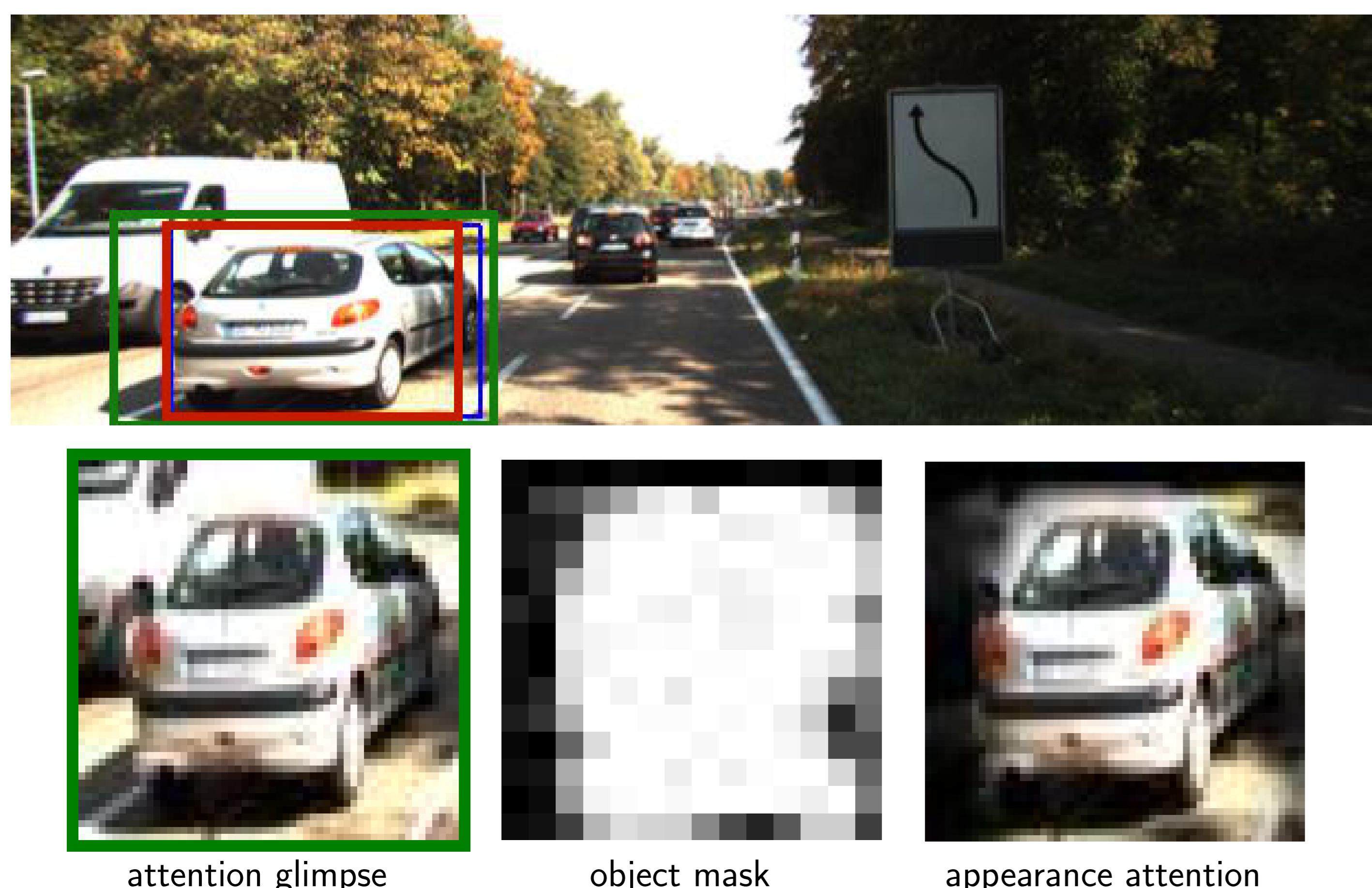


Problem Statement

- What:** Class-agnostic single object tracking in real-world videos with camera motion
- Difficulties:** No target-specific discriminative models
Cluttered backgrounds with many distractors
- How:** Discard uninformative background features
Learn arbitrary motion models
Anticipate appearance changes
- Approach:** Recurrent Neural Network with Hierarchical Attention Mechanism

