

Modular Decomposition and Analysis of Registration based Trackers

Abhineet Singh, Ankush Roy, Xi Zhang
and Martin Jagersand

Reading

Abhineet Singh and Martin Jagersand, "**Modular Tracking Framework: A Fast Library for High Precision Tracking**",
IEEE/RSJ International Conference on Intelligent Robots and
Systems (IROS), September 2017 [\[pdf\]](#) [\[video\]](#) [\[code\]](#)

MTF is available at: <http://webdocs.cs.ualberta.ca/~vis/mf/>
along with all datasets and papers



UNIVERSITY OF
ALBERTA

Registration based Tracking

- Find the optimal warp or geometric transformation that registers each image in a sequence with the template

$$\mathbf{p}_t = \underset{\mathbf{p}}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$

$$\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N], \mathbf{x}_k = [x_k, y_k]^T \in \mathbb{R}^2$$

$$\mathbf{I}(\mathbf{x}) = [I(x_1, y_1), I(x_2, y_2), \dots, I(x_N, y_N)]^T \in \mathbb{R}^N, I(x, y) : \mathbb{R}^2 \mapsto \mathbb{R}$$

$$\mathbf{p} = [p_1, p_2, \dots, p_S], S : \text{DOF of image motion}$$

$$\mathbf{w} : \mathbb{R}^2 \times \mathbb{R}^S \mapsto \mathbb{R}^2$$

$$f : \mathbb{R}^N \times \mathbb{R}^N \mapsto \mathbb{R}$$

Tracking Lucas-Kanade algorithm

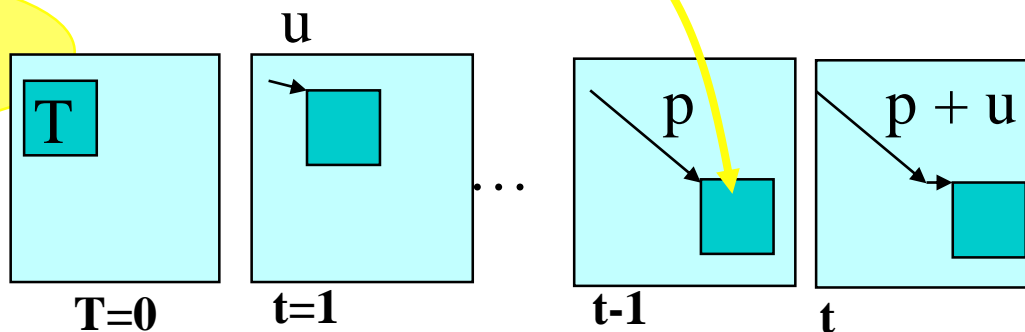
- Create tracking loop, iterate for each new image

Init $p=0$, Template T

For $t = 1 \dots$

1. Receive $I(t+1)$
2. Compute $dIm = I(t+1, x+p) - T$
3. Solve $-Im_t = M u$
 - Use $u = M \backslash Im_t$
4. Update $p = p + u$

Template sourced from pixel window shifted by the state vector p

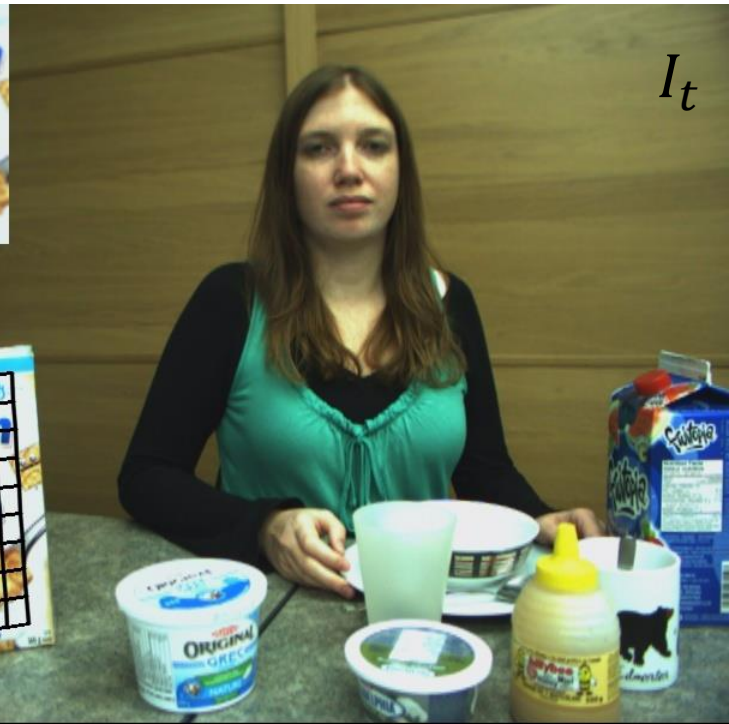


Registration based Tracking

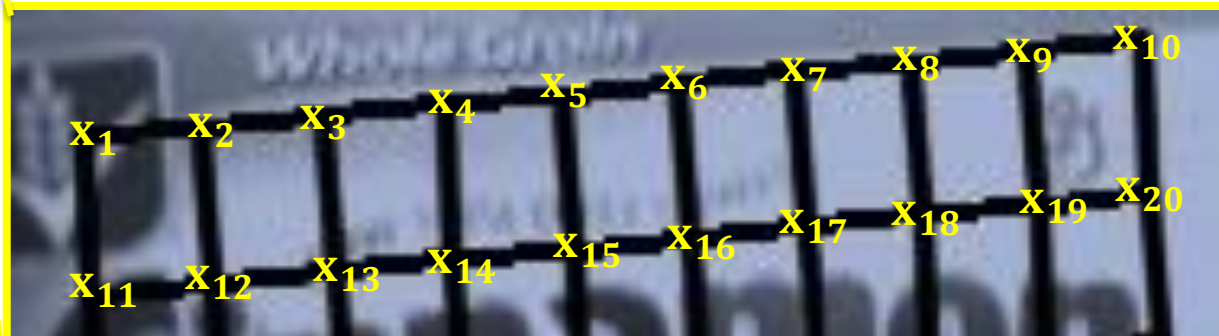
$$\mathbf{p}_t = \underset{\mathbf{p}}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$



I_0



I_t



Motivation

- Learning/detection based trackers are not suitable for tasks requiring **fast** and **high precision** tracking
 - Visual Servoing
 - Virtual reality
 - SLAM

Registration (8DOF)



Learning (3DOF)



MTF Usage Example – Multi Target
Tracking

UAV Trajectory Estimation

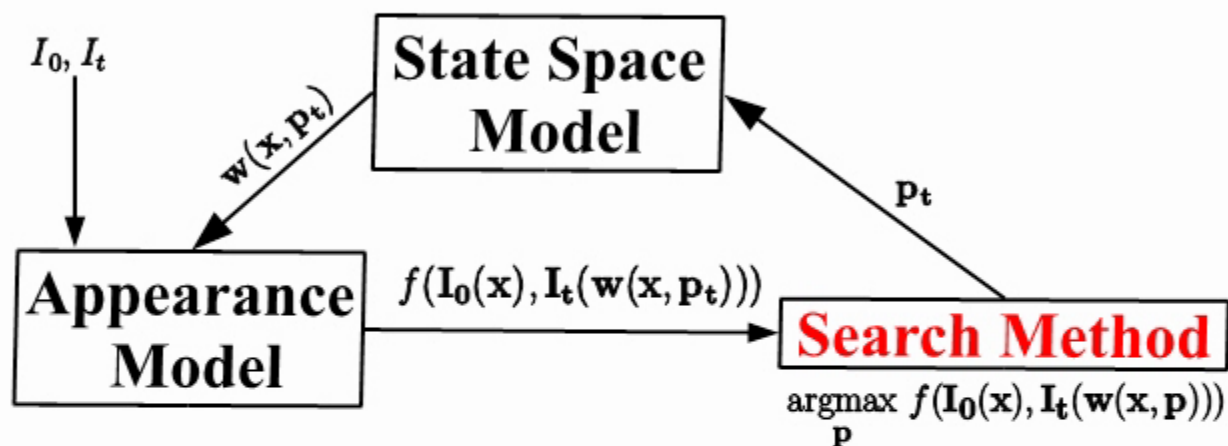
Online Image Mosaicing

Motivation

- Progress in registration based tracking has become **fragmented** since Lucas Kanade^[Lucas81]
 - myriad of contributions that are not well connected
- An intuitive way exists to **relate** these by decomposing the tracking task into three modules
 - most contributions are confined to only one or two of these modules
- Modular Tracking Framework (**MTF**)^[Singh16] to easily plug in new methods

Modular Decomposition

$$\mathbf{p}_t = \underset{\mathbf{p}}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$



- Appearance Model (**AM**)
 - Measures the **similarity** between a warped patch and the template
- State Space Model (**SSM**)
 - Defines the possible ways to **warp** the object patch
- Search Method (**SM**)
 - **Finds** the warp that maximizes the similarity measure

State Space Model

$$\mathbf{p}_t = \underset{\mathbf{p}}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$

- A warping function or geometric transformation that represents the set of allowable image motions of the object
 - embodies any constraints placed on the warp parameter space
 - search efficiency
 - alignment precision
 - includes
 - degrees of freedom (DOF) of allowed motion
 - actual parameterization of the warping function

Registration: from trans \mathbf{u} to warp $\mathbf{w}(\mathbf{x}, \mathbf{p})$

Find parameters of a warping function such that:

$$\mathcal{I}(\mathbf{w}(\bar{\mathbf{x}}; \mathbf{p}_i)) = \mathcal{I}_T(\mathbf{p}_i)$$

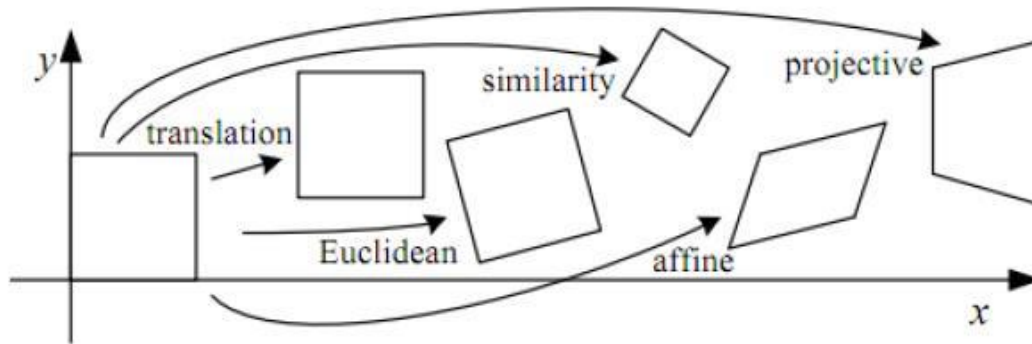
for all template points \mathbf{p}_i



$\mathbf{w}(\bar{\mathbf{x}}; \mathbf{p}_i)$



State Space Model – Examples



Translation : $S = 2$

$$- \mathbf{w}(\mathbf{x}_k, \mathbf{p}) = \begin{bmatrix} x_k + p_1 \\ y_k + p_2 \end{bmatrix}$$

Isometry/Euclidean : $S = 3$

$$- \mathbf{w}(\mathbf{x}_k, \mathbf{p}) = \begin{bmatrix} x_k \cos p_1 - y_k \sin p_1 + p_2 \\ x_k \sin p_1 + y_k \cos p_1 + p_3 \end{bmatrix}$$

• **Similitude/Similarity**: $S = 4$

$$- \mathbf{w}(\mathbf{x}_k, \mathbf{p}) = p_4 \begin{bmatrix} x_k \cos p_1 - y_k \sin p_1 + p_2 \\ x_k \sin p_1 + y_k \cos p_1 + p_3 \end{bmatrix}$$

• **Affine** : $S = 6$

$$- \mathbf{w}(\mathbf{x}_k, \mathbf{p}) = \begin{bmatrix} (1 + p_1)x_k - p_2y_k + p_3 \\ (1 + p_1)y_k + p_2x_k + p_4 \end{bmatrix}$$

State Space Model – Examples (cont'd)

- **Homography** : $S = 8$

$$- \mathbf{w}(\mathbf{x}_k, \mathbf{p}) = \left[\frac{(1+p_1)x_k + p_2y_k + p_3}{(1+p_7)x_k + p_8y_k + 1}, \frac{(1+p_4)y_k + p_5x_k + p_6}{(1+p_7)x_k + p_8y_k + 1} \right]^T$$

- **SL3 Homography**^[Benhimane04] : $S = 8$

$$- \mathbf{w}(\mathbf{x}_k, \mathbf{p}) = \mathbf{G} \hat{*} \begin{bmatrix} x_k \\ y_k \end{bmatrix}$$

$$\bullet \mathbf{G} = \exp(\sum_{i=1}^8 p_i \mathbf{A}_i) \in \mathbb{SL}(3), \mathbf{A}_i : \mathfrak{sl}(3) \text{ basis}$$

- **Corner Homography** : $S = 8$

$$- \mathbf{w}(\mathbf{x}_k, \mathbf{p}) = \mathbf{G} \hat{*} \begin{bmatrix} x_k \\ y_k \end{bmatrix}$$

