

# Object Detection using SIFT

Yuan Tjiam

August 15, 2022



## Abstract

Surgical teams use instrument counts to prevent leaving unintended objects in patients. This is done manually, but could potentially be done through computer vision software. This paper presents a proof of concept for detecting instruments in the operating room with the Scale Invariant Feature Transform (SIFT). The SIFT algorithm is explored and tested on a variety of household appliances to substitute medical instruments. The algorithm responds differently to metal objects compared to matte objects and has room for many improvements. Further research on run time and multi object images is necessary. The proof of concept is considered successful when not taking run time into account.

**Keywords**— Scale Invariant Feature Transform, Computer vision, Instrument detection, Operating room

# Acknowledgements

I want to thank my Professor John van den Dobbelen for giving me another chance. And I want to thank my supervisor Rick Butler for our illuminating weekly discussions. Without those I would not have been able to finish. Finally I want to show my appreciation to my parents and my girlfriend for their endless support.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Method</b>	<b>4</b>
2.1	Multiple images per instrument . . . . .	4
2.2	Generating results . . . . .	4
2.3	Camera setup . . . . .	4
2.4	Running SIFT . . . . .	4
<b>3</b>	<b>Theoretical Framework</b>	<b>6</b>
3.1	Overview . . . . .	6
3.2	Scale space extrema detection . . . . .	8
3.3	Keypoint localization and removal . . . . .	8
3.4	Keypoint orientation . . . . .	9
3.5	Local Image descriptor . . . . .	9
3.6	Keypoint matching . . . . .	10
3.7	Keypair clustering . . . . .	10
3.8	Keypair parameters . . . . .	11
3.9	Hough Transform . . . . .	11
3.10	Computing models . . . . .	12
3.11	Probability of model . . . . .	13
3.12	Duplicate removal by NMS . . . . .	13
<b>4</b>	<b>Results</b>	<b>14</b>
<b>5</b>	<b>Discussion</b>	<b>16</b>
5.1	Matte and metal . . . . .	16
5.2	Similarity and affine transform . . . . .	16
5.3	Mirrored images . . . . .	16
5.4	Hyper parameters . . . . .	16
5.5	Precision and Recall . . . . .	17
5.6	Requirements and limitations . . . . .	18
5.7	Future research . . . . .	18
<b>6</b>	<b>Conclusion</b>	<b>19</b>
<b>A</b>	<b>Variables and constants</b>	<b>20</b>
<b>B</b>	<b>Affine</b>	<b>21</b>
<b>C</b>	<b>Probability</b>	<b>21</b>
<b>D</b>	<b>Code</b>	<b>21</b>

# 1 Introduction

The operating room (OR) is by its nature a hazardous place. Many things can and do go wrong, resulting in increased cost and reduced quality of care. One method of preventing mistakes in the OR is to have checks and redundancy checks. These can vary from pre-procedure anaesthetic and allergy checks to instrument counts [9].

Leaving an unintended sponge, needle or instrument in a patient is a mistake that occurs with a frequency ranging from 1 in 5000 to 1 in 7000 [1, 3]. The instrument count is a standard procedure that takes place pre- and post- surgery. It has been standardized by the World Health Organization [9]. Sponge, needle and instrument counts are done at all procedures that risk the possibility of an unintended foreign object being retained in the patient. These objects can lead to inflammation, obstruction, perforation sepsis and death [11]. The risk of foreign object retention increases for emergency operations, a sudden change in surgical procedure and patients with a high body-mass index [4].

The count takes several minutes and is done multiple times, often by multiple members for redundancy. Reconciling miscounts is done via X-rays and results in increased OR time and costs.

If the count could be automated in such a way that it is both fast and reliable, it frees up multiple members of the team and OR time decreases along with cost. Automating the count could be done using Radio Frequency IDentification (RFID) tracking [11], which requires trackers on all objects entering the body. Another method would be to use computer vision software to track objects without adjusting instruments and with minimal change to operating room procedure.

There are different options for computer vision software. In this paper we chose to work with the Scale Invariant Feature Transform (SIFT) by Lowe [7]. Simi-

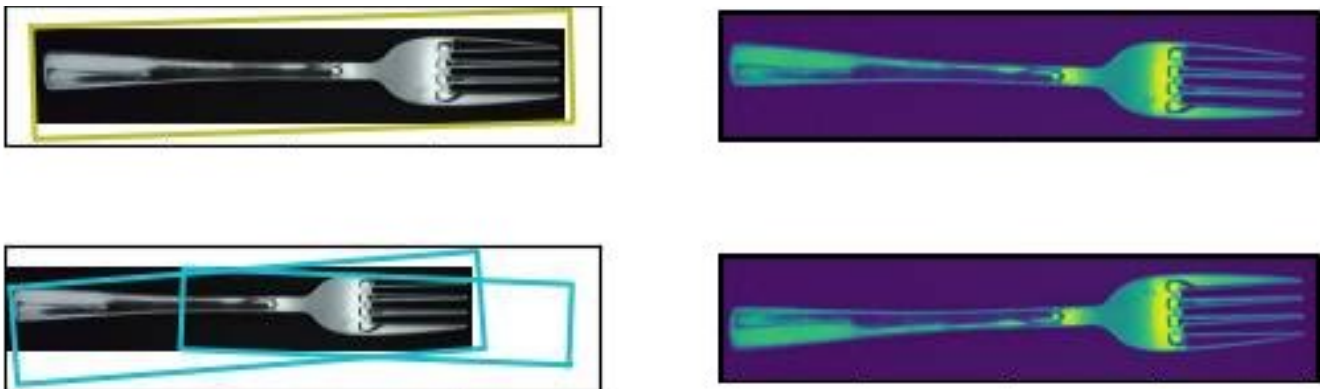
lar software such as Oriented Features from Accelerated Segment Test and Rotated Binary Robust Independent Elementary Features (ORB) is faster but focuses on the center of the image. Software Speeded Up Robust Features (SURF) outperforms SIFT only in noisy images, but is less accurate when it comes to clear images [2]. Another option would be to use deep learning. Deep learning emulates the working of the human brain. It requires a lot of data to train the software. After it has been trained the inner mechanisms are often difficult to understand. In this thesis, we focus on accuracy and reliability over speed and we use sharp and clear images. We also want to be able to understand what effect different variables have on our outcome. For these reasons, we chose to work with SIFT [5, 6, 7, 8].

SIFT compares keypoints in images and matches these to a database. This software is invariant to scale, translation, rotation and illumination. Additionally, it can detect objects rotated in depth up to 20 degrees. These properties make SIFT suitable as instruments are not necessarily placed neatly on a tray under the same lighting conditions pre- and post- operation.

The research question of this thesis is:

”Can the Scale Invariant Feature Transform (SIFT) algorithm reliably detect objects in an operating room setting?”

This thesis will serve as a proof of concept.



**Figure 1:** An example of a true positive identification (upper) and a false positive identification(lower). A false negative identification is the same but without any colored boxes.

## 2 Method

In an operating room setting, a health practitioner would take a picture of all instruments on a surface at various stages of surgery, but at least before and after. These are the scene images. The algorithm would identify the number and type of all objects in a scene image from a previously defined database of instrument images. Non-detections and false positives increase the likelihood of harm to a patient. Which is why we define a correct identification as the detection of the correct type of instrument with no false positives and a 70% overlap of area between the scene and instrument photo. This 70% was chosen because it seemed reasonable.

Positive identifications are shown by colored boxes around the detected instrument in the scene image as in figure 1, upper left. In the case of comparing single instrument scene images, if there are multiple boxes per image, it's determined to be a false positive. All detections are verified by eye.

In an operating room setting, run time would be essential to consider as well. In that setting, the run time of the program should be no more than the current count procedure time. However, there are too many variables influencing run time of the program to consider the run time with regards to the research question.

