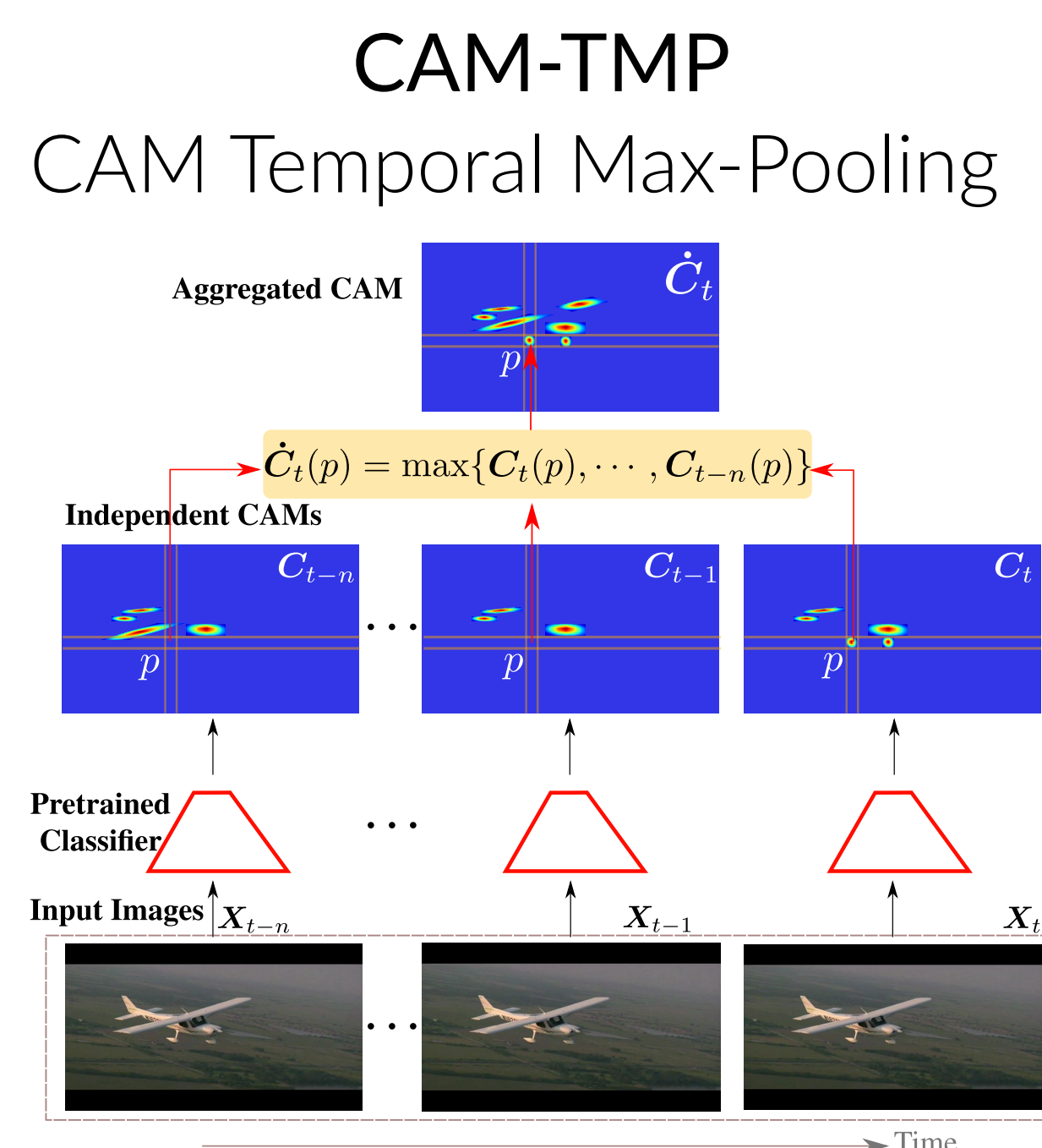
**Weak supervision:** global video tags (classes) are available

State-of-the-art methods in weakly-supervised video object localization have good performance, but:

- multiple sequential and independent stages
- video tags (labels) are used only to cluster video
- ROI are not necessarily discriminative
- motion cues (optical flow) are noisy, not always discriminative, and need post-processing
- requires solving an optimization problem at inference time: slow inference (build a model per class).