Hierarchical Attention



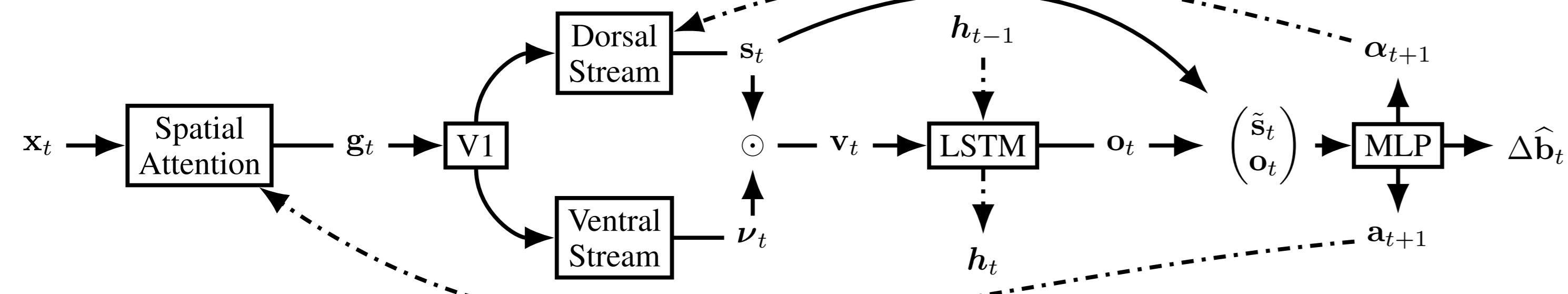
attention glimpse

object mask

appearance attention

- Bio-inspired:** Two-stream processing pathway and attention mechanisms adapted from human visual cortex.
- Interpretable:** Important features selected by Spatial Attention and Object Segmentation mechanisms.
- Scalable:** Applicable to real-world data due to distractor suppression and auxiliary loss terms.
- Efficient:** Attention quickly discards irrelevant features
> 120 fps on a laptop!

Two-Stream Attentive Model



x_t input image

g_t attention glimpse

v_t appearance-based features

s_t object segmentation

ν_t masked features

h_t hidden state

o_t LSTM output

α_{t+1} appearance

$\Delta \hat{b}_t$ bounding-box update

a_{t+1} spatial attention

V1: Shared low-level convolutional neural network for feature extraction.

Dorsal Stream: Attention driven by high-level appearance features in a *top-down* fashion.

Ventral Stream: Appearance features computed in a *bottom-up* fashion.

Learning to Attend for Tracking

The model is optimised for the following three objectives.

Main Tracking Objective: Maximise the overlap of predicted box with the true object.

Spatial Attention: The glimpse should contain the object, but shouldn't be too big.

Appearance Attention: The dynamic appearance features should respond specifically to the target object.

