

Combining Deep Learning and Verification for Precise Object Instance Detection

Junyu Nan

CMU-RI-TR-19-NN

November 20, 2019



The Robotics Institute
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA

Thesis Committee:

David Held

Kris Kitani

Xiaofang Wang

*Submitted in partial fulfillment of the requirements
for the degree of Master of Science in Robotics.*

Copyright © 2019 Junyu Nan. All rights reserved.

Abstract

Deep learning object detectors often return false positives with very high confidence. Although they optimize generic detection performance, such as mean average precision (mAP), they are not designed for *reliability*. For a reliable detection system, if a high confidence detection is made, we would want high certainty that the object has indeed been detected. To achieve this, we have developed a set of verification tests which a proposed detection must pass to be accepted. We develop a theoretical framework which proves that, under certain assumptions, our verification tests will not accept any false positives. Based on an approximation to this framework, we present a practical detection system that can verify, with *high precision*, whether each detection of a machine-learning based object detector is correct. We show that these tests can improve the overall accuracy of a base detector and that accepted examples are highly likely to be correct. This allows the detector to operate in a high precision regime and can thus be used for robotic perception systems as a reliable instance detection method. Code is available at <https://github.com/siddancha/FlowVerify>.

Acknowledgments

I would like to thank my advisor Professor David Held for his advice to this project and overall support of my Master's study. I've learned a lot from his insights about the research field and methodologies.

I would also like to thank Siddharth Ancha for collaborating on this project. It has been a very pleasant and inspiring experience working with him. Besides, I would also thank Professor Kris Kitani and Xiaofang Wang for being my committee members.

I want to also express my appreciations to all the amazing people I met at the Robotics Institute. Finally, great thanks to my family for their support and understanding.

Funding

This material is based upon work supported by the United States Air Force and DARPA under Contract No. FA8750-18-C-0092 as well as by the NSF under Grant No. IIS-1849154.

Contents

1	Introduction	1
2	Related Work	5
2.1	Object Instance Detection	5
2.2	Verification	6
2.3	Dense Pixel Correspondence.	7
3	Theoretical Framework	9
4	Approach	17
4.1	Verification Tests	17
5	Experiments	21
5.1	Implementation Details	21
5.1.1	Base Instance Detector	21
5.1.2	Dense Pixel-wise Correspondence	22
5.1.3	Network Architecture	22
5.2	Results	25
5.3	SIFTVERIFY Baseline	26
5.4	Ablation	27
5.5	Qualitative Analysis	28
5.6	Qualitative Comparison with SIFTVERIFY	29
5.6.1	Visualization Format	29
5.6.2	SIFT Successful Cases	30
5.6.3	SIFT Unsuccessful Cases	31
5.7	Failure Cases	36
6	Future Work	39
7	Conclusion	41

List of Figures

1.1	Pipeline of our instance detection verification system. <i>Base Detector</i> : generates proposed object instance detections; FLOWMATCHNET: computes dense pixelwise correspondences between template images of an object and a proposed detection; FLOWVERIFY: a suite of verification tests is applied to the detection using the estimated correspondences. The detection is accepted only if all verification tests pass.	2
5.1	(a), (b): Target and scene images from MS-COCO synthetic dataset that applies affine transformations. (c), (d): Target and scene images from BigBIRD-RGBD synthetic dataset that applies homography transformations.	24
5.2	Precision-Recall (PR) curves for the base instance detector, the detector improved by FLOWVERIFY, and the SIFTVERIFY baseline, evaluated on the GMU (a) and W-RGBD (b) datasets.	25
5.3	(a, b): False positives from the base detector filtered by FLOWVERIFY; (c, d): True positives from the base detector accepted by FLOWVERIFY.	28
5.4	False detections which are successfully filtered out by both FLOWVERIFY and SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.	32
5.5	True detections that are successfully retained out by both FLOWVERIFY and SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.	33
5.6	False detections which are incorrectly retained by SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.	34
5.7	True detections that are incorrectly filtered by SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.	35
5.8	Failure cases of FLOWVERIFY. (a) is an incorrect detection accepted by our method; (b) is a correct detection rejected by our method; Blue box in each example denotes the bounding box prediction from instance detector.	36

List of Tables

- 5.1 Overall performance (mAP) and max Precision results. 26
- 5.2 Tradeoff of performance vs speed as we vary the number of viewpoints. 26
- 5.3 Ablation analysis. 28

Chapter 1

Introduction

Instance detection is the task of detecting instances of a particular object in a scene. Here (as in previous work [18, 24, 53]), the term “instance” refers to a specific sub-type of object (e.g. “coke can” rather than just “can”). For example, a robot may be shown an example image of a cereal box, and it may be required to detect it in a scene in order to fetch it, even if it is in a slightly different viewpoint than the example image. A user may wish to give a robot instructions that refer to a specific object, e.g. to bring the user’s coffee mug (as opposed to a random coffee mug). Unfortunately, current systems for object instance detection are not sufficiently reliable for use in real world applications; most methods will fail in real-world scenarios due to occlusions, lighting changes, viewpoint variation, and other difficulties.

Traditional non-parametric methods for object instance detection rely on keypoint-matching [5, 41, 45] or template matching [25, 28] to a set of template images. However, these methods are not robust to large changes in object viewpoint (i.e. greater than 25 degrees [24]) and often fail for significant lighting changes. For a robust home perception system, we cannot always guarantee that the objects being observed will be viewed from the same conditions as in the templates. Thus we desire to have an object detection system that is robust to significant changes in the object viewpoint, as well as lighting, occlusions, and other variations. Recently, a number of machine learning approaches have been used for object instance detection [21, 24, 34]. However, machine-learning based approaches often produce a large number of false positive detections. These false detections can prevent deployment of robots in

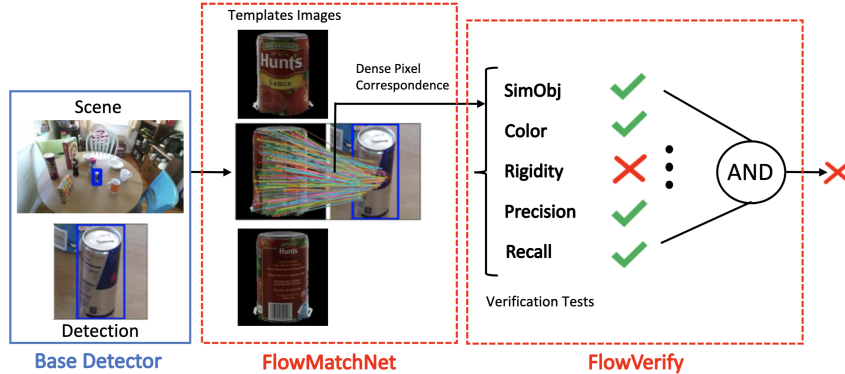


Figure 1.1: Pipeline of our instance detection verification system. *Base Detector*: generates proposed object instance detections; *FLOWMATCHNET*: computes dense pixelwise correspondences between template images of an object and a proposed detection; *FLOWVERIFY*: a suite of verification tests is applied to the detection using the estimated correspondences. The detection is accepted only if all verification tests pass.

real-world applications.

Our first insight is that a combination of parametric and non-parametric approaches can be used to obtain both robust and reliable object detection. While machine learning methods enable us to leverage large amounts of data to learn necessary invariances like lighting and viewpoint changes, non-parametric methods that match candidate detections to template images can ensure that these detections are reliable. In sections 3 and 4.1 we propose verification tests that operate on top of existing machine learning detectors; these tests take proposed detections as input and verify them by matching against template images. Thus, by not throwing away training data (in contrast to the parametric machine learning approach) and using them to verify detections at test time, we show that we can increase the accuracy of the detector, especially in the high-precision regime. This makes our system more precise, reliable and interpretable, making it more suitable for robotics applications.

Our approach is based on a novel theoretical framework for verified object detection. Specifically, under our theoretical framework, we prove that under certain conditions, a series of verification tests will reject all false positive detections. Under these assumptions, any detection that is verified by our tests is guaranteed to be correct. These verification tests are based on estimated dense correspondences between the proposed detection and a set of template images. It is important to note that our

theoretical framework does not make assumptions on the accuracy of the base detector or the accuracy of the estimated correspondences; the tests can fail due to either a false positive proposed by the detector or due to inaccurately estimated correspondences. However, the tests will pass only if the detection is a true positive.

Based on our theoretical framework, we implement an approximate but more practically suitable version of these tests. In our practical framework, if any of the tests fail, rather than rejecting the detection, we reduce its initial confidence score. We show that our method improves the performance of the detector in the high precision regime and performs no worse in the low-precision regime, thereby improving the overall accuracy of the detector. Figure 1.1 shows a diagram of our method.

Because our system matches proposed detections to template images, the runtime is a function of the number of template images used. However, the number of template images can be varied to obtain an accuracy-speed tradeoff; we demonstrate that the user of our system can use more templates from different viewpoints to achieve higher accuracy, or they can choose to use fewer templates to increase speed. We summarize our contributions as:

- A theoretical framework for verified object instance detection, which guarantees, under certain assumptions, no false positive detections.
- An approximate implementation of this framework which leads to a practical instance detection method
- A demonstration that our method leads to significantly improved detector performance, especially in the high-precision regime.

CHAPTER 1. INTRODUCTION

Chapter 2

Related Work

2.1 Object Instance Detection

General Instance Detection. Traditional methods for object instance detection rely on keypoint-matching [5, 41, 45] or template matching [25, 28]. Recently, detectors based on deep neural networks have shown improved performance over these traditional approaches. Many general object category detectors such as SSD [38] or Faster R-CNN [46] are still used to solve the instance detection problem. Other methods [12, 20] leverage synthetic data to attack the instance problem. Target Driven Instance Detection (TDID) [3] is a state of the art instance detection method based on a Siamese style neural architecture [32]. We compare to TDID in our work and show improved performance over this state-of-the-art baseline.