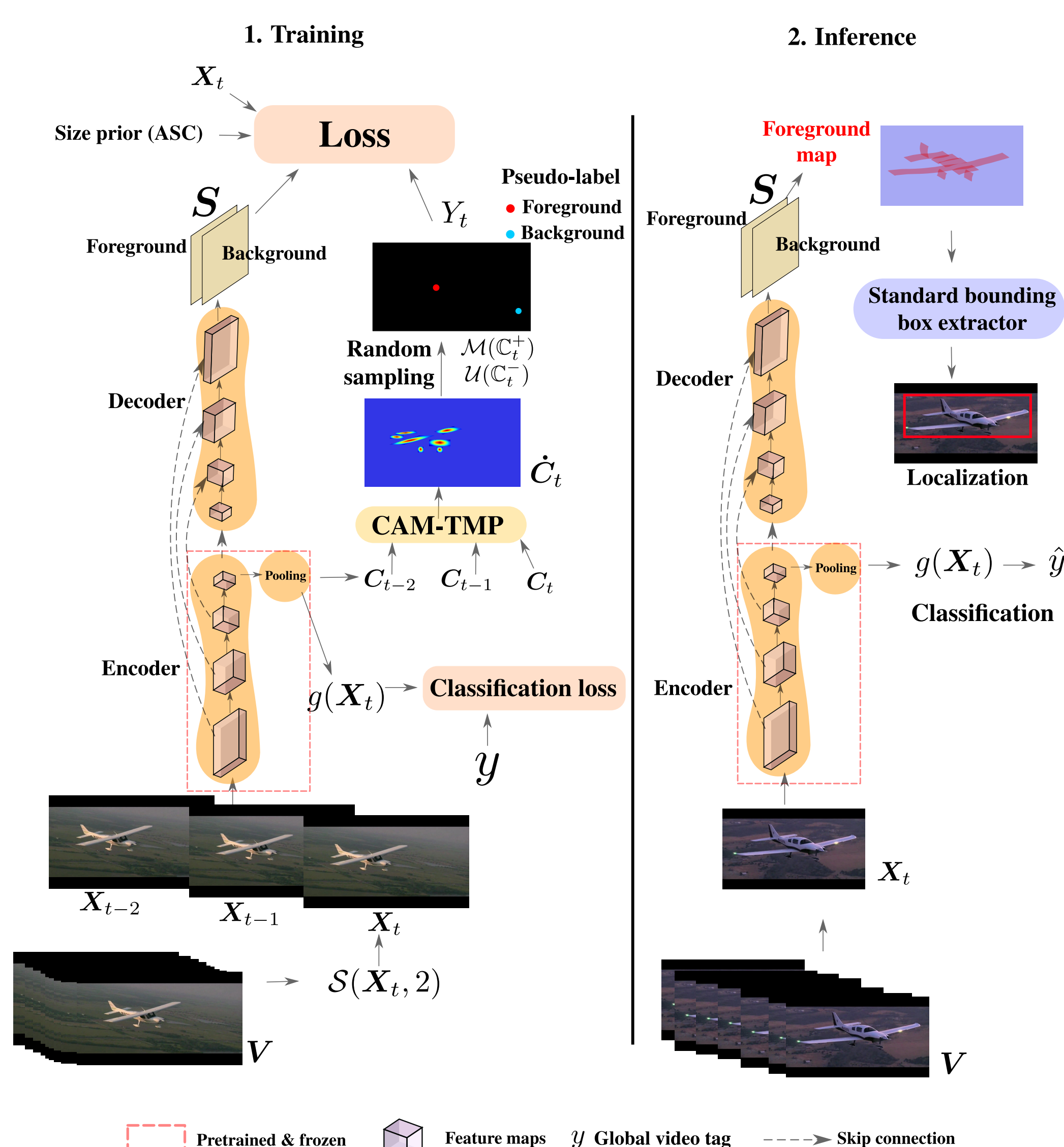
## Proposal: Use CAMs with Temporal Dependency

- CAM methods are successful in developing discriminative WSOL models for CNNs
- We adapt CAM methods to exploit the spatiotemporal dependency in videos
- Leverage the slight variations in sets of consecutive frames
- The **CAM-TMP** module aggregates diversified CAMs from  $n$  frames



