

Deep Tracking

Slides by: David Held

This session will be recorded
for educational use by other
students in this course

Types of Tracking

- Single object tracking
- Joint segmentation and tracking
- Multi-object tracking
- 3D tracking
- Dense optical flow
- 3D scene flow
- Self-supervised tracking

Single-Target Tracking

Input: Label in first frame



Output:



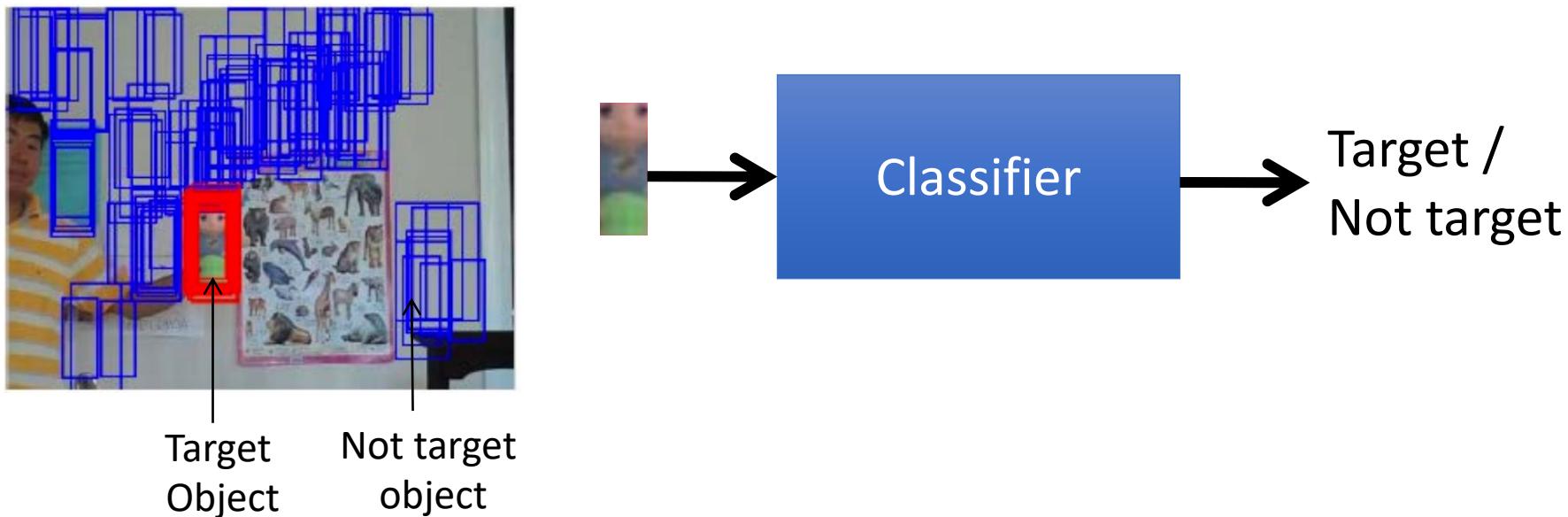
Changes:

- Object motion
- Camera motion
- Deformations
- Lighting changes
- Occlusions

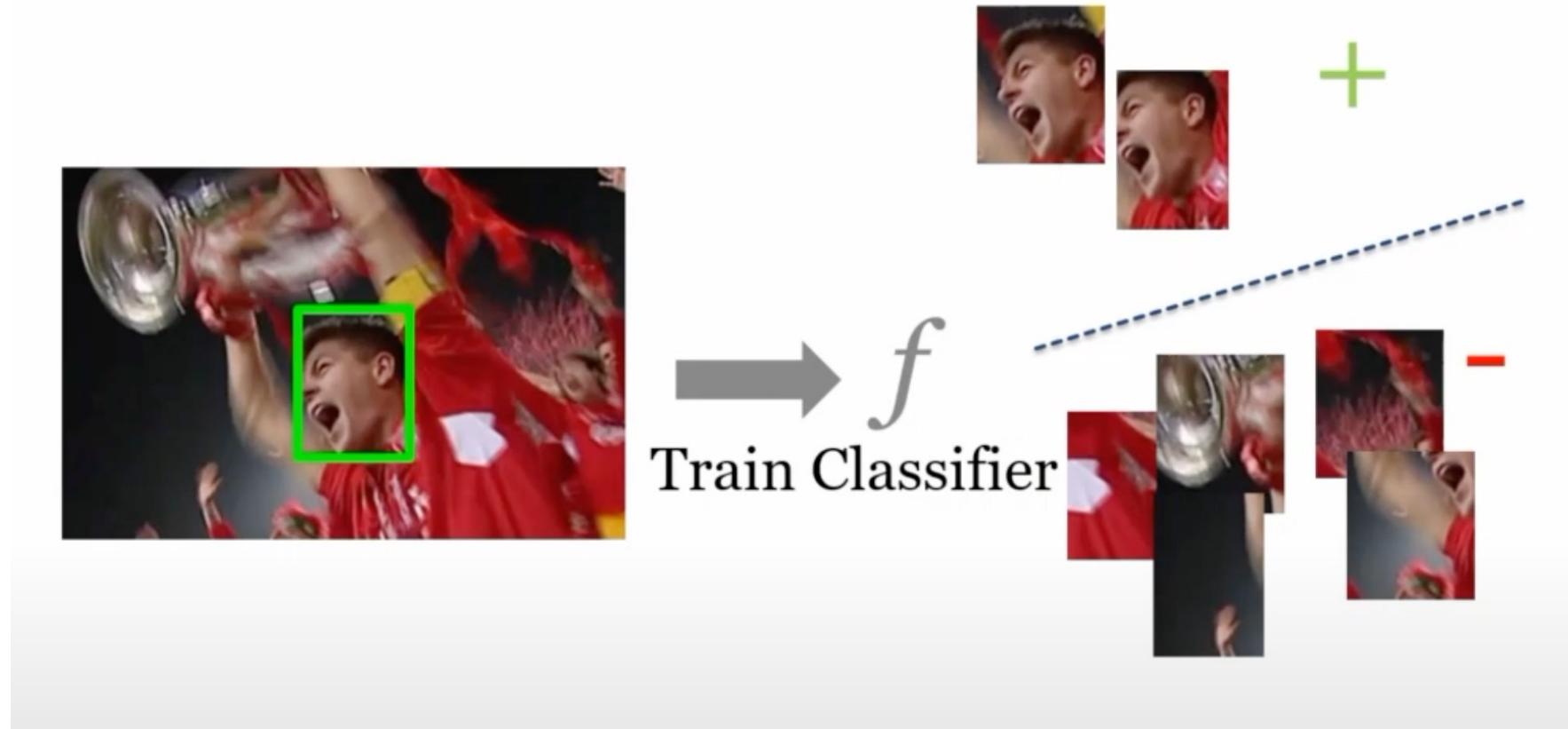
Tracking Novel Objects



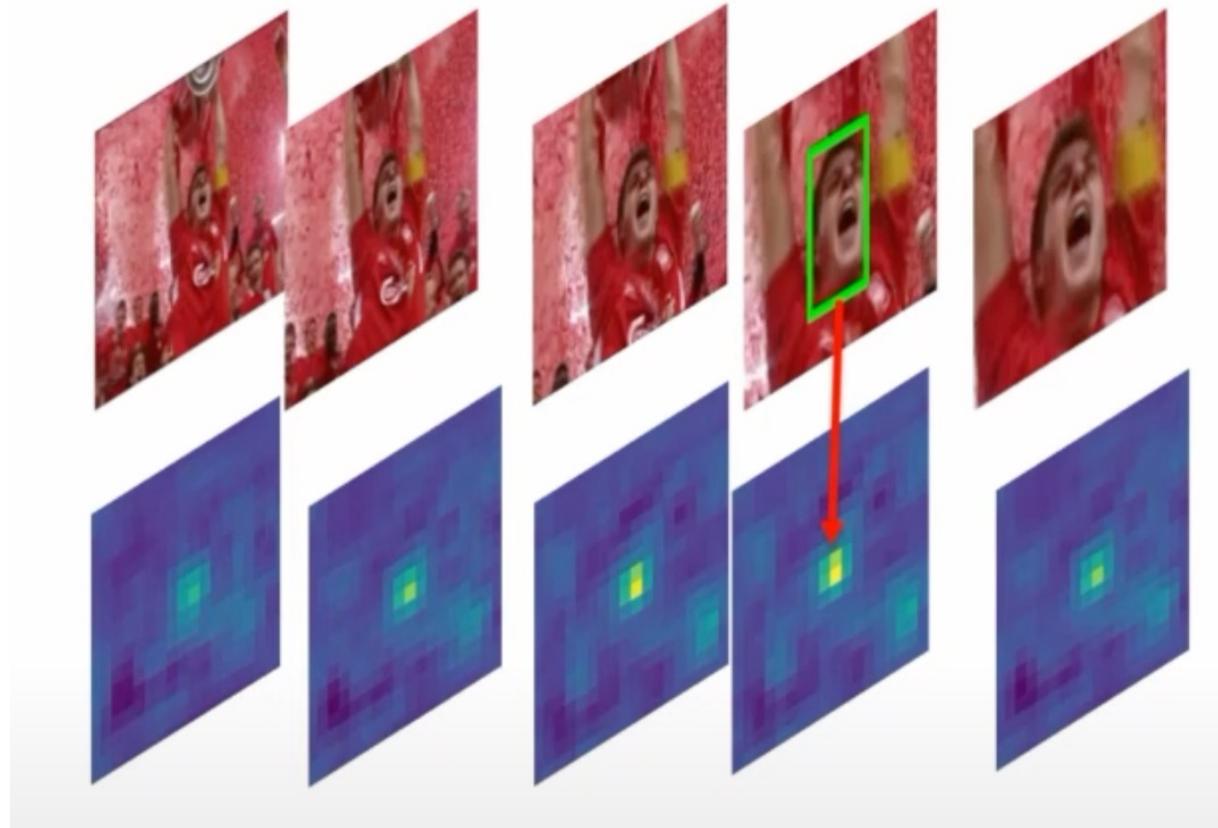
Initial Approaches: Tracking as Classification



Train a “Classifier” to Distinguish Target from Background



On Test Frame, Run Detector over Different Positions and Scales



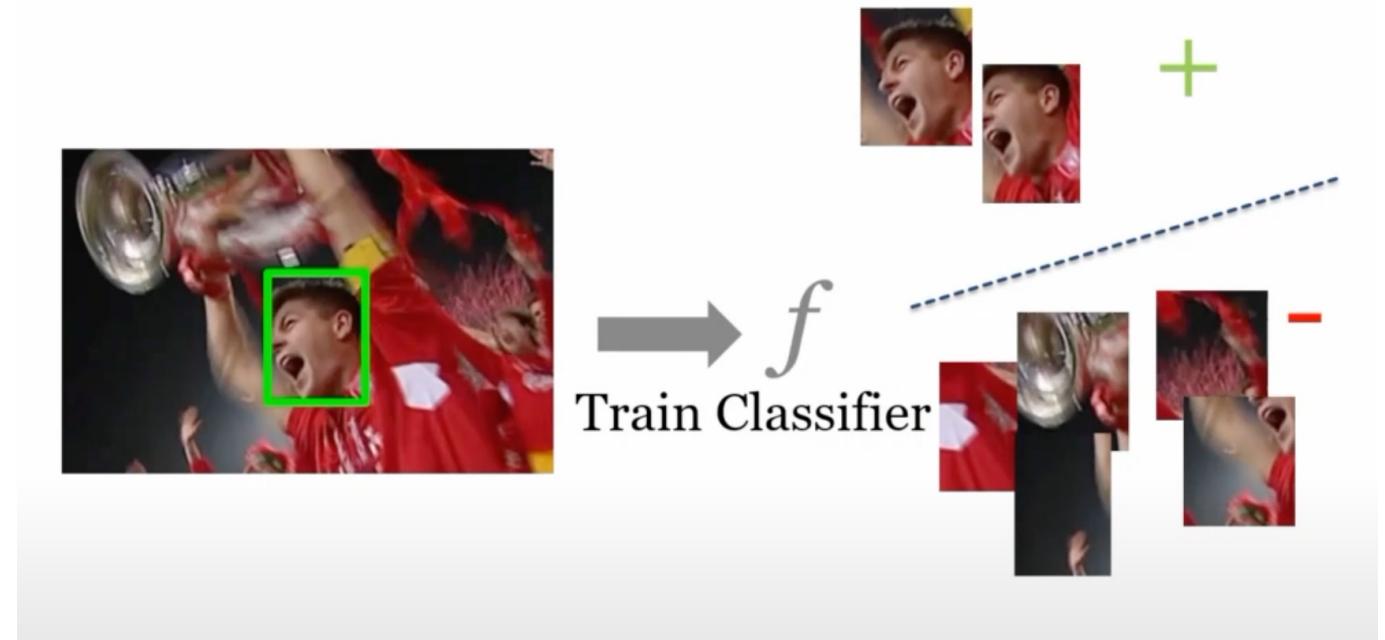
Location and scale with the highest classifier score is the coarse prediction

VOT 2015

Tracker	A	R	$\hat{\Phi}$
Deep DCF → DeepSRDCF*	0.56	1.05	0.32
EBT	0.47	1.02	0.31
SRDCF*	0.56	1.24	0.29
LDP*	0.51	1.84	0.28
sPST*	0.55	1.48	0.28
SC-EBT	0.55	1.86	0.25
NSAMF*	0.53	1.29	0.25
Struck*	0.47	1.61	0.25

Online Train a “Classifier” to Distinguish Target from Background

- What are the limitations of this approach?



VOT 2015

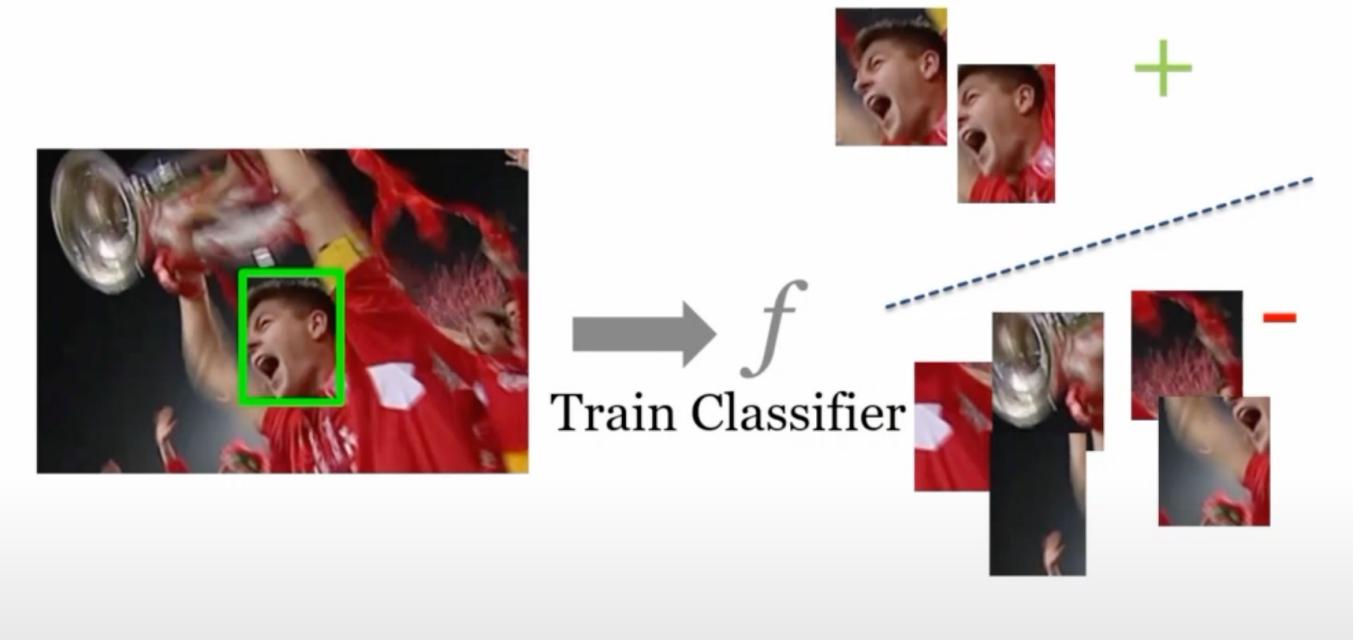
Tracker	A	R	$\hat{\Phi}$	Speed	
Deep DCF → DeepSRDCF*	0.56	1.05	0.32	0.38	(< 1 FPS)
EBT	0.47	1.02	0.31	1.76	
SRDCF*	0.56	1.24	0.29	1.99	
LDP*	0.51	1.84	0.28	4.36	
sPST*	0.55	1.48	0.28	1.01	
SC-EBT	0.55	1.86	0.25	0.80	
NSAMF*	0.53	1.29	0.25	5.47	
Struck*	0.47	1.61	0.25	2.44	

This works well but:

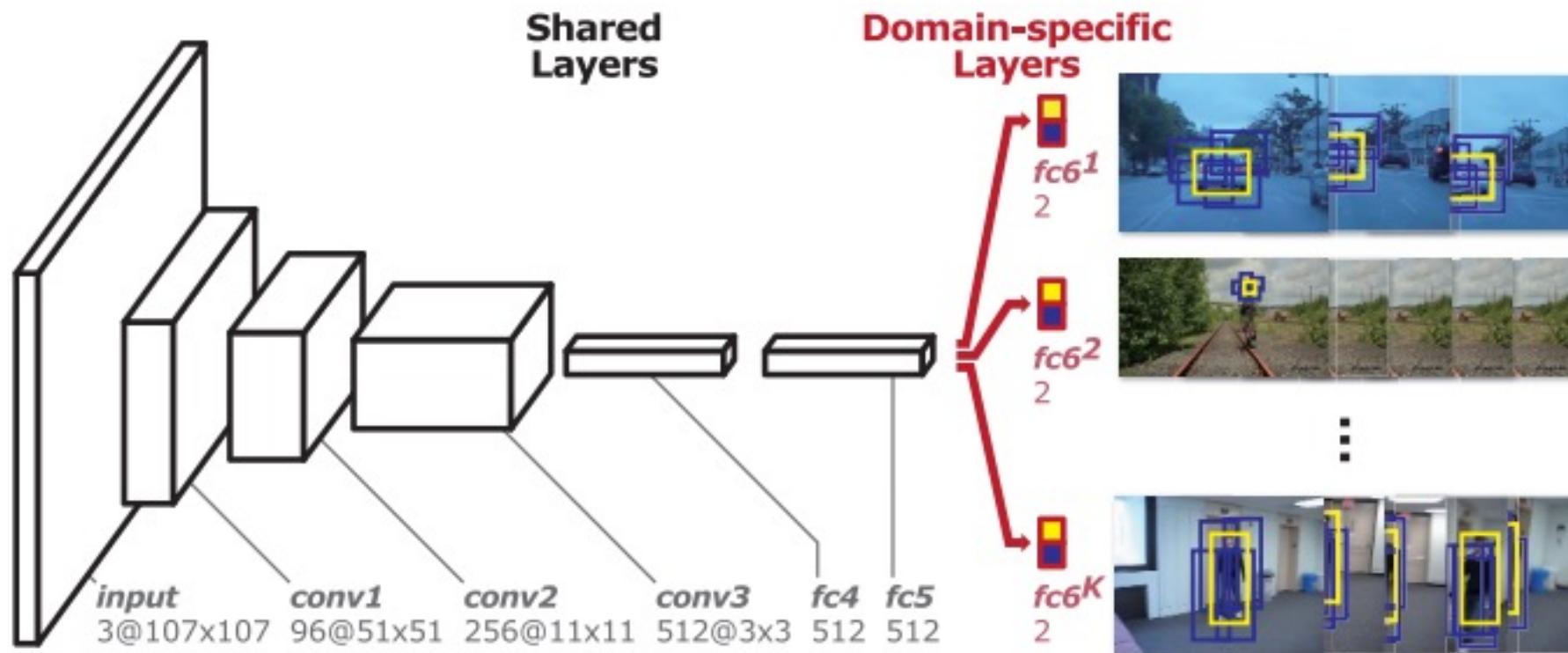
- It is slow (< 1 fps) since the target vs background classifier must be trained online
- Performance is limited since the deep features are trained on classification but then used for tracking

Questions:

- **How can we speed this up?**



MDNet: Train early layers offline, later layers online



Nam, Hyeonseob, and Bohyung Han. "Learning multi-domain convolutional neural networks for visual tracking." CVPR. 2016.

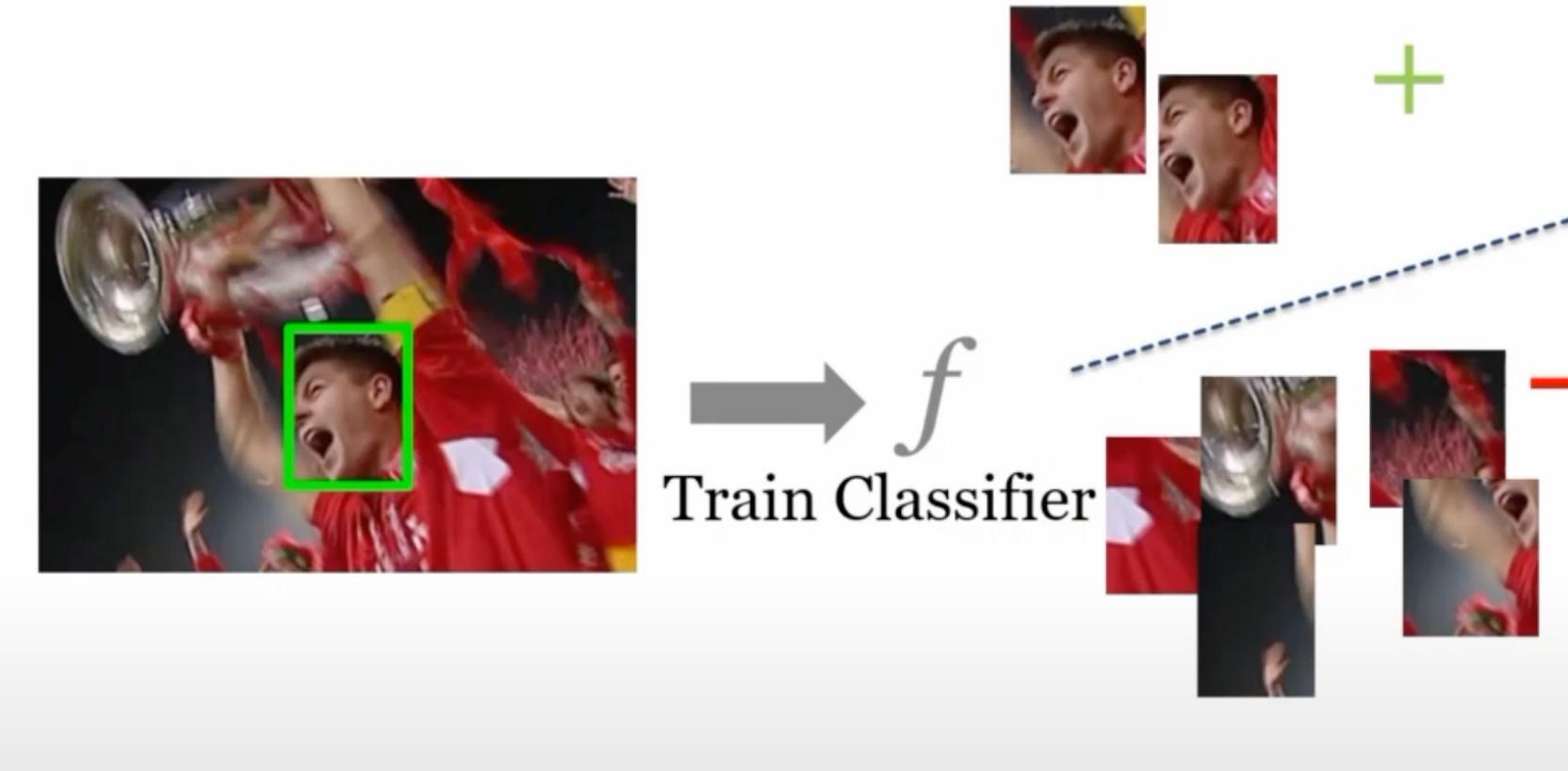
MDNet is the winner of VOT 2015

Tracker	A	R	$\hat{\Phi}$	Speed	
MDNet*	0.60	0.69	0.38	0.87	(1 FPS)
Deep DCF → DeepSRDCF*	0.56	1.05	0.32	0.38	(< 1 FPS)
EBT	0.47	1.02	0.31	1.76	
SRDCF*	0.56	1.24	0.29	1.99	
LDP*	0.51	1.84	0.28	4.36	
sPST*	0.55	1.48	0.28	1.01	
SC-EBT	0.55	1.86	0.25	0.80	
NSAMF*	0.53	1.29	0.25	5.47	
Struck*	0.47	1.61	0.25	2.44	

MDNet has better performance but is still too slow!

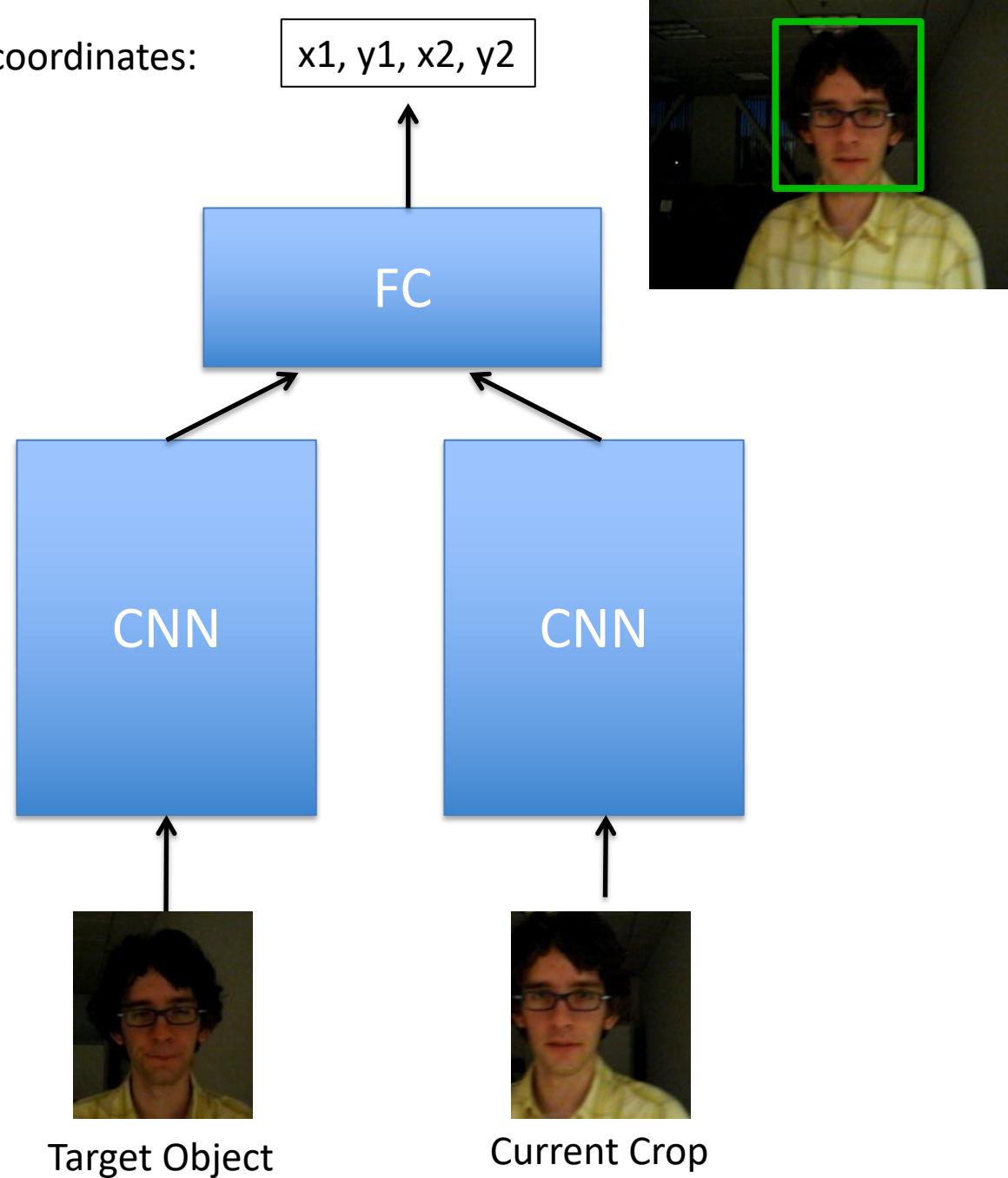
Early approach: Online Train a “Classifier” to Distinguish Target from Background

Can we train a tracker entirely offline?

