Modular Decomposition and Analysis of Registration based Trackers

Abhineet Singh, Ankush Roy, Xi Zhang and Martin Jagersand



Registration based Tracking

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

$$\mathbf{p_t} = \operatorname*{argmax}_{\mathbf{p}} f(\mathbf{I_0}(\mathbf{x}), \mathbf{I_t}(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$

$$\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N], \mathbf{x}_k = [x_k, y_k]^T \in \mathbb{R}^2$$
$$\mathbf{I}(\mathbf{x}) = [I(x_1, y_1), I(x_2, y_2), ..., I(x_N, y_N)]^T \in \mathbb{R}^N, I(x, y) : \mathbb{R}^2 \mapsto \mathbb{R}$$
$$\mathbf{p} = [p_1, p_2, ..., p_S], S : \text{DOF of image motion}$$
$$\mathbf{w} : \mathbb{R}^2 \times \mathbb{R}^S \mapsto \mathbb{R}^2$$
$$f : \mathbb{R}^N \times \mathbb{R}^N \mapsto \mathbb{R}$$

Registration based Tracking

 $\mathbf{p_t} = \operatorname{argmax} f(\mathbf{I_0}(\mathbf{x}), \mathbf{I_t}(\mathbf{w}(\mathbf{x}, \mathbf{p})))$





Motivation

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



<u>MTF Usage Example – Multi Target</u> <u>Tracking</u>

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

B. Lucas, T. Kanade, "An iterative image registration technique with an application to stereo vision", 1981
A. Singh, M. Jagersand, "Modular Tracking Framework: A Unified Approach to Registration based Tracking", 2016, available at: <u>http://webdocs.cs.ualberta.ca/~vis/mtf/</u>

- 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

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

Search Method – Examples (Gradient Descent)

- Variants of Lucas Kanade (LK)^[Baker01] method
 - Forward Additive (FALK)
 - $\Delta \mathbf{p}_{t} = \underset{\Delta \mathbf{p}_{t}}{\operatorname{argmax}} f(\mathbf{I}_{0}(\mathbf{x}), \mathbf{I}_{t}(\mathbf{w}(\mathbf{x}, \mathbf{p}_{t-1} + \Delta \mathbf{p}_{t})))$
 - $\ p_t = p_{t-1} + \Delta p_t$
 - Inverse Additive (IALK)
 - uses constant approximation of ∇I_t computed from I_0
 - Forward Compositional (FCLK)
 - $\Delta \mathbf{p}_{t} = \underset{\Delta \mathbf{p}_{t}}{\operatorname{argmax}} f(\mathbf{I}_{0}(\mathbf{x}), \mathbf{I}_{t}(\mathbf{w}(\mathbf{w}(\mathbf{x}, \Delta \mathbf{p}_{t}), \mathbf{p}_{t-1})))$
 - $p_t = p_{t-1} \circ \Delta p_t$
 - Inverse Compositional (ICLK)
 - $\Delta \mathbf{p}_{t} = \underset{\Delta \mathbf{p}_{t}}{\operatorname{argmax}} f(\mathbf{I}_{0}(\mathbf{w}(\mathbf{x}, \Delta \mathbf{p}_{t})), \mathbf{I}_{t}(\mathbf{w}(\mathbf{x}, \mathbf{p}_{t-1})))$
 - $p_t = p_{t-1} \circ \Delta p_t^{-1}$
- Efficient Second Order Minimization (ESM)^[Benhimane04]
 - combines FCLK and ICLK

S. Baker, I. Matthews, "Equivalence and Efficiency of Image Alignment Algorithms", 2001

S. Benhimane, E. Malis, "Real-time image-based tracking of planes using efficient second-order minimization", 2004

Search Method – Examples (Stochastic)

- Nearest Neighbor Search (NN) [Dick13]
 - generate samples by warping $I_0(x)$
 - find the nearest neighbor to $I_t(w(x,p_{t-1}))$ and update p_{t-1} with the inverse of the corresponding Δp_t
 - combined with ICLK for stability (NNIC)
- Particle Filter (PF) [Kwon14]
 - generate samples by warping $I_t(w(x, p_{t-1}))$
 - compute weight for each and estimate Δp_t as weighted average of samples

9

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]

10

Appearance Model – Examples

• Sum of Squared Differences (SSD)^[Baker01]

$$-f(\mathbf{I_0}, \mathbf{I_t}) = -\frac{1}{2} \parallel \mathbf{I_0} - \mathbf{I_t} \parallel^2$$

• Sum of Conditional Variance (SCV)[Richa11]

$$-f(\mathbf{I_0}, \mathbf{I_t}) = -\frac{1}{2} \parallel E[\mathbf{I_t} | \mathbf{I_0}] - \mathbf{I_t} \parallel^2$$

- Using several joint distributions computed from corresponding sub regions of I_t and I_0 gives a variant called ${\sf LSCV}^{[{\sf Richa14}]}$
- Reversed Sum of Conditional Variance (RSCV)^[Dick13]

$$-f(\mathbf{I_0}, \mathbf{I_t}) = -\frac{1}{2} \| \mathbf{I_0} - E[\mathbf{I_0} | \mathbf{I_t}] \|^2$$