### 2.1 Multiple images per instrument

Single image comparison often would not yield positive identifications. If objects are rotated too much along the wrong axis, the algorithm loses the ability to detect it. As instruments are not always neatly placed in real life situations, there was a need for a comprehensive way of detecting instruments regardless of lighting or angle. We decided to make a map of images per instrument that vary in several parameters.

These parameters are instrument rotation, camera angle, lighting angle and lighting intensity. All images per instrument were added to an instrument folder. Then every instrument in a scene has multiple different instrument images to compare to.

The program is also not mirror-invariant. To solve this, images of instruments were mirrored and both versions of the instrument image were used.

### 2.2 Generating results

The program was tested on household appliances that are similar to surgical tools in shape and size shown in table 1. Preliminary tests of SIFT showed a large difference in outcome between metal and matte objects. This is why we divided the instruments into two categories shown in table 1.

**Table 1:** Objects that were used to test the program.

Glare	Matte
Fork	Wooden spatula
Butter knife	Black spatula
Spoon	Black soup spatula
Ice cream scoop	Blue spatula
Whisk	

### 2.3 Camera setup

The objects were photographed in a dark room where light sources were controlled. Each object was placed individually underneath the camera on either a white or black background, depending on the colour of the object. The light never shone directly on the objects, but always on a white wall to imitate omnidirectional homogeneous lighting. True homogeneous lighting would decrease the variability of the results and we expect it would decrease number of false positives, as even slight errors in light direction might lead to different keypoints. The light we used was a Dörr SLR-16 Bi-Color Selfie Ring Light, which allows for multiple angles on lightning and camera angle. The camera is from a Huawei P20 Lite phone, with 16 megapixel camera, f/2.2 lens and a 2 megapixel depth sensor. In figure 2 we see the default configuration of the set up. The camera is parallel with the horizontal of the table. The object is placed directly beneath it and parallel with the alignment of the camera. The light points up toward the white ceiling.

For each object, depending on the size of the object, we varied the distance of the camera to the surface to capture the entire object. We then cycled through all light intensities. We took the lowest intensity that was still clear and sharp because over-lighting the image could result in a surplus of keypoints. From this default position we varied light intensities, the object rotation, the orientation of the light source and the tilt of the camera. For each variation in orientation of light and camera tilt, we rotated the object. Table 2 shows the object rotation, light orientation and camera tilt variations. Per object there are 54 photos, 62 including the light intensity variations.

### 2.4 Running SIFT

We took the default configuration of the camera and position of the object as the scene image. All other images from an instrument were put into a folder and served as the instrument images. We ran the program for three variables in two settings resulting in eight configurations. Table 3 shows the type of transform, the cluster threshold and the probability threshold. These are found in section 3 The cluster threshold was chosen to be three and four instead of the minimum because

**Table 2:** Rotation angles of the object, light and camera.

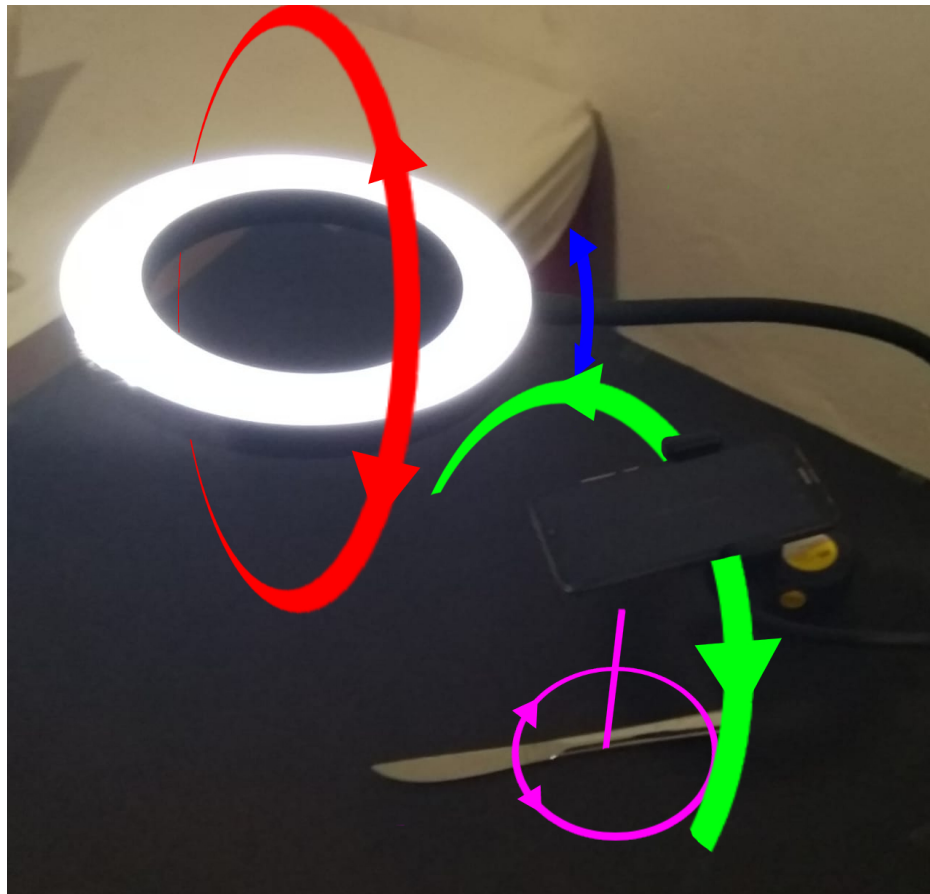
Object rotation	Light orientation	Camera tilt
0°	Up	0°
7.5°	180°	10°
17.5°	90°	27.5°
45°		-12.5°
90°		
135°		
180°		
-7.5°		
-17.5°		

**Table 3:** Settings for SIFT

Transform	Cluster threshold	Probability threshold
Similarity	3	0.98
Similarity	3	0.99
Similarity	4	0.98
Similarity	4	0.99
Affine	3	0.98
Affine	3	0.99
Affine	4	0.98
Affine	4	0.99

preliminary tests showed an overabundance of false positives when the cluster threshold was two.

All instrument images of every object in table 3 were compared to their scene images. Photos were manually checked whether it was a true positive, a false negative or a false positive. False positives were defined as having multiple identifications within the same image.



**Figure 2:** The setup used to take pictures of instruments. The ring light can rotate (red) and tilt (blue). The camera can be tilted so perspective (skewed) photos can be taken (green). Finally the object itself can be rotated (pink). In this local reference system the direction of the knife is 0° for the object rotation and light orientation. A light orientation towards the ceiling is defined as up. The angle for the tilt is defined as 0° directly above the object, with -90° and 90° when the camera hits the table along the green arrows.

### 3 Theoretical Framework

The algorithm used in this paper is the SIFT algorithm based on Lowe’s work [7][6]. Generally, SIFT compares 2D images and detects objects they share. In this application of SIFT, one scene image is compared to many instrument images. An instrument image contains the object that must be detected. The scene image potentially contains it as well. For each unique instrument we use a range of reference images to increase the probability of detection of the instrument in the scene image. The algorithm is coded in Python using integrated development environment Spyder version 5.

Many of the values chosen for the parameters in this section are experimentally determined by Lowe [7, 6]. Many of these could be chosen to have a different value. An oversight of these with more in depth explanations can be found in appendix A. We decided to focus on cluster threshold and probability threshold as the parameters to tune the algorithm with. We chose the cluster threshold because it has a straightforward effect and because the minimum step difference of the cluster threshold is already significant on the performance of the algorithm. We chose the probability threshold because it is one of the final tuning opportunities. It also has a straightforward effect.

