

# DINO-Tracker: Taming DINO for Self-Supervised Point Tracking in a Single Video

*Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel*

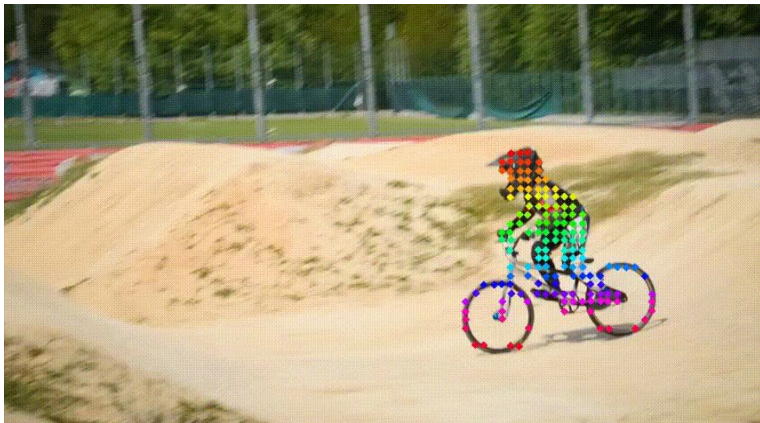
*Weizmann Institute of Science*

- **Problem/Objective**
  - tracking in video
- **Contribution/Key Idea**
  - tracking model