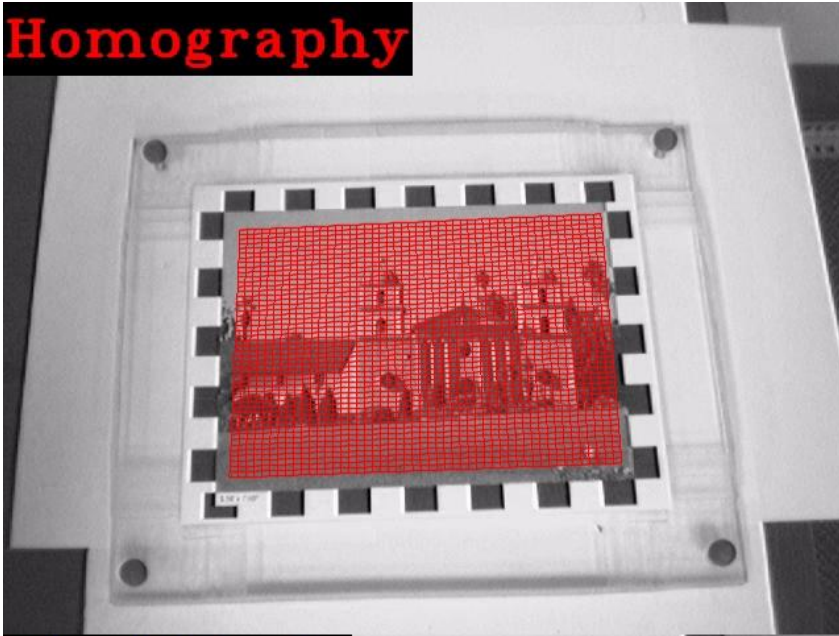
$$\bullet \mathbf{G} = \underset{\mathbf{M}}{\operatorname{argmin}} \sum_{i=1}^4 \left\| \mathbf{M} \hat{*} \begin{bmatrix} c_{ix} \\ c_{iy} \end{bmatrix} - \begin{bmatrix} c_{ix} + p_{2i-1} \\ c_{iy} + p_{2i} \end{bmatrix} \right\|^2$$

$$\bullet \left\{ c_i = \begin{bmatrix} c_{ix} \\ c_{iy} \end{bmatrix} \mid 1 \leq i \leq 4 \right\} : \text{bounding box corners}$$

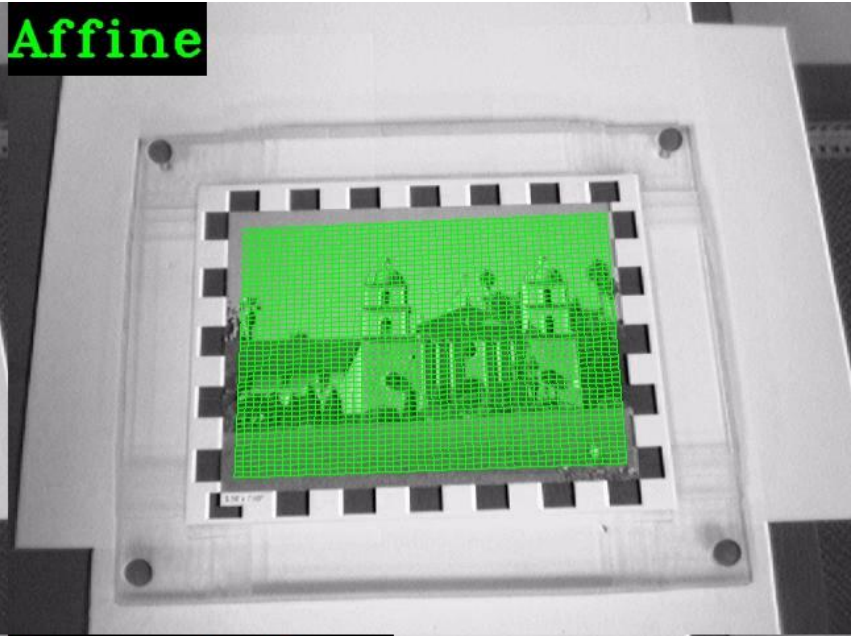
$$\mathbf{G} \hat{*} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} g_{00}x + g_{01}y + g_{02} \\ g_{10}x + g_{11}y + g_{12} \\ g_{20}x + g_{21}y + g_{22} \end{bmatrix}^T \text{ with } \mathbf{G} = \begin{bmatrix} g_{00} & g_{01} & g_{02} \\ g_{10} & g_{11} & g_{12} \\ g_{20} & g_{21} & g_{22} \end{bmatrix}$$

State Space Model – Examples (Demo)

Homography



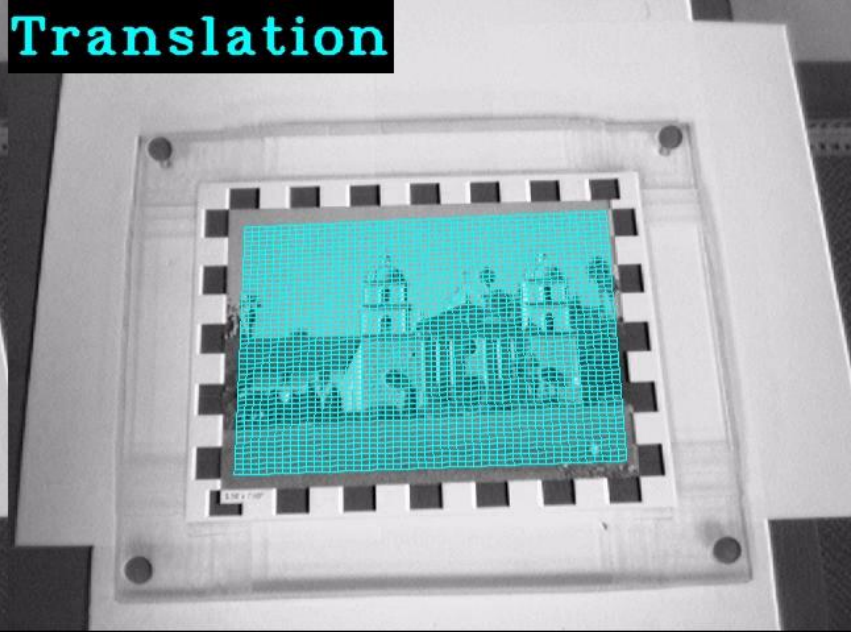
Affine



Similitude



Translation



(in 2 weeks) Homography = Planar Projective Warping



$$x_i' = Hx_i$$
$$i = 1 \dots 4$$

A novel view rendered via four
points with known structure

Results – State Space Models (Demo)

nl_juice_s5: 121

8DOF



6DOF



4DOF



2DOF



Search Method

$$\mathbf{p}_t = \underset{\mathbf{p}}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$

- Optimization method that finds the SSM parameters corresponding to the warped patch that maximizes the AM similarity function.
- Two main categories:
 - Gradient descent
 - Newton or Gauss Newton method
 - Stochastic Search
 - Sampling based

Simple image registration algorithm

SSD error norm

$$E(u, v) = \sum_{x, y} (I(x + u, y + v) - T(x, y))^2$$

Exhaustive search:

For each offset (u, v)

compute $E(u, v)$;

Choose (u, v) which minimizes $E(u, v)$;

(Gauss) Newton optimization:

Solve

$$\mathbf{u} = \mathbf{M} \backslash \mathbf{Im_t} \quad \begin{pmatrix} \vdots \\ -\frac{\partial \text{Im}}{\partial t} \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots & \vdots \\ \frac{\partial \text{Im}}{\partial x} & \frac{\partial \text{Im}}{\partial y} \\ \vdots & \vdots \end{pmatrix} \begin{pmatrix} u_x \\ u_y \end{pmatrix}$$

Search Method – Examples (Gradient Descent)

- Variants of Lucas Kanade (LK)^[Baker01] method
 - Forward Additive (**FALK**)
 - $\Delta \mathbf{p}_t = \underset{\Delta \mathbf{p}_t}{\operatorname{argmax}} f \left(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p}_{t-1} + \Delta \mathbf{p}_t)) \right)$
 - $\mathbf{p}_t = \mathbf{p}_{t-1} + \Delta \mathbf{p}_t$
 - Inverse Additive (**IALK**)
 - uses constant approximation of $\nabla \mathbf{I}_t$ computed from \mathbf{I}_0
 - Forward Compositional (**FCLK**)
 - $\Delta \mathbf{p}_t = \underset{\Delta \mathbf{p}_t}{\operatorname{argmax}} f \left(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{w}(\mathbf{x}, \Delta \mathbf{p}_t), \mathbf{p}_{t-1})) \right)$
 - $\mathbf{p}_t = \mathbf{p}_{t-1} \circ \Delta \mathbf{p}_t$
 - Inverse Compositional (**ICLK**)
 - $\Delta \mathbf{p}_t = \underset{\Delta \mathbf{p}_t}{\operatorname{argmax}} f \left(\mathbf{I}_0(\mathbf{w}(\mathbf{x}, \Delta \mathbf{p}_t)), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p}_{t-1})) \right)$
 - $\mathbf{p}_t = \mathbf{p}_{t-1} \circ \Delta \mathbf{p}_t^{-1}$
- Efficient Second Order Minimization (**ESM**)^[Benhimane04]
 - combines FCLK and ICLK

Homogenous coordinates: How to translate a 2D point:

- Old way: $x' = x + dx$
$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix}$$

- New way: $x' = M * dx = M \circ dx$

$$\begin{pmatrix} x' \\ y' \\ 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & \Delta x \\ 0 & 1 & \Delta y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}$$

- Can chain many transf: $x' = M1 * M2 * dx$

- Euclidean transform SE2:
$$p' = \begin{pmatrix} R & T \\ 0 & 0 & 0 & 1 \end{pmatrix} p$$

Search Method – Examples (Stochastic)

- Nearest Neighbor Search (**NN**) [Dick13]
 - generate samples by warping $\mathbf{I}_0(\mathbf{x})$
 - find the nearest neighbor to $\mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p}_{t-1}))$ and update \mathbf{p}_{t-1} with the *inverse* of the corresponding $\Delta\mathbf{p}_t$
 - combined with ICLK for stability (**NNIC**)
- Particle Filter (**PF**) [Kwon14]
 - generate samples by warping $\mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p}_{t-1}))$
 - compute weight for each and estimate $\Delta\mathbf{p}_t$ as weighted average of samples

Appearance Model

$$\mathbf{p}_t = \underset{\mathbf{p}}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$

- A similarity measure between two image patches:
 - candidate warped patch from the current image
 - template extracted from the initial image
- Two main categories:
 - SSD like
 - Robust^[Richa12]

Appearance Model – Examples

- Sum of Squared Differences (**SSD**)^[Baker01]
 - $f(I_0, I_t) = -\frac{1}{2} \| I_0 - I_t \|^2$
- Sum of Conditional Variance (**SCV**)^[Richa11]
 - $f(I_0, I_t) = -\frac{1}{2} \| E[I_t | I_0] - I_t \|^2$
 - Using several joint distributions computed from corresponding sub regions of I_t and I_0 gives a variant called **LSCV**^[Richa14]
- Reversed Sum of Conditional Variance (**RSCV**)^[Dick13]
 - $f(I_0, I_t) = -\frac{1}{2} \| I_0 - E[I_0 | I_t] \|^2$
- Zero mean Normalized Cross Correlation (**ZNCC**)^[Ruthotto10]
 - $f(I_0, I_t) = -\frac{1}{2} \left\| \frac{I_0 - \mu_0}{\sigma_0} - \frac{I_t - \mu_t}{\sigma_t} \right\|^2$

Appearance Model – Examples (cont'd)

- Mutual Information (**MI**)^[Dame10]

$$-f(\mathbf{I}_0, \mathbf{I}_t) = \sum_{ij} P_{I_t I_0}(i, j) \log \left(\frac{P_{I_t I_0}(i, j)}{P_{I_t}(i) P_{I_0}(j)} \right)$$

