

# BİL-722 ADVANCED TOPICS IN COMPUTER VISION

Çağdaş Baş, N10266943

Paper: Robust Object Tracking with Online Multi-lifespan Dictionary Learning

Authors: Junliang Xing, Jin Gao, Bing Li, Weiming Hu and Shuicheng Yan

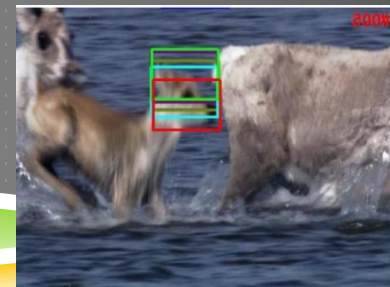
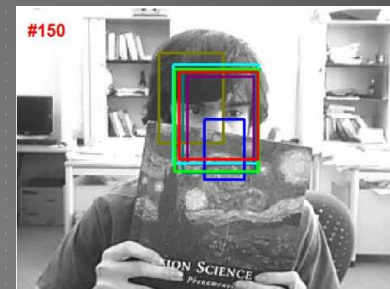
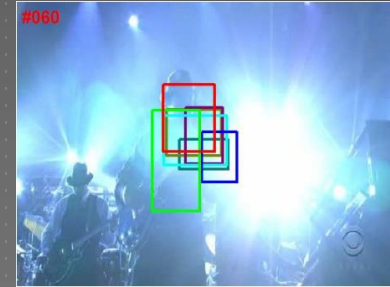
# TRACKING



Images from: SUN Dataset

# WHY IS TRACKING A DIFFICULT PROBLEM?

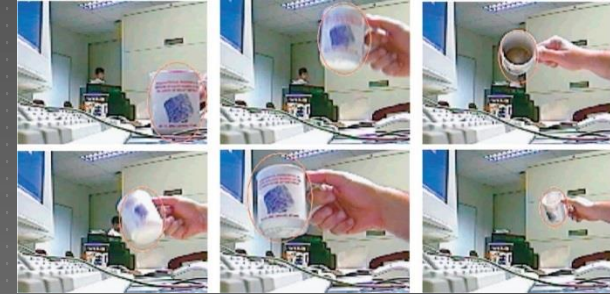
- ❖ Image noise and background clutter
- ❖ Illumination changes
- ❖ Clutter
- ❖ Motion



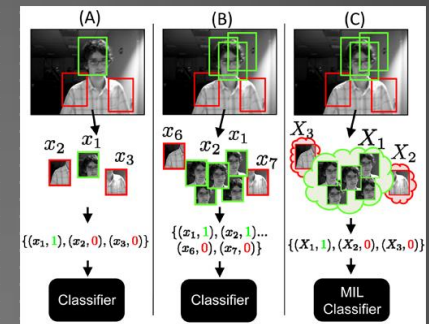
# RELATED WORK

❖ There are three types of methods in general:

1. Generative Methods  
Eigentrapper, meanshift tracker etc.
2. Discriminative Methods  
Ensemble tracker, MIL tracker etc.
3. Sparse Learning Methods  
This method



MeanShift Tracker [1]



MIL Tracker [2]

# DIFFERENT STAGES OF TRACKING

1. Designing good templates
2. Solving optimization problem efficiently
3. Updating the object template.

Many of the papers focus mainly on these two parts.

Present methods use fixed templates or incremental template update.