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Tracking



Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

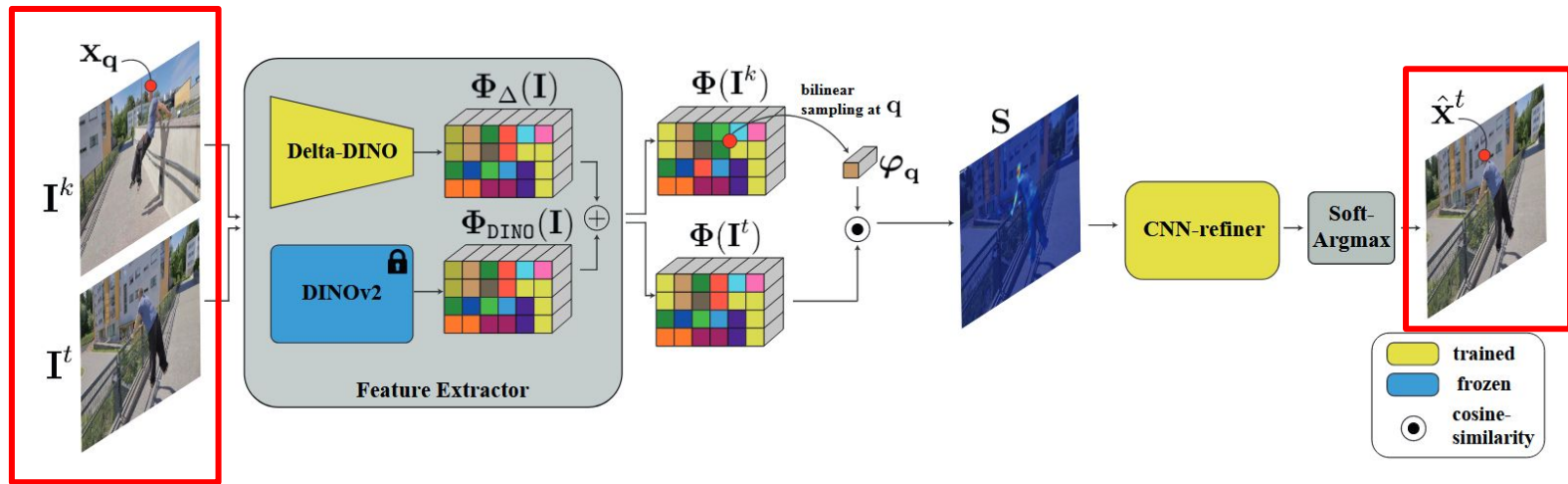
## - Tracking



Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Overview



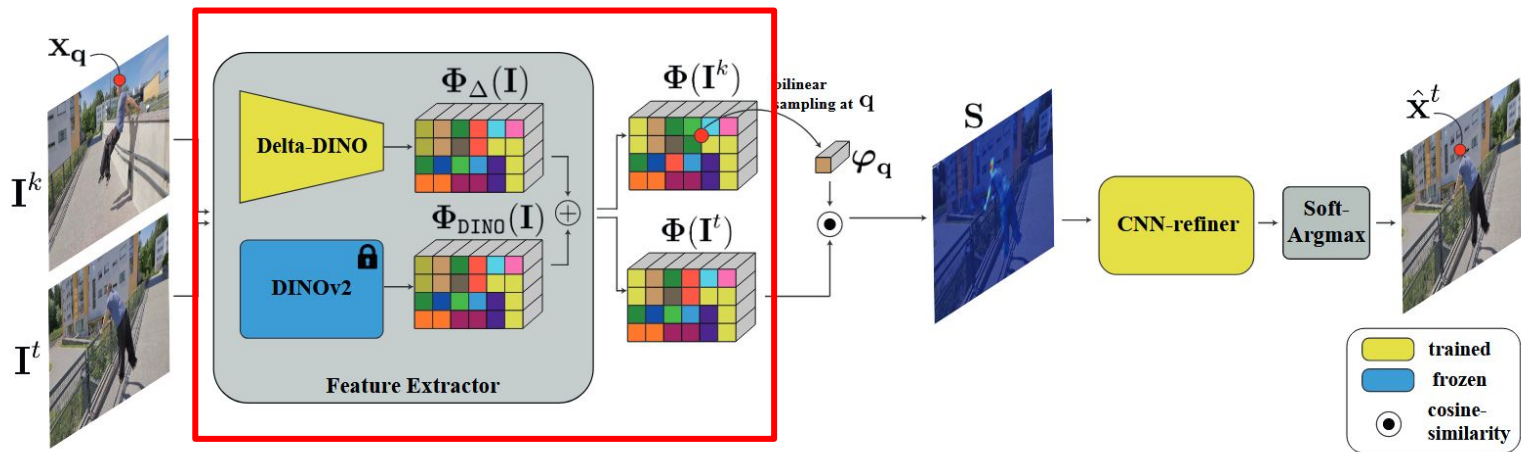
Input : query point  $\mathbf{x}_q$  , video  $\{\mathbf{I}^t\}_{t=1}^T$

Output : position estimates  $\{\hat{\mathbf{x}}^t\}_{t=1}^T$

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - DINO-Tracker



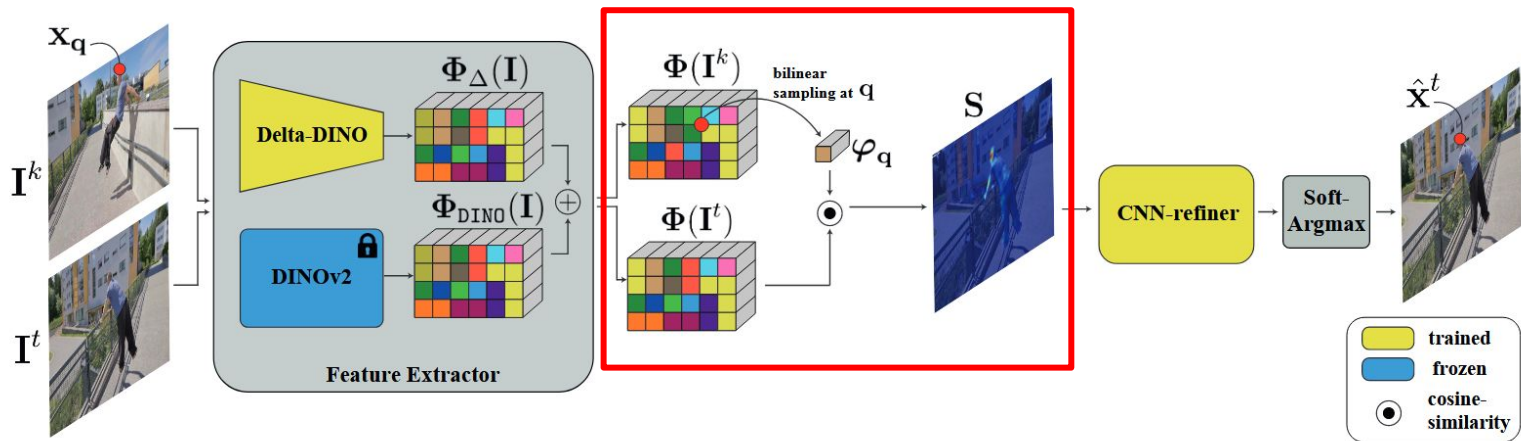
$$\Phi(I) = \Phi_{DINO}(I) + \Phi_{\Delta}(I) \quad (1)$$



Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - DINO-Tracker



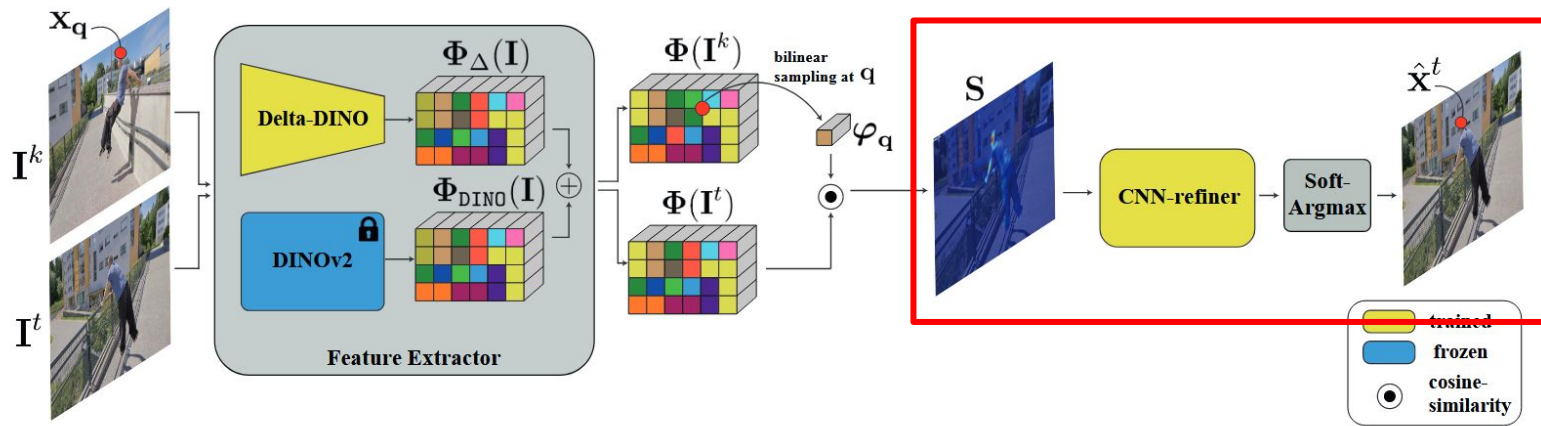
$$\varphi_q = \Phi(I^k)[q] \quad \Phi^t = \Phi(I^t)$$

$$S(p) = \text{cos-sim}(\varphi_q, \Phi^t(p)) \quad \text{where} \quad \text{cos-sim}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a}^T \cdot \mathbf{b}}{\|\mathbf{a}\|_2 \cdot \|\mathbf{b}\|_2}$$