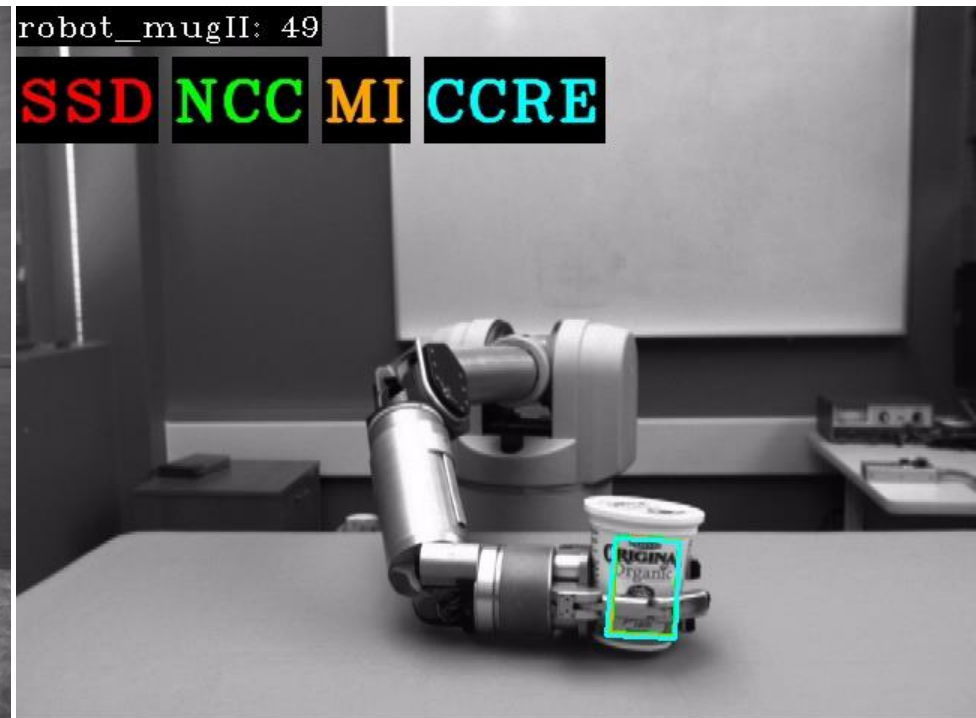
- Cross Cumulative Residual Entropy (**CCRE**)^[Richa12]

$$-f(\mathbf{I}_0, \mathbf{I}_t) = \sum_{ij} P_{I_t I_0}^*(i, j) \log \left(\frac{P_{I_t I_0}^*(i, j)}{P_{I_t}^*(i) P_{I_0}(j)} \right)$$

- Normalized Cross Correlation (**NCC**)^[Scandaroli12]

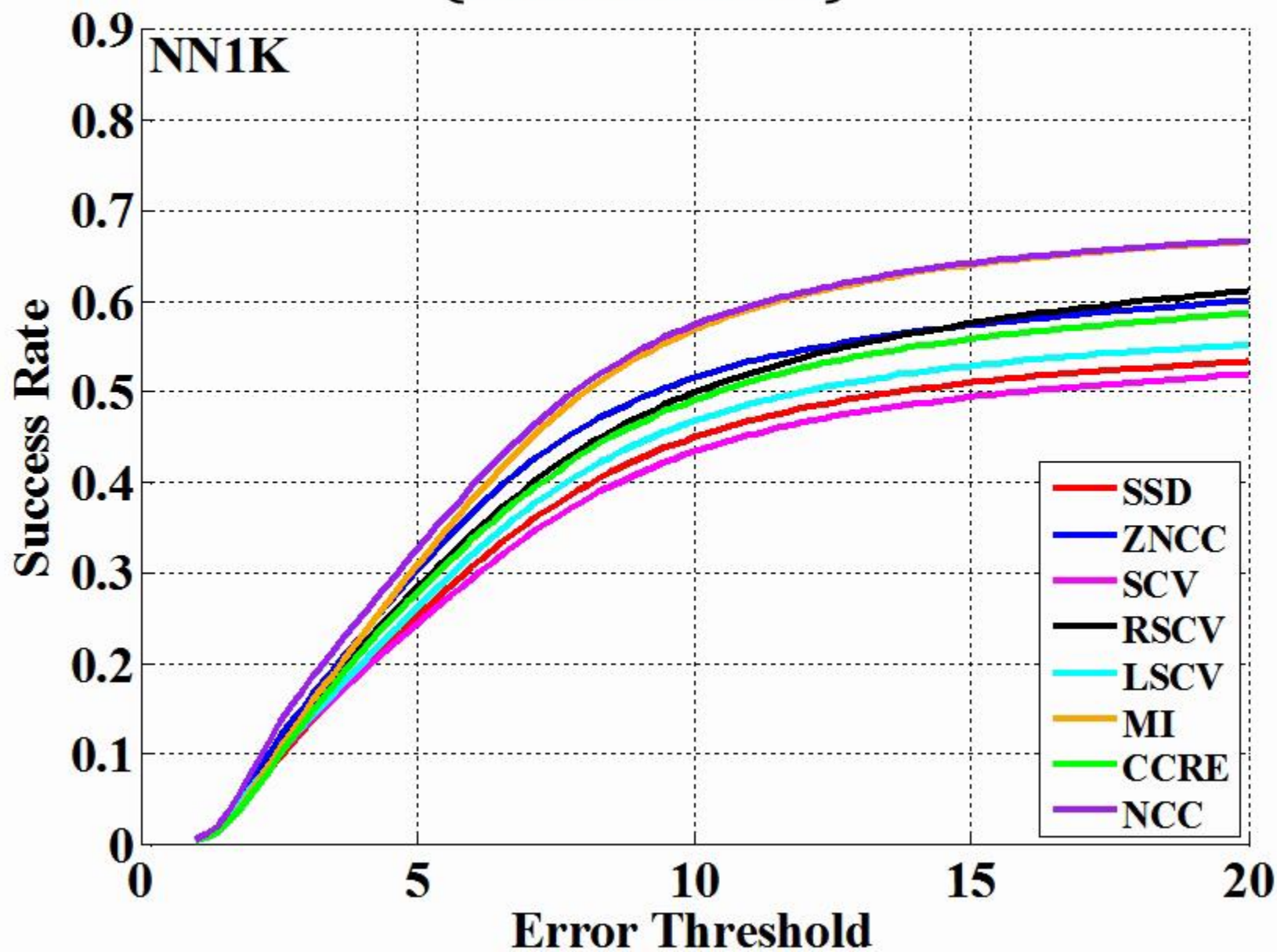
$$-f(\mathbf{I}_0, \mathbf{I}_t) = \frac{\mathbf{I}_0 - \mu_0}{\sigma_0} \cdot \frac{\mathbf{I}_t - \mu_t}{\sigma_t}$$

Results – Appearance Models (Demo)

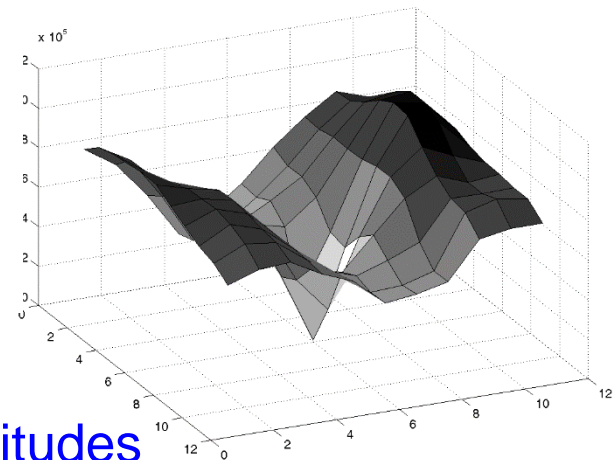
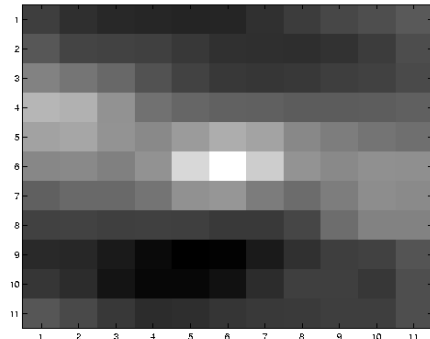


FCLK with Homography

Results – Appearance Models (Stochastic)



High textured region

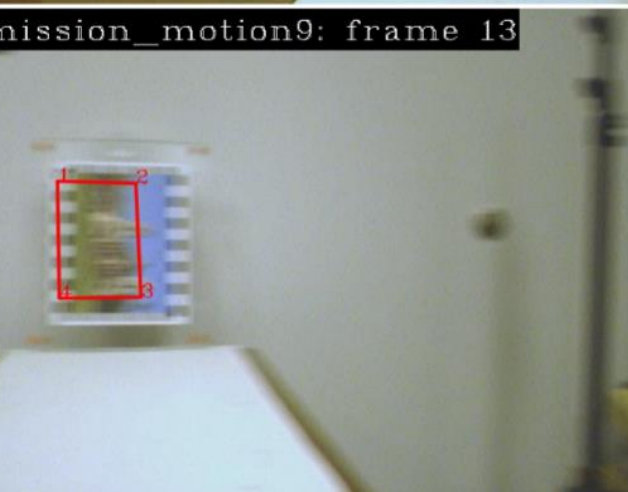
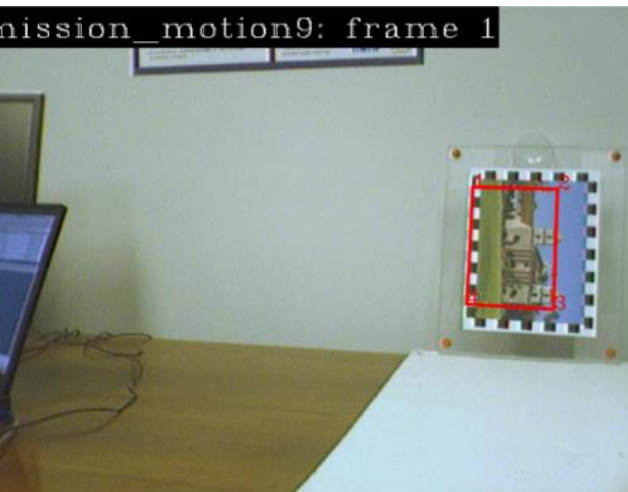


$$\sum \nabla I (\nabla I)^T$$

- gradients are different, large magnitudes
- large λ_1 , large λ_2

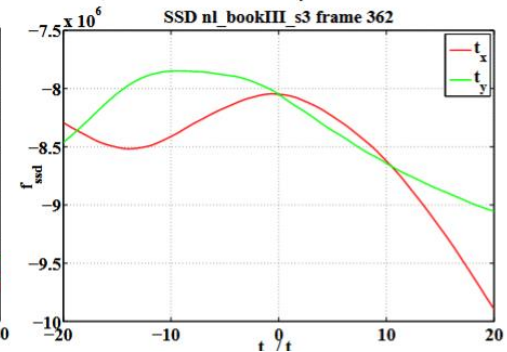
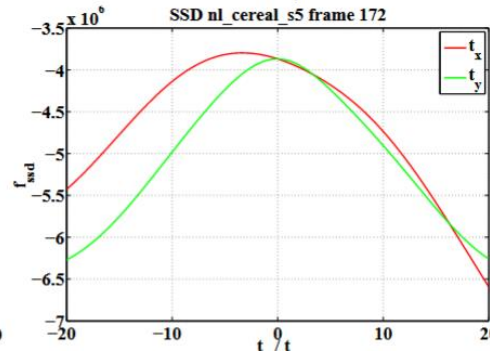
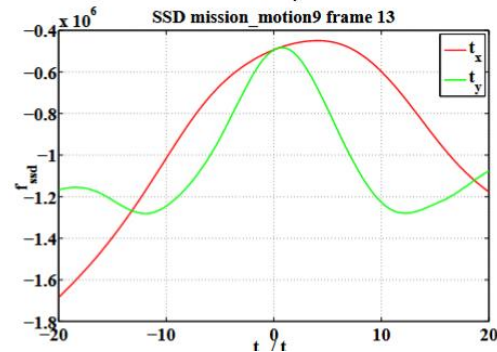
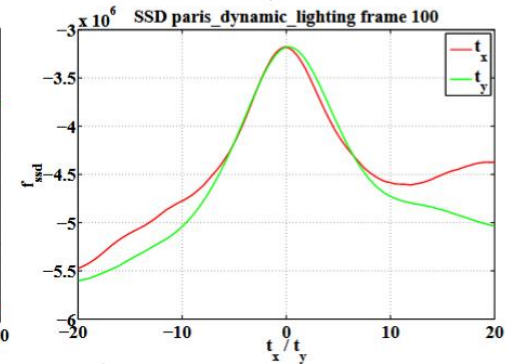
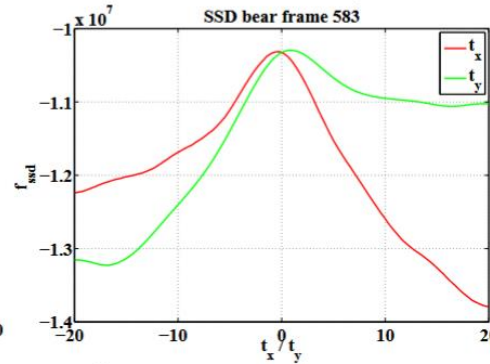
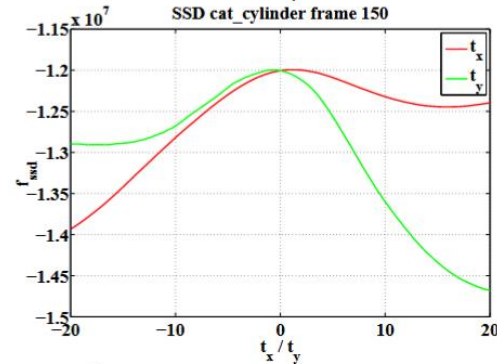
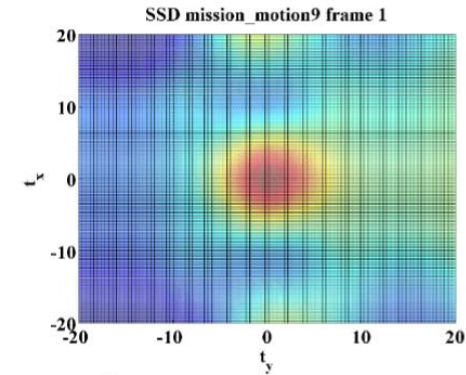
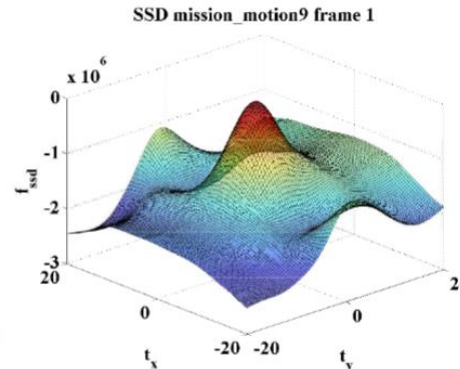
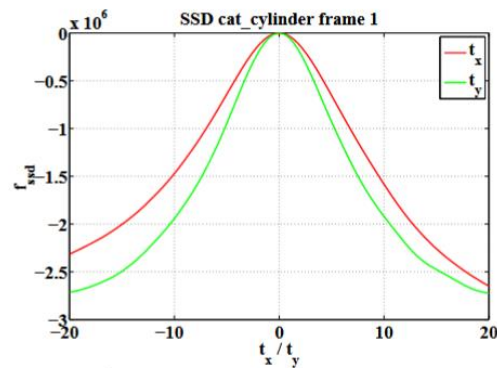
Appearance model