SIFT consists of a number of steps, to be repeated for every instrument and scene image. The goal of the algorithm is to create a singular model for each instrument present in the scene that transforms the coordinates of the instrument image into scene image coordinates.

To clearly explain SIFT, example instrument image of figure 3 will be matched onto example scene image of figure 4.

### 3.1 Overview

This is a very short overview of how the algorithm works. Sections 3.2 through 3.12 explain it more in depth.

First the scene and instrument images are loaded. After loading, every subsequent step builds on the previous one. First keypoints are evaluated for both the instrument and scene image based on their position. Data based on the immediate area around every keypoint is added to these keypoints. Keypoints between the instrument and scene are matched into keypairs. Similar keypairs are clustered. From every cluster, a model is computed. A model maps the instrument onto the scene image. The models are evaluated and filtered. The remaining models are the detections in the scene.

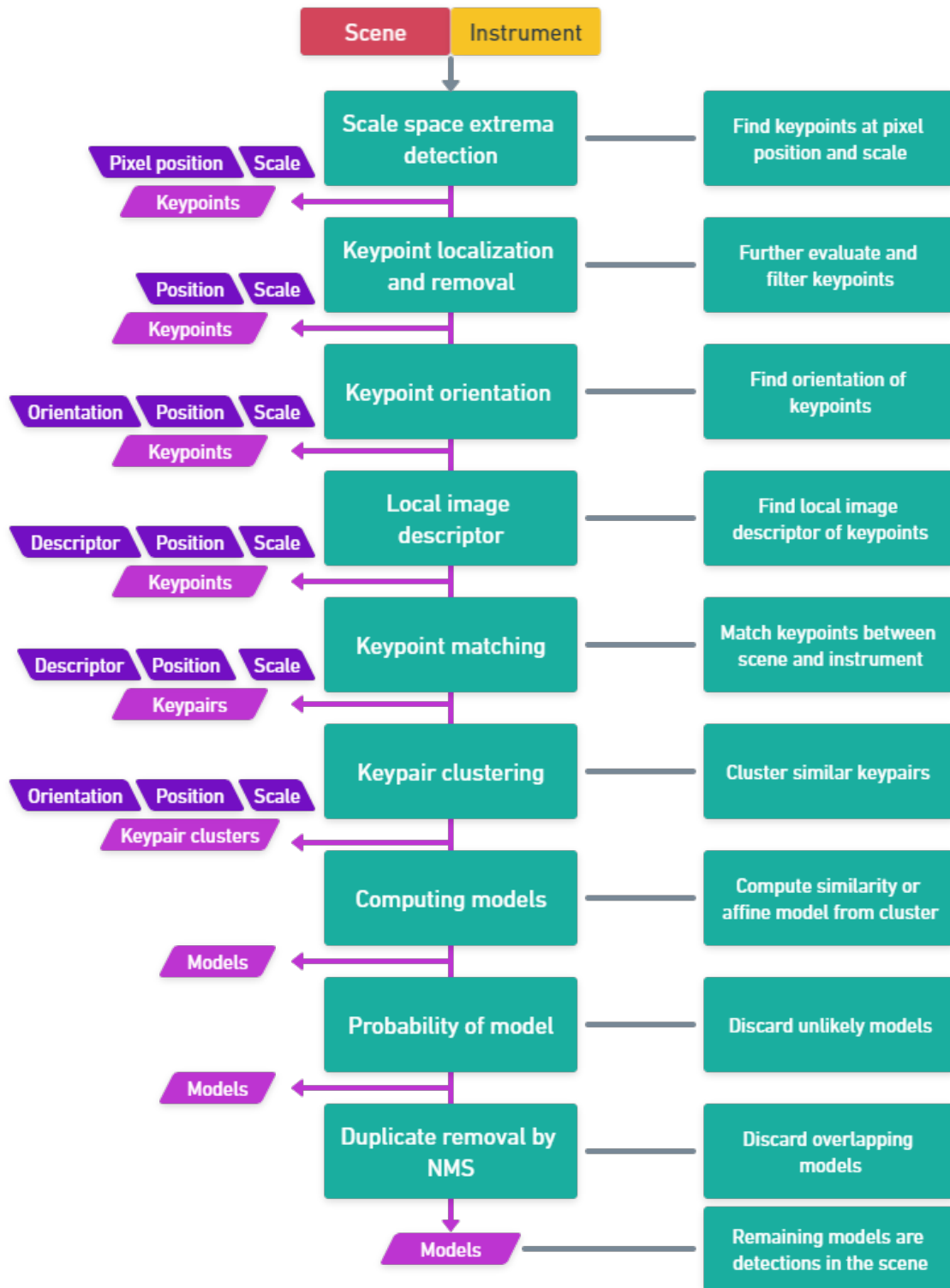
Figure 5 shows the steps in order. The middle blocks are the steps. On the right these are summarized. On the left the information that’s passed on from step to step is shown. The purple blocks are the contents of the pink blocks.



Figure 3: This image of a fork will be the example instrument image to illustrate the steps of SIFT.



Figure 4: This image of a fork will be the example scene image to illustrate the steps of SIFT



**Figure 5:** A general overview of SIFT. A scene image and instrument image are loaded into the program at the top. If done correctly the final product is a model that maps the instrument image onto the scene image.

### 3.2 Scale space extrema detection

SIFT uses extrema detection to define keypoints because this method is rotation invariant and is highly efficient [5]. The image is blurred by convolving it with Gaussians of different standard deviation values  $\sigma$ . The blurred images are called scales. The different  $\sigma$  values are separated by a constant factor  $k$ . The Gaussians have standard deviations  $\sigma_0, \sqrt{2}\sigma_0, 2\sigma_0$ , etc.

As this trend continues, more computational time is needed as a higher  $\sigma$  means that a larger area in the image must be evaluated for blurring per individual pixel. However,  $k$  is chosen to be  $\sqrt{2}$ . Because of this, every fourth scale can be replaced by a downsampling of the image by a factor of 2 (i.e. taking every second pixel of an image). As shown on the left side in figure 6, this downsample is original of the second scale. This saves on computational time and has the same accuracy [7]. A set of scales before downsampling is called an octave

because it is done with a factor of 2. Lowe recommends 4 octaves [7], which is what we will be using.

In each octave the scales are stacked on top of another. Adjacent scales are then subtracted from each other, with 5 scales giving 4 Difference of Gaussians (DoGs) as shown on the right side of figure 6. This is done for each scale. Of these 4 DoGs, two are the middle layers. Every pixel on these two middle layers except for the edge pixels have 26 neighbors, 9 on the layer above them, 8 on their own layer and 9 on the layer below. Figure 7 shows the pixel marked 'x' as an example. If the value of this pixel is larger or smaller than all its neighbours, it is a local maximum or local minimum and defined as a keypoint.