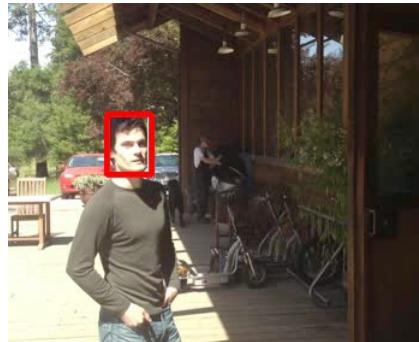
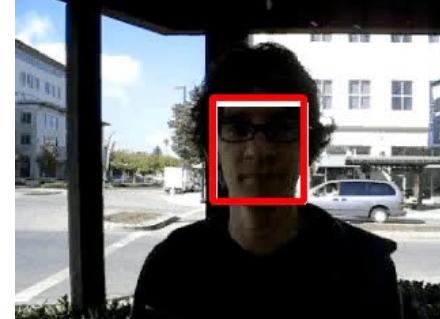


Regress to BB coordinates:

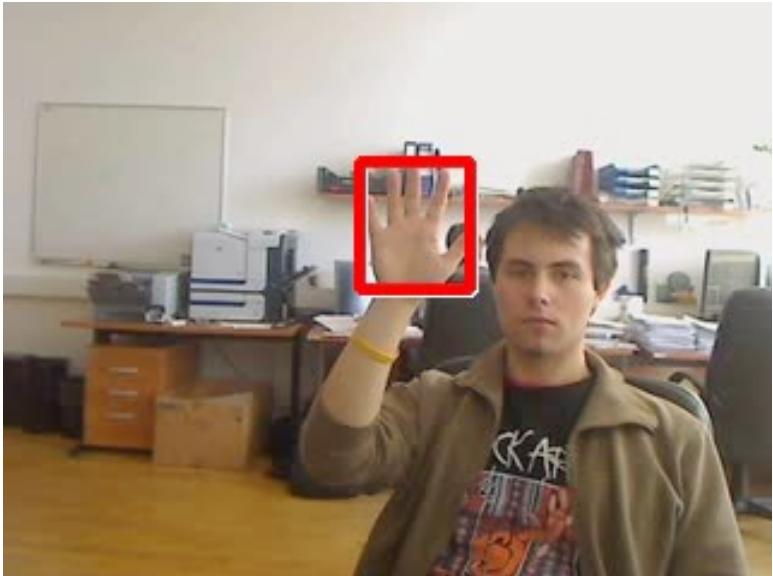
x_1, y_1, x_2, y_2



Tracking at > 100 FPS



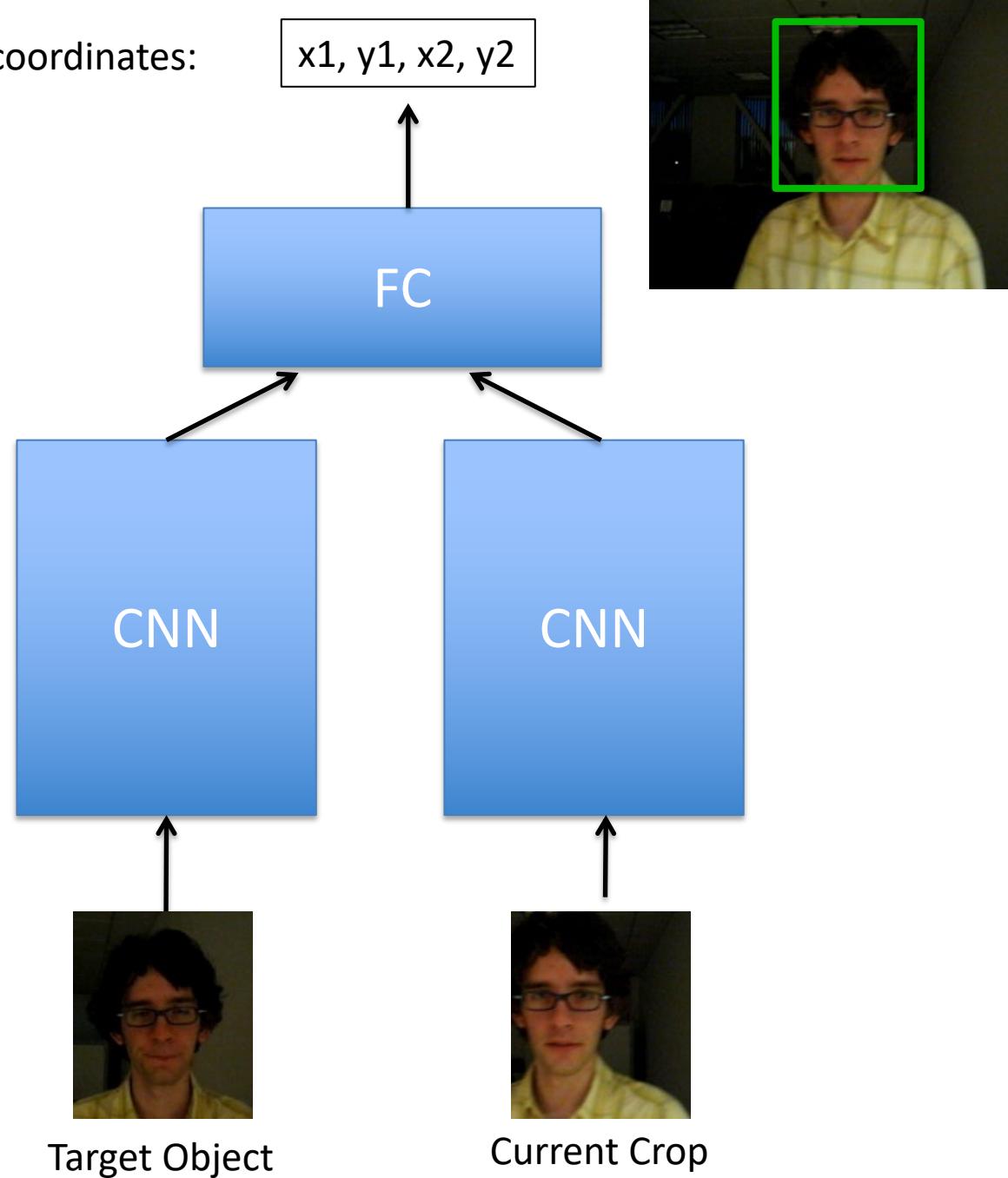
Why does this happen:



Issue: Network does not always pay “attention” to the target,
just learns objectness

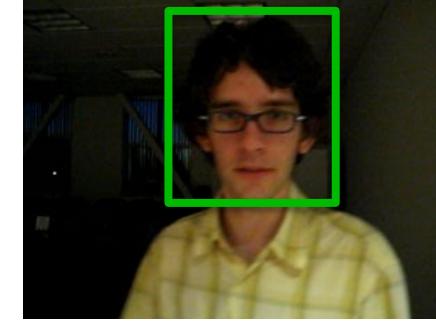
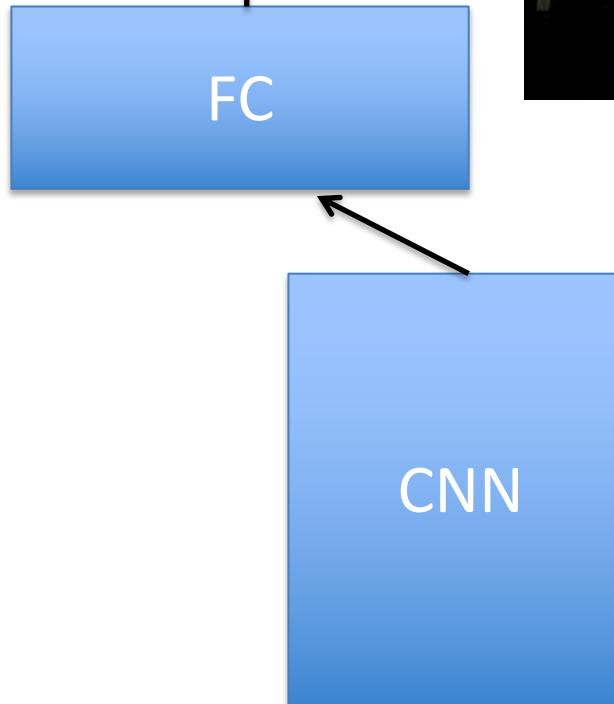
Regress to BB coordinates:

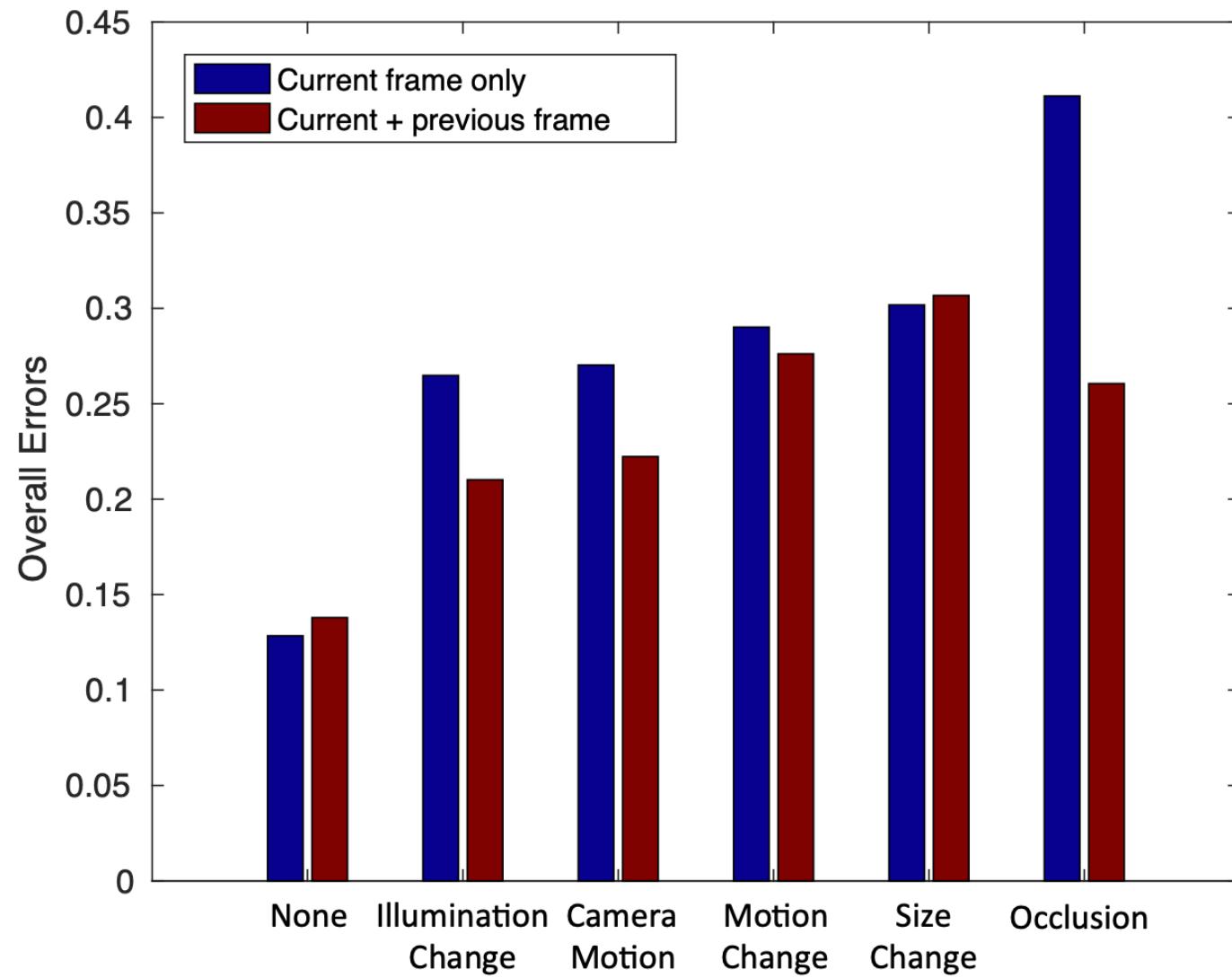
x_1, y_1, x_2, y_2



Regress to BB coordinates:

x_1, y_1, x_2, y_2





Regress to BB coordinates:

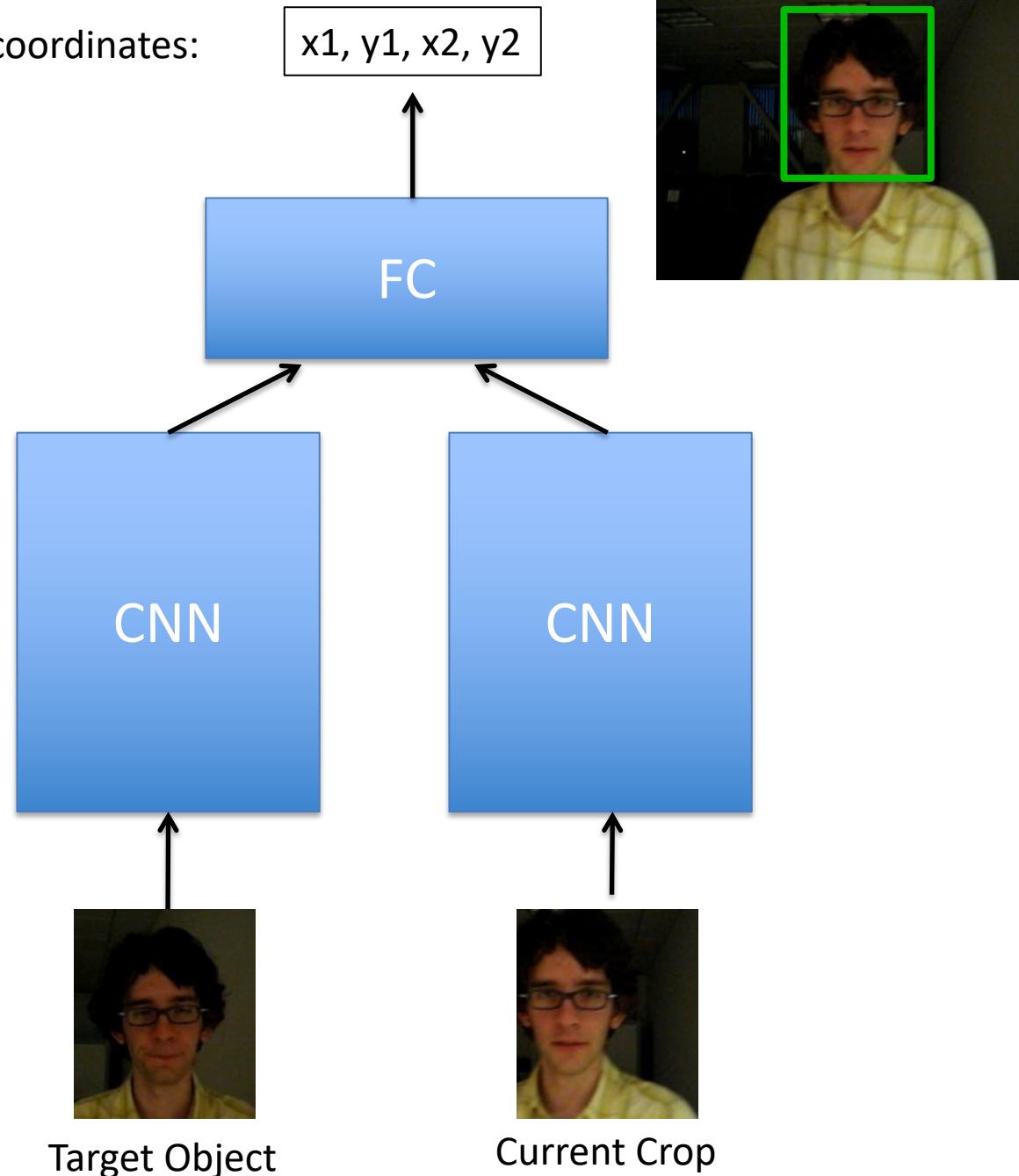
x_1, y_1, x_2, y_2

How can we design a network that pays more “attention” to the target object?

Real problem:

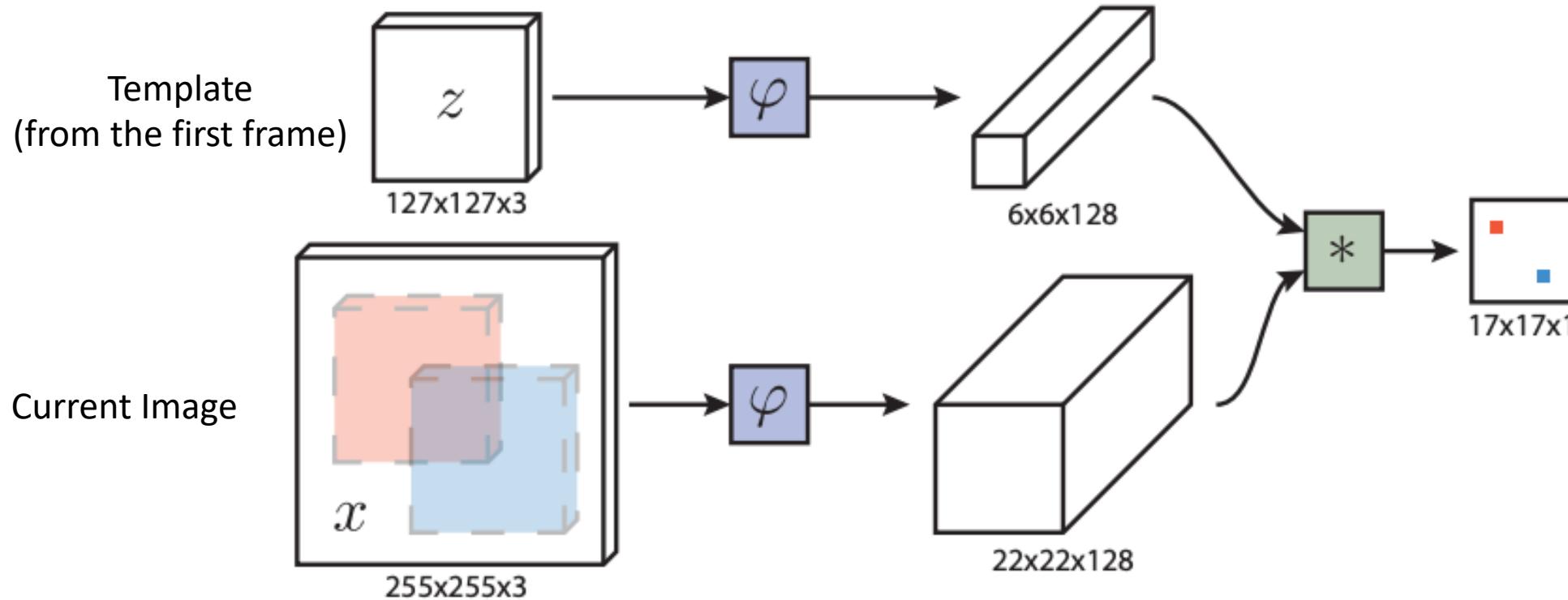
We concatenate the target features with the current image features

Ideas?



Fully-Convolutional Siamese Networks for Tracking

Don't just condition on the target: use a dot product



The network cannot just ignore the template features!

Bertinetto, Luca, et al. "Fully-convolutional siamese networks for object tracking." *European conference on computer vision (ECCV)*, 2016.

Training Data (ImageNet Video)



Results

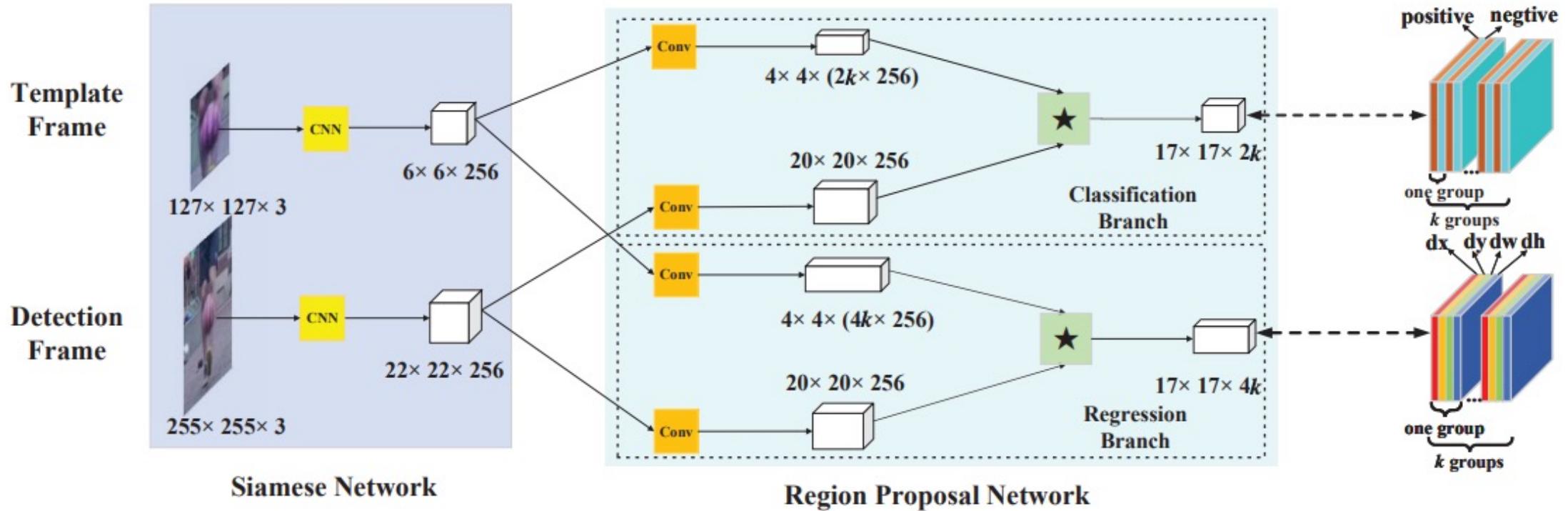


<https://www.youtube.com/watch?v=V7Fab9y3ASk>

Much faster than previous approaches due to no online training but performs almost as well!

Tracker	accuracy	# failures	overlap	speed (fps)
MDNet [9]	0.5620	46	0.3575	1
EBT [41]	0.4481	49	0.3042	5
DeepSRDCF [6]	0.5350	60	0.3033	< 1 *
SiamFC-3s (ours)	0.5335	84	0.2889	86
SiamFC (ours)	0.5240	87	0.2743	58
SRDCF [42]	0.5260	71	0.2743	5
sPST [43]	0.5230	85	0.2668	2
LDP [12]	0.4688	78	0.2625	4 *
SC-EBT [44]	0.5171	103	0.2412	—
NSAMF [45]	0.5027	87	0.2376	5 *
StruckMK [3]	0.4442	90	0.2341	2
S3Tracker [46]	0.5031	100	0.2292	14 *
RAJSSC [12]	0.5301	105	0.2262	2 *
SumShift [46]	0.4888	97	0.2233	17 *
DAT [47]	0.4705	113	0.2195	15
SO-DLT [7]	0.5233	108	0.2190	5

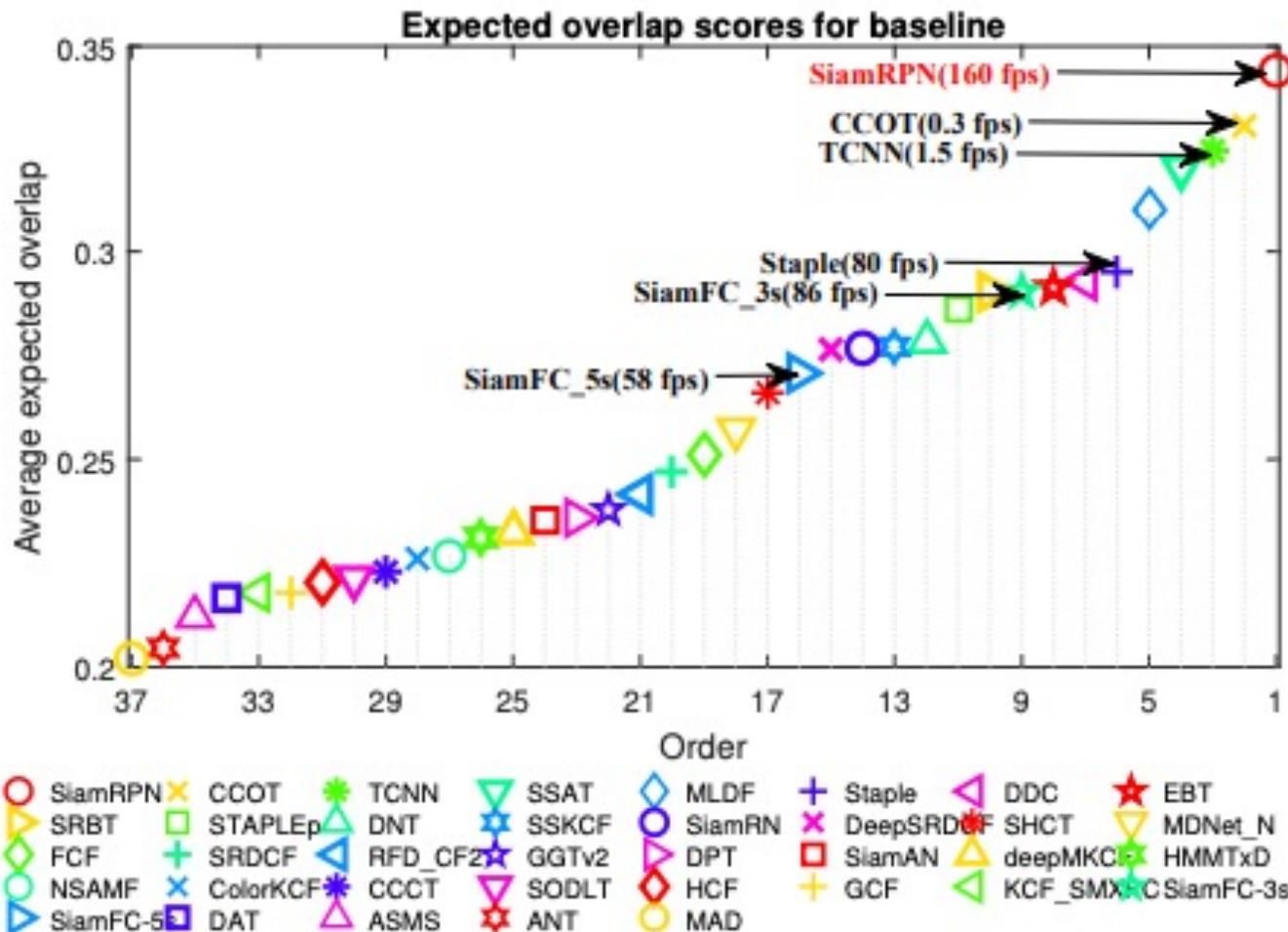
SiamRPN: Tracking as One-shot Detection



Trained with detection-style losses
Classification (coarse localization) and Regression (fine localization)

Li, Bo, et al. "High performance visual tracking with siamese region proposal network." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.

Results (VOT 2016)



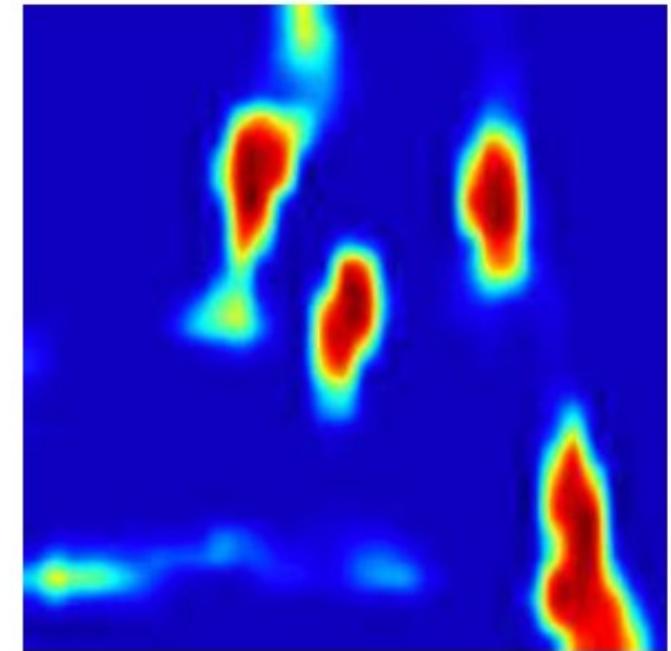
Winner of real-time challenge in VOT 2018 (>20 fps)

Problem: Distractors

Why does this happen?



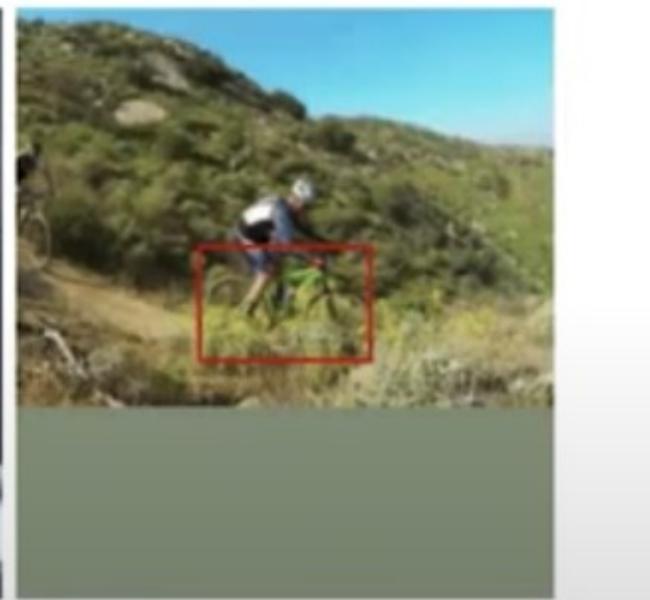
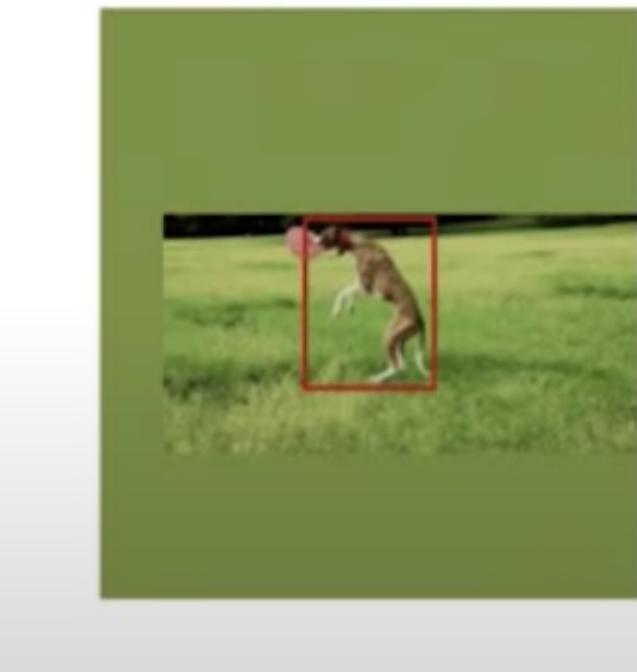
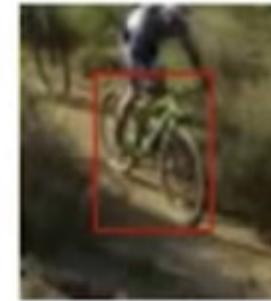
ROI



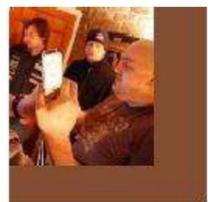
SiamRPN

Problem: Training Data

How can we fix this?



How to fix this: harder negatives?



negative pairs from the
same categories

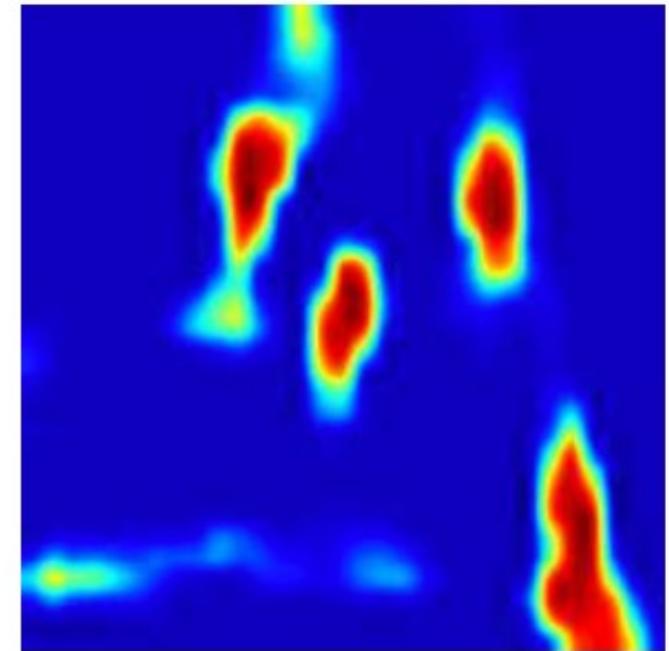
negative pairs from
different categories

Will this fix the issue?