Test images



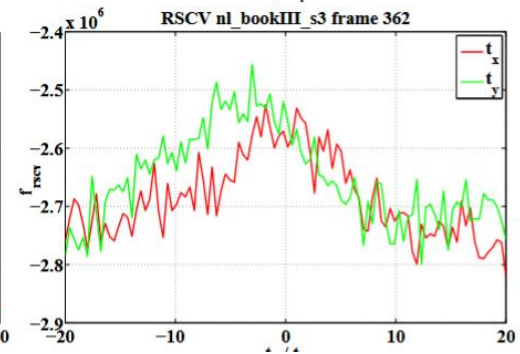
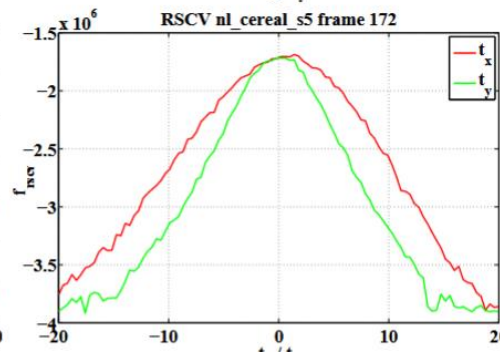
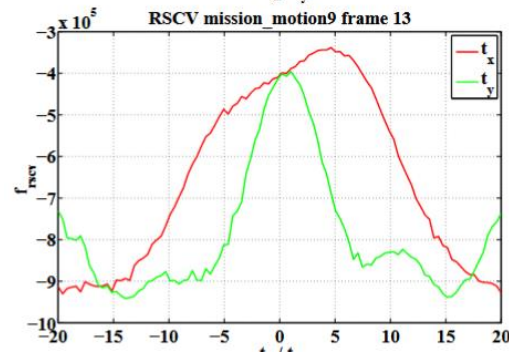
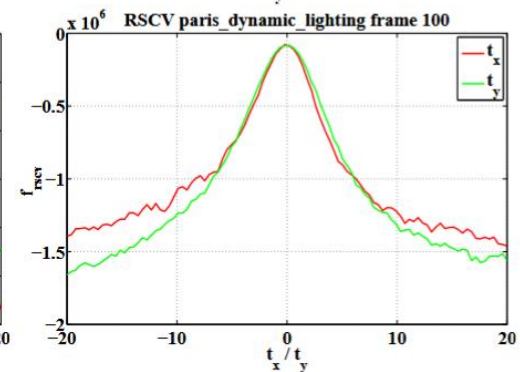
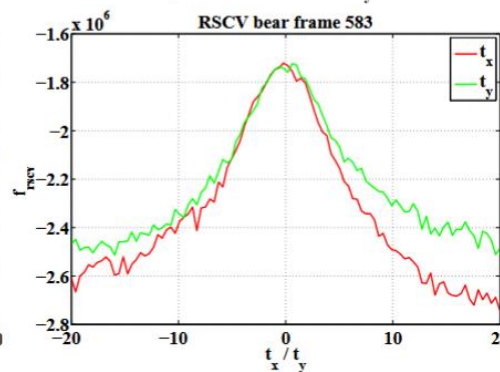
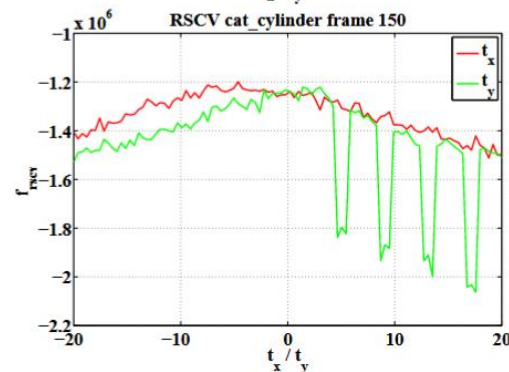
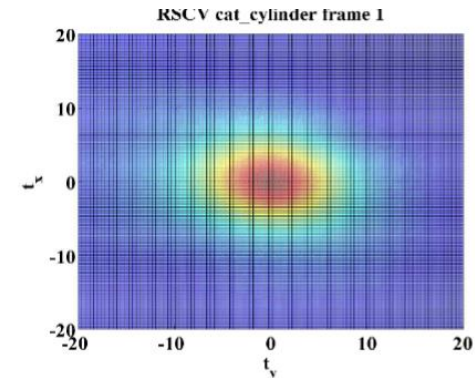
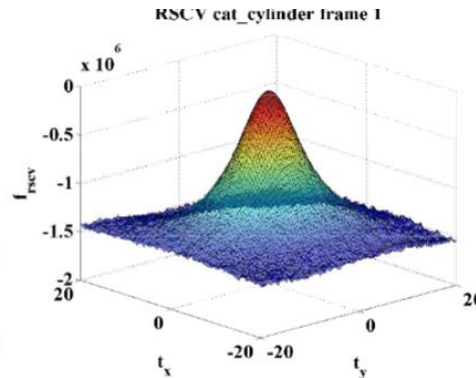
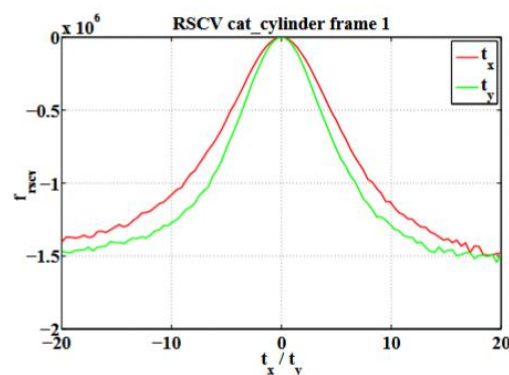
Appearance model

$L^2 ||T-I||^2$ aka “SSD”



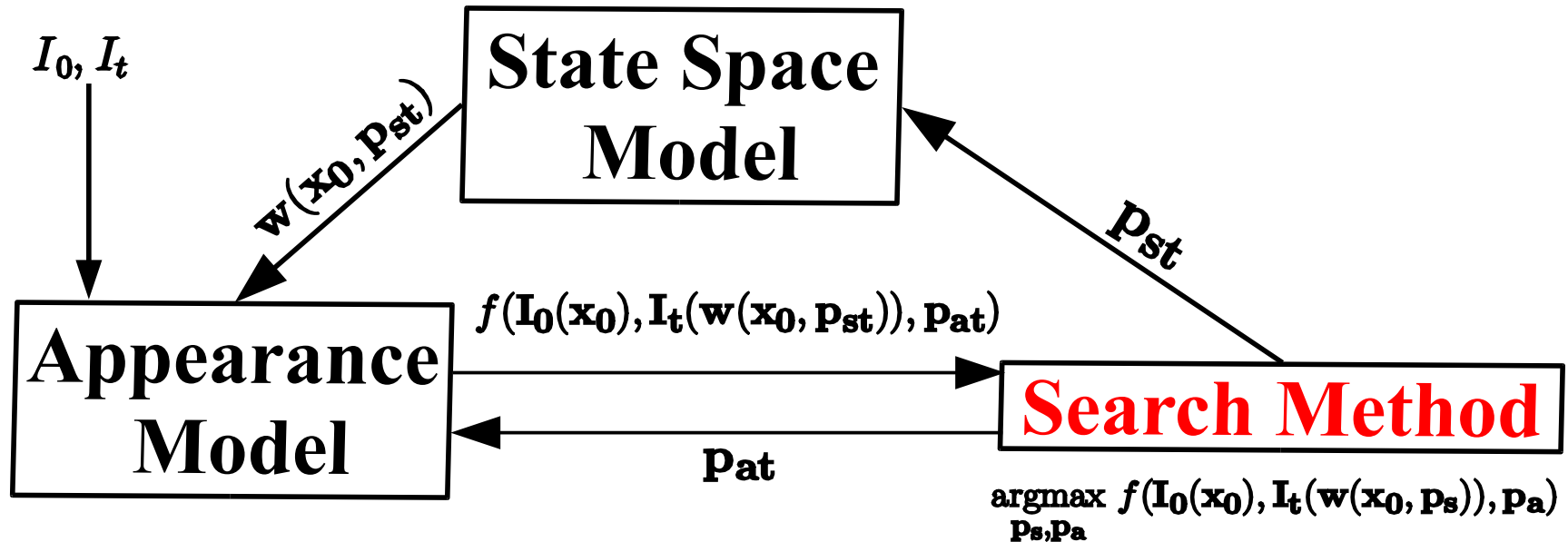
Appearance model

RSCV – Reversed Sum of Conditional Variance



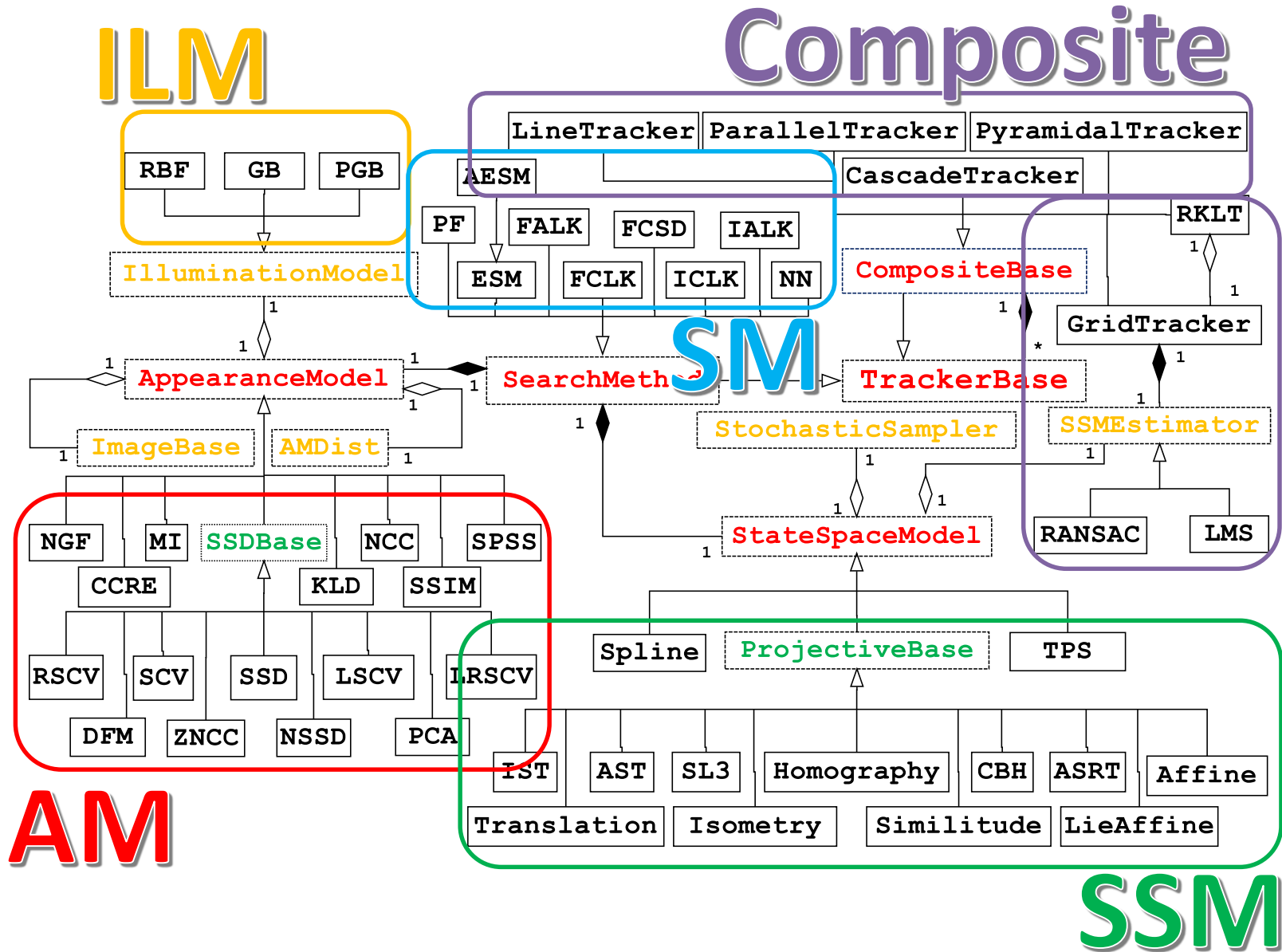
System design

$$\mathbf{p}_t = \underset{\mathbf{p}_s, \mathbf{p}_a}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}_0), \mathbf{I}_t(\mathbf{w}(\mathbf{x}_0, \mathbf{p}_s)), \mathbf{p}_a)$$