### 3.3 Keypoint localization and removal

Now that we have a set of keypoints, we want to make them more descriptive to be able to match them later

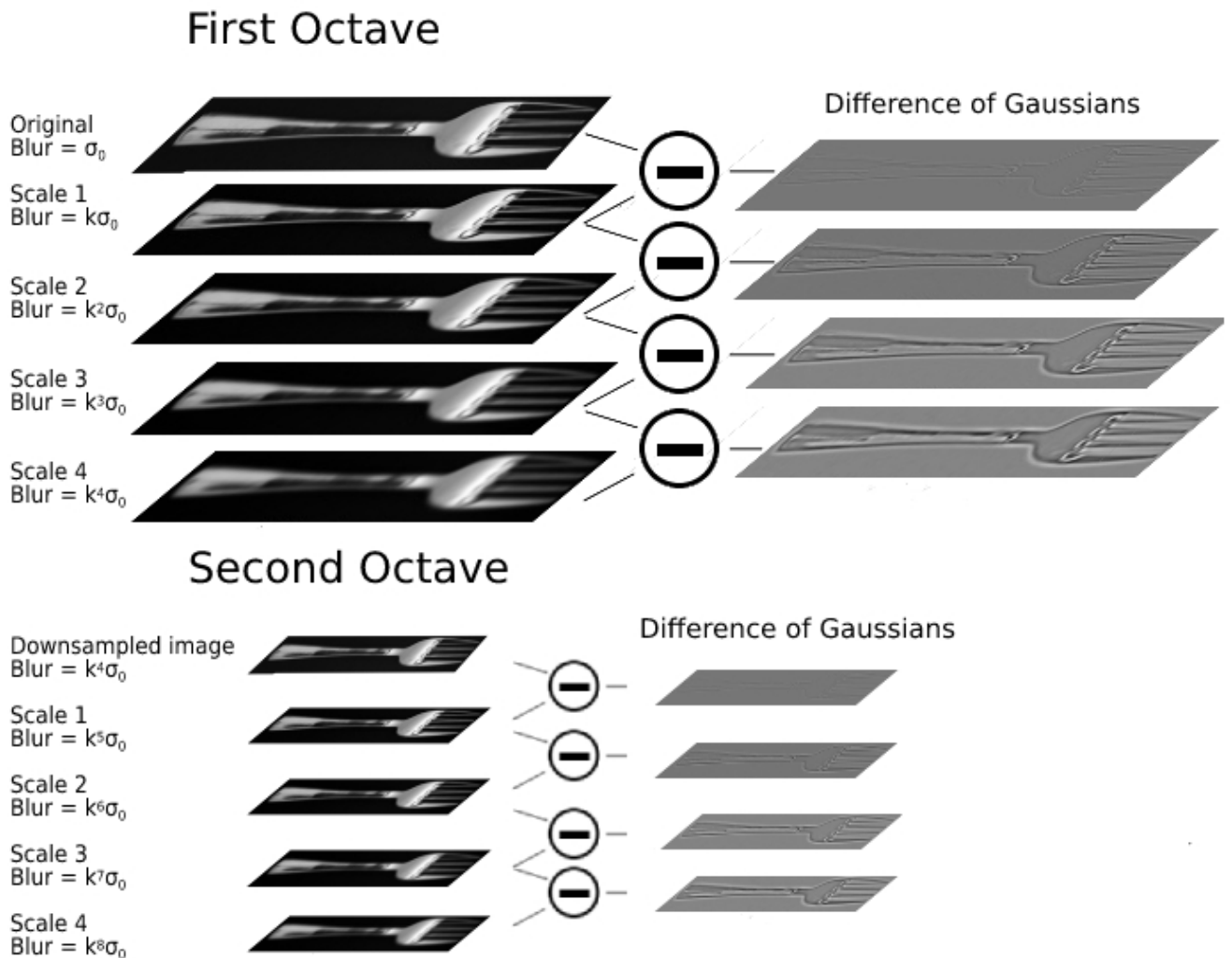
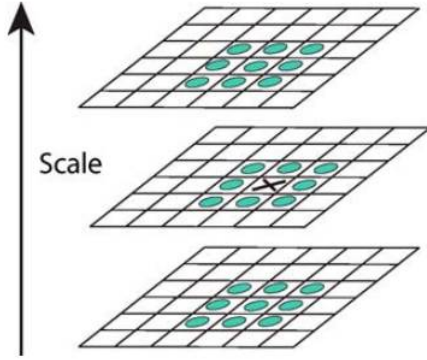


Figure 6: The scales are images convolved with  $k\sigma_0$ . Every second scale is downsampled to become the first of the next octave.





**Figure 7:** The middle of the middle layer is evaluated. If all 26 surrounding pixels are collectively darker or brighter, the pixel is defined as a keypoint. Image taken from Lowe, 2004 [7].

on whilst removing keypoints that are not robust across images due to, for example, noise. Now, the keypoints possess a discrete location and a scale.

Pixels exist in discrete space. However, the positions of local maxima and minima of the pixels can be evaluated in continuous space. This makes them more descriptive. We can find these positions by evaluating the derivative in every direction, where a higher derivative means that the extrema leans toward that direction. Additionally, it allows for a value at the place of the extrema, giving an opportunity to filter keypoints based on this value compared to a threshold value. Low values are extrema with low contrast which are sensitive to noise [8][7]. This value is discarded for all keypoints as it will no longer be of use.

Keypoints at the edges of the image are common. But as it is the edge of the image, local extrema are much more likely and much less reliable. These are eliminated by evaluating the derivative in the x direction and contrasting it with the y direction. At the edge, one of these is going to be big whilst the other will be small. If the ratio between the two is smaller than the edge ratio (appendix A), the keypoint is discarded.

### 3.4 Keypoint orientation

So far, each keypoint has a scale, a continuous location and a value. Now we assign an orientation. The closest

Gaussian smoothed image to the scale of each keypoint is taken. Selecting for scale in this manner approaches scale-invariance. With this scale, the magnitude  $m$  and orientation  $\theta$  are computed using pixel differences as in formula 1 and 2.

$$m_{x,y} = \sqrt{(L_{x+1,y} - L_{x-1,y})^2 + (L_{x,y+1} - L_{x,y-1})^2} \quad (1)$$

$$\theta(x,y) = \arctan \left( \frac{L_{x,y+1} - L_{x,y-1}}{L_{x+1,y} - L_{x-1,y}} \right) \quad (2)$$

These formula are applied to each pixel in the image for this scale. Then, for every keypoint, we create an orientation histogram with 36 bins over 360 degrees. The histogram is filled by taking the orientation and magnitude of the pixels within a  $1.5\sigma$  area around the keypoint and adding these to the histogram. The farther away a pixel, the less weight it carries in the histogram. The peak position is then approximated more accurately by fitting a parabola on the two histogram values closest to each peak and taking the position of the maximum. If the second highest peak is within 80% of the highest peak, the keypoint is duplicated with this orientation. Figure 8 shows all 3134 keypoints for our instrument image. In figure 9, the largest circle clearly shows two orientations.