Grocery product recognition. There has been significant effort in recognizing products on shelves of retail stores [15, 17, 22, 42, 50, 52]. Many works investigate the features to use for production detection and recognition, such as SIFT [17, 22], color histograms [17], and HoG [52]. Other works exploit the statistical correlation between nearby products on shelves [1, 6]. This problem is simpler than the general object instance recognition problem that we are aiming to solve due to the structured environment. For example, objects in such scenes are typically placed in front-facing canonical viewpoints with minimal occlusions and good lighting conditions. Furthermore, there is often a great deal of contextual information, such as the location of the product on the shelf. Our method is designed for the more general usage in

the home or other unstructured locations.

2.2 Verification

Verification for Deep Neural Networks. Due to popularity of deep neural networks, many efforts have been made to verify certain properties about DNNs [27, 31, 43, 44]. These works generally attempt to verify consistencies of network output under small perturbations of inputs [27, 31, 43, 44], which is not directly related to verifying the accuracy of an detection network. Additionally, they often have scalability issues [27, 31, 43, 44], and therefore have not been applied to instance detection due to the large-scale input and network structure. In contrast, we design a verification framework to directly verify the accuracy of detection network and prove that it outputs no false positives under certain assumptions; moreover, our method can apply generally to any deep-learning based instance detectors.

Classification with reject option. The notion of verification is related to the literature on classification with rejection. Some of these previous approaches assume a cost function for rejection is provided and find the optimum rejection rule [4, 7, 9] or learn classification and rejection functions simultaneously [10, 16]. Other works [23, 48] treat the problem as conformal prediction, where the classification system can predict any subset of classes, including the null set which stands for rejection.

In contrast, in our work, we verify whether the object class output by the detector is correct, and we re-score the confidence of detections based on the verification tests.

Verification in grocery product recognition. Tonioni *et al.* [50] proposes a method for grocery product recognition, which involves a refinement stage to remove false positive detections. The method uses a geometric consistency criterion to rerank detections and remove false positives, which computes the similarity between the features in the query image and features in the reference image using their spatial locations relative to the center of the image. This criterion assumes query and the reference images are at the same canonical viewpoint, which is generally not applicable to the general instance detection setting. In comparison, we aim to design verification tests for the more general instance detection problem.

2.3 Dense Pixel Correspondence.

There has been many works on predicting dense pixel correspondences. Specifically, the task of estimating pixel correspondences between consecutive frames in the video, or estimating optical flow, has been a traditional computer vision problem and has wide applications in video analysis and processing. Traditional methods for optical flow estimation includes variational approaches [26], and variational approaches combined with combinatorial matching [47]. Recently, deep learning based methods [8, 13, 29] have obtained state-of-the-art performance.

Many methods use dense pixel correspondences to solve other tasks, such as tracking [54], video segmentation [51] and robot manipulation [11]. In this paper, we build upon previous works a matching network to predict dense pixel correspondences. However, unlike past work, we use these correspondences to improve performance of object instance recognition.

CHAPTER 2. RELATED WORK

Chapter 3

Theoretical Framework

In this section, we lay out a formal framework for high-precision instance detection. We specify a set of assumptions and prove that under these assumptions, the verification-based detector will have no false positives. We then discuss a modification of this method that we implement in practice which improves overall detection accuracy.

Problem Statement. We assume there are O categories of objects that we are trying to recognize. We also assume access to a dataset which consists of a variety of images per object, recorded from a set of different viewpoints and lighting conditions. We refer to these images as ‘template images’ (see Figure 1.1). At test time, we are given ‘scene images’ – these are real world images that contain (multiple) objects that need to be detected (see Figure 1.1). We also assume that we have an object instance detector trained to detect objects of interest; this detector will return detections of the target objects, along with a proposed object class for each detection. However, many of these proposed object classes will be incorrect, i.e. false positives. The goal of our framework is to filter these out.

Notation. We denote by O_i the i th object from the O categories of objects. For each object O_i , we assume access to a dataset of template images recorded from a variety of viewpoints and lighting conditions; we denote M_i as the index set for these template images, and $I_{i,m}$ as the m th template image for object O_i .

We denote by $T : (u_1, v_1) \rightarrow (u_2, v_2)$ a 2D mapping from pixels of one image to the pixels of another image. We overload notation such that $T(I)$ is an application of the 2D mapping T to all pixels in the image I to produce a new image $T(I)$. Let $r(T)$

CHAPTER 3. THEORETICAL FRAMEWORK

denote the “rigidity” of such a transformation, measured by the fraction of inlier pixel matches under the best approximating rigid transformation. In other words, for any given transformation T , we find

$$r(T) = \max_{\bar{T} \in T^R} \frac{\sum_{u,v} \mathbb{1}\{\|T(u,v) - \bar{T}(u,v)\| \leq \epsilon\}}{\sum_{u,v} 1} \quad (3.1)$$

where $\mathbb{1}\{\cdot\}$ is the indicator function and ϵ is a constant. We denote by T^R the set of 2D mappings which are perfectly rigid, such that $r(T) = 1$. Let D denote a detection in a scene image (e.g. a crop of a scene image).

We assume access to a similarity classifier which returns a score $c(I_1, I_2)$ indicating the confidence that images I_1 and I_2 each contain the same object class. We also assume access to a distance metric $d(I_1, I_2)$ that measures the the distance between images I_1 and I_2 . Any distance function $d : \mathcal{I}^2 \rightarrow \mathbb{R}^+$ that satisfies the following two properties can be used (the properties will be necessary for the proof later):

1. Triangle Inequality: $d(I_1, I_2) \leq d(I_1, I_3) + d(I_2, I_3)$.
2. Permutational Invariance: $d(I_1, I_2) = d(\sigma(I_1), \sigma(I_2))$, where σ is any permutation of image pixels.

Note that many distance metrics satisfy the above properties, such as the maximum difference in pixel intensities $d(I_1, I_2) = \|I_1 - I_2\|_\infty$, any L_p norm, or normalized cross-correlation.

We denote $F(I_1, I_2)$ as a function that computes a 2D mapping between images I_1 and I_2 . Typically, the objective of this function is to find a rigid 2D mapping that minimizes some distance metric d , i.e.

$$\arg \min_{T \in T^R} d(T(I_1), I_2) \quad (3.2)$$

However, in our framework, we make no assumptions about the output of $F(I_1, I_2)$; we will prove that our system will have no false positives, regardless of the output of $F(I_1, I_2)$.

Assumptions. We make the below assumptions; we will prove that, with these assumptions, our algorithm will lead to no false positives. In practice, we will need to relax these assumptions for a practical implementation; nonetheless, this theoretical

framework provides the basis for our approach.

Assumption 1. (Dense Dataset)

For each detection D of ground truth class O_i , $\exists m \in M_i, T \in T^R$ such that $d(T(I_{i,m}), D) \leq \gamma$.

Assumption 2. (Similarity Classifier Smoothness)

We have a similarity classifier c with a corresponding constant δ that satisfies the following property: for any two images I_1 and I_2 , and for any detection D , if $\exists T \in T^R$ such that $d(T(I_1), I_2) \leq 2\gamma$, then $|c(I_1, D) - c(I_2, D)| < \delta$.

The first assumption states that our dataset is dense enough such that every detection can be constructed as a rigid transformation of some image in the dataset, with a bounded lighting change applied. The second assumption states that the similarity classifier c is a smooth function: if two images I_1 and I_2 are sufficiently similar such that I_2 can be created by a rigid transformation and a small lighting change applied to I_1 , then the similarity classifier c will output a similar score when comparing I_1 and D as when comparing I_2 and D .

TheoreticalFlowVerify. Our object detection pipeline proceeds as follows: we assume that an initial detector finds detections D in an image and proposes an object class O_i for each detection. We then pass the detection and its corresponding proposed class to our verification system `THEORETICALFLOWVERIFY(i, D)` which returns True if it can verify that D contains an image of class O_i and False otherwise. We present `THEORETICALFLOWVERIFY` in Algorithm ??; this method is a theoretical version of our practical algorithm `FLOWVERIFY` described in Section 4.1. The algorithm proceeds as follows: For any detection D and proposed object class O_i , `THEORETICALFLOWVERIFY` iterates over all images $I_{i,m}$ of the proposed object class O_i and compares each template image $I_{i,m}$ to the detection D using `VERIFYMATCH`. `VERIFYMATCH` then estimates a set of dense pixel-wise correspondences \hat{T} between the detection D and the image $I_{i,m}$, and determines whether it can validate the detection using the following verification tests:

1. Similar Object Comparison: Tests whether an image $I_{j,n}$ of another object class O_j looks similar to D according to similarity function c ; formally: $\forall j \neq i, \forall n \in M_j$, if $c(I_{i,m}, D) < c(I_{j,n}, D) + \delta$, return False
2. Color Comparison: Tests whether the image distance between $I_{i,m}$ transformed

CHAPTER 3. THEORETICAL FRAMEWORK

using \hat{T} and the detection D is sufficiently small: if $d(\hat{T}(I_{i,m}), D) > \gamma$, return False

3. Flow Rigidity: Tests whether the the correspondences \hat{T} could have been derived from a rigid object transformation: if $r(\hat{T}) < 1$, return False

This algorithm is designed to return False for all false positive detections. Here we state and prove the main theorem of our approach:

Theorem 1. (No False Positive Theorem)

Under assumptions 1 and 2, THEORETICALFLOWVERIFY does not produce any false positives. That is, the following statement always holds: $\text{THEORETICALFLOWVERIFY}(i, D)$ returns False whenever the ground-truth class for detection D is different from O_i .

