# Attentional Correlation Filter Network for Adaptive Visual Tracking

Jongwon Choi [1], Hyung Jin Chang [2], Sangdoo Yun [1], Tobias Fischer [2], Yiannis Demiris [2], Jin Young Choi [1]

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

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

CVPR July 21-26 2017

Codes & Results are available.
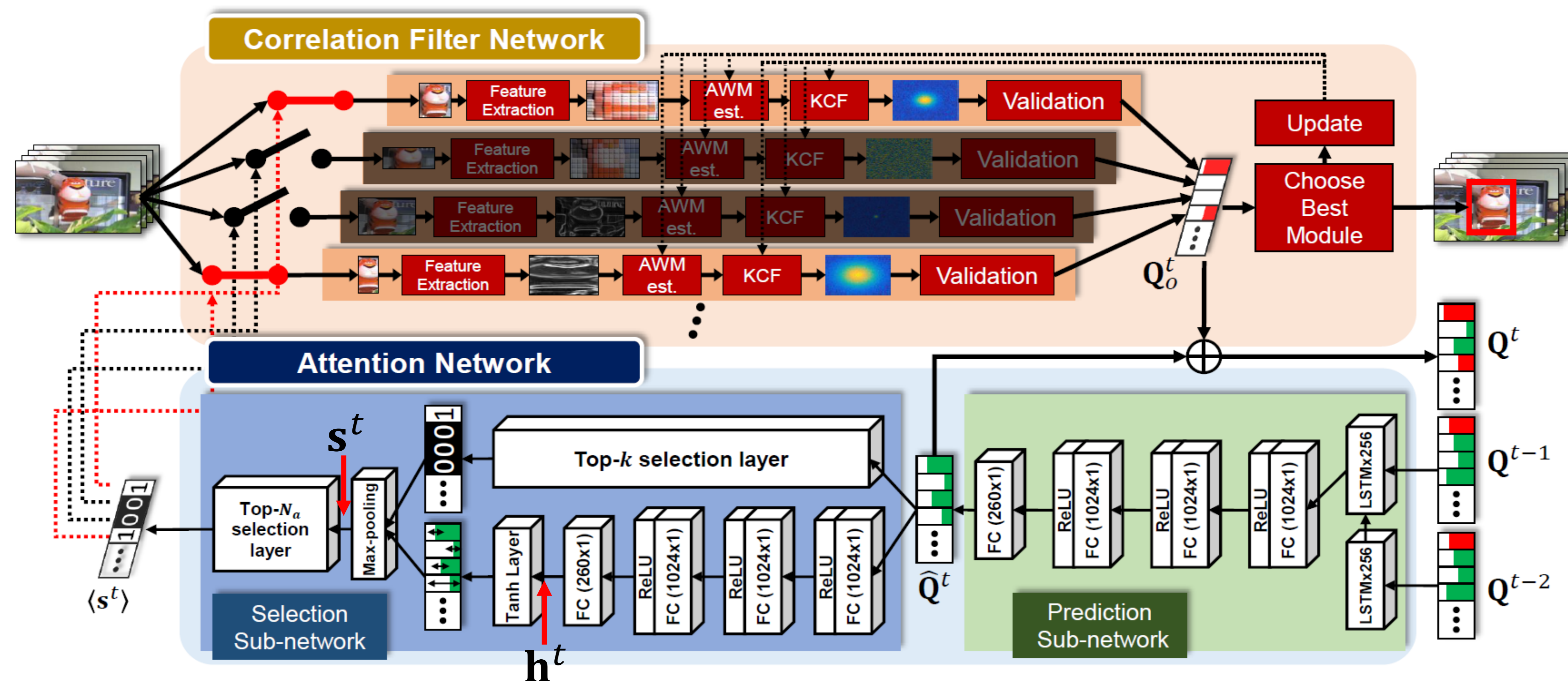**homepage**: https://sites.google.com/site/jwchoivision

## Target Problems

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

## Approach & Contribution
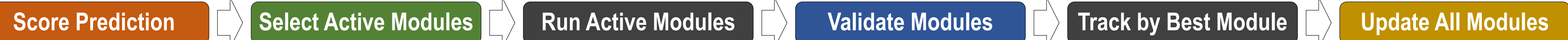
Shrinking Target

Shape-deformed Target

- **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.)