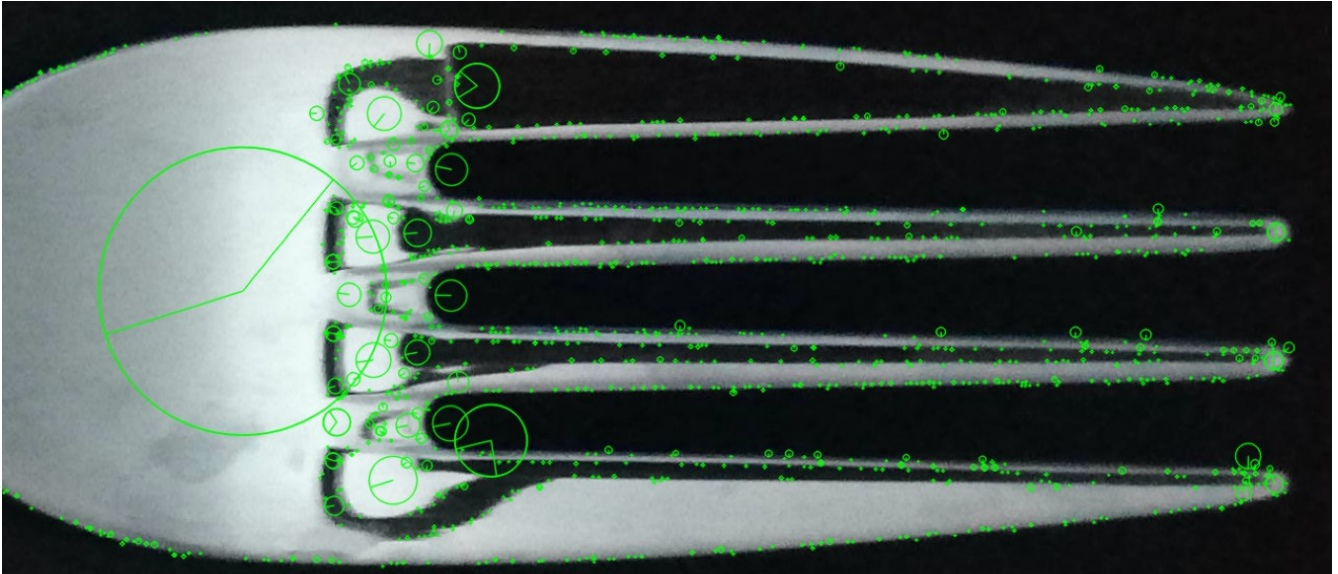
### 3.5 Local Image descriptor

Each keypoint now has an continuous location, a scale and an orientation. Around every keypoint location, for the scale of that keypoint, the orientations and gradient magnitudes of all sample points in that region have been computed in section 3.4. These are rotated relative to the keypoint orientation and weighted according to distance from the keypoint location with a Gaussian weighting function  $\sigma_w$  equal to one half the width of the descriptor window (appendix A). The sample points come from a  $16 \times 16$  area around the keypoint with respect to the rotation. This  $16 \times 16$  area is divided into  $4 \times 4$  regions each with a histogram in 8 directions as shown in figure 10.

To avoid boundary effects, the orientations are tri-linearly interpolated across adjacent histogram bins to ensure continuity across bins. The histograms can be

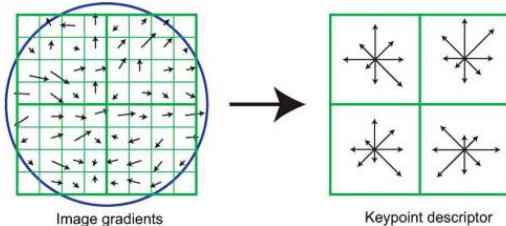


**Figure 8:** The keypoints of the instrument image 3. There are 3134 keypoints in this image with 3134 corresponding descriptors. Every green circle is a keypoint



**Figure 9:** The tip of the fork of image 8 is enlarged to show the keypoints. Every circle is a keypoint where the center indicates position and the radii the orientation. Note that there can be up to 2 radii in a keypoint, depending on whether they satisfy the 80% histogram requirement. Higher octaves and scales (i.e. blurrier images) produce less keypoints.

summarized in a  $4 \times 4 \times 8 = 128$  vector for each keypoint. This is the local image descriptor. This vector is then normalized to unit length, which keeps its value if all pixels' illumination increase by the same amount. This makes the local image descriptor illumination invariant, provided the illumination is constant over all pixels in the region around the keypoint. If that is not the case, gradient magnitudes are disproportionately affected with respect to gradient orientations. By capping the gradient magnitudes and renormalizing the vector, emphasis is placed on orientation.



**Figure 10:** This example image uses a  $8 \times 8$  set of samples from which a  $2 \times 2$  descriptor array is computed. Image taken from Lowe, 2004 [7].

See the following example in equation 3, where the vector is 5 dimensional. The vector in this example is normalized, capped at 0.5 and renormalized. The actual value of the cap is experimentally determined to be 0.2 (appendix A[7]).

$$\begin{bmatrix} 7 \\ 3 \\ 2 \\ 1 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 7/8 \\ 3/8 \\ 2/8 \\ 1/8 \\ 1/8 \end{bmatrix} \rightarrow \begin{bmatrix} 4/8 \\ 3/8 \\ 2/8 \\ 1/8 \\ 1/8 \end{bmatrix} \rightarrow \frac{1}{\sqrt{31}} \begin{bmatrix} 4/8 \\ 3/8 \\ 2/8 \\ 1/8 \\ 1/8 \end{bmatrix} \quad (3)$$

### 3.6 Keypoint matching

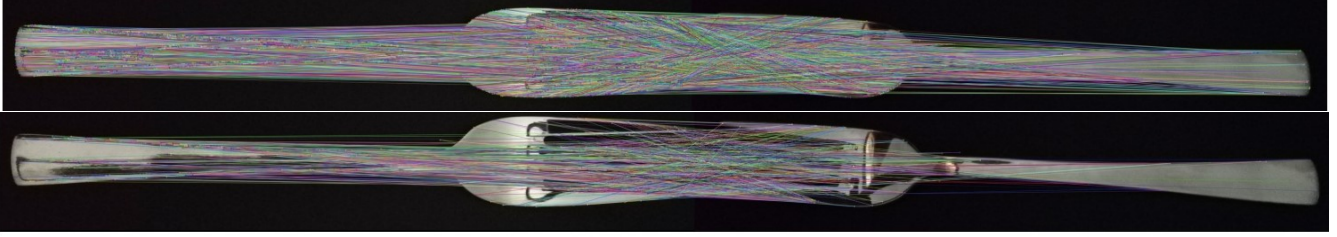
Each keypoint possesses a continuous location, a scale, an orientation and a local image descriptor.

The local image descriptor is used to match keypoints of one picture to another. A match of keypoints is a keypair. The minimum Euclidian distance of the descriptor vector is found by brute-force comparisons. These distances can be globally thresholded, but this method performs poorly because Euclidian distance isn't a great predictor of distinctiveness of the descriptor. For that, relative Euclidian distance between the closest neighbor and second-closest is used. If the closest match is much closer than the second-closest match (appendix A), then the descriptor is distinctive. We use a brute force approach to match the keypoints. Lowe [7] uses the Best-Bin-First algorithm, which is faster. We use the brute force method because it is more accurate. Figure 11 shows keypoint matching between images without and with the Euclidian distance requirement.

### 3.7 Keypair clustering

We are now working with keypairs. Every keypair has a pair of continuous position, a pair of orientations and a pair of local image descriptors. The local image descriptor is no longer needed and is discarded.

Keypairs indicate where the instrument should be placed in the scene. It can be likened to using thumbtacks to overlay images on a scrap board. The more thumbtacks are present, the surer one can be that the image is correctly placed. However, this does not mean



**Figure 11:** The keypoint matching algorithm applied to the instrument (upper left) 3 and scene (upper]right) 4 image. Every line matching these is a keypair. The bottom is the same but with the Euclidian distance requirement

all keypairs can be used in this manner because not all keypairs indicate true positive instrument detections within the scene. This is why the keypairs are grouped together based on how much they are alike.

A keypair has the following information in both the instrument and scene image: the position, orientation and the scale. Grouping keypairs together must be done based on invariant parameters. In this case the difference between the positions, orientations and scales. This gives a translation, rotation and scaling. These are computed in following section 3.8. They can then be grouped according to these differences.