Empirically it helps a little but doesn't fix it

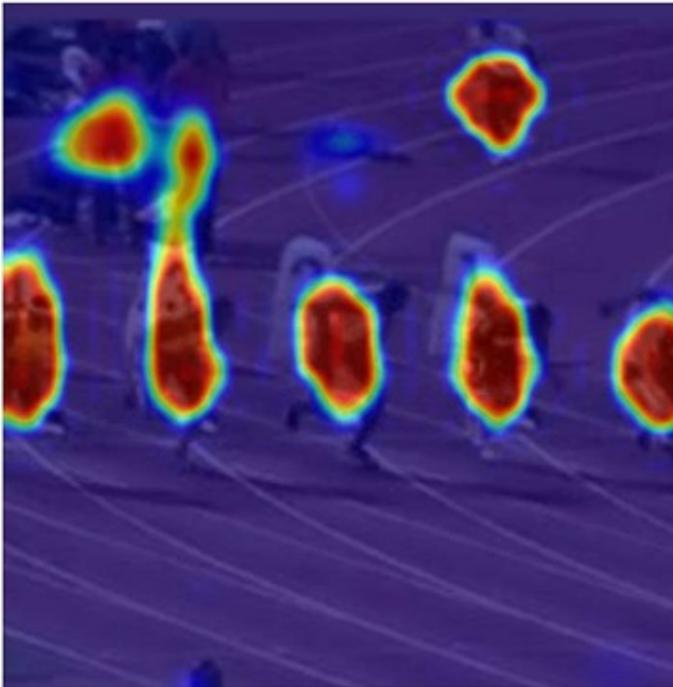


ROI

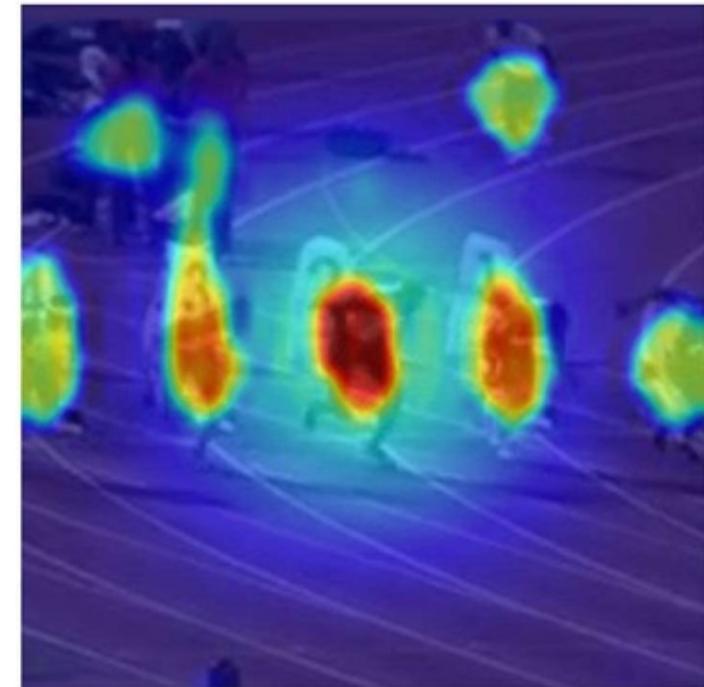


SiamRPN

How to fix this: cosine window?



SiamRPN

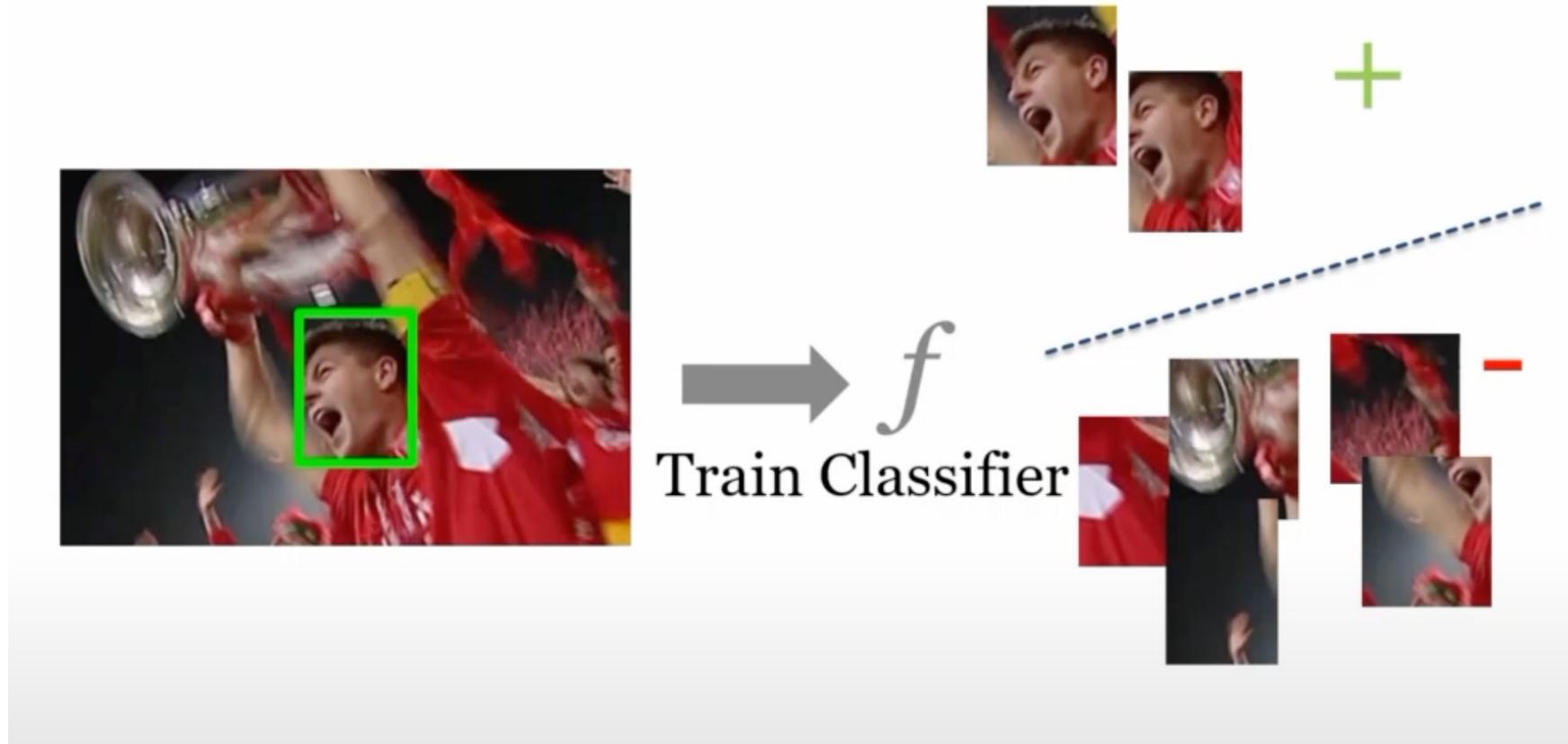


Cosine Window

Problem: Cannot handle large object or camera motion

Solution: Train the network online!

Online Train a “Classifier” to Distinguish Target from Background



Online Train a “Classifier” to Distinguish Target from Background



This helps but it is too slow – what do we do?

<https://www.youtube.com/watch?v=j7A83F6PRAE&t=3630s>

View tracking from the optimization perspective

How to minimize $f(x)$?

Gradient descent: $x_{n+1} = x_n - \alpha \nabla_x f(x)$

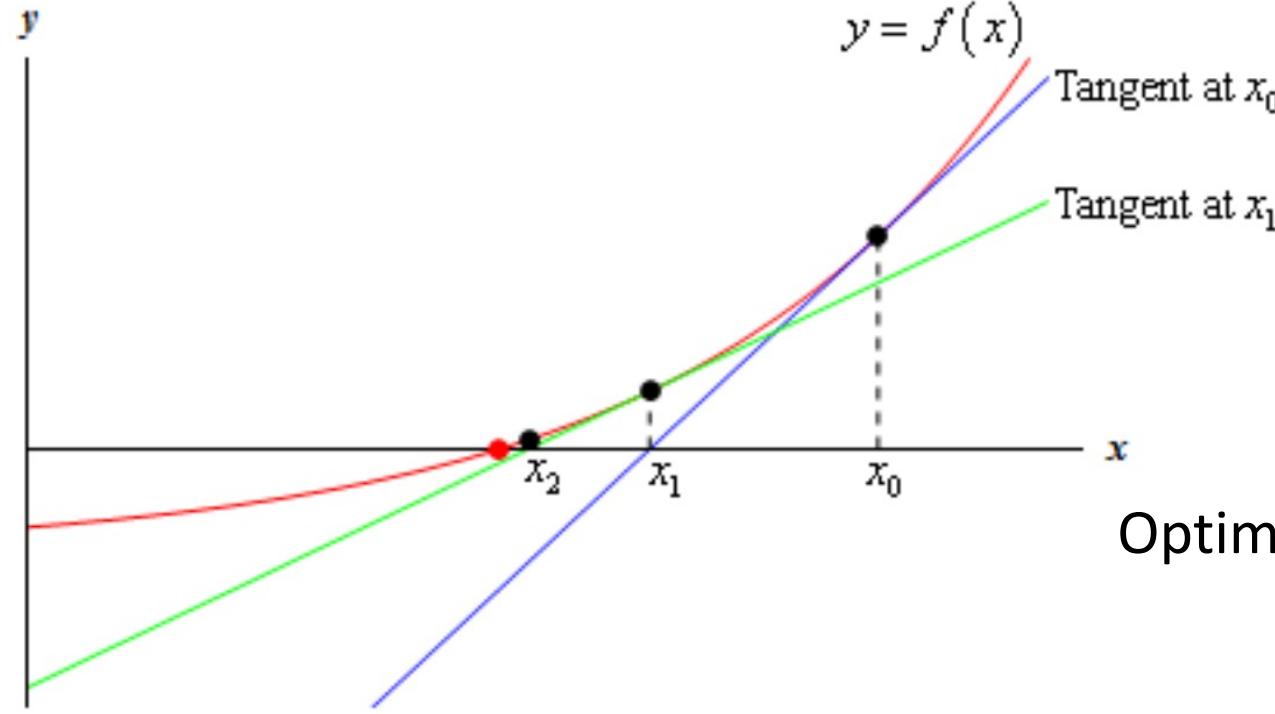
Can we design a better method that converges more quickly?

What if we assume that the object doesn't change appearance much from frame to frame?

Credits: Alfonso Croeze, Lindsey Pittman, Winnie Reynolds:

<https://www.math.lsu.edu/system/files/MunozGroup1%20-%20Presentation.pdf>

Newton's Method: Iterative method to find roots



1. Linearize the function $f(x)$
2. Solve for the roots of the linear function

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

3. Repeat

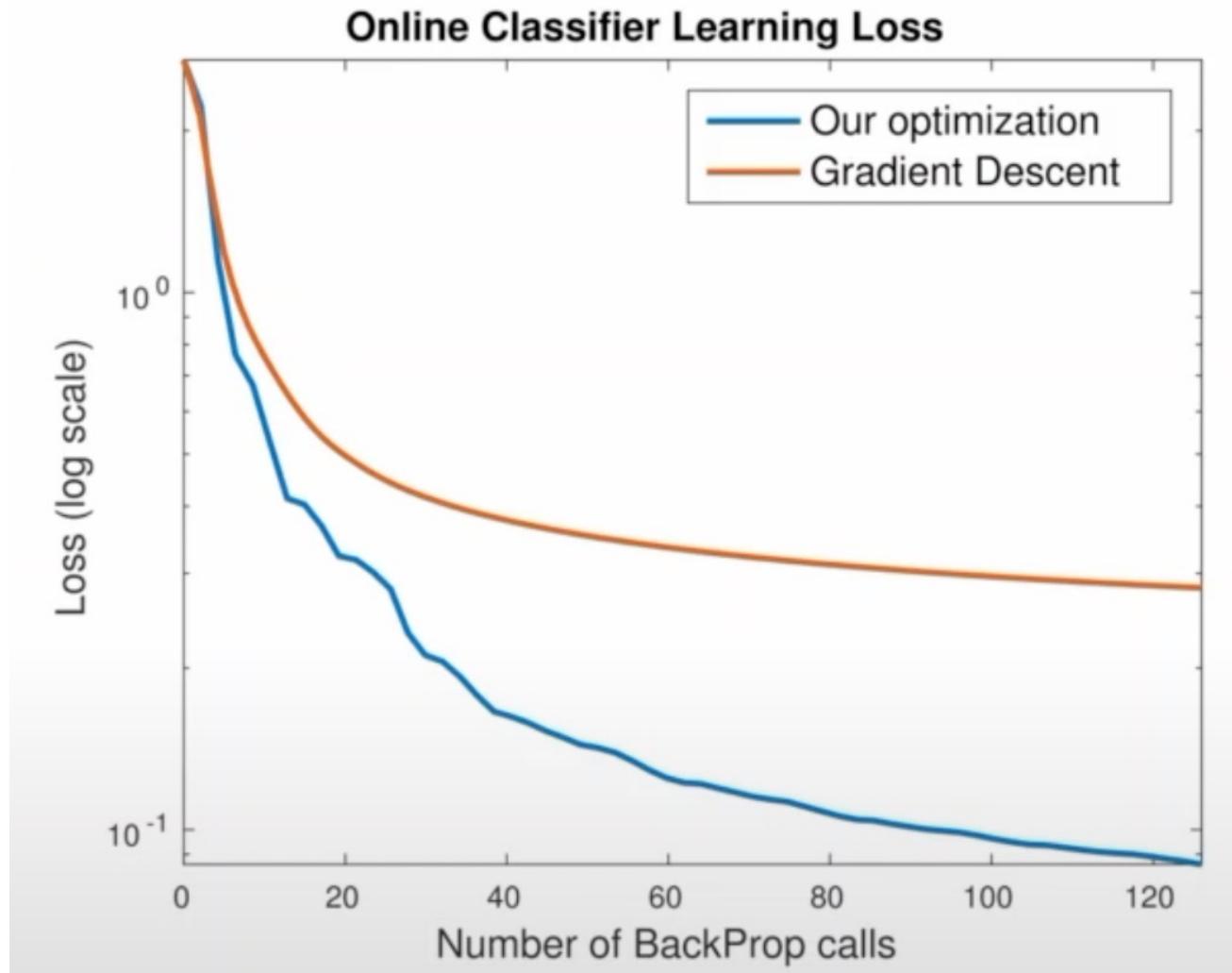
Optimization: we want to find the roots of $f'(x)$

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

Credits: Alfonso Croeze, Lindsey Pittman, Winnie Reynolds:

<https://www.math.lsu.edu/system/files/MunozGroup1%20-%20Presentation.pdf>

How to train this quickly? Gauss-Newton!



Trackers trained online are robust (stick roughly to the target) but are not sufficiently accurate



[1] G. Bhat, J. Johnander, M. Danelljan, F. S. Khan, and M. Felsberg. Unveiling the power of deep tracking. In ECCV, 2018

- How do we (as humans) know how to adjust the bounding box for these cases?
- We cannot determine this from a single image
- We need to learn this offline from data!
- But if we learn offline then we cannot be robust to distractors! What to do?
- Idea: Learn to roughly localize from online training; learn bounding box refinement from offline data

First frame

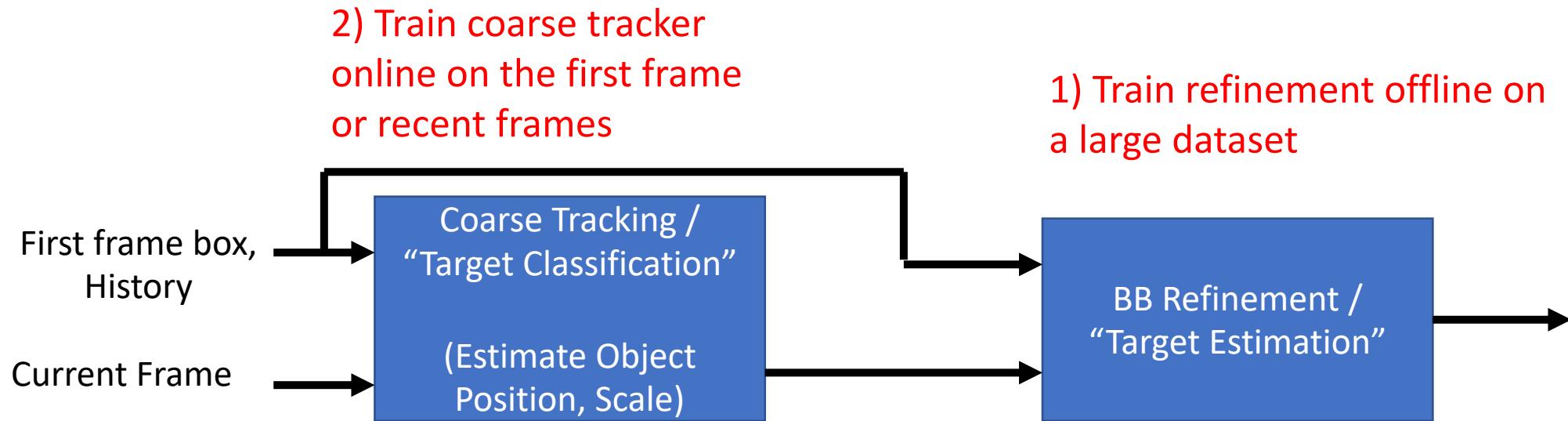


Frame T

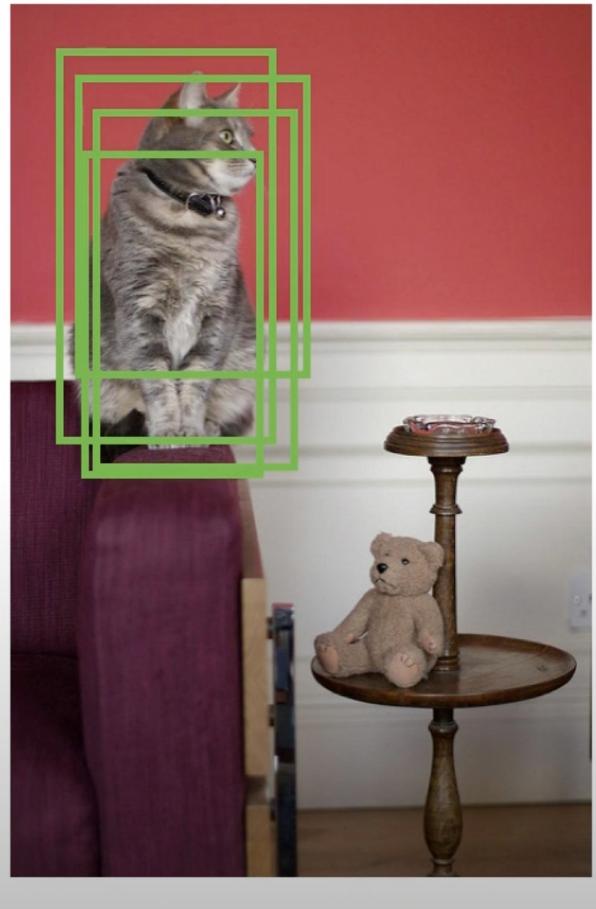


[1] G. Bhat, J. Johnander, M. Danelljan, F. S. Khan, and M. Felsberg. Unveiling the power of deep tracking. In ECCV, 2018

General Tracking Pipeline



How can we localize objects in detection?



Region Proposal Network

Generate box proposals.

Classification.

Classification score

Assign class labels: cat/dog...

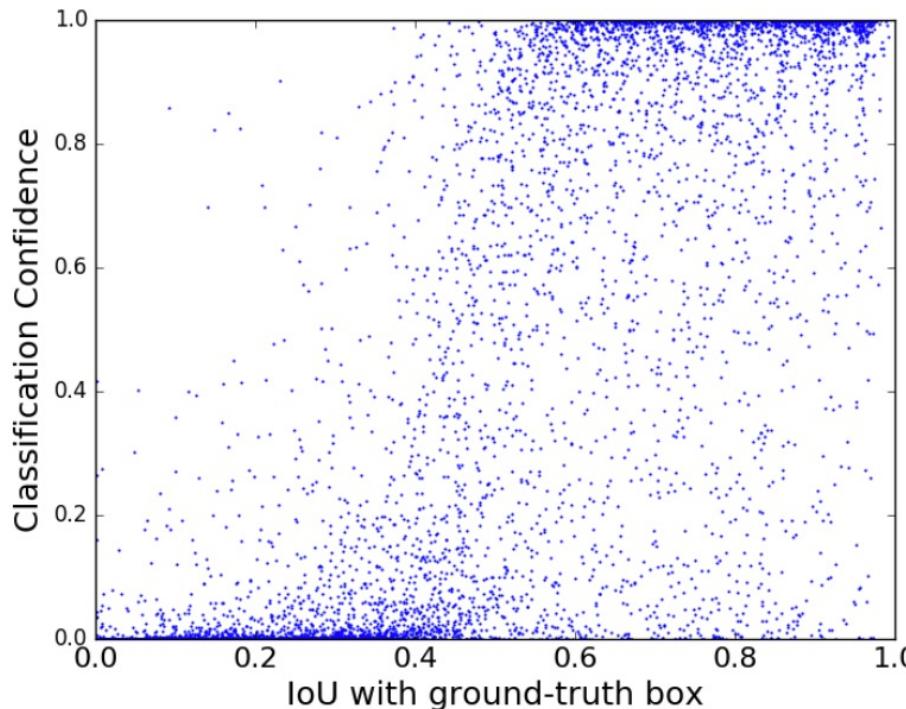
Bounding Box Regression.

Make the box accurate.

Non-Maximum Suppression

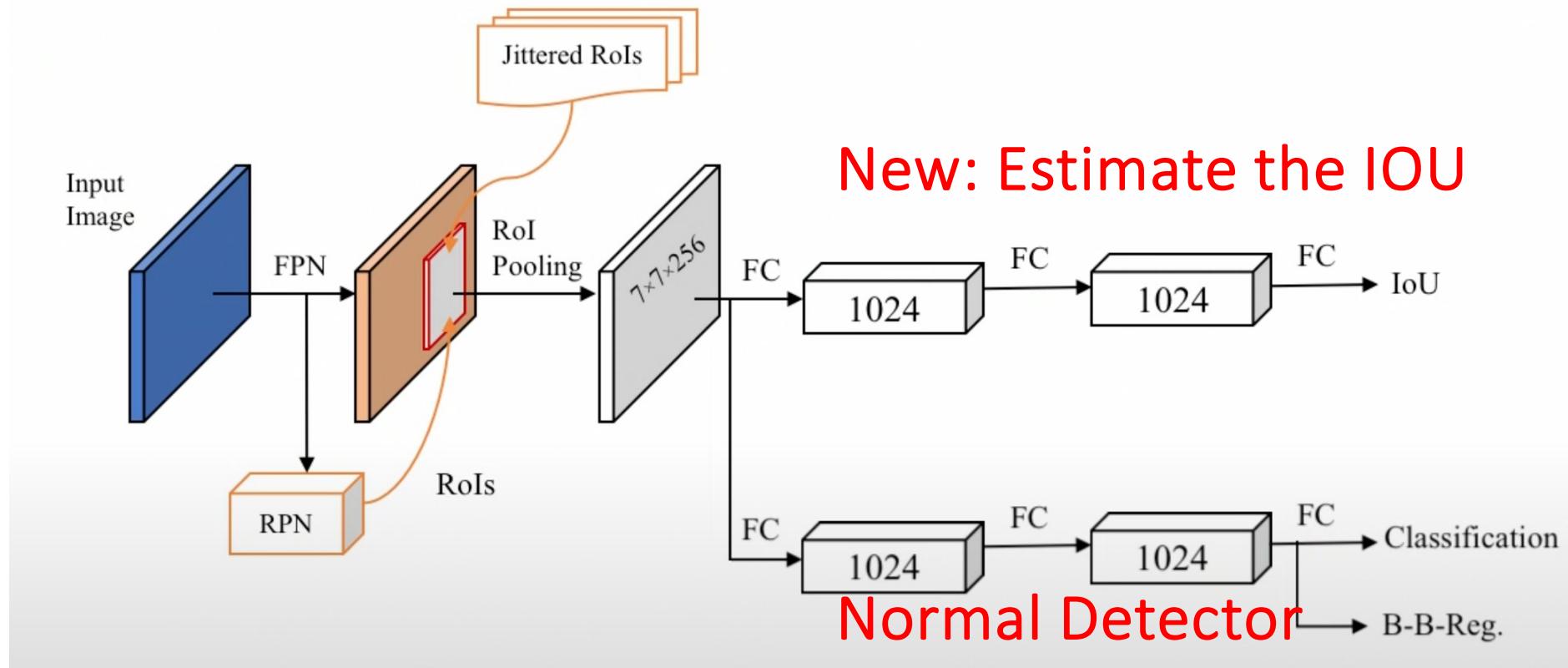
Based on
classification score

Remove duplicated detections.

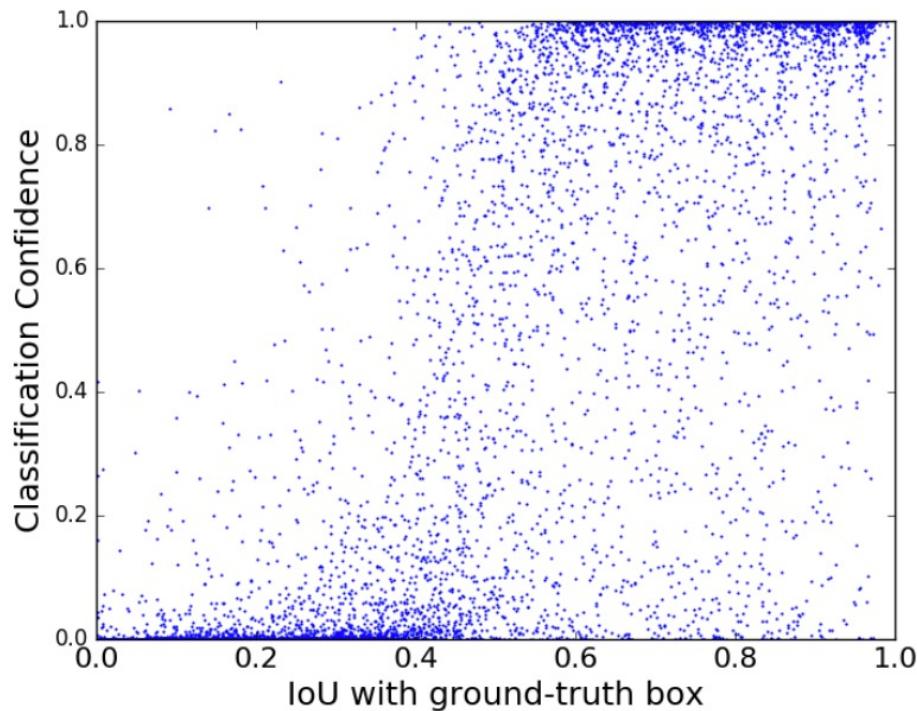


(a) IoU vs. Classification Confidence

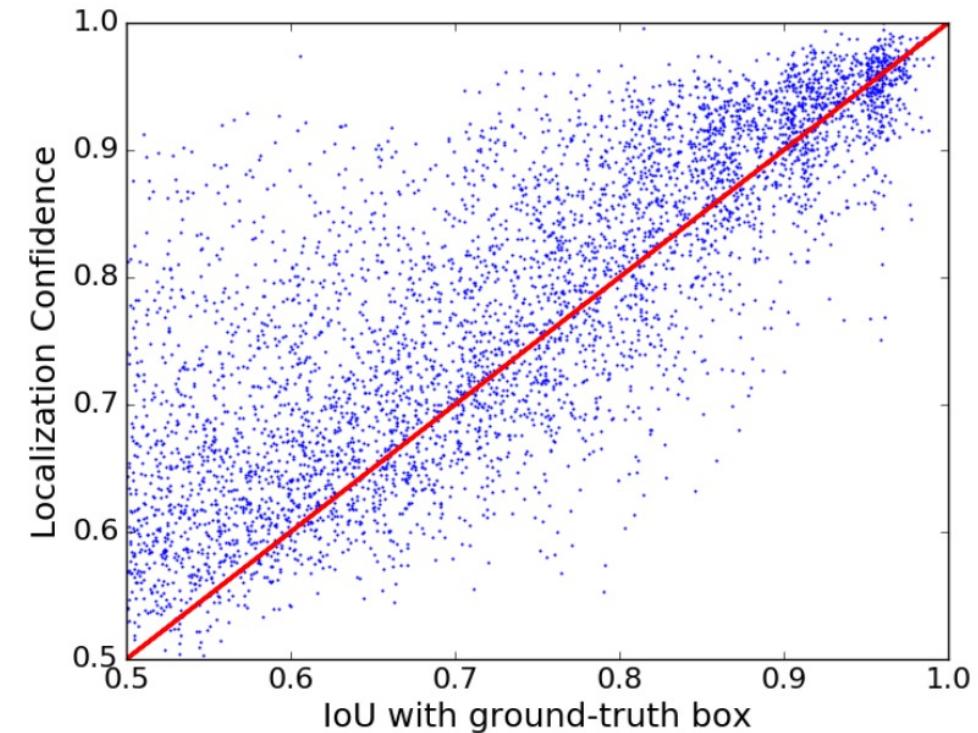
IOUNet: B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang. Acquisition of localization confidence for accurate object detection. In ECCV, 2018



IOUNet: B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang. Acquisition of localization confidence for accurate object detection. In ECCV, 2018



