Action-Decision Networks for Visual Tracking with Deep Reinforcement Learning

Sangdoo Yun, Jongwon Choi, Youngjoon Yoo, Kimin Yun, and Jin Young Choi

Motivation

Visual tracking

- Find the target position in a new frame.
- Deep CNN-based tracking method (Tracking-by-detection)



Problem

- Inefficient search strategy.
- Need lots of labeled video frames to train CNNs.

Approach

Action-driven tracking

• Dynamically capture the target by selecting sequential actions.



Previous frame

CNN-based tracker [1]

Our method

[1] H. Nam and B. Han. Learning multi-domain convolutional neural networks for visual tracking., CVPR 2016.





Reward:

$$(s_t) = \begin{cases} 1, \text{ if } IOU(p_t, G) > 0 \\ -1, \text{ otherwise} \end{cases}$$

$$\Delta W \propto \sum_{t}^{T} \frac{\partial \log(p(a_t|s_t; W))}{\partial W} r_t$$

Experiment setting

- Trained on <u>VOT dataset</u> & evaluated on <u>OTB-100 dataset</u>.
- ADNet $(N_I: 3000, N_O: 250, I: 10) \rightarrow 3$ fps (Prec: 88%)
- ADNet-fast $(N_I: 300, N_O: 50, I: 30) \rightarrow 15$ fps (Prec: 85%)

Analysis on action

93% of total frames have smaller than 5 actions

Self-comparison

- SL: supervised learning
- SS: uses 1/10 gt annotations
- RL: reinforcement learning

OTB-100 test results

	Algorithm	Prec.(20px)
	ADNet	88.0%
	ADNet-fast	85.1%
Non real-time	MDNet [24]	90.9%
	C-COT [9]	90.3%
	DeepSRDCF [8]	85.1%
	HDT [25]	84.8%
	MUSTer [15]	76.7%
Real-time	MEEM [42]	77.1%
	SCT [5]	76.8%
	KCF [13]	69.7%
	DSST [7]	69.3%
	GOTURN [12]	56.5%

Sequential actions selected by ADNet



SEOUL NATIONA UNIVERSITY

Experiments

MatConvNet toolbox, i7-4790K CPU, Nvidia Titan X GPU.