### 3.8 Keypair parameters

A keypair has the  $x$  position, the  $y$  position, the orientation  $\theta$  and the the scale  $s$  for the instrument and the scene image. From these the following parameters can be computed: the difference in x coordinates  $d_x$ , the difference in y coordinates  $d_y$ , the scale ratio  $d_{Scale}$  and the difference in orientation  $d_\theta$ . The parameters  $d_{Scale}$  and  $d_\theta$  are straightforwardly computed in equations 4 and 5.

$$d_{Scale} = \frac{scale_{Scene}}{scale_{Instr}} \quad (4)$$

$$d_{Angle} = \theta_{Scene} - \theta_{Instr} \quad (5)$$

Parameters  $d_x$  and  $d_y$  are more complicated. The coordinates from the instrument image must be mapped to the scene image first. Instrument image coordinates must be adjusted for scale and rotation. The order is relevant, but all orders are possible. In this paper we do scaling first, rotation second. The scaling is relative to the center of the instrument image. See equations 6 and 7.

$$x_{scaledInstr} = (x_{Instr} - \frac{x_{length_{Instr}}}{2})d_{Scale} \quad (6)$$

$$y_{scaledInstr} = (y_{Instr} - \frac{y_{length_{Instr}}}{2})d_{Scale} \quad (7)$$

After scaling, the scaled instrument positions are rotated based on the difference in orientation  $d_\theta$  to find the x and y positions in the scene  $x_{InstrScene}$  and  $y_{InstrScene}$ . This is done with a rotation matrix shown in equation 8.

$$\begin{bmatrix} x_{InstrScene} \\ y_{InstrScene} \end{bmatrix} = \begin{bmatrix} \cos d_\theta & -\sin d_\theta \\ \sin d_\theta & \cos d_\theta \end{bmatrix} \begin{bmatrix} x_{scaledInstr} \\ y_{scaledInstr} \end{bmatrix} \quad (8)$$

Now the position of the instrument keypoint can be located in the scene. Finding parameters  $d_x$  and  $d_y$  is a simple subtraction. See equations 9 and 10.

$$d_x = x_{Scene} - x_{InstrScene} \quad (9)$$

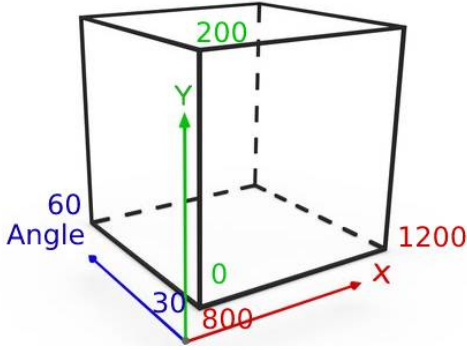
$$d_y = y_{Scene} - y_{InstrScene} \quad (10)$$

### 3.9 Hough Transform

Each keypair now has the four known parameters  $d_x$ ,  $d_y$ ,  $d_{Scale}$  and  $d_\theta$ . These describe how a keypoint in the instrument image is placed in the scene image. From these keypairs a model can be computed that maps all the pixels from the instrument image onto the scene image.

One keypair is sufficient to transform the instrument into the scene with the proper rotation. However, the second keypair determines the size of the image. The model that can be computed from two keypairs is called the similarity transform. The instrument or the scene image can also be taken at an angle, making it necessary to skew the instrument image. For this a third keypair is needed. The model computed from three keypairs is called the affine transform. Not just any three keypairs can be considered to find an instrument in the scene. The keypairs must be sufficiently similar. We use the Hough transform to cluster the keypairs.

The Hough transform considers the keypoints in parameter space, creating a 4 dimensional box known as a bin where the 4 parameters are enclosed by the 4 dimensions. Figure 12 shows the principle in 3 dimensions.



**Figure 12:** This example box is illustrated in 3 dimensions. Every keypair with an angle difference between 30-60 degrees, a y difference between 0-200 pixels and an x difference between 800-1200 pixels would fall in this box. For illustrative purposes the boundary problems are not shown, but can be viewed as one point falling into two adjacent boxes.

Each parameter is one dimension of the bin. The bin sizes per dimension are taken from Lowe [7]. They are 30 degrees for orientation difference, factor 2 for scale ratio and the maximum projected training image dimension for  $d_x$  and  $d_y$  is 0.25 (appendix A). As the entire parameter space must be accounted for, larger bins mean less bins and vice versa.

The voting is done by assigning a unique number to each box. Table 4 is an example of a keypair with four parameter values. For angles, 30 degrees lead to 12 bins, values 0-11 in python. The scale ratio of 2 include a number of scales, and can be assigned a number by taking the  $\log_2(\text{scale ratio})$ . For  $d_x$  and  $d_y$ , the ratio between  $d_x$  and the size of the scene image in the x dimension times the number of bins assigns a number. Rounding up and down yields whole numbers, two for every dimension for each keypair.

**Table 4:** Table of example parameter values for difference in orientation, scale ratio, difference in x and difference in y. Parameters  $d_x$  and  $d_y$  are dependent on image size, so normalized values are arbitrary. These are normalized and given a lowerbound and upperbound number.

	$d_\theta$	$d_{Scale}$	$d_x$	$d_y$
Parameter value	45	8.0	489	9312
Normalized value	1.5	3	8.6	4.5
Lowerbound number	1	3	8	4
Upperbound number	2	4	9	5

In this example, there are four sets of two values. These are added together like letters, to assign a number to a bin. The bin numbers for this example would be 1384, 1385, 1394, 1395, 1484, 1485, 1494, 1495, 2384, 2385, 2394, 2395, 2484, 2485, 2494 and 2495. Each of these bins gets one vote from this keypair. Other keypairs might vote for different bins, which will be generated. Multiple votes end up in the same bin only if their parameter values fall in the same bin. A bin needs at least two votes to compute a similarity model and

three votes for an affine model. These are explained in section ???. Every filled bin is called a cluster.

Boundary problems are prevented by adding the keypair match to the two closest bins for each parameter. For four parameters, this leads to 16 boxes with a vote for each keypair. This also increases total number of votes for bins by a factor 16, increasing the chance that a bin receives multiple votes.

### 3.10 Computing models

Both the similarity and affine transform are models that transform instrument image coordinates into scene image coordinates. They both have the same basic formula 11[7].

$$\begin{bmatrix} x_{Scene} \\ y_{Scene} \end{bmatrix} = \begin{bmatrix} m_1 & m_2 \\ m_3 & m_4 \end{bmatrix} \begin{bmatrix} x_{Instr} \\ y_{Instr} \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (11)$$

The similarity transform model takes into account scale, rotation and translation. The affine model takes into account scale, rotation, translation and skew. These are hidden in the  $m$  parameters from formula 11.

The  $m$  parameters from the similarity transform are found by multiplying the rotation with the scale. The scale is a scalar  $s$ . The rotation is defined by a standard rotation matrix  $R$  of rotation  $\theta$ . The  $m$  parameter matrix for the similarity transform is found in equation 12.

$$m = sR = s \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \quad (12)$$

The affine transform adds another element to the  $m$  parameters: the shear matrix  $C$  13.

$$m = s R C = s \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} 1 & c_x \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ c_y & 1 \end{bmatrix} \quad (13)$$

The  $m$  parameters can be found by taking formula 11 and using a least squares method. The similarity transform angle  $\theta$  and scale  $s$  can then be calculated from the  $m$  parameters in a straightforward manner as in equations 14 and 15.

$$\theta = \arctan \frac{m_2}{m_1} = \arctan \frac{-s \sin \theta}{s \cos \theta} \quad (14)$$

$$s = \sqrt{m_1^2 + m_2^2} = \sqrt{(s \cos \theta)^2 + (-s \sin \theta)^2} \quad (15)$$

The affine  $m$  parameters are more difficult to decompose into angle, scale and shear. We used Maple for the calculations. These can be found in appendix B.

The translation, scale and rotation of the model is then checked against the original keypairs. The model must be closer to the keypairs than half the original bin sizes of the Hough transform (appendix A). If this is not the case, the offending keypairs are removed from the cluster and the model is recomputed without those keypairs until the model is accepted. If the number of keypairs is insufficient, the entire cluster is discarded.

### 3.11 Probability of model

The candidate models can be rank ordered and thresholded by assigning a probability [6]. This is the probability that model  $m$  exists given the keypair set  $f$ . Which is the probability that the keypairs are present if the model is present divided by the probability that the keypairs are present in general. This probability can be approximated using equation 16 (appendix C).

$$P(m|f) \approx \frac{P(m)}{P(m)+P(f|-m)} \quad (16)$$

If this probability is larger than a certain threshold, the model is accepted. Lowe recommends a probability threshold of 0.98 (appendix A).