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - DINO-Tracker



$$\hat{\mathbf{x}}^t = \frac{\sum_{\mathbf{p} \in \Omega} \mathbf{H}(\mathbf{p}) \cdot \mathbf{x}_{\mathbf{p}}}{\sum_{\mathbf{p} \in \Omega} \mathbf{H}(\mathbf{p})} \quad (2)$$

$$\Omega = \{\mathbf{p} : \|\mathbf{x}_{\mathbf{p}} - \mathbf{x}_{\mathbf{p}_{max}}\|_2 \leq R\}$$

$$\Pi(\mathbf{x}_{\mathbf{q}}, t) = \hat{\mathbf{x}}^t \quad \mathcal{T}_q = \{\hat{\mathbf{x}}^t : \hat{\mathbf{x}}^t = \Pi(\mathbf{x}_{\mathbf{q}}, t), t = 1 \dots T\}$$

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Self-Supervision

### • Optical flow

$$\Omega_{\text{flow}} = \{(\mathbf{x}^i, \mathbf{x}^j) \text{ cycle-consistent}\}$$

1) t 프레임에서 트랙 종료

$$\|\mathbf{x}^t - (\mathbf{x}^{t+1} + \mathbf{f}_{t+1 \rightarrow t}(\mathbf{x}^{t+1}))\| \geq \gamma_{\text{of}}$$

2) optical flow 와 일관되지 않는 correspondence  $\mathbf{x}^j$  제거

$$\|\mathbf{x}^j - \mathbf{x}^{i \rightarrow j}\|_2 \geq \gamma_{\text{of-lng}}$$

$$\|\mathbf{x}^i - (\mathbf{x}^{i \rightarrow j} + \mathbf{f}_{j \rightarrow i}(\mathbf{x}^{i \rightarrow j}))\|_2 \leq \gamma_{\text{of}}$$

$$\mathbf{x}^{i \rightarrow j} = \mathbf{x}^i + \mathbf{f}_{i \rightarrow j}(\mathbf{x}^i), \gamma_{\text{of-lng}} = 2\text{px}$$



Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

- **Self-Supervision**
- Feature correspondences

1) best-buddies

$$\Omega_{\text{dino-bb}} = \{(\mathbf{p}^i, \mathbf{p}^j) \text{ DINO bb}\}$$
$$NN(\varphi_{\text{DINO}}^i, \Phi_{\text{DINO}}(\mathbf{I}^j)) = \varphi_{\text{DINO}}^j \wedge NN(\varphi_{\text{DINO}}^j, \Phi_{\text{DINO}}(\mathbf{I}^i)) = \varphi_{\text{DINO}}^i \quad (3)$$

