

Attentional Correlation Filter Network for Adaptive Visual Tracking

Jongwon Choi¹, Hyung Jin Chang², Sangdoo Yun¹, Tobias Fischer², Yiannis Demiris², Jin Young Choi¹

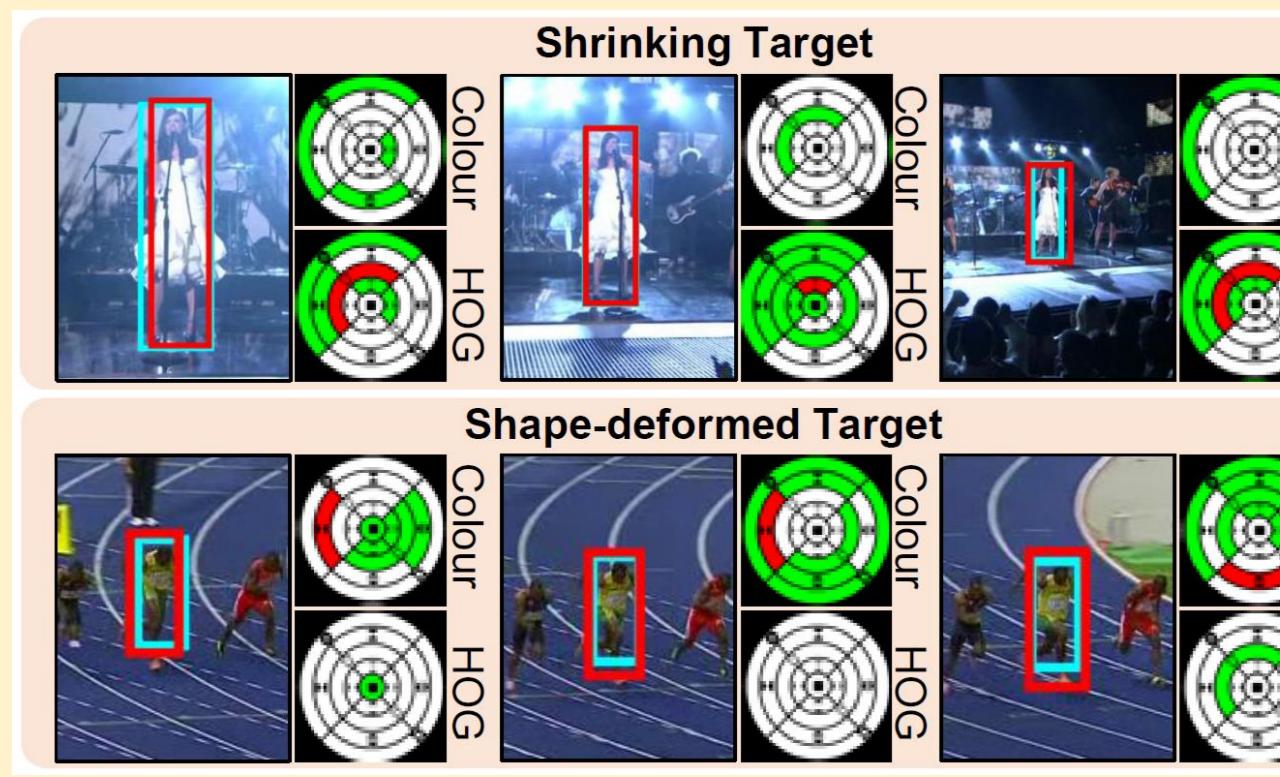
jwchoi.pil@gmail.com, {yunsd101, jychoi}@snu.ac.kr, {hj.chang, t.fischer, y.demiris}@imperial.ac.uk

¹Dept. of EC. Eng., ASRI, Seoul National Univ., South Korea. ²Dept. of EE. Eng., Imperial College London, UK.

