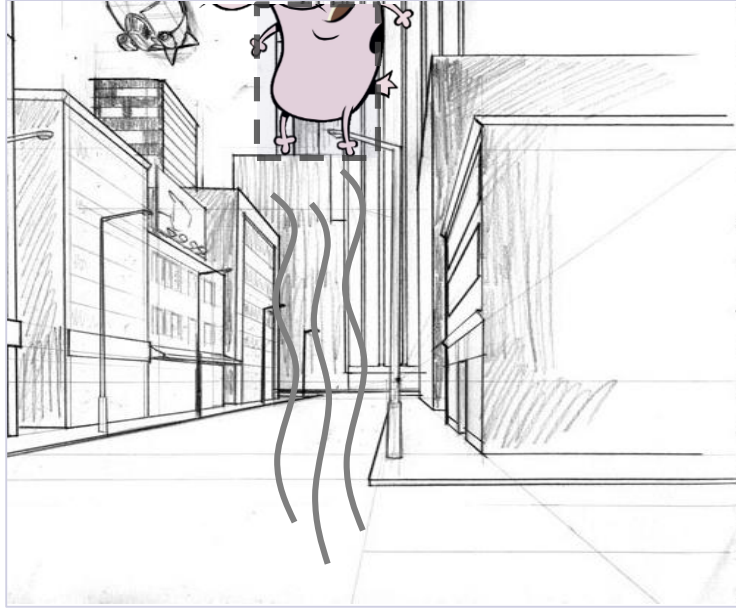# Advanced CV methods
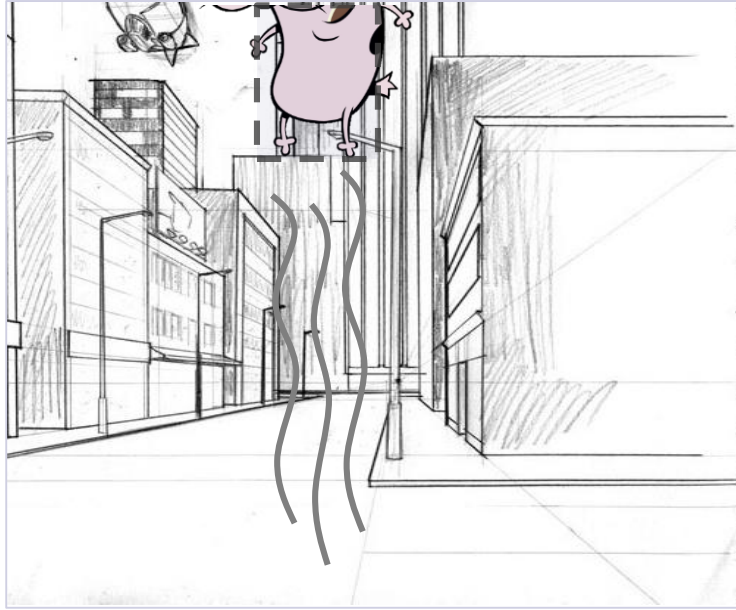# Long-Term tracking

## Matej Kristan

Laboratorij za Umetne Vizualne Spoznavne Sisteme,
Fakulteta za računalništvo in informatiko,
Univerza v Ljubljani

# Long-term tracking (LTT)



- Regardless of how well the visual model is designed, any short-term tracker will eventually fail

- Disappears from the field of view, gets fully occluded, etc.

# Long-term tracking (LTT)



- The general LT tracking properties:

  - Determine when the target has been lost (or disappeared)

  - Re-detect the target after losing the target

  - Update the visual model very carefully to minimize drifting

# Taxonomy: Short-term/long-term spectrum[1]

| | Position reported | Tracking failure detection | Target re-detection |
|---|---|---|---|
| $ST_0$: Basic ST | each frame | no | no |
| $ST_1$: Basic ST with conservative updating | each frame | not explicitly, selective update of visual model | no |
| $LT_0$: Pseudo LT | only when visible | yes | no |
| $LT_1$: Re-detecting LT | only when visible | yes | yes |

- $ST_0$ (e.g., vanilla DCF, MS); $ST_1$ (e.g., MDNet) -> easily converted to $LT_0$