(a) IoU vs. Classification Confidence

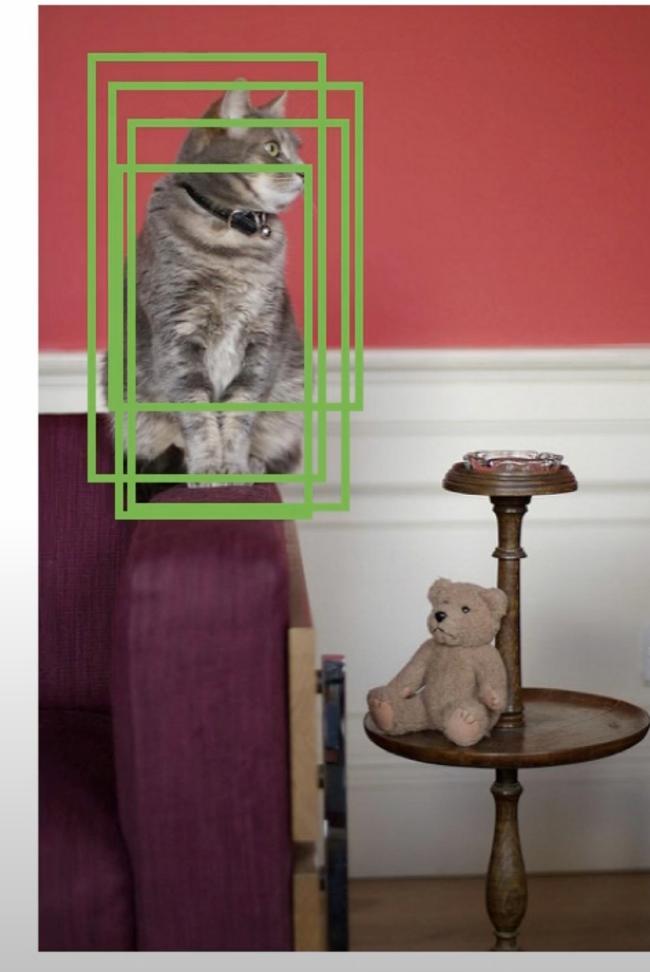


(b) IoU vs. Localization Confidence

IOUNet: B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang. Acquisition of localization confidence for accurate object detection. In ECCV, 2018

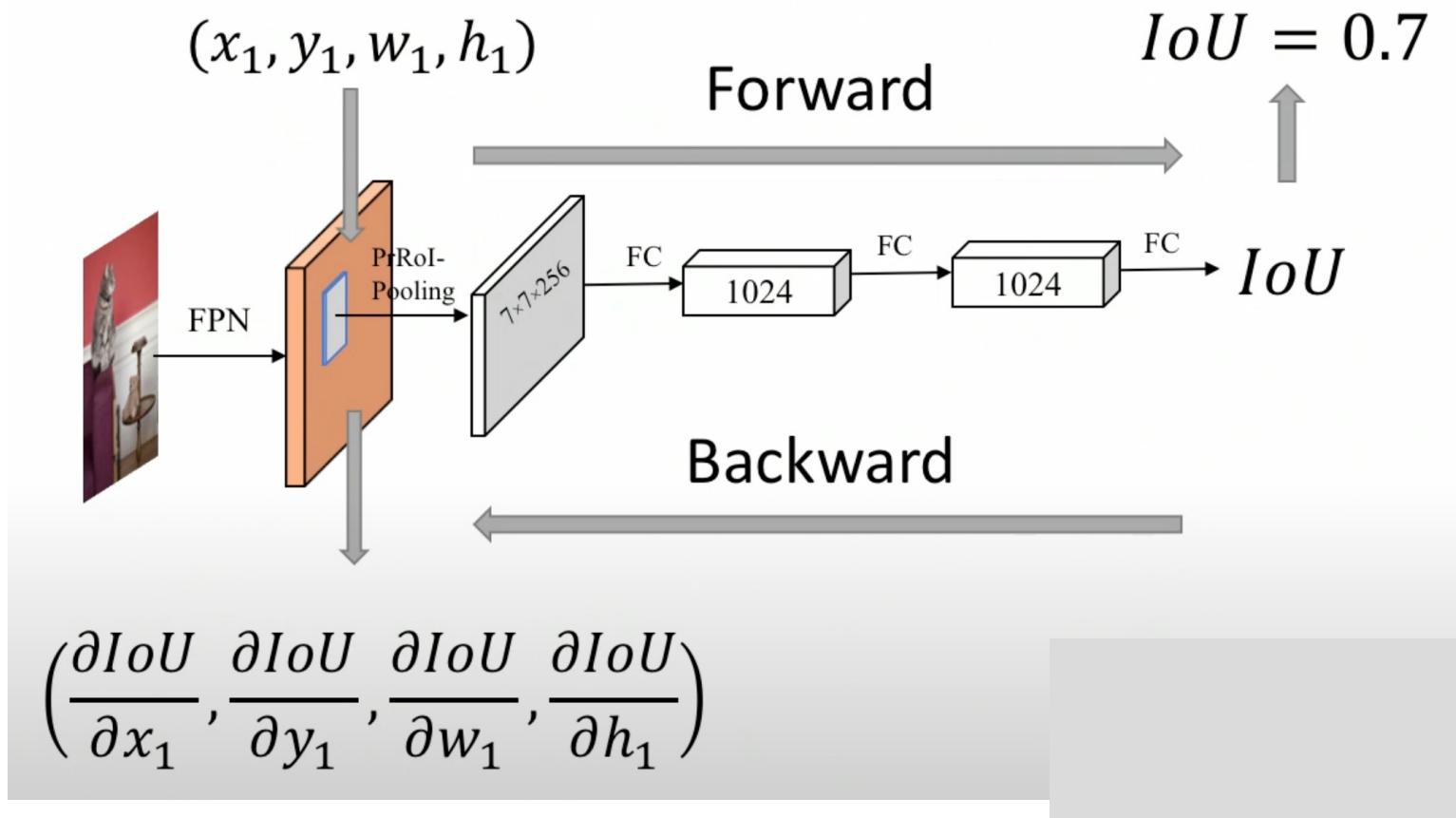
Improve NMS: IOU-Guided NMS

- First keep the box with the highest IOU (objectness)
- Next, find all overlapping boxes and choose the label with the most confident (highest) classification score



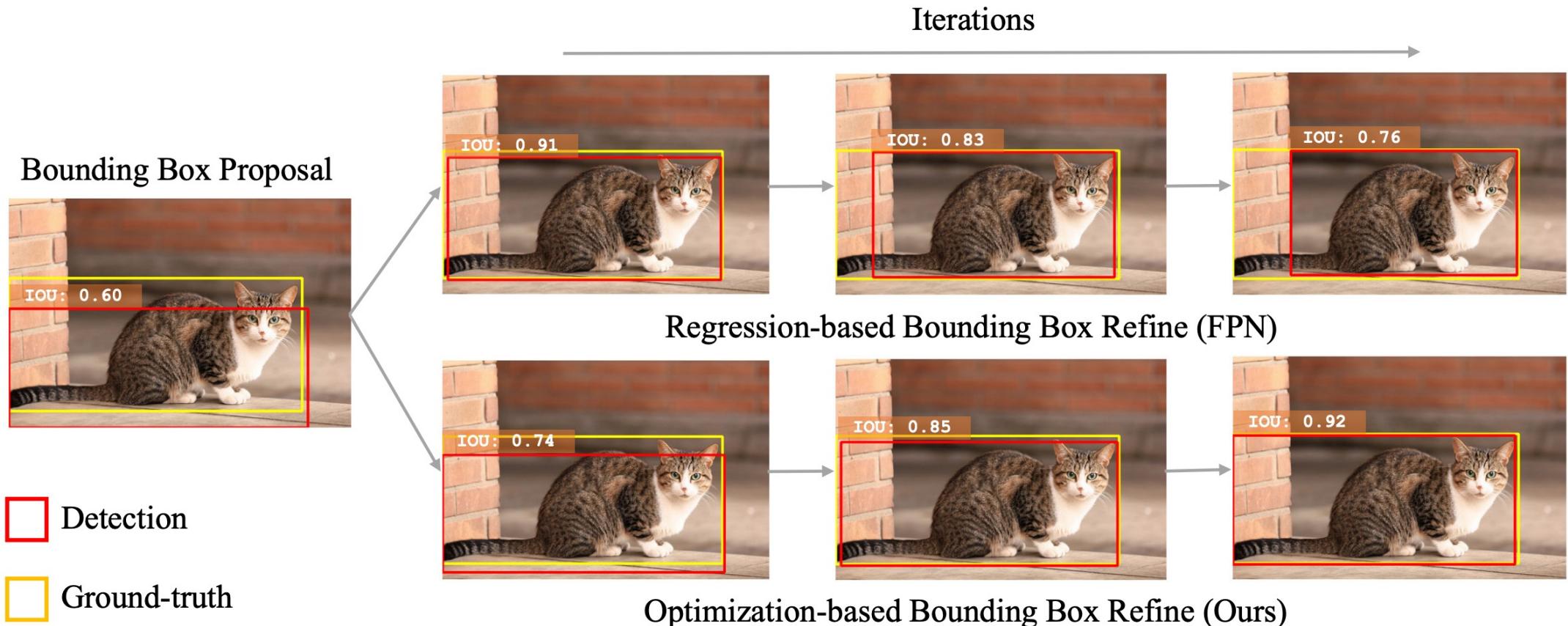
	Cls. Conf.	Loc. Conf.
Box 1	0.85	0.95
Box 2	0.9	0.82
Box 3	0.96	0.74
Box 4	0.85	0.70
.....		

Alternative: Refinement by IoU Optimization

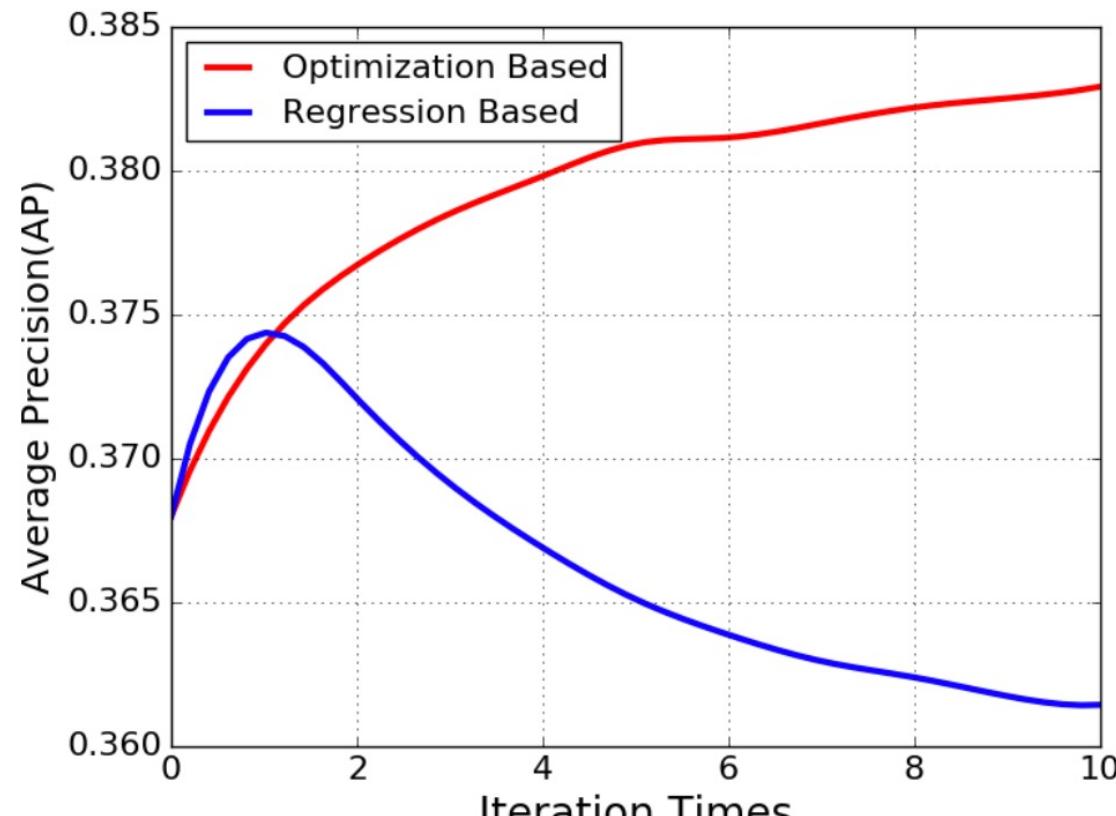


IOUNet: B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang. Acquisition of localization confidence for accurate object detection. In ECCV, 2018

Regression vs Optimization



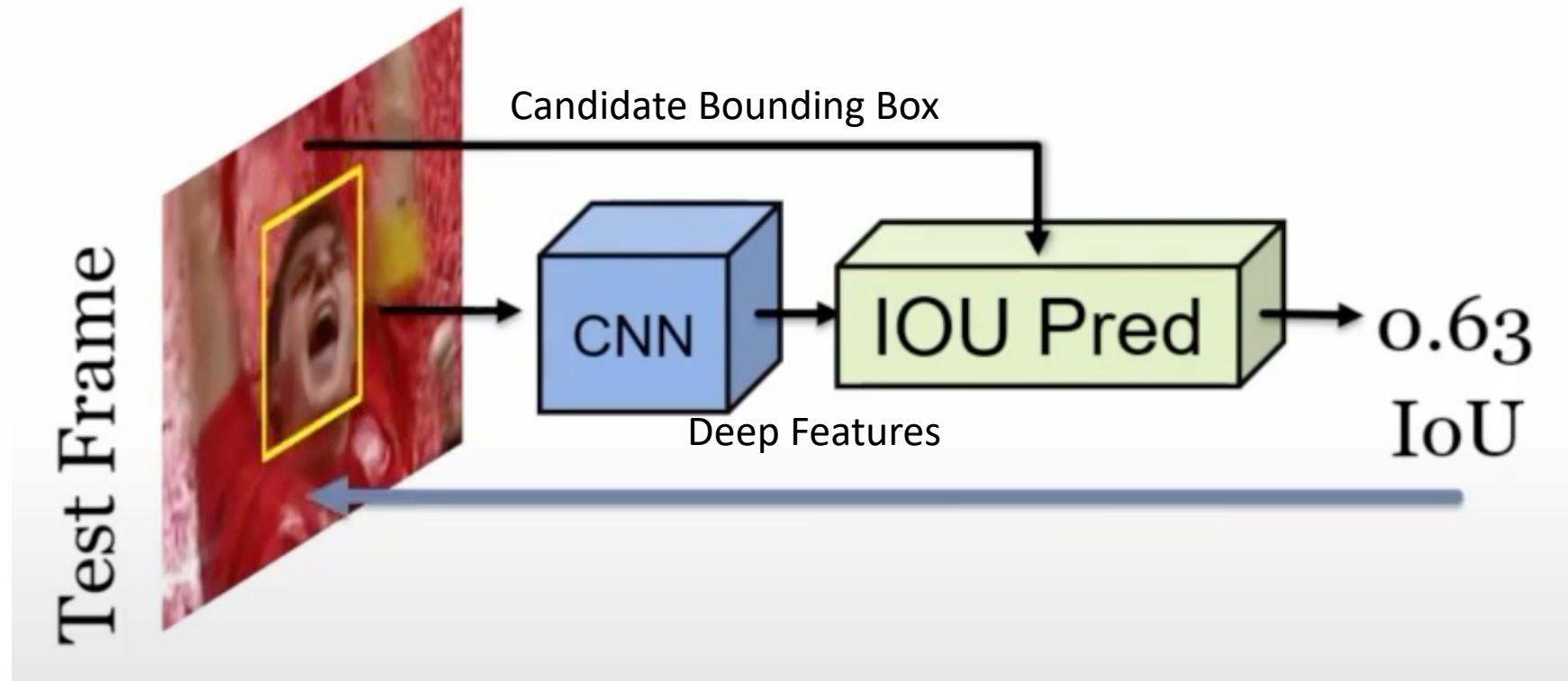
Regression vs Optimization



(a) FPN

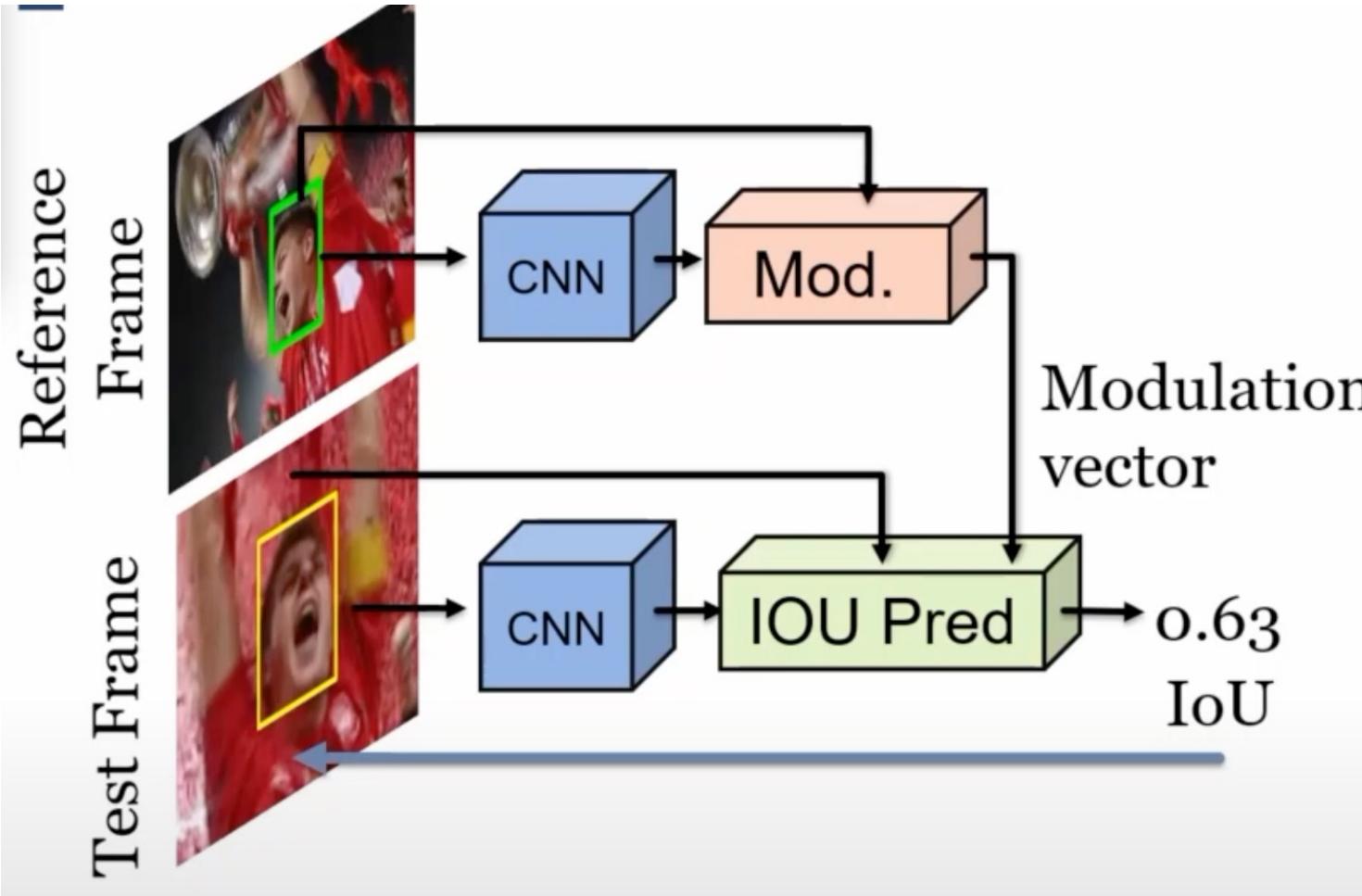
IOUNet: B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang. Acquisition of localization confidence for accurate object detection. In ECCV, 2018

IoU Net – designed for detection, shifts the box to the nearest “object”



IOUNet: B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang. Acquisition of localization confidence for accurate object detection. In ECCV, 2018

Tracking BB Refinement (ATOM) – shifts the box to the location that matches the target



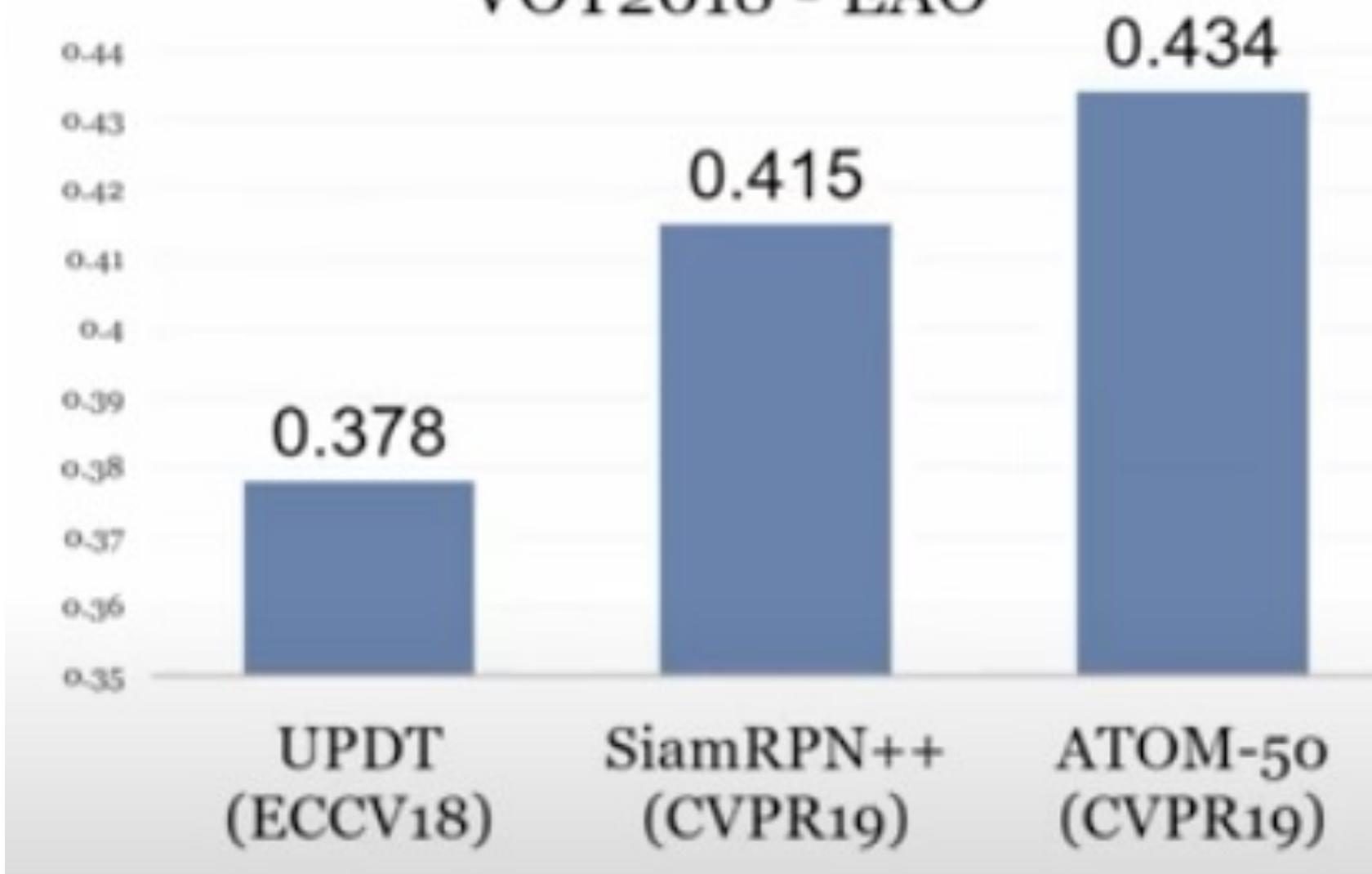
Danelljan, Martin, et al. "Atom: Accurate tracking by overlap maximization." CVPR. 2019.

ATOM: Results



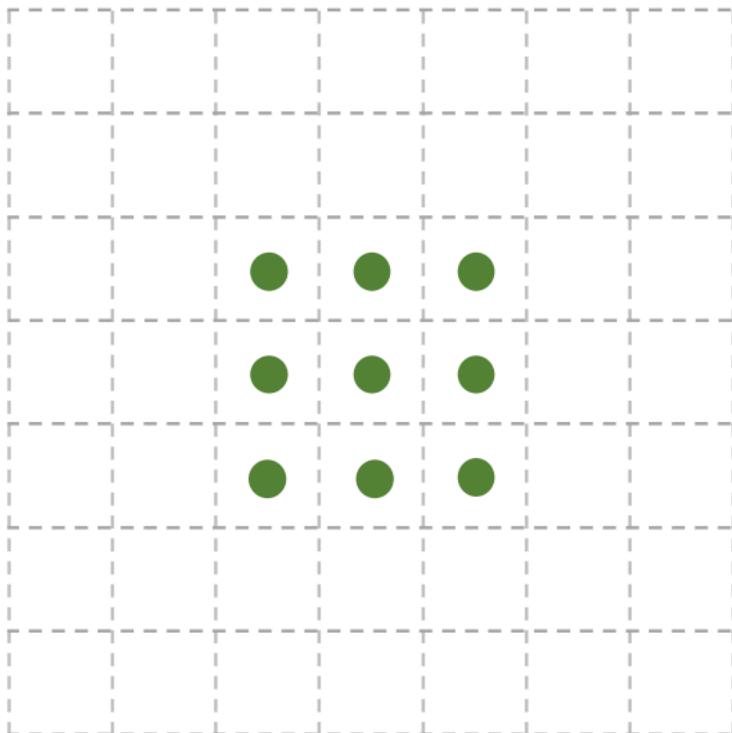
Danelljan, Martin, et al. "Atom: Accurate tracking by overlap maximization." CVPR. 2019.

VOT2018 - EAO

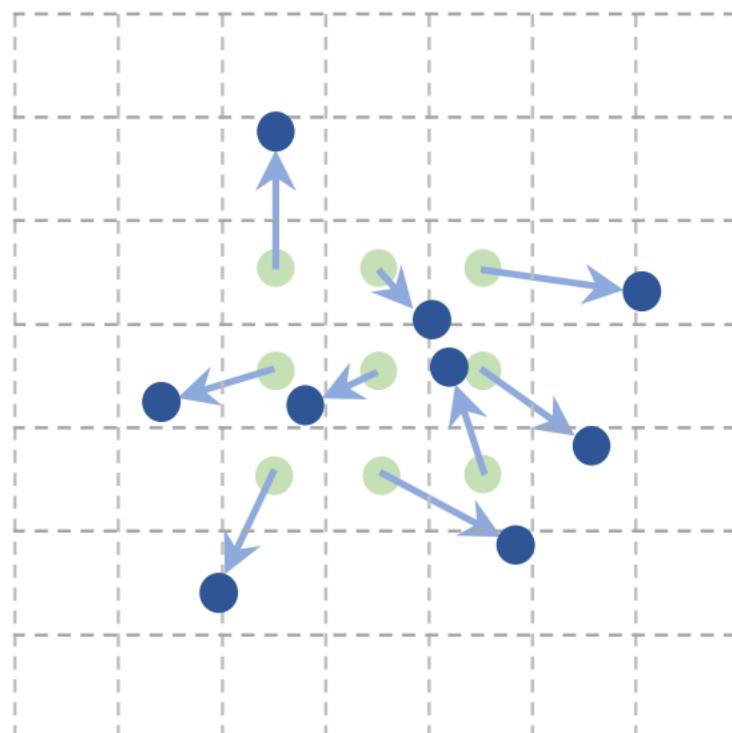


<https://www.youtube.com/watch?v=j7A83F6PRAE&t=3629s>

Further improvement: Deformable Convolution

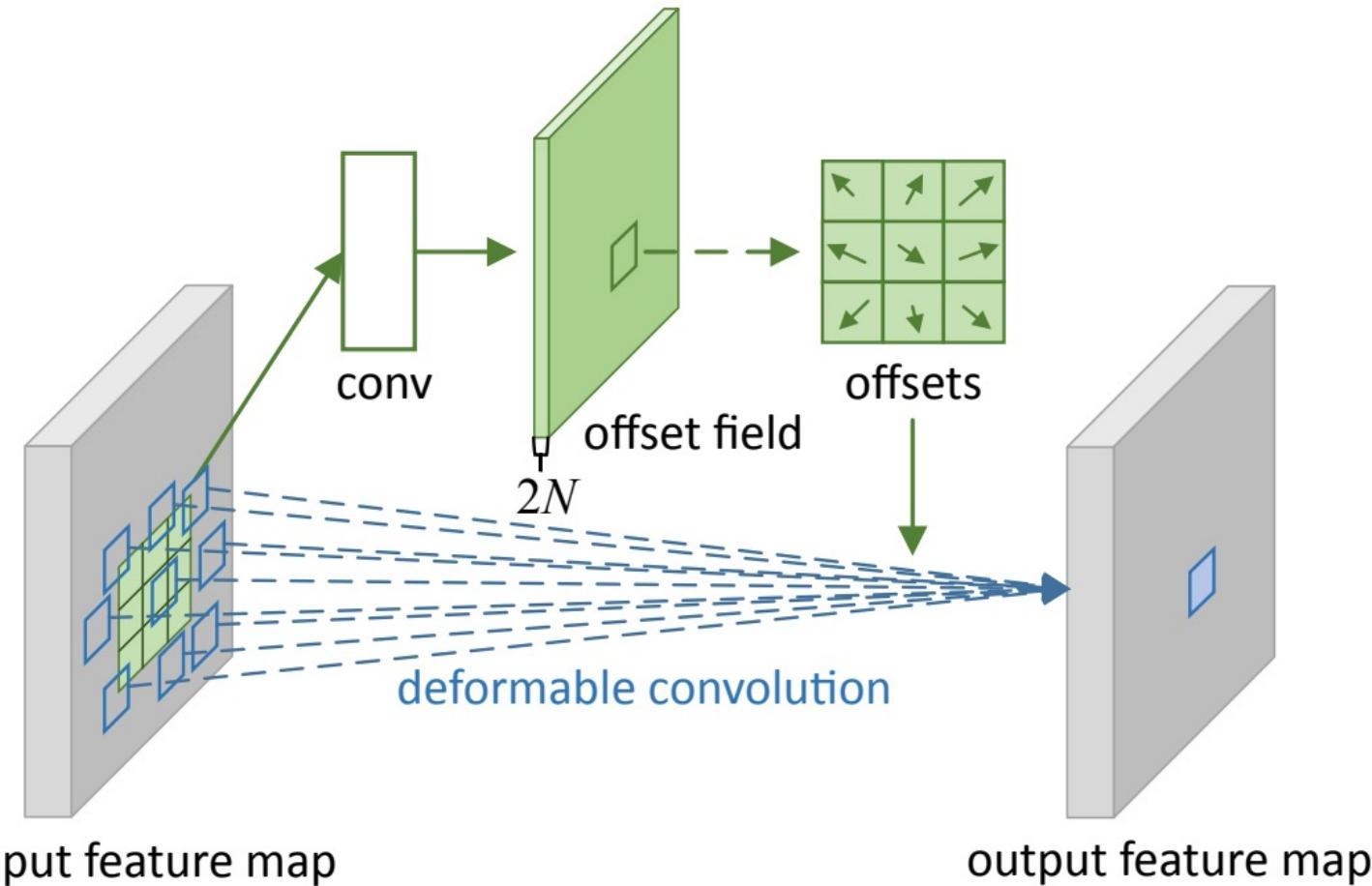


Regular (3x3) Convolution



Deformable Convolution

Deformable Convolution



Offsets and convolution filters are trained end-to-end!