Note that $\text{THEORETICALFLOWVERIFY}(i, D)$ may sometimes return False even if i is the correct ground-truth class of D if it cannot verify this proposal. In such cases, $\text{THEORETICALFLOWVERIFY}$ will return false for all object classes, meaning that we cannot verify the category of detection D using our approach. Still, the benefit of our approach is that, when $\text{THEORETICALFLOWVERIFY}(i, D)$ returns true, we can be assured that detection D is an object of class O_i .

Proof. We need to show that $\text{THEORETICALFLOWVERIFY}(i, D)$ returns False whenever the ground-truth class of object D is not O_i . Note that from line 17, we need to prove that for any object class O_i that is not equal to the ground truth class, every template image $I_{i,m}$, for all $m \in M_i$ of object O_i needs to be rejected by VERIFYMATCH . That is, $\text{VERIFYMATCH}(I_{i,m}, D) = \text{False}, \forall m \in M_i$, if $gt(D) \neq O_i$. Assume that the ground-truth object class of D is O_j . We can partition M_i into three classes:

1. $M_{i,1} := \{m : \exists T \in T^R \text{ such that } d(T(I_{i,m}), D) \leq \gamma\}$.
2. $M_{i,2} := \{m \in M_i \mid d(\hat{T}(I_{i,m}), D) > \gamma\} \setminus M_{i,1}$.
3. $M_{i,3} := M_i \setminus (M_{i,1} \cup M_{i,2})$.

Now we show $\text{VERIFYMATCH}(I_{i,m}, D) = \text{False}$, for each of the three cases above.

1. Case 1: $m \in M_{i,1}$. VERIFYMATCH will return False at line 6.

Since $m \in M_{i,1}$, we know that $\exists T_i \in T^R$ such that $d(T_i(I_{i,m}), D) \leq \gamma$. In other words, $M_{i,1}$ is a set of images of object class O_i that can be described by a rigid

Algorithm 1 Theoretical Flow Verifier

```

1: procedure VERIFYMATCH( $I_{i,m}, D$ )      ▷ returns True if verification succeeds
2:   // Test 1: Similar Object Comparison
3:   for  $j \in \mathcal{I} \setminus \{i\}$  do
4:     for  $n \in M_j$  do
5:       if  $c(I_{i,m}, D) < c(I_{j,n}, D) + \delta$  then
6:         return False
7:       end if
8:     end for
9:   end for
10:  // Test 2: Color Comparison
11:   $\hat{T} \leftarrow F(I_{i,m}, D)$                                 ▷ Estimated flow
12:  if  $d(\hat{T}(I_{i,m}), D) > \gamma$  then
13:    return False
14:  end if
15:  // Test 3: Flow Rigidity
16:  if  $r(\hat{T}) < 1$  then
17:    return False
18:  end if
19:  return True
20: end procedure
21: procedure THEORETICALFLOWVERIFY( $i, D$ )
22:   for  $m \in M_i$  do
23:     if VERIFYMATCH( $I_{i,m}, D$ ) then
24:       return True
25:     end if
26:   end for
27:   return False
28: end procedure

```

CHAPTER 3. THEORETICAL FRAMEWORK

transformation and a small lighting change applied to D (which is an object of class O_j). Thus there are two object classes, O_i and O_j , that both appear similar to detection D , although detection D is an object of class O_j .

From assumption 1, $\exists n \in M_j, T_j \in T^R$ such that $d(T_j(I_{j,n}), D) \leq \gamma$. Let us define a new operator $T_{ij} = (T_j)^{-1} \circ T_i$, i.e. first apply T_i and then apply inverse of T_j . First, note that, if $T_i \in T^R$ and $T_j \in T^R$, then $T_{ij} \in T^R$ since the composition of these rigid transforms is a rigid transform. Secondly,

$$\begin{aligned} d(T_{ij}(I_{i,m}), I_{j,n}) &= d(T_j \circ T_{ij}(I_{i,m}), T_j(I_{j,n})) \\ &= d(T_i(I_{i,m}), T_j(I_{j,n})) \\ &\leq d(T_i(I_{i,m}), D) + d(T_j(I_{j,n}), D) \\ &\leq 2\gamma \end{aligned}$$

The first line holds by the permutation invariance property of d , since T_j is a permutation of pixels. The second line holds by definition of T_{ij} . The third line holds by triangle inequality of d . The last line holds because of the definition of m and T_i that $d(T_i(I_{i,m}), D) \leq \gamma$, and by definition of n and T_j that $d(T_j(I_{j,n}), D) \leq \gamma$.

By assumption 2, since $d(T_{ij}(I_{i,m}), I_{j,n}) \leq 2\gamma$, then $|c(I_{i,m}, D) - c(I_{j,n}, D)| < \delta$. This then implies that

$$\begin{aligned} c(I_{i,m}, D) &< c(I_{j,n}, D) + \delta \\ c(I_{i,m}, D) &< \max_n c(I_{j,n}, D) + \delta \end{aligned}$$

Hence, the conditional at line 5 will be true and so `VERIFYMATCH`($I_{i,m}, D$) will return False at line 6.

2. Case 2: $m \in M_{i,2}$. `VERIFYMATCH` will return False at line 10.

By definition of $M_{i,2}$, we know that $d(\hat{T}(I_{i,m}), D) > \gamma$. Hence the conditional at line 9 will be true, so `VERIFYMATCH`($I_{i,m}, D$) will return False at line 10.

3. Case 3: $m \in M_{i,3}$. `VERIFYMATCH` will return False at line 13.

If $m \in M_{i,3}$, then by definition, $m \notin M_{i,1}$ and $m \notin M_{i,2}$. Since $m \notin M_{i,1}$, then $\nexists T \in T^R$ such that $d(T(I_{i,m}), D) \leq \gamma$. Since $m \notin M_{i,2}$, then $d(\hat{T}(I_{i,m}), D) \leq \gamma$. Hence we know that $\hat{T} \notin T^R$, so $r(\hat{T}) < 1$. Hence the conditional at line 12 will be true, so `VERIFYMATCH`($I_{i,m}, D$) will return False at line 13.

Hence we have shown that $\forall i \neq j, \forall m \in M_i, \text{VERIFYMATCH}(I_{i,m}, D) = \text{False}$.

This implies that `THEORETICALFLOWVERIFY`(i, D) = False whenever $i \neq j$. \square

Note that `THEORETICALFLOWVERIFY`(i, D) returns False whenever it cannot verify that the ground-truth class of D is O_i . This can sometimes occur even if O_i represents the ground-truth class of D , if the estimated correspondences \hat{T} are not accurate. It will also return false whenever the ground-truth class of D is not O_i . Thus `THEORETICALFLOWVERIFY` has 100% precision; every time it returns True, the proposed object class is the same as the ground-truth. On the other hand, the recall of the algorithm is not guaranteed; it may return False for an arbitrary proportion of examples, even if the proposed object class is correct. Note that the No False Positive Theorem does not make any assumptions on the quality of the set of predicted correspondences \hat{T} ; if the correspondences are not accurate, then `THEORETICALFLOWVERIFY` will also reject the detection. We will now describe an approximation of this framework that leads to a practical implementation of this algorithm.

CHAPTER 3. THEORETICAL FRAMEWORK

Chapter 4

Approach

We will now describe an approximation of the theoretical framework that leads to a practical implementation of our method. The pipeline for high-precision detection consists of three stages:

1. *Base Instance Detector*: We first run an instance detector trained on the objects of interest. This stage provides candidate bounding box detections for target objects, as well as a proposed object class O_i for each box. The focus of our work is verifying these detections, as described below.
2. *Dense Pixel-wise Correspondence* (FLOWMATCHNET): We next predict dense pixel-wise correspondences between template images and each proposed detection, cropped from the scene. See Appendix 5.1.2 for more details on the network architecture of FLOWMATCHNET.
3. *Verification tests* (FLOWVERIFY): Given the proposed detection and template images and estimated pixel-wise correspondences between the two, we conduct a set of *verification tests* to ascertain whether the proposed class of the detection is correct.

4.1 Verification Tests

We design verification tests that are modifications to our theoretical framework (Section 3). These tests, which we call ‘FLOWVERIFY’ tests (since they make use of

predicted flow correspondences), are intended to be stringent; any detection that does not pass these tests is deemed ‘rejected’ and will receive a lower detection score than the detections that pass the tests. We can be highly confident that the detections that pass our verification tests are likely to be true detections, boosting the performance of our detector in the high-precision regime, as our results demonstrate. The first three verification tests are derived from our theoretical framework of Section 3:

1. SimObj: This test corresponds to the similar object comparison test in our theoretical framework. For each detection, there should be only one target object being matched to it with high confidence. The similar object comparison test is the following: for each detection D we compute a confidence score $c(D)$ using the initial detector (such as TDID [3]). We then search whether there is another detection D' of another object which has an IoU of at least η_{iou} with D . If no such detection is found, $\text{SIMOBJ} = 1$. Otherwise, the similarity score for D is the minimum of the confidence difference across all other detections D' : $\text{SIMOBJ} = \min_{D'} \max(0, c(D) - c(D'))$. We define the *similar object* test using a boolean variable given by $T_{\text{SIMOBJ}} = (\text{SIMOBJ} > \eta_{diff})$.

