# Towards efficient search methods in object tracking: An evaluation and application to precise tracking

by

Ankush Roy

A thesis submitted in partial fulfillment of the requirements for the degree of

Master of Science

Department of Computing Science

University of Alberta

# Abstract

Object tracking is a much researched subject in the computer vision community. With more and more tracking algorithms reported every year, standard benchmarking and evaluation methods are reported for long term tracking systems.

In this thesis we present a public dataset to evaluate trackers used for human and robot manipulation tasks. For these tasks high degrees of freedom (DOF) motion of the object is to be tracked with high accuracy. We describe in detail, both the process of recording the sequences and how ground truth data was generated for the videos. As an initial example, we evaluate the performance of seven published trackers [41, 66, 56, 22, 6, 9, 80] and analyze the results. We describe a new evaluation metric to test sensitivity of trackers to speed. A total of 100 annotated and tagged sequences are reported. All the videos, ground truth data, tagged image frames, original implementation of trackers and evaluation scripts are made publicly available.

We also introduce a new search method in tracking. Sequential Graph based Approximate Nearest Neighbour Search algorithm [69, 32] or SGANNS. It uses overlapping image features in videos to build a connected graph, offline. This graph is then searched efficiently during tracking to predict the best warp parameters. We test this algorithm on the dataset reported and further analyze the results.

Finally we show that using a detection module, registration based trackers can be made more robust. We address tracking challenges of occlussion and varying appearance which a regular registration based tracker fails to track.

*Imagination is more important than knowledge. For knowledge is limited to all we now know and understand, while imagination embraces the entire world, and all there ever will be to know and understand.*

– Albert Einstien

# Acknowledgements

v

# Contents

# List of Tables

# List of Figures

# Chapter 1

# Introduction

An important component of system(s) that rely on estimating the trajectory of an object in a continuous stream of images is *Object Tracking*. Applications like automated surveillance [38], targeting systems [59], robot manipulation [22], augmented reality [48], user interface design [15], image stabilization [78], medical analysis [71], video compression [21] and many others directly benefit from advances in tracking research.

2D Object tracking is challenging for several reasons. One may need to deal with complexities such as noise in images from camera sensors, varying illumination, low texture of the object, loss of information in two dimensional (2D) images of original three dimensional (3D) scenes, occlusion of the object, blur at high speed object motion and complex object motion. An ideal tracker should be both robust to these challenges and precise in predicting the pose (position and orientation, as the case might be) of the object. Two distinct categories are observed in tracking research. (1) Trackers that are very precise [51, 31, 6, 22, 9], that are used in manipulation and servoing tasks, (2) trackers that adapt to the changing appearance of the object over time [41, 66, 56], thus focussing more on robustness and are used in surveillance tasks. The former is referred to as registration based trackers as the core algorithm works around registering the image to the target template to find the best match. The latter is known as Online Learned Trackers (OLT), since the trackers learn and update the appearance of the objects online as tracking goes on.

Figure 1.1 shows the two categories discussed, on a video sequence from [9]

Figure 1.1: [TOP ROW] Example of an Online Learned Tracker, tracking (position, x,y centroid co-ordinates) in TLD [41] [BOTTOM ROW] Precise registration based tracker tracking full pose (4 corners) in IC [6]

[1]. Inverse Compositional tracker (IC) [6] tracks the full pose (precise bounding box represented by the green rectangle) of the object. Tracking Learning Detection (TLD) [41] tracker on the other hand roughly follows the centroid of the object shown by the red dot.

Object tracking is defined as a sequential state prediction problem. State of an object ($\mathbf{p_t}$) at time $t$ is a vector that represents pose (location and orientation) of the object. Based on how precisely a tracker tracks, size of the state vector (DOF - Degrees Of Freedom) varies. Precise trackers in most cases track full pose (8 DOF, [51, 31, 6, 9, 22]) transformation of the object compared to OLT, 3 in [41, 5], 6 in [66, 56]. OLT updates the appearance model of the object, taking into account newer templates, as more appearances are exposed. This process is computationally expensive, hence in-spite of tracking a low DOF state, they are mostly subpar with real time requirements.

In this thesis I focus on precise tracking, trackers that converge within sub-pixel accuracy and can be used in manipulation tasks.

---

[1]Images taken from [9] Available at: http://esm.gforge.inria.fr/ESMdownloads.html

Figure 1.2: Example sequences show the two different sets of challenges that each of the tracking modalities (OLT and Registration based) test in a tracker. IVT [66]) an OLT tracker is shown to track videos [TOP] that has high occlusion, unstructured motion and rapid change of appearance, but they need not follow the object precisely. Registration tracker in ESM [9] [BOTTOM] precisely tracks the high DOF motion of the object.

## 1.1 Motivation

Any algorithm when reported, needs proper evaluation to bring out the full contribution and/or limits of the algorithm. Tracking research being divided, has different challenge sets that each of the sub groups (OLT and registration based) are interested in. This makes homogenous evaluation of tracking difficult. There are no common grounds to evaluate trackers from the two categories, in terms of the videos and evaluation metrics. Figure 1.2 show few example sequences used to evaluate OLT trackers [2] and registration based trackers [3] and how they are different in their focus.

Lately, there has been much focus on evaluating OLT trackers [76, 45, 72], which over a period of time has evolved into making the evaluation process reach a set standard on a wide range of videos. This has been mostly missing (Metaio [47] attempted, but had limitations discussed in Chapter 2) in registration based trackers. Most published methods have been tested only on a

---

[2]Available at: http://cvlab.hanyang.ac.kr/tracker_benchmark/datasets.html
[3]Available at: http://webdocs.cs.ualberta.ca/vis/trackDB/

3

few anecdotal videos. They are however limited (Chapter 2) in the number of challenges they present. To the best of our knowledge Metaio [47] is the only work that reports a structured dataset and how to use it. We show in Chapter 2 how their setup is limited. This motivated us [4] to first benchmark the evaluation process for registration based tracker by recording a diverse dataset and carefully choose the evaluation metrics. This would be relevant for all research groups working on registration based trackers.

One of the difficult problems in the domain of registration based is to track fast object motion maintaining the high precision required. Most search methods [51, 6] assume very little change in object pose, which is critical to their formulation. This restrict their application to a small set of object motion. Travis et al. in [22] showed that with Approximate Nearest Neighbour search using randomized KD-Trees this problem can be handled in a better way than existing techniques. We provide an improvement over this in terms of using the sequential property of video data with a Graph based Approximate Nearest Neighbour search [32] method. We evaluate this method on the dataset recorded, and compare it with existing trackers [6, 9, 22, 41, 66, 56, 80].

## 1.2 Contribution and Thesis Organization

We report a dataset to benchmark 2D object tracking algorithms. Furthermore, a new search method in tracking is reported with detailed results reported and analysed. The salient contributions in this thesis are:

- A publicly available dataset to evaluate registration based trackers. Evaluation scripts are made available.

- A new evaluation metric, that tests a tracker's robustness against large object motions.

- A new search method in registration based tracking that uses the temporal coherency of video data.

_____

[4]In collaboration with Xi Zhang, Nina Wolleb and Camilo Perez

First, we focus on identifying the challenges that a tracker would face when used in a manipulation setup. We mimic the same by recording videos with a human user. The user had no prior knowledge of tracking algorithms, to eliminate any bias while recording the videos. We further augment the dataset with videos of a robot arm (WAM arm with a Barrett hand) performing the same tasks. Robot motion is normally planned with smooth velocity profile and acceleration. Before grasping the object, the robot comes to a smooth stop, then smooth start. Human motions are often less smooth, as a human may grasp the object on the fly, causing the tracked object to accelerate quickly. However with the robot, there is the problem of partial occlusion from the end effector.

As an initial example we evaluate and rank eight trackers, [6, 9, 22, 66, 56, 41, 80] and a tracker we develop (Chapter 3). We experimentally show, that OLT trackers are not generic enough to be used in precise target tracking applications. Here we introduce a new evaluation metric that evaluates trackers on how they perform against large object motion. We argue that previous methods [6, 9, 22] that simulated large motions were not fully accurate. Chapter 5 has detailed discussions on this.

We also introduce a new search method [32] in tracking. This builds on the idea of using Approximate Nearest Neighbour search in tracking, proposed in [22]. It uses the sequential coherency of video data to build a graph of spatially related image templates. This connected graph makes search faster, as the number of computations to predict pose at time $t$ decreases with a smart choice of starting position, based on pose at time $t-1$. Results and analysis of the algorithm is done in Chapter 5.

Finally, we address the problem of occlusion in videos. This is a pressing issue when objects becomes occluded by the end effector of a robot or any other object. We use a fast detection system coupled with registration based tracker. The detection system reinitializes the object tracker whenever it drifts substantially and using a registration based tracker maintains the high convergence. The content of this thesis is organized as follows.

- Chapter 2 lists existing tracking methods in literature. The second half of the Chapter provide a detailed summary of existing datasets and benchmarking methods.

- Chapter 3 discusses how a tracking algorithm can be broken down into subsequent components and individual components studied. Here we present a new search method that can be applied to registration based tracking and show how it's suited for tracking application, given the sequential nature of a video sequence.

- Chapter 4 describes in detail our setup for recording the videos and how we generate ground truth data for all videos. Furthermore, we describe our error metric and evaluation methods, highlighting their relevance in evaluating a registration tracker.

- Chapter 5 presents the evaluation results of the tracker we develop in light of the dataset reported. Finally, we present a global ranking of all trackers we study.

- Chapter 6 shows an application of using a simple detection technique on top of the registration trackers to make trackers robust to occlusion without compromising convergence.

- Chapter 7 concludes the work highlighting the significant contributions in the thesis and how it can be further extended to cover more challenges in the dataset.

# Chapter 2

# Background

In this chapter we discuss different existing tracking algorithms. We summarize years of tracking literature and group them into different categories. The second half of this Chapter deals with details of existing benchmarking techniques for Object Trackers, their usability and limitations to evaluate registration based trackers.

## 2.1 Tracking Algorithms

Visual tracking is the process of repeated estimation of the state of an object (position and orientation) in an image, given state(s) in previous frame(s). Approaches from different domains of computer vison have been adapted to track objects. Some of the existing works can be grouped into four major categories which are (i) **Feature Based**, (ii) **Segmentation Based**, (iii) **Detection Based** and (iv) **Registration Based**.

### 2.1.1 Feature Based

In feature based tracker, the object to be tracked is represented as a collection of features. Some of the popular features used in tracking are, transform features in intensity [51, 6, 9], color [55, 18, 55], steerable filter responses in [39] (Figure 2.1 (a)) and [44], optical flow features in [7], eigen features [75, 11], integral histogram in [2]. Invariant features like SIFT [50] and SURF [36] (Figure 2.1 (b)) have also been shown to work in tracking. Sometimes, a combination of these features is also utilized to improve tracking performance.

(a)



(b)

Figure 2.1: (a) Tracking is done using filter responses from a steerable pyramid using mixture model to fit data. Image taken from [39] show the tracked region and the model's mixing probability, mean, and ownership (left to right). (b) Tracking using SURF features. Images taken from [36] show the extracted features on the image frames and how the tracker follow the object in the video.

### 2.1.2 Segmentation Based

Segmentation algorithms partition the image into similar regions. This can be adapted to be used in tracking, by segmenting the object from the background. Naturally, object clustering approaches [81] translate well into video tracking. Approaches such as, mean shift [17] (Figure 2.2 (c)), normalised cross cut [77] (Figure 2.2 (a)) active contour [64] (Figure 2.2 (b)) and [79, 10] have been shown to work well in tracking. Segmentation based trackers in most cases provide pixel level accuracy [42] of the moving object by marking the exact pixel co-ordinates of the object boundary.

(a)



(b)



(c)

Figure 2.2: (a) Leukocytes are tracked in [77] (b) Snake based tracking as shown in [64]. The white boundary is the starting position of the snake where it evolves and the final black boundary is the new snake in the current frame (c) Mean Shift tracker in [17] is shown tracking pedestrians in a sub way system

## 2.1.3 Detection Based



(a)



(b)

Figure 2.3: (a) Tracking using Support Vector Machine is done in [4]. The best fitting bounding box is searched close to the last known co-ordinates which eventually converge as shown in the right image. (b) Detected features are used to vote for the object position in [29]. The blue lines shows the suppotring features used to predict position.

Detection based trackers, use an object detection module together with a tracker to aid long term robust tracking. The detection module localises the object when the tracker fails to converge [5, 41] and subsequently reinitialises the tracker. Another way to use detection, is to use frame to frame detection [54]. Discriminative classifiers in SVM [4] (Figure 2.3 (a)), graph based pruning [33], kernelized structured output support vector machine [34], online boosting [28, 29] (Figure 2.3 (b)) are used to give a binary decision of object/background are such examples. Such systems rely on having object models that are either static or dynamic. Detection is done based on this model.