Target Problems

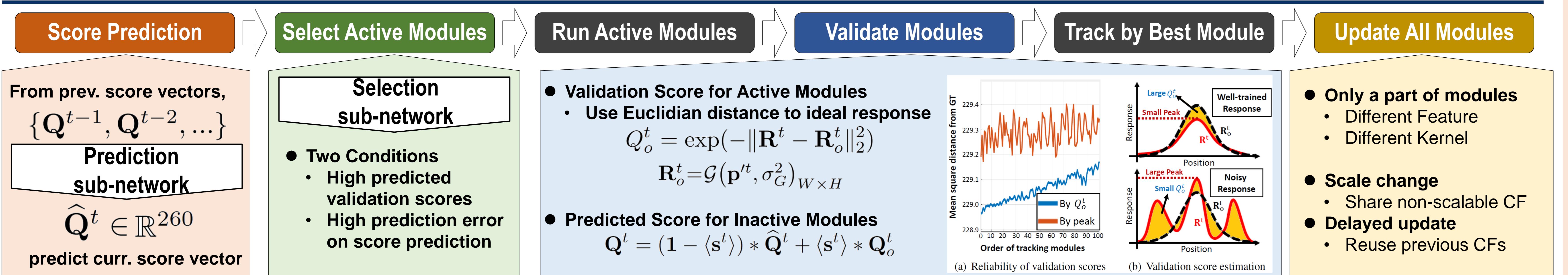
- By using many properties, tracking performance can be improved
- But, needs much time to consider various properties of target

Approach & Contribution



- **Attentional Correlation Filter Network**
 - Attention Network
 - >> Predict the module-wise performance
 - >> Select the attentional modules
 - Correlation Filter Network
 - >> A lot of tracking modules with different properties
 - >> Novel properties (flexible aspect ratio, delay etc.)