2. FColor: This test corresponds to the color comparison test in our theoretical framework. We estimate the pixel-wise correspondences \hat{T} between a template image I_i of the proposed object class O_i and the detection D . We then check the image similarity between D and the template image transformed using the predicted correspondences $\hat{T}(I_i)$. We use normalized cross correlation (“ncc”) to measure this similarity, which lies in $[-1, 1]$. We define $\text{FCOLOR} = \frac{1}{2}(\text{ncc}(\hat{T}(I_i), D) + 1) \in [0, 1]$ and we define the *flow color* test using a boolean variable given by $T_{color} = (\text{FCOLOR} > \alpha_{color})$.

3. FRigidity: This test corresponds to the rigidity test in our theoretical framework. Assuming objects are rigid, if the detection box contains the target object, then we would expect an ideal mapping between the detection and a corresponding template image from a similar viewpoint to describe a rigid body transformation. We use RANSAC [14] to find the best-fit fundamental matrix (using the 8-point algorithm [39]) that maximizes the number of corresponding inlier pairs. The proportion of inliers under the best-fit mapping is a measure of how rigid the flow is. We define flow

rigidity as

$$\text{FRIGIDITY} = \frac{\#\text{inliers}}{\#\text{total correspondences}}$$

and we define a *flow rigidity* test using a boolean variable $T_{rig} = (\text{FRIGIDITY} > \alpha_{rig})$. Besides the verification tests motivated by our theoretical framework, we design two additional tests to evaluate the extent of each bounding box prediction:

4. FPrecision: If we expect the detection box to contain the entire object (and not be too small), we would expect all pixels in the template image to be mapped to pixels *inside* the bounding box. We define *flow precision* as

$$\text{FPRECISION} = \frac{\#\text{target correspondences mapped to inside bbox}}{\#\text{total correspondences}}$$

and the *flow precision* test is defined using a boolean variable $T_{prec} = (\text{FPRECISION} > \alpha_{prec})$.

5. FRecall: Complementary to precision, if we expect the detection box to contain only the detected object (and not be too large), we would expect the flow mapping to cover all pixels of the detection box. To measure this, we can look at the bounding box that tightly fits the extent of the mapped pixels in the scene. We define *flow recall* as the IoU between the box suggested by flow to the bounding box output by the detector, and compute a *flow recall* score as

$$\text{FRECALL} = \text{IoU}(\text{detector bbox}, \text{bbox suggested by flow}).$$

The *flow recall* test is defined using a boolean variable $T_{rec} = (\text{FRECALL} > \alpha_{rec})$. Our overall verification test is

$$\text{FLOWVERIFY} = T_{\text{SIMOBJ}} \wedge T_{color} \wedge T_{rig} \wedge T_{prec} \wedge T_{rec},$$

with given threshold values and parameters $\eta_{iou}, \eta_{diff}, \alpha_{color}, \alpha_{rig}, \alpha_{prec}, \alpha_{rec} \in [0, 1]$, tuned with a validation set as described in Section 5. For a given object of interest, we could have multiple template images from different viewpoints and lighting conditions. In this case, we say that the set of template images for a given class passes FLOWVERIFY if any one of the template images for that class passes these tests.

While FLOWVERIFY tests are designed to improve performance in high precision

CHAPTER 4. APPROACH

regime, they might worsen performance in the low-precision regime as they could reject true positives. As there are cases where performance in low-precision regime is also important, instead of completely rejecting detections not passing FLOWVERIFY tests, we rerank all detections from our base detector based on 1) whether a detection passes FLOWVERIFY tests and 2) the score predicted by base detector. All detections that pass FLOWVERIFY are ranked higher than all detections that don't. The reranking procedure can be viewed as reducing the confidence of the detections that do not pass the FLOWVERIFY tests. As we will see, this improves performance in the high-precision regime while at the same time maintaining performance in the low-precision regime.

Chapter 5

Experiments

In this chapter, we describe a practical implementation of our framework, and evaluate its performance on several instance detection benchmarks.

Datasets. We evaluate our framework on two tests sets – the GMU Kitchens test split [19] and W-RGBD scenes v1 Dataset [35]. Both datasets contain images of objects placed in indoor scenes such as kitchen surfaces, table tops, and living rooms. For FLOWMATCHNET and FLOWVERIFY, at test-time we use 15 equally spaced viewpoints per object taken from [49] or [33] to form our “template images”. We also vary the number of template images to explore the speed-accuracy tradeoff of our system in section 5.2. For GMU, we evaluate on the 11 BigBIRD objects; for W-RGBD, we evaluate on 9 textured objects.

5.1 Implementation Details

We explain in this section the implementation details of each of the three stages in our pipeline, as listed in Chapter 5.1.2.

5.1.1 Base Instance Detector

Our method focuses on improving the precision of an existing instance detector. One could use any instance detector with our approach; we use *target driven instance detection* (TDID) [3], as this method produces state-of-the-art results for instance

detection, including one-shot instance detection that we evaluate on in our work. In the one-shot scenario, the training and validation objects are separate from the objects evaluated on at test time; the detector must generalize to novel objects. TDID showed state-of-the-art performance for this task [3]. Similarly, we train other components of our system, such as FLOWMATCHNET, in a one-shot manner, such that our entire system can be used to detect novel objects.

The base detector TDID is trained on the Active Vision Dataset (AVD) [2], W-RGBD scenes Dataset [35], and synthetic images from Cut-Paste-Learn [12]. GMU objects common to the AVD dataset are removed from training so that we can evaluate this detector in a one-shot manner. In order to evaluate on W-RGBD Scenes dataset in a one-shot manner as well, for this evaluation we use the TDID model released by [3] that is trained only on AVD.

FlowMatchNet.

5.1.2 Dense Pixel-wise Correspondence

5.1.3 Network Architecture

We call our network for predicting dense pixel-wise correspondences FLOWMATCHNET. For every proposed detection, we pad the detection to square, crop it (which we refer to as the ‘cropped scene image’), and feed it to FLOWMATCHNET. FLOWMATCHNET also takes an image from the dataset and computes a mapping from every pixel in the template image to some pixel in the cropped scene image.

We compute pixel correspondence from each template image of the object to the cropped scene image. We train a deep neural network using a modification of the FlowNetC neural network architecture [11]. Specifically for FlowNetC, a series of convolutional layers are applied separately to input images to extract features from each image. Then, a *cross-correlation* layer is used to combine features coming from each image branch, to compare feature in every location in image1 with every location in image2. However, FlowNetC was designed for optical flow, which is typically applied to consecutive frames of a video with the assumption that pixels have small displacements between consecutive frames. Therefore, in the original FlowNetC, cross-correlation is approximated by considering only a 21×21 pixel neighborhood window [11, 30].

However, we cannot make such a ‘small displacement’ assumption in our case because the object in the scene crop might have undergone a fairly large rotation and translation relative to the template image. Hence, we perform a full cross correlation between all pixel pairs in the scene and template images. Furthermore, flow prediction is traditionally made in terms of displacement vectors. Since a full cross-correlations is implemented as unrolling the entire 2D feature map, we reparameterize flow in terms of global coordinates of pixels being mapped to in the second image. Hence we concatenate (x, y) pixel coordinates as additional feature maps, following [37].

Training details of FlowMatchNet

Optical flow models are trained using synthetic datasets that have ground truth flow [11, 30]. Flownet based models have been shown to generalize well when trained on such synthetic datasets.

For training FLOWMATCHNET, we train our model successively on synthetic datasets with increasing difficulty –

1. *MS-COCO with affine transformations*: we take images from the MS-COCO dataset [36], crop an area around the object using its annotated bounding box to simulate detections from a base detector, and apply a random rotation and translation to create a simulated scene image. These transformations help the flow model to learn simple affine transformations. Although the transformations are relatively simple, MS-COCO has a huge variety of objects; training on such a diverse set of objects can help the network to generalize to different objects types. Examples of transformations are in Figure 5.1.
2. *BigBIRD + RGBD scenes with homographies*: We take objects from the BigBIRD dataset [49], crop them out using segmentation masks, and paste them onto background scenes in the W-RGBD scenes dataset [35]. Since we want FLOWMATCHNET to generalize to novel objects, we exclude training on the 11 objects that are part of GMU Kitchens [19], our test set. We apply random homographies to the cropped images and blend them into scenes using Cut-Paste-Learn [12], while simulating randomized lighting conditions and image blurs. We limit the random homography to ensure that it is not too unrealistic and does not distort the object too much; specifically we ensure that

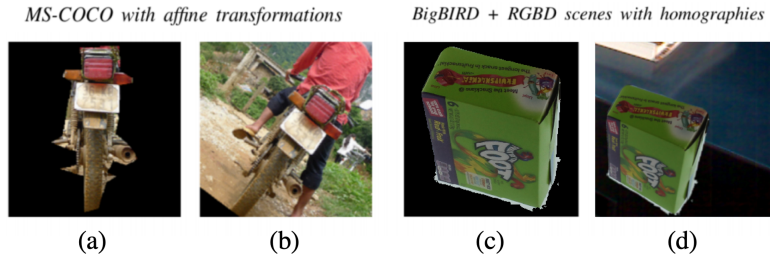


Figure 5.1: (a), (b): Target and scene images from MS-COCO synthetic dataset that applies affine transformations. (c), (d): Target and scene images from BigBIRD-RGBD synthetic dataset that applies homography transformations.

the top-left corner of the bounding box is still the top-left corner after applying the homograph, the bottom-right corner is still at the bottom-right, and so on. Importantly, since these are controlled synthetic transformations, ground truth flow can be computed. We train on L2 loss for optical flow similar to [11, 30].