## Proposal: Overall Architecture

- Includes a learnable decoder to produce accurate CAMs
- Training: accounts for spatio-temporal dependency at a CAM level – it leverages sequences of  $n$  frames
- Uses aggregated CAMs to sample pixel pseudo-labels for training the decoder
- Fast inference: CAMs produced from independent frames



Pretrained & frozen Feature maps  $y$  Global video tag Skip connection

## Proposal: Training Loss

Total loss: pseudo labels, CRF, unsupervised size constraint.

$$\min_{\theta} \sum_{p \in \Omega_t^l} H_p(Y_t, S_t) + \lambda \mathcal{R}(S_t, X_t), \quad \text{s.t.} \quad \sum S_t^r \geq 0, \quad r \in \{0, 1\} \quad (1)$$

where the classifier is frozen,  $\sum_{p \in \Omega_t^l} H_p(Y_t, S_t)$  is partial cross-entropy loss over pixel pseudo-labels,  $\mathcal{R}(S_t, X_t)$  is a CRF loss, and  $\sum S_t^r \geq 0$  is unsupervised size constraint (the Absolute Size Constraints, ASC).

## Empirical Results

**Datasets:** Localization accuracy (CorLoc) on YouTube-Object v1.0 and v2.2 datasets

### Comparison with State-of-Art Methods

Dataset	Method (venue)	Aero	Bird	Boat	Car	Cat	Cow	Dog	Horse	Mbike	Train	Avg	Time/Frame
YT0v1	Prest et al. (cvpr,2012)	51.7	17.5	34.4	34.7	22.3	17.9	13.5	26.7	41.2	25.0	28.5	N/A
	Papazoglou et al. (iccv,2013)	65.4	67.3	38.9	65.2	46.3	40.2	65.3	48.4	39.0	25.0	50.1	4s
	Joulin et al. (eccv,2014)	25.1	31.2	27.8	38.5	41.2	28.4	33.9	35.6	23.1	25.0	31.0	N/A
	Kwak et al. (iccv,2015)	56.5	66.4	58.0	76.8	39.9	69.3	50.4	56.3	53.0	31.0	55.7	N/A
	Rochan et al. (icv,2016)	60.8	54.6	34.7	57.4	19.2	42.1	35.8	30.4	11.7	11.4	35.8	N/A
	Tokmakov et al. (eccv,2016)	71.5	74.0	44.8	72.3	52.0	46.4	71.9	54.6	45.9	32.1	56.6	N/A
	POD (cvpr,2016)	64.3	63.2	73.3	68.9	44.4	62.5	71.4	52.3	78.6	23.1	60.2	N/A
	Tsai et al. (eccv,2016)	66.1	59.8	63.1	72.5	54.0	64.9	66.2	50.6	39.3	42.5	57.9	N/A
	Haller et al. (iccv,2017)	76.3	71.4	65.0	58.9	68.0	55.9	70.6	33.3	69.7	42.4	61.1	0.35s
	Croitoru et al. (LowRes-Net <sub>iter1</sub> ) (icv,2019)	77.0	67.5	77.2	68.4	54.5	68.3	72.0	56.7	44.1	34.9	62.1	0.02s
	Croitoru et al. (LowRes-Net <sub>iter2</sub> ) (icv,2019)	79.7	67.5	68.3	69.6	59.4	75.0	78.7	48.3	48.5	39.5	63.5	0.02s
	Croitoru et al. (DilateU-Net <sub>iter2</sub> ) (icv,2019)	85.1	72.7	76.2	68.4	59.4	76.7	77.3	46.7	48.5	46.5	65.8	0.02s