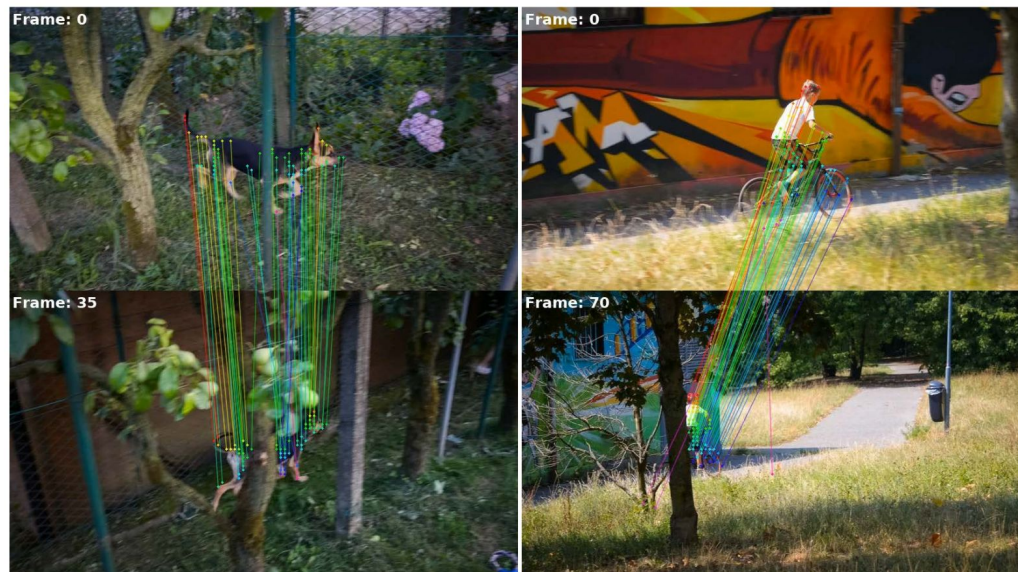
2) refined best buddies

$$\Omega_{\text{rfn-bb}} = \{(\mathbf{p}^i, \mathbf{p}^j) \text{ refined bb}\}$$

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

- **Self-Supervision**
- Feature correspondences



**Fig. 8:** *DINO best-buddies*. We visualize best-buddy pairs between distant frames. DINO best-buddies provide localized semantic correspondences, allowing the model to recover the object past repeating occlusions.

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Objective

$$\mathcal{L} = \mathcal{L}_{\text{flow}} + \lambda_1 \mathcal{L}_{\text{dino-bb}} + \lambda_2 \mathcal{L}_{\text{rfn-bb}} + \lambda_3 \mathcal{L}_{\text{rfn-cc}} + \lambda_4 \mathcal{L}_{\text{prior}} \quad (5)$$

$$\mathcal{L}_{\text{flow}} = \sum_{(\mathbf{x}^i, \mathbf{x}^j) \in \Omega_{\text{flow}}} L_H(\Pi(\mathbf{x}^i, j), \mathbf{x}^j) + L_H(\Pi(\mathbf{x}^j, i), \mathbf{x}^i)$$

$$\mathcal{L}_{\text{dino-bb}} = \frac{1}{|\Omega_{\text{dino-bb}}|} \sum_{(\varphi^i, \varphi^j) \in \Omega_{\text{dino-bb}}} \frac{1}{2} w_{\text{dino-bb}}^{ij} (l(\varphi^i, \varphi^j) + l(\varphi^j, \varphi^i))$$

$$l(\varphi^i, \varphi^j) = -\log \frac{\exp(\cos\text{-sim}(\varphi^i, \varphi^j)/\tau)}{\sum_{\mathbf{p}} \exp(\cos\text{-sim}(\varphi^i, \Phi^j(\mathbf{p}))/\tau)}$$

$$\mathcal{L}_{\text{rfn-bb}} = \frac{1}{|\Omega_{\text{rfn-bb}}|} \sum_{(\varphi^i, \varphi^j) \in \Omega_{\text{rfn-bb}}} \frac{1}{2} w_{\text{rfn-bb}}^{ij} (l(\varphi^i, \varphi^j) + l(\varphi^j, \varphi^i))$$

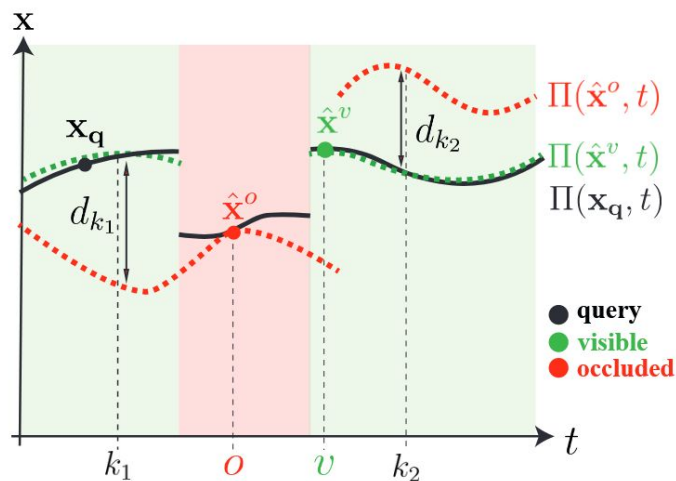
$$\mathcal{L}_{\text{rfn-cc}} = \sum_{(\mathbf{x}^i, \mathbf{x}^j) \in \Omega_{\text{rfn-cc}}} \frac{1}{2} w_{\text{rfn-cc}}^{ij} (L_H(\Pi(\mathbf{x}^i, j), \mathbf{x}^j) + L_H(\Pi(\mathbf{x}^j, i), \mathbf{x}^i)) \quad (4)$$