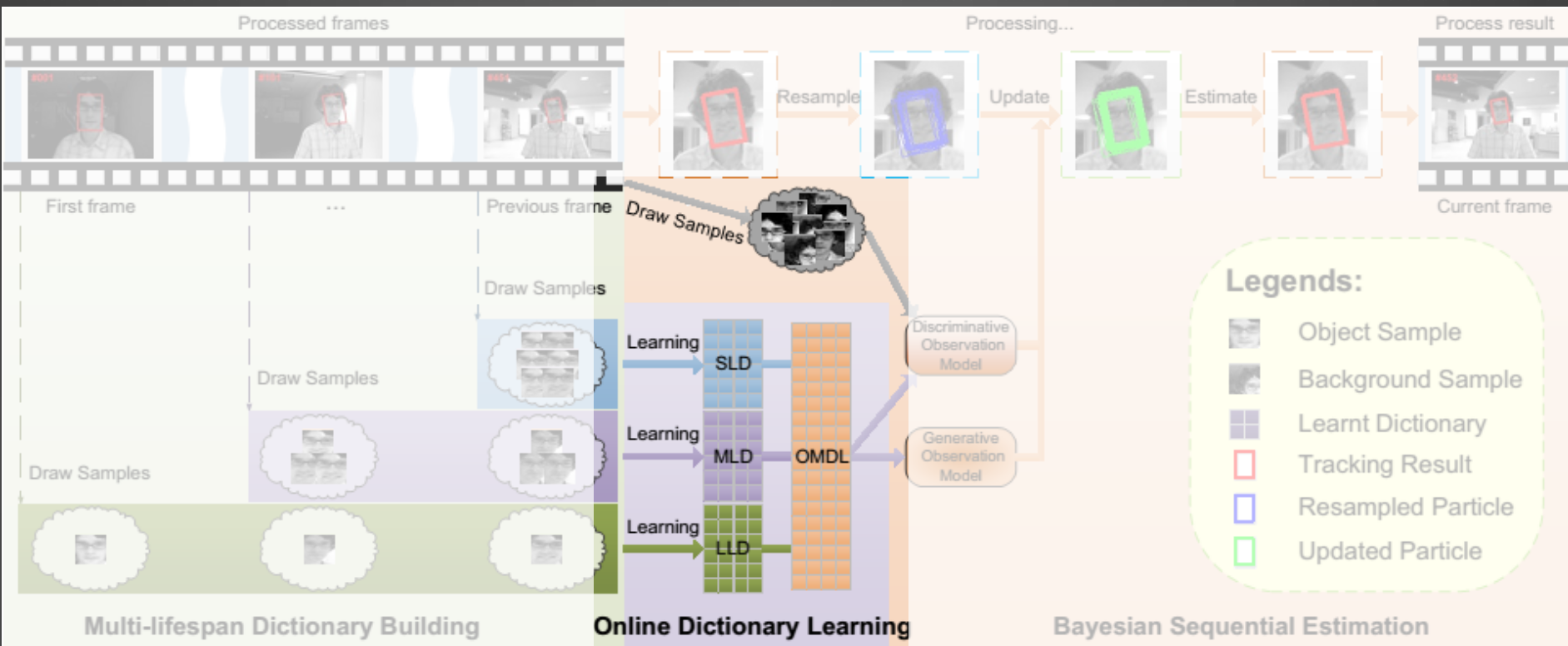
# WHAT IS SPARSE LEARNING?

- ❖ Represents an instance with a minimum set of dictionary elements.

$$\min_c \|Tc - y\|_2^2 + \lambda|c|_1$$

- ❖ Find the sparse representation  $c$  of  $y$  in dictionary  $T$

# OVERALL APPROACH



# TRACKING AS ONLINE DICTIONARY LEARNING

- ❖ Extract candidate regions  $\mathcal{Y}$
- ❖ Optimize a new object template (dictionary) by minimizing:

$$D^* = \underset{c}{\operatorname{argmin}} \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} l(y, D)$$

- ❖ Template set is not predefined but learned in time.



# ONLINE LEARNING ALGORITHM

**Algorithm 1** Online dictionary learning for template update

**Input:** frame data  $\mathbf{I}_t$ , tracking results  $\mathbf{x}_t$ , learned dictionary  $\mathbf{D}_{t-1}$ ,  $\mathbf{C}_{t-1}$ ,  $\mathbf{Y}_{t-1}$  in the previous frame,  $\lambda$  (regularization parameter),  $M$  (sample drawing number).

**Output:** learned dictionary  $\mathbf{D}_t$  in the current frame.

1: **Initialization:**  $\mathbf{D}_t \leftarrow \mathbf{D}_{t-1}, \mathbf{C}_t \leftarrow \mathbf{C}_{t-1}, \mathbf{Y}_t \leftarrow \mathbf{Y}_{t-1}$ .