- $LT_1$ most sophisticated, typical composition:

  - Short-term tracker (ST) for frame-to-frame localization

  - Detector for target re-detection

  - Algorithm for interaction between ST and detector

[1] Lukežič, et al., *Now you see me: evaluating performance in long-term visual tracking*, TCyb 2020

# LT1 trackers origin

- Most of the $LT_1$ originate from two main paradigms introduced by *TLD*[1] (aka Predator) and *Alien*[2]

- In the following we will overview the TLD

[1]Kalal, Mikolajczyk, Matas, Tracking-Learning-Detection, TPAMI2010

[2]Pernici, F. and Del Bimbo, A., Object Tracking by Oversampling Local Features, TPAMI2013

Advanced computer vision methods

# TRACKING BY TRACKING, LEARNING, DETECTION (PREDATOR)

# Tracking learning detection: TLD aka Predator[1]

- Detector is the main component

- It's all about robust detector updating

- Run Detector and ST tracker in parallel

- Use the ST and Detector output to construct training samples for Detector

Short-term:

A flock of flows

Detector:

Grayscale NCC

[1]Kalal, Mikolajczyk, Matas, Tracking-Learning-Detection, TPAMI2010

# Fast-forward… "TLD in action"



Kalal, Mikolajczyk, Matas, Tracking-Learning-Detection, TPAMI2010

# The short-term tracker



- A "cell" grid of Lucas-Kanade trackers

- Each LK tracker has a reliability estimate

- Robustly estimates motion from 50% of most reliable displacements

  (could also use a robust estimator, e.g., RANSAC)

- 2 layers of Pyramidal LK tracker

  with $10 \times 10$ pixels patches.

- Fairly robust frame-to-frame localization

  in absence of severe occlusion

Z. Kalal, K. Mikolajczyk, and J. Matas. Forward-Backward Error: Automatic Detection of Tracking Failures. ICPR, 2010
Improved version:
T. Vojir and J. Matas. Robustifying the flock of trackers. CVWW2011

# The detector visual model

- Appearance model: a grayscale patch

- Bounding box with fixed aspect
  (only scale changes, proportions constant)

- Patch resampled into 15x15 size

- Object model is a collection of multiple positive and negative patches!

- Forget patches (randomly) to keep the number of patches low enough
  (memory and speed efficiency)



15pixels

15pixels

Model:

Positive exemplar patches:

...

Negative exemplar patches:

...

# The detector application

- A scanning window

- Compare patches using a normalized cross correlation (NCC)

- A nearest-neighbor classifier using the NCC score

- Problem: A brute force would require comparing all locations with all patches in the model!

- Solution: Apply cascaded approach that quickly rejects many potential image locations by using simple and fast features.



Positive exemplar patches:

Negative exemplar patches:

1-NN classifier

Fast classifiers with low FP/FN, high TP



Patch variance

Ensemble classifier

1-NN classifier

Accepted patches

Rejected patches

# The ST-Detector interaction algorithm

- PN learning: Responsible for training the Detector

- PN (semi-supervised) learning assumptions:

  - Two classes of labelling processes are available: P and N

  - "P" proposes positive, the "N" proposes negative examples only.

  - Both processes are noisy and can make mistakes

  - By carefully addressing the conflicts between the two labelling processes,
    a long-term stability is achieved.

# Interaction algorithm P-event: "Loop"

- Guideline: *Do not trust the learning examples until you are absolutely sure about their labels!*

- Exploits temporal structure

- Assumption: If an adaptive tracker fails, it is unlikely to recover.

- Rule: Patches from a track starting and ending in the current model (red), i.e. are validated by the detector, are added to the model.