## 2.1.4 Registration Based

First introduced by Lucas and Kanade in [51], registration based trackers use image alignment, to find the best possible warp parameters, that when applied

to the object image ($\mathbf{I}$) closely resembles the template ($\mathbf{I}^*$). Gradient descent is the most commonly used method to acheive this among others like difference decomposition [27] (Figure 2.4 (a)) and linear regression [19]. Different image comparision metrics in sum of squared difference [51, 6], mutual information [20, 23], sum of conditional variance in [65] (Figure 2.4 (b)) and [22], normalised cross correlation [77] is used to compare the object image with the original template.



(a)



(b)

Figure 2.4: (a) Piecewise projective transformation is done in [27]. This allows the algorithm to track flexible objects as shown in the figure above (b) Sum of Conditional Variance (SCV) is used with ESM tracker in [65]. SCV corrects for illumination variation which allows it to track objects even in the presence of specular reflection.

In this thesis we focus mainly on registration based trackers. We show how they are more suitable to use in fine alignment tasks compared to some of the other trackers. To do this, we need to define an evaluation method and record videos that simulate fine alignment tasks.

## 2.2   Dataset and Evaluation

An important part of tracking research is to evaluate trackers using a standard set of videos that are carefully recorded to present a wide spectrum of challenges. Some of the existing datasets in literature are summarised in Table

2.1.

Table 2.1: Comparision of Different Object Tracking Datasets

| Dataset | TR | RO | IL | PR | SR | SC | OC | TX | DOF |
|---|---|---|---|---|---|---|---|---|---|
| Metaio [47] | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | 8 |
| Static [22, 6, 22] | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | 8 |
| Zimmermann [82] | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | 6 |
| OLT Benchmark [76] | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | 2 |
| Gauglitz [26] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | 8 |
| ESM [9] | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | 8 |
| BoBot [43] | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | 2 |
| VOT Challenge [45] | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | 2 |
| MOT Challenge [46] | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | 2 |
| [72] | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | 2 |
| KTH Dataset [3] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 6 |
| TMT [67] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 8 |

TR = Translation; RO = Rotation; IL = Illumination; PR = Perspective; SR = Specular; OC = Occlusion; TX = Texture; DOF = Degrees of Freedom

## 2.2.1 Anecdotal Video Datasets

Previous works either use synthetically generated data by applying random warps [40, 6] to a template followed by tracking it, or test algorithms on a small set of videos that are either recorded or pooled from the Internet. Zimmerman et al. [82] tested their algorithm on 3 grayscale videos (linTrack) that covers 6 DOF transformation of the object. They manually annotate the bounding box using either a special marker or visual texture based cues. [49] and [62]

showed application of template tracking in augmented reality, but restricted themselves to a small test set when evaluating their tracking system.

Other popular sources of videos in tracking literature are by Babenko et al. in [5] [1], Ross et al. in [66] [2], Kalal et al. [41] [3], Collins et al. in PETS2005 [16] [4], Klein et al. in [43] [5] and Fisher et al. in CAVIAR [25] [6]. These datasets have sequences to test tracking algorithms that are mainly developed for surveillance applications were loosely following ($> 50\%$ area overlap of the target with ground truth) the object is good enough for tracking to succeed. The main focus is on testing trackers for occlusion, out of frame motion and large appearance variations.



Figure 2.5: Sequence showing both the robot and human user performing identical tasks of pouring cereal in a bowl under normal light settings. The red rectangle shows the tracking result on the sequence using ESM [9]

---

[1] Available at: `http://vision.ucsd.edu/~bbabenko/project_miltrack.html`

[2] Available at: `http://www.cs.toronto.edu/~dross/ivt/`

[3] Available at: `http://personal.ee.surrey.ac.uk/Personal/Z.Kalal/TLD/TLD_dataset.ZIP`

[4] Available at: `http://vision.cse.psu.edu/data/vividEval/main.html`

[5] Available at: `http://www.iai.uni-bonn.de/~kleind/tracking/`

[6] Availabe at: `http://homepages.inf.ed.ac.uk/rbf/CAVIAR/`

### 2.2.2 Full Video Datasets

PETS2005 [16] made initial efforts of creating a comprehensive dataset having wide range of challenges, but as tracking research evolved it needed more diverse set of videos to test trackers. In a recent paper Yang et al. 2013 [76] [7] categorized some of the above mentioned video sequences to publish a common benchmark [8]. Any new tracker can be evaluated and compared with other state of the art trackers in the literature using this benchmark. VOT [9] challenge [45] has a running catalog of videos and evaluation metrics by which they benchmark any new tracker. They reported an unique way to globally rank trackers based on how each tracker perform on individual metrices of rubustness, accuracy etc. Each frame in VOT is labelled with the challenges it presents. This makes it easy for an user to asses performance of the tracker, which might be suseptible to some specific challenges. MOT [46] host a similar competition for multiple object tracking in association with Winter conference on Application of Computer Vision (WACV). [45] proposes an unique way of globally ranking trackers. However, these video sequences are more suited to surveillance tracking. Furthermore, we elaborate why the error metrics used and the evaluation thresholds set [76, 45] are too coarse and ill suited for manipulation tasks.

Table 2.1 shows how diverse the popular datasets are in the set of challenges they provide. It also brings out how little attention has been paid to benchmark registration based trackers with recent works mostly focussing on evaluating Online Learned Trackers (OLT).

Closest to our aim use to be Metaio dataset [10] meant to evaluate planar homography tracking by Lieberknecht et al. [47]. They collected videos under various motions, illuminations and textures. They used a camera mounted on a Faro arm (a passive robot arm with very high precision joint encoders) that was calibrated and all transformations of pose were stored. The stored

---

[7]Available at: https://sites.google.com/site/trackerbenchmark/benchmarks/v10
[8]Available at: http://cvlab.hanyang.ac.kr/tracker_benchmark/datasets.html
[9]Available at: http://www.votchallenge.net
[10]Available at: http://www.metaio.com/research/

values were used to calculate ground truth data. Instead of a 3D scene the Metaio benchmark like [57], records a printed (2D planar) poster, making the benchmark somewhat artificial. Furthermore, the setup with a moving camera on an arm means that the motions are not from natural tasks and were restricted by the arm workspace and mass. Metaio, haven't released the ground truth data which makes the dataset un-useable now that it's offline. An extensive dataset for 3D tracking is reported in [3]. [11] It reports 6 DOF ground turth, 3 degrees in rotation and 3 in translation.

To this effect, we report a public video dataset [12] to evaluate $2D$ marker-less single object tracking algorithms. The tasks are natural table top manipulations, frequently performed in our daily life. We provide videos and annotated ground truth of both a human and a robot arm performing the same tasks. All videos are tagged with the task performed and the challenges that a tracker would face while tracking. Using these tags, susceptibility of the tracking algorithm can be properly narrowed down and subsequently improved. Each task is repeated with different speeds and under two different lighting conditions, of normal office light and diffused light. Together with the videos, we present detailed analysis of seven trackers. We also report a new metric that is suited for evaluating registration based trackers, and finally rank all the trackers.

---

[11]Available at: http://robocoffee.org/datasets/
[12]Available at: http://webdocs.cs.ualberta.ca/~vis/trackDB/

# Chapter 3

# Modular Decomposition of Tracking

A video stream is modeled as a temporal sequence of images, $\mathbf{F} = \{\mathbf{I_0}, \mathbf{I_1}, ..., \mathbf{I_n} : \mathbb{R}^2\}$. The pixel locations in an image patch with $N$ pixels are denoted by $\mathbf{x} = (\mathbf{x_1}, \mathbf{x_2}, ..., \mathbf{x_N})^T$ where $\mathbf{x_k} = (x_k, y_k)$ are the corresponding pixel intensities in image, denoted by $\mathbf{I(x)}$. When a geometric transform $\mathbf{W}$ with parameters $\mathbf{p} = (p_1, p_2, ..., p_l)$ is applied to an image patch $\mathbf{x}$, the transformed patch is denoted by $\mathbf{x'} = \mathbf{W(x; p)}$ and the corresponding pixel values transform to $\mathbf{I(W(x; p))}$. Object tracking can then be formulated (Eq 3.1) as a state prediction problem where we need to predict the optimal transform parameters $\mathbf{p_t}$ for an image $\mathbf{I_t}$ that minimizes the difference, measured by a suitable distance metric $\mathbf{d}$, between the target patch $\mathbf{I^*} = \mathbf{I_0(W(x; p_0))}$ and the warped image patch $\mathbf{I_t(W(x; p_t))}$. Common choices of distance functions are Eucledian distance ($l_2$) [51] and Manhattan distance ($l_1$) [13].

$$\mathbf{p_t} = \underset{\mathbf{p}}{\text{argmin }} \mathbf{d}(I(\mathbf{W(x; p)}), \mathbf{I^*}) \tag{3.1}$$

This can be formulated more clearly by exemplifying how higher dimensional warps (affine and homography), larger value of $l$, both give more precise state information (important in e.g. robot manipulation) and allows a better template alignment and thus lower residual error in the search (improves convergence of many methods e.g. Gauss Newton). This determines the state space for search. The algorithm that minimizes Eq 3.1 is the search method.

In registration based tracking, the three submodules (i) Appearance ($\mathbf{d}$),

(ii) Search Method and (iii) State Space ($\mathbf{p}$) combine together as shown in Figure 3.1. This process is repeated over several iterations till the image ($\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}))$) is similar to the template ($\mathbf{I}^*$). This assumes static appearance models wherein no updates are done to the appearance model with newer templates.



Figure 3.1: Modular breakdown of a tracking algorithm assuming there is no dynamic update to the appearance model. This shows how different components work. The state vector $\mathbf{p}$ parameterizes the pose of the object.

## 3.1 Appearance Model

Appearance model is the transform space (generated by the function $\mathbf{d}$) wherein the search method compares different warped patches from the current image to get the closest match with the original image template $\mathbf{I}^*$.

### 3.1.1 Sum of Squared Difference (SSD)

Sum of squared difference ($\mathcal{SSD}$) is a common measure [51, 6, 9, 22] where raw pixel intensities are compared. $\mathcal{SSD}$ is expressed as an Eucleadian ($l_2$)

norm in the image space.

$$\mathcal{SSD} = ||\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})) - \mathbf{I}^*||_2 \tag{3.2}$$

$\mathcal{SSD}$ is susceptible to lighting variations, as it does not correct for illumination bias. Sometimes smoothing of the image is done prior to tracking to remove the effect of local noise.

## 3.1.2 Sum of Conditional Variance (SCV)

Sum of Conditional Variance ($\mathcal{SCV}$) was originally used in medical imaging [63]. The authors show how $\mathcal{SCV}$ is non-invariant to non-linear illumination variations. This motivated researchers to use it in object tracking [22, 65], where changing appearance of the object is a common challenge. $\mathcal{SCV}$ corrects for illumination changes by taking into account the original template ($\mathbf{I}^*$) as reference.

$$\mathcal{SCV} = \sum_{\mathbf{x}}(\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})) - \mathbb{E}[\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}))|\mathbf{I}^*])^2 \tag{3.3}$$

$\mathbb{E}$ is the Expectation Operator. We use the same formulation of $\mathcal{SCV}$ as done in [65, 22]. $\mathbb{E}[\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}))|\mathbf{I}^*]$ is computed from the joint intensity distribution between $\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}))$ and $\mathbf{I}^*$.

## 3.1.3 Normalised Cross Correlation (NCC)

Normalized Cross Correlation ($\mathcal{NCC}$) is another measure that accounts for illumination changes by centering and normalising the image data before comparing the two images.

$$\mathcal{NCC} = \frac{\sum_{\mathbf{x}}(\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})) - \bar{\mathbf{I}})(\mathbf{I}^*(\mathbf{x}) - \bar{\mathbf{I}}^*)}{\sqrt{\sum_{\mathbf{x}}(\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})) - \bar{\mathbf{I}})^2}\sqrt{\sum_{\mathbf{x}}(\mathbf{I}^*(\mathbf{x}) - \bar{\mathbf{I}}^*)^2}} \tag{3.4}$$

It has further been shown in [68] that maximizing $\mathcal{NCC}$ between two images $\mathbf{I}^*$ and $\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}))$ is equivalent to minimizing the $\mathcal{SSD}$ score between two Z-score [37] normalized images, which makes it easier to implement.