• Zero mean Normalized Cross Correlation (ZNCC)[Ruthotto10]

$$-f(\mathbf{I_0}, \mathbf{I_t}) = -\frac{1}{2} \| \frac{\mathbf{I_0} - \mu_0}{\sigma_0} - \frac{\mathbf{I_t} - \mu_t}{\sigma_t} \|^2$$

R. Richa, R. Sznitman, R. Taylor, G. Hager, "Visual Tracking Using the Sum of Conditional Variance", 2011 R. Richa, et. al, "Direct visual tracking under extreme illumination variations using the sum of conditional variance", 2014 L. Ruthotto, "Mass-preserving registration of medical images", 2010

11

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]

$$-f(\mathbf{I_0},\mathbf{I_t}) = \frac{\mathbf{I_0} - \mu_0}{\sigma_0} \cdot \frac{\mathbf{I_t} - \mu_t}{\sigma_t}$$

A. Dame, E. Marchand, "Accurate Real-time Tracking Using Mutual Information", 2010G. G. Scandaroli, M. Meilland, R. Richa, "Improving NCC-Based Direct Visual Tracking", 2012

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

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

$$-\mathbf{w}(\mathbf{x}_{k},\mathbf{p}) = \begin{bmatrix} (1+p_{1})x_{k} - p_{2}y_{k} + p_{3} \\ (1+p_{1})y_{k} + p_{2}x_{k} + p_{4} \end{bmatrix}$$

State Space Model – Examples (cont'd)

• Homography : S = 8

$$-\mathbf{w}(\mathbf{x}_{k},\mathbf{p}) = \left[\frac{(1+p_{1})x_{k}+p_{2}y_{k}+p_{3}}{(1+p_{7})x_{k}+p_{8}y_{k}+1}, \frac{(1+p_{4})y_{k}+p_{5}x_{k}+p_{6}}{(1+p_{7})x_{k}+p_{8}y_{k}+1}\right]^{T}$$

• **SL3 Homography**^[Benhimane04]: S = 8

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

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

• Corner Homography : S = 8

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

• $\mathbf{G} = \underset{\mathbf{M}}{\operatorname{argmin}} \sum_{i=1}^4 \left\| \mathbf{M} \cdot \begin{bmatrix} c_{ix} \\ c_{iy} \end{bmatrix} - \begin{bmatrix} c_{ix} + p_{2i-1} \\ c_{iy} + p_{2i} \end{bmatrix} \right\|^2$
• $\left\{ c_i = \begin{bmatrix} c_{ix} \\ c_{iy} \end{bmatrix} \middle| 1 \le i \le 4 \right\}$: bounding box corners

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

State Space Model – Examples (Demo)

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

System Design

Evaluation Benchmarks

UCSB

LinTrack

MJ Note CRV16 Results

- Two datasets
 - Tracking for Manipulation Tasks (TMT)^[Roy15]
 - 109 sequences, 70592 frames
 - Visual Tracking Dataset provided by UCSB^[Gauglitz11]
 - 96 sequences, 6889 frames
- Alignment Error (E_{AL}) used to compare tracking result with ground truth

$$-E_{AL} = \frac{1}{4} \left\| C_{track} - C_{gt} \right\|$$

- Success Rate (SR) plot used to measure tracking performance
 - **x** axis : error threshold $t_p \in [0, 20]$
 - **y** axis : fraction of frames with $E_{AL} < t_p$

A. Roy, X. Zhang, N. Wolleb, C. P. Quintero, M. Jagersand, *"Tracking Benchmark and Evaluation for Manipulation Tasks"*, 2015 S. Gauglitz, T. Hollerer, M. Turk. *"Evaluation of interest point detectors and feature descriptors for visual tracking"*, 2011

Note CRV 16 results

Evaluation Methodology - Datasets

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

		Without Subsequences			With Subsequences		
— I IVI I	Dataset	Soguoncos	Total	Trackable	Sub-	Total	Trackable
		Sequences	Frames	Frames	sequences	Frames	Frames
-0CSB	TMT	109	70592	70483	1090	390470	389380
 LinTrack 	UCSB	96	6889	6793	960	41170	40210
	LinTrack	3	12477	12474	30	68700	68670
- ΡΔΜΙ	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})

$$-E_{AL} = \frac{1}{4} \left\| C_{track} - C_{gt} \right\|$$

- 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)

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

MTF vs. ViSP Speeds (FPS)

Results – Learning Based Trackers

ZNCC with **Translation**

Results -Learning Based Trackers (Demo)

ZNCC with Translation

Results – Search Methods

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

Results – Appearance Models

Results – Appearance Models (Demo)

FCLK with Homography

Results - Appearance Models (Demo)

Results – State Space Models

ESM with ZNCC

Results – State Space Models (Demo)

Results - State Space Models (Demo)

- 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: <u>http://webdocs.cs.ualberta.ca/~vis/mtf/</u> along with all datasets and this presentation

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: <u>http://webdocs.cs.ualberta.ca/~vis/mtf/</u> 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 -Learning Based Trackers (Demo)