	Croitoru et al. (MultiSelect-Net <sub>iter2</sub> ) (icv,2019)	84.7	72.7	78.2	69.6	60.4	80.0	78.7	51.7	50.0	46.5	67.3	0.15s
	SPFTN (M) (tpami,2020)	66.4	73.8	63.3	83.4	54.5	58.9	61.3	45.4	55.5	30.1	59.3	N/A
	SPFTN (P) (tpami,2020)	<b>97.3</b>	27.8	<b>81.1</b>	65.1	56.6	72.5	59.5	<b>81.8</b>	79.4	22.1	64.3	N/A
	FPVOS (cvpr,2021)	77.0	72.3	64.7	67.4	79.2	58.3	74.7	45.2	<b>80.4</b>	42.6	65.8	0.29s
YT0v2.2	CAM (cvpr,2016)	75.0	55.5	43.2	69.7	33.3	52.4	32.4	74.2	14.8	50.0	50.1	0.2ms
	GradCAM (iccv,2017)	86.9	63.0	51.3	81.8	45.4	62.0	37.8	67.7	18.5	50.0	56.4	27.8ms
	GradCAM++ (wacv,2018)	79.8	85.1	37.8	81.8	75.7	52.4	64.9	64.5	33.3	<b>56.2</b>	63.2	28.0ms
	Smooth-GradCAM++ (corr,2019)	78.6	59.2	56.7	60.6	42.4	61.9	56.7	64.5	40.7	50.0	57.1	136.2ms
	XGradCAM (bmvc,2020)	79.8	70.4	54.0	<b>87.8</b>	33.3	52.4	37.8	64.5	37.0	50.0	56.7	14.2ms
	LayerCAM (ieee,2021)	85.7	<b>88.9</b>	45.9	78.8	75.5	61.9	64.9	64.5	33.3	<b>56.2</b>	65.6	17.9ms
	TCAM (ours)	90.5	70.4	62.2	75.7	<b>84.8</b>	<b>81.0</b>	<b>81.0</b>	64.5	70.4	50.0	<b>73.0</b>	18.5ms
	Haller et al. (iccv,2017)	76.3	68.5	54.5	50.4	59.8	42.4	53.5	30.0	53.5	<b>60.7</b>	54.9	0.35s
	Croitoru et al. (LowRes-Net <sub>iter1</sub> ) (icv,2019)	75.7	56.0	52.7	57.3	46.9	57.0	48.9	44.0	27.2	56.2	52.2	0.02s
	Croitoru et al. (LowRes-Net <sub>iter2</sub> ) (icv,2019)	78.1	51.8	49.0	60.5	44.8	62.3	52.9	48.9	30.6	54.6	53.4	0.02s
	Croitoru et al. (DilateU-Net <sub>iter2</sub> ) (icv,2019)	74.9	50.7	50.7	60.9	45.7	60.1	54.4	42.9	30.6	57.8	52.9	0.02s
	Croitoru et al. (BasicU-Net <sub>iter2</sub> ) (icv,2019)	<b>82.2</b>	51.8	51.5	62.0	50.9	64.8	55.5	45.7	35.3	55.9	56.5	0.02s
	Croitoru et al. (MultiSelect-Net <sub>iter2</sub> ) (icv,2019)	81.7	51.5	54.1	62.5	49.7	68.8	55.9	50.4	33.3	57.0	55.6	0.15s
	CAM (cvpr,2016)	52.3	66.4	25.0	66.4	39.7	<b>87.8</b>	34.7	53.6	45.4	43.7	51.5	0.2ms
	GradCAM (iccv,2017)	44.1	68.4	50.0	61.1	51.8	79.3	56.0	47.0	44.8	42.4	54.5	27.8ms
	GradCAM++ (wacv,2018)	74.7	78.1	38.2	69.7	56.7	84.3	61.6	61.9	43.0	44.3	61.2	28.0ms
	Smooth-GradCAM++ (corr,2019)	74.1	83.2	38.2	64.2	49.6	82.1	57.3	52.0	51.1	42.4	59.5	136.2ms
	XGradCAM (bmvc,2020)	68.2	44.5	45.8	64.0	46.8	86.4	44.0	57.0	44.9	45.0	54.6	14.2ms
	LayerCAM (ieee,2021)	80.0	84.5	47.2	<b>73.5</b>	55.3	83.6	71.3	60.8	55.7	48.1	66.0	17.9ms
	TCAM (ours)	79.4	<b>94.9</b>	<b>75.7</b>	61.7	<b>68.8</b>	87.1	<b>75.0</b>	<b>62.4</b>	<b>72.1</b>	45.0	<b>72.2</b>	18.5ms