We first trained on the MS-COCO synthetic flow dataset for 1.6M iterations using Adam optimizer with a learning rate of 10^{-4} with a batch size of 2. We then finetuned on the BigBIRD-RGBD synthetic flow dataset for 50,500 iterations using Adam optimizer with a learning rate of 10^{-6} and batch size 2.

Verification Tests

We use a validation set to tune the parameters for the tests in FLOWVERIFY. We set $\eta_{iou} = 0.5$, which is common used for mAP evaluation. As our validation set, we use the YCB-Video training set, on which we tune the values for $\alpha_{rig}, \alpha_{color}, \alpha_{prec}, \alpha_{rec}, \eta_{diff} \in [0, 1]$ by optimizing mAP. We find the optimal values to be $\alpha_{rig} = 0.9, \alpha_{color} = 0.5, \alpha_{prec} = 0.9, \alpha_{rec} = 0.3$ and $\eta_{diff} = 0.0$ (meaning that all other object classes should have strictly lesser confidence scores if detected).

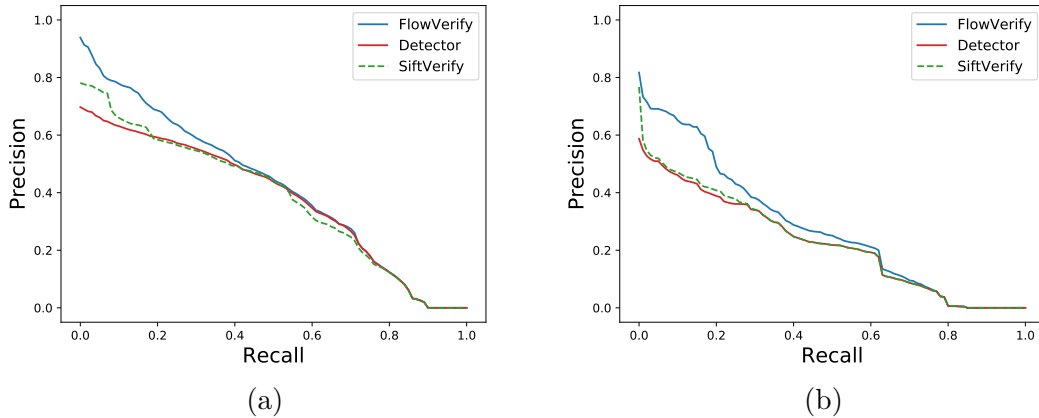


Figure 5.2: Precision-Recall (PR) curves for the base instance detector, the detector improved by FLOWVERIFY, and the SIFTVERIFY baseline, evaluated on the GMU (a) and W-RGBD (b) datasets.

5.2 Results

Figure 5.2 shows precision-recall curves of the base instance detector (TDID) and the improved detector using FLOWVERIFY, on GMU (a) and W-RGBD (b) datasets. On both datasets, FLOWVERIFY significantly improves detection performance in the high-precision regime to the left end of the PR curve. For high scoring detections, it makes significantly fewer errors than TDID as it is able to filter out many false positive detections using the verification tests. Table 5.1 summarizes mAP and maximum precision of baseline detectors and FLOWVERIFY on GMU and W-RGBD. It is worth noting that although FLOWVERIFY is designed to reject false positives and improve performance in the high-precision regime, the PR curves show that, with the re-ranking procedure, the performance does not degrade in the low-precision regime (right side of the curve). In the low-precision regime, FLOWVERIFY nearly matches the base detector in performance. This shows that verification tests can make an instance detector more reliable and precise while simultaneously improving overall performance.

We note that our method presents a tradeoff between precision/accuracy and timing performance. As more template images are used, the precision/accuracy of our system increases, but it will take longer to run. Hence, the user of our system can adjust

Method	GMU-Test		RGBD		time(s)
	mAP	max Precision	mAP	max Precision	
Base Detector	0.379	0.697	0.222	0.587	0.037
SIFTVERIFY	0.382	0.780	0.228	0.767	0.940
FLOWVERIFY	0.422	0.939	0.278	0.822	0.872

Table 5.1: Overall performance (mAP) and max Precision results.

the number of template images based on their specific needs. We report the running times of our method for varying number of template viewpoints in Table 5.2.

Method	GMU-Test		RGBD		time(s)
	mAP	max Precision	mAP	max Precision	
FLOWVERIFY(15 vp)	0.422	0.939	0.278	0.822	0.872
FLOWVERIFY-12vp	0.413	0.946	0.278	0.838	0.696
FLOWVERIFY-9vp	0.411	0.962	0.275	0.847	0.540
FLOWVERIFY-6vp	0.408	0.898	0.263	0.766	0.379
FLOWVERIFY-3vp	0.405	0.932	0.252	0.783	0.201

Table 5.2: Tradeoff of performance vs speed as we vary the number of viewpoints.

5.3 SIFTVerify Baseline

We implemented and evaluated another verification method using SIFT-based [40, 41] keypoint correspondences as a baseline for our learning-based dense-correspondences computed by FLOWMATCHNET. This baseline is designed to test our hypothesis that verification tests should make use of both machine learning and non-machine learning approaches; we use machine learning for computing the correspondences (FLOWMATCHNET) but we use the non-learning based verification tests of FLOWVERIFY to verify detections.

In contrast, SIFT [40, 41] is a non-learning based approach for computing correspondences, and we combine it with a non-learning based approach for verification. We implement verification tests on top of SIFT that are analogous to FLOWVERIFY. In order to estimate rigidity using the fundamental matrix, we need at least 8 matches.

However, SIFT matches are often sparse, and have as few as 2-3 matches. Hence, we simply compute the raw number of SIFT keypoint matches to replace the rigidity test, which we refer to as the SMATCHES test. For the precision test, we compute the proportion of matched pixels that lie inside the predicted bounding box, and refer to this test as SPRECISION. Since SIFT matches are few in number, it is also not possible to estimate recall the way we do for FRECALL, so we omit that test. We find that these settings give the best performance for this baseline. We call this baseline SIFTVERIFY. We run grid-search to find the best thresholds for SIFTVERIFY on YCB following the same procedure as for FLOWVERIFY.

Figure 5.2 and Table 5.1 show the performance on the GMU and W-RGBD datasets. For both datasets, SIFTVERIFY slightly improves performance over the base detector, but performs consistently worse than FLOWVERIFY. This demonstrates the value of using machine learning for computing correspondences, even if the final verification tests are not learned.

5.4 Ablation

This section lists detailed results of our ablation analysis in terms of mAP and maximum precision for FLOWVERIFY, versus dropping each verification test one at a time. As shown in Table 5.3, FLOWVERIFY has the best performance in terms of maximum precision and mAP on GMU-Test, as well the second highest maximum precision on W-RGBD. Dropping FRIGIDITY results in the largest drop in maximum precision and mAP on both GMU-Test and W-RGBD. This confirms the importance of *flow rigidity* test as suggested by our theoretical framework.

We also observe that removing the FCOLOR test leads to very little change in performance. We believe that this is because FLOWVERIFY is approximately optimizing Supplementary Equation 3.2, with a priority on matching colors:

$$F(I_{i,m}, D) = \hat{T} \approx \arg \min_T d(T(I_{i,m}), D). \quad (5.1)$$

In other words, FLOWVERIFY = $F(I_{i,m}, D)$ returns a transform \hat{T} that tries to match colors in $I_{i,m}$ with similar colors in D , thereby minimizing $d(\hat{T}(I_{i,m}), D)$. Because the FCOLOR test rejects the detection if $d(\hat{T}(I_{i,m}), D) > \gamma$, our system will very rarely



Figure 5.3: (a, b): False positives from the base detector filtered by FLOWVERIFY; (c, d): True positives from the base detector accepted by FLOWVERIFY.

reject a detection based on this criteria.

Method	GMU-Test		W-RGBD	
	mAP	max Precision	mAP	max Precision
FLOWVERIFY	0.422	0.939	0.278	0.822
FLOWVERIFY-SIMOBJ	0.412	0.906	0.316	0.732
FLOWVERIFY-FRIGIDITY	0.394	0.807	0.235	0.700
FLOWVERIFY-FCOLOR	0.422	0.933	0.278	0.823
FLOWVERIFY-FPRECISION	0.415	0.913	0.280	0.811
FLOWVERIFY-FRECALL	0.413	0.887	0.277	0.755

Table 5.3: Ablation analysis.

5.5 Qualitative Analysis

Figure 5.3 (a, b) shows examples of false positive detections output by the base detector (TDID) but filtered out using FLOWVERIFY. Example (a) is a detection with high FRIGIDITY, FPRECISION, and FRECALL scores, but with a low FCOLOR score. In such cases, the detection is of an object with very different color and texture from the target object, resulting in low color similarity. Example (b) shows a detection with high FCOLOR and FPRECISION scores but with a low FRIGIDITY score. Here, the matched colors are fairly similar; however, the correspondences cannot be derived from a rigid body transformation, and hence the detection has a low FRIGIDITY score. Examples (c) and (d) in Figure 5.3 show true positive predictions output by the base detector which pass FLOWVERIFY. We can see that the flow quality is also quite good by the correct matching of features across the target and the bounding

box area.

5.6 Qualitative Comparison with SIFTVerify

In this section, we perform a comparative and qualitative analysis of our method FLOWVERIFY with the SIFTVERIFY baseline. We will first show examples where SIFTVERIFY succeeds in filtering out false positives from a base detector and where it is successful in retaining true positives. Then we will analyze some failure cases of SIFTVERIFY, where our model is able to filter false positives and retain correct detections whereas SIFTVERIFY is not.