### 3.1.4 Illumination Invariant Appearance (IVA)

Finally, we use another Illumination Invariant Appearance ($\mathcal{IVA}$) model in [60]. Colour image is decomposed into the consitituent channels of Red ($\mathbf{R}$), Green ($\mathbf{G}$) and Blue ($\mathbf{B}$). This independent channels are then used to re-compute the illumination invariant image using two coefficients, $\alpha$ for Blue channel and $\beta$ for Red channel. For a 2D image the transformed $\mathcal{IVA}$ image is calculated as shown in Equation 3.5.

$$\mathcal{IVA} = \ln(\mathbf{G}[i,j]) - \alpha * \ln(\mathbf{B}[i,j]) - \beta * \ln(\mathbf{R}[i,j]) \tag{3.5}$$

The parameters $\alpha$ and $\beta$ are subjected to $\frac{1}{\lambda_1} = \frac{\alpha}{\lambda_2} + \frac{\beta}{\lambda_3}$, where $\lambda_1, \lambda_2, \lambda_3$ are the peak sensitivity (wavelength) for each sensor. Parameters $\alpha$ and $\beta$ for the camera sensors we use (SONY ICX274), are 0.482 and 0.518 respectively.



| (a) | (b) | (c) | (d) | (e) |

Figure 3.2: Example of different transforms, (a) Original Image (b) Original Image smoothed using a Gaussian Filter of window $5 \times 5$ and $\sigma = 3$ (c) Reverse Sum of Conditional Variance (d) Normalised Cross Correlation and (e) Illumination Invariant Appearance

## 3.2 Search Method

Search step finds the best transform parameters (state) $\mathbf{p}$ at time $t$, that minimize the difference between the target image ($\boldsymbol{I}^* = I(\mathbf{W}(\mathbf{x}; \mathbf{p}_0))$) and the transformed template $I(\mathbf{W}(\mathbf{x}; \mathbf{p}))$. $\mathbf{W}$ is the transformation function and $\mathbf{p}$ the transform parameters. The difference is expressed as a distance score, $\mathbf{d}$ as shown in Eq 3.6.

$$\mathbf{p_t} = \underset{\mathbf{p}}{\operatorname{argmin}} \ \mathbf{d}(\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})), \mathbf{I}^*) \tag{3.6}$$

Eucleadian distance or $l_2$ norm is a popular distance function [51, 6, 12]. The corresponding optimization is formulated as Eq 3.7.

$$\mathbf{p_t} = \underset{\mathbf{p}}{\operatorname{argmin}} \ ||\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})) - \mathbf{I}^*||^2 \tag{3.7}$$

Gradient based search is the most popular [51, 6, 9, 22] method used for optimization. Recent works include use of supervised learning [22]. In the following sub-sections, I elaborate these methods and describe how improvements can be done.

## 3.2.1 Gradient based search

One way to solve the above non-linear optimization problem is to iteratively update $\mathbf{p}$ till it converges to $\mathbf{p}^*$. Typically, iterations continue till some norm of the update ($\triangle \mathbf{p}$) is below a constant $\epsilon$ ($\triangle \mathbf{p} \le \epsilon$). Lucas and Kanade [51] approximated the objective function ($||\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})) - \mathbf{I}^*||^2$), using a first order expansion of taylor series. This gives a closed form solution of the update $\triangle \mathbf{p} = \mathbf{H}^{-1} \sum_{\mathbf{x}} [\nabla \mathbf{I} \frac{\delta \mathbf{W}}{\delta \mathbf{p}}]^T [I^* - \mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}))]$. $\mathbf{H}$, the Hessian matrix, is computed as $\mathbf{H} = \sum_{\mathbf{x}} [\nabla \mathbf{I} \frac{\delta \mathbf{W}}{\delta \mathbf{p}}]^T [\nabla \mathbf{I} \frac{\delta \mathbf{W}}{\delta \mathbf{p}}] + \sum_{\mathbf{x}} r_x R_x$. Here $r_x$ is the residual and $R_x$ the second order differential of the residual. Residuals are normamlly very low, and almost zero near the optimum. This makes the second term negligeble. Hence an approximate Hessian ($\mathbf{H} \approx \sum_{\mathbf{x}} [\nabla \mathbf{I} \frac{\delta \mathbf{W}}{\delta \mathbf{p}}]^T [\nabla \mathbf{I} \frac{\delta \mathbf{W}}{\delta \mathbf{p}}]$) is used [58]. The authors [51] used an additive update, $\mathbf{p_t^k} \leftarrow \mathbf{p_t^k} + \triangle \mathbf{p_t^k}$, for $k^{th}$ iteration. $t$ is dropped since the iteration is to find the best parameter for the same time instance. With an assumption that on each iteration the template image ($\mathbf{I}^*$) can be approximated very closely by the warped image $I(\mathbf{W}(\mathbf{x}; \mathbf{p}))$ [31], avoided the computationally expensive step of calculating $\mathbf{H}$ in every iteration. There warp update was $\mathbf{p^k} \leftarrow \mathbf{p^k} - \triangle \mathbf{p^k}$. The general form of additive parameter update is shown in Eq. 3.8

$$\mathbf{p}^* = \underset{\mathbf{p}}{\operatorname{argmin}} \ ||\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p} + \triangle \mathbf{p})) - \mathbf{I}^*||^2 \tag{3.8}$$

The formulation in [31] assumes the image intensities are linearly dependent on the motion vector, i.e. they are consistent, to solve the over determined system. This assumption doesn't always hold true for real images. To avoid this, Jurie and Dhome [40] learnt the jacobian by fitting hyperplanes. This was done by randomly sampling data in the neighbourhood of $\mathbf{p}$

Instead of using an additive update, Baker et. al [6] used an inverse compositional update, $\mathbf{W}(\mathbf{x}; \mathbf{p^k}) \leftarrow \mathbf{W}(\mathbf{x}; \mathbf{p^k}) \circ \mathbf{W}(\mathbf{x}; \triangle\mathbf{p^k})^{-1}$. The resultant minimization is as shown in Eq 3.9.

$$\mathbf{p}^* = \underset{\mathbf{p}}{\operatorname{argmin}} \ ||\mathbf{I}(\mathbf{W}(\mathbf{W}(\mathbf{x}; \triangle\mathbf{p}); \mathbf{p})) - \boldsymbol{I}^*||^2 \tag{3.9}$$

The parameter update ($\triangle\mathbf{p}$) is expressed as, ($\triangle\mathbf{p} = \mathbf{H}^{-1}\sum_{\mathbf{x}}[\nabla\mathbf{T}\frac{\delta\mathbf{W}}{\delta p}]^T[\mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p})) - \mathbf{I}^*]$. Hessian matrix ($\mathbf{H}$) is calculates once before tracking commences. $\mathbf{H} = \sum_{\mathbf{x}}[\nabla\mathbf{I}^*\frac{\delta\mathbf{W}}{\delta\mathbf{p}}]^T[\nabla\mathbf{I}^*\frac{\delta\mathbf{W}}{\delta\mathbf{p}}]$. This speeds up the optimization since both the gradient ($\nabla\mathbf{I}^*$) and the Jacobian ($\frac{\delta\mathbf{W}}{\delta\mathbf{p}}$) can be pre computed.

Benhimane and Mallis [9] presented a compositional update scheme for the warp parameter $\mathbf{p}$. They use an Efficient Second order Minimization (ESM) [53] approach. The Hessian matrix ($\mathbf{H}$) is approximated using two Jacobians. One of which is computed once, when tracking starts (on the template, $\mathbf{J}(\mathbf{e})$). However, the other Jacobian ($\mathbf{J}(\mathbf{x_c})$), needs to be updated on every iteration. This though comparitively slow, on account of computing the Jacobian on every iteration has it's own benefits. The authors [9] show that the algorithm converges using lesser number of iterations and also avoids the problem of not converging unless search starts close to $\mathbf{p}^*$.

### 3.2.2 Nearest neighbour based search

As an alternative to gradient descent search, Dick et. al [22] in a recent study (NNIC) used approximate nearest neighbor (ANN) search to find the best matching warp from a large list of pre-computed warps. They used randomized KD-trees [8] in Flann [1]. KD-Trees was used previously in [30] for tracking, however the authors used it as an object detector.

---

[1] Available at: `https://github.com/mariusmuja/flann`

**Randomised KD-Tree Based**

NNIC [22] use a cascade tracking system, with a coarse tracker in ANN (randomised KD-Tree) and fine alignment done using IC [6]. The authors randomly sample warps from a Gaussian distribution $\mathbf{p_i}$. Warped template images are generated as shown in Algorithm 1. The warped images are then used to build KD-Tree offline before tracking starts. While tracking, the best warp ($\mathbf{p}^*$) is searched using randomised KD-Tree search [8]. This warp $\mathbf{p}^*$ is the corresponding warp for which $\mathbf{I}(\mathbf{W}(x; \mathbf{p}))$ is the most similar to the Object Template ($\mathbf{I}^*$). Update to the parameters are done as, $\mathbf{p} \leftarrow \mathbf{p} \circ \mathbf{p}^{*-1}$. Motivated by how intuitive this algorithm is, we adapt a Sequential Graph Based Approximate Nearest Neighbour Search (SGANNS) [32] in tracking in GNNIC.

---

**Algorithm 1** Pre-Computation of Sampled Images

---

1: **procedure** WARPIMAGE($I^*, \sigma$)                    $\triangleright I^*$ Template Image
2:     $warpedData = dict\{\}$
3:     $N =$ No of Samples
4:     **while** $I < N$ **do**
5:         $\mathbf{p} = sample\ \mathcal{N}(0, \sigma)$
6:         $\hat{I} = \mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}))$
7:         $warpedData[\mathbf{p}] = \hat{I}$
8:         $I = I + 1$
9:     **end while**
10:     **return** $warpedData$
11: **end procedure**

---

**Sequential Graph based Approximate Nearest Neighbour Based**

Randomised KD-Tree based approach doesn't take into account the temporal coherency of video data, an assumption that is implicit in videos. The original graph based approximate nearest neighbour algorithm proposed in [69], though similar in performance to KD-Tree, also starts search from a random node. The authors in [32] showed how the temporal coherency of data can be used to speed up search by having an intelligent guess of the start node. We use the same pre-computed sampled images from Algorithm 1 to construct the graph structure needed in GNNIC.

**Algorithm 2** Construction of the graph **G**

---

1: **procedure** BUILD GRAPH(**warpedData, k**)  $\qquad$ ▷ k = Number of
   Neighbours
2: $\quad$ **G** = {}
3: $\quad$ **while** $I < len(\textbf{warpedData})$ **do**
4: $\qquad$ $Q = \textbf{warpedData}[I]$
5: $\qquad$ $(indexes) = knnsearch(Q, \textbf{warpedData, k})$
6: $\qquad$ $\textbf{G}[Q] = indexes$
7: $\qquad$ $I = I + 1$
8: $\quad$ **end while**
9: $\quad$ **return G**
10: **end procedure**

---

Each image template corresponds to a node on the graph **G**. The graph **G** is built by each node being connected with **k** nearest nodes. We explore both $l_2$ and $l_1$ distances while finding nearest nodes. The parameter **k** decides sparsity of the resultant graph **G**. Higher the value of **k** the more connected the graph is. We set **k** to 400, to have a fairly connected graph. *knnsearch* which calculates the $k$ nearest neighbours, is implemented using quick-select [52], that has a computational complexity of $O(n)$, where $n$ is the number of samples. The cost of constructing the graph **G** is $O(dn^2)$, where $d$ is the resolution of the template ($d = (m \times n)$). Search during tracking is done as explained in Algortihm 3.

When tracking starts, at every iteration, the algorithm searches for the best **E** possible matches within the **k** connected nodes. $N(\textbf{Y}, \textbf{E}, \textbf{G})$ searches for **E** nearest neighbour of **Y** in the connected graph **G** where $\textbf{E} \leq \textbf{k}$. This is similar to *knnsearch*, but within the **k** connected child nodes of a node, **Y**. The current node is replaced by the neighbor that is closest to the query and each of these **E** child nodes are then further explored. This exploration continues for $T$ iterations to find the best match. Termination can also be done when residual of the best match with the template is less than a set thershold. Note the difference between **K** the number of nearest neighbours to be returned and **k**, degree of connectivity of the graph. In tracking we use **K** = 1, that returns the closest match. For the first frame, the algorithm initialises from a

---
**Algorithm 3** SGANNS Search
---
1: **procedure** SGANNS($\mathbf{G}, Q, T, E$)   ▷ T= No of greedy steps, E = No of possibilities
2:    $\mathbf{S} = \{\}$
3:    $\mathbf{U} = \{\}$
4:    $Y_0$: a point drawn randomly from $\mathbf{G}$.
5:    **while** $I < T$ **do**
6:       $Y_t = \text{argmin}_{Y \in N(Y_{t-1}, E, \mathbf{S})}(Y, Q)$
7:       $\mathbf{S} = \mathbf{S} \cup N(Y_{t-1}, E, \mathbf{S})$
8:       $\mathbf{U} = \mathbf{U} \cup \{\mathbf{d}(Y, Q) : Y \in N(Y_{t-1}, E, \mathbf{S})\}$
9:       $I = I + 1$
10:   **end while**
11: Sort $\mathbf{U}$, pick the first K elements, and return the corresponding elements in $\mathbf{S}$
12: **end procedure**
---

random node. From there on, the best node of the last search is used as the starting point. Hence the name Sequential Graph based Approximate Nearest Neighbour Search (SGANNS).