Further improvement: Deformable Convolution

Convolution
Filter (bilinear interpolation)

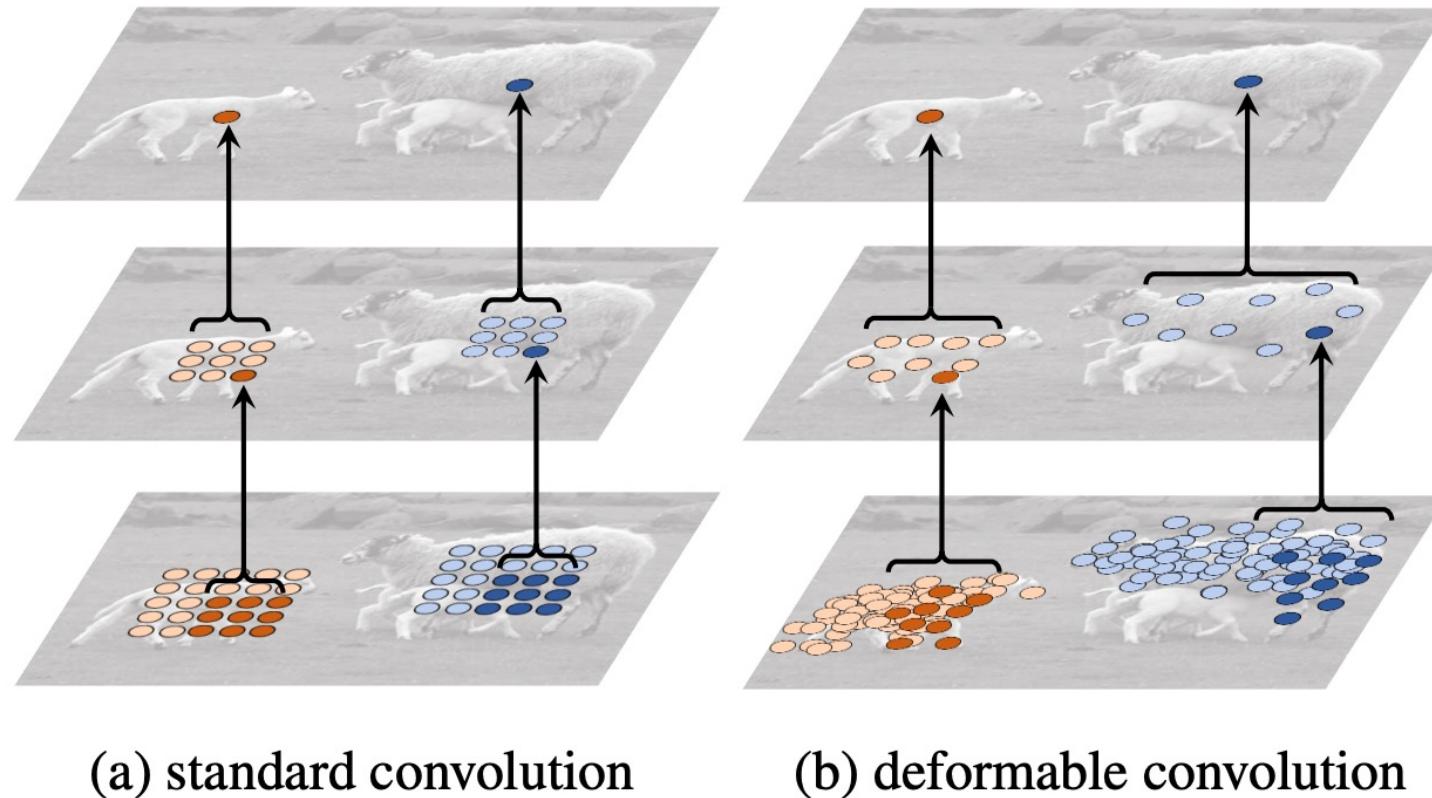
$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n + \Delta p_n)$$

Regular Grid Points Center Point Regular Grid Points Offsets

The diagram illustrates the calculation of a deformable convolution output $y(p_0)$. It shows a 'Center Point' at the origin, surrounded by a grid of 'Regular Grid Points'. 'Offsets' are applied to these grid points to calculate shifted points $p_0 + p_n + \Delta p_n$. The output is computed as a weighted sum of the input x at these shifted points, using a 'Filter' $w(p_n)$ and bilinear interpolation.

Offsets and convolution filters are trained end-to-end!

Adaptive Receptive Field



(a) standard convolution

(b) deformable convolution

Adaptive Receptive Field



Adaptive receptive field for three layers of 3x3 deformable filters for activation units at the g
Total points = $9^3 = 729$ points

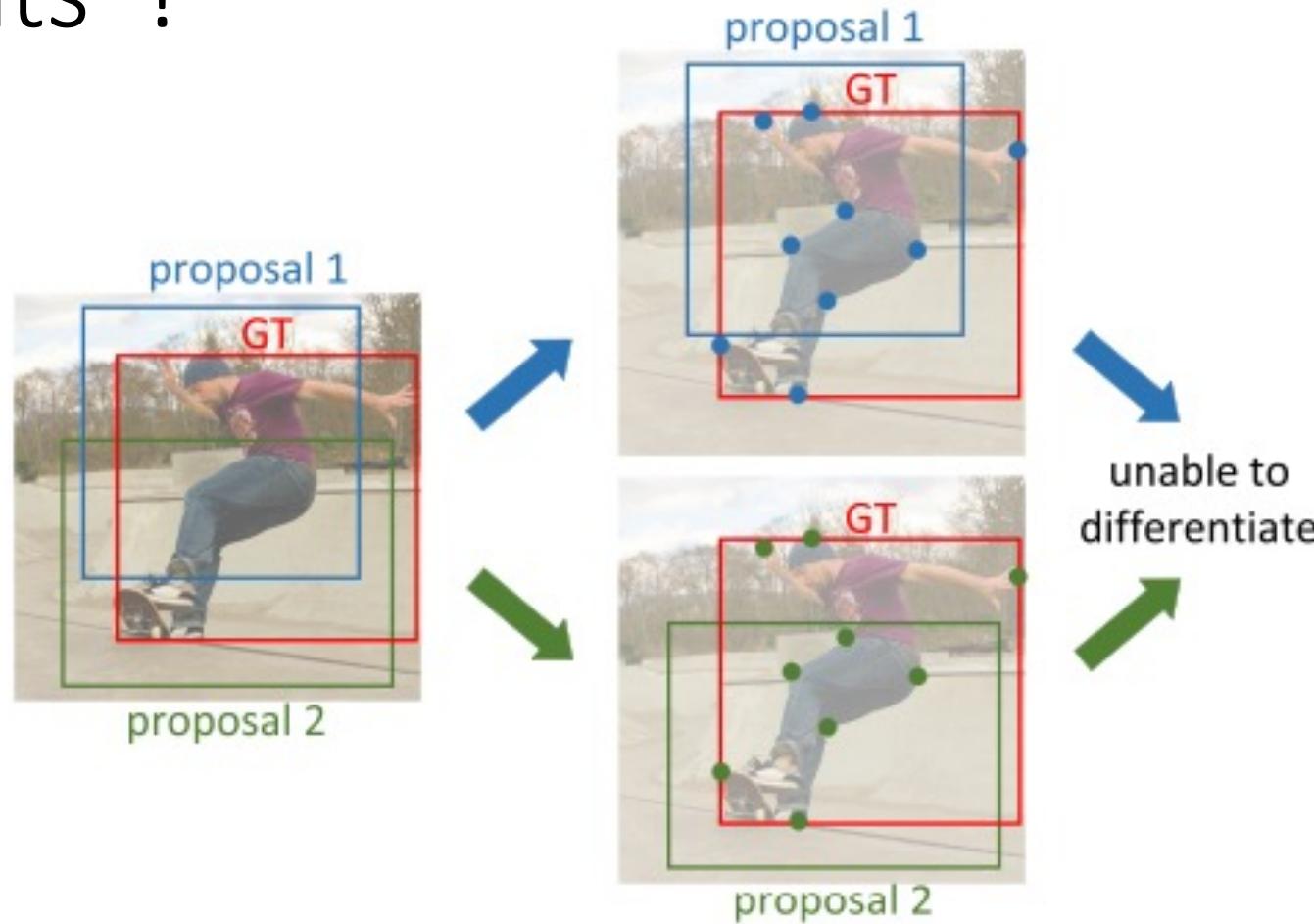
Deformable Convolutions are also very helpful for tracking!



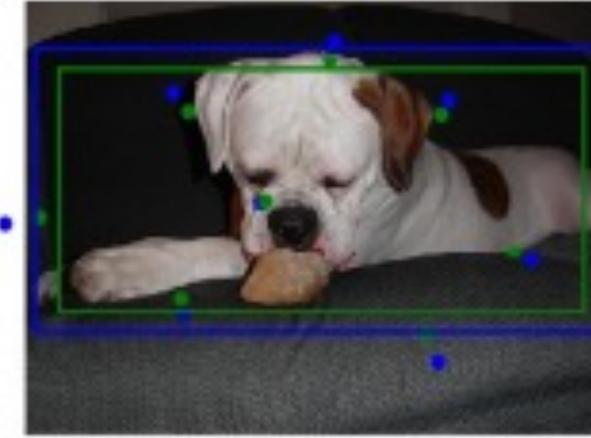
Distracting background region in the bounding box

Receptive field can deform to focus on the target object

Do Deformable Convolutions Capture Object “Keypoints”?



Do Deformable Convolutions Capture Object “Keypoints”?

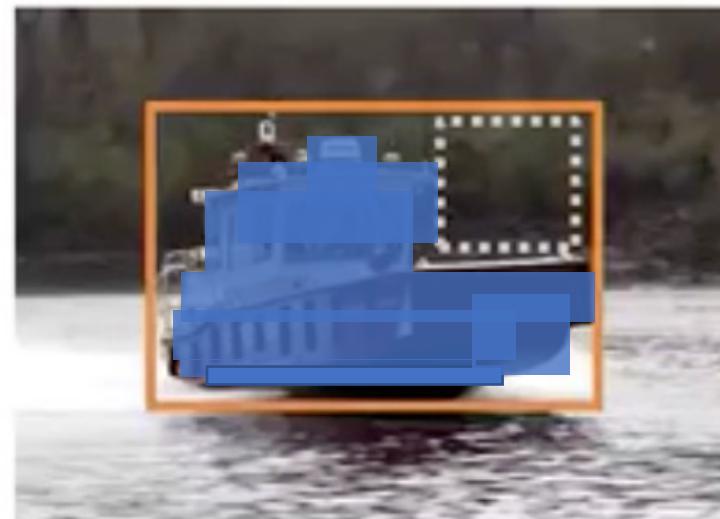


What is the right object representation?

Bounding box?



Segmentation?



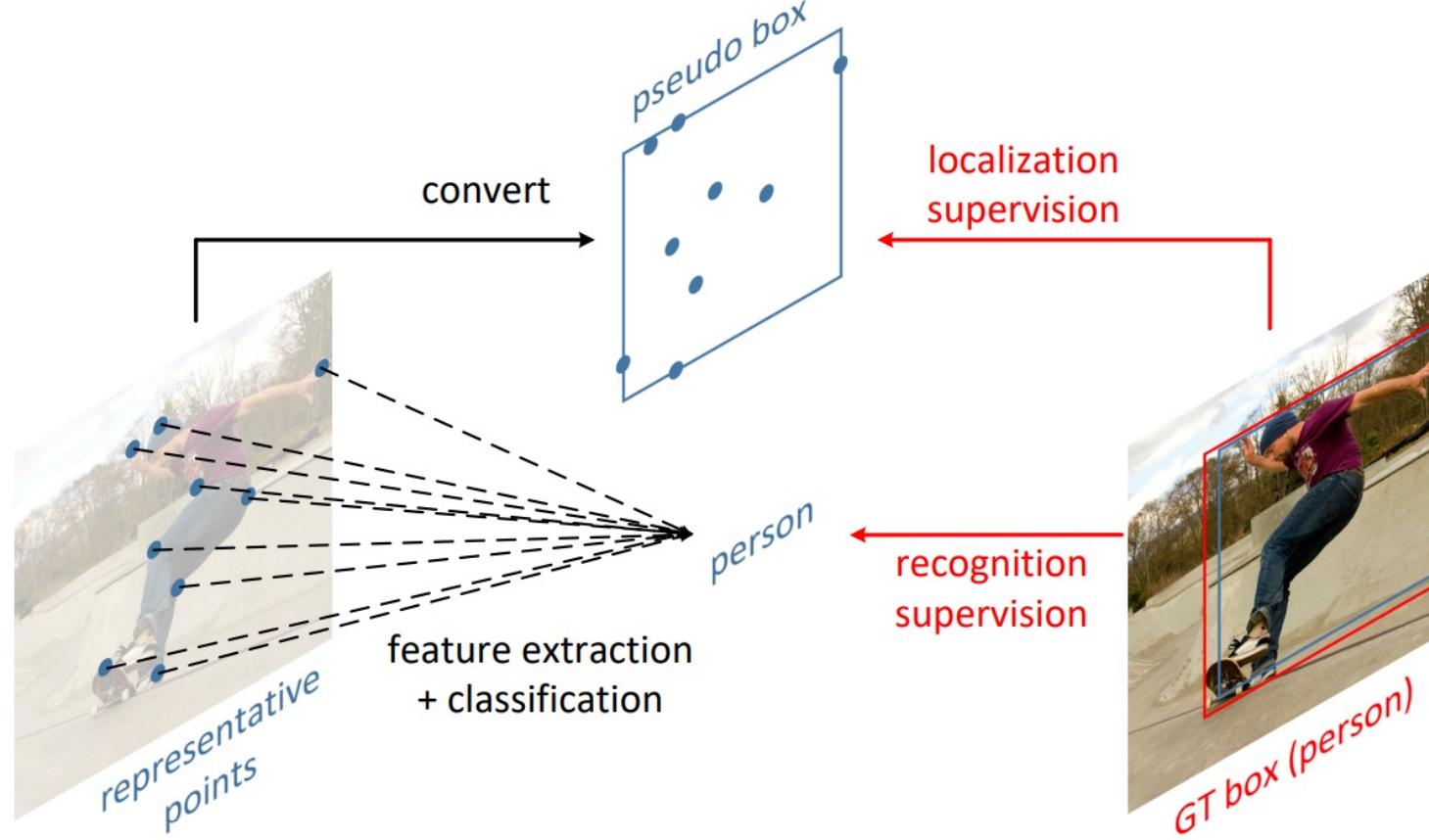
Set of Points



But this includes
the background

Requires segmentation
labels

RepPoints



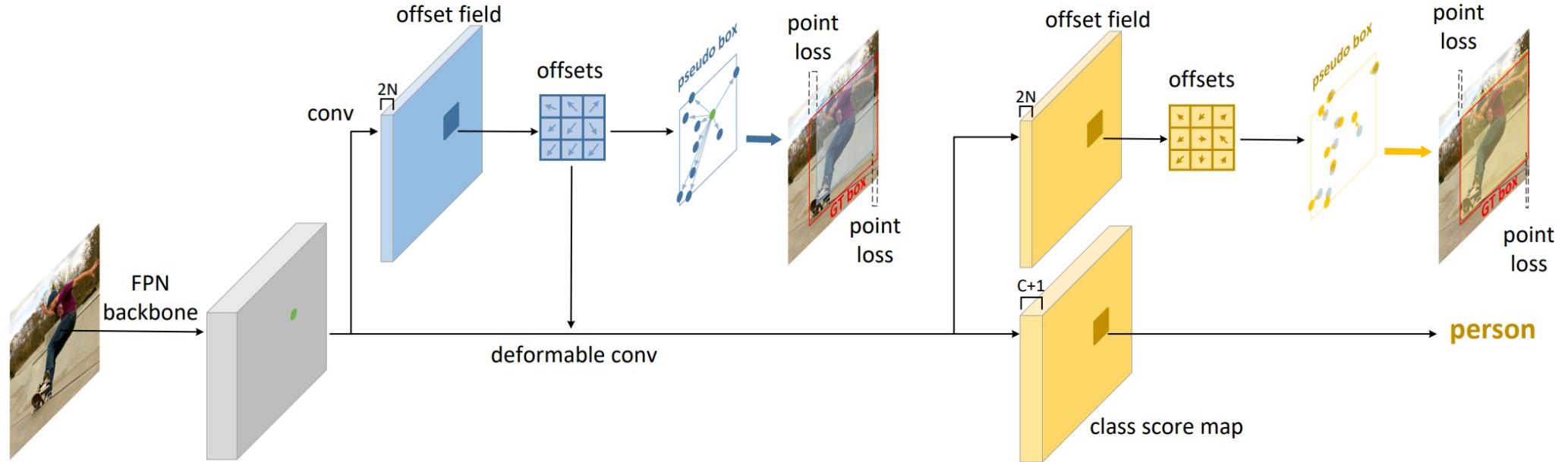
Convert from set of points to a pseudo-box using differentiable min-max:

$$B_p = (\min \{x_k + \Delta x_k\}, \min \{y_k + \Delta y_k\}, \max \{x_k + \Delta x_k\}, \max \{y_k + \Delta y_k\})$$

RepPoints



RepPoints



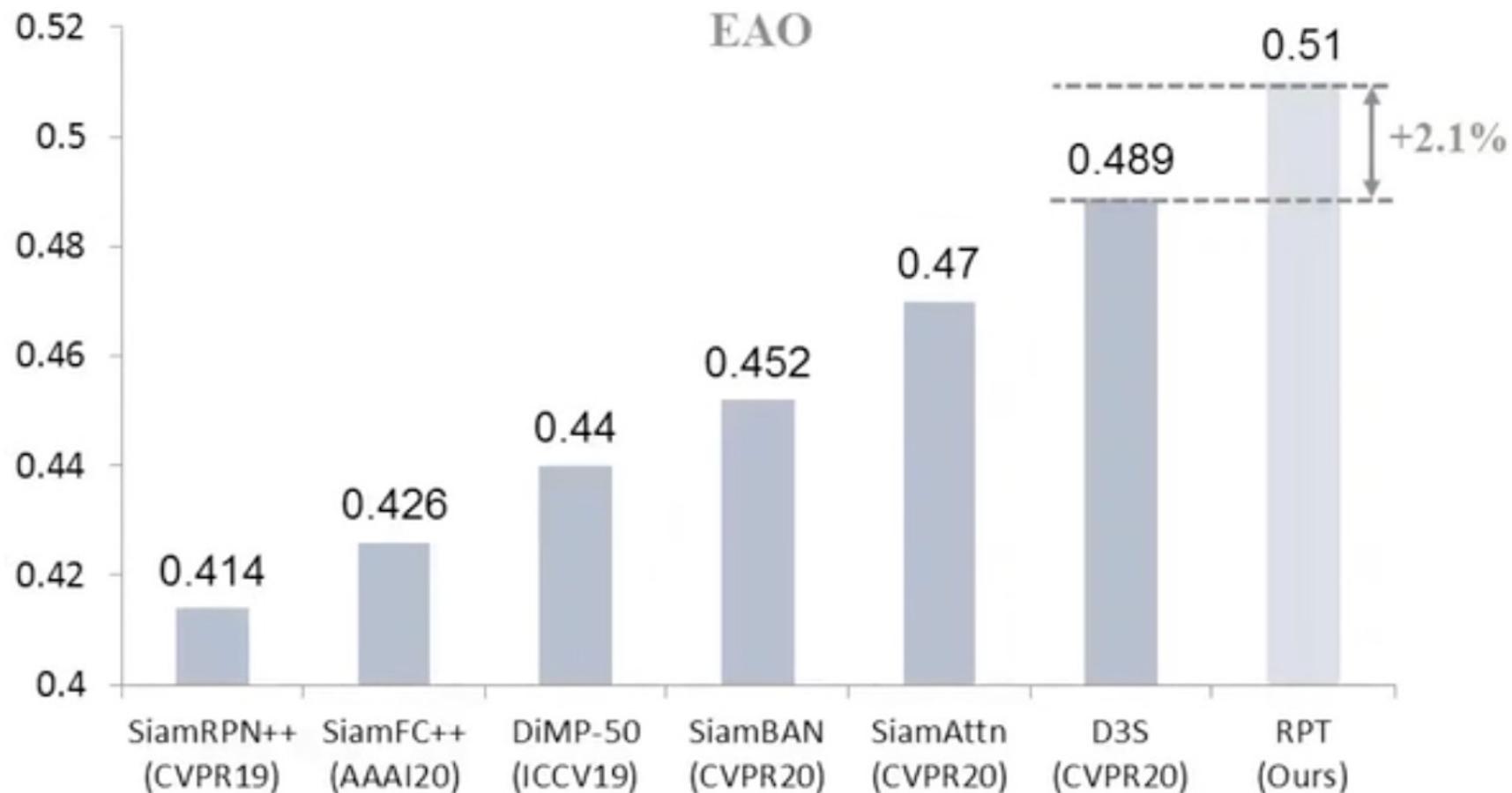
Extension: Dense RepPoints



Yang, Ze, et al. "Dense reppoints: Representing visual objects with dense point sets." ECCV 2020.

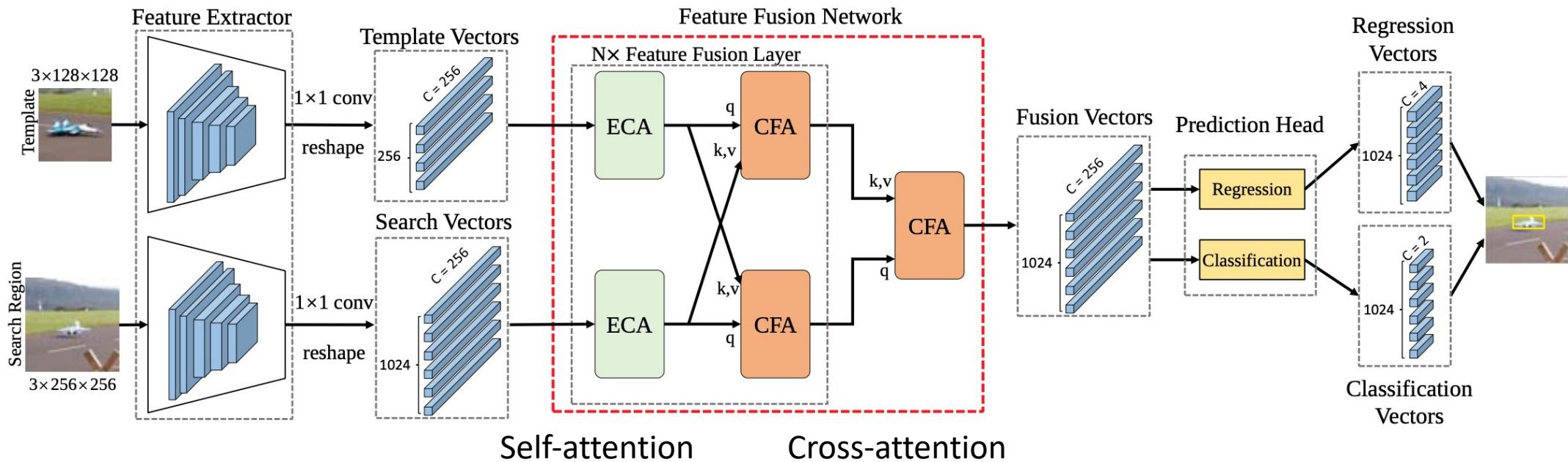
RPT: Uses RepPoints for Tracking

VOT2018



The Future of Tracking: Transformers?

Transformer Tracking (TransT)



Search Region
Cross -Attention
Map



Chen, Xin, et al. "Transformer tracking." *CVPR*. 2021.

Single-object Tracking: A summary

- Online training: Better robustness to distractors
- Offline training: Leverages prior knowledge to “learn to track”
- Best: Combine both, but especially online for coarse tracking and offline for BB refinement
- Fast online optimization: Gauss-Newton
- Learn to output localization confidence and optimize online
- Bounding box is not the best representation; we need to reason about object vs background (e.g. set of points, segmentation mask)
- Network needs to learn to compare the query to the search
 - Convolution, dot-product attention (transformers)
- Watch latest “tricks” in detection and think about how to adapt to tracking!

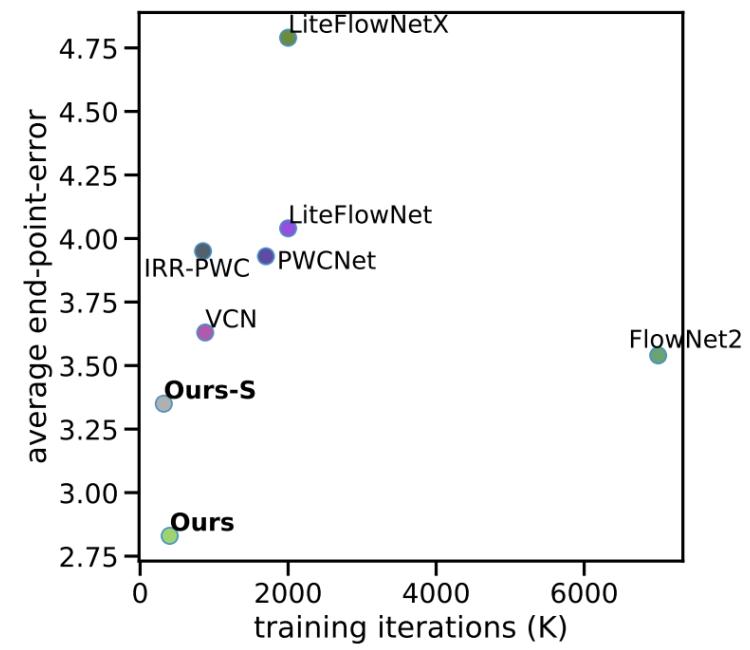
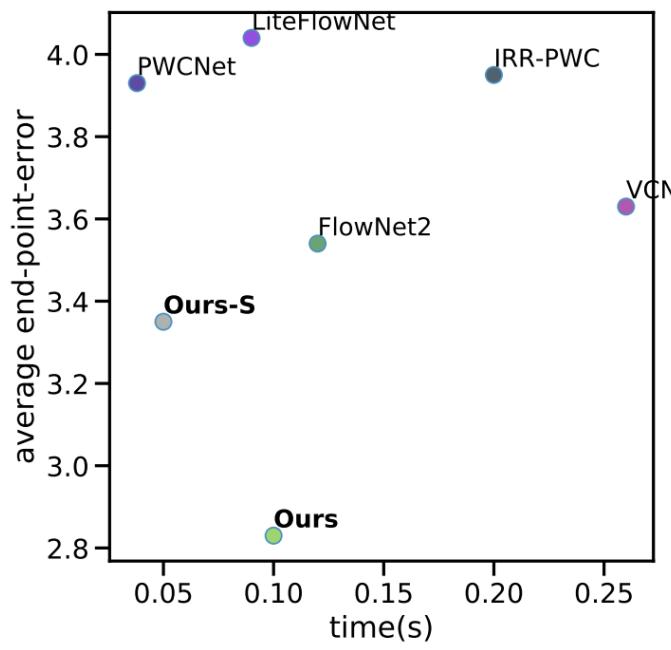
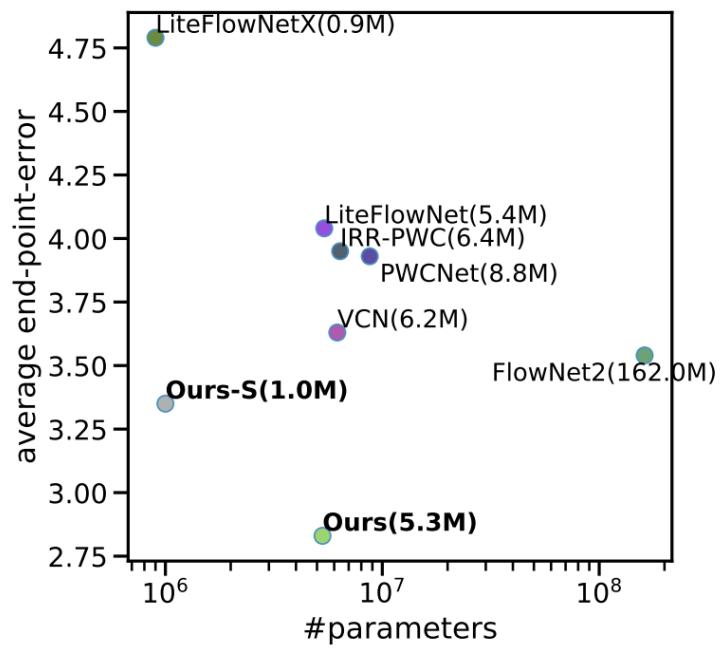
Types of Tracking

2D

- **Single Object Tracking**
- Joint Segmentation and Tracking
- Multi-object tracking
- Dense Optical Flow
- Self-supervised tracking

3D

- 3D Multi-object Tracking
- 3D Scene Flow
- Self-supervised 3D Scene Flow



Teed, Zachary, and Jia Deng. "Raft: Recurrent all-pairs field transforms for optical flow." *ECCV*, 2020.

RAFT Summary

- Compute features on previous and current frame
- Compute a correlation volume
- Based on initial flow estimate, look up features in correlation volume
- Iteratively update the flow

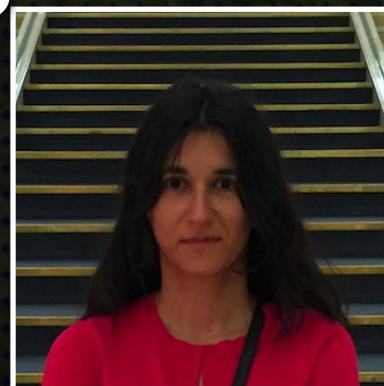
Particle Video Revisited: Tracking Through Occlusions Using Point Trajectories



Adam W. Harley



Zhaoyuan Fang



Katerina Fragkiadaki

Modern object trackers survive occlusions using temporal priors



Rajasegaran et al., Tracking People by Predicting 3D Appearance, Location & Pose, CVPR 2022.