It is important to note that since we are interested in high-precision detection, we will analyze detections which were given a very high confidence score (> 0.99) by the original detector on the GMU Kitchens [19] test set. In order to improve detector performance in the high precision regime (the leftmost parts of the PR-curve), the system must be able to filter out high-confidence false positives whilst not rejecting too many true detections with high confidence.

As a reminder, FLOWVERIFY uses five tests – SIMOBJ, FCOLOR, FRIGIDITY, FPRECISION and FRECALL . In order to pass a test, there must exist at least one viewpoint of the target object whose SIMOBJ, FCOLOR, FRIGIDITY, FPRECISION and FRECALL scores are all above their respective thresholds. The thresholds that were tuned on YCB-Video train for FCOLOR, FRIGIDITY, FPRECISION and FRECALL turn out to be 0.5, 0.9, 0.9, and 0.3 respectively, with parameters for SIMOBJ being $\eta_{iou} = 0.5$, $\eta_{diff} = 0.0$.

Similarly, SIFT filters detections by computing the number of keypoint matches between the target image and cropped scene image. If the SMATCHES and SPRECISION are above their corresponding thresholds, the detection is deemed to be true by the SIFT baseline. When tuned on YCB-Video Train, the thresholds for SMATCHES and SPRECISION turns out to be 30 and 0.9 respectively.

5.6.1 Visualization Format

The visualization format is consistent across all Figure 5.4, Figure 5.5, Figure 5.6, and Figure 5.7 – the top image shows a visualization of matching by FLOWMATCHNET

and filtering by FLOWVERIFY, whereas the bottom image shows matching using SIFTVERIFY. Each visualization consists of two images. On the left is a ‘target image’ – it is an image of the object being detected. If the detection passes the FLOWVERIFY or SIFT verification tests, the displayed viewpoint is the one that passes all tests and has the largest product of scores, which denotes the “best” viewpoint of the object that is able to be matched by each of the methods; otherwise, the canonical viewpoint of the object is displayed in the visualization.

In each visualization, 20 randomly-chosen pixel-wise mappings are depicted. In the case of SIFT, sometimes there are less than 20 keypoint matches found. In such cases, all keypoint matches are depicted.

5.6.2 SIFT Successful Cases

In this section, we analyze cases where SIFT matching successfully filters false positives while retaining true detections.

False Positives rejected by SIFT and FlowVerify

Figure 5.4 contains four examples that denote false detections that are successfully filtered out by both methods. In these examples, the homography based on flow computed by FLOWMATCHNET produces many outliers, and hence are rejected by the FRIGIDITY test. This happens because when the object are distinct and do not match, in which case the predicted flow can be arbitrary. SIFT also successfully rejects all these examples, as it produces very few (< 30) keypoint matches.

True Positives accepted by SIFT and FlowVerify

Figure 5.5 contains four examples that represent true detections successfully retained by both methods. In these examples, FLOWMATCHNET and SIFT matches seem nearly perfect. This section demonstrates that SIFT features can be effective in object matching, and our implementation of SIFT is a reasonably strong baseline. However, we show below that SIFT can also fail on many cases.

5.6.3 SIFT Unsuccessful Cases

In this section, we analyze cases where SIFT matching fails to filter false positives or retain high-confidence true detections. We illustrate how FLOWVERIFY successfully handles these cases, to highlight the strengths of our method over the SIFT baseline.

False Positives accepted by SIFT but rejected by FlowVerify

In Figure 5.6, the four examples represent false detections that SIFT fails to filter out, but are successfully filtered out our method. In Figure 5.6(a), the target object is different from the object in the bounding box, but they have the same text in their logos. SIFT fails in this case since it can find enough keypoint matches just in the logo area. FLOWVERIFY handles this case successfully since it relies on dense pixelwise correspondence. As the object in the cropped scene is different from that in the target image, FLOWMATCHNET cannot find a rigid transformation that matches color for all pixels, thereby not passing the FRIGIDITY test.

Figure 5.6(b, c, d) illustrate another difference between dense and sparse pixel-wise matching. In these examples, the bounding boxes are incomplete and the cropped scene only contains a part of the target object. As SIFT only needs to match a few keypoints, even if small parts of two objects seem to match, it is still enough for SIFT to accept it as a correct detection. FLOWVERIFY combined with FLOWMATCHNET successfully rejects these as false detections since FLOWMATCHNET computes dense matches across the entire object, resulting in flow fields that have a low FRIGIDITY score.

True Positives rejected by SIFT but retained by FlowVerify

Figure 5.7 contains four examples that denote true positives, which are successfully retained by FLOWVERIFY but are incorrectly *rejected* by SIFTVERIFY. In all examples (a, b, c, d) in Figure 5.7, FLOWMATCHNET produces good-quality pixel-wise correspondence, leading to high FCOLOR, FRIGIDITY, FPRECISION and FRECALL scores. However, for most examples, SIFT can produce very few keypoint matches. Our experiments show that SIFT struggles to perform reliable keypoint matching in our setting where there can be drastic changes in viewpoints, lighting conditions

CHAPTER 5. EXPERIMENTS

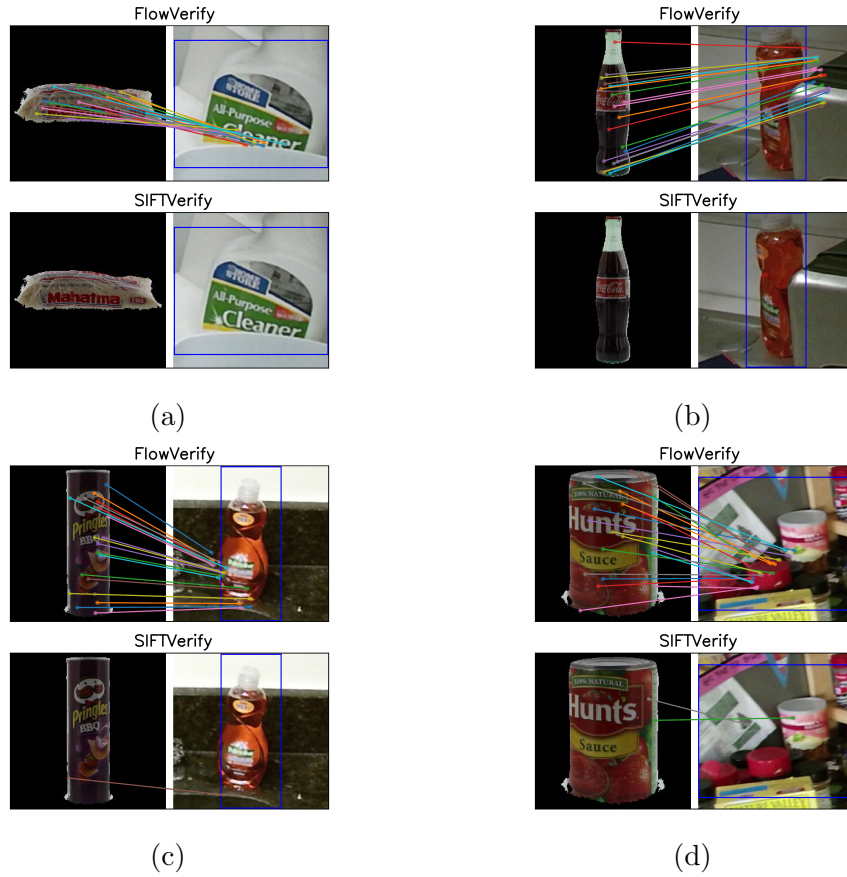


Figure 5.4: False detections which are successfully filtered out by both FLOWVERIFY and SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.

and even occlusions. We show that in such scenarios, predicting a dense pixel-wise correspondence and designing subsequent verification tests can improve instance detectors towards high-precision and verifiable recognition.

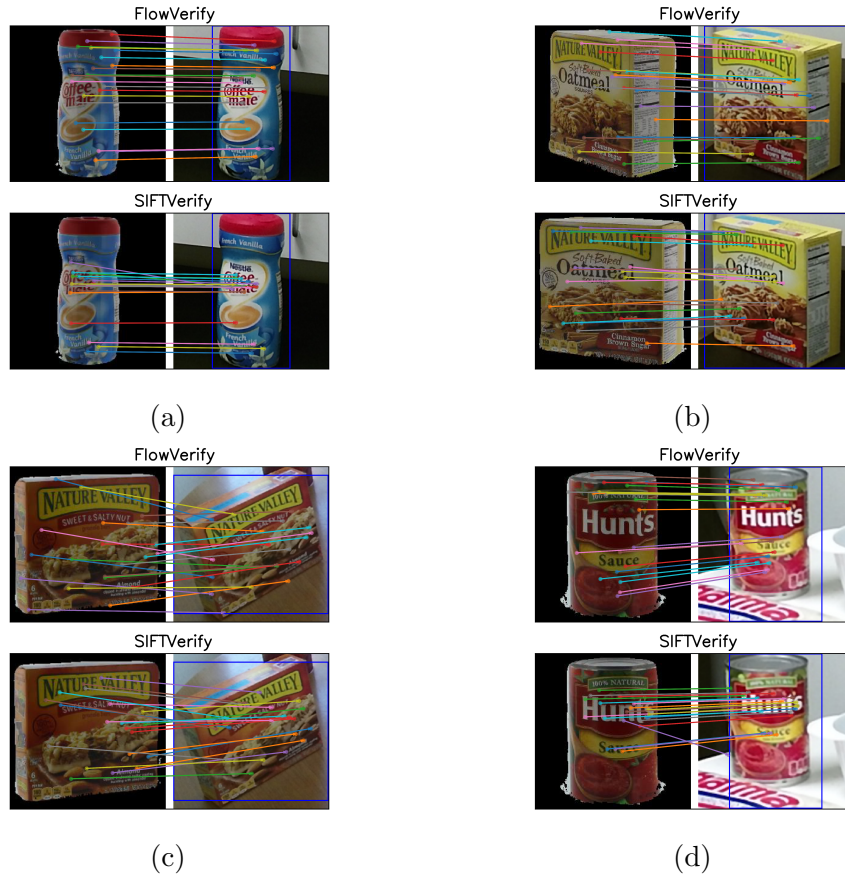


Figure 5.5: True detections that are successfully retained out by both FLOWVERIFY and SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.