This allows us to remove the random restart loop in TABLE I [32]. The assumption here is, features are repeatable in consequitive frames, since video being a temporal sequence there is lot of overlap between two frames. We use the same setup, by using SGANNS with a cascade of IC [6] to further improve alignment. Detailed analysis and results of this tracker GNNIC is shown in Chapter 5.

## 3.3   State Space Model

The state space model embodies any constraints that are placed on the search space of warping parameters to make the search process more efficient. This includes both the DOF of allowed motion, as well as the actual parameterization of the warping function. For instance, the ESM tracker [9] can be considered to search over a different parameterization of the state space from conventional Lucas Kanade type trackers [51, 6] since it uses an $\mathcal{SL}3$ parameterization of homography rather than the actual entries of the corresponding matrix as used by the others.

The advantage of allowing higher DOF in the warping function is achieving greater precision in the aligned warp since transforms that are higher up in the hierarchy [35, sec. 2.4] can better approximate the projective transformation process that captures the relative motion between the camera and the object in the 3D world into the 2D images. However, there are two issues with having to estimate more free parameters. Firstly, the iterative search takes longer to converge making the tracker slower. Secondly, the search process becomes more likely to either diverge or end up in a local minimum thus causing the tracker to be less stable and more likely to lose track. The latter is a well known phenomenon with Lucas Kanade type trackers [14] whose higher DOF (e.g. affine) variants are much less robust than simple 2 DOF versions. Since we focus on high DOF trackers we restrict the study to full pose 8 DOF trackers.

# Chapter 4

# Tracking and Manipulation Tasks Dataset and Evaluation

Designing tracking systems raises the question of how to best evaluate them. Tracking research being divided into two focus group, it is (i) too varied to claim generality, (ii) each interesting enough to have its own set of evaluation standards. A longstanding problem of evaluating registration based trackers is unavailability of a standard dataset and evaluation methodology. We report a **Tracking and Manipulation Tasks (TMT)** [67] dataset to address this.

**TMT** is both a video dataset and an evaluation benchmark. The dataset consists of more than **100** videos of manipulation tasks carefully recorded by a human user and a robot arm. The tasks are chosen to mimic routine tasks one would perform on a daily basis. The videos are structured to systematically evaluate a tracker with evaluation scripts, challenge tags, sample results and ground truth data are made public. [1]

## 4.1   Video Capture Set Up

All videos (Fig  4.3) were recorded using a GRAS-20S4C-C fire-wire camera equipped with a Kowa LM6NCM F1.2/6mm lens. For the human recorded videos the camera stood on a fixed position on a tabletop 90 cm from the edge of the table. The human user had no involvement in designing tracking algorithm. This removes bias in the way the motion was planned while ma-

---

[1]Available at: `http://webdocs.cs.ualberta.ca/~vis/trackDB/`

Table 4.1: Parameter Setting for Video Capture

| Parameters | Diffuse Light | Normal Light |
|---|---|---|
| Exposure (in IL) | 0.73 | 1.06 |
| Gain (in dB) | 6.02 | 0 |
| Shutter (in s) | 0.04 | 0.07 |
| White Balance | Blue/U927 Red/V493 | Blue/U757 Red/V490 |
| Saturation (in %) | 95.22 | 119.04 |

nipulating the objects. Camera settings were implemented through coriander 2.0.1 [61]. Lighting was varied between normal office lighting and diffused light by placing an umbrella to avoid direct source. The videos were recorded at 30 frames per second (F.P.S.) at a resolution of 600x800 in YUV colour space. Exact parameters of recording can be found in Table 4.1.

For the robot recorded videos, a 7 DOF WAM arm [1] with a Barrett hand was used. Fig. 4.1 shows the set up. The camera was fixed on a tripod stand at a distance of 120 cm from the WAM arm. Other parameters of exposure, gain, shutter, white balance and saturation were kept same as in Table 4.1.



Figure 4.1: Set up for recording videos under normal light using a WAM arm [1] and a Barret hand. The camera is set 120 cm from the base of the robot arm.

### 4.1.1 Description of Videos

A tracking algorithm is considered to be robust if it successfully tracks in the presence of "*Translation (TR)*", "*Rotation (RO)*", "*Illumination (IL)*", "*Perspective (PR)*", "*Specularity (SR)*", "*Occlusion (OC)*", "*Texture (TX)*" and "*Speed (SP)*" variations.

We structure the dataset into two categories of videos. *Oriented Motion* tasks (upper block in Table 4.2) and *Composite Motion* Tasks (lower block in Table 4.2). Each video consits of one or more of the challenges mentioned above.

Table 4.2: Description of Videos in TMT Dataset

| Seqs | Object | Challenges | # of Seqs |
|---:|---|---|---|
| Juice | Juice Box | TR,RO,SP,IL | 6 NL 6 DL |
| Cereal | Cereal Box | TR,RO,SP,IL | 6 NL 6 DL |
| Book I | Hardcover Book | PR,SP,IL,SR | 6 NL 5 DL |
| Book II | Hardcover Book | SC,SP,IL | 6 NL 6 DL |
| Book III | Hardcover Book | TR,OC,SP,IL | 5 NL 5 DL |
| Mug I | Coffee Mug | TR,SP,IL,SR | 6 NL 6 DL |
| Mug II | Coffee Mug | TR,PR,SP,IL,SR | 6 NL 6 DL |
| Mug III | Coffee Mug | RO,SP,IL,SR | 6 NL 6 DL |
| Bus | Toy Bus | TR,SC,PR,SP,IL | 1 NL 1 DL |
| Highlighting | Newspaper | OC,SP,IL,TR | 1 NL 1 DL |
| Letter | Envelope | PR,SP,IL | 1 NL 1 DL |
| Newspaper | Newspaper Page | PR,SC,SP,IL | 1 NL 1DL |

TR = Translation; RO = Rotation; IL = Illumination; PR = Perspective; SR = Specular; OC = Occlusion; TX = Texture; SP = Speed; NL = Normal Light; DL = Diffuse Light

**Oriented Motion Tasks**

Oriented or single motion tasks refer to highly structured motion of the object, where the human user or the robot hand performs one simple action in manipulating an object. Details of the tasks are described below.

Figure 4.2: Image frames of the video sequences show sample tasks performed by both a human user and a robot hand with two different objects.

**Juice**    Juice from a juice box is poured onto a container. The goal is to track the juice box in the presence of both translational (TR) and rotational (RO) motion of the object.

**Cereal**    This task is similar to *Juice*, with an added challenge. There is large jerky motion when the cereal is poured onto the bowl. Also, the cereal box has higher texture compared to the juice box.

**Book I**    The object (book) is tilted from a vertical upright position to finally lie flat on the table and back up again. The challenge is to handle perspective (PR) transformation of the object in the image plane. During this motion there is Specular Reflection (SR) from the cover of the book.

**Book II**    Here a book is brought near the camera, parallel to the camera axis and away. The goal is to capture scale (SC) variation of the object.

Juice            Cereal Box            Book            Mug

(a)



Newspaper            Bus            Page            Envelope

(b)

Figure 4.3: Planar projection of the objects used in the dataset are shown, both planar and curvilinear objects are chosen with varying texture, lambertian/specular to have a full spectrum of challenges. (a) Objects used to record *Oriented Motion Tasks*, (b) Objects used to record *Composite Motion Tasks*

**Book III**   The object (book) is placed inside a book holder. The challenge in this video is to track the translational (TR) and Rotational (RO) motion of the object, despite varying occlusion (OC) from the book holder.

**Mug I**   Translational (TR) motion of a coffee mug is recorded when it's lifted up. The cup being a small shiny object, specular reflection (SR) and low texture (TX) of the mug present additional challenges in this case.

**Mug II**   This sequence, although similar to *Mug I*, is more difficult to track because of high perspective (PR) deformation of the mug when it is tilted to drink the contents.

**Mug III**   A low texture (TX) coffee mug is rotated (RO) from its initial position and back along the axis perpendicular to the principal camera axis.

**Composite Motion Tasks**

Composite motion tasks consists of videos that can be decomposed into simpler *Oriented Motion* tasks. Some tasks involve the use of both the hands as against the simple motions in *Oriented Motion* tasks. For the composite motion taks only human recorded videos are reported.

**Bus**   The planar surface in front of a toy bus, that undergoes scale change (SC) and perspective deformation (PR) at varying speed (SP) is to be tracked.

**Highlighting**   A portion of the newspaper is highlighted with a yellow marker pen. The challenge is to track the object in the presence of changing texture (TX) caused by highlighting and partial occlusion (OC) by the hand of the user and the pen she uses for highlighting.

**Letter**   The sequence records a letter being put inside an envelope and out again. The challenge in largely to tackle the perspective deformation (PR) of the object. This is a complex motion where both hands are used, one manipulating the envelope, the other taking out the letter from the envelope.

**Newspaper**  A textured portion of the newspaper is to be tracked in the presence of perspective (PR) and scale (SC) changes under varying speed.

Most of the *Oriented Motion* tasks have 6 videos recorded in both Natural Light (NL) and Diffused Light (DL) settings. The six videos are of the same task performed at varying speeds. Speed vary from very slow to very fast with an added option of "increasing speed", in which the speed is varied while performing the task. *Complex Motion* tasks, on the other hand report one video with varying speed. Having videos of different speeds, correctly captures blurr as a challenge at higher speeds.

A total of 8 different objects (Fig 4.3) were used to record the videos. To have a diverse selection, we choose objects that are planar (book, cereal box, juice box), curvillinear (mug, yogurt can), lambertian (newspaper), specular (book cover, mug, envelope), large (book, cereal box, juice box) and small (bus, mug). Each frame is tagged with the different challenges it presents. The distribution of frames over the challenges is shown in Figure 4.4. Sometimes more than one challenge is present simultaneously in one frame. All such challenges are tagged. The number on top of each bar in Figure 4.4 represent the number of frames having the corresponding challenge. TR or translation is the most common challenge since in most videos it is part of the initial motion that the object goes through. Only the highlighting sequence present change in Texture (TX). However, we use different objects to cover a varied set of textures.

## Frame distribution of Challenges

Figure 4.4: The distribution of number of frames over the different challenges for normal light is shown. The number on top of each bar represent the number of frames with the corresponding challenge. Abbreviations for challenges are, TR = Translation; RO = Rotation; PR = Perspective; SR = Specular; OC = Occlusion; TX = Texture, BL = Blurr casued at high speeds, SC = Scale change.

## 4.2 Ground Truth

Ground Truth (GT) refers to the correct position and orientation of the object in the image frame. A tracker's performance is evaluated against this data. Video datasets that evaluate surveillance trackers manually annotate a loose bounding box around the object of either 3 DOF (x, y, scale) [76] or 4 DOF (x, y, scale, rotation) [45]. They also benefit from the fact that the trackers are required to roughly follow the object with an allowable error limit of 50% overlap. Frame to frame correspondence of the target object is not very rigid in this case.

With registration based trackers the thresholds are much more stringent. Metaio [47] uses an absolute error threshold of 10 pixels to decide if a tracker has failed on not. They used a camera mounted on a Faro arm that was calibrated and all transformations of pose were stored. The stored values were used to calculate ground truth data. They also used fiducial markers to double check the ground truth hence generated. This setup is restricted by the workspace and mass of the robot arm.

We report sub-pixel level co-ordinate positions based on tracking data. Three trackers [22, 6, 9] are initiated on the first frame. Ground truth is registered only when bounding box co-ordinates reported by all three of them lie within $\pm 1$ pixel. The reason for choosing the trackers are because of their high convergence which is further shown in Chapter 5.

The stringent convergence criteria ($\pm 1$ pixel) sometimes result in the trackers failing to converge to a common bounding box. We reinitialize the three trackers using positional information from previous frames every time it fails to converge. Number of reinitializations varied from 0 to 12, for more difficult sequences. As a final step we also verify all ground truth data manually.

## 4.3 Error Metric

Error in single object tracking is a quantitative estimate of how closely the tracker follows the object. Tracked output ($X_T$) is compared with Ground

Truth ($X_{GT}$) to calculate this error. Any further analysis of a tracker first needs an error measure. Some of the common error measures in literature are, *Center Distance* [76] and *Area Overlap* [45].