## Overall Framework



Correlation Filter Network

Feature Extraction → AWM est. → KCF → Validation

Attention Network

Top-k selection layer

$s^t$ 0001

Top-$N_a$ selection layer / Tanh Layer / FC (260x1) / ReLU FC (1024x1) ... Max-pooling

⟨$s^t$⟩  1001

Selection Sub-network

$h^t$

$\hat{Q}^t$

FC (260x1) / ReLU FC (1024x1) / ReLU FC (1024x1) / ReLU FC (1024x1) / LSTMx256

Prediction Sub-network

Update → Choose Best Module

$Q_o^t$

$Q^t$
$Q^{t-1}$
$Q^{t-2}$

LSTM256

## Tracking Step

Score Prediction → Select Active Modules → Run Active Modules → Validate Modules → Track by Best Module → Update All Modules

**Score Prediction**

From prev. score vectors,

$\{\mathbf{Q}^{t-1}, \mathbf{Q}^{t-2}, \ldots\}$

Prediction sub-network

$\hat{\mathbf{Q}}^t \in \mathbb{R}^{260}$

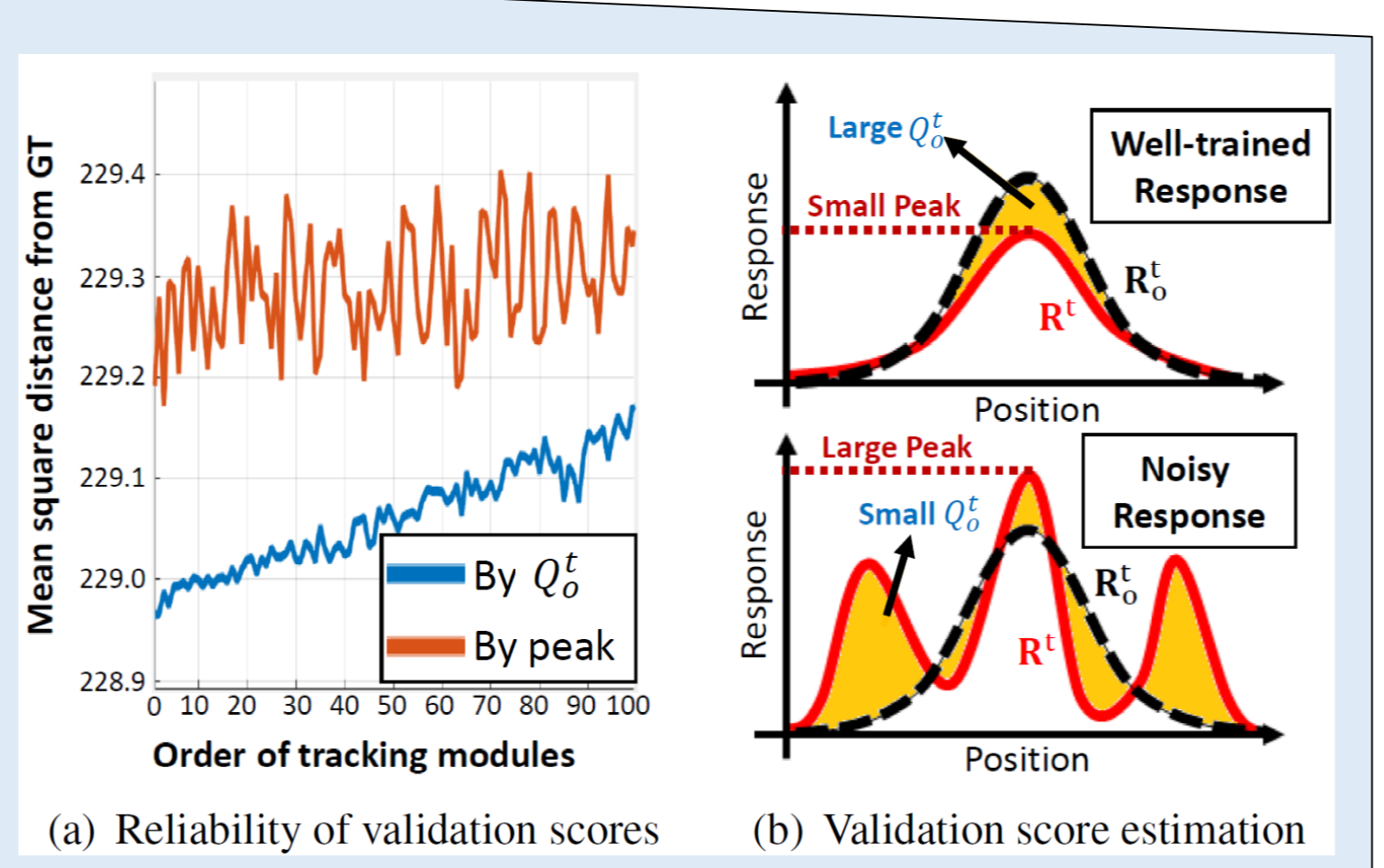predict curr. score vector

**Select Active Modules**

Selection sub-network

- Two Conditions
  - High predicted validation scores
  - High prediction error on score prediction