$$\mathcal{L}_{\text{prior}} = \frac{1}{H' \cdot W'} \cdot \sum_{\mathbf{p}} \underbrace{\left| 1 - \frac{\|\Phi(\mathbf{I})[\mathbf{p}]\|_2}{\|\Phi_{\text{DINO}}(\mathbf{I})[\mathbf{p}]\|_2} \right|}_{\mathcal{L}_{\text{norm}}} + \underbrace{|1 - \cos\text{-sim}(\Phi(\mathbf{I})[\mathbf{p}], \Phi_{\text{DINO}}(\mathbf{I})[\mathbf{p}])|}_{\mathcal{L}_{\text{angle}}}$$

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Occlusion prediction



**Fig. 3:** *Visibility via trajectory agreement.* To determine the visibility of  $\mathbf{x}_q$  at time  $t=o$ , we track  $\hat{\mathbf{x}}^o$  across time and check the agreement between  $\Pi(\hat{\mathbf{x}}^o, t)$  and  $\Pi(\mathbf{x}, t)$ . This is done by measuring  $d_{k_1}, d_{k_2}$  – displacements between the (black and red) tracks for anchor time steps  $k_1, k_2$ . Since these displacements are large, we classify  $\mathbf{x}_q$  as occluded for  $t=o$ . For  $t=v$ , the track  $\Pi(\hat{\mathbf{x}}^v, t)$  (green) agrees with  $\Pi(\mathbf{x}, t)$ , thus  $\mathbf{x}_q$  is classified as visible for  $t=v$ .

Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Experiments

Method	DAVIS-256			DAVIS-480			Kinetics-256			Kinetics-480			BADJA	
	$\delta_{avg}^x$	OA	AJ	$\delta_{avg}^x$	OA	AJ	$\delta_{avg}^x$	OA	AJ	$\delta_{avg}^x$	OA	AJ	$\delta^{seg}$	$\delta^{3px}$
RAFT [47]	56.7	—	—	66.7	—	—	50.4	—	—	60.5	—	—	45.0	5.8
DINOv2 [38]	61.4	—	—	64.7	—	—	60.3	—	—	61.0	—	—	62.8	8.4
TAP-Net* [12]	53.4	81.4	38.4	66.4	79.0	46.0	61.7	86.6	48.5	67.1	81.5	47.7	45.4	9.6
PIPs++* [63]	71.5	—	—	73.6	—	—	68.2	—	—	70.8	—	—	59.0	9.8
TAPIR* [13]	74.7	<b>89.4</b>	<u>62.8</u>	77.3	<b>89.5</b>	<b>65.7</b>	69.5	<u>89.1</u>	57.3	69.8	86.7	57.5	<u>68.7</u>	10.5
Co-Tracker* [26]	<b>79.2</b>	<u>89.3</u>	<b>65.1</b>	<u>79.4</u>	<b>89.5</b>	<u>65.6</u>	<u>72.9</u>	88.9	<b>59.9</b>	<u>72.8</u>	<u>88.9</u>	<u>59.8</u>	64.0	<u>11.2</u>
Omnimotion <sup>†</sup> [51]	67.5	85.3	51.7	74.1	84.5	58.4	69.2	<b>89.2</b>	55.0	—	—	—	45.2	6.9
Ours <sup>†</sup>	<u>78.2</u>	87.5	62.3	<b>80.4</b>	<u>88.1</u>	64.6	<b>73.3</b>	88.5	<u>59.7</u>	<b>74.3</b>	<b>89.2</b>	<b>60.9</b>	<b>72.4</b>	<b>14.3</b>