2: **for**  $i = 1 \rightarrow M$  **do**

3:   **Step 1:** fix  $\mathbf{D}_t$  to find the best coefficients,

$$\mathbf{c}_t^{(i)} = \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y}_t^{(i)} - \mathbf{D}_t \mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_1.$$

4:   **Step 2:** fix  $\{\mathbf{c}_t^{(i)}\}$  to update the dictionary,

$$\mathbf{C}_t \leftarrow \mathbf{C}_t - \frac{\mathbf{c}_t \mathbf{c}_t^{(i)\top} \mathbf{c}_t^{(i)\top} \mathbf{C}_t}{1 + \mathbf{c}_t^{(i)\top} \mathbf{C}_t \mathbf{c}_t^{(i)}}, \mathbf{Y}_t \leftarrow \mathbf{Y}_t + \mathbf{y}_t^{(i)} \mathbf{c}_t^{(i)\top},$$

$$\mathbf{D}_t = \operatorname{argmin}_{\mathbf{D} \in \mathcal{D}} \sum_{j=1}^i \frac{1}{2} \|\mathbf{y}_t^{(j)} - \mathbf{D} \mathbf{c}_t^{(j)}\|_2^2 + \lambda \|\mathbf{c}_t^{(j)}\|_1,$$

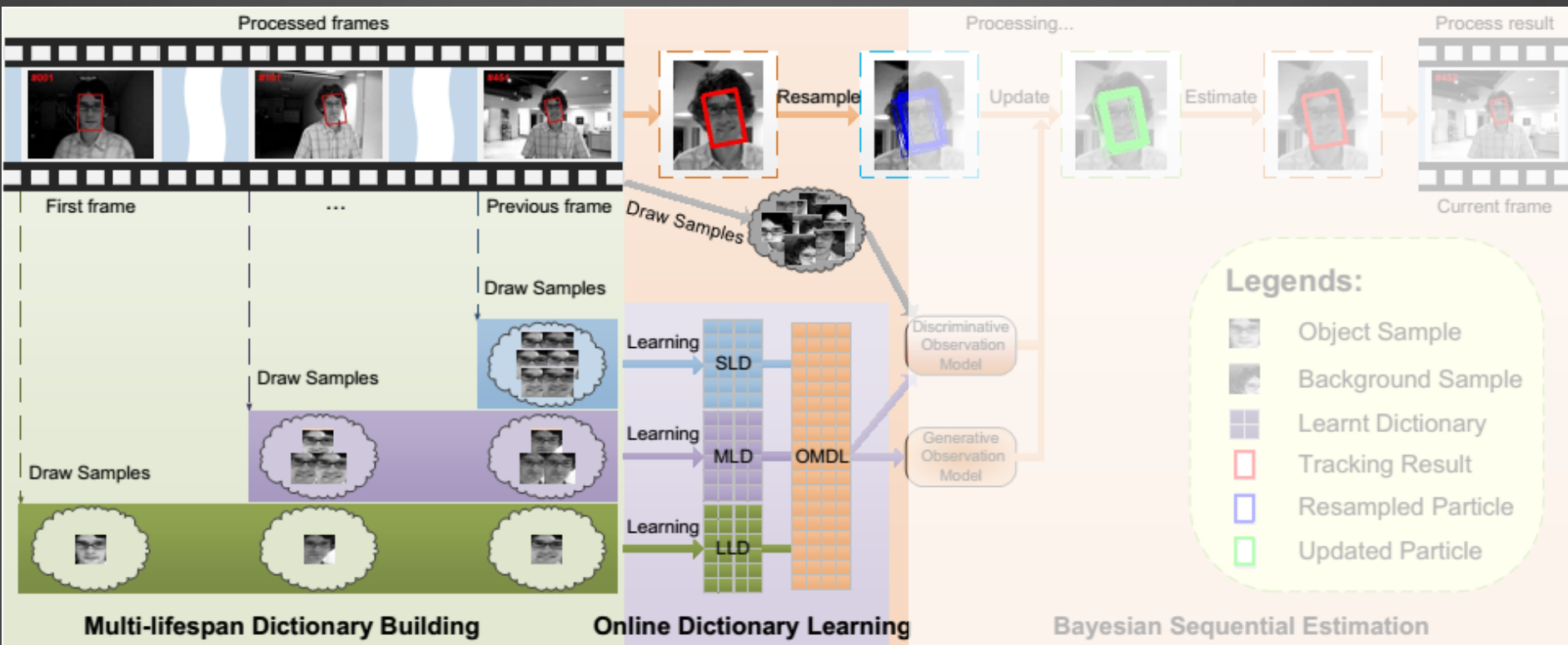
$$= \left( \sum_{j=1}^i \mathbf{c}_t^{(j)} \mathbf{c}_t^{(j)\top} \right)^{-1} \left( \sum_{j=1}^i \mathbf{y}_t^{(j)} \mathbf{c}_t^{(j)\top} \right),$$

$$= \mathbf{C}_t \mathbf{Y}_t.$$

5: **end for**

6: Save dictionary  $\mathbf{D}_t$ , intermediate variable  $\mathbf{C}_t$  and  $\mathbf{Y}_t$ .