**Validate Modules**

- Validation Score for Active Modules
  - Use Euclidian distance to ideal response

$Q_o^t = \exp(-\|\mathbf{R}^t - \mathbf{R}_o^t\|_2^2)$

$\mathbf{R}_o^t = \mathcal{G}(\mathbf{p}'^t, \sigma_G^2)_{W \times H}$

- Predicted Score for Inactive Modules

$\mathbf{Q}^t = (\mathbf{1} - \langle\mathbf{s}^t\rangle) * \hat{\mathbf{Q}}^t + \langle\mathbf{s}^t\rangle * Q_o^t$

(a) Reliability of validation scores — By $Q_o^t$ / By peak — Mean square distance from GT, Order of tracking modules

(b) Validation score estimation — Large $Q_o^t$ Small Peak, Well-trained Response, $\mathbf{R}_o^t$; Large Peak Small $Q_o^t$, Noisy Response, $\mathbf{R}_o^t$

**Track by Best Module**

- Only a part of modules
  - Different Feature
  - Different Kernel
- Scale change
  - Share non-scalable CF
- Delayed update
  - Reuse previous CFs

## 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)

## Pre-training of Attention Network

- **Loss Function**

$E = \sum_{i=1}^{N} \{\|\mathbf{Q}(i) - \mathbf{Q}_{GT}(i)\|_2^2 + \lambda\|\mathbf{s}(i)\|_0\}$
$\mathbf{Q}(i) = (1 - \langle\mathbf{s}(i)\rangle) * \mathbf{Q}(i) + \langle\mathbf{s}(i)\rangle * \mathbf{Q}_{GT}(i)$

**Relaxation**

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

- Prediction sub-network

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

- Selection sub-network

$E = \sum_{i=1}^{N} \{\|(\mathbf{1}-\mathbf{s}(i))*(\hat{\mathbf{Q}}(i)-\mathbf{Q}_{GT}(i))\|_2^2 + \lambda\ln(1+\|\mathbf{h}(i)\|_1)\}$
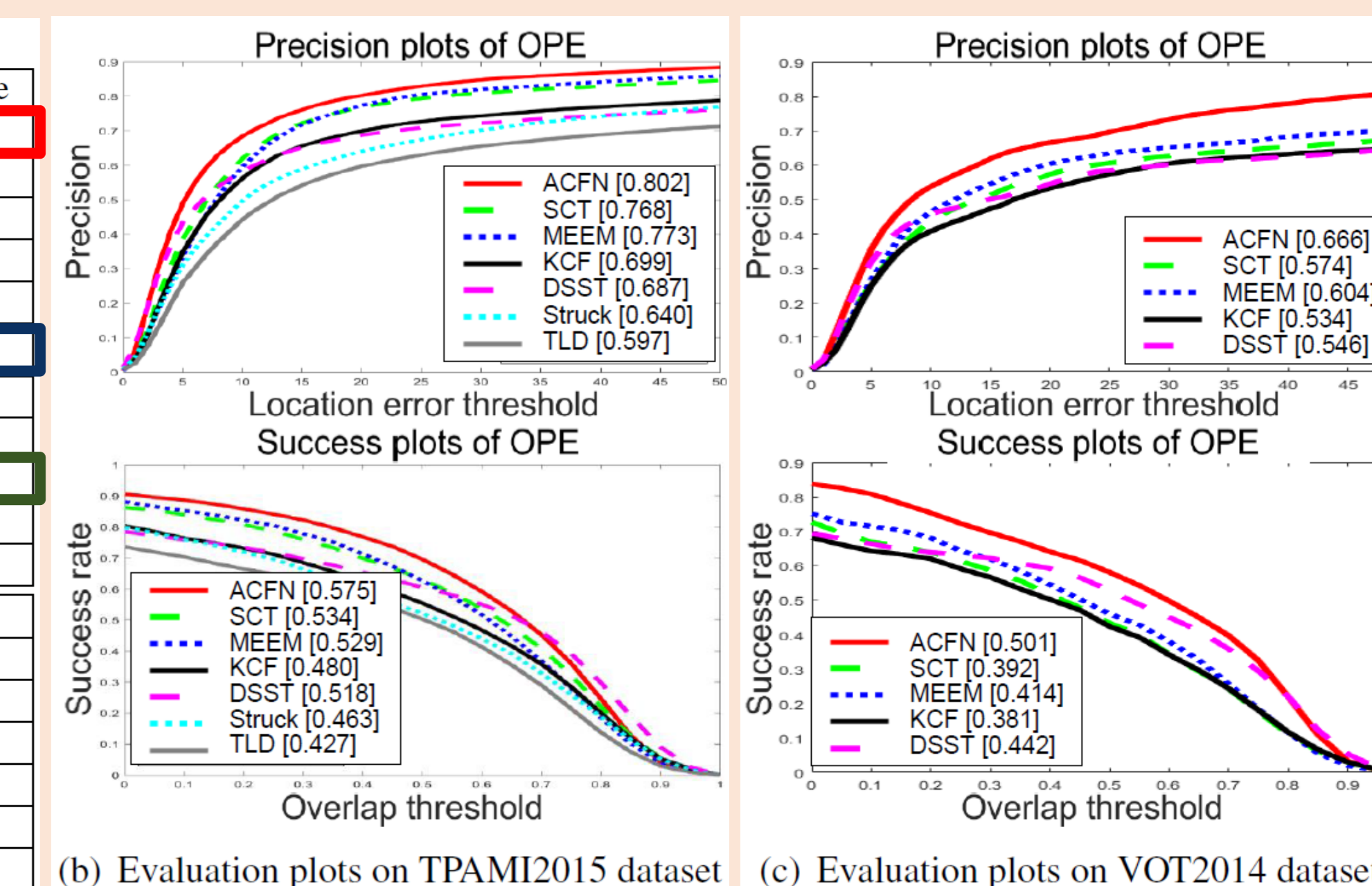
## Experiment

