

Hierarchical Attentive Recurrent Tracking

Adam R. Kosiorek, Alexy Bewley, Ingmar Posner

Applied AI Lab, Department of Engineering Science, University of Oxford, UK

Hierarchical Attention



Ground-truth and predicted bounding boxes and an attention glimpse. The bottom row shows the three layers of attention:
 1st layer extracts an attention glimpse (left)
 2nd layer uses appearance attention to build a location map (middle)
 3rd layer uses the location map to suppress distractors (right)

Loss

Directly optimise Intersection-over-Union (IoU) and guide attention mechanisms.

$$\mathcal{L}_{\text{HART}}(\cdot) = \lambda_t \mathcal{L}_t(\cdot) + \lambda_s \mathcal{L}_s(\cdot) + \lambda_a \mathcal{L}_a(\cdot) + \beta R(\cdot)$$

Tracking: Negative log of Intersection-over-Union.

$$\mathcal{L}_t(\mathcal{D}, \theta) = \mathbb{E}_{p(\hat{\mathbf{b}}_{1:T}|\mathbf{x}_{1:T}, \mathbf{b}_1)} \left[-\log \text{IoU}(\hat{\mathbf{b}}_t, \mathbf{b}_t) \right]$$

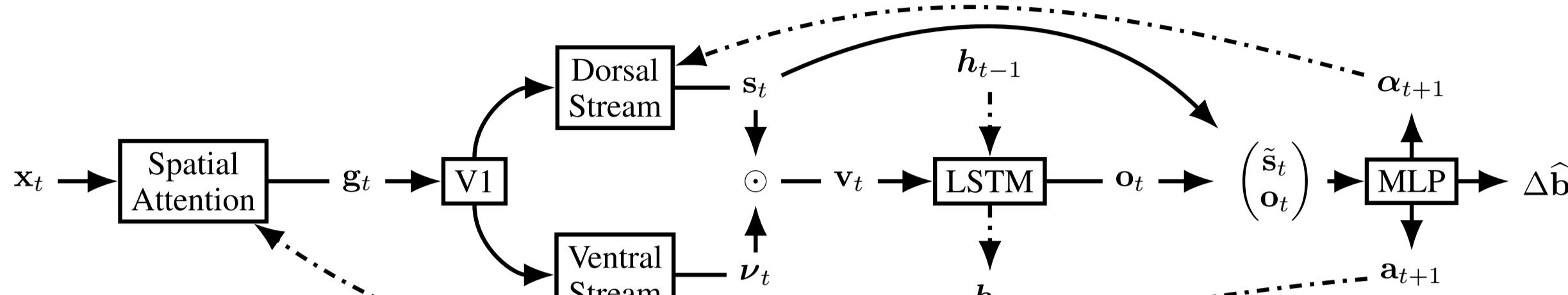
Spatial Attention: It follows the object, but shouldn't be too big.

$$\mathcal{L}_s(\mathcal{D}, \theta) = \mathbb{E}_{p(\mathbf{a}_{1:T}|\mathbf{x}_{1:T}, \mathbf{b}_1)} \left[-\log \left(\frac{\mathbf{a}_t \cap \mathbf{b}_t}{\text{area}(\mathbf{b}_t)} \right) - \log(1 - \text{IoU}(\mathbf{a}_t, \mathbf{x}_t)) \right]$$

Appearance Attention: Cross-entropy with dynamically created target mask $\tau(\mathbf{a}_t, \mathbf{b}_t)$.

$$\mathcal{L}_a(\mathcal{D}, \theta) = \mathbb{E}_{p(\mathbf{a}_{1:T}, \mathbf{s}_{1:T}|\mathbf{x}_{1:T}, \mathbf{b}_1)} [H(\tau(\mathbf{a}_t, \mathbf{b}_t), \mathbf{s}_t)]$$