Goal: Track any pixel across occlusions



Goal: Track any pixel across occlusions



Existing structured trackers do not solve this problem

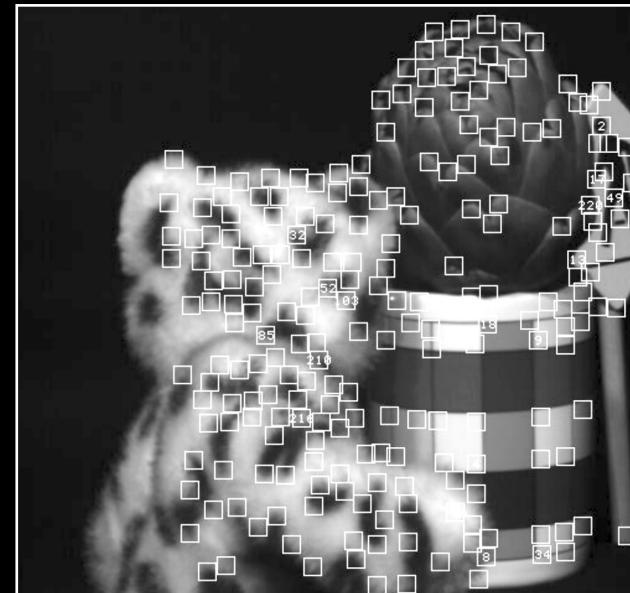
Semantic keypoints



Provides tracks for faces, hands, joints
Not general-purpose

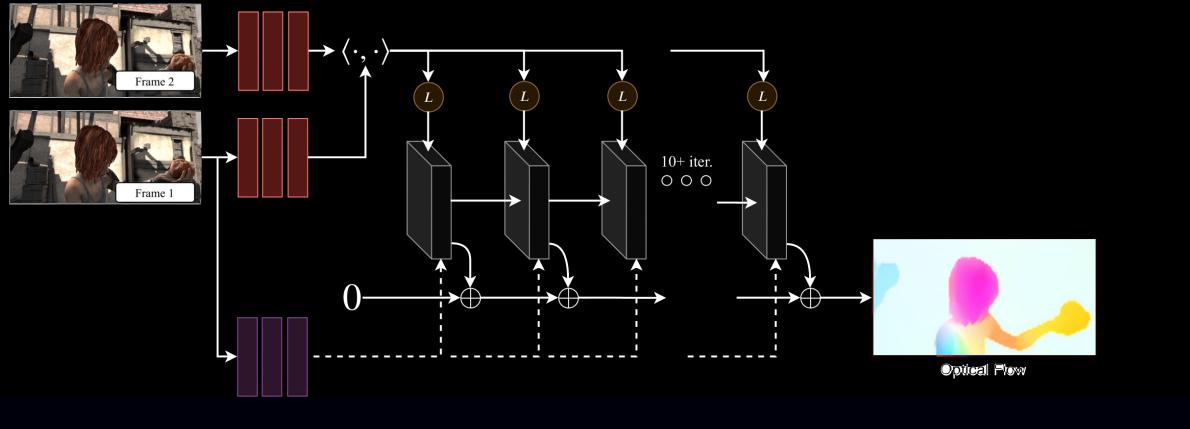
Jhuang et al. "Towards understanding action recognition." ICCV 2013.

Structure from Motion

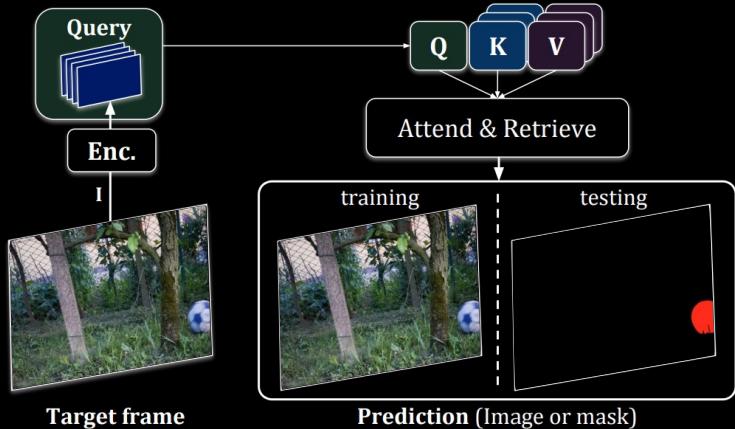


Existing general-purpose trackers do not solve this problem

Optical Flow



Feature matching



Describes motion of visible pixels
Propagates labels by chaining flows
Narrow timespan

Teed & Deng. RAFT: Recurrent all-pairs field transforms for optical flow. ECCV 2020.

Provides correspondable features
Propagates labels using feature similarity
Weak/absent temporal prior

Lai et al. MAST: A memory-augmented self-supervised tracker, CVPR 2020

Inspiring work: “Particle Video”



Sand & Teller. "Particle video: Long-range motion estimation using point trajectories." IJCV 80.1 (2008): 72.

Modeling time can help us track through occlusions



Input



video



target

Initialize a zero-velocity trajectory

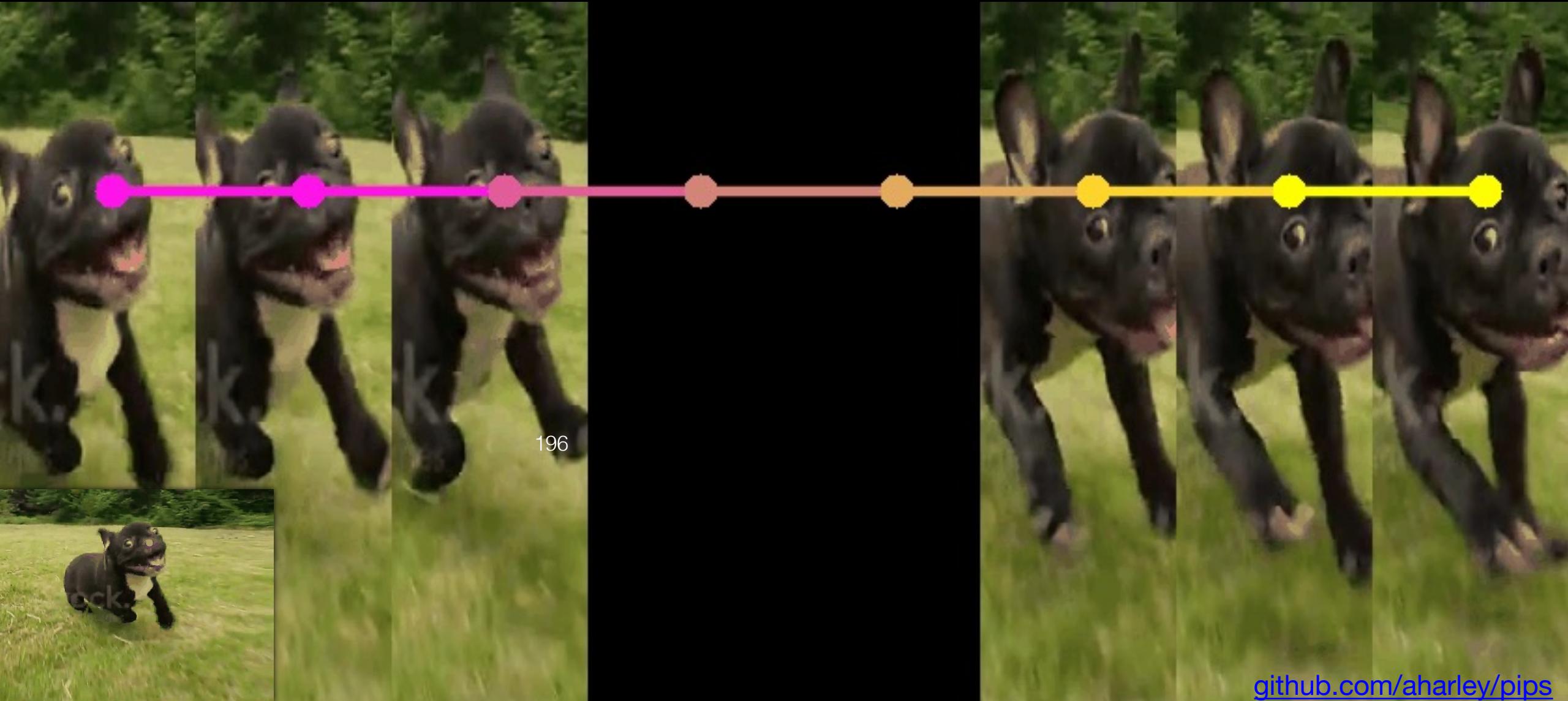


Iteratively refine the trajectory

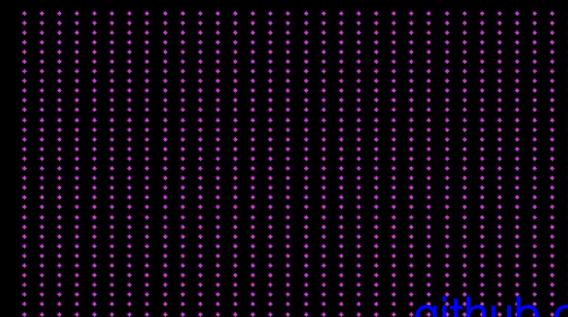
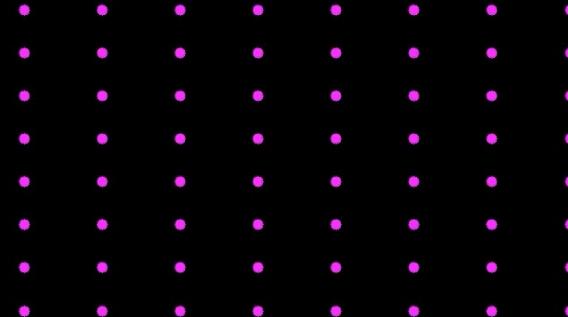


github.com/aharley/pips

Consider the full sequence simultaneously

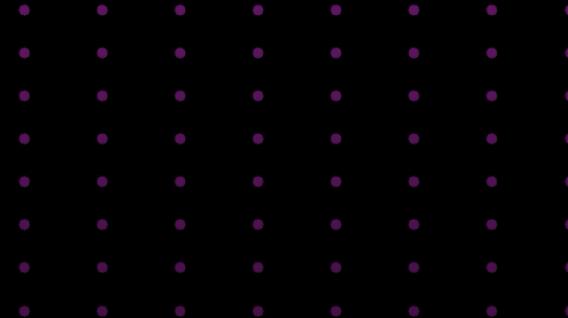


Track every particle independently



github.com/aharley/pips

Track every particle independently



Persistent Independent Particles (PIPs)