- Search Method (**SM**)
 - **Finds** the warp that maximizes the similarity measure

System Design



Evaluation Methodology - Datasets

- 4 large publicly available datasets with a total of over **100K** frames

- TMT
- UCSB
- LinTrack
- PAMI

Dataset	Without Subsequences			With Subsequences		
	Sequences	Total Frames	Trackable Frames	Sub-sequences	Total Frames	Trackable Frames
TMT	109	70592	70483	1090	390470	389380
UCSB	96	6889	6793	960	41170	40210
LinTrack	3	12477	12474	30	68700	68670
PAMI	28	16511	16483	280	91400	91120
Total	236	106469	106233	2360	591740	589380

- Each sequence tested from 10 different starting points for an effective total of nearly **600K** frames

Evaluation Methodology – Performance Metric

- **Alignment Error (E_{AL})**

- $E_{AL} = \frac{1}{4} \| \mathbf{C}_{track} - \mathbf{C}_{gt} \|$

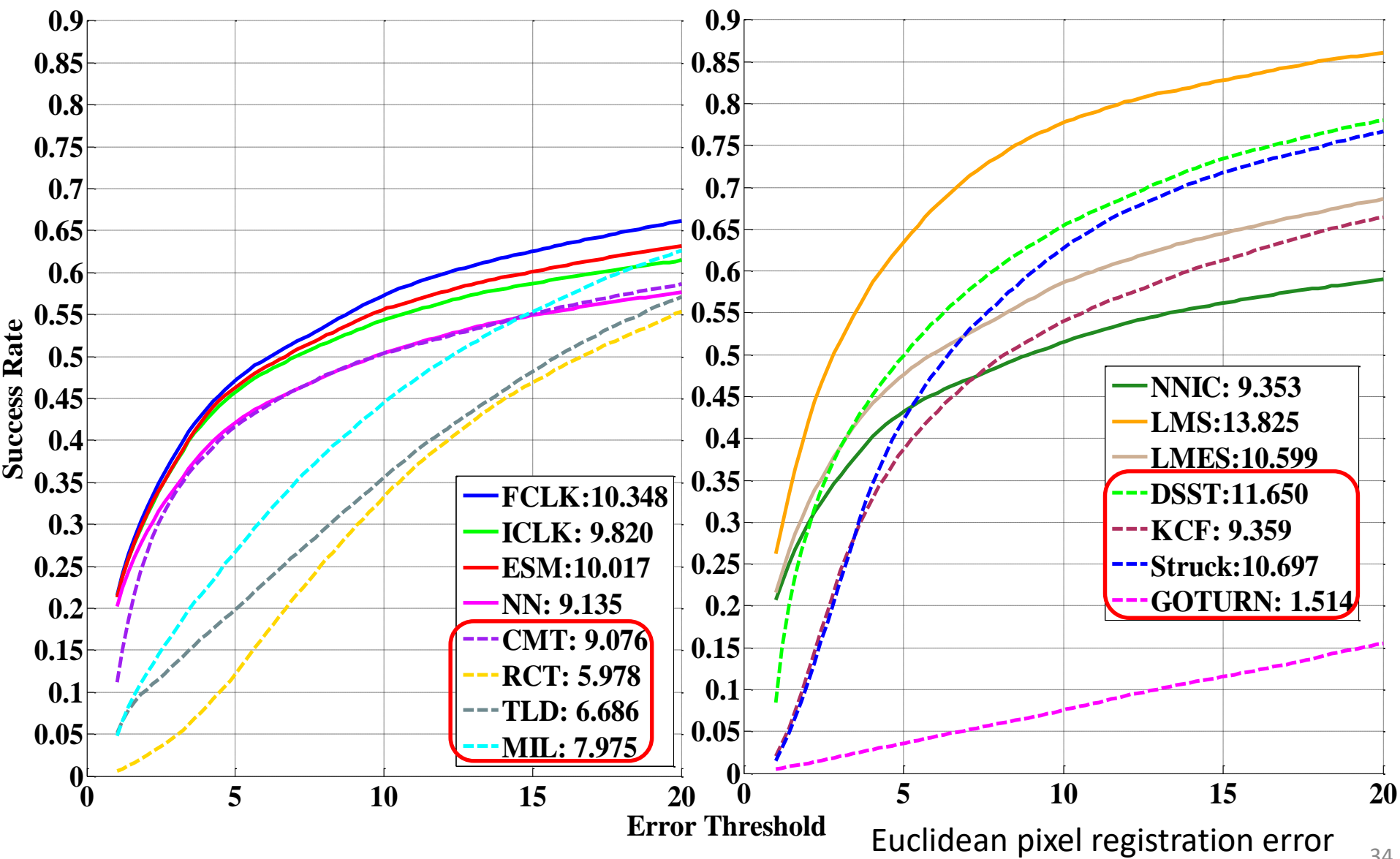
- **Success Rate (SR)**

- **x** axis : error threshold $t_p \in [0, 20]$
 - **y** axis : fraction of frames with $E_{AL} < t_p$
 - each sequence tracked from 10 different starting points
 - measures both accuracy and robustness

- **Failure Rate (FR)**

- reinitialize whenever E_{AL} exceeds 20
 - count the number of such failures
 - additional metric for tracking robustness

Results: Learning vs. 2DOF Registration Based Trackers (Accuracy)



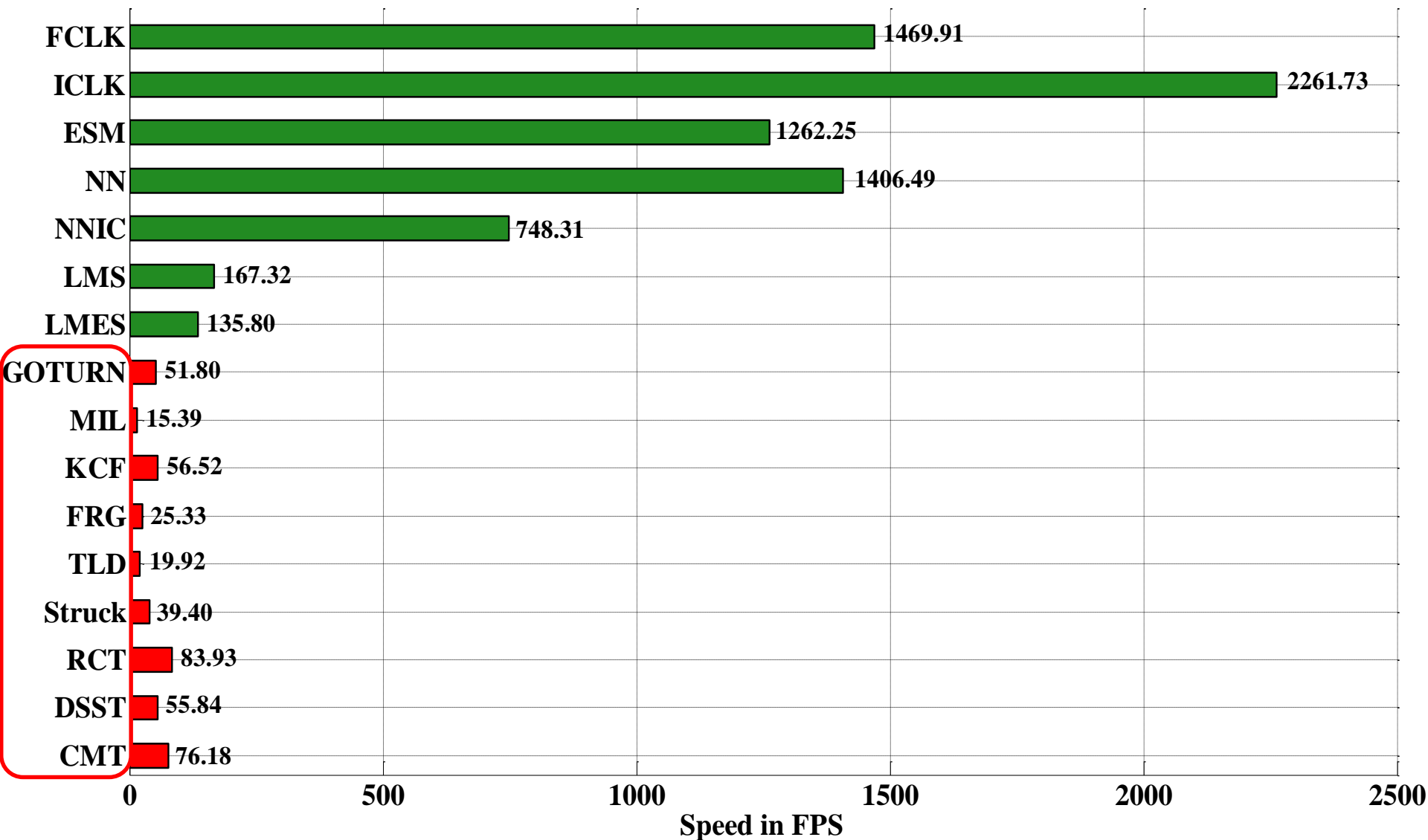
RBT vs Learn vs GoTurn

nl_cereal_s3: frame 1

GOTURN DSST ESMLMES



Results: Learning vs. 2DOF Registration Based Trackers (Speed)

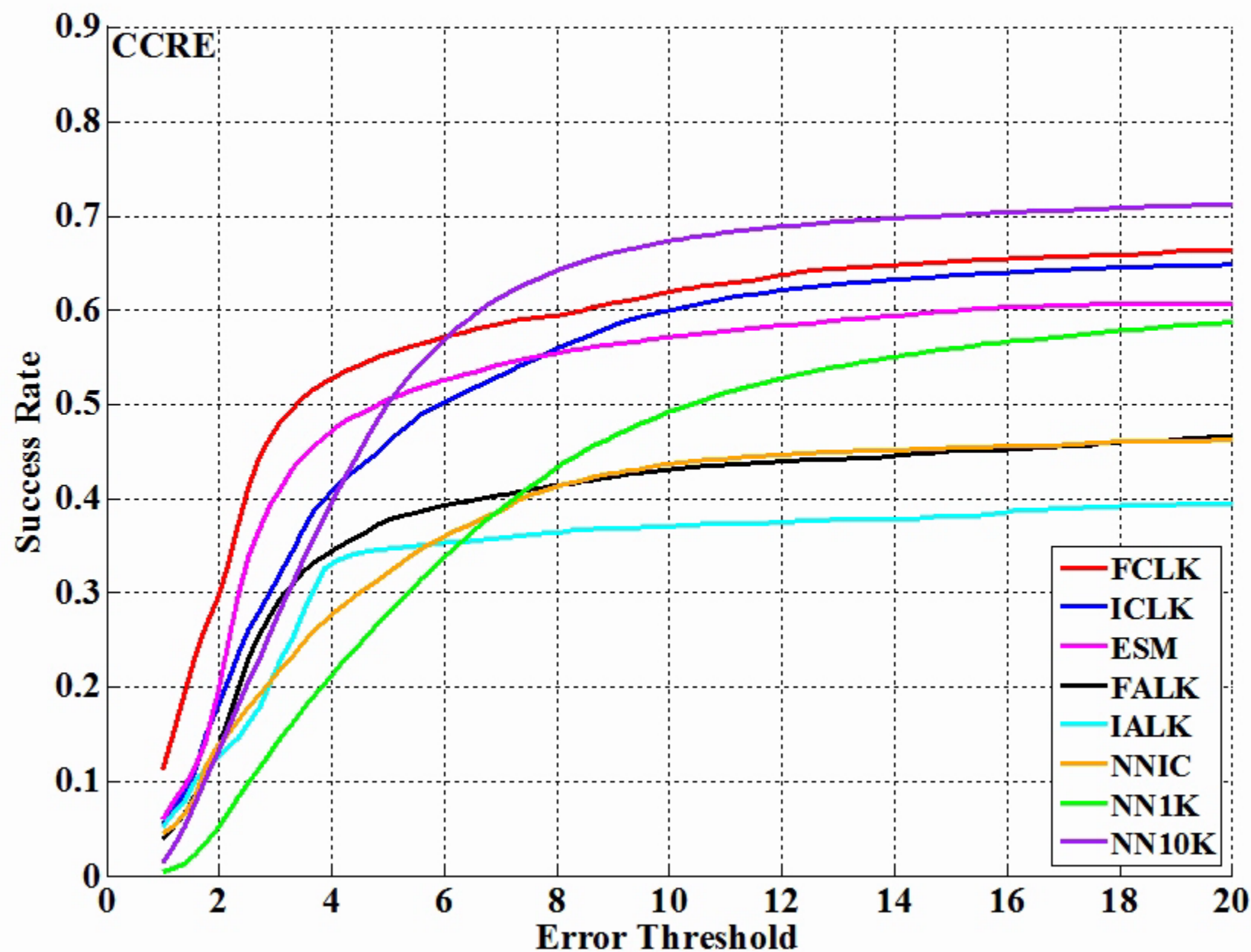


Results –Learning Based Trackers (Demo)