Loop example

Failure example

Detector

Short-term component

# Interaction algorithm N-event: "Uniqueness"

- Exploits spatial structure

- Assumption:

  Object is unique in a single frame

  (no other object looks alike)

- Rule: If the tracked patch is

  in the model, all other detections

  within the current frame (red) are

  assumed wrong

  → *are pruned from the model*

# Interaction algorithm: Model learning

Defined by:

- P-events, N-events, detector learning method

- P and N events are defined in terms of tracker and detector outputs

# TLD tracking-learning example



Detector templates (positives)

# TLD tracking example

# TLD summary

- **PN Learning trains a robust detector** by observing the object of interest
  (no a priori labelled training data, no constraints on the video)

- Detector **improves over time** (experimentally validated)

- A stable semi-supervised learning algorithm

- Matlab/C++ implementation runs at > 20 fps (back in 2010)


- Code is available online:

  http://personal.ee.surrey.ac.uk/Personal/Z.Kalal/

  Kalal, Mikolajczyk, Matas, Tracking-Learning-Detection, TPAMI2010

# Long-Term Architecture Implementation Issues

| Tracker | Short-term tracker | Detector | Interaction |
|---------|-------------------|----------|-------------|
| Alien [6] | Keypoints (SIFT) | Keypoints (SIFT) | F-B, Ransac |
| TLD [1] | Optical flow | Random forest | P-N learning |
| MUSTER [2] | Correlation filter | Keypoints (SIFT) | F-B, Ransac |
| LCT [3] | Correlation filter | Random fern | K-NN, response thresh. |
| CMT [4] | Keypoints (flow) | Keypoints (static) | F-B, clustering |
| PTAV [5] | Correlation filter | CNN (Siam. Net.) | CNN confidence score |

Approaches from different methodologies

- Prohibits tight interaction e.g., feature/model sharing
- Leads to complicated implementation

[1] Kalal et al., Tracking-Learning-detection, TPAMI 2010
[2] Ma et al., Long-Term Correlation Tracking, CVPR 2015
[3] Hong et al., MUlti-Store Tracker (MUSTer): a Cognitive Psychology Inspired Approach to Object Tracking, CVPR 2015
[4] Nebehay et al., Clustering of Static-Adaptive Correspondences for Deformable Object Tracking, CVPR 2015
[5] Fan et al., Parallel Tracking and Verifying: A Framework for Real-Time and High Accuracy Visual Tracking, ICCV 2017
[6] Pernici, F. and Del Bimbo, A., Object Tracking by Oversampling Local Features, TPAMI2013
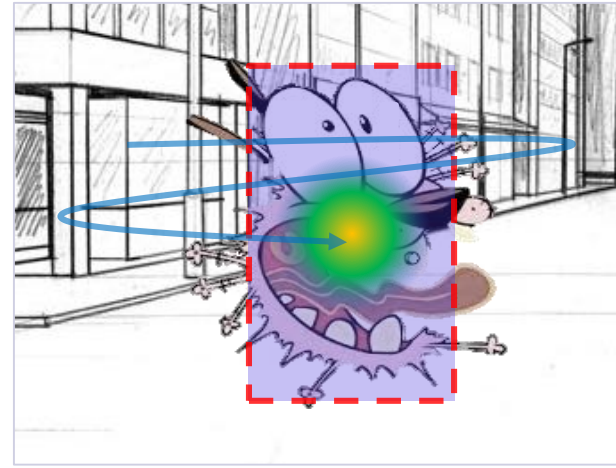
# Long-Term Architecture Implementation Issues

| Tracker | Short-term tracker | Detector | Interaction |
|---|---|---|---|
| Alien [6] | Keypoints (SIFT) | Keypoints (SIFT) | F-B, Ransac |
| TLD [1] | Optical flow | Random forest | P-N learning |
| MUSTER [2] | Correlation filter | Keypoints (SIFT) | F-B, Ransac |
| LCT [3] | Correlation filter | Random fern | K-NN, response thresh. |
| CMT [4] | Keypoints (flow) | Keypoints (static) | F-B, clustering |
| PTAV [5] | Correlation filter | CNN (Siam. Net.) | CNN confidence score |
| FCLT [7] | Correlation filter | Correlation filter | Correlation uncertainty |

Shared target representation: tight interaction, efficient implementation

- Short-term tracker and a detector within a single methodology
- A single DCF learner, two interacting models

[7] Lukežič, Čehovin, Vojir, Matas, Kristan, *FuCoLoT -- A Fully-Correlational Long-Term Tracker*, ACCV 2018

# Fully Correlational Long-term Tracker (FCLT)



- Discriminative correlation filter in two separate components.

- Detector activated when ST not confident.