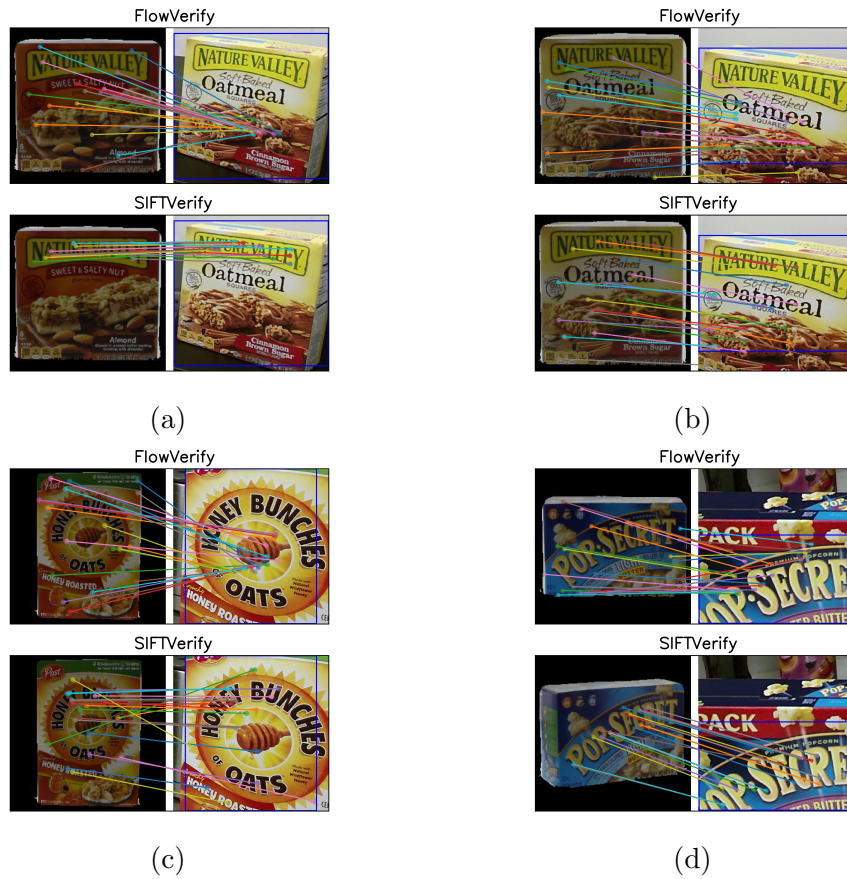


Figure 5.6: False detections which are incorrectly retained by SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.

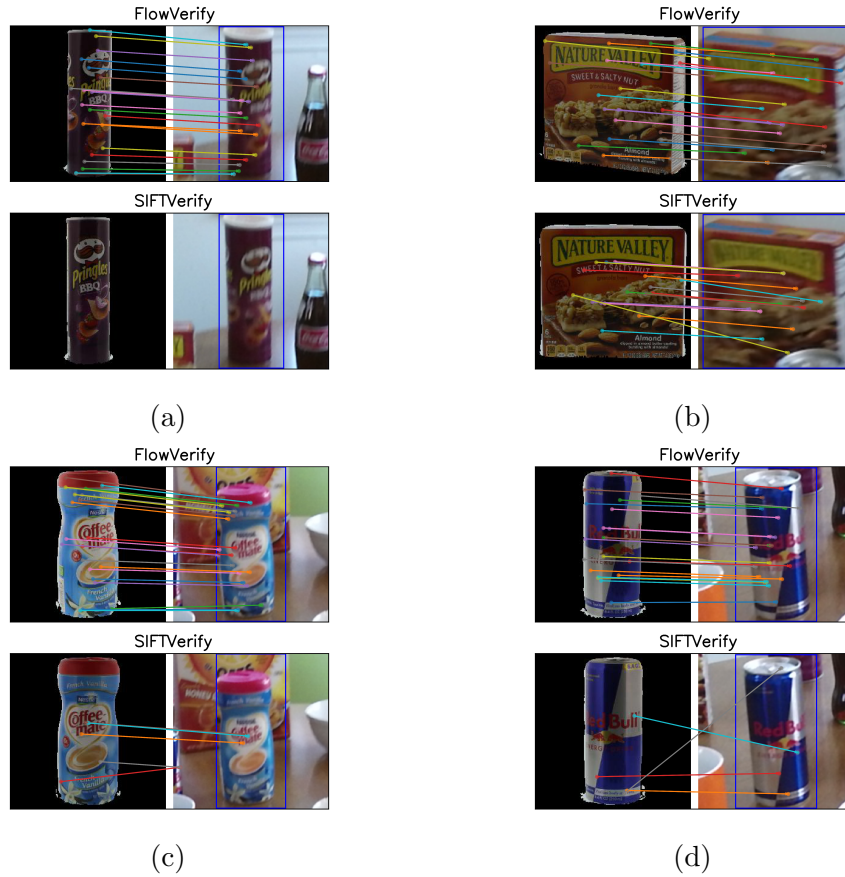


Figure 5.7: True detections that are incorrectly filtered by SIFTVERIFY. In each image, top: FLOWVERIFY, bottom: SIFTVERIFY.



(a) A false positive of FLOWVERIFY. (b) A false negative of FLOWVERIFY.

Figure 5.8: Failure cases of FLOWVERIFY. (a) is an incorrect detection accepted by our method; (b) is a correct detection rejected by our method; Blue box in each example denotes the bounding box prediction from instance detector.

5.7 Failure Cases

In this section, we present some failure cases. We visualize two different kinds of failure cases in Figure 5.8:

Figure 5.8(a) is a false positive example of our method, meaning that it is an incorrect detection accepted by FLOWVERIFY, which means that it passes all of our verification tests. This is actually an example where Assumption 1 is broken, since the detection does not include any object in our dataset. Unfortunately, our theoretical framework is currently unable to handle this case, where the detection does not contain any object in our dataset, but has similar appearance to some object in the dataset after applying a rigid transformation. Nevertheless, we note that such scenario is intrinsically challenging to any purely vision based 2d instance detection tasks, and perhaps need to be handled using extra information or interactions with the environment. For example, if a robot is allowed to take multiple images of the area of interest from multiple viewpoints, we will be more likely to identify the example in Figure 5.8(a) as a false positive. One possible future work along this line is to apply our method to these scenarios, and verify using detections from multiple perspectives of the same area.

Figure 5.8(b) is a false negative example of our method, which is a correct detection rejected by our method. Currently we have no guarantees about true positives, so this failure is sort of expected. One of a possible future work is to extend the theoretical framework to true positive detections. Intuitively, if there is some guarantee in the

accuracy of the predicted correspondences, our method should be able to benefit from it in terms of recall. We can explore this theoretical extension in the future to further improve in our verification system. We also observe in Figure 5.8(b) that the instance is rejected since a significant part of the instance is occluded in the detection. Currently, our framework does not take account of occlusions. However, since occlusions are common in many real world scenarios, we can extend our work to handle it in the future.

CHAPTER 5. EXPERIMENTS

Chapter 6

Future Work

To address failure case discussed in Section 5.7, we can extend the framework to true positives, handle occlusions in our framework, and apply our framework to an interactive environment to collect and verify multiple detections of the same object. Additionally, we can further improve the network for predicting dense pixel correspondences. On one hand, increasing accuracy of the correspondence predictions is related to retaining more true positives; on the other hand, we can extend FLOWMATCHNET to predict occlusions in addition to correspondences as one way to handle occlusions. Furthermore, our current framework only considers rigid objects. If we want to apply similar verification approaches to deformable objects, we need to devise other metric to check validity of correspondences instead of the rigidity test we are currently using. We can further extend the work to rigid multibodies like robot, which has rigid parts but also non-rigid joints. We can also extend the current method to class-level detection verification.

CHAPTER 6. FUTURE WORK

Chapter 7

Conclusion

We have proposed a method to combine machine learning based detection and correspondence matching with non-learning based verification tests to increase the accuracy of an existing instance-detection system in the high-precision regime (without reducing overall detection performance). The verification tests are based on dense-pixel correspondences computed between the detection and template images; we reduce the confidence of any detection that does not pass these tests, thereby rejecting many false positives. Our system is grounded in a novel theoretical framework that we prove leads to no false positives, under certain assumptions. Furthermore, we use our method in a one-shot fashion, applying our approach to a novel set of objects at test time without finetuning. For future work, we would like to extend the current framework to handle occlusions, and introduce guarantees about true positives with respect to accuracy of matching predictions. We hope that our system will be useful for robotic systems that need reliable performance for high confidence detections.