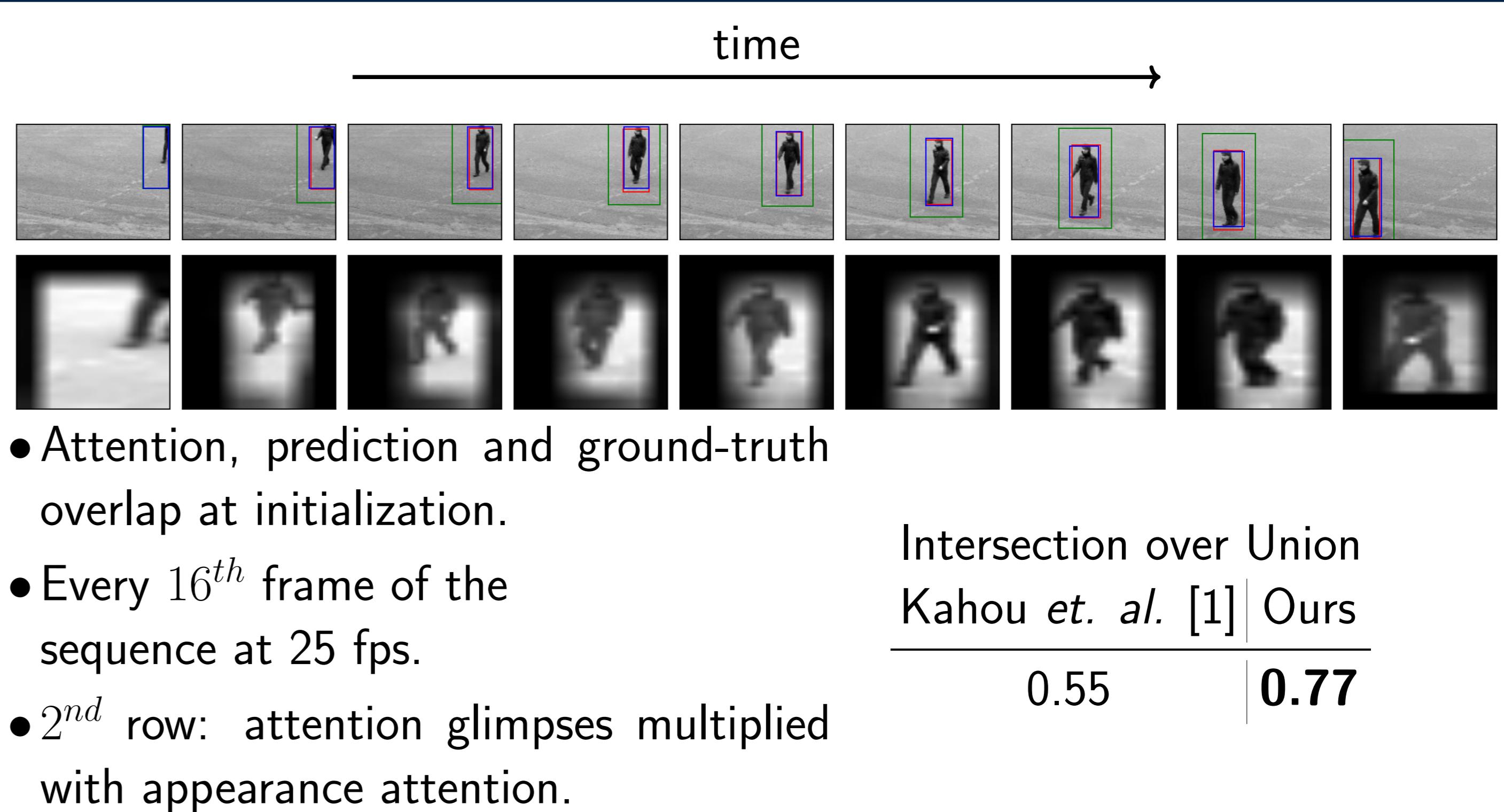


Appearance attention loss (top) prevents an ID swap when a pedestrian is occluded by another one (bottom).

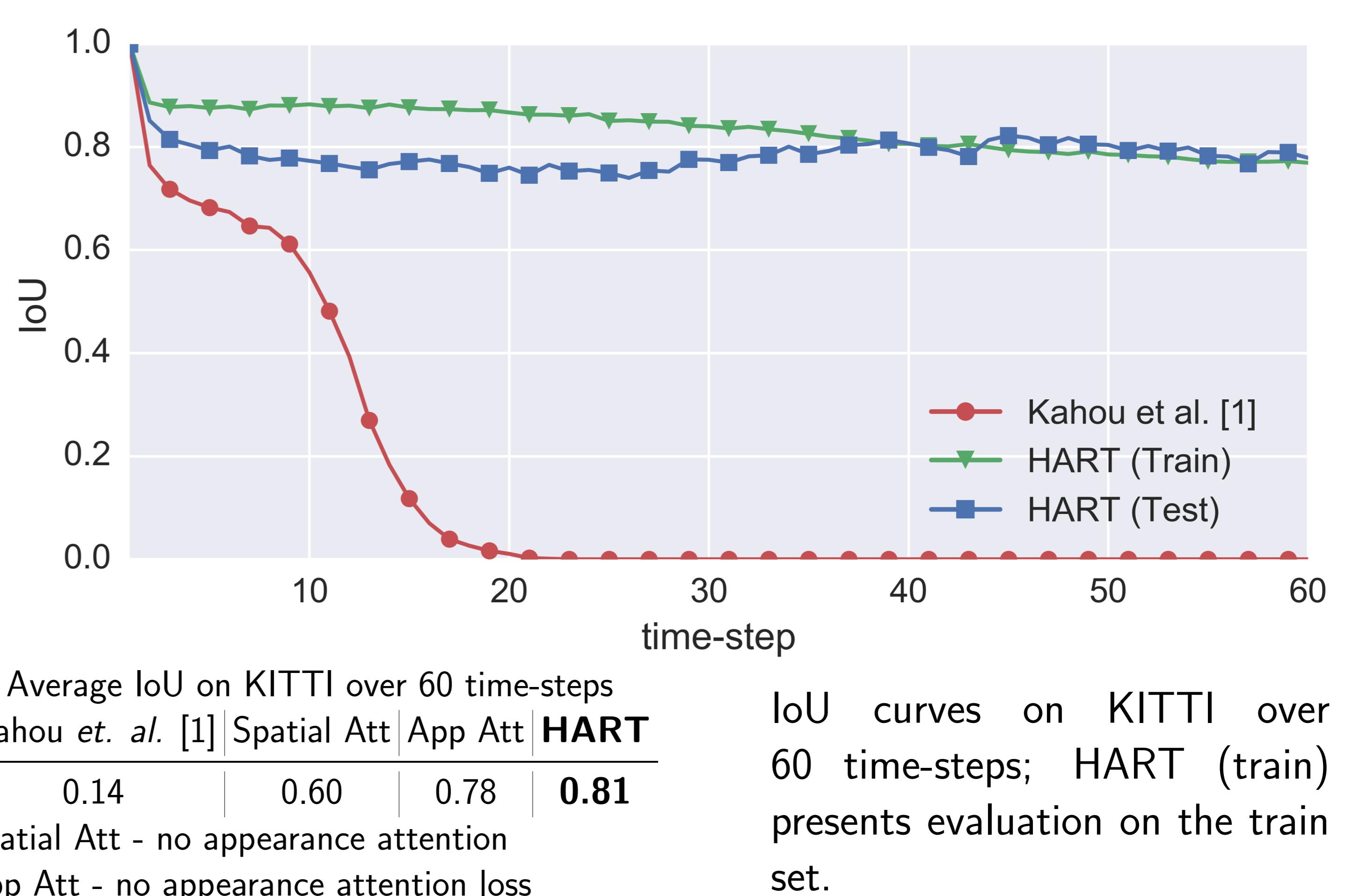


Left to right: glimpses and segmentations learned with and without appearance loss. Attention loss leads to distractor suppression.

Pedestrian Tracking: KTH Dataset [2]



Scaling to Real-World Data: KITTI [3]



References

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