# MULTI-LIFESPAN DICTIONARY LEARNING



# MULTI-LIFESPAN DICTIONARY LEARNING

## ❖ Sample starting frame changes the lifespan

1. SLD, Short Life Span Dictionary is learned the samples extracted from only previous frame. (Starting frame:  $t-1$ , Ending frame:  $t$ )
  - ❖ Learned for best adaptation
2. LLD, Short Life Span Dictionary is learned the samples extracted from all the frames before current (Starting frame:  $1$ , Ending frame:  $t$ )
  - ❖ Learned for robustness
3. MLD, Short Life Span Dictionary is learned to balance short life span and long life span (Starting frame:  $t/2$ , Ending frame:  $t$ )
  - ❖ Balances trade-off between short term adaptive model and long term robust model.

## ❖ Final Template is:

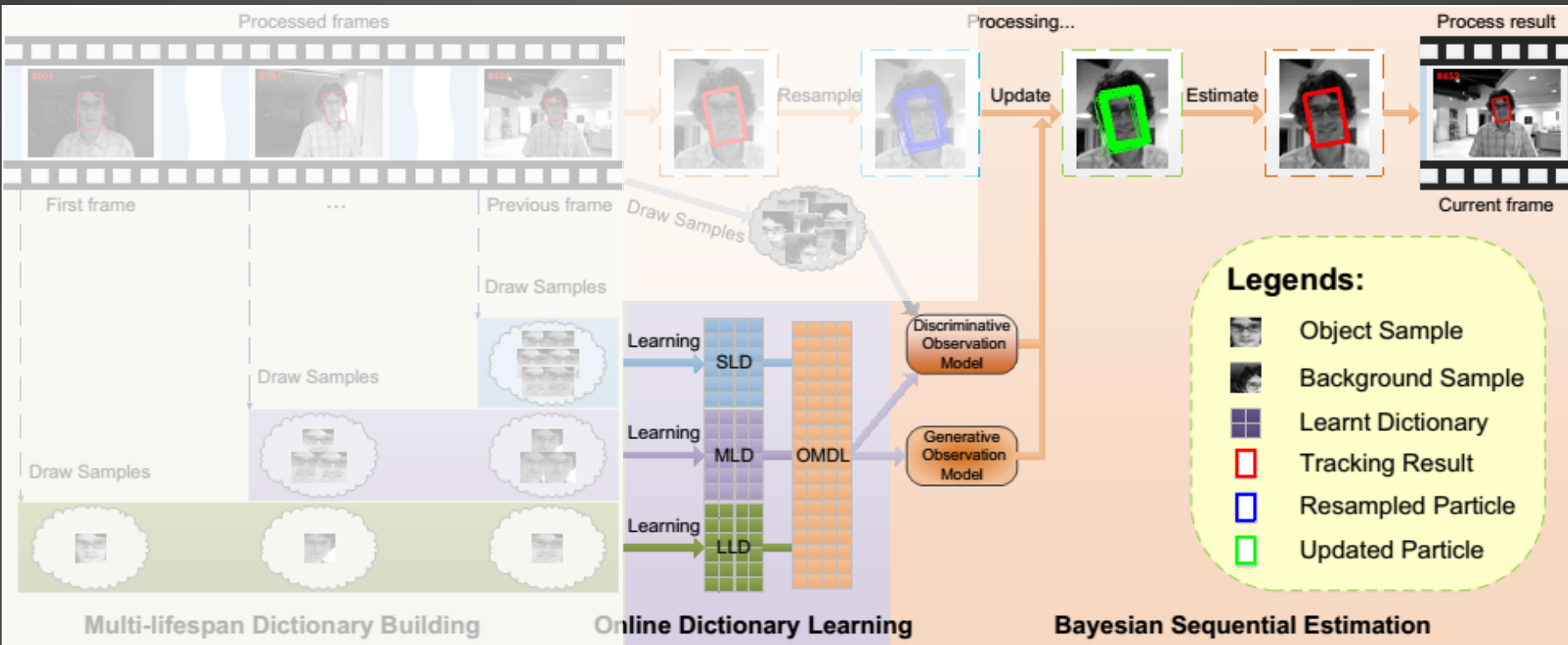
$$D^* = \{D^S, D^M, D^L\}$$

# MULTI-LIFESPAN DICTIONARY LEARNING

- ❖ Learned examples of different life spanned dictionaries.



# BAYESIAN SEQUENTIAL ESTIMATION



# PARTICLE FILTER

- ❖ Solving maximum posterior problem:

$$\hat{x}_t = \underset{x_t}{\operatorname{argmax}} p(x_t | y_{1:t})$$

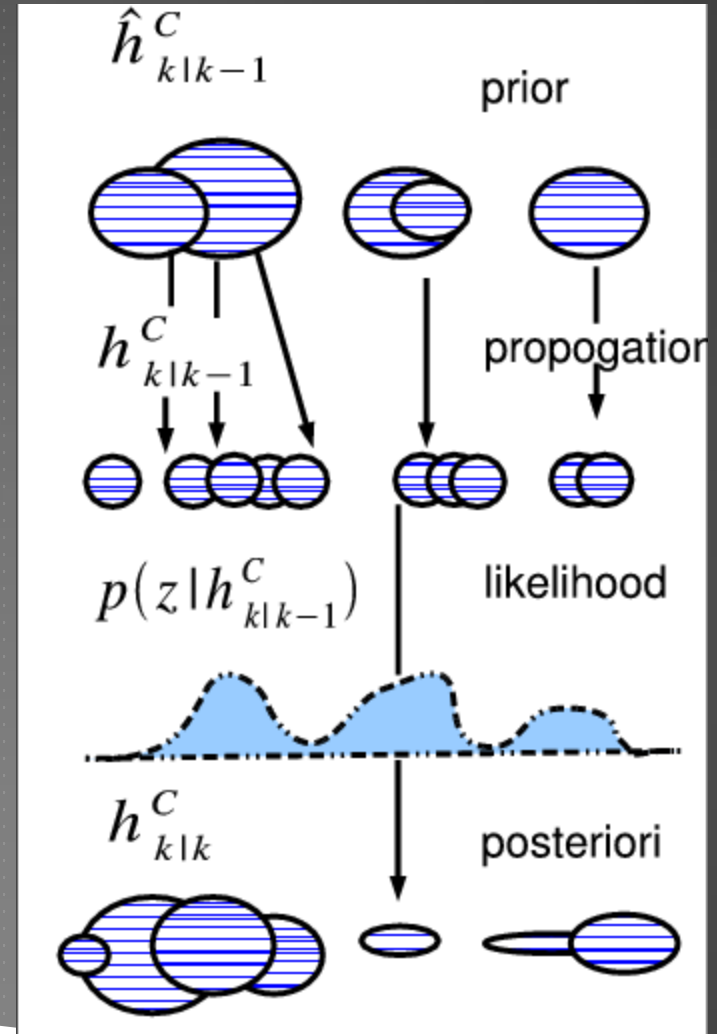
means tracking.

- ❖ Use learned OMDL as observation model in Particle Filter:

$$p(y_t | x_t) \propto \underbrace{g(y_t | x_t)}_{\text{Generative: Fix dictionary and optimize sparsity}} \underbrace{d(y_t | x_t)}_{\text{Discriminative: Extract negative samples and optimize with labels}}$$