ZNCC with Translation

Results – Search Methods (**Robust**)



Results – Search Methods (Demo)

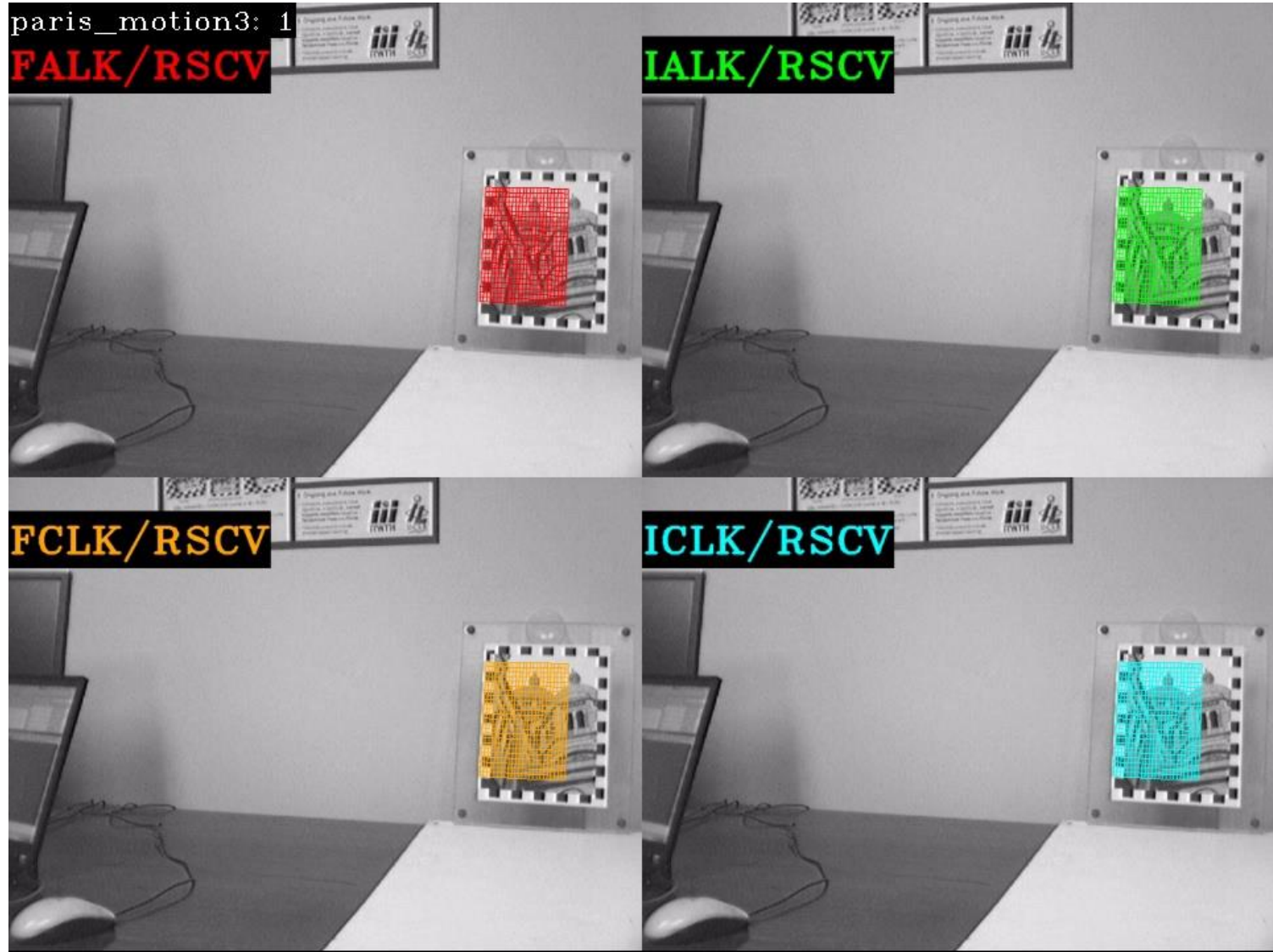
- The four variants of Lucas Kanade fail at different times
- Sequences from TMT



RSCV with Homography

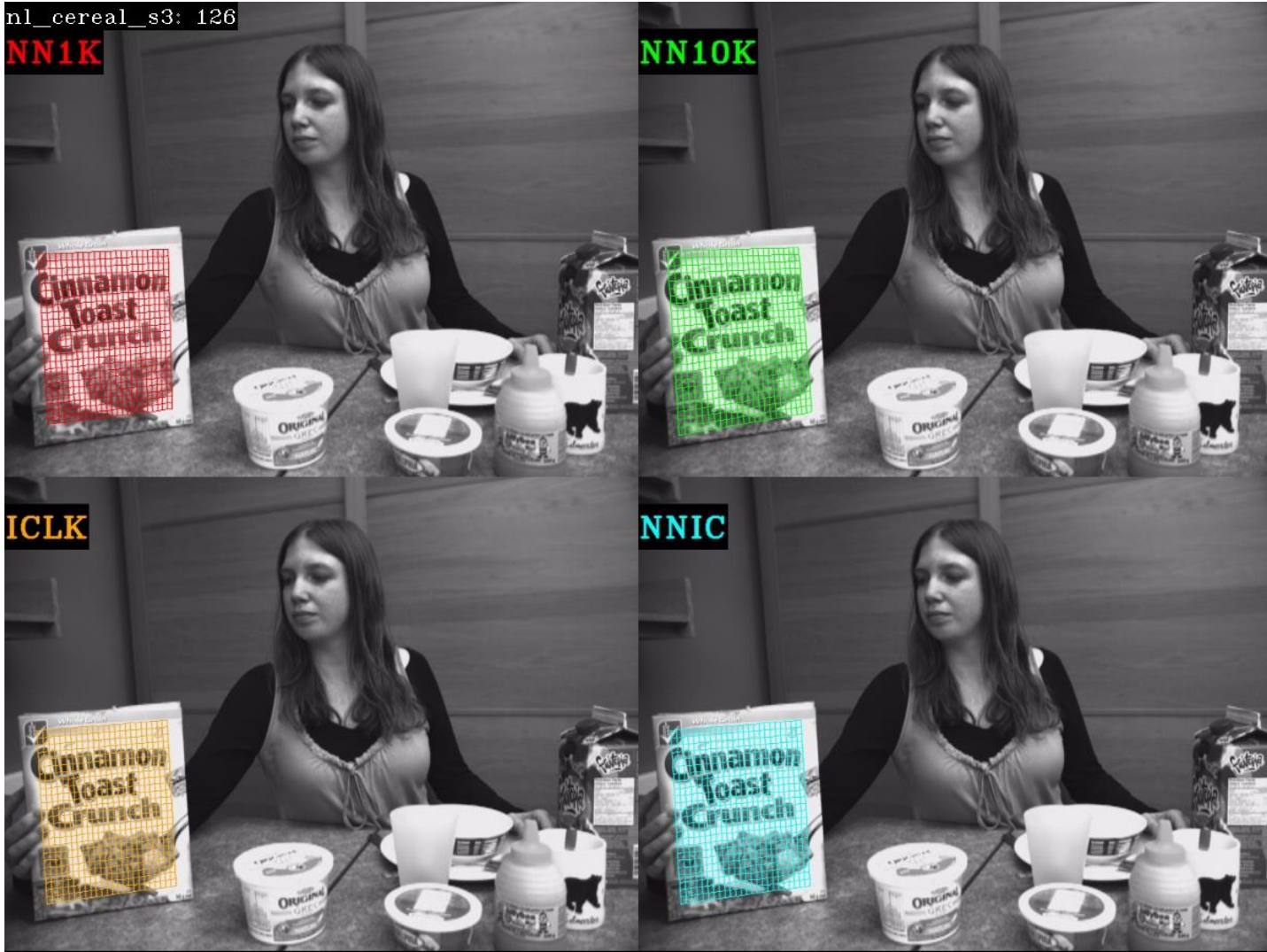
Results – Search Methods (Demo)

- The four variants of Lucas Kanade fail at different times
- Sequence from UCSB



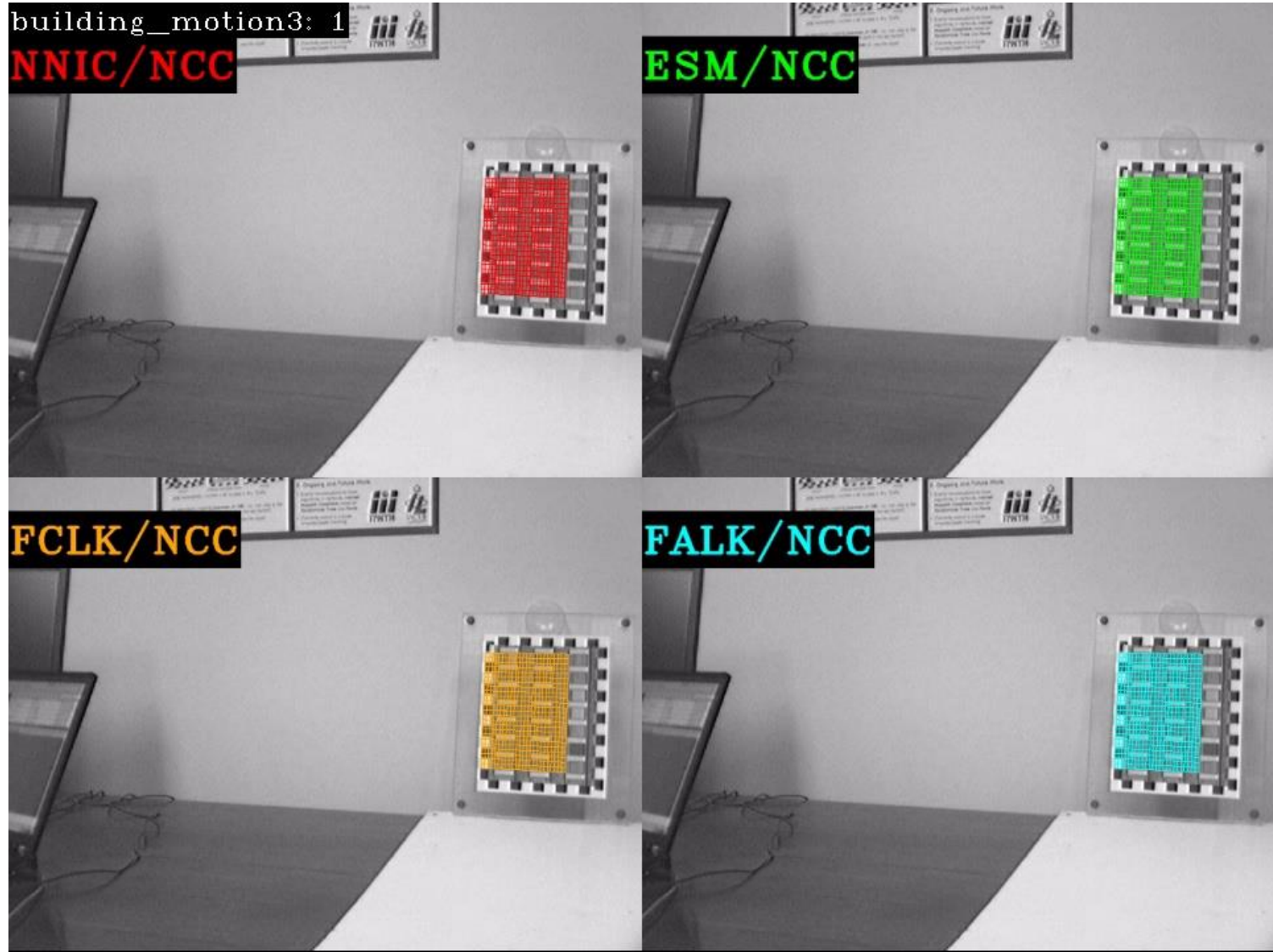
Results – Search Methods (Demo)

- NN has more jitter than LK type SMs
 - decreases with more samples



Results – Search Methods (Demo)

- NNIC is more robust to motion blur
- Sequence from UCSB



Conclusions?