CHAPTER 7. CONCLUSION

Bibliography

- [1] Siddharth Advani, Brigid Smith, Yasuki Tanabe, Kevin Irick, Matthew Cotter, Jack Sampson, and Vijaykrishnan Narayanan. Visual co-occurrence network: using context for large-scale object recognition in retail. In *2015 13th IEEE Symposium on Embedded Systems for Real-time Multimedia (ESTIMedia)*, pages 1–10. IEEE, 2015. [2.1](#)
- [2] Phil Ammirato, Patrick Poirson, Eunbyung Park, Jana Košecká, and Alexander C Berg. A dataset for developing and benchmarking active vision. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 1378–1385. IEEE, 2017. [5.1.1](#)
- [3] Phil Ammirato, Cheng-Yang Fu, Mykhailo Shvets, Jana Kosecka, and Alexander C Berg. Target driven instance detection. *arXiv preprint arXiv:1803.04610*, 2018. [2.1](#), [4.1](#), [5.1.1](#)
- [4] Peter L Bartlett and Marten H Wegkamp. Classification with a reject option using a hinge loss. *Journal of Machine Learning Research*, 9(Aug):1823–1840, 2008. [2.2](#)
- [5] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *ECCV*, pages 404–417. Springer, 2006. [1](#), [2.1](#)
- [6] Ipek Baz, Erdem Yoruk, and Mujdat Cetin. Context-aware hybrid classification system for fine-grained retail product recognition. In *2016 IEEE 12th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP)*, pages 1–5. IEEE, 2016. [2.1](#)
- [7] C Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on information theory*, 16(1):41–46, 1970. [2.2](#)
- [8] Christopher B Choy, JunYoung Gwak, Silvio Savarese, and Manmohan Chandraker. Universal correspondence network. In *Advances in Neural Information Processing Systems*, pages 2414–2422, 2016. [2.3](#)
- [9] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Boosting with abstention. In *Advances in Neural Information Processing Systems*, pages 1660–1668, 2016. [2.2](#)

- [10] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with rejection. In *International Conference on Algorithmic Learning Theory*, pages 67–82. Springer, 2016. [2.2](#)
- [11] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2758–2766, 2015. [2.3](#), [5.1.3](#), [5.1.3](#), [2](#)
- [12] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *The IEEE international conference on computer vision (ICCV)*, 2017. [2.1](#), [5.1.1](#), [2](#)
- [13] Philipp Fischer, Alexey Dosovitskiy, Eddy Ilg, Philip Häusser, Caner Hazirbas, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. *CoRR*, abs/1504.06852, 2015. URL <http://arxiv.org/abs/1504.06852>. [2.3](#)
- [14] 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. [4.1](#)
- [15] Annalisa Franco, Davide Maltoni, and Serena Papi. Grocery product detection and recognition. *Expert Systems with Applications*, 81:163–176, 2017. [2.1](#)
- [16] Yonatan Geifman and Ran El-Yaniv. Selectivenet: A deep neural network with an integrated reject option. *arXiv preprint arXiv:1901.09192*, 2019. [2.2](#)
- [17] Weidong Geng, Feilin Han, Jiangke Lin, Liuyi Zhu, Jieming Bai, Suzhen Wang, Lin He, Qiang Xiao, and Zhangjiong Lai. Fine-grained grocery product recognition by one-shot learning. In *2018 ACM Multimedia Conference on Multimedia Conference*, pages 1706–1714. ACM, 2018. [2.1](#)
- [18] Georgios Georgakis, Md Alimoor Reza, Arsalan Mousavian, Phi-Hung Le, and Jana Kosecka. Multiview rgb-d dataset for object instance detection. *arXiv preprint arXiv:1609.07826*, 2016. [1](#)
- [19] Georgios Georgakis, Md Alimoor Reza, Arsalan Mousavian, Phi-Hung Le, and Jana Kosecka. Multiview rgb-d dataset for object instance detection. *arXiv preprint arXiv:1609.07826*, 2016. [5](#), [2](#), [5.6](#)
- [20] Georgios Georgakis, Arsalan Mousavian, Alexander C. Berg, and Jana Kosecka. Synthesizing training data for object detection in indoor scenes. *CoRR*, abs/1702.07836, 2017. URL <http://arxiv.org/abs/1702.07836>. [2.1](#)
- [21] Georgios Georgakis, Arsalan Mousavian, Alexander C Berg, and Jana Kosecka. Synthesizing training data for object detection in indoor scenes. *arXiv preprint*

- arXiv:1702.07836*, 2017. 1
- [22] Marian George and Christian Floerkemeier. Recognizing products: A per-exemplar multi-label image classification approach. In *European Conference on Computer Vision*, pages 440–455. Springer, 2014. 2.1
- [23] Yotam Hechtlinger, Barnabás Póczos, and Larry Wasserman. Cautious deep learning. *arXiv preprint arXiv:1805.09460*, 2018. 2.2
- [24] David Held, Sebastian Thrun, and Silvio Savarese. Robust single-view instance recognition. In *Robotics and Automation (ICRA), 2016 IEEE International Conference on*, pages 2152–2159. IEEE, 2016. 1, 1
- [25] Stefan Hinterstoisser, Cedric Cagniart, Slobodan Ilic, Peter Sturm, Nassir Navab, Pascal Fua, and Vincent Lepetit. Gradient response maps for real-time detection of textureless objects. *PAMI*, 34(5):876–888, 2012. 1, 2.1
- [26] B.K. Horn and B.G. Schunck. Determining optical flow. In *Artificial Intelligence*, 1981. 2.3
- [27] Xiaowei Huang, Marta Kwiatkowska, Sen Wang, and Min Wu. Safety verification of deep neural networks. *CoRR*, abs/1610.06940, 2016. URL <http://arxiv.org/abs/1610.06940>. 2.2
- [28] Daniel P. Huttenlocher, Gregory A. Klanderman, and William J Rucklidge. Comparing images using the hausdorff distance. *PAMI*, 15(9):850–863, 1993. 1, 2.1
- [29] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. *CoRR*, abs/1612.01925, 2016. URL <http://arxiv.org/abs/1612.01925>. 2.3
- [30] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *IEEE conference on computer vision and pattern recognition (CVPR)*, volume 2, page 6, 2017. 5.1.3, 5.1.3, 2
- [31] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Reluplex: An efficient SMT solver for verifying deep neural networks. *CoRR*, abs/1702.01135, 2017. URL <http://arxiv.org/abs/1702.01135>. 2.2
- [32] Gregory Koch. Siamese neural networks for one-shot image recognition. 2015. 2.1
- [33] K. Lai, L. Bo, X. Ren, and D. Fox. A large-scale hierarchical multi-view rgb-d object dataset. In *2011 IEEE International Conference on Robotics and Automation*, pages 1817–1824, May 2011. doi: 10.1109/ICRA.2011.5980382. 5
- [34] Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox. A large-scale hierarchical

- multi-view rgb-d object dataset. In *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pages 1817–1824. IEEE, 2011. 1
- [35] Kevin Lai, Liefeng Bo, and Dieter Fox. Unsupervised feature learning for 3d scene labeling. In *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, pages 3050–3057. IEEE, 2014. 5, 5.1.1, 2
- [36] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1
- [37] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. In *Advances in Neural Information Processing Systems*, pages 9628–9639, 2018. 5.1.3
- [38] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: single shot multibox detector. *CoRR*, abs/1512.02325, 2015. URL <http://arxiv.org/abs/1512.02325>. 2.1
- [39] H Christopher Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293(5828):133, 1981. 4.1
- [40] David G Lowe. Object recognition from local scale-invariant features. In *ICCV*, volume 2, pages 1150–1157. IEEE, 1999. 5.3
- [41] David G Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004. 1, 2.1, 5.3
- [42] Michele Merler, Carolina Galleguillos, and Serge Belongie. Recognizing groceries in situ using in vitro training data. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007. 2.1
- [43] Luca Pulina and Armando Tacchella. An abstraction-refinement approach to verification of artificial neural networks. In Tayssir Touili, Byron Cook, and Paul Jackson, editors, *Computer Aided Verification*, pages 243–257, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg. ISBN 978-3-642-14295-6. 2.2
- [44] Luca Pulina and Armando Tacchella. Challenging smt solvers to verify neural networks. *AI Commun.*, 25:117–135, 2012. 2.2
- [45] A Quadros, James Patrick Underwood, and Bertrand Douillard. An occlusion-aware feature for range images. In *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, pages 4428–4435. IEEE, 2012. 1, 2.1
- [46] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *CoRR*, abs/1506.01497, 2015. URL <http://arxiv.org/abs/1506.01497>. 2.1

- [47] Jérôme Revaud, Philippe Weinzaepfel, Zaïd Harchaoui, and Cordelia Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. *CoRR*, abs/1501.02565, 2015. URL <http://arxiv.org/abs/1501.02565>. 2.3
- [48] Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(Mar):371–421, 2008. 2.2
- [49] Arjun Singh, James Sha, Karthik S Narayan, Tudor Achim, and Pieter Abbeel. Bigbird: A large-scale 3d database of object instances. In *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, pages 509–516. IEEE, 2014. 5, 2
- [50] Alessio Tonioni, Eugenio Serra, and Luigi di Stefano. A deep learning pipeline for product recognition on store shelves. *CoRR*, abs/1810.01733, 2018. URL <http://arxiv.org/abs/1810.01733>. 2.1, 2.2
- [51] Yi-Hsuan Tsai, Ming-Hsuan Yang, and Michael J. Black. Video segmentation via object flow. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 2.3
- [52] Gül Varol and Rıdvan Salih. Toward retail product recognition on grocery shelves. page 944309, 03 2015. doi: 10.1117/12.2179127. 2.1
- [53] Ziang Xie, Arjun Singh, Justin Uang, Karthik S Narayan, and Pieter Abbeel. Multimodal blending for high-accuracy instance recognition. In *IROS*, pages 2214–2221. IEEE, 2013. 1
- [54] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. Flow-guided feature aggregation for video object detection. *CoRR*, abs/1703.10025, 2017. URL <http://arxiv.org/abs/1703.10025>. 2.3