### 4.3.1   Center Distance

Center distance ($E_C$) is expressed as the Eucledian distance between the tracked data and ground truth data. Mathematically it is expressed as shown in Eq. 4.1

$$E_C = \sqrt{(X_t - X_{GT})^2} \tag{4.1}$$

$X_t$ is centroid of a tracked frame, $X_{GT}$ the corresponding ground truth. This measure is more suitable for 2 DOF trackers that track translational motion of the object.

### 4.3.2   Area Overlap

Area overlap ($E_A$), a measure inspired from segmentation scores [70], is used in tracking in [76, 45]. Area overlap is expressed as a ratio of the area of intersection of the target area $A_T$ and ground truth area $A_{GT}$ with that of their union. Mathematically it is expressed as Eq. 4.2.

$$E_A = \frac{|A_T \cap A_{GT}|}{|A_T \cup A_{GT}|} \tag{4.2}$$

Both the above mentioned errors fail to properly capture the misalignment of pose. Fig 4.5 show a case where the Ground Truth ($X_{GT}$) and Tracked result ($X_T$) are 180° out of phase. *Center Distance ($E_C$)* Sec. 4.3.1 and *Area Overlap ($E_A$)* Sec. 4.3.2 fail to account for this out of phase alignment, and have very low error scores. To address this, we describe a measure called Alignment Error, Sec. 4.3.3. This error measure is similar to the one used in [47].

Figure 4.5: Alignment error $E_{AL}$ accounts for the displacement of each corner (colour coded) from the corresponding GT data. It is expressed as a root mean square (RMSE) score to account for the misalignment of pose.

### 4.3.3 Alignment Error

Alignment error quantifies the change in pose between the Target image $(X_T)$ and Ground Truth image $(X_{GT})$. Alignment Error $(E_{AL})$ is expressed as a root mean square distance of misalignment of the target image $(X_T)$ with ground truth $(X_{GT})$. Mathematically it is expressed as shown in Eq. 4.3.

$$Error(E_{AL}) = \sqrt{\frac{\sum_{i=1}^{4}(X_T - X_{GT})^2}{4}} \qquad (4.3)$$

The Root Mean Square Error (RMSE) is calculated between the corresponding four co-ordinates of the bounding box. For all of our experiments we use Alignment Error $(E_{AL})$.

## 4.4 Evaluation Measures

Using Alignment Error, Sec. 4.3.3, we evaluate trackers based on four different measures. (i) **Overall Success** (Robustness), (ii) **Average Drift** (Convergence), a measure that we define (iii) **Speed Sensitivity** and (iv) **Overall Rank** similar to [45].

### 4.4.1 Overall Success

Overall Success (OC) is a commonly used measure [45, 76]. It is defined as the fraction of total frames that a tracker tracks within an error $(E_{AL})$ threshold

of $t_p$ pixels. It is expressed as

$$\mathbf{OC} = \frac{|\mathbf{S}|}{|\mathbf{F}|}$$

$$\mathbf{S} = \{f^i \in \mathbf{F} : E^i_{AL} < t_p\}$$

(4.4)

$\mathbf{S}$ is the set of successfully tracked frames, $\mathbf{F}$ set of all frames, $E_{AL}$ error of frame $(f^i)$. Manipulation tasks calls for high precision in tracking. We set threshold $t_p$ to be 5 pixels in our experiments.

## 4.4.2 Average Drift

Average Drift (AD) computes the expected error of the tracker when it operates within a allowed drift of $t_p$ pixels. This is similar to the one proposed in [74].

$$\mathbf{AD} = E[e|e < t_p] = \frac{\sum\limits_{f_i \in \mathbf{F}} E_{AL}}{|\mathbf{S}|}$$

$$\text{subject to } e^i_{AL} < t_p$$

(4.5)

Average Drift checks a tracker, for how precisely it aligns itself with the target. Yang et al. [76] used a pixel threshold $t_p$ of 20 pixels. This is too high for manipulation tasks. As used in OC above, we set $t_p$ to be 5 pixels. It is important to calculate AD when the tracker is within 5 pixel drift or else it would be biased by large drifts when the tracker no longer tracks the object.

## 4.4.3 Speed Sensitivity

A common way [6, 22, 9] to quantify tracking convergence in dealing with large motion is to synthetically warp the standard "Lena" image and allow the tracker to recover this warp. The random warps are sampled from a Normal distribution having zero mean and standard deviation $\sigma$, $(\mathcal{N}(0, \sigma))$. However, this experimental methodology overestimates convergence. Warps sampled using a large sigma, with a certain probability still have smaller warps in it.

This is further illustrate in Figure 4.6. Overall Success (OC) is plotted against both interframe object motion ($m_i$ as shown in Eq. 4.6 and sigma ($\sigma$). Trackers perform poorly on large motions of 12 pixels (0.6 Success Rate)

Figure 4.6: Success rate ($t_p = 2$) plotted against varying inter frame motion and sigma. It shows that even though the success rate for higher sigma is high it's lower for the corresponding inter frame motion.

compared to the corresponding sigma (0.9 Success Rate) for NNIC tracker.

$$m_i = rmse(gt^{i+1}, gt^i) = \sqrt{\frac{\sum_{k=1}^{4}(gt_k^{(i+1)} - gt_k^{(i)})^2}{4}} \qquad (4.6)$$

Algorithm 4 shows how *Speed Sensitivity* is computed. Since focus in given on precise trackers, threshold $t_p$ (TH here) is 5 pixels. $\#tf$ is the number of frames that have error threshold below $t_p$ with inter frame motion $m$ in Algorithm 4. $\#f$ is the total number of frames having inter frame motion $m$.

When used with real videos from the TMT dataset, the number of samples having fast object motion is less. We skip frames (track every second, third and fourth frame) to simulate fast motion. We also track both forward and backward, initialise tracker on $k^{th}$ frame and track $k + 1^{th}$ frame, initialise on $k+1^{th}$ and track $k^{th}$ frame. This gives us enough samples to have a statistically significant result.

**Algorithm 4** Speed Sensitivity of Trackers

---

1: **procedure** SPEEDSENSITIVITY($GT, T$)                  ▷ T, tracked Result
2:     $r = dict\{\}$                                        ▷ {[m]:(#tf , #f)}
3:     $TH = 5$                          ▷ TH = Threshold for successful tracking
4:     **while** $I < len(GT) - 1$ **do**
5:         $m = rmse(GT[I + 1], GT[I])$
6:         $err = rmse(GT[I + 1], T[I + 1])$
7:         **if** err $<$ TH **then**
8:             dict[m] = (#tf+= 1, #f+= 1)
9:         **else**
10:            dict[m] = ( #tf, #f+= 1)
11:        **end if**
12:        $I = I + 1$
13:    **end while**
14:    **return** $dict.items()$
15: **end procedure**

---

### 4.4.4   Overall Rank

Measures discussed above rank trackers on isolated attributes of success and robustness. They fail to provide a global view of how the trackers compare. We use an idea similar to AR plots in [45]. Trackers are ranked based on Overall Success (OC) and Robustness. Finally a 2D plot of Overall Rank shows their relative performance taking into account both attributes of OC (Y-axis) and Robustness (X-axis). Since we are interested in evaluating precise trackers we use Alignment Error ($E_{AL}$) instead of $E_A$ used in [45].

Robustness is quantified by the number of times a tracker fails, and needs to be re-initialised. A tracker is considered to have failed when $E_{AL} > t_p = 5$ pixels. Once the tracker fails, the tracker is re-initialised 10 frames from the frame it has failed to remove any bias and account for the same challenge twice. Total number of re-initialisations give an indication of how robust a tracker is in a manipulation setup. Two graphs are plotted. A rank plot graph where X-axis has the robustness rank and Y-axis the success rank with both axises flipped. Trackers on the top right are the better performing trackers. Another graph, un-normalised graph plot, shows the actual number of re-initializations on the X-axis and the overall success on the Y-axis.

# Chapter 5

# Results and Evaluation

In this Chapter we evaluate and analyse both registration based trackers and OLT. Since we are interested in precise tracking, we use videos from TMT dataset (Chapter 4) [67] which is a publicly available dataset and reports 8 DOF ground truth. Static experiments using the standard "Lena Image" as done in [6, 9, 22, 80] is also done with results reported and analysed.

## 5.1    Baseline Trackers

All trackers were initialised on the first frame of the video with bounding box co-ordinates of the object. Original implementation of RKLT [80], NNIC [22], IVT [66], TLD [41] and L1-APG [56] were used. One pass evaluation is done on the videos, with the tracker running the full length of the video.

**Registration Based**    Among registration based trackers, we test ESM [9], NNIC [22], GNNIC (Chapter 3), IC [6] and RKLT [80]. 30 iterations with a resolution of $100 \times 100$ was used in IC for convergence. NNIC used a look up table of $(0.06 \times 0.04)$, $(0.03 \times 0.02)$ and $(0.015 \times 0.01)$ sigmas, each using 4000 samples and a resolution of $50 \times 50$. The corresponding IC used in the cascade literature, used 10 iterations with a resolution of $50 \times 50$. ESM used 30 iterations with a resolution of $50 \times 50$. RKLT used a template size of $40 \times 40$, with 10 iterations in IC for refinement. The parameters were chosen to give each tracker approximately the same execution time. Finally, we test SGANNS agorithm in GNNIC, using the same parameters as done

for NNIC, with degree of connectivity ($\mathbf{k}$) for graph construction set to 400. This gives a common baseline to compare the two related search methods of randomized kD-Trees and SGANNS. Two distance measures are used with GNNIC, Eucleadian distance ($l_2$) and Manhattan distance ($l_1$). Additionally all four apperance models discussed in Chapter 3 are studied in conjugation with the registration based trackers.

**Online Learning Based**   Three online learning based trackers (OLT) were studied, IVT [66] , L1-APG [56] and TLD [41]. IVT and L1-APG uses 6 DOF particle filter as the search method. TLD uses median flow motion model. Yang et al. in [76] evaluated some of the state of the art trackers that uses online learning to adapt the appearance model. We choose three of the trackers that performed well on most of the videos. Both IVT [66] and L1-APG [56] in our experiments used 600 particles for the particle filter search. Additionally, L1-APG [56] used a template subspace of 10 templates of the target image. Original parameters as specified in [41] were used for TLD.

A diverse selection, such as this, would show the usability of some of the popular trackers in precise tracking requirements like manipulation and visual servoing. This is missing in the literature, since OLT benchmarks don't impose high convergence.

## 5.2   Result and Analysis

**Overall Success**   Overall Success, averaged over three speeds (very slow, slow, normal) is reported in Table 5.1. It uses SCV as the appearance model. Later in Chapter we study the effect of appearance variation on tracking performance. The best performing tracker for each sequence is shown in bold.

Table 5.1 clearly shows that registration based trackers have better overall performance compared to OLT. OLT are ill suited for manipulation tasks. They build an appearance model of the object (lower dimensional subspace of the initial template(s)) to begin with. This model is then updated as tracking

Table 5.1: Overall Success Expressed as Fraction of Frames Tracked

| Video | TLD | L1APG | IVT | ESM | NNIC | IC | RKLT | GNNIC |
|---|---|---|---|---|---|---|---|---|
| Cereal | 0.00 | 0.24 | 0.99 | **1.00** | **1.00** | **1.00** | **1.00** | **1.00** |
| Book I | 0.00 | 0.10 | 0.48 | **1.00** | **1.00** | **1.00** | **1.00** | **1.00** |
| Book II | 0.00 | 0.79 | 0.30 | **1.00** | **1.00** | **1.00** | **1.00** | **1.00** |
| Book III | 0.00 | 0.42 | **0.72** | 0.34 | 0.32 | 0.32 | 0.38 | 0.34 |
| Juice | 0.00 | 0.16 | 0.98 | **1.00** | 0.41 | **1.00** | **1.00** | **1.00** |
| Mug I | 0.84 | 0.10 | 0.91 | **1.00** | **1.00** | **1.00** | **1.00** | **1.00** |
| Mug II | 0.41 | 0.30 | 0.72 | 0.89 | 0.89 | 0.89 | 0.31 | **0.92** |
| Mug III | 0.51 | 0.54 | 0.68 | **1.00** | 0.59 | 0.65 | **1.00** | 0.68 |
| Bus | 0.37 | 0.57 | 0.94 | **1.00** | 0.99 | 0.96 | 0.81 | **1.00** |
| Highlighting | 0.00 | 0.67 | **0.95** | 0.76 | 0.70 | 0.33 | 0.72 | 0.60 |
| Letter | 0.11 | 0.19 | 0.25 | **1.00** | **1.00** | **1.00** | **1.00** | 0.89 |
| Newspaper | 0.00 | 0.61 | 0.92 | **1.00** | 0.51 | 0.43 | **1.00** | 0.24 |

progresses. Newer templates are accounted for in the appearance of the object. Though this makes the trackers more robust, this comes at the expense of convergence. Newer templates that are shifted, scaled from the original, are not exactly the same as the target image.