- Motion model used with detector.

Short-term:       Detector:

correlation filter

[1]Lukežič et al., *Discriminative Correlation Filter Tracker with Channel and Spatial Reliability*, IJCV 2018

# FCLT: ST and Detector learning

- Short-term (ST) model is a CSRDCF[1] with standard update

- Detector:

  - Standard DCF cannot be used for image-wide detection

  - Utilize constrained learning from CSRDCF[1] from a wider region

  - Several object models updated at various time scales

Detector 1:        Detector 2:        Detector 3:        Detector N:



Never update      Update every 250th      Update every 50th      Update every frame

[1]Lukežič, Vojir, Čehovin Zajc, Matas and Kristan, *Discriminative Correlation Filter Tracker with Channel and Spatial Reliability*, IJCV 2018

# FCLT: Detector application



Correlation response

Motion consistency

Final response

Final target candidate position

Low values          High values

Detector 1:  Detector 2:  Detector 3:  Detector N:

Target last seen here.

If not detected: Cycle through N detectors and scales in subsequent frames.
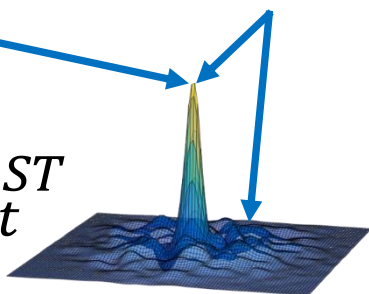
# FCLT : ST tracking failure detecton

- Reliability score $q_t$
  on correlation response $R_t^{ST}$

$$q_t = MAX(R_t^{ST}) \times PSR(R_t^{ST})$$

$$* H_t^{ST} = R_t^{ST}$$

- Threshold on the ratio: $\dfrac{\overline{q_t}}{q_t}$

  $\overline{q_t}$ is mean over past frames

- When failure detected:
  - Activate detector
  - Stop updating visual model

# Example: Tracking with FCLT



Short-term tracker

Detector

Tracking uncertainty

Lukežič, Čehovin, Vojir, Matas, Kristan, *FuCoLoT -- A Fully-Correlational Long-Term Tracker*, ACCV 2018

# Redetection capability (LT$_0$ vs LT$_1$)

FCLT[1]

MDNet[2]



Re-detects after target re-appears

Never recovers after drift

[1] Lukežič, Čehovin, Vojir, Matas, Kristan, *FuCoLoT -- A Fully-Correlational Long-Term Tracker*, ACCV2018
[2] Nam, Han, Learning, Multi-Domain Convolutional Neural Networks for Visual Tracking, CVPR2016
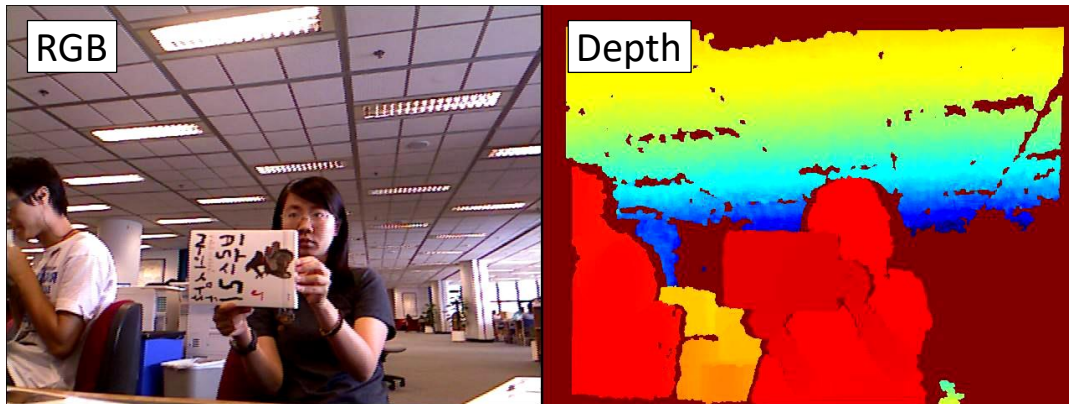
# Extension of D3S to LT setup