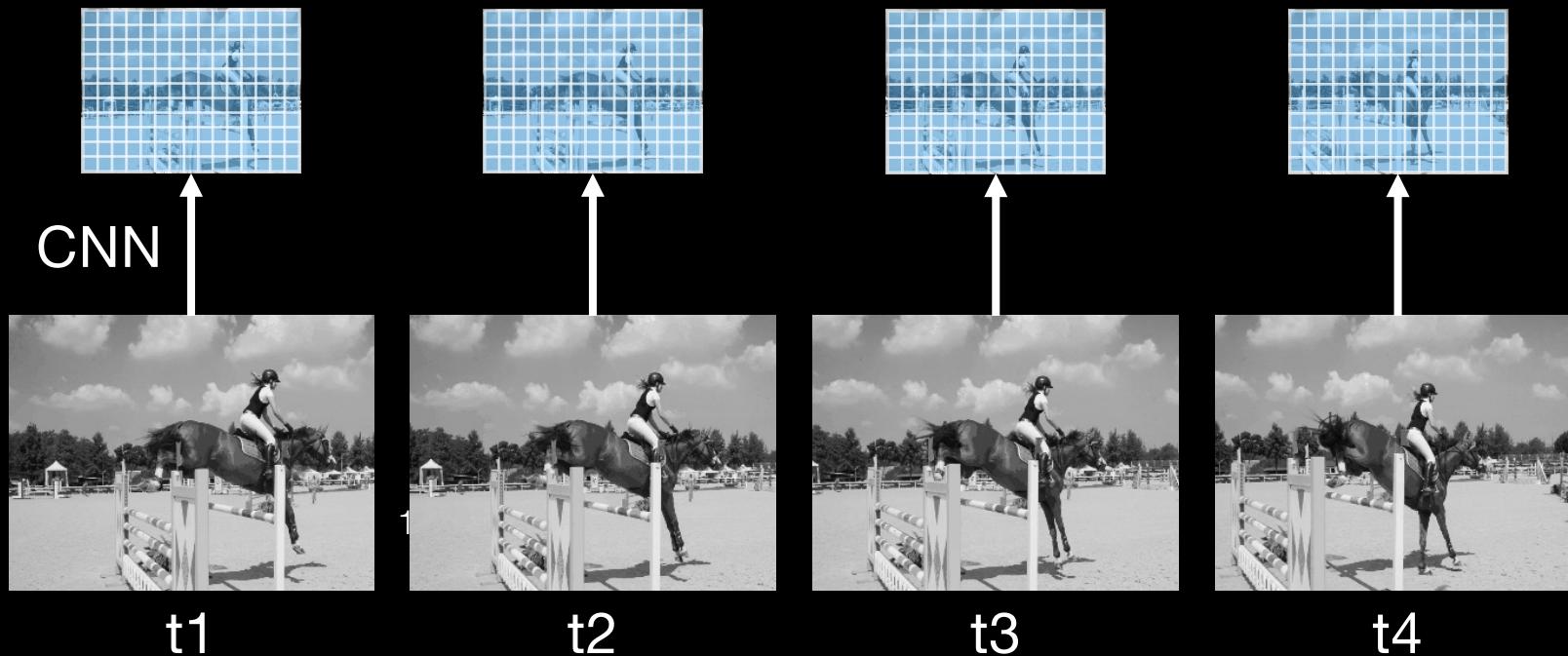


198

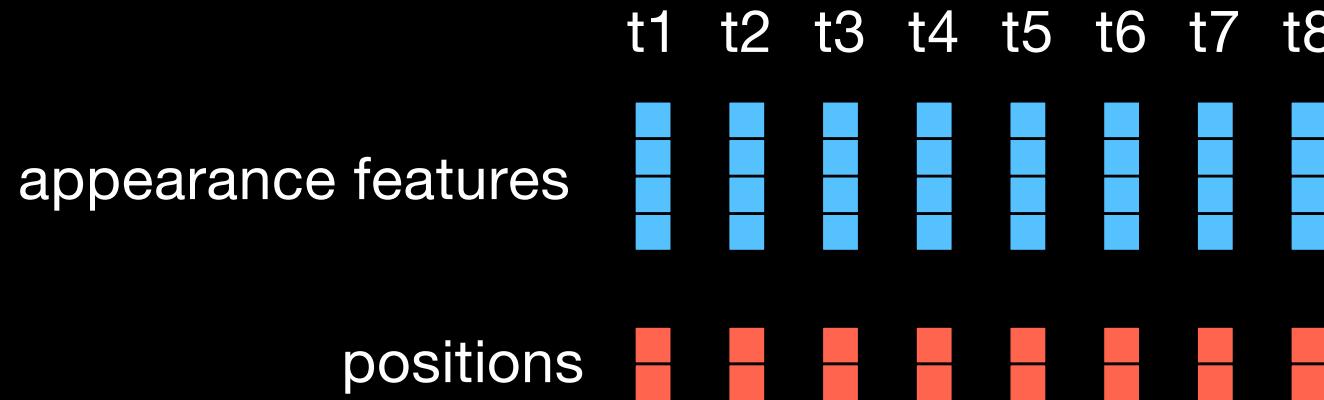


github.com/aharley/pips

Before tracking begins, compute features for all frames

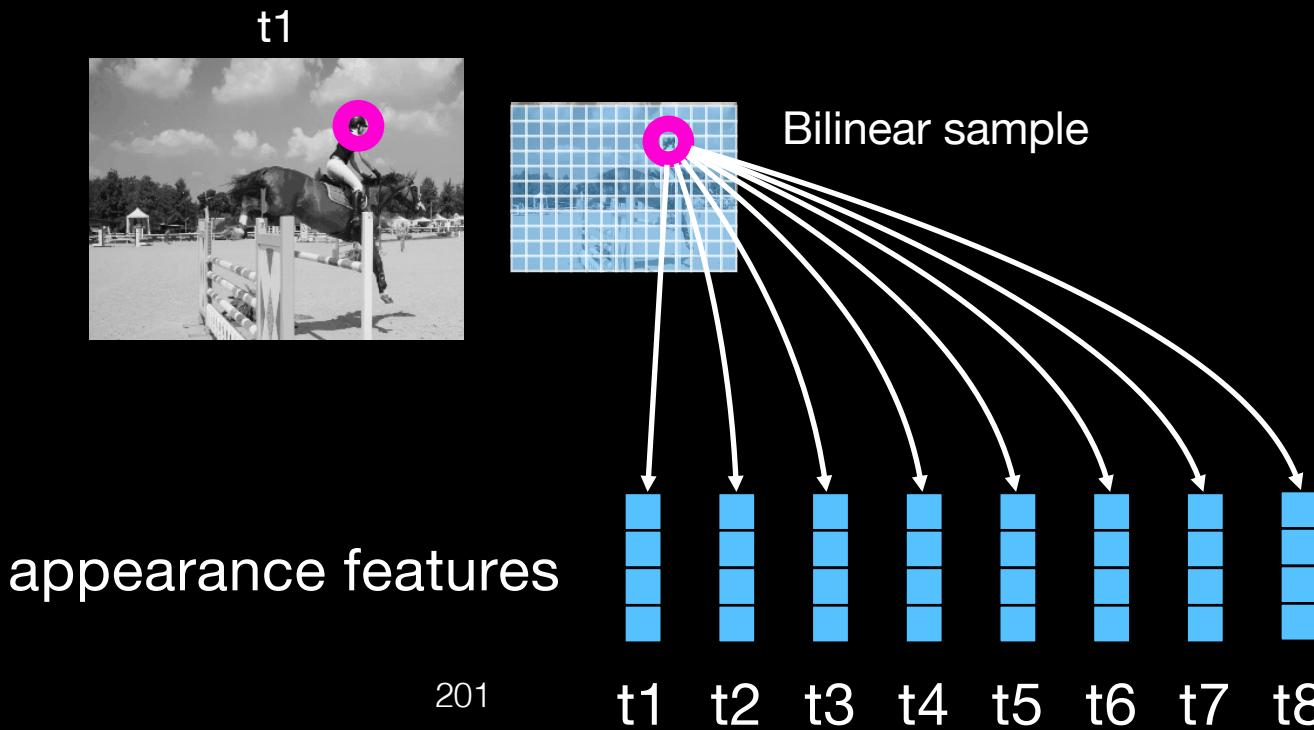


Represent the target with a list of features and positions across time

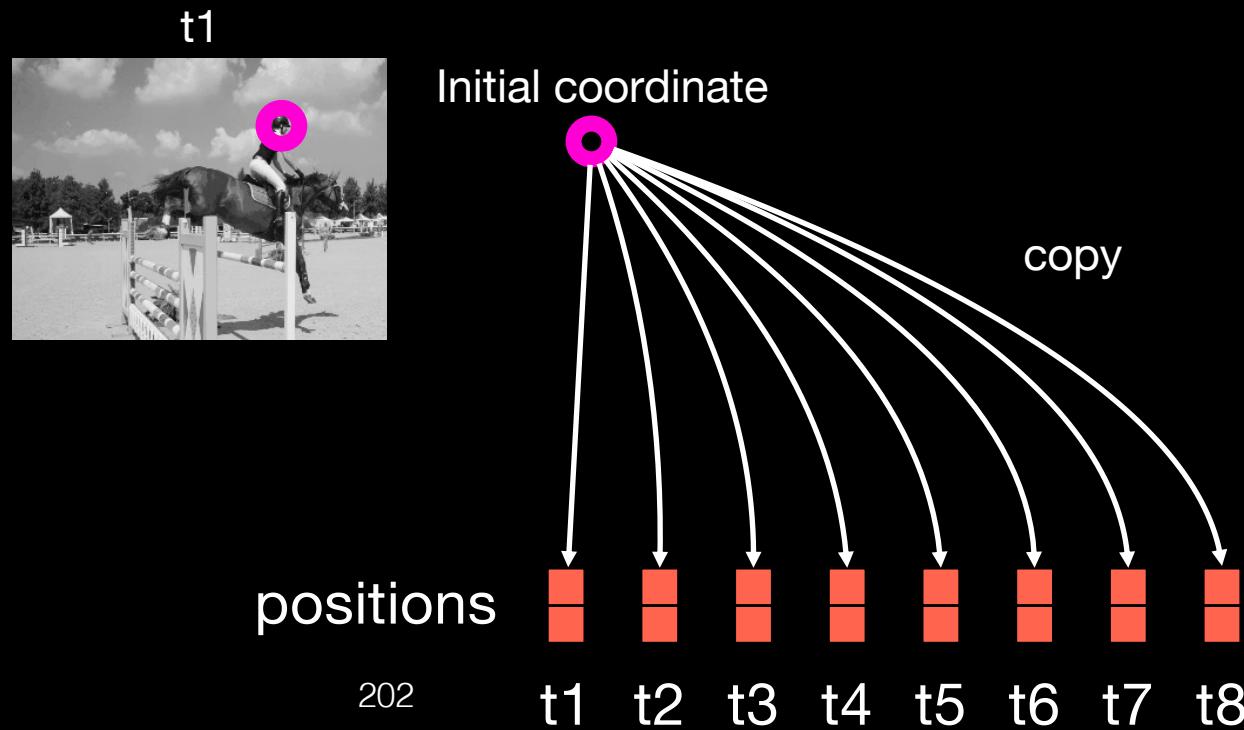


200

Initialize the feature list using a bilinear sample

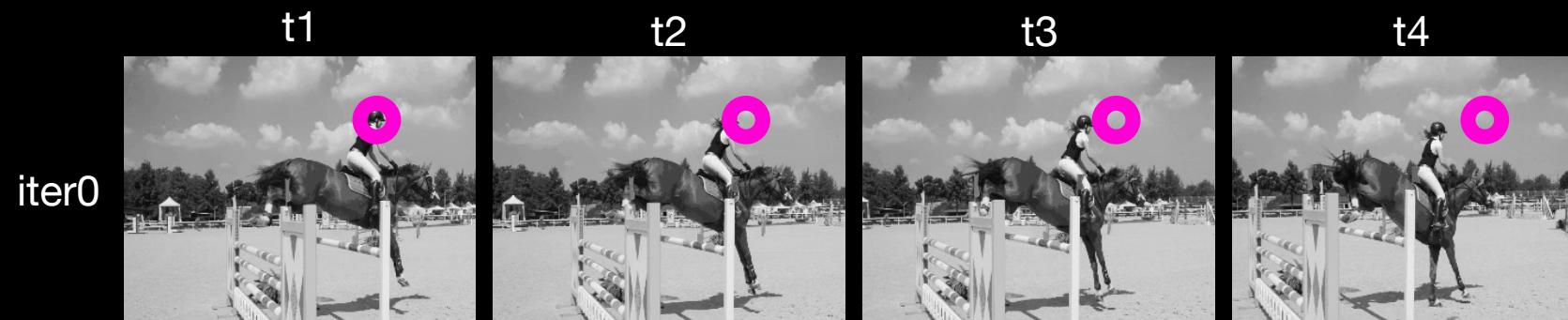


Initialize the position list with zero velocity



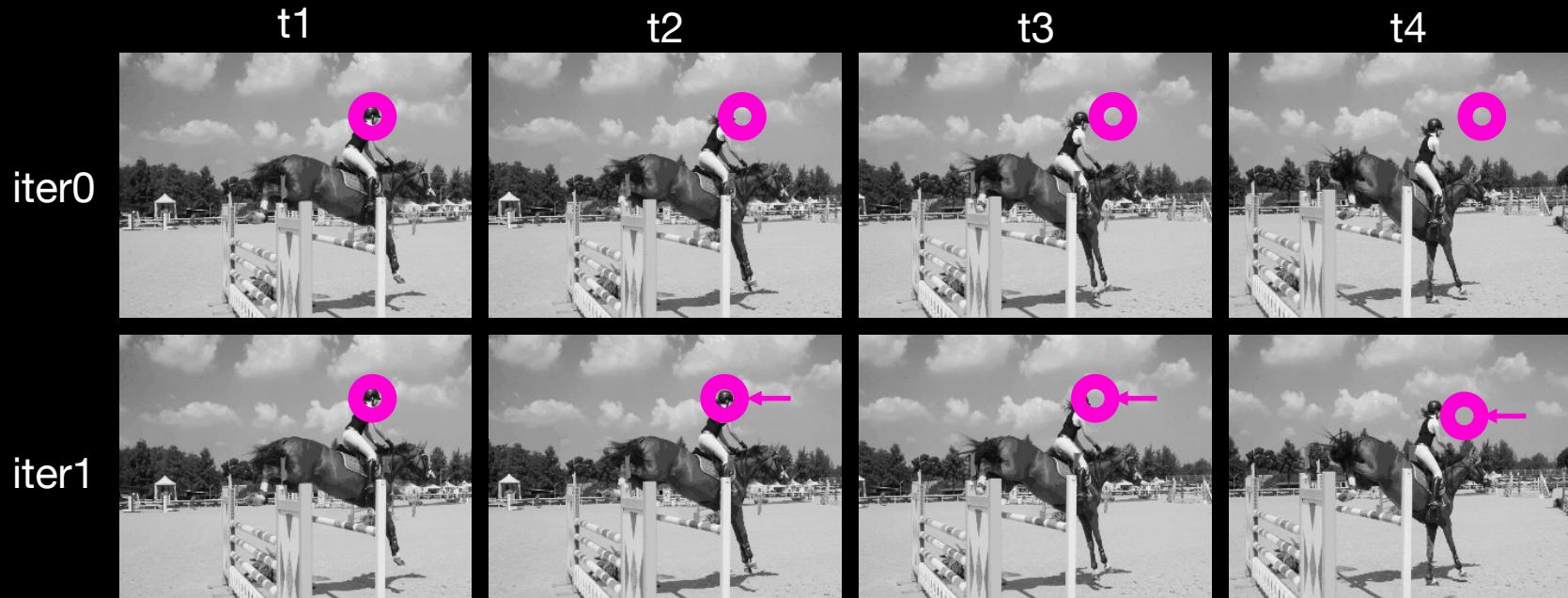
202

Inference: iteratively update positions and features



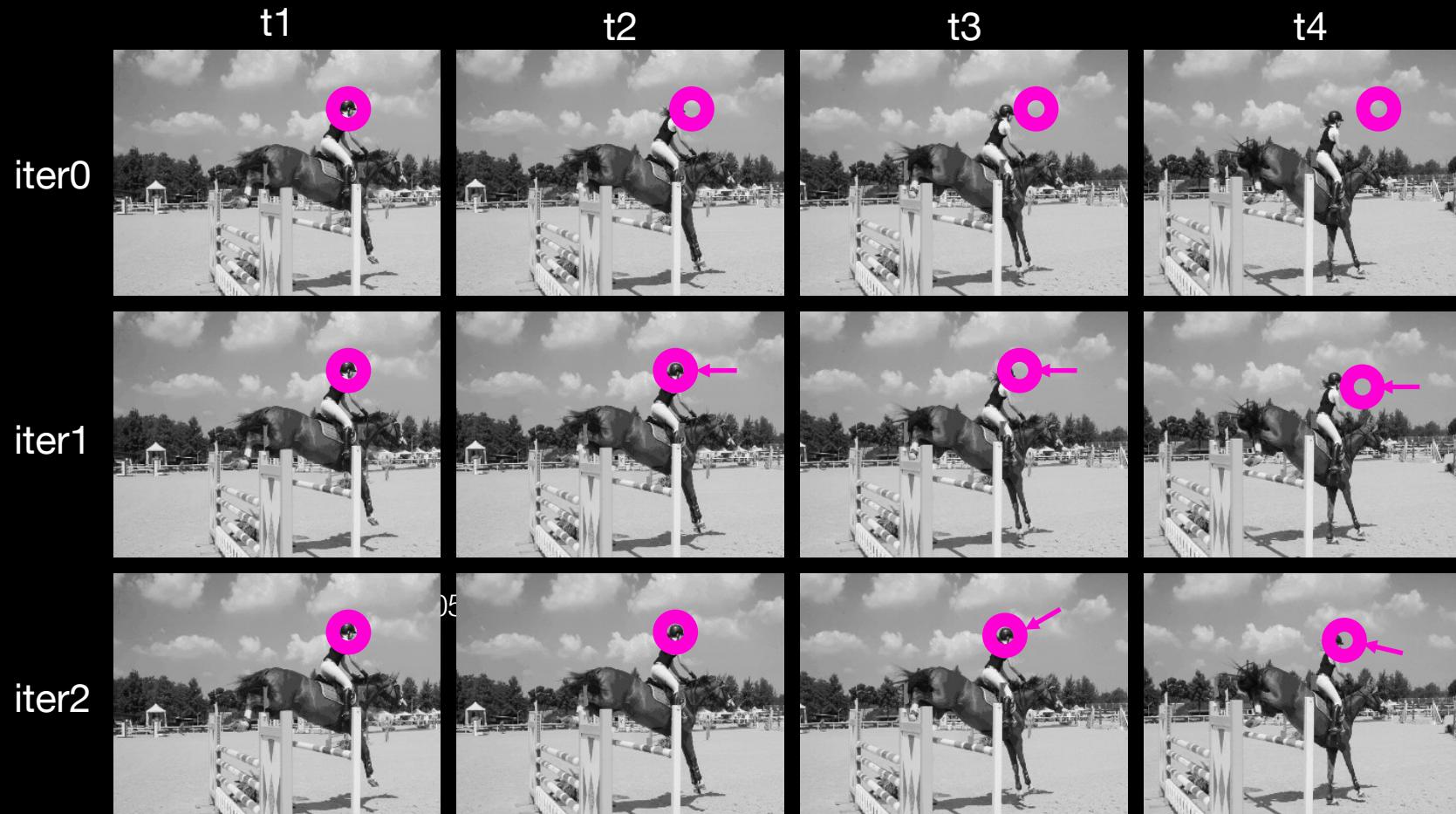
203

Inference: iteratively update positions and features

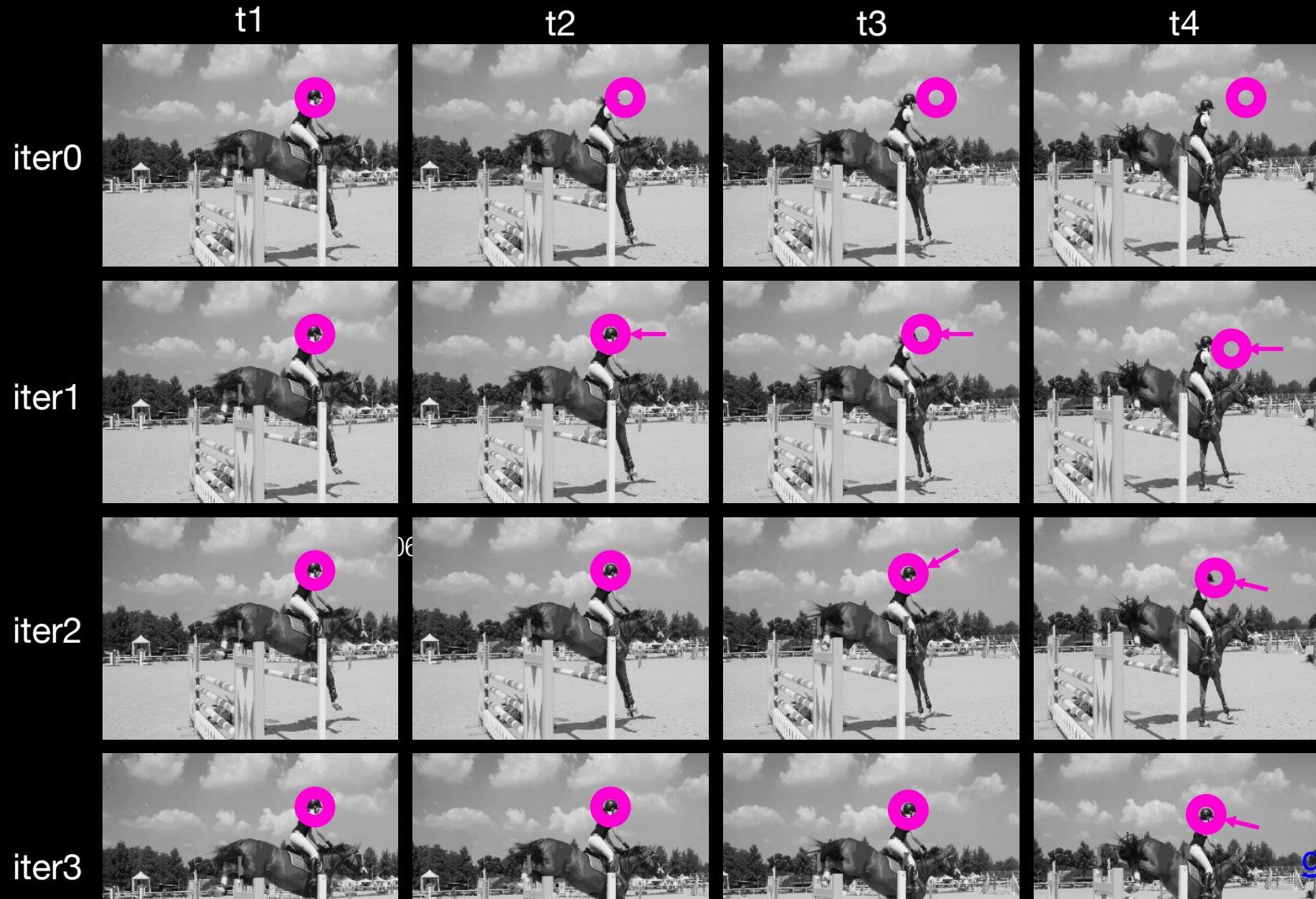


204

Inference: iteratively update positions and features



Inference: iteratively update positions and features



github.com/aharley/pips

Show the updater the scores of local deltas

feature_t



featuremap_t



score map

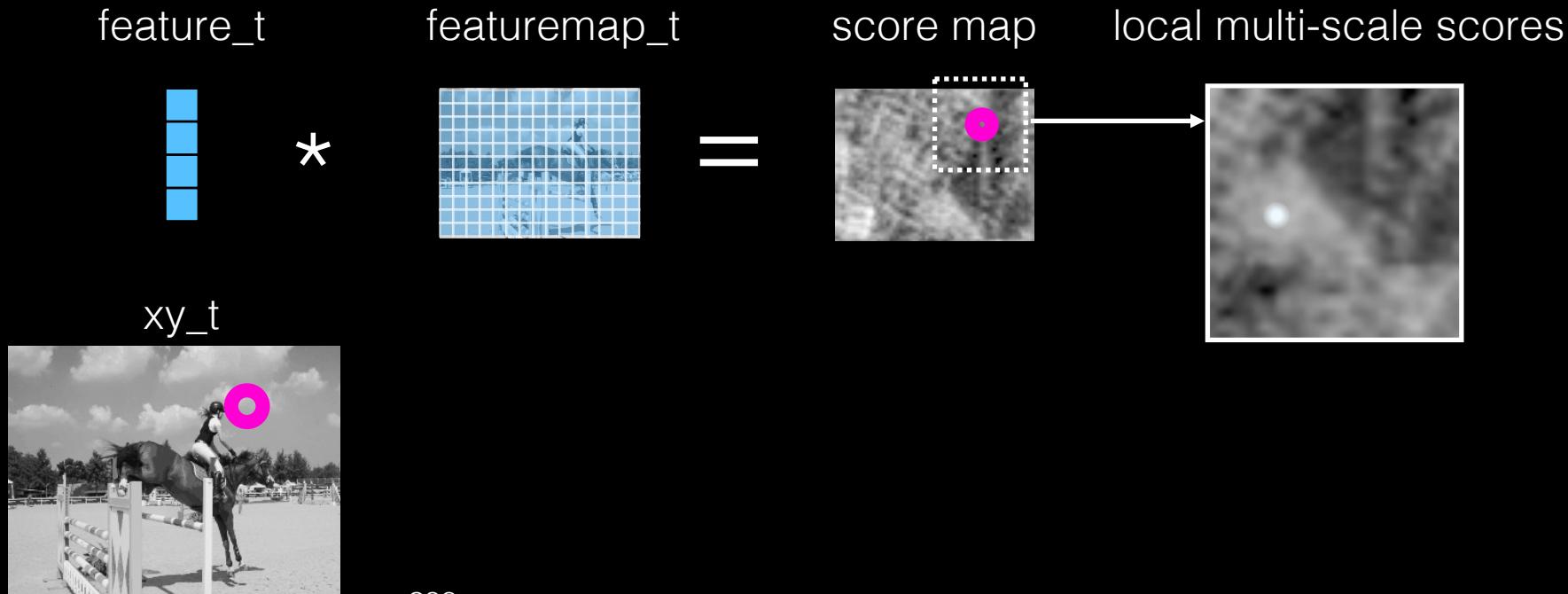
=



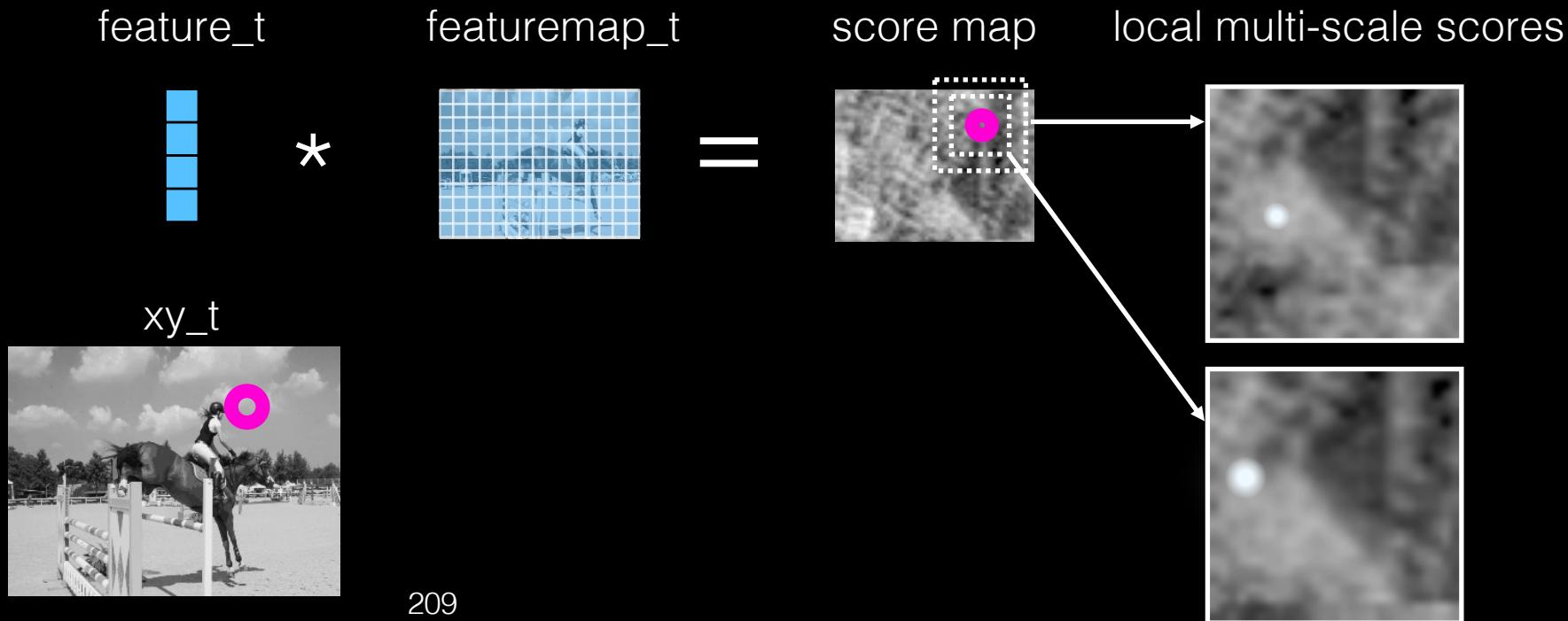
xy_t



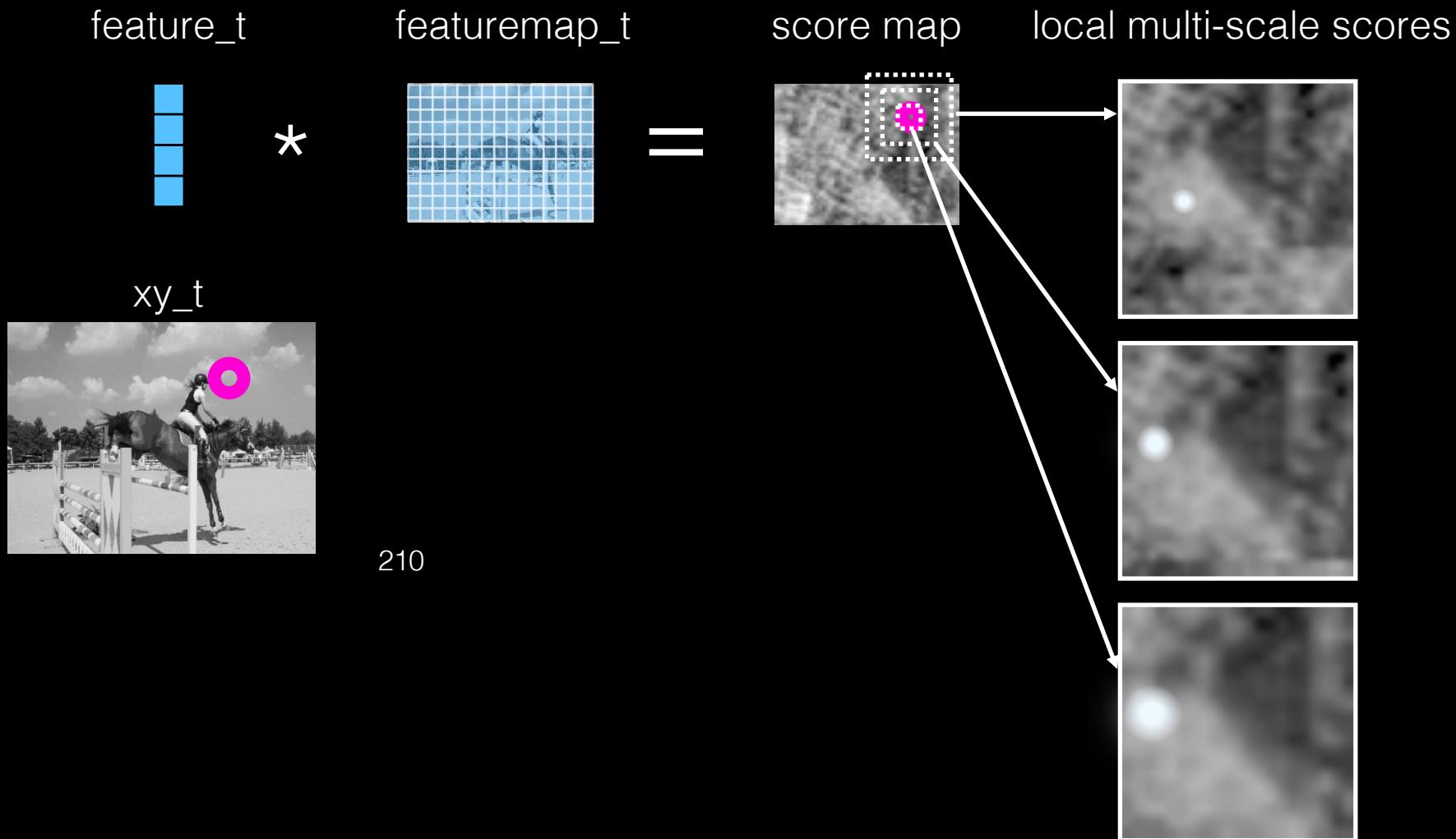
Show the updater the scores of local deltas



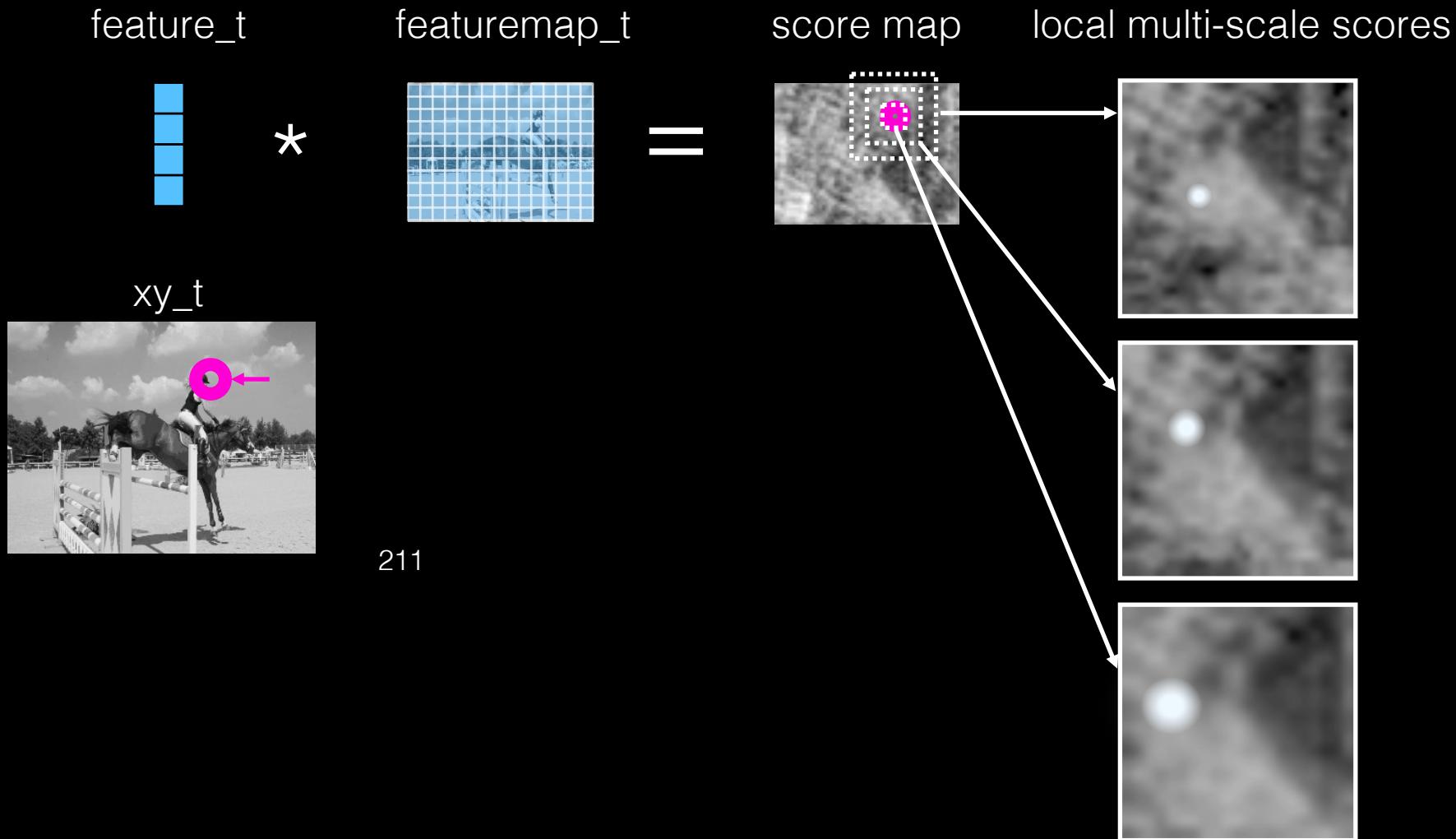
Show the updater the scores of local deltas



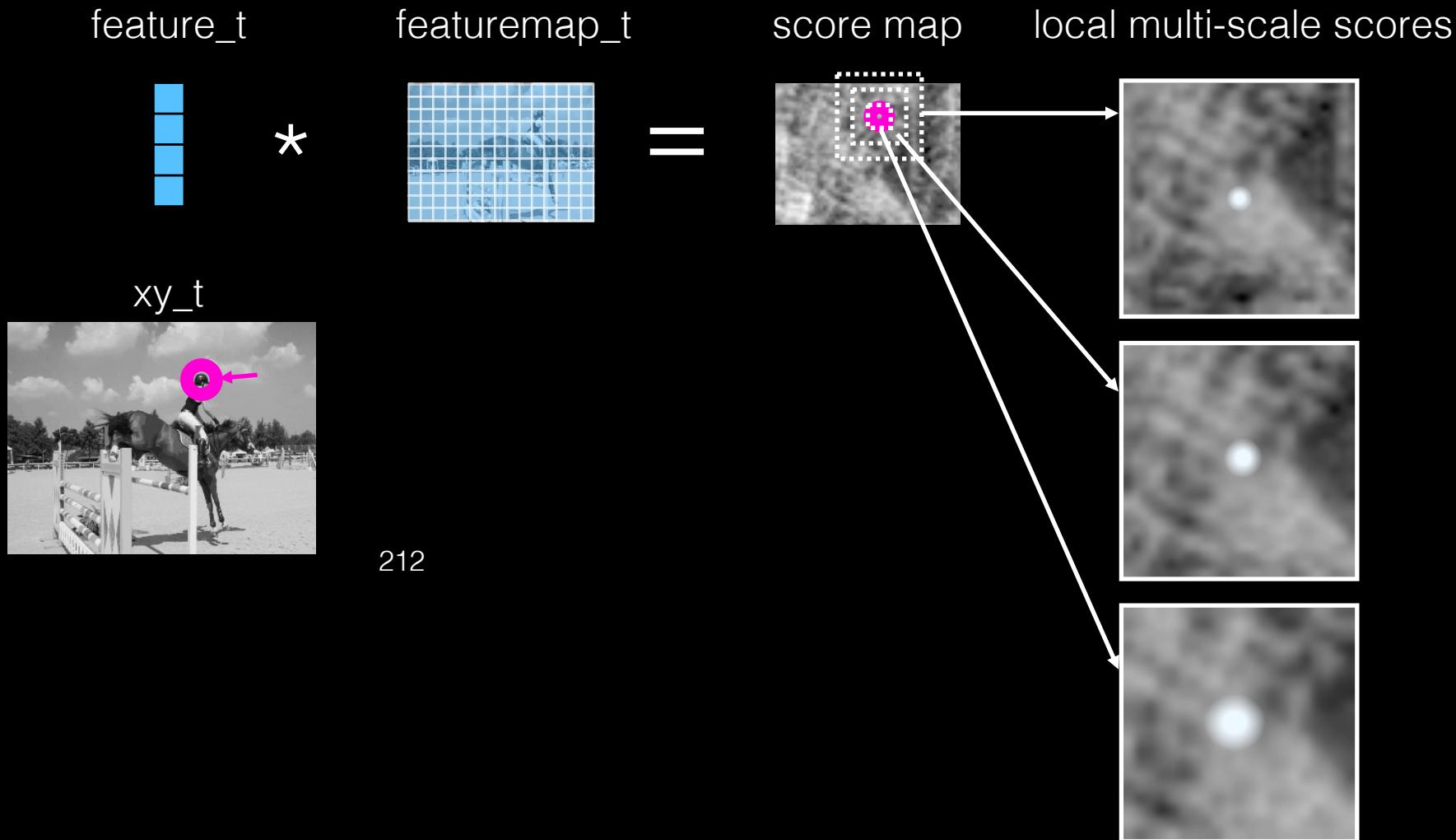
Show the updater the scores of local deltas



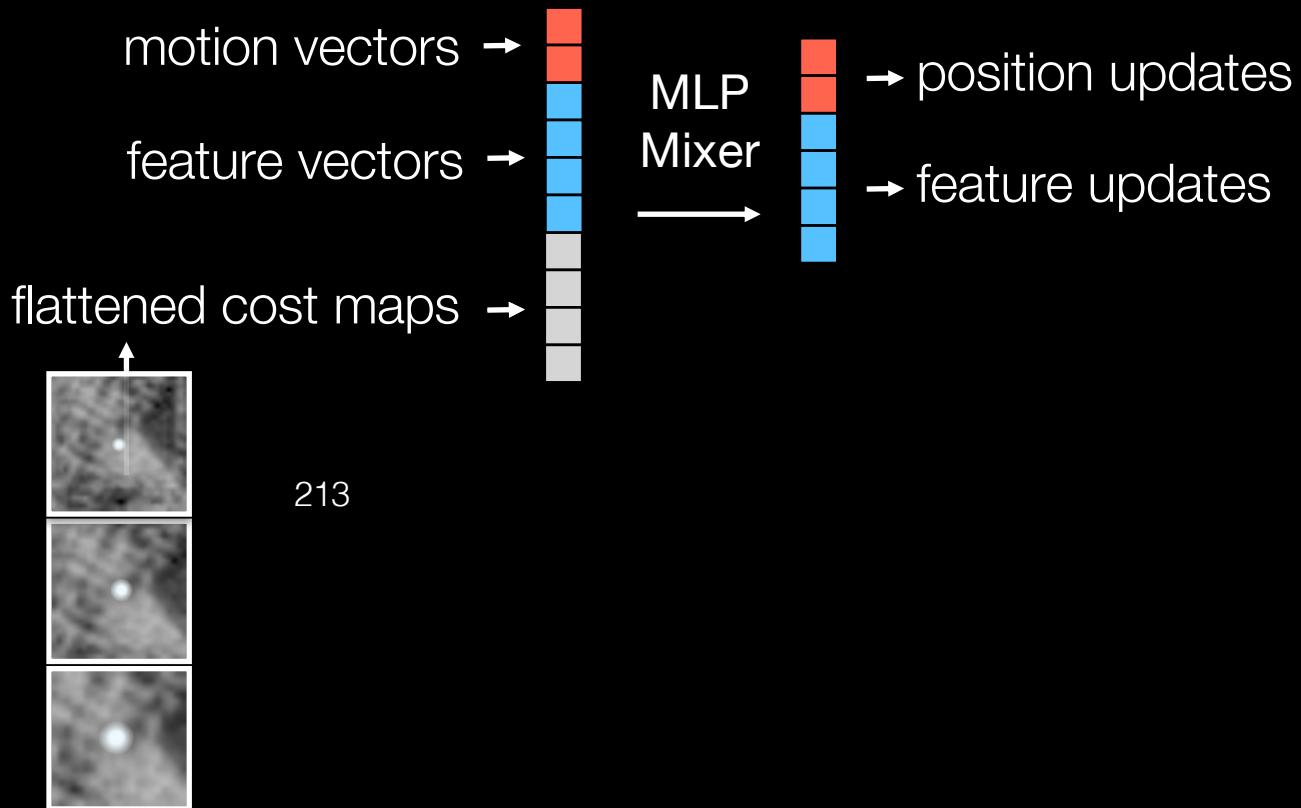
Show the updater the scores of local deltas



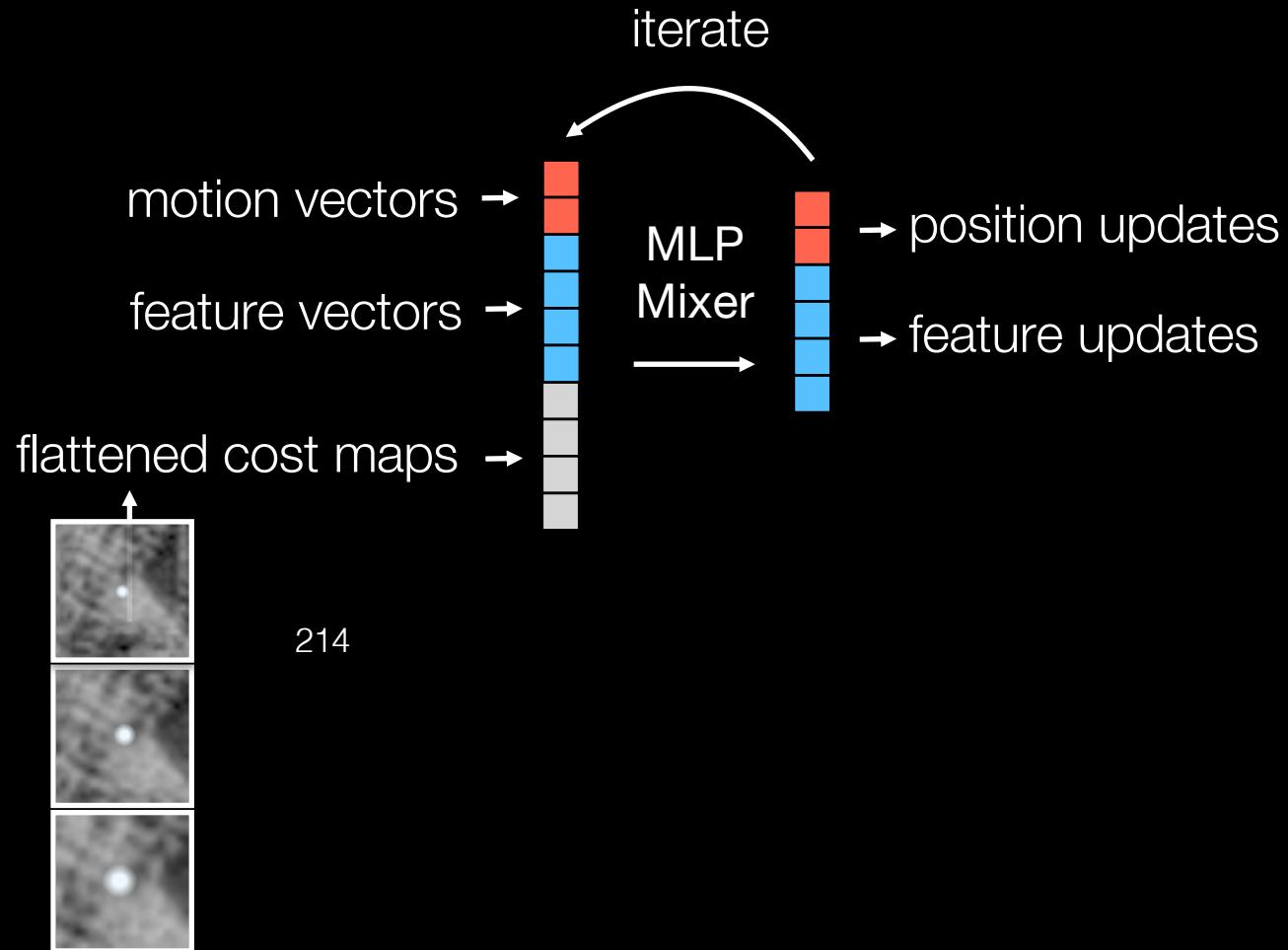
Show the updater the scores of local deltas



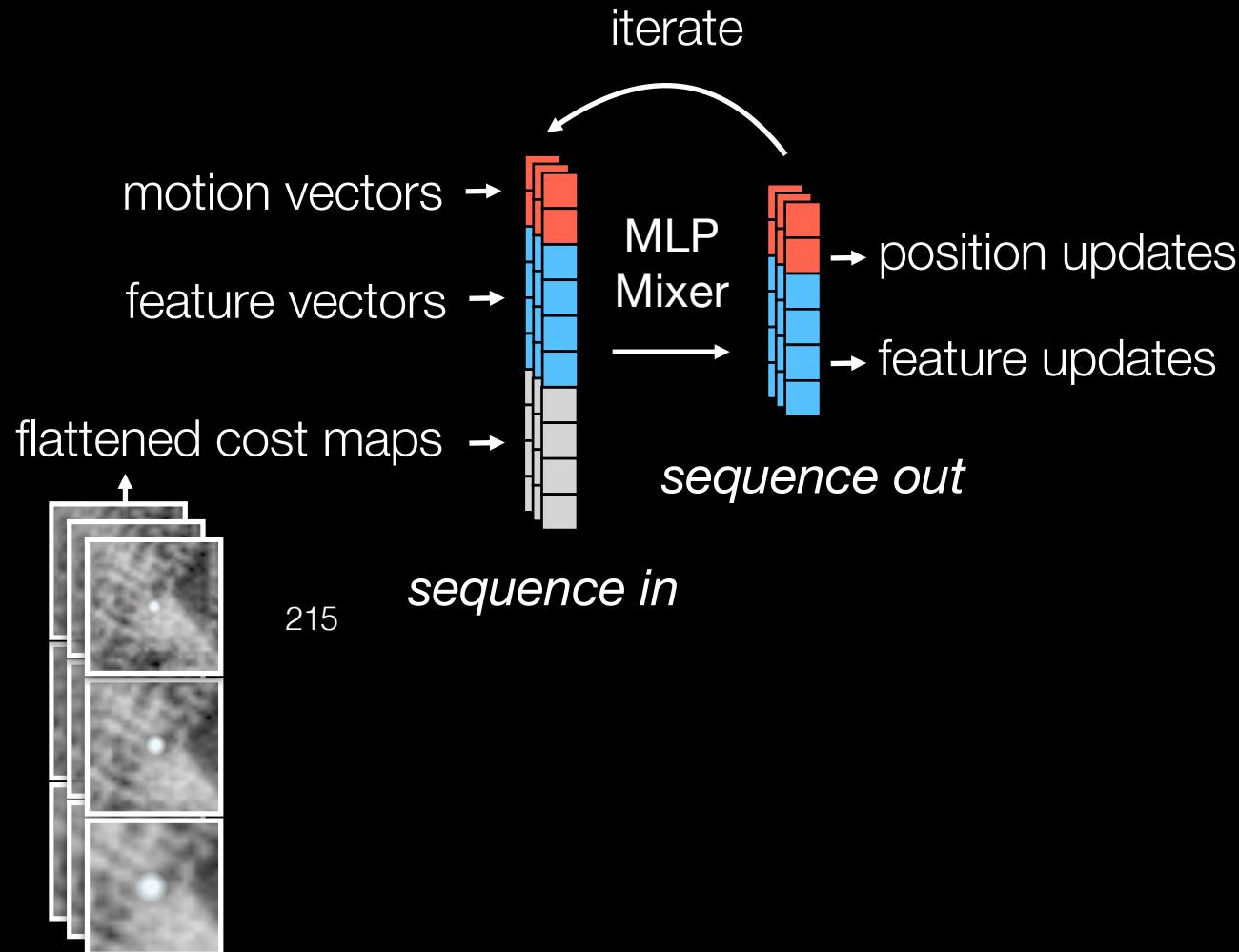
Input: score maps, features, and motions. Output: residual updates



Iterate 8x



Consider the full sequence simultaneously



Correspondence is allowed to fail!