Tracking Step

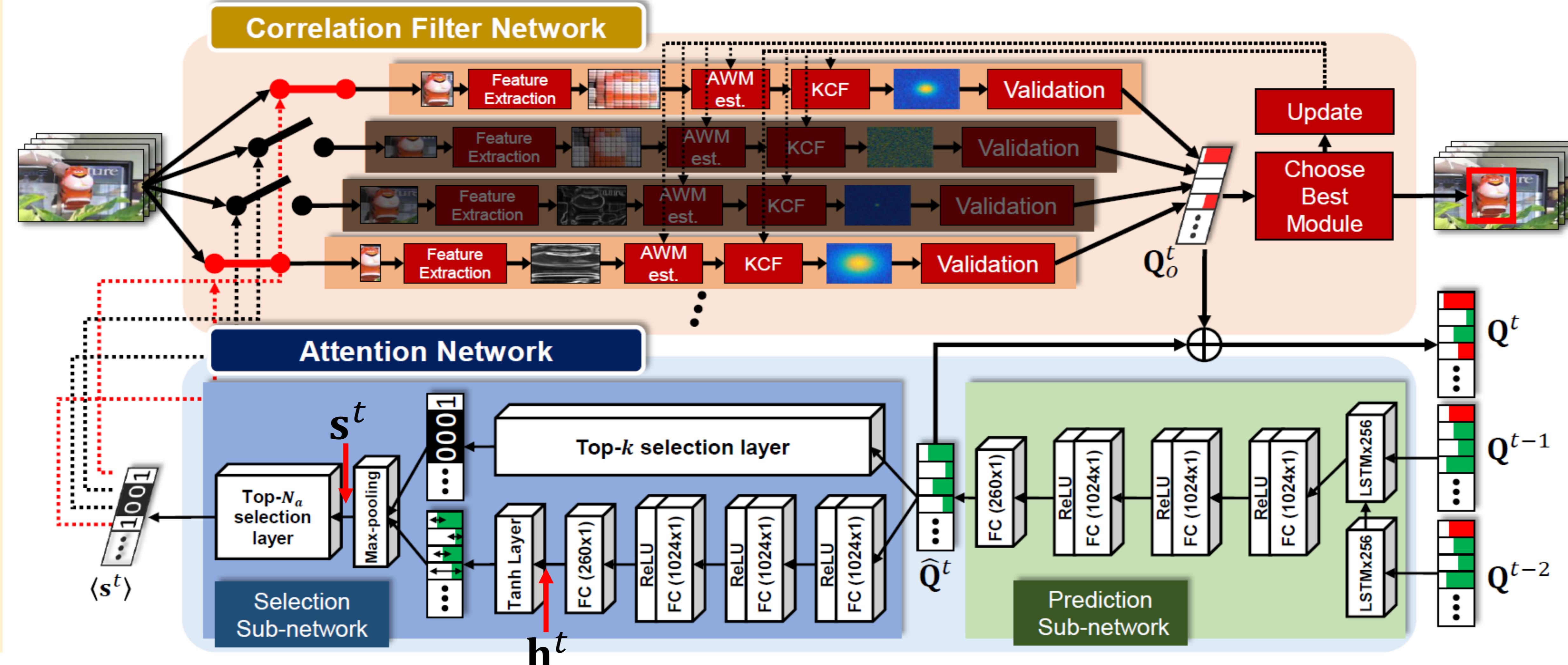


Correlation Filter Network

260 Tracking Modules

- Each tracking module is AtCF [1]
- 2 Features (Color intensity, HOG)
- 2 Kernel types (Gaussian, Polynomial)
- 13 Flexible scale changes (-2x, -x, +x, +2x, -2y, -y, +y, +2y, +xy, +2xy, 0)
- 5 Delayed updates (0, -1, -2, -3, -4 frames)

Overall Framework



Pre-training of Attention Network

Loss Function

$$E = \sum_{i=1}^N \left\{ \|Q(i) - Q_{GT}(i)\|_2^2 + \lambda \|s(i)\|_0 \right\}$$

$$Q(i) = (1 - \langle s(i) \rangle) * Q(i) + \langle s(i) \rangle * Q_{GT}(i)$$

- Prediction sub-network

$$E = \sum_{i=1}^N \left\{ \|\hat{Q}(i) - Q_{GT}(i)\|_2^2 \right\}$$

Relaxation

$$E = \sum_{i=1}^N \left\{ \|(1 - s(i)) * (\hat{Q}(i) - Q_{GT}(i))\|_2^2 + \lambda \|s(i)\|_0 \right\}$$

$$Q(i) = (1 - s(i)) * \hat{Q}(i) + s(i) * Q_{GT}(i)$$

- Selection sub-network

$$E = \sum_{i=1}^N \left\{ \|(1 - s(i)) * (\hat{Q}(i) - Q_{GT}(i))\|_2^2 + \lambda \ln(1 + \|h(i)\|_1) \right\}$$