- Tested **different combinations** of sub modules leading to several interesting observations that were missing in the original papers.
 - used two large datasets with over 77,000 frames in all to ensure statistical significance.
- Compared **robust similarity metrics** with traditional SSD type measures.
- Compared formulations against **online learning based trackers** to validate their usability for precise tracking
- Provided **an open source tracking framework** called **MTF** using which all results can be reproduced
 - can also address practical tracking requirements with its efficient C++ implementation

MTF is available at: <http://webdocs.cs.ualberta.ca/~vis/mtf/> along with all datasets and this presentation

Project and Research opportunities in video tracking

- Combine learning and registration tracking search methods
 - Direct deep network methods not precise
 - Predict with deep network, refine with registration
- New appearance models
 - Kullback-Liebler divergence
 - AM based on deep features
- More detailed experimental evaluation
 - Study failure causes in individual frames, solve those
 - Our TMT benchmark data marked up for this (per-frame annotation of blur, motion etc.)

Questions ?

- Tested **different combinations** of sub modules leading to several interesting observations that were missing in the original papers.
 - used two large datasets with over 77,000 frames in all to ensure statistical significance.
- Compared **robust similarity metrics** with traditional SSD type measures.
- Compared formulations against **online learning based trackers** to validate their usability for precise tracking
- Provided **an open source tracking framework** called **MTF** using which all results can be reproduced
 - can also address practical tracking requirements with its efficient C++ implementation

MTF is available at: <http://webdocs.cs.ualberta.ca/~vis/mtf/> along with all datasets and this presentation

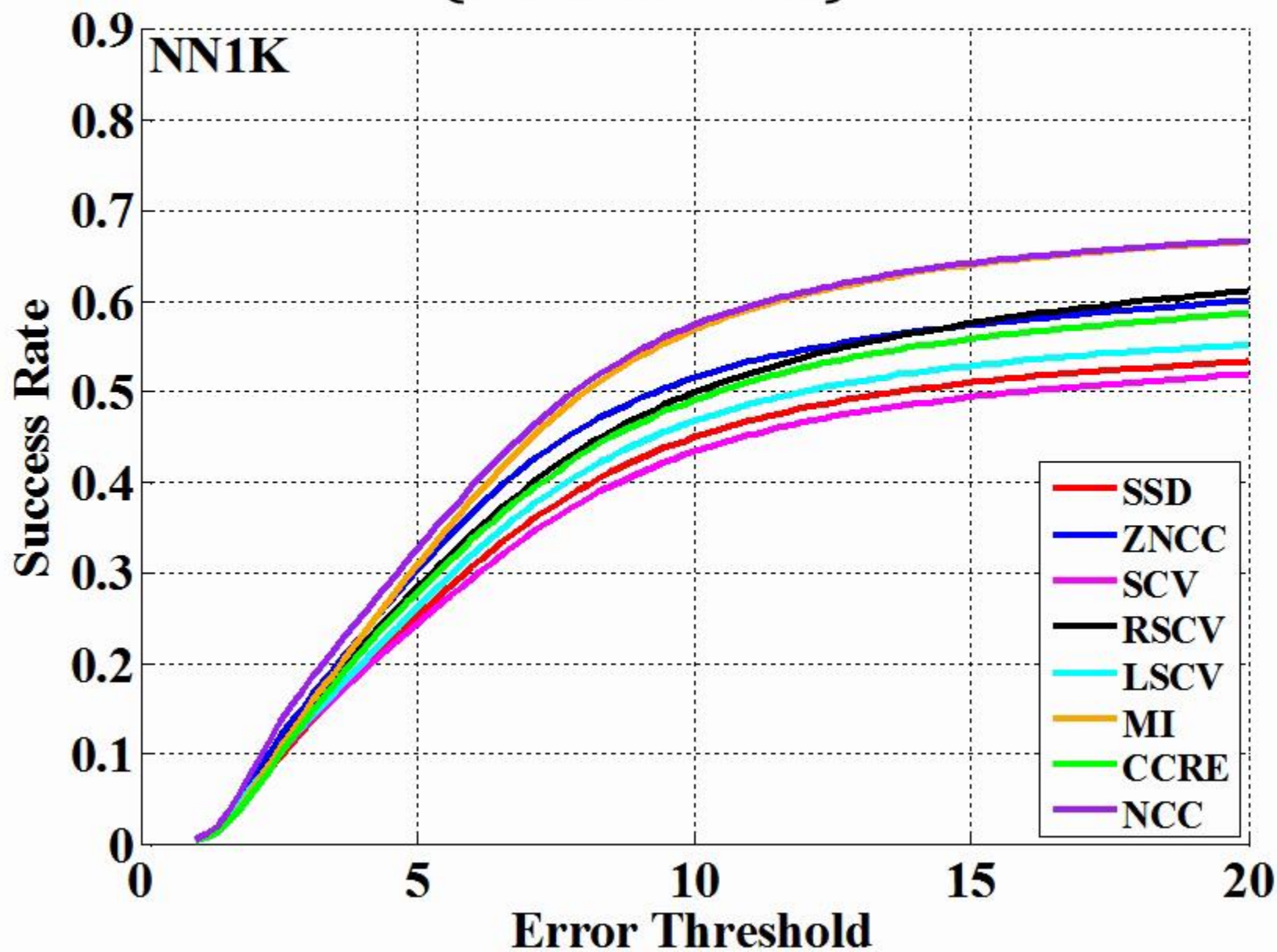
References

- S. Baker and I. Matthews, “*Equivalence and Efficiency of Image Alignment Algorithms*”, CVPR 2001
- S. Benhimane and E. Malis, “*Real-time image-based tracking of planes using efficient second-order minimization*”, IROS 2004
- A. Dame and E. Marchand, “*Accurate Real-time Tracking Using Mutual Information*”, ISMAR 2010
- T. Dick, C. Perez, A. Shademan and M. Jagersand, “*Realtime Registration-Based Tracking via Approximate Nearest Neighbor Search*”, RSS 2013
- S. Gauglitz, T. Hollerer and M. Turk. “*Evaluation of interest point detectors and feature descriptors for visual tracking*”, IJCV 2011
- J. Kwon, H. S. Lee, F. C. Park, and K. M. Lee, “*A Geometric Particle Filter for Template-Based Visual Tracking*”, TPAMI 2014
- R. Richa, R. Sznitman, R. Taylor and G. Hager, “*Visual Tracking Using the Sum of Conditional Variance*”, IROS 2011

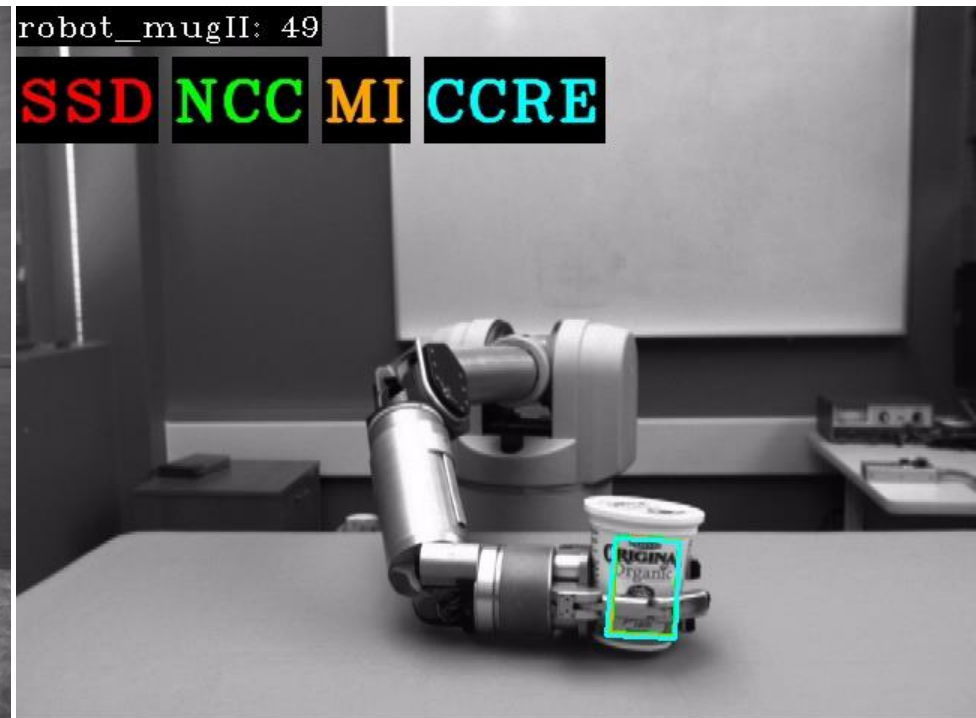
References (cont'd)

- R. Richa, R. Sznitman and G. Hager, “*Robust Similarity Measures for Gradient-based Direct Visual Tracking*”, CIRL Technical Report 2012
- R. Richa, M. Souza, G. Scandaroli, E. Comunello and A. Wangenheim, “*Direct visual tracking under extreme illumination variations using the sum of conditional variance*”, ICIP 2014
- A. Roy, X. Zhang, N. Wolleb, C. P. Quintero and M. Jagersand, “*Tracking Benchmark and Evaluation for Manipulation Tasks*”, ICRA 2015
- L. Ruthotto, “*Mass-preserving registration of medical images*”, Thesis 2010
- G. G. Scandaroli, M. Meilland and R. Richa, “*Improving NCC-Based Direct Visual Tracking*”, ECCV 2012
- A. Singh and M. Jagersand, “*Modular Tracking Framework: A Unified Approach to Registration based Tracking*”, arXiv:1602.09130, 2016

Results – Appearance Models (Stochastic)



Results – Appearance Models (Demo)

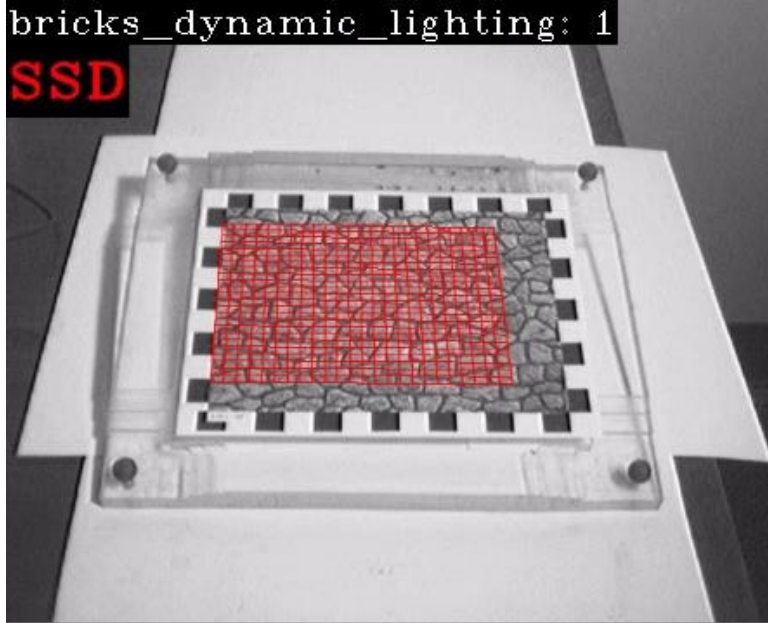


FCLK with Homography

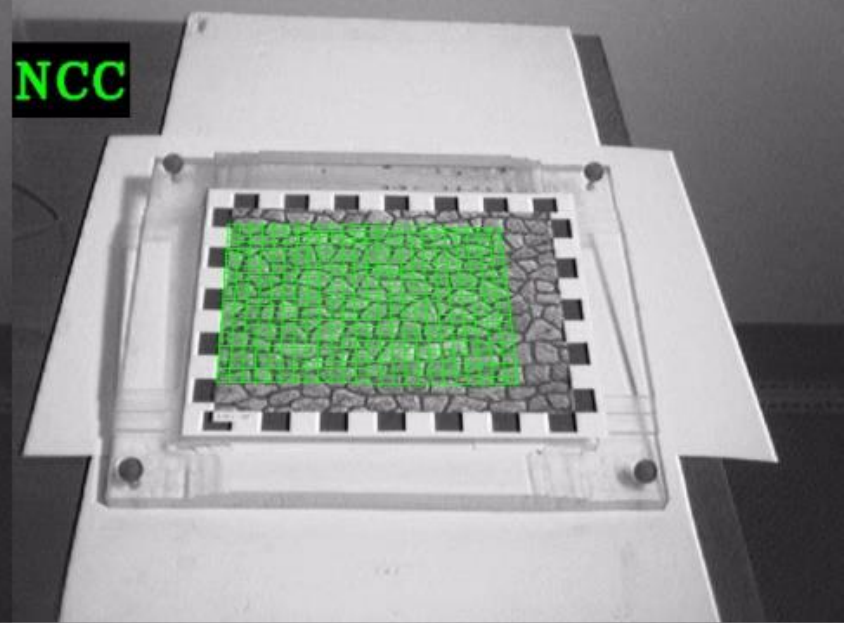
Results – Appearance Models (Demo)