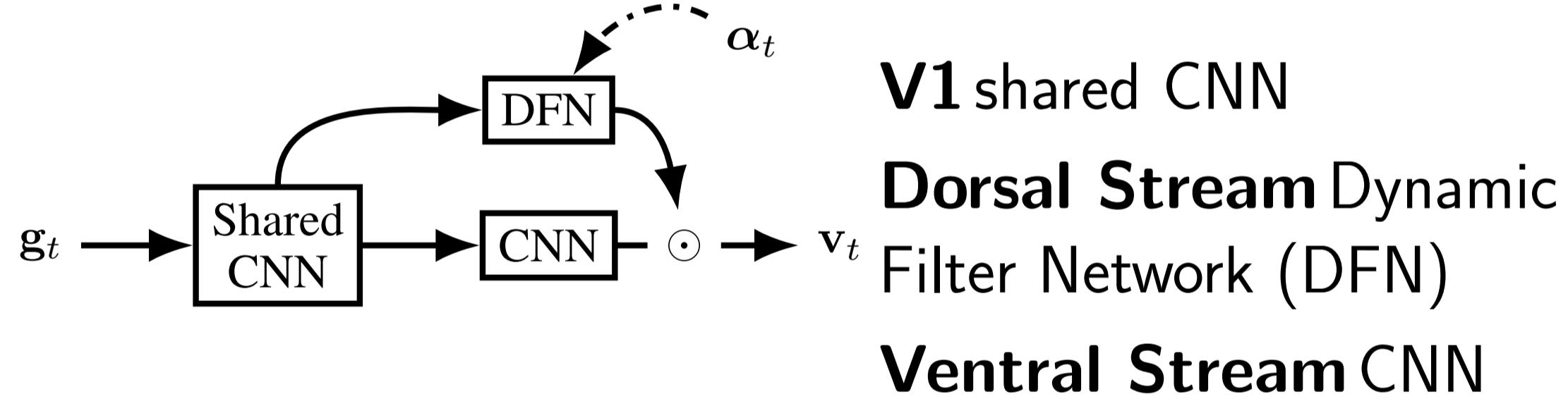
Two-stream Attentive Model



x_t input image
 g_t attention glimpse
 ν_t appearance-based features
 s_t object segmentation
 v_t masked features

h_t hidden state
 o_t LSTM output
 α_{t+1} appearance
 $\Delta \hat{\mathbf{b}}_t$ bounding-box update
 a_{t+1} spatial attention

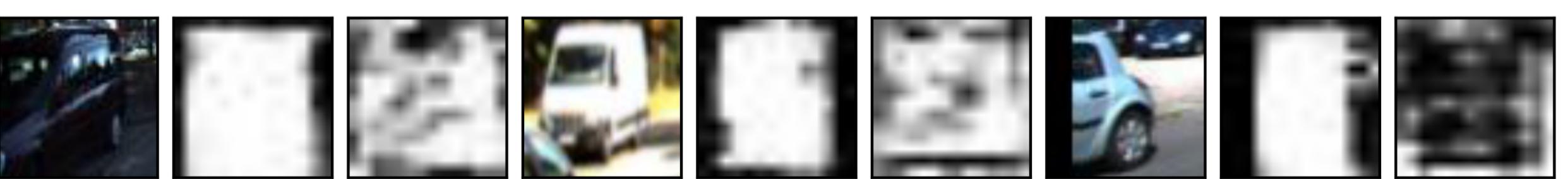
Appearance attention architecture:



Is Attention Loss Important?



Appearance attention loss (top) prevents an ID swap when a pedestrian is occluded by another one (bottom).

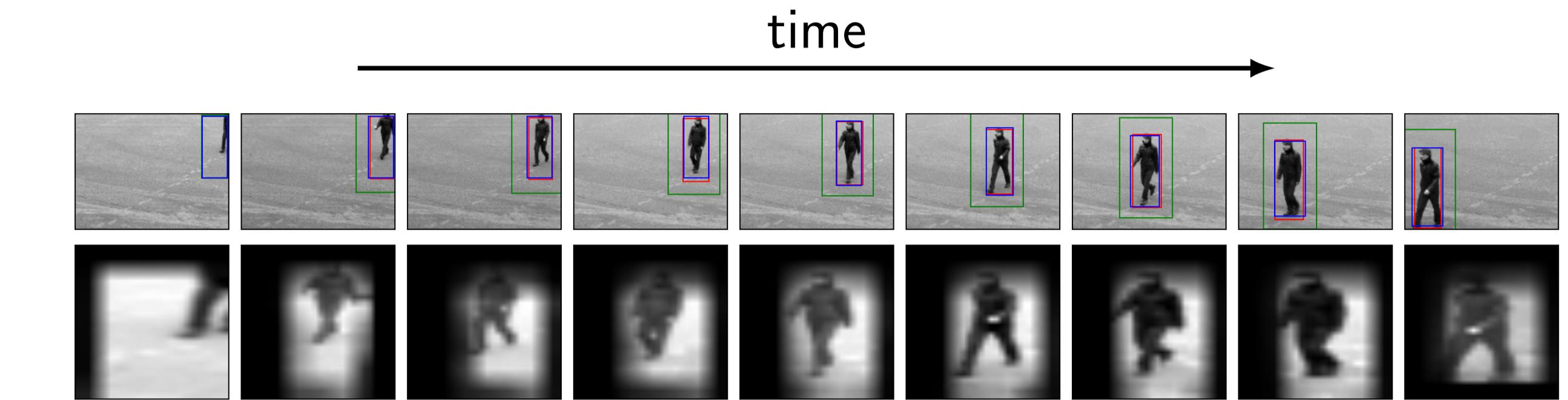


Left to right: glimpses and segmentations learned with and without appearance loss. Attention loss leads to distractor suppression.