## Visual Results

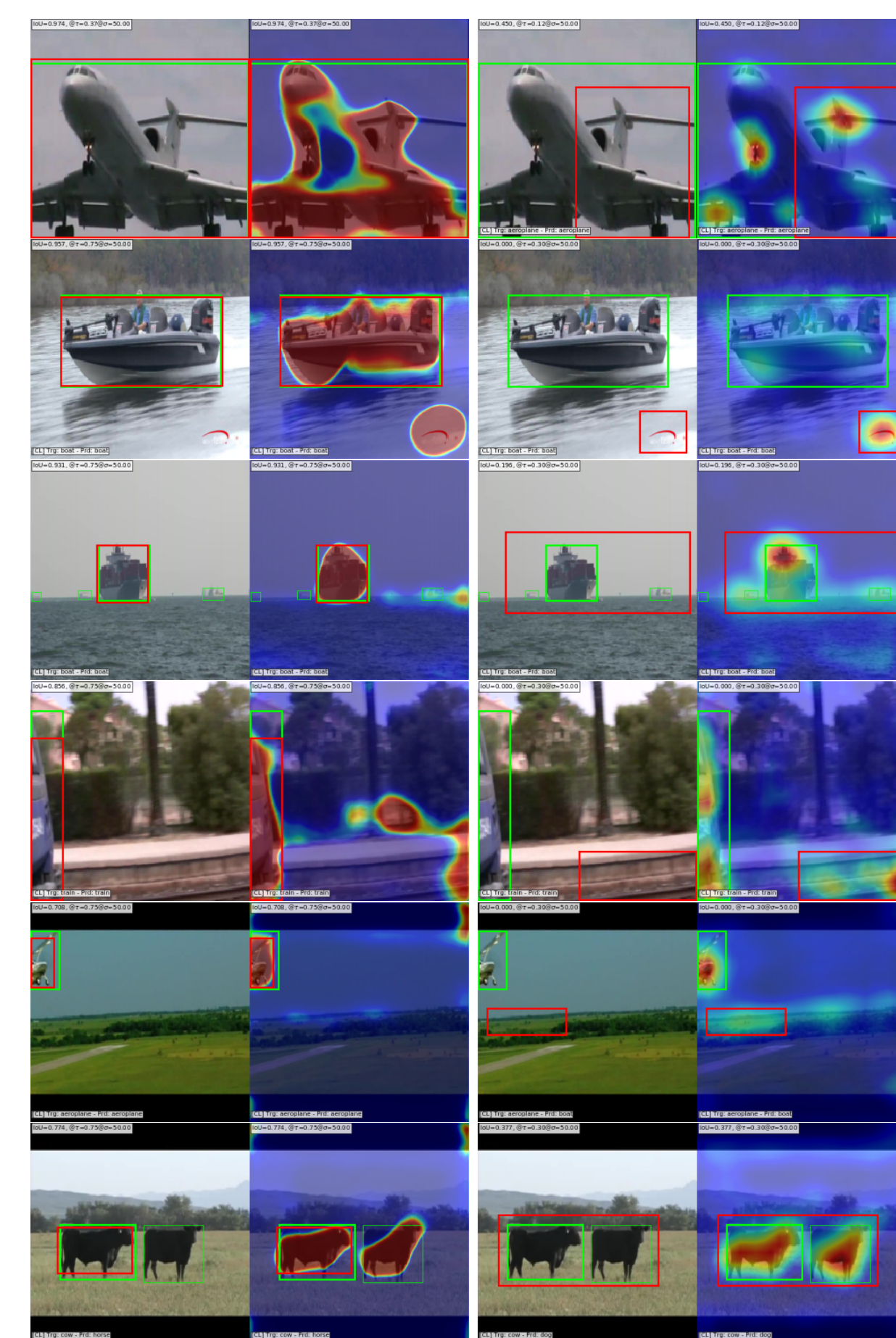


Figure 1. Prediction examples of test sets frames. Left: TCAM (ours). Right: baseline CAM method, LayerCAM. Bounding box: ground truth (green), prediction (red). The second column is predicted CAM on images.