Experiment

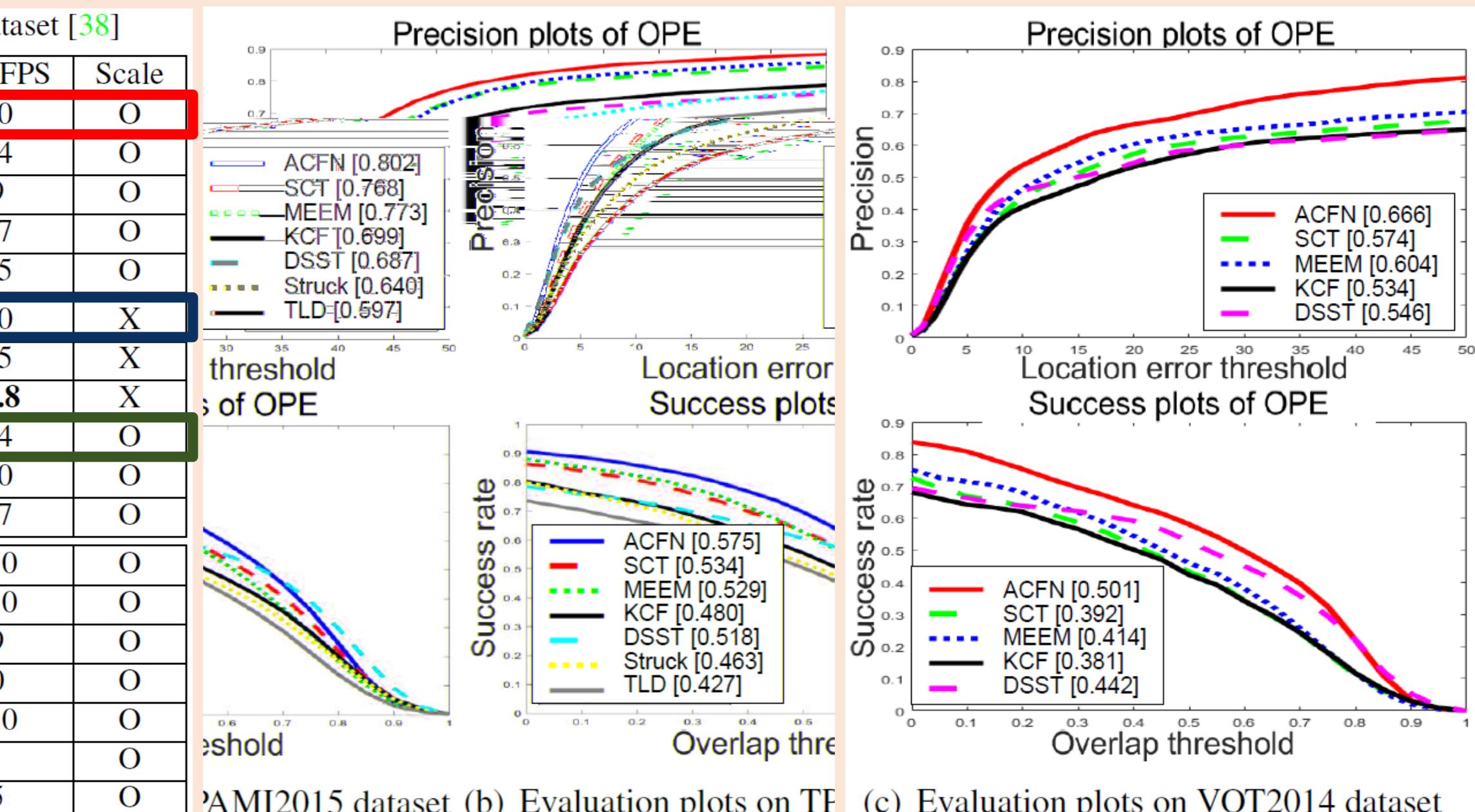
Implementation

- Tensorflow (CF-Net) + MATLAB (At-Net) (By socket communication)
- i7-6900K CPU, 32GB RAM, NVIDIA GTX1070 GPU

Quantitative Results

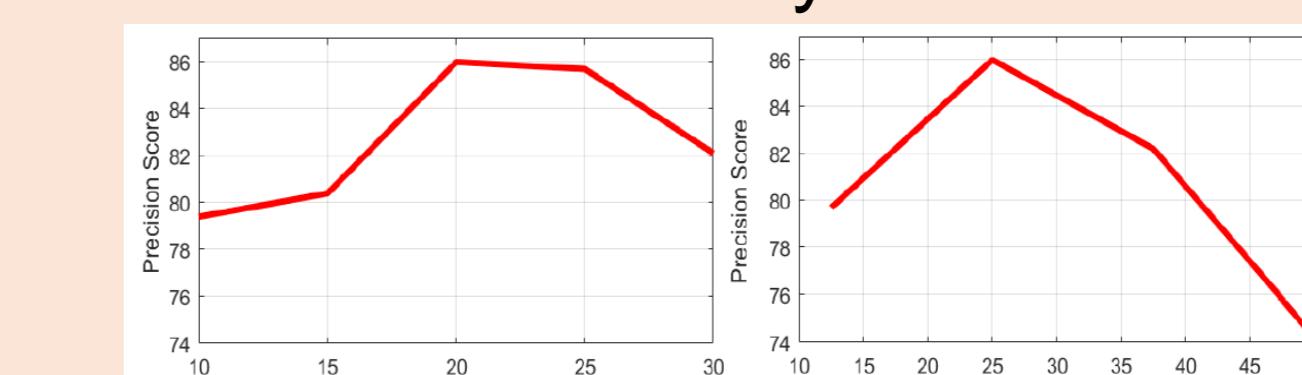
Table 1. Quantitative results on the CVPR2013 dataset [38]