References:

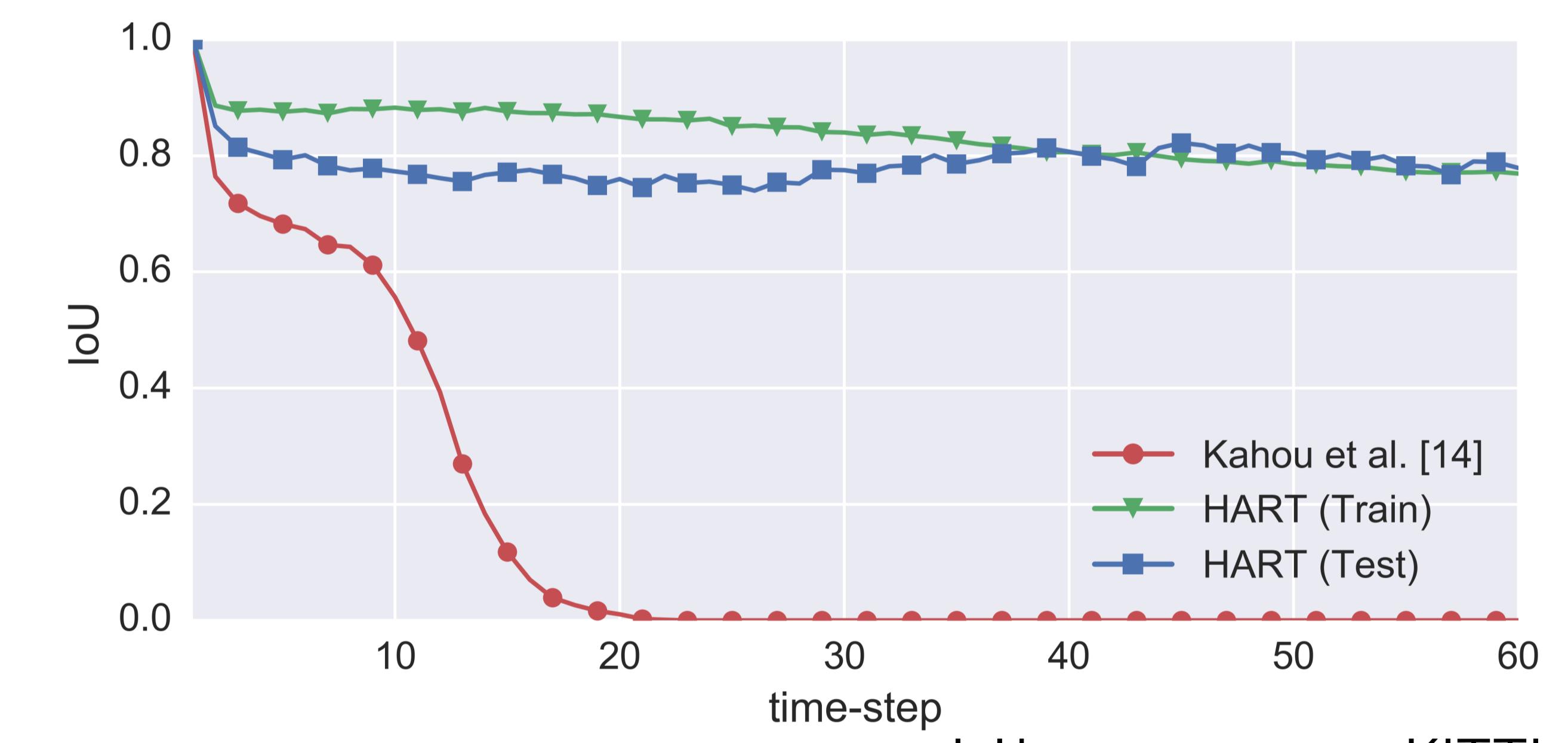
- [1] Samira Ebrahimi Kahou, Vincent Michalski, and Roland Memisevic. RATM: Recurrent Attentive Tracking Model. CVPR Work., 2017.
- [2] Christian Schuldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: A local SVM approach. In ICPR. IEEE, 2004.
- [3] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The KITTI dataset. Int. J. Rob. Res., 32(11):12311237, sep 2013.

Pedestrian Tracking: KTH Dataset [2]



- Attention, prediction and ground-truth overlap at initialization
 - Every 16th frame of the sequence at 25 fps.
 - 2nd row shows attention glimpses multiplied with appearance attention
- | | |
|-------------------|-------------------------|
| | Intersection over Union |
| Kahou et. al. [1] | 55.03% |
| Ours | 77.11% |

Scaling to Real-world Data: KITTI [3]



Average IoU on KITTI over 60 time-steps
 Kahou et. al. [1] | Spatial Att | App Att | HART
 0.14 | 0.60 | 0.78 | 0.81
 HART (train) presents evaluation on the train set.

Conclusion

- Bio-inspired:** Neural Recurrent tracking with Attention Mechanisms.
- Interpretable:** Important features selected by spatial attention and object segmentation mechanisms.
- Scalable:** Auxiliary loss terms allow scaling to complex real-world datasets.
- Efficient:** > 120 fps on a laptop!
- Future Work:** Multi-object tracking.