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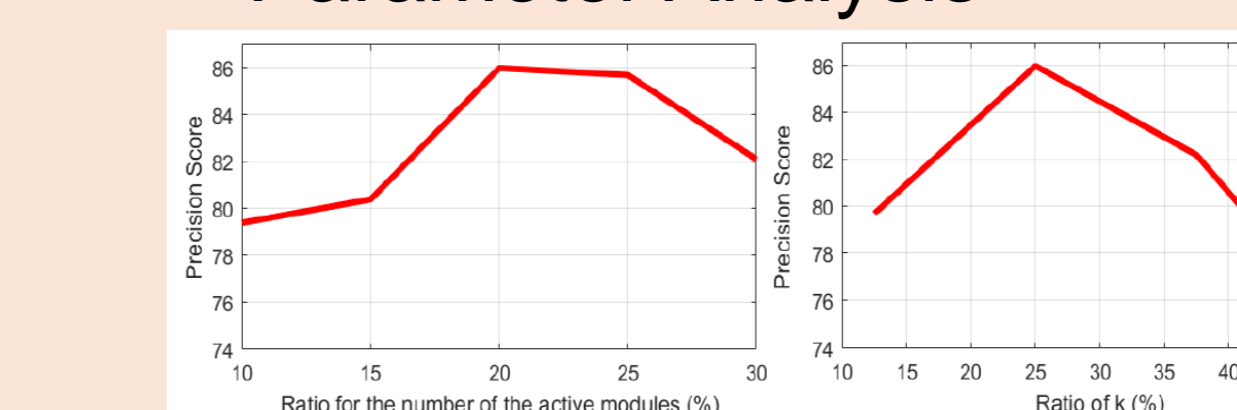  - i7-6900K CPU, 32GB RAM, NVIDIA GTX1070 CPU

- **Quantitative Results**

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

| | Algorithm | Pre. score | Mean FPS | Scale |
|---|---|---|---|---|
| Proposed | 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 |
| Real-time | 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 |
| Non 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 |

Precision plots of OPE
ACFN [0.802], SCT [0.768], MEEM [0.773], KCF [0.699], DSST [0.687], Struck [0.640], TLD [0.597]

Success plots of OPE
ACFN [0.575], SCT [0.534], MEEM [0.529], KCF [0.480], DSST [0.518], Struck [0.463], TLD [0.427]

(b) Evaluation plots on TPAMI2015 dataset

Precision plots of OPE
ACFN [0.666], SCT [0.574], MEEM [0.604], KCF [0.534], DSST [0.546]

Success plots of OPE
ACFN [0.501], SCT [0.392], MEEM [0.414], KCF [0.381], DSST [0.442]

(c) Evaluation plots on VOT2014 dataset

- **Analysis**
  - Parameter Analysis
    Ratio for the number of the active modules (%) / Ratio of k (%)
  - Attention Map Definition

(a) Attention map  (b) Global attention map

- Frequency map for various cases

All Frames / Enlarging Frames / Shrinking Frames / Failure Scenes

Freq. Map of Active Modules / Freq. Map of Best Modules

- **Qualitative Results**



Freeman4 / Singer1 / Couple / Lemming / Walking2 / Skiing

ACFN / SCT / MEEM / KCF / DSST

## Reference

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

Perception and Intelligence Laboratory / imperial.ac.uk/PersonalRobotics / Seoul National University