We train our model with ground truth trajectories in synthetic data



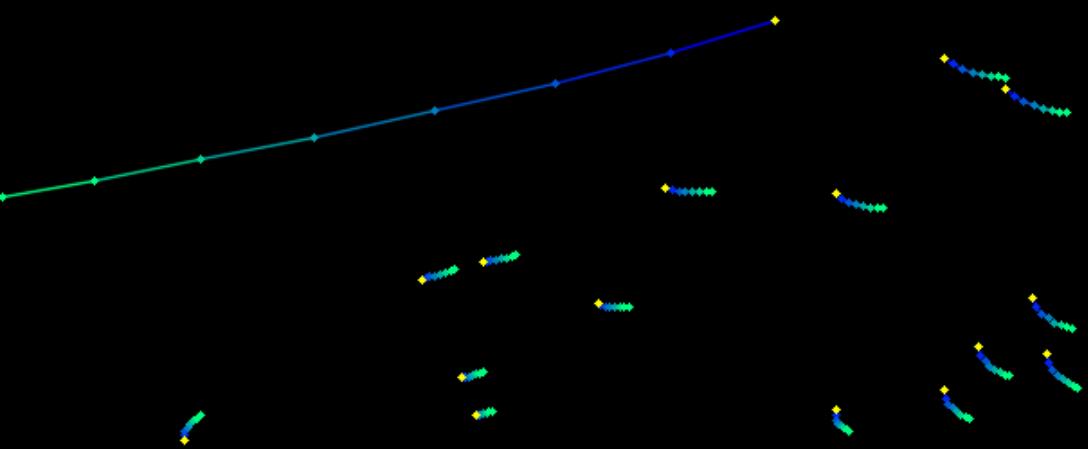
RGBs



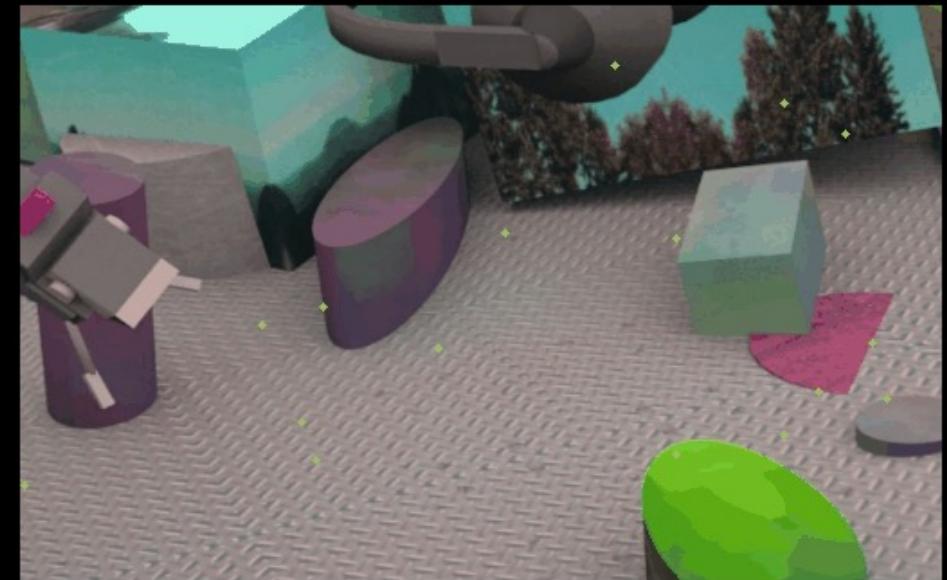
Annotations

Qualitative results in FlyingThings++

1.3557



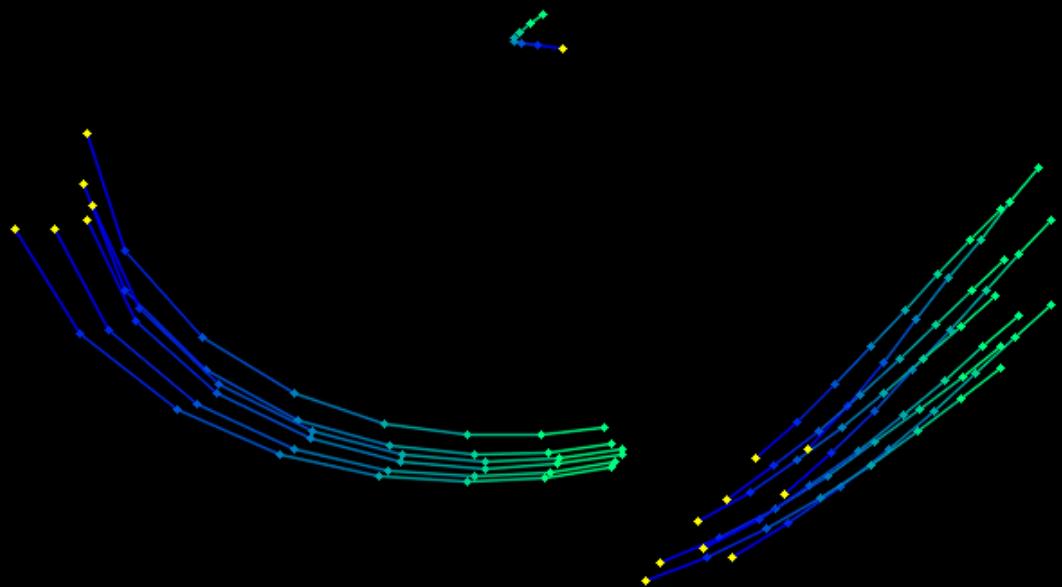
Iterative updates



Final result

Qualitative results in FlyingThings++

5.16689



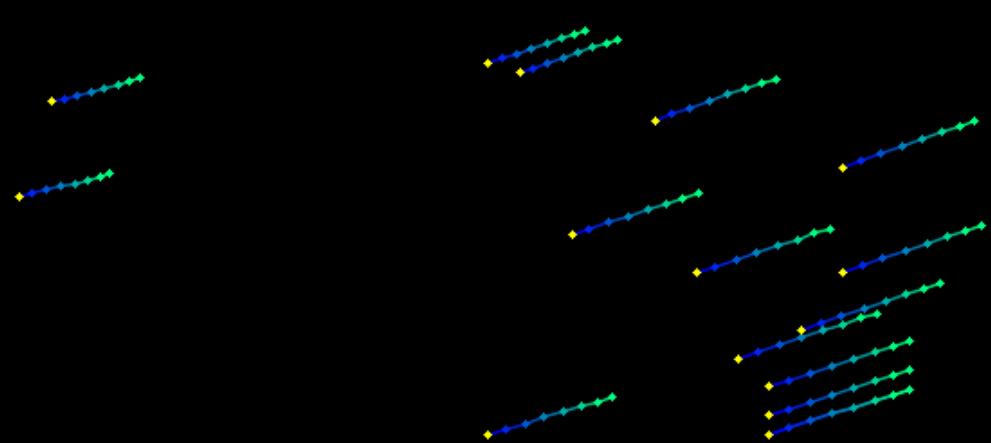
Iterative updates



Final result

Qualitative results in FlyingThings++

20.9623



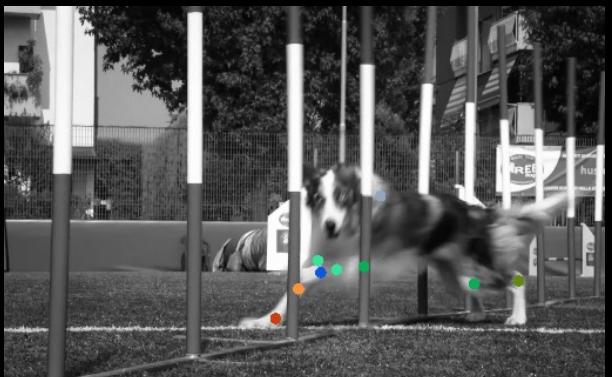
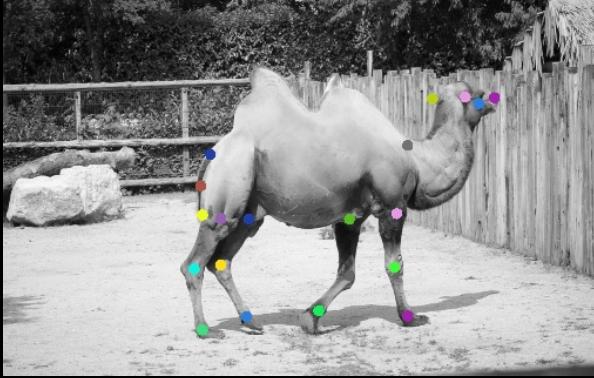
Iterative updates



Final result

Long-range keypoint propagation in BADJA

Biggs et al. Creatures great and SMAL: Recovering the shape and motion of animals from video. ACCV 2018.



Goal: propagate each keypoint from frame 1 all the way to the end

PIPs outperforms all methods on keypoint propagation

Method	Mean PCK
DINO [1]	53.5
CRW [2]	42.6
MAST [3]	30.2
VFS [4]	44.6 222
RAFT [5]	45.6
Ours	62.3

[1] Caron et al. Emerging properties in self-supervised vision transformers. ICCV, 2021

[2] Wang et al. Learning correspondence from the cycle-consistency of time. CVPR, 2019.

[3] Lai et al. MAST: A memory-augmented self-supervised tracker, CVPR 2020.

[4] Xu & Wang. Rethinking self-supervised correspondence learning: A video frame-level similarity perspective. ICCV 2021.

[5] Teed & Deng. RAFT: Recurrent all-pairs field transforms for optical flow. ECCV 2020.

BADJA trajectory visualizations



DINO (baseline)



223

RAFT (baseline)



PIPs (ours)

BADJA trajectory visualizations



DINO (baseline)



224

RAFT (baseline)



PIPs (ours)

BADJA trajectory visualizations



DINO (baseline)



225

RAFT (baseline)



PIPs (ours)

BADJA trajectory visualizations



DINO (baseline)



226

RAFT (baseline)



PIPs (ours)

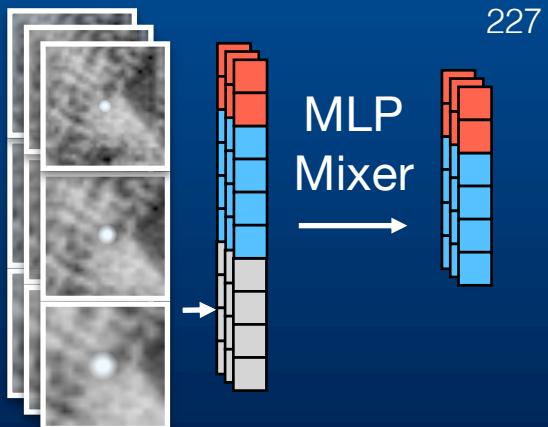
Particles in **small / boundary** regions are difficult for all models

PIPs

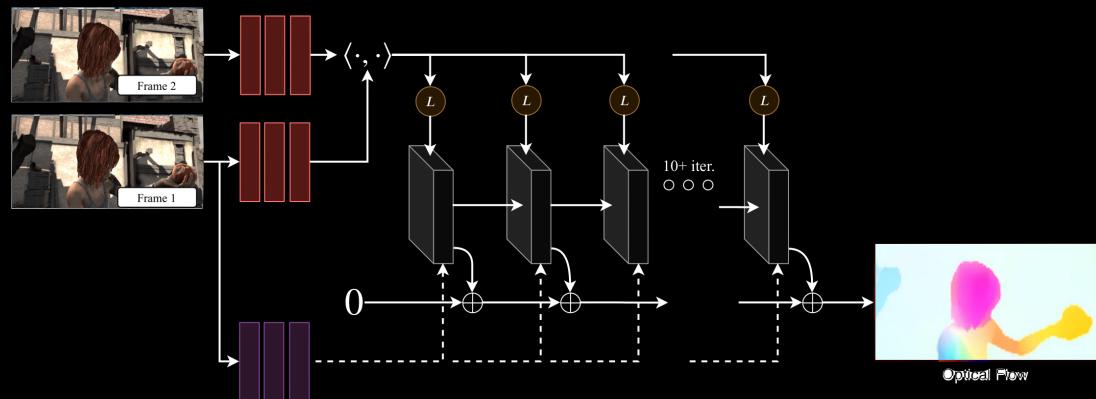
vs

RAFT

- **T-frame temporal window (e.g., T=8)**
- Does not share motion information between targets
- $H \times W \times (T-1)$ cost volume
- captures temporal context with iterated 12-block MLP-Mixer
- **runtime: 300ms per batch of particles per 8-frame clip (e.g., 128 particles)**

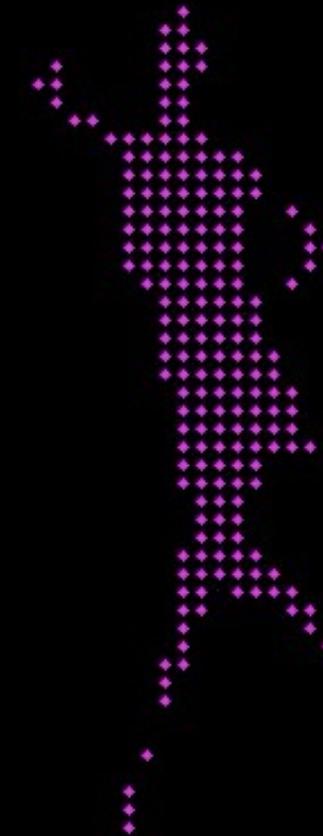


- **2-frame temporal window**
- Shares motion information locally
- $H \times W \times H \times W$ cost volume
- captures spatial context with iterated 2-layer convGRU
- **runtime: 250ms per frame (2000ms for 8)**

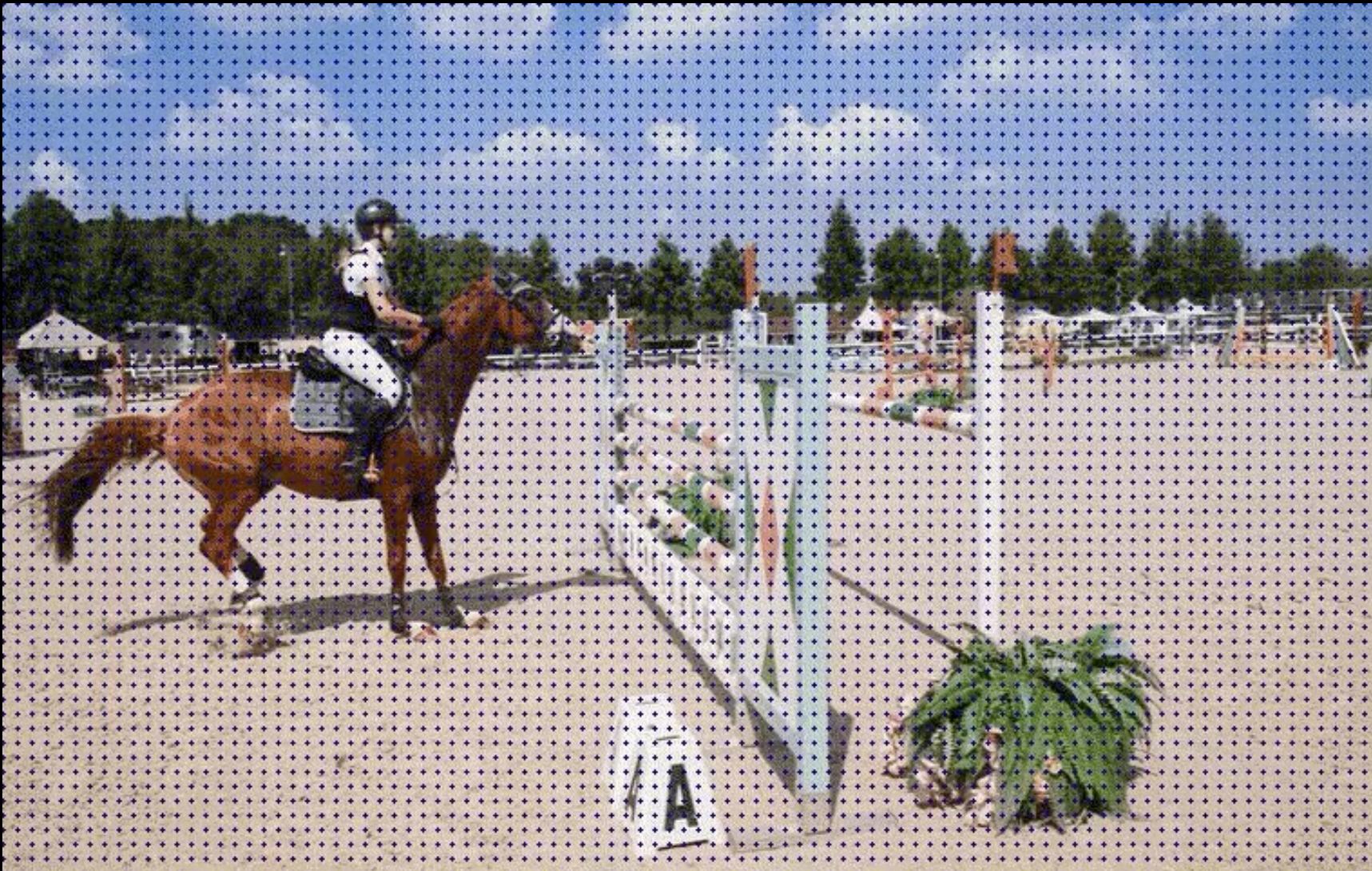


Teed & Deng. RAFT: Recurrent all-pairs field transforms for optical flow. ECCV 2020.

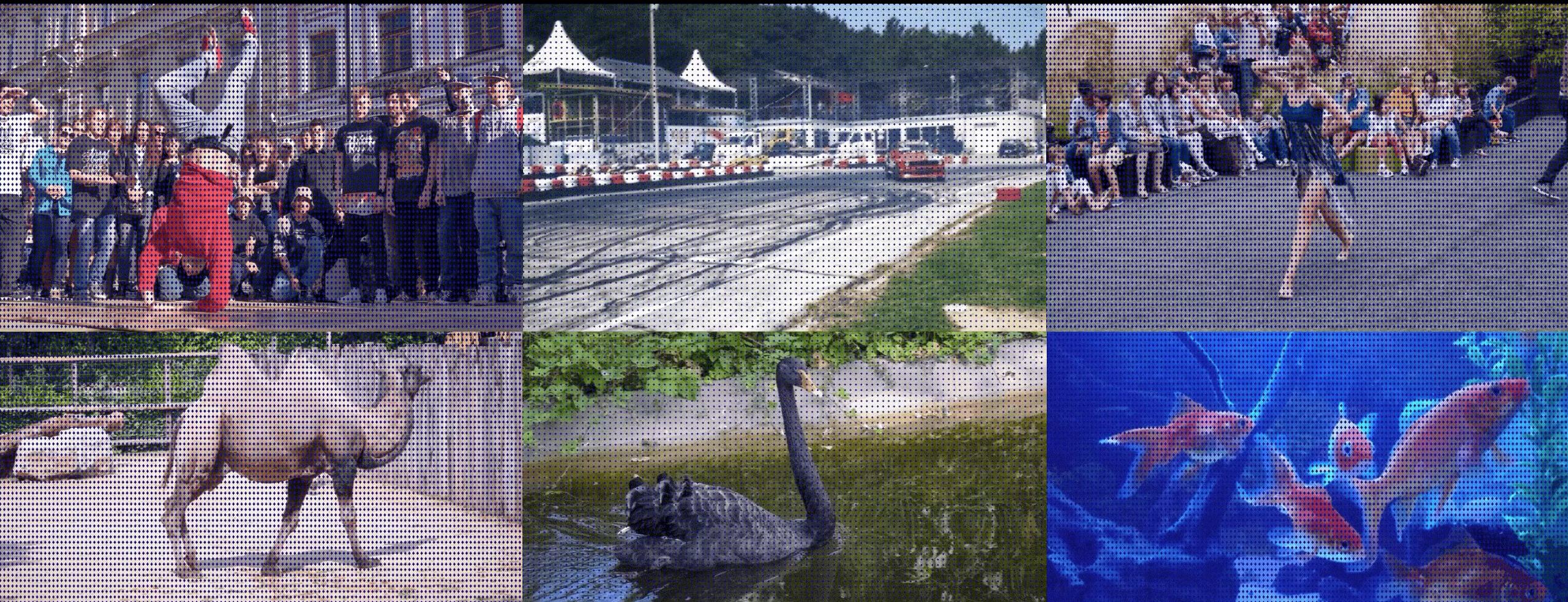
Generate trajectories for any irregular subset of pixels



Generalizes well to YouTube videos



Persistent Independent Particles (PIPs)



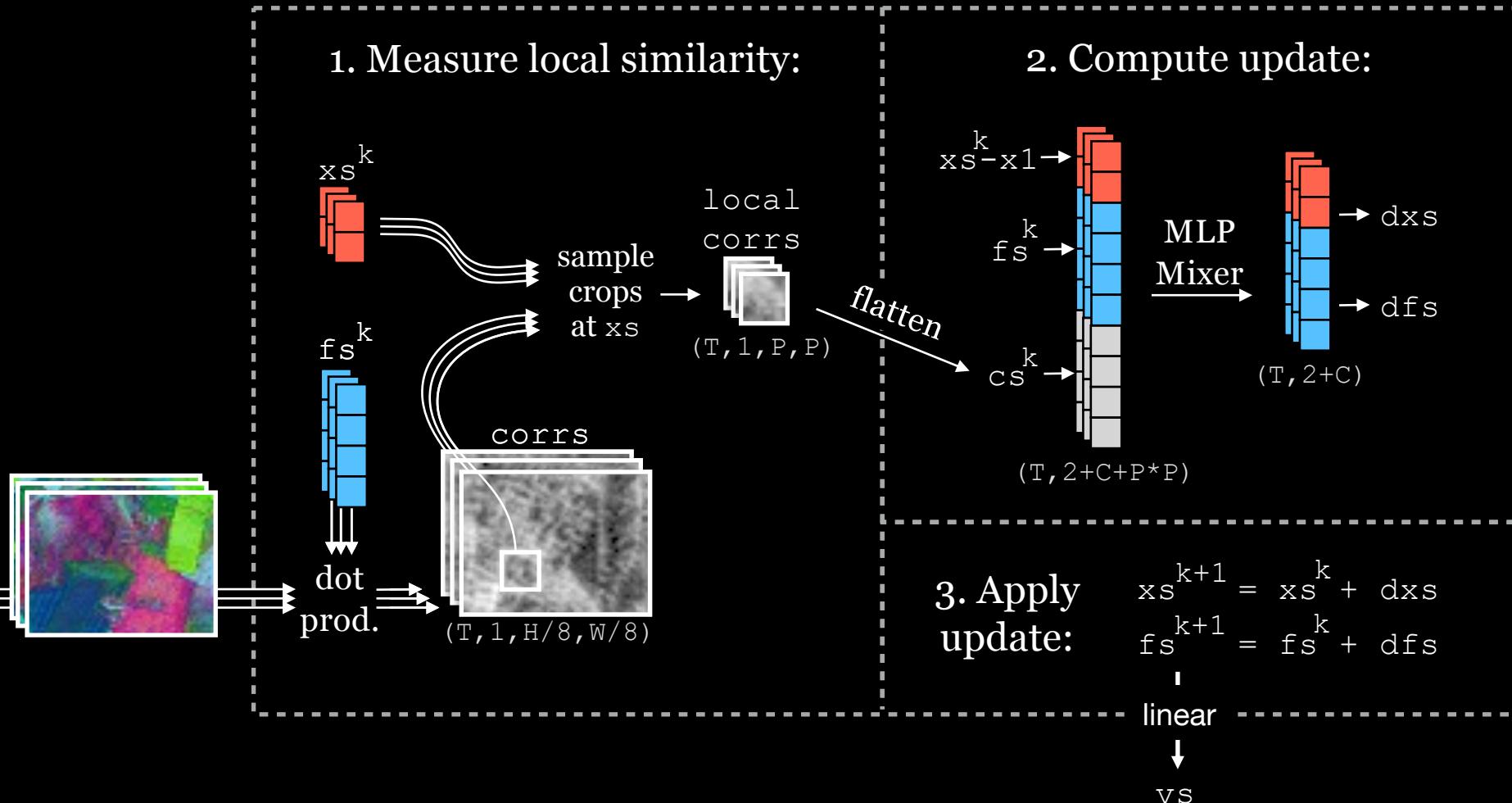
particle-video-revisited.github.io

github.com/aharley/pips

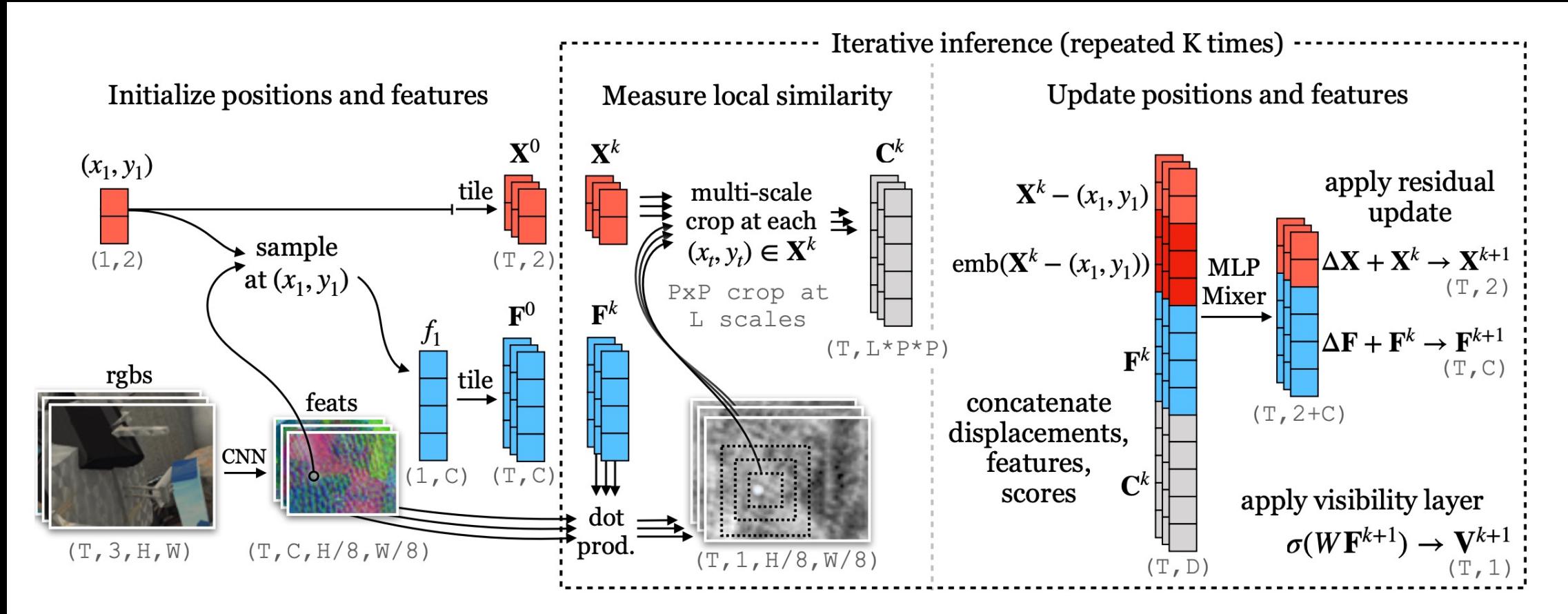
Limitation 1: No message-passing

Limitation 2: No recurrence

Iterate $\times K$



Architecture



Visibility-aware trajectory linking



$v=0.99$

$v=0.25$

$v=0.23$

$v=0.13$

$v=0.90$

init —————

236

↓
re-init —————→

Better on occluded points in FlyingThings++

Method	Visible pixels EPE	Occluded pixels EPE
DINO [1]	40.7	77.8
RAFT [2]	24.3	46.7
Ours	15.5 ²³⁷	36.7



[1] Caron et al. Emerging properties in self-supervised vision transformers. ICCV, 2021.

[2] Teed & Deng. RAFT: Recurrent all-pairs field transforms for optical flow. ECCV 2020.

[3] Mayer et al. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation, CVPR 2016.

Tied with RAFT on KITTI

Method	Vehicles	Pedestrians
DINO [1]	13.3	13.5
RAFT [2]	4.0	6.8
Ours	4.4	5.3



[1] Caron et al. Emerging properties in self-supervised vision transformers. ICCV, 2021.

[2] Teed & Deng. RAFT: Recurrent all-pairs field transforms for optical flow. ECCV 2020.

[3] Geiger et al. Vision meets robotics: The KITTI dataset. IJRR, 2013.

Better on head tracking in CroHD

Method	Visible heads	Occluded heads
DINO [1]	22.5	26.1
RAFT [2]	7.9	13.0
Ours	5.2	7.6



[1] Caron et al. Emerging properties in self-supervised vision transformers. ICCV, 2021.

[2] Teed & Deng. RAFT: Recurrent all-pairs field transforms for optical flow. ECCV 2020.

[3] Sundararaman et al. Tracking pedestrian heads in dense crowd. CVPR 2021.

Summary

- Dense optical flow can be computed efficiently and accurately with neural networks
- Learn a pixel-wise comparison function
- Train using diverse data with a curriculum
- Learn an iterative procedure
- Flow can be used for manipulation (cloth folding)
- We can learn flow without labels (self-supervised)