We observe that L1-APG tracker fails to track in plane rotation (Fig 5.1). Although Mei et. al [56] use a 6 d.o.f motion model, the state space update mostly accounts for two DOF. The motion model use a velocity component $\tilde{\mathbf{v}} = (v_1, v_2)$ (horizontal and vertical velocities), which are the last two terms of the state vector $\mathbf{p}$. Particles are updated based on average of the last few frames tracked. This accounts for the transitional motion of the object but fails to capture the full pose. It is not reported in any existing work that we know of. TLD with it's 3 d.o.f. parametrization understandably fails to precisely estimate the object pose.

We further analyse registration based trackers. We test it on higher speed videos (fast and very fast from TMT [67] Dataset). RKLT [80] and GNNIC are the two best performing trackers on higher speed object motion. RKLT [80] tracker uses several point trackers (KLT trackers), that it initialises on each frame. This is then followed by RANSAC to estimate pose. This allows the

Figure 5.1: L1-APG (orange) [56] and TLD (black) [41] fails to track in plane rotation. Registration based trackers, ESM (red) [9], IC (yellow) [6] and NNIC (cyan) [22] and IVT (green) [66] tracks without considerable drift

tracker to handle large pose variation till some inliers are detected by RANSAC to estimate homography (Algorithm I in [80]). NNIC [22] and GNNIC use a two step search process. It models the possible pose variations prior to tracking. With sufficient number of samples, large motion variations can be modelled. Fine alignment following the Approximate Nearest Neighbour step is done using Inverse Compositional [6] search in NNIC [22] and GNNIC. The two search methods are complementary to each other, thus improving robustness maintaining high convergence.

Both NNIC [22] and GNNIC use Approximate Nearest Neighbour, which are computed using two different algorithms. We compare these two algorithms on high speed videos since for slower speeds both these algorithms track the full sequence, which makes it hard to distinguish. Figure 5.3 show the result. The red boxes indicate those sequences where GNNIC is better than or equal to NNIC [22]. Apart from the one exception of "mugI" on very fast speed ("mugI-s5"), GNNIC is better than NNIC in terms of overall success. Further experiments on high speed object motions are shown and analysed in the next Section 5.2.

Both the algorithms (NNIC and GNNIC) map warps to sampled images to model possible transformations before tracking starts. The efficiency of the system depends on how many samples are needed to properly model the object.

Figure 5.2: Overall Success of six registration based tracker are plotted. The six trackers being ESM [9], IC [6], RKLT [80], NNIC [22] and two variants of GNNIC one using Eucleadian distance (GNNIC-L2) and the other using Manhattan distance (GNNIC-L1) on 12 video sequences from TMT Dataset [67].



Figure 5.3: Overall Succes (OC) of NNIC [22] and the two variants of GNNIC on the higher speed videos from TMT Dataset [67].

We further test the two algorithms to verify this. KD-Tree [8] based NNIC and SGANNS based GNNIC are compared without using the cascade framework IC [6], to test the contribution of the approximate nearest neighbour part of the algorithm. Overall success is plotted against varying threshold $t_p$. SGANNS using 5000 samples perform comparable to ANN (for low pixel thresholds used in manipulation) using much higher particles ($> 10{,}000$ particles). Figure 5.4 shows the comparison. Number of samples needed to appropiately represent data is much less in SGANNS which in turn makes the system more efficient. The result validates that keeping the same setting (number of particles and resolution of templates) SGANNS has better overall accuracy for most videos

compared to KD-Tree in Figure 5.3.



Figure 5.4: Overall Success is plotted with varying threshold $t_p$ for different number of samples on "BookI" sequence. Results show that GNNIC performs well with much lesser samples compared to NNIC [22].

Not only does GNNIC require less number of samples to model the object, it speeds up search as pointed out in [32]. The speed up using SGANNS is measured based on the number of distance computations during search, to find the optimal parameters.

Table 5.2: Speed Comparison expressed as number of distance computations.

| Video | BookI | BookII | MugI | MugII | Juice | Cereal |
|---|---|---|---|---|---|---|
| NNIC | 1412 | 1216 | 811 | 821 | 2107 | 1781 |
| GNNIC | 922 | 651 | 408 | 411 | 732 | 691 |
| Speed up | 1.53 | 1.86 | 1.98 | 1.99 | 2.87 | 2.57 |

The reason for using such a method to compare speed is because KD-Tree is implemented using Flann [1], which is a more optimised version. The

---
[1]Available at: https://github.com/mariusmuja/flann

number of distance computations are shown in the Table 5.2 on six videos. All results are reported on medium speed. **Speed up** is the ratio of the number of distance computations needed in NNIC compared to that of GNNIC. The speed is almost double in most cases. This is explained by the fact that having a smart criteria to initialise search (start traversing the graph **G**, from the best node in the previous search) reduces the number of distance computations.

**Speed Sensitivity** Evaluation measure *Speed Sensitivity* introduced in Section 4.4.3 studies tracking performance on fast object motion. This is important for registration based trackers, which are often used in manipulation tasks, that have a wide range of motion. Originally, fast object motion was simulated by warping the Standard "Lena" image (Figure 5.5) using a random warp, followed by the tracker retrieving this warp. Success Rate (OC) over a large number of trials (5000) for each sigma is plotted with increasing sigma.

As reasoned in Section 4.4.3 and shown in Figure 4.6 , this is not an accurate estimate of Speed Sensitivity. Speed Sensitivity if lower, trackers can track higher inter frame motion. This is further shown in Figure 5.5 (b). All five registraton trackers are tested using the original Static Experimental setup in [6, 9, 22] and the reported metric *Speed Sensitivity* in Chapter 4. NNIC [22], RKLT [80] and GNNIC are the three best performing trackers. Note that the success plots are very much similar when plotted against sigma, but shows a high degree of variability when plotted against interframe motion. RKLT [80] has the best performance with GNNIC close second. The trackers are tested further on real motion sequences. Speed sensitivity is reported taking into account all five speeds. We have made ground truth data publicly available, users can define their own evaluation criteria. *Speed Sensitivity* results are grouped under two object categories i.e. *Planar Object* and *Curvillinear Object.*

**Planar Objects:** "BookII" sequence is the easiest to track. The only challenge is scale (SC) variation. Compared to "BookII", "BookI" (Fig 5.6) has both perspective transformation (PR) and specularity (SR).

(a)



(b)

Figure 5.5: (a) Original "Lena" image that is morphed and tracked and (b) Speed Sensitivity experiments on the same image.

(a)



(b)

Figure 5.6: Speed Sensitivity result on "bookI" sequence. Frame 1 and Frame 190, in the "bookI" sequence, recorded at medium speed. The considerable change in appearance due to specular reflection is visually shown on top right corner

In this sequence when the book is tilted, (Fig 5.6) the appearance changes considerably because of specular reflection from the surface of the book as shown in Figure 5.6 (b). This makes it hard to track. Another interesting observation can be made (Fig: 5.7), though "Cereal" and "Juice" sequences have same sets of challenges, the "Juice" sequence is more sensitive to large motion.



Figure 5.7: Speed sensitivity is plotted against inter frame motion for "Cereal" (top) and "Juice" (bottom) sequences

The only difference between them being texture, we investigate further.

Static image experiments [22, 6, 9] with optimal parameters are done. The results in Fig 5.8 show the effect of texture playing a crucial role, making the two sequences distinct in their own way. Cereal box having richer texture, is easier to track.



Figure 5.8: Static experiment [22, 6, 9] with 5000 trials for each sigma ($\sigma \in \{1, ..14\}$), is done on the cereal box (top) and juice box (bottom), cereal box having a rich texture is easier to track

*Curvilinear Object:* The mug sequences are on an average more difficult to track. With higher inter frame motion the success rate falls steeply. The specularity (SR) on the surface makes it even hard to track. "MugII" is the most difficult sequence in the entire set, with the highest number of challenges

Table 5.3: Average Drift Expressed as an Expectation Value

| Video | TLD | L1-APG | IVT | ESM | NNIC | IC | RKLT | GNNIC |
|---|---|---|---|---|---|---|---|---|
| Juice | N/A | 2.39 | 3.27 | 0.79 | 0.89 | 0.94 | **0.73** | 1.57 |
| Cereal | N/A | 2.67 | 2.80 | 0.61 | 0.61 | **0.15** | 0.63 | 0.64 |
| Book I | N/A | 2.87 | 2.96 | 0.73 | 0.73 | **0.70** | 0.92 | 1.57 |
| Book II | N/A | 3.32 | 3.37 | **0.70** | **0.70** | 0.73 | 0.82 | 1.67 |
| Book III | 4.97 | 1.21 | 1.14 | **1.38** | 1.41 | 1.49 | 1.45 | 1.41 |
| Mug I | 3.48 | 3.3 | 1.34 | 0.49 | 0.49 | **0.46** | 0.50 | 0.90 |
| Mug II | 3.64 | 2.47 | 3.09 | 1.87 | 1.88 | 1.870 | 1.87 | **0.74** |
| Mug III | 3.20 | 2.51 | 2.28 | 0.74 | 0.97 | 0.75 | 0.75 | **0.42** |
| Bus | 3.25 | 0.68 | 2.18 | 0.63 | **0.59** | 0.63 | 1.20 | 1.20 |
| Highlighting | N/A | 3.71 | 1.61 | 1.23 | 1.21 | **0.47** | 1.06 | 1.21 |
| Letter | 4.53 | 1.79 | 1.68 | **0.36** | 0.505 | **0.36** | 1.1 | 0.89 |
| Newspaper | N/A | 2.49 | 3.18 | 0.42 | 0.53 | **0.31** | 1.63 | 0.92 |

N/A = There were no frames that had an error $(E_{AL})$ less than 5 pixels

(Table 4.2).

**Average Drift**   Average Drift as described in Chapter 4 is calculated on all the sequences and reported in Table 5.3. Registration based trackers have low average drift compared to the online learned trackers for most of the sequences. The reason being, particle filter uses random samples to get an initial state estimate. This is unlikely to hit the best alignment when iterations are limited, while the Gauss Newton method of the original registration trackers provide fast and higher convergence when they do converge. IC tracker has the least drift in most of the sequences. This prompted us to use IC as the precise tracker in conjugation with the Approximate Nearest Neighbour based trackers which in itself does a coarse search are often don't converge well.

**Overall Rank:**   We rank the eight trackers in Figure 5.9, similar to AR plots in [45] . This gives a global overview of how the trackers compare. The methodology is described in Chapter 4. RKLT [80] ESM [9] and GNNIC are the 3 top performing trackers. Approximate Nearest Neighbour based trackers have higher success rates compared to ESM due of it's ability to track larger

object motion (Fig 4.6). RKLT [80] initialises KLT trackers on every frame, which makes it robust to fast motion and thus has the best overall success rate. TLD is the worst among the eight trackers, when used for precise tracking. It tracks 3 DOF state space and re-initialises often. IVT is the best among the Online Learned Trackers because of the 6 DOF motion model.
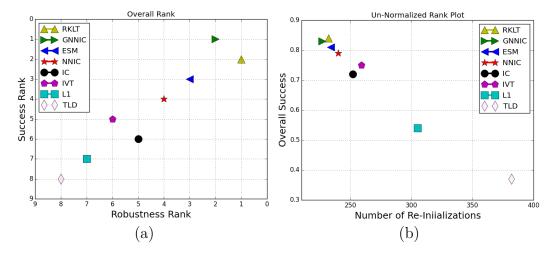


Figure 5.9: (a) Trackers are ranked [45] based on robustness and overall success. The error measure used is $E_{AL}$ with $t_p = 5$ pixels. (b) Un-Normalised rank plot shows the number of re-initializations and overall sucess. Apart from TLD [41] and L1 [56] the trackers perform similarly.

Figure 5.9 shows the rank plot and the un-normalised rank plots. The un-normalised plot (Figure 5.9 (b)) indicates that the top performing trackers (upper right corner of graph Figure 5.9 (a)) are similar in performance (Figure 5.9 (b) top left corner).

An important reason for the trackers to perform robustly is attributed to the appearance model of the object. Previously we tested using SCV as the appearance model. It corrects for illumination and hence tracks in the presence of specular reflection. In the next subsection we study the effect of appearance models and how it affects tracking performance.

### 5.2.1 Effect of Appearance Models in Tracking

Trackers are often subjected to changing appearance of the object over the course of a video. Hence different appearance models are used in registration

based trackers, to correct for illumination variation. Here we study four illumination correcting appearance models described in Chapter 3. We show average success rate on six videos, averaged over all the different speeds and search methods in Figure 5.10. IVA [60] is only slightly better than SSD. IVA [60] was originally designed to eliminate the effect of shadows in autonomous navigation. A downside of this is it removes considerable texture information from the object to be tracked. This makes tracking difficult, particularly in the case of low textured objects like "Juice" and "Mug" sequences. Figure 5.11 show sample transformed images. This is not the case for "Book" sequences. Image of the book still retain lot of the original texture after the transform.