 $\delta_{avg}^x$  : position accuracy

OA : occlusion accuracy

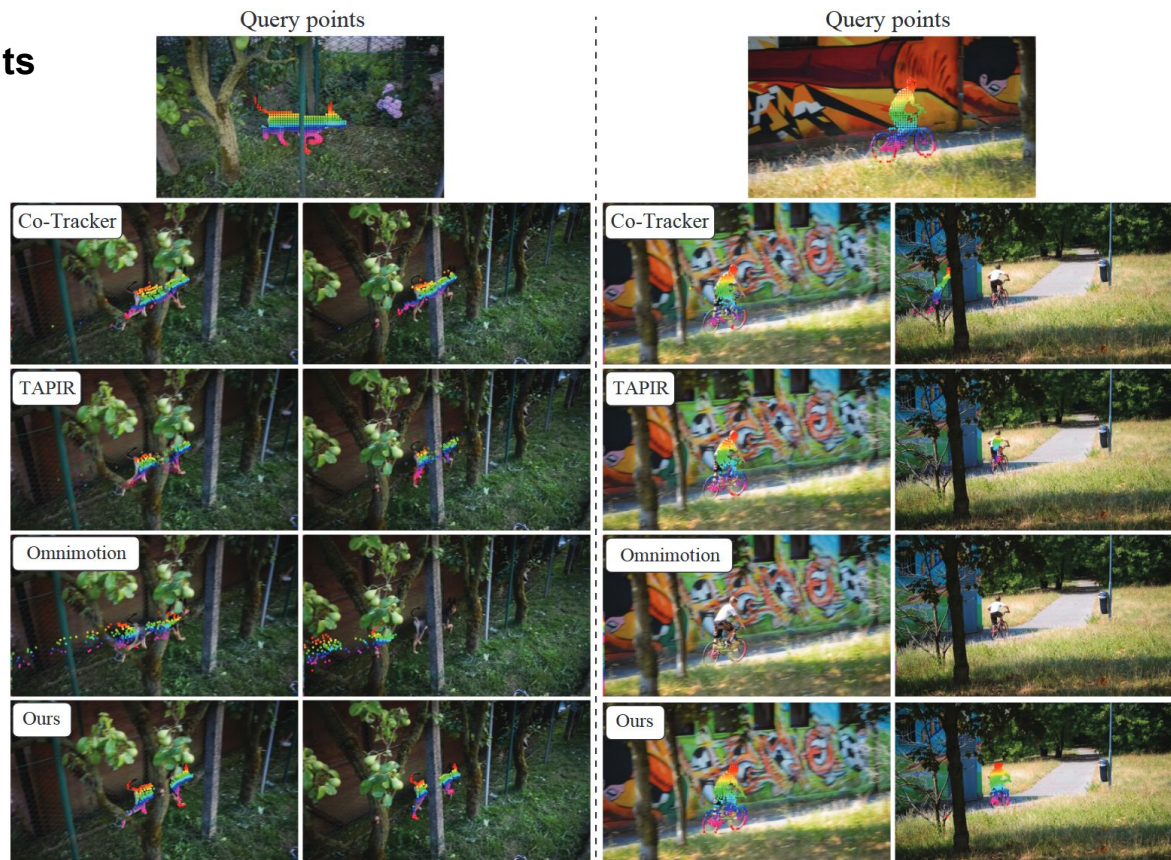
AJ : average jaccard : position + occlusion accuracy 둘다 고려



Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Experiments



Narek Tumanyan\*, Assaf Singer\*, Shai Bagon, Tali Dekel

Weizmann Institute of Science

## - Experiments