- Similar to FCLT, only using DCF from GEM for global re-detection (and few additional upgrades, such as MDNet verifier)



Džubur et al., A Long-Term Discriminative Single Shot Segmentation Tracker, ERK2022
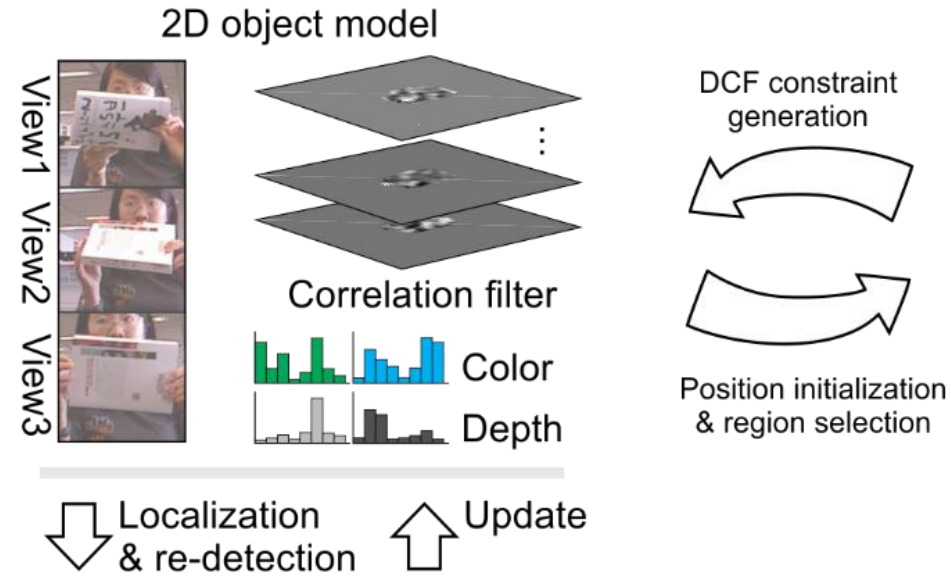
# A 2D Object Assumption in Standard Trackers
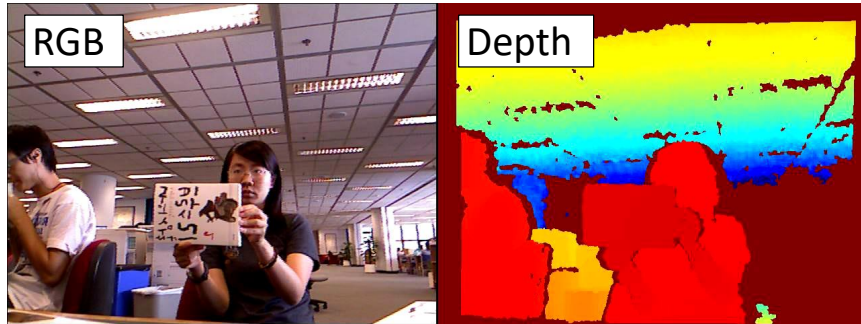
- Existing tracking methods treat a tracked object as a 2D structure

- Problem: Cannot distinguish between pose change and (self)occlusion

# Extension to RGBD tracking

- Extend FCLT by 3D reconstruction to improve occlusion detection



Kart, Lukezic, Kristan, Kämäräinen, Matas, Object Tracking by Reconstruction with View-Specific Discriminative Correlation Filters CVPR2019

# Object tracking by reconstruction (OTR)

- Top performance among all RGBD trackers on PTB [Song et al., ICCV2013] and STC [Xiao et al.] benchmarks.



Kart, Lukezic, Kristan, Kämäräinen, Matas, Object Tracking by Reconstruction with View-Specific Discriminative Correlation Filters CVPR2019

# Recent deep LT developments (2018)

- Region proposal network akin to SSD[1] and SiamRPN[2]