Figure 5.10: Average Success Rate plotted for the four appearance models discussed in Chapter 3 on six sequences from TMT dataset [67]

NCC and SCV have similar performances in terms of Overall Success. However NCC is faster than SCV. Using NCC as the appearance model is equivalent to using SSD on normalised image [68] as opposed to calculating

the expected intensity value in SCV. The expectation operator in SCV has a complexity of $O(n^2m^2)$, where $(n \times m)$ is the dimension of the template. SSD, NCC and IVA, all three have a complexity of $O(nm)$.



Figure 5.11: Sample images of tranformed template image using IVA [60]. Very little texture is retained in most cases which makes tracking using this appearance model difficult.

NCC being both fast and accurate is the preferred choice of appearance model.

# Chapter 6

# An application to track Occlusion

Unlike TLD [41], the trackers discussed in the thesis don't have any detection module. Which means, once it looses track, it never re-initailises. This is a bottleneck when faced with occlussion or objects that have very little texture exposed at the beginning of the sequence. Note that, even though OLT adapts for changing appearances, the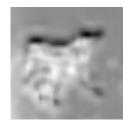y rarely re-initialize once tracking fails. To counter this, we propose a detection system that works together with the tracker and re-initialises when it has drifted substancially. Workflow of the full system with the tracker is shown in Figure 6.1. Alternatively the detection module can also be used to initialize tracking, if the object to be tracked is known ahead of time.

Object template ($\mathbf{I}^*$) initialised by the user is used to build the initial model of the object $\mathcal{O}$. `keypoints` (positional information) and the corresponding `feature descriptors` are extracted from $\mathbf{I}^*$. We use SIFT [50] features. This forms the initial object model, $\mathcal{O}$. As shown in Figure 6.1 we check for tracker failure (blue box) at every image frame. The tracker is considered to have failed when residual error ($||\mathbf{I}^* - \mathbf{I}(\mathbf{W}(\mathbf{x}; \mathbf{p}*))||_2 > TH$) is above a certain threshold $TH$. We choose $TH$ to be 20. If the tracker is judged to have failed, the detection system (green box) detects the object and re-initialises the tracker.

Feature descriptors from the new image (frame where the tracker drifted) are matched with the existing object model using Approximate Nearest Neigh-
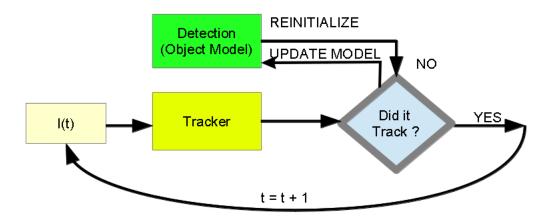
Figure 6.1: Worflow of the detection system together with the tracking module. When the residual is higher than a threshold, the object is detected and tracker re-initialised. The object module is updated with every re-initialization taking into account he new object template.

bour [8]. The matched keypoints are further filtered using RANSAC [24]. RANSAC checks every combination for the best fit with the existing object model, $\mathcal{O}$. RANSAC is computaionally expensive. Approximate Nearest Neighbour applied prior to RANSAC, narrows down the number of potential points to be checked. Finally, we compute homography between the keypoints in the object model and the matched keypoints from the image frame. Homography parameters give the new bounding box co-ordinates of the object.

This is used to re-initialise the tracker. Note that if the object is out of focus, there are no matched keypoints, and the detection system checks the next frame. With every re-initialization the object model is updated with the matched keypoints (*goodKeypoints* in Algorithm 5). Figure 6.2 shows the steps when the tracker fails. Figure 6.2 (a) shows the SIFT keypoints. These keypoints are matched with the object model, (b) shows the keypoints within a certain threshold $t$ (distance ratio). Finally Figure 6.2 (c) shows RANSAC selecting a handful of keypoints to estimate homography parameters to give the new bounding box shown in red.

**Algorithm 5** Detection Module

---

1: **procedure** DETECTION($\mathcal{O}, I$)  ▷ $\mathcal{O}$=(`reference keypoints,`
   `feature descriptors`), I = Image
2:  $goodKeypoints = \{\}$
3:  $(keypoints, descriptors) = SIFT(I)$
4:  **for** `keypoint` *in* `keypoints` **do**
5:     $dist = nnsearch(\mathcal{O}, keypoint)$  ▷ Approximate Nearest Neighbour
6:     **if** $dist < t$ **then**
7:        $goodKeypoints.append($`reference keypoint, keypoint`$)$
8:     **end if**
9:  **end for**
10:  $RANSAC(goodKeypoints)$  ▷ Filters $goodKeypoints$
11:  $POSE = HOMOGRAPHY(goodKeypoints)$
12:  $UPDATE(\mathcal{O})$
13:  **return** $POSE$
14: **end procedure**

---
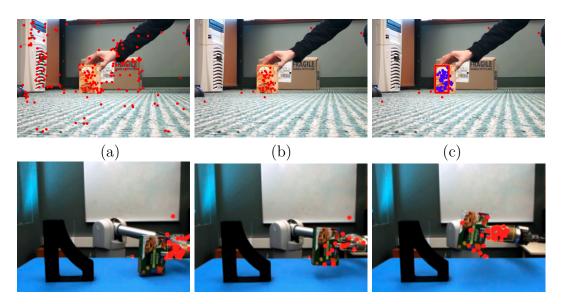


|  (a)  |  (b)  |  (c)  |

Figure 6.2: [TOP IMAGE] (a) Detected keypoints on an Image Frame (red points), (b) Keypoints that match closely with the object model $\mathcal{O}$, (c) Filtered keypoints after RANSAC in purple, used to calculate homography parameters. The resultant bounding box is shown in red. [BOTTOM IMAGE] End effector of the robot arm is tracked with the filtered keypoints.

Two examples are shown when using a detection system makes the tracking system more robust. The first example in Figure 6.3 shows occlussion from the book holder. Registration based trackers fail to track, but NNIC [22] coupled with the detection system re-initialises the tracker and it successfully tracks the object. The second example shows tracking very small objects. In this case a fish head. The position of the fish head is used as a reference to find the hook which is threaded. Being a very small object, it lacks significant texture. Other trackers fail very fast. The updating model of the object enables tracking the fish head for a long period of time.



Figure 6.3: Examples frames where the tracker was re-initialized using the detection module. Seven trackers are shown, DNNIC refers to NNIC with the detection module (red), GNNIC, NNIC [22], ESM [9], RKLT [80], IC [6] and TLD [41]

TLD [41] is shown to track both the "Book" sequence and the small fish head in Figure 6.3. However, in the "Book" sequence, it fails to track full pose. This is becasue of the 3 DOF motion model used in TLD. The proposed detection system tracks full pose yet maintaining high robustness by re-initializing the tracker when it fails, hence is better suited in robotic applications.

# Chapter 7

# Conclusion and Discussion

## 7.1 Contributions

In this thesis we address registration based tracking. While tracking evaluation methodologies exist in literature, little attention has been paid to evaluate registration based trackers. We report an extensive dataset for that purpose on which we evaluate a diverse set of trackers and finally rank them. A new evaluation measure, *Speed Sensitivity*, is introduced that test a tracker on large object motion.

We adapt a new search method, SGANNS [32], in tracking. It uses the sequential property of video data to efficiently search for the next best warp parameters, using the current warp as the starting node. The algorithm is explained in Chapter 3. Detailed results and analysis of the algorithm is presented in Chapter 5. Here we also study four appearance models in relation to registration based trackers. This is useful in making the tracker robust to illumination changes.

Interesting observations are revealed in our study. Firstly, results show L1-APG [56] tracker fails to track in plane rotation. This is attributed to the motion model used that primarily updates the translational parameters. Secondly, we observe that registration based trackers overall perform better than OLT. This explains the popular choice of registration based trackers in manipulation and visual servoing tasks. Thirdly, we show SGANNS to be more efficient compared to KD-Tree in tracking, with acheiving better accuracy with less number of samples to model motion. Fourthly, we study different

appearance models. NCC is shown to be have the best trade off in terms of speed and accuracy. It corrects for illumination by centering the data.

Finally we describe a detection system in Chapter 6. This coupled with a registration based trackers re-initialises the tracker when it drifts. We show applications in tracking occlussion, track small objects and objects that change appearance.

## 7.2 Future Work

### 7.2.1 Dataset Augmentation and Tracker Evaluation

One key area, is to augment the dataset with more videos. This can be done by including eye in hand camera or first person view camera [73] videos, which would include large scale and perspective changes in appearance of the object. Figure 7.1 shows a typical setup from three view points. The camera is attached to the palm of the robot hand as it grabs the box and puts the content in to the blue container. Tracking in such videos is more challenging due to considerable change in appearance of the object.



| (a) | (b) | (c) |

Figure 7.1: Three view points of the setup is shown. (a) An external camera shows the overall setup used for recording. (b) A camera attached to the palm of the robot hand provides egocentric view. (c) First person view is shown from a camera placed on top of the base of the robot.

We evaluate eight trackers and analyze performance. However the dataset being publicly downloadable with ground truth data, others are encouraged to use the dataset to evaluate their trackers.

### 7.2.2 Approximate Nearest Neighbour Search using Graph Structure

In this thesis we show SGANNS to work well with sequential data in tracking videos. Application in SLAM is shown in [32]. Future research on this is to work on ways to efficiently add and delete nodes to the connected graph which forms the model prior to tracking. This is important, as it would allow for the model to adapt to newer appearances of the object (newer observations added) without substancial memory usage (prune the directed graph and remove old observations) which are no longer useful.

# Bibliography

[1] Wam arm by barett technologies.

[2] Amit Adam, Ehud Rivlin, and Ilan Shimshoni. Robust fragments-based tracking using the integral histogram. In *Computer vision and pattern recognition, 2006 IEEE Computer Society Conference on*, volume 1, pages 798–805. IEEE, 2006.

[3] K.Pauwels A.Pieropan, G.Salvi and H.Kjellstrøm. A dataset of human manipulation actions, 2014.

[4] Shai Avidan. Support vector tracking. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 26(8):1064–1072, 2004.

[5] Boris Babenko, Ming-Hsuan Yang, and Serge Belongie. Visual tracking with online multiple instance learning. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 983–990. IEEE, 2009.

[6] Simon Baker and Iain Matthews. Lucas-kanade 20 years on: A unifying framework. *International journal of computer vision*, 56(3):221–255, 2004.

[7] John L Barron, David J Fleet, and Steven S Beauchemin. Performance of optical flow techniques. *International journal of computer vision*, 12(1):43–77, 1994.

[8] Jeffrey S Beis and David G Lowe. Shape indexing using approximate nearest-neighbour search in high-dimensional spaces. In *Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on*, pages 1000–1006. IEEE, 1997.

[9] Selim Benhimane and Ezio Malis. Real-time image-based tracking of planes using efficient second-order minimization. In *Intelligent Robots and Systems, 2004.(IROS 2004). Proceedings. 2004 IEEE/RSJ International Conference on*, volume 1, pages 943–948. IEEE, 2004.

[10] Marcelo Bertalmio, Guillermo Sapiro, and Gregory Randall. Morphing active contours. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(7):733–737, 2000.

[11] Michael J Black and Allan D Jepson. Eigentracking: Robust matching and tracking of articulated objects using a view-based representation. *International Journal of Computer Vision*, 26(1):63–84, 1998.

[12] Michael J Black and Allan D Jepson. Eigentracking: Robust matching and tracking of articulated objects using a view-based representation. *International Journal of Computer Vision*, 26(1):63–84, 1998.

[13] Gunilla Borgefors. Distance transformations in digital images. *Computer vision, graphics, and image processing*, 34(3):344–371, 1986.

[14] Jean-Yves Bouguet. Pyramidal Implementation of the Affine Lucas Kanade Feature Tracker: Description of the algorithm. Technical report, Intel Corporation Microprocessor Research Labs, 2001.

[15] Gary R Bradski. Computer vision face tracking for use in a perceptual user interface. 1998.

[16] Robert Collins, Xuhui Zhou, and Seng Keat Teh. An open source tracking testbed and evaluation web site. In *IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, pages 17–24, 2005.

[17] Dorin Comaniciu, Visvanathan Ramesh, and Peter Meer. Real-time tracking of non-rigid objects using mean shift. In *Computer Vision and Pattern Recognition, 2000. Proceedings. IEEE Conference on*, volume 2, pages 142–149. IEEE, 2000.

[18] Dorin Comaniciu, Visvanathan Ramesh, and Peter Meer. Kernel-based object tracking. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 25(5):564–577, 2003.