### 3.12 Duplicate removal by NMS

The method of applying the SIFT algorithm as described in the previous section 2.1 leads to more positive identifications. However, as each instrument photo is treated by the algorithm as a new possible instrument, the same instrument can be found in the same position multiple times as shown in figure 13 (left).

These duplicates in combination with the possibility of multiple detections per instrument image distorted the results of the algorithm. Every duplicate is a false positive and throws the count off. To remove the duplicates we used the Non-Maximum Suppression algorithm (NMS).

The NMS algorithm was used to identify duplicates and remove them [12]. It functions by ranking the proposed models by probability, taking the model with the highest probability of being a true positive. This model

was then compared with all other models by calculating the intersection of the bounding boxes and dividing by the union of the bounding boxes. This is the Intersection over Union (IoU). If this is larger than 0.9 (appendix A, that model is removed. This threshold is chosen to be high to ensure that only very similar models are removed. If the model is not removed, it is retained as a separate detection and matched to all remaining models.

The remaining models are the detected instruments. Figure 13 (right) shows a successful removal of duplicates leaving only one positive identification.

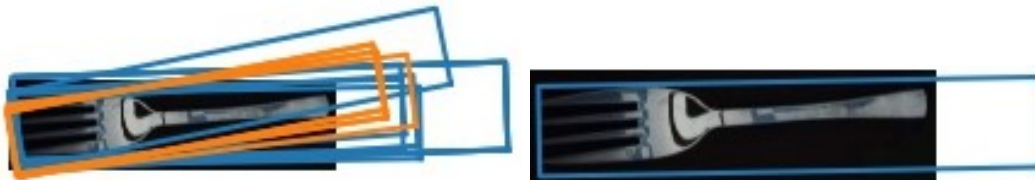


Figure 13: There are 14 detections before the NMS algorithm (left). Afterwards, only 1 detection remains (right).

## 4 Results

The results displaying the differences between glare and matte objects and the differences between the affine and similarity transform can be found in table 5. Table 6 shows the effect of increasing the cluster threshold for both glare and matte objects. The difference between regular and mirrored images for glare, matte and their average is found in table 7. Finally, table 8 shows the difference between the similarity and affine transform for both regular images and images taken by a tilted camera along the green axis of figure 2.

Additionally, as an example we applied the program on a scene with multiple instruments with an instrument map of 28 instrument images. This took 40 seconds. Figure 14 shows what a real life application of the program could look like. Of the 4 objects, only two were detected. These were detected with no false positives.

Tables 5, 6, 7 and 8 show a significant difference between the performance of the glare and matte objects in terms of both precision and recall. Tables 5 and 8 show

that affine transform performs worse than the similarity transform in terms of both precision and recall.



**Figure 14:** An example of a scene image with multiple instruments detected. There are 2 false negatives and 0 false positives. The instrument map here contains 28 images, 7 per instrument.

**Table 5:** Comprehensive results displaying the true positive, false positive, false negative, precision and recall of glare and matte objects for similarity (Sim) and affine (Aff) transform and their average (Avg).

	Glare Sim	Glare Aff	Mat Sim	Matte Aff	Glare Avg	Mat Avg
True positive	20%	14%	52%	39%	17%	46%
False positive	73%	67%	26%	26%	70%	26%
False negative	7%	20%	22%	34%	13%	28%
Precision	0.21	0.17	0.67	0.60	0.19	0.64
Recall	0.74	0.41	0.70	0.53	0.55	0.62

**Table 6:** Difference between glare and matte for cluster thresholds 3 and 4.

	Glare Cluster 3	Glare Cluster 4	Matte Cluster 3	Matte Cluster 4
Precision	0.14	0.23	0.61	0.67
Recall	0.51	0.56	0.67	0.59

**Table 7:** Difference of precision and recall between regular (Reg) and mirrored (Mir) images for the glare, matte and their average.

	Glare Reg	Glare Mir	Matte Reg	Matte Mir	Avg Reg	Avg Mir
Precision	0.24	0.18	0.81	0.74	0.57	0.49
Recall	0.66	0.50	0.57	0.47	0.59	0.47

**Table 8:** Difference between tilt and regular for affine and similarity, for glare and matte objects.

Tilt	Sim Glare	Aff Glare	Sim Mat	Aff Mat	Avg Sim	Avg Aff
Precision	0.13	0.16	0.67	0.58	0.34	0.32
Recall	0.87	0.54	0.79	0.61	0.81	0.59
Regular	Sim Glare	Aff Glare	Sim Mat	Aff Mat	Avg Sim	Avg Aff
Precision	0.20	0.16	0.71	0.70	0.41	0.39
Recall	0.49	0.25	0.72	0.60	0.59	0.40

## 5 Discussion

### 5.1 Matte and metal

In columns 5 and 6 of table 5 the difference between glare and matte objects is clear. Glare objects have a much higher false positive percentage and a much lower true positive percentage than matte objects. Tables 7 and 8 show the difference between in terms of glare and precision.

Table 6 concretely shows a lower precision and recall for glare objects across cluster thresholds 3 and 4. Most notably the precision for glare objects is much lower. This is because the number of false positives is higher in glare objects.

This difference is probably due to the unpredictability of light reflection in metals. In the first step of SIFT, scale space extrema detection, keypoints are generated by local minima and maxima. Metal surfaces have specular reflection, while matte surfaces have diffuse reflection. A specular reflection can have many extrema and therefor many keypoints, depending on what it reflects. Additionally, the variability of these keypoints between images can differ vastly due to the specular reflection of metal surfaces.

On the other hand, diffuse reflections give a surface a constant color. This uniform reflection has fewer local extrema, and therefor fewer keypoints than a metal surface. For matte objects, this means non-detections are more prevalent.

Variables such as image quality, lighting, background and object size and complexity being equal, metal objects generate more keypoints and these have a higher variability.

The SIFT program should then account for what type of instrument it is detecting, and adjust parameters accordingly. For glare objects, the number of accepted models should be reduced by increasing parameter thresholds.

### 5.2 Similarity and affine transform

Columns 1 through 4 of table 5 show the difference between the similarity and affine transform. The similarity transform performs better in both true positives and false negatives. False positives remain about equal, with some variance below 10%. This indicates that the similarity transform performs better than the affine transform across the board. The nature of the affine transform isn't congruent with these results: it should have fewer false negatives at least because it should be able to fit models onto the scene more easily.

This not being the case indicates that there is some problem with the affine transform. The affine transform uses a complicated equation (appendix B) to de-

compose the  $m$  parameters into scale, skew and rotation. However, for some affine model computations, these formula attempt to take the square root of a negative number. These potential models were skipped as candidate models, which makes it likely that an indeterminate number of true and false positives were lost. Possible solutions can be to find a different mathematical way of decomposing the  $m$  parameters or to use the similarity transform for these instances.

However, the affine transform is comparatively better when looking at skewed images. Comparing the similarity transform with the affine transform from table 8, the average precision for glare and matte stays quite even but recall drops significantly. For regular photos the recall drops 8% more than for tilted photos. When comparing skewed and regular images, the skewed images have an increase in false negatives for both the affine and similarity transform. However, the increase of the affine transform is smaller than the increase of the similarity transform. This indicates that the affine transform performs comparatively better than the similarity transform for skewed images.

So there is evidence that the affine transform can be more effective at comparing depth rotated objects. This is only relevant if the affine transform can be made to work properly, as it now performs worse in both precision and recall than the similarity transform.

### 5.3 Mirrored images

Table 7 shows that the regular images perform better. This is expected because SIFT is not a mirror-invariant program. This means that copying and mirroring images contributes to a higher positive detection rate, as objects that would have been mirrored physically can now be successfully detected.