Generative: Fix dictionary and optimize sparsity

Discriminative: Extract negative samples and optimize with labels

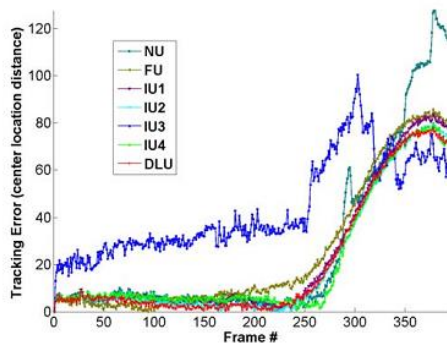


# EXPERIMENTS

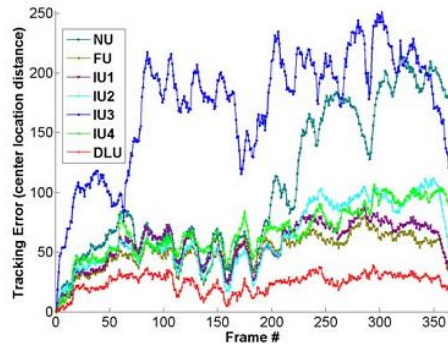
- ❖ Results compared with two different metric:
  - ❖ Center location distance
  - ❖ Overlap Ratio

# EXPERIMENTS

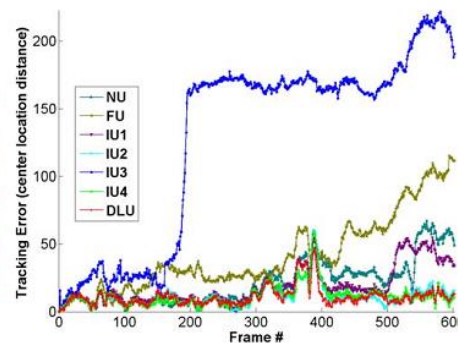
❖ The template update method is evaluated firstly:



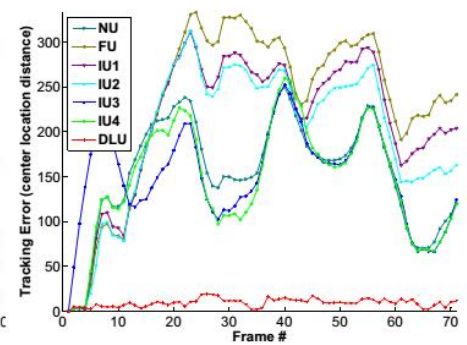
(a) Tracking Error on *car11*



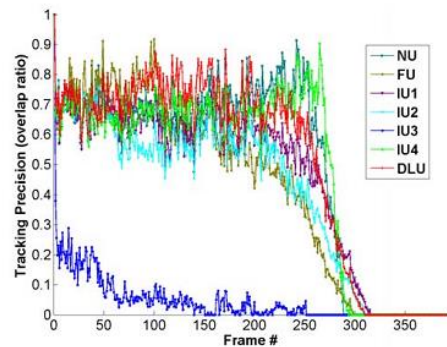
(b) Tracking Error on *shaking*



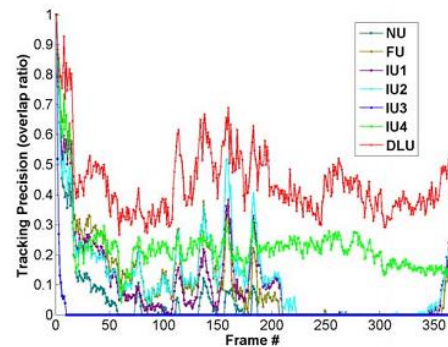
(c) Tracking Error on *faceocc2*



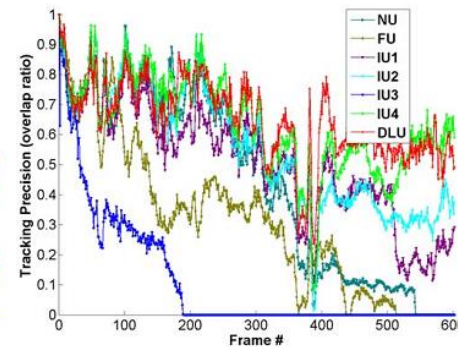
(d) Tracking Error on *animal*



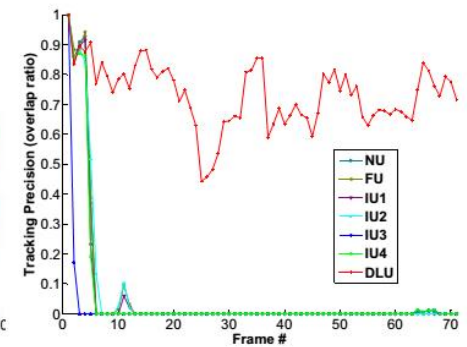
(e) Tracking Precision on *car11*



(f) Tracking Precision on *shaking*



(g) Tracking Precision on *faceocc2*



(h) Tracking Precision on *animal*



# EXPERIMENTS

## ❖ Overall tracking error and precision:

Table 2. Average tracking errors (in pixels). The best and second best results are respectively shown in **red** and **blue** colors.