## Ablation Studies

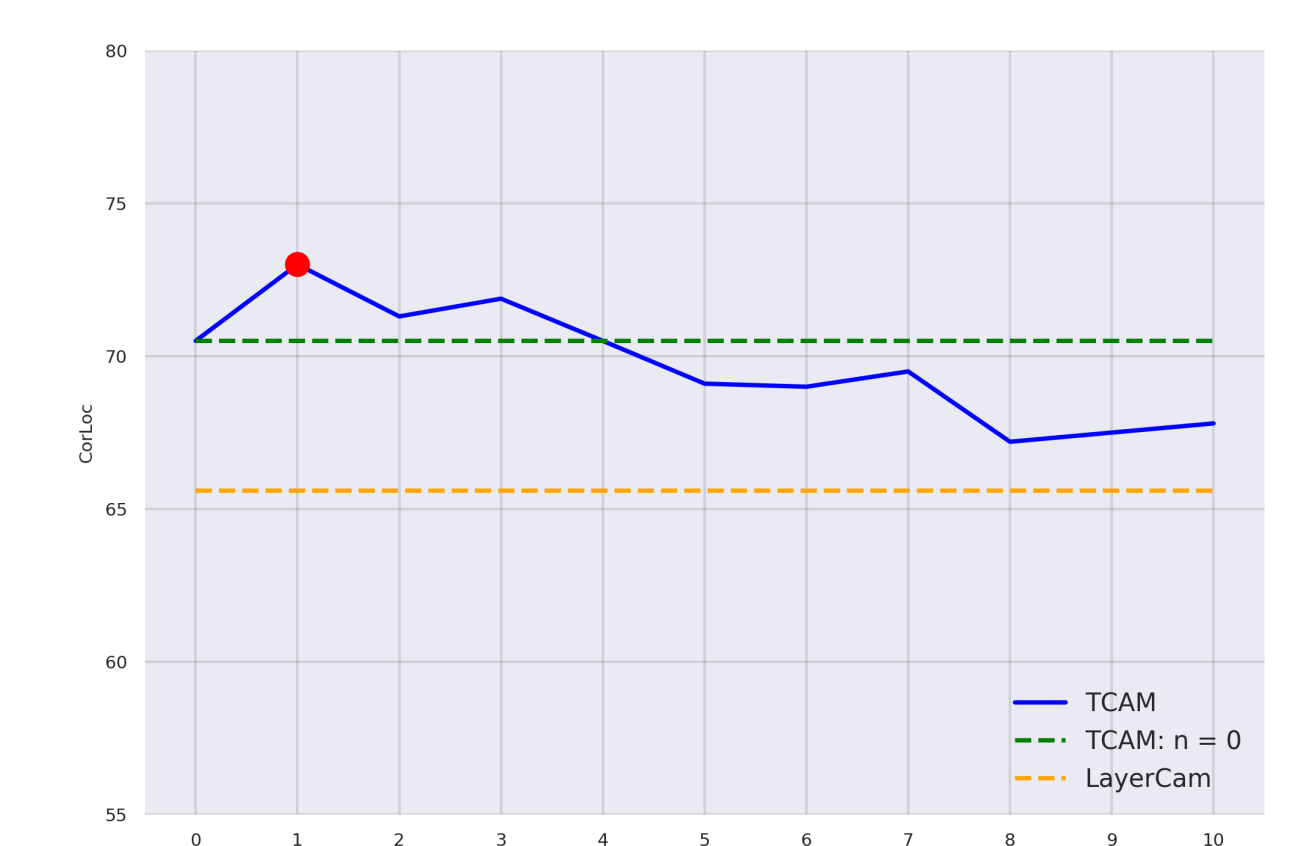


Figure 2. Localization accuracy of TCAM with different temporal dependencies  $n$  on the YT0v1 test set.

Methods	CorLoc
Layer-CAM (ieee,2021)	65.6
Ours + $C_t^+$ + $C_t^-$	68.5
$n = 0$ Ours + $C_t^+$ + $C_t^-$ + CRF	69.6
Ours + $C_t^+$ + $C_t^-$ + ASC	66.2
Ours + $C_t^+$ + $C_t^-$ + CRF + ASC	70.5
$n > 0$ Ours + $C_t^+$ + $C_t^-$ + CRF + ASC + CAM-TMP	73.0
Improvement	+7.4

Table 1. Localization accuracy of TCAM with different losses on the YT0v1 test set.

## Main Conclusions

- Standard CAM methods: can yield discriminative CNNs with accurate localization
- With TCAM: leveraging temporal information during training yielded new state-of-the-art results