Proposed	Algorithm	Pre. score	Mean FPS	Scale
	ACFN	86.0%	15.0	O
	CFN+predNet	82.3%	14.4	O
	CFN	81.3%	6.9	O
	CFN+simpleSel.	79.4%	15.7	O
	CFN	78.4%	15.5	O
	SCT [3]	84.5%	40.0	X
	MEEM [42]	81.4%	19.5	X
	KCF [16]	74.2%	223.8	X
	DSST [5]	74.0%	25.4	O
	Struck [15]	65.6%	10.0	O
	TLD [19]	60.8%	21.7	O
Real-time	C-COT [8]	89.9%	<1.0	O
	MDNet-N [29]	87.7%	<1.0	O
	MUSTER [18]	86.5%	3.9	O
	FCNT [35]	85.6%	3.0	O
	D-SRDCF [6]	84.9%	<1.0	O
	SRDCF [7]	83.8%	5	O
	STCT [36]	78.0%	2.5	O
Non Real-time	All Frames	86	Ratio for the number of active modules (%)	Precision Score
	Enlarging Frames	86	Ratio for the number of active modules (%)	Precision Score
	Shrinking Frames	86	Ratio for the number of active modules (%)	Precision Score
	Failure Scenes	86	Ratio for the number of active modules (%)	Precision Score

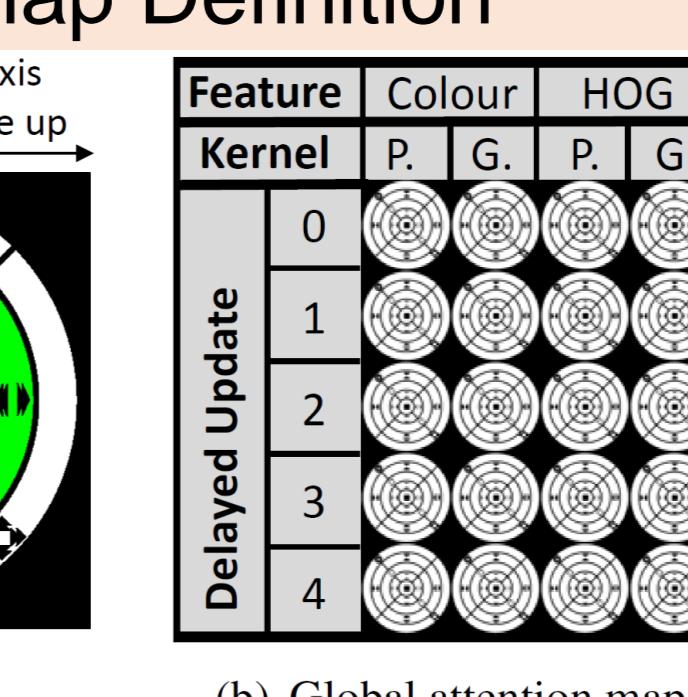


Analysis

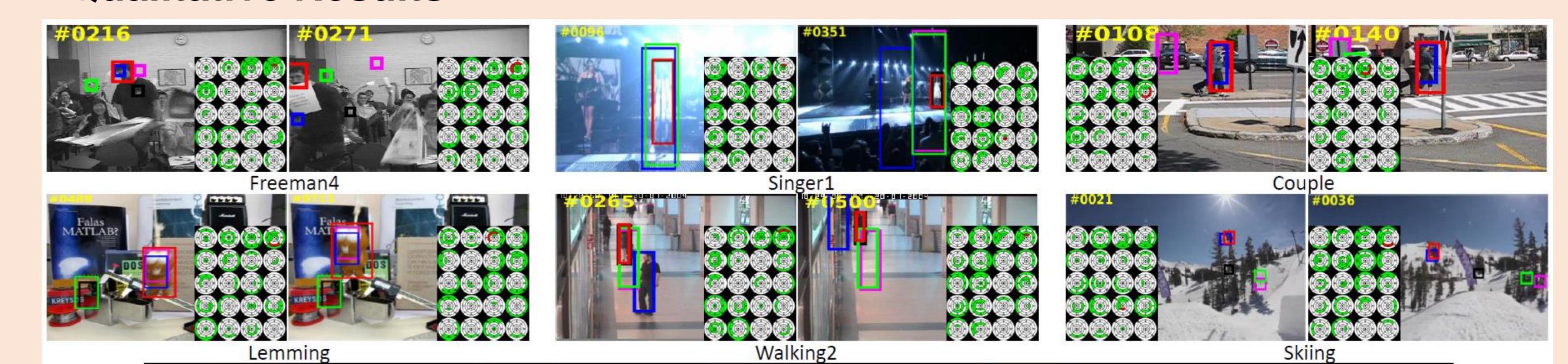
Parameter Analysis



Attention Map Definition



Qualitative Results



Reference

[1] Choi et al., "Visual tracking using attention-modulated disintegration and integration", CVPR2016