[19] Timothy F Cootes, Gareth J Edwards, and Christopher J Taylor. Active appearance models. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (6):681–685, 2001.

[20] Amaury Dame and Eric Marchand. Accurate real-time tracking using mutual information. In *Mixed and Augmented Reality (ISMAR), 2010 9th IEEE International Symposium on*, pages 47–56. IEEE, 2010.

[21] Alessio Del Bue, Dorin Comaniciu, Visvanathan Ramesh, and Carlo Regazzoni. Smart cameras with real-time video object generation. In *Image Processing. 2002. Proceedings. 2002 International Conference on*, volume 3, pages III–429. IEEE, 2002.

[22] Travis Dick, Camilo Perez Quintero, Martin Jägersand, and Azad Shademan. Realtime registration-based tracking via approximate nearest neighbour search. In *Robotics: Science and Systems*. Citeseer, 2013.

[23] Nicholas Dowson and Richard Bowden. Mutual information for lucas-kanade tracking (milk): An inverse compositional formulation. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (1):180–185, 2007.

[24] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.

[25] Robert Fisher, Jose Victor, and James Crowley. Caviar: Context aware vision using image-based active recognition dataset.

[26] Steffen Gauglitz, Tobias Höllerer, and Matthew Turk. Evaluation of interest point detectors and feature descriptors for visual tracking. *International journal of computer vision*, 94(3):335–360, 2011.

[27] Michael Gleicher. Projective registration with difference decomposition. In *Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on*, pages 331–337. IEEE, 1997.

[28] Helmut Grabner, Michael Grabner, and Horst Bischof. Real-time tracking via on-line boosting. In *BMVC*, volume 1, page 6, 2006.

[29] Helmut Grabner, Jiri Matas, Luc Van Gool, and Philippe Cattin. Tracking the invisible: Learning where the object might be. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 1285–1292. IEEE, 2010.

[30] Steve Gu, Ying Zheng, and Carlo Tomasi. Efficient visual object tracking with online nearest neighbor classifier. In *Computer vision–ACCV 2010*, pages 271–282. Springer, 2011.

[31] Gregory D Hager and Peter N Belhumeur. Efficient region tracking with parametric models of geometry and illumination. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 20(10):1025–1039, 1998.

[32] Kiana Hajebi and Hong Zhang. An efficient index for visual search in appearance-based slam. In *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, pages 353–358. IEEE, 2014.

[33] Mei Han, Amit Sethi, Wei Hua, and Yihong Gong. A detection-based multiple object tracking method. In *Image Processing, 2004. ICIP'04. 2004 International Conference on*, volume 5, pages 3065–3068. IEEE, 2004.

[34] Sam Hare, Amir Saffari, and Philip HS Torr. Struck: Structured output tracking with kernels. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 263–270. IEEE, 2011.

[35] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, New York, NY, USA, 2 edition, 2003.

[36] Wei He, Takayoshi Yamashita, Hongtao Lu, and Shihong Lao. Surf tracking. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 1586–1592. IEEE, 2009.

[37] Anil Jain, Karthik Nandakumar, and Arun Ross. Score normalization in multimodal biometric systems. *Pattern recognition*, 38(12):2270–2285, 2005.

[38] Omar Javed and Mubarak Shah. Tracking and object classification for automated surveillance. In *Computer Vision?ECCV 2002*, pages 343–357. Springer, 2002.

[39] Allan D Jepson, David J Fleet, and Thomas F El-Maraghi. Robust online appearance models for visual tracking. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 25(10):1296–1311, 2003.

[40] Frédéric Jurie and Michel Dhome. Hyperplane approximation for template matching. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 24(7):996–1000, 2002.

[41] Zdenek Kalal, Krystian Mikolajczyk, and Jiri Matas. Tracking-learning-detection. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 34(7):1409–1422, 2012.

[42] Changick Kim and Jenq-Neng Hwang. Fast and automatic video object segmentation and tracking for content-based applications. *Circuits and Systems for Video Technology, IEEE Transactions on*, 12(2):122–129, 2002.

[43] Dominik Alexander Klein, Dirk Schulz, Simone Frintrop, and Armin B Cremers. Adaptive real-time video-tracking for arbitrary objects. In *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*, pages 772–777. IEEE, 2010.

[44] Bianca Kochner, Dietrich Schuhmann, Markus Michaelis, Gerd Mann, and Karl-Hans Englmeier. Course tracking and contour extraction of retinal vessels from color fundus photographs: Most efficient use of steerable filters for model-based image analysis. In *Medical Imaging'98*, pages 755–761. International Society for Optics and Photonics, 1998.

[45] Matej Kristan, Roman Pflugfelder, Ale Leonardis, Jiri Matas, Fatih Porikli, Luka Cehovin, Georg Nebehay, Gustavo Fernandez, Toma Vojir, Adam Gatt, et al. The visual object tracking vot2013 challenge results. In *Computer Vision Workshops (ICCVW), 2013 IEEE International Conference on*, pages 98–111. IEEE, 2013.

[46] Laura Leal-Taixé, Anton Milan, Ian D. Reid, Stefan Roth, and Konrad Schindler. Motchallenge 2015: Towards a benchmark for multi-target tracking. *CoRR*, abs/1504.01942, 2015.

[47] Sebastian Lieberknecht, Selim Benhimane, Peter Meier, and Nassir Navab. A dataset and evaluation methodology for template-based tracking algorithms. In *Mixed and Augmented Reality, 2009. ISMAR 2009. 8th IEEE International Symposium on*, pages 145–151. IEEE, 2009.

[48] Liang Lin, Yongtian Wang, Yue Liu, Caiming Xiong, and Kun Zeng. Marker-less registration based on template tracking for augmented reality. *Multimedia Tools and Applications*, 41(2):235–252, 2009.

[49] Liang Lin, Yongtian Wang, Yue Liu, Caiming Xiong, and Kun Zeng. Marker-less registration based on template tracking for augmented reality. *Multimedia Tools and Applications*, 41(2):235–252, 2009.

[50] David G Lowe. Object recognition from local scale-invariant features. In *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, volume 2, pages 1150–1157. Ieee, 1999.

[51] Bruce D Lucas, Takeo Kanade, et al. An iterative image registration technique with an application to stereo vision. In *IJCAI*, volume 81, pages 674–679, 1981.

[52] Hosam M Mahmoud, Reza Modarres, and Robert T Smythe. Analysis of quickselect: An algorithm for order statistics. *Informatique théorique et applications*, 29(4):255–276, 1995.

[53] Ezio Malis. Improving vision-based control using efficient second-order minimization techniques. In *Robotics and Automation, 2004. Proceedings. ICRA'04. 2004 IEEE International Conference on*, volume 2, pages 1843–1848. IEEE, 2004.

[54] Mario Edoardo Maresca and Alfredo Petrosino. Matrioska: A multi-level approach to fast tracking by learning. In *Image Analysis and Processing–ICIAP 2013*, pages 419–428. Springer, 2013.

[55] Stephen J McKenna, Yogesh Raja, and Shaogang Gong. Tracking colour objects using adaptive mixture models. *Image and vision computing*, 17(3):225–231, 1999.

[56] Xue Mei and Haibin Ling. Robust visual tracking using l1 minimization. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 1436–1443. IEEE, 2009.

[57] Xue Mei and Haibin Ling. Robust visual tracking and vehicle classification via sparse representation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(11):2259–2272, 2011.

[58] A. Zisserman Michaelmas. Optimization lectures. In *B1 Optimization*, volume 2.

[59] Rainer Mobus and Uli Kolbe. Multi-target multi-object tracking, sensor fusion of radar and infrared. In *Intelligent Vehicles Symposium, 2004 IEEE*, pages 732–737. IEEE, 2004.

[60] Michael Paton, Kirk MacTavish, Chris J Ostafew, and Timothy D Barfoot. It's not easy seeing green: Lighting-resistant stereo visual teach & repeat using color-constant images. In *Robotics and Automation (ICRA), 2015 IEEE International Conference on*, pages 1519–1526. IEEE, 2015.

[61] Gord Peters. Corainder: Ieee1394 camera gui. *Linux Magazine*, 69, 1999.

[62] Antoine Petit, Guillaume Caron, Hideaki Uchiyama, Eric Marchand, et al. Evaluation of model based tracking with trakmark dataset. In *2nd Int. Workshop on AR/MR Registration, Tracking and Benchmarking*, 2011.

[63] Mark R Pickering, Abdullah Muhit, Jennie M Scarvell, Paul N Smith, et al. A new multi-modal similarity measure for fast gradient-based 2d-3d image registration. In *Engineering in Medicine and Biology Society, 2009. EMBC 2009. Annual International Conference of the IEEE*, pages 5821–5824. IEEE, 2009.

[64] Nilanjan Ray, Scott T Acton, and Klaus Ley. Tracking leukocytes in vivo with shape and size constrained active contours. *Medical Imaging, IEEE Transactions on*, 21(10):1222–1235, 2002.

[65] Rogério Richa, Raphael Sznitman, Russell Taylor, and Gregory Hager. Visual tracking using the sum of conditional variance. In *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*, pages 2953–2958. IEEE, 2011.

[66] David A Ross, Jongwoo Lim, Ruei-Sung Lin, and Ming-Hsuan Yang. Incremental learning for robust visual tracking. *International Journal of Computer Vision*, 77(1-3):125–141, 2008.

[67] Ankush Roy, Xi Zhang, Nina Wolleb, Camilo Perez, Quenterio, and Martin Jagersand. Tracking benchmark and evaluation for manipulation tasks. In *International Conference on Robotics and Automation*. IEEE, 2015.

[68] Lars Ruthotto. Mass-preserving registration of medical images. *German Diploma Thesis (Mathematics), Institute for Computational and Applied Mathematics, University of Münster*, 2010.

[69] Thomas B Sebastian and Benjamin B Kimia. Metric-based shape retrieval in large databases. In *Pattern Recognition, 2002. Proceedings. 16th International Conference on*, volume 3, pages 291–296. IEEE, 2002.

[70] David W Shattuck, Gautam Prasad, Mubeena Mirza, Katherine L Narr, and Arthur W Toga. Online resource for validation of brain segmentation methods. *Neuroimage*, 45(2):431–439, 2009.

[71] Ihor Smal, Katharina Draegestein, Niels Galjart, Wiro Niessen, and Erik Meijering. Particle filtering for multiple object tracking in dynamic fluorescence microscopy images: Application to microtubule growth analysis. *Medical Imaging, IEEE Transactions on*, 27(6):789–804, 2008.

[72] Arnold WM Smeulders, Dung M Chu, Rita Cucchiara, Simone Calderara, Afshin Dehghan, and Mubarak Shah. Visual tracking: An experimental survey. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 36(7):1442–1468, 2014.

[73] Siddhartha S Srinivasa, Dave Ferguson, Casey J Helfrich, Dmitry Berenson, Alvaro Collet, Rosen Diankov, Garratt Gallagher, Geoffrey Hollinger, James Kuffner, and Michael Vande Weghe. Herb: a home exploring robotic butler. *Autonomous Robots*, 28(1):5–20, 2010.

[74] Björn Stenger, Thomas Woodley, and Roberto Cipolla. Learning to track with multiple observers. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 2647–2654. IEEE, 2009.

[75] Tiesheng Wang, Irene YH Gu, and Pengfei Shi. Object tracking using incremental 2d-pca learning and ml estimation. In *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*, volume 1, pages I–933. IEEE, 2007.

[76] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Online object tracking: A benchmark. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 2411–2418. IEEE, 2013.

[77] Ning Xu, Narendra Ahuja, and Ravi Bansal. Automated lung nodule segmentation using dynamic programming and em-based classification. In *Medical Imaging 2002*, pages 666–676. International Society for Optics and Photonics, 2002.

[78] Junlan Yang, Dan Schonfeld, and Magdi Mohamed. Robust video stabilization based on particle filter tracking of projected camera motion. *Circuits and Systems for Video Technology, IEEE Transactions on*, 19(7):945–954, 2009.

[79] Alper Yilmaz, Xin Li, and Mubarak Shah. Contour-based object tracking with occlusion handling in video acquired using mobile cameras. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 26(11):1531–1536, 2004.

[80] Xi Zhang, Abhineet Singh, and Martin Jagersand. Rklt: 8 dof real-time robust video tracking combining coarse ransac features and accurate fast template registration. In *Canadian Conference on Robot Vision, 2015. CRV 2015*, 2015.

[81] Tao Zhao and Ram Nevatia. Tracking multiple humans in crowded environment. In *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, volume 2, pages II–406. IEEE, 2004.

[82] Karel Zimmermann, Jiri Matas, and Tomas Svoboda. Tracking by an optimal sequence of linear predictors. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 31(4):677–692, 2009.