- Verification network, essentially MDNet[3]
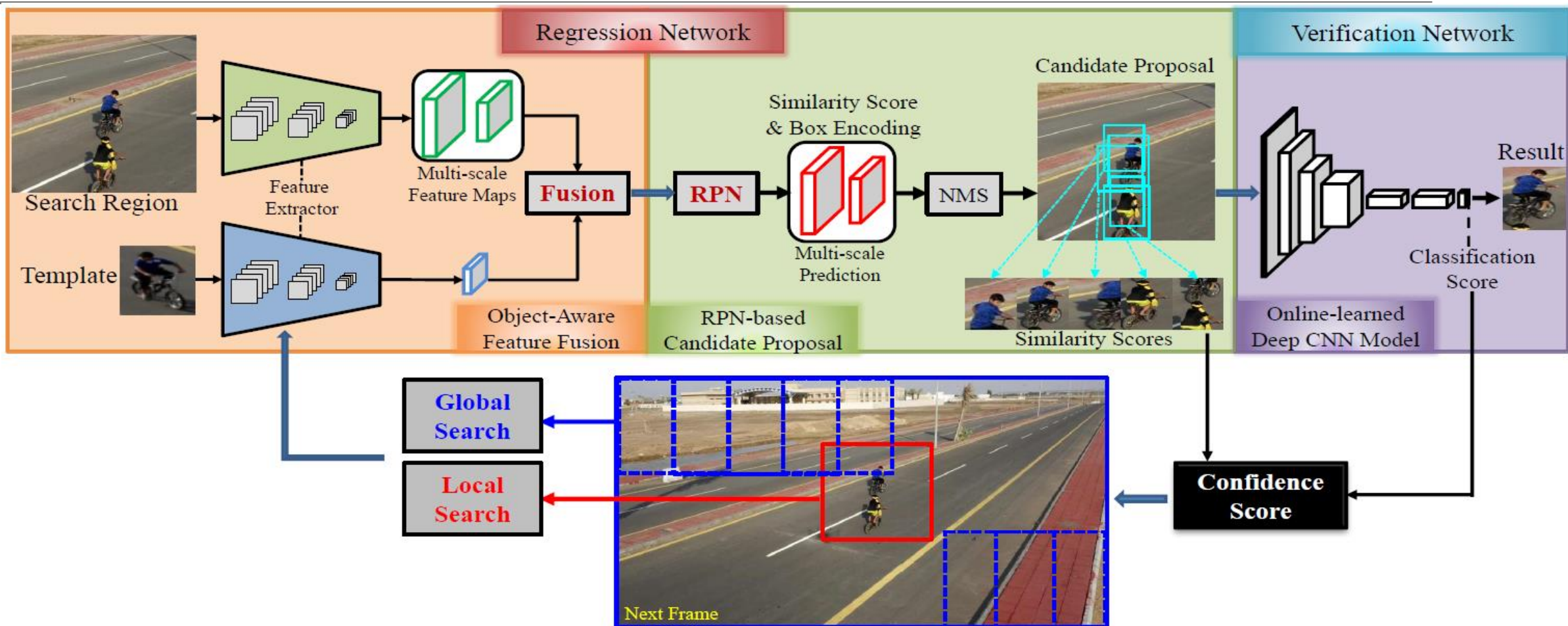
- Interaction akin to FCLT

[1]Liu et al., SSD: Single shot multibox detector, ECCV2016
[2]Li et al., High Performance Visual Tracking with Siamese Region Proposal Network, CVPR2018
[3]Nam et al., Learning multi–domain convolutional neural networks for visual tracking, CVPR2016

Zhang et al., Learning regression and verification networks for long-term visual tracking, ArXiv 2018

# MBMD deep long-term tracker



- Modern state-of-the-art trackers are based on transformers (e.g., STARK-like) with a large localization range + a discriminator like Dimp

Zhang et al., Learning regression and verification networks for long-term visual tracking, ArXiv 2018   https://github.com/xiaobai1217/MBMD

# References

- TLD:
  - Kalal, Z., Mikolajczyk, K. and Matas, J., Tracking-Learning-Detection, IEEE TPAMI2010
  - Page + code: http://personal.ee.surrey.ac.uk/Personal/Z.Kalal/
- FCLT:
  - Lukežič, Čehovin, Vojir, Matas, Kristan, *FuCoLoT -- A Fully-Correlational Long-Term Tracker*, ACCV 2018

# Acknowledgment

- Thanks to Jiri Matas for kindly sharing some of their slides that I used in preparation of this lecture.