bricks_dynamic_lighting: 1

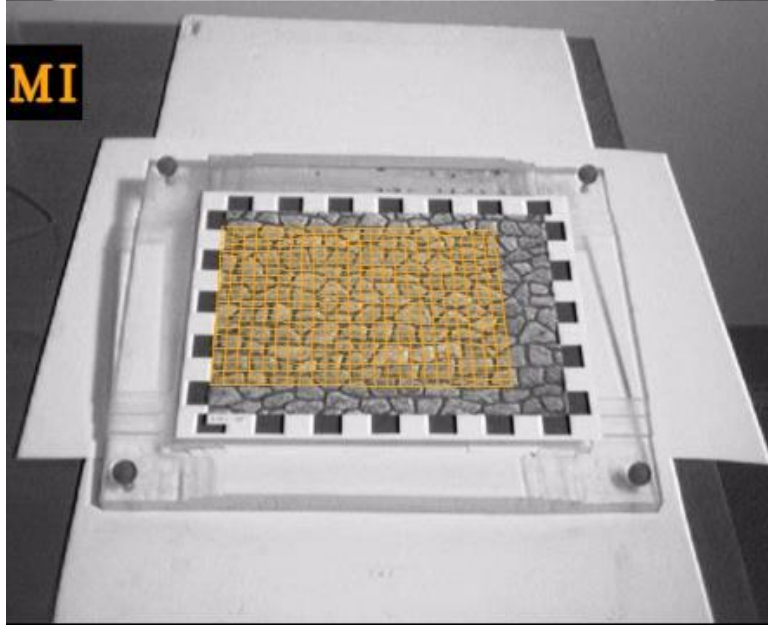
SSD



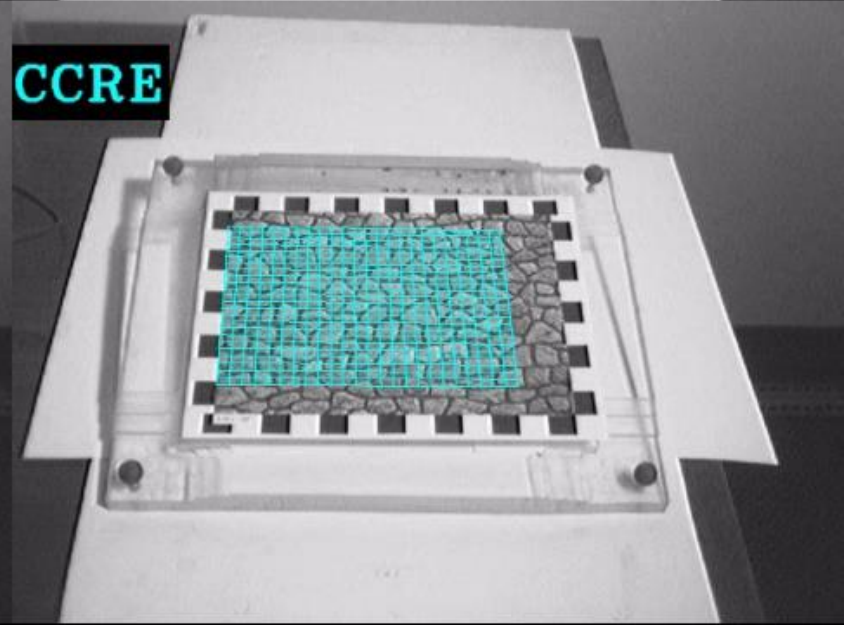
NCC



MI



CCRE



Results – State Space Models (Demo)

nl_juice_s5: 121

8DOF



6DOF



4DOF



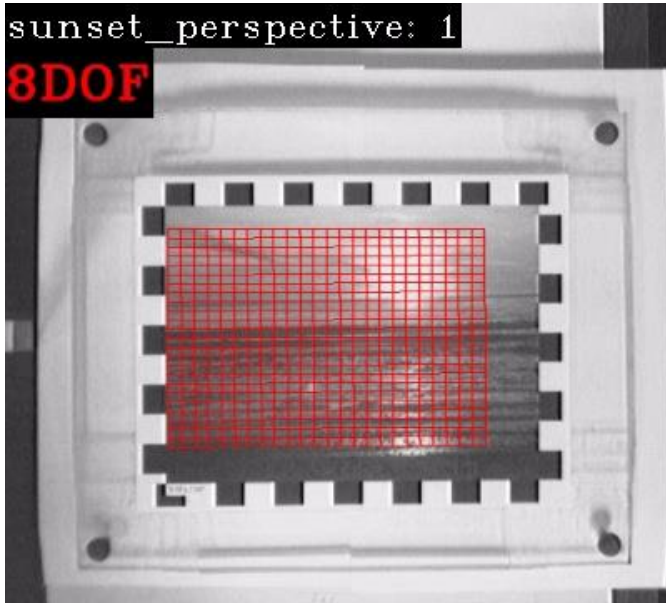
2DOF



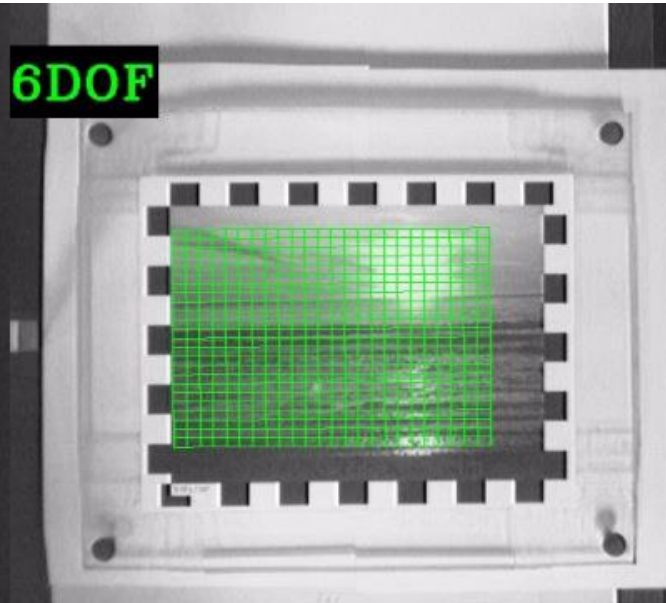
Results – State Space Models (Demo)

`sunset_perspective: 1`

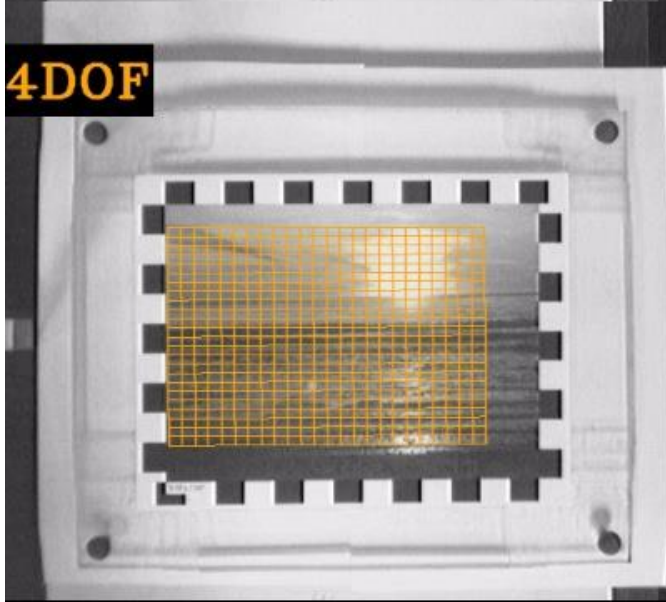
8DOF



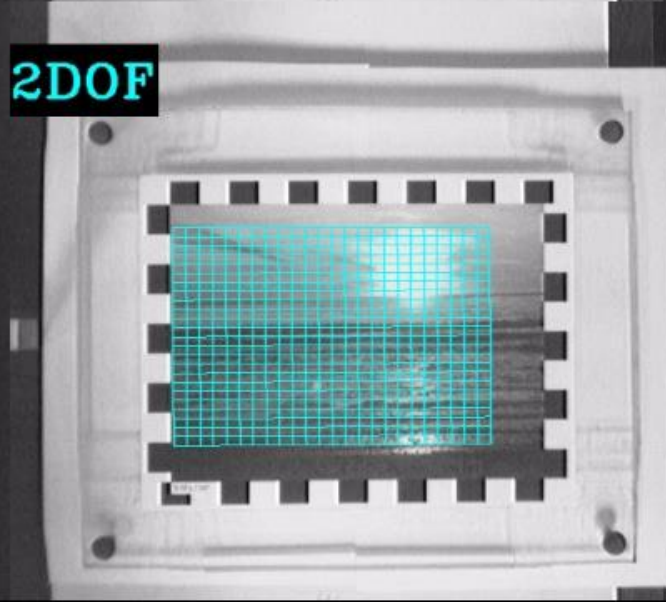
6DOF



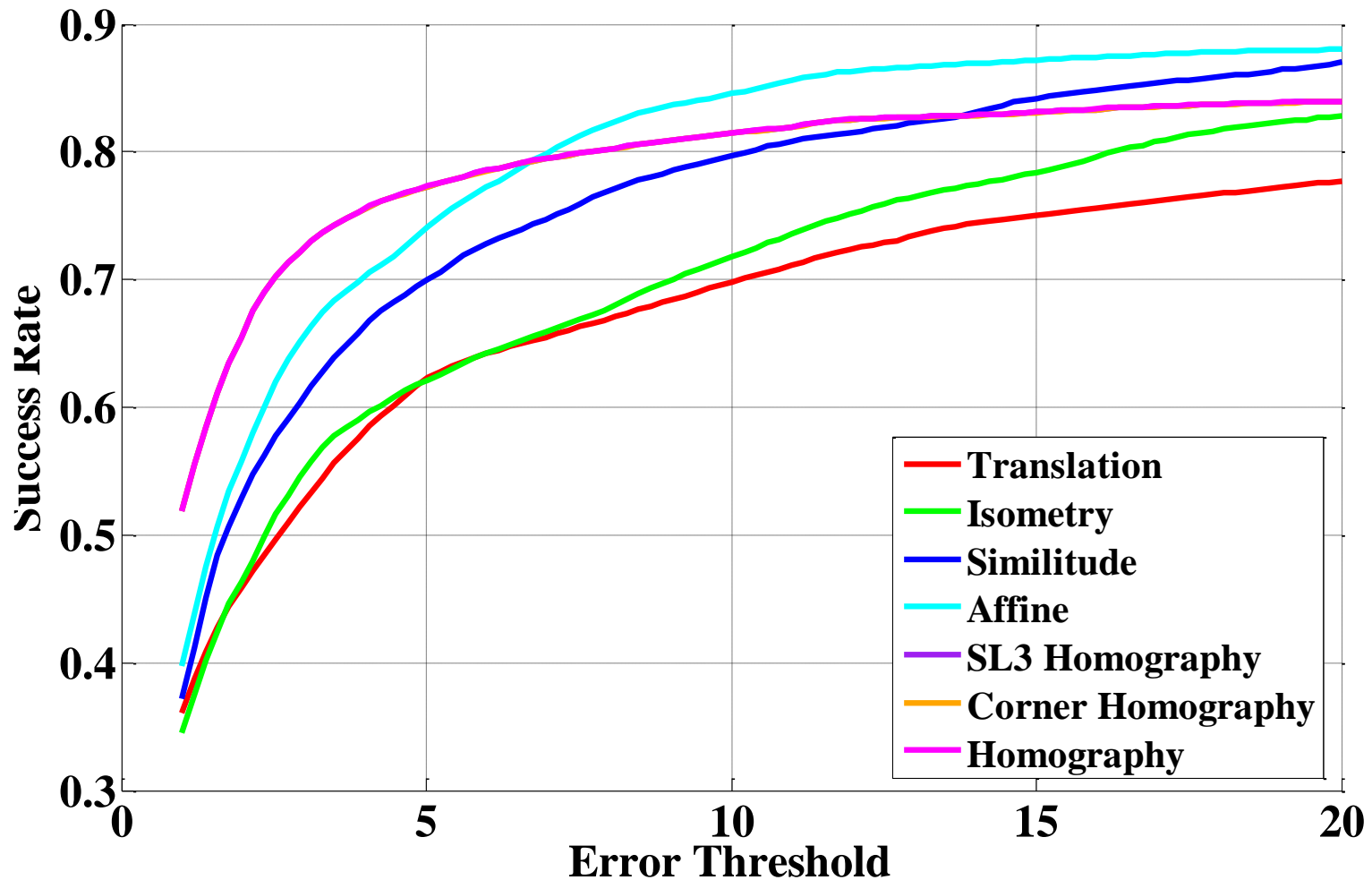
4DOF



2DOF



Results – State Space Models



ESM with ZNCC