Sequence	Frag	IVT	MIL	VTD	$\ell_1$	MTT	Ours
<i>sylv</i>	0.245	0.875	<b>0.156</b>	0.220	0.961	0.260	<b>0.139</b>
<i>bike</i>	2.109	0.075	0.083	0.086	0.082	<b>0.070</b>	<b>0.054</b>
<i>car11</i>	1.436	<b>0.062</b>	0.848	<b>0.065</b>	0.378	0.403	0.161
<i>david</i>	0.946	<b>0.057</b>	0.194	0.351	0.210	1.103	<b>0.110</b>
<i>woman</i>	1.302	1.590	1.351	<b>1.126</b>	1.305	2.281	<b>0.135</b>
<i>animal</i>	0.934	0.101	0.182	<b>0.056</b>	0.059	<b>0.047</b>	<b>0.047</b>
<i>coke11</i>	1.247	0.894	0.381	0.759	0.954	<b>0.338</b>	<b>0.178</b>
<i>shaking</i>	0.704	1.005	<b>0.222</b>	0.279	1.286	0.336	<b>0.161</b>
<i>jumping</i>	0.169	<b>0.094</b>	0.245	1.121	1.417	0.666	<b>0.081</b>
<i>faceocc2</i>	0.137	<b>0.101</b>	0.252	0.117	0.149	<b>0.097</b>	0.113
Average	0.923	0.485	<b>0.391</b>	0.418	0.680	0.560	<b>0.118</b>

Table 3. Average tracking precision. The best and second best results are respectively shown in **red** and **blue** colors.

Sequence	Frag	IVT	MIL	VTD	$\ell_1$	MTT	Ours
<i>sylv</i>	0.617	0.450	0.751	<b>0.810</b>	0.323	0.770	<b>0.833</b>
<i>bike</i>	0.136	<b>0.983</b>	0.917	<b>1.000</b>	0.908	<b>1.000</b>	<b>1.000</b>
<i>car11</i>	0.097	<b>1.000</b>	0.102	<b>0.972</b>	0.682	0.687	0.781
<i>david</i>	0.089	<b>0.905</b>	0.537	0.615	0.435	0.320	<b>0.779</b>
<i>woman</i>	0.256	0.204	0.209	<b>0.309</b>	0.215	0.198	<b>0.440</b>
<i>animal</i>	0.099	0.887	0.747	<b>0.972</b>	<b>0.972</b>	<b>1.000</b>	<b>1.000</b>
<i>coke11</i>	0.051	0.119	0.271	0.068	0.085	<b>0.559</b>	<b>0.678</b>
<i>shaking</i>	0.222	0.025	0.414	<b>0.784</b>	0.011	0.099	<b>0.578</b>
<i>jumping</i>	0.690	<b>0.959</b>	0.233	0.230	0.118	0.198	<b>0.984</b>
<i>faceocc2</i>	0.767	0.772	0.537	0.743	0.419	<b>0.929</b>	<b>0.826</b>
Average	0.302	0.630	0.472	<b>0.650</b>	0.417	0.576	<b>0.790</b>

# EXPERIMENTS

## ❖ Speed Analysis:

Algorithm	[19]	[20]	[5]	[27]	[28]	[13]	Ours
Speed	0.01fps	0.05fps	1fps	2fps	2.5fps	2.5fps	2.5fps

# VISUAL RESULTS



(e) *woman, occlusions and viewpoint changes*



(i) *jumping, fast motions and blurs*



(h) *shaking, dynamic illumination changes and scale variations*

# THANKS

1. D. Comaniciu and P. Meer. Kernel-based object tracking. TPAMI, 25(5):564–77, 2003.
2. B. Babenko, M. Yang, and S. Belongie. Robust object tracking with online multiple instance learning. TPAMI, 33(8):1619–32, 2011.