The average difference found in columns 5 and 6 show a difference of 14% in precision and 20% in recall. However, copying and mirroring every instrument image also nearly doubles computing time.

### 5.4 Hyper parameters

The two hyper parameters that we varied were the cluster threshold and the probability threshold.

The first two columns of table 6 show that the precision and recall of glare objects increases with increased cluster threshold. The stricter threshold predictably reduces false positives and increases false negatives. This is paired with an increase in true positives, as it reduces false positives more than it reduces false negatives. For the final two columns, matte objects show the same effect on false positives and negatives, though false negatives increase more than false positives decrease, which



means there are less true positives. This indicates, all other things being equal, that the cluster threshold of 4 is more suited to glare objects while a threshold of 3 to matte objects.

For matte and glare objects, the best cluster thresholds seem to be 3 and 4 respectively. A cluster threshold of 3 means that there are 3 keypairs required to compute a model. Going from 3 to 4 means that every model that would have been accepted with 3 keypairs is now discarded. However, these are often models that are similar to the best model.

Changing the cluster threshold is more impactful than it appears when solely looking at precision and recall. When increasing the threshold, the number of accepted models decreases. However, many of these models are similar and would have been removed anyway due to the NMS algorithm. This means that the effect of the cluster threshold change on precision and recall is not as noticeable. However, it means that any effect on precision and recall is based solely on the models the NMS algorithm could not catch.

The cluster threshold of 2 in this configuration of SIFT generated too many false positives and a cluster threshold of 5 generated many false negatives. However, depending variables such as the size of the instrument map, higher cluster thresholds could be useful.

Varying the probability threshold from 0.98 to 0.99 proved too small of a step. We still believe that the probability threshold is a logical choice to tune the algorithm. It is near the end of the algorithm, which makes it a very intuitive variable to tune. Model probabilities often reached 0.999999, so further testing should be done keeping a such a threshold in mind compared to the 0.98 recommended by Lowe.

In the SIFT algorithm, there are many other variables that can be used to tune the program. An overview of these and what would happen to the program if you changed them can be found in appendix A. Many of these were experimentally determined by Lowe [7]. These can be tuned as well but we believe that the cluster and probability threshold parameters should be researched more thoroughly. The variables should be considered first are the cluster threshold, the probability threshold and the Hough transform bin size. These are straightforward in their effect: reduce the number of positive detections by increasing thresholds or reducing the bin size. A more fundamental approach would be to reduce the number of keypoints detected in any given image. However, it is hard to predict what discarding information at such an early stage would do.

## 5.5 Precision and Recall

The results are largely measured by precision and recall. These are measured over categories. The precision shows the ratio of true positives and total positives, that is what fraction of all detections is true. Recall tells us the ratio between the true positives and the sum of the true positives and false negatives, that is what fraction of all relevant detections is detected. In this subsection we explain how results should be interpreted when taking into account recall and precision as separate.

We want to test the program in such a way that that quantitative data is produced that sheds light on the functioning of the program. This was done through the matching of individual photos, which isn't the way the program would function in an OR setting. The program takes an entire map of photos of an instrument and tries to locate that instrument in the photo, for each photo.

Here we have tested individual images. We do this because testing the entire map on a picture can give qualitative results, but would generate little data that indicate how the program would react to different configurations. This means that the results must be interpreted whilst taking the individual nature of the test into consideration.

For the precision this has little bearing; the higher the precision, the better the program functions. However, recall is another matter. A large number of false negatives will decrease recall, but a that isn't necessarily a bad thing. Each false negative is just one undetected instrument from an instrument map with potentially hundreds of photo's, which means false negatives aren't a problem so long as one of the other images in the map finds a true positive, thereby detecting the instrument.

F score is meant to be a measure of accuracy and often goes hand in hand with precision and recall [10]. In this thesis we disregard F score as its relationship with accuracy is skewed due to the individual nature of the tests contrasted with the wholesale nature of the intended program.

A high recall means only that there are few false negatives. It must be evaluated in concert with the precision. Table 9 shows how the configuration of the program should be altered for different situations. If there's a high precision and a high recall, then the program is working as intended. If there is high precision and low recall, the program might be too severe in its boundaries. It could also be just fine, depending on whether the map of instruments can generate at least one true positive. For low precision and high recall, the program needs to be stricter to increase the precision, even though recall might fall. And for low precision and low recall, the program needs to be stricter to increase precision and perhaps more photos are to be added to

the map to increase the chance of finding a true positive.

Table 9 shows how the algorithm should be configured when

**Table 9:** Precision and recall for a map of images.

	Low precision	High precision
Low recall	↑ Threshold & images	↓ Threshold
High recall	↑ Threshold	Good

Using an entire map of photos allows us to prioritize precision over recall, which means making the program stricter is preferable to reduce false positives, as false negatives can be compensated by adding more photos of instruments to maps. This will however increase computational time.

## 5.6 Requirements and limitations

The main limitations are practical. The program cannot detect what it cannot see. Therefore, the instruments have to be placed on a table in such a way that every object is at least partly visible. Also, the program would work best when the objects are placed on a contrasting background. These requirements would probably require a human touch or at least a deviation on current procedure.

Pictures have to be taken. This can be done automatically with some sort of camera set up or by a medical practitioner. Additionally, a good user interface is essential. If the medical team can easily input what instrument the program missed and draw a box around the instrument, the database can be filled quickly. On the other hand, if some instrument is supposed to be missing, that too should be easily entered in the user interface.

## 5.7 Future research

This algorithm needs more testing. Individual image tests were promising, but more cluttered images need to be tested. The NMS algorithm has not been tested in a setting where multiple objects are stacked on top of each other. The threshold for removing duplicates should be strict enough to ensure that doesn't happen.

Lowe [6] also integrates multiple views of an object. He takes an initial view of the object and uses the Hough transform [6] to add subsequent views from different angles. This creates a larger set of keypoints with which to match an image. This increases the chance of matching the model to an image. This method uses the similarity transform instead of the affine because the latter "provides a poor approximation for rotation in depth of more complex 3D objects" [6]. We believe this approach is suited to this application. Instead of

comparing multiple models per training image for multiple images in our image map, the relevant keypoints can be combined into one model. Additional training image keypoints can then be compared to that current model. If it matches, those keypoints can be added to the model. If it doesn't match, a new model can be made. This increases robustness of the program [6]. More specific to our implementation, not just 3D depth rotation could be varied. If the lighting changes per training image, unique keypoints that are generated by arbitrary specular reflection can be compared to the model and be discarded. Additionally, keypoints from the model can be weighted according to how often they appear in training images.

Run time has not been the focus of this research. For figure 14, the run time was 40 seconds on my laptop. and had 28 instrument images across 4 objects. In an OR setting, more images should be used per face of the instrument, increasing run time. However, the computer on which the program is run could be much faster instrument images can be loaded before every operation. Additionally, in section 3.6 we choose to use a brute force approach to matching, instead of the Best-Bin-First algorithm. This would make the program faster. There are many factors which can influence run time, and future research should make evaluate the run time thoroughly in an operating room setting.

Cluttered image tests should be done with images from actual surgeries, where the instrument count is done with SIFT and compared to a hand count. As with any practical application, testing in the field is preferable. In this case it is necessary, as it concerns technology where errors might result in harm to patients.

## 6 Conclusion

Glare objects generate more keypoints and more false positives than matte objects. The SIFT algorithm should be run on different settings for both. The cluster threshold should be 3 for matte objects and 4 for glare objects in the current configuration.

Positive detections will increase when instrument images and their mirrors are used. It is unclear whether that is worth the increased run time.

The similarity transform performs better across all metrics, probably because the affine transform does not function as expected. The affine transform, when working correctly, is expected to be able to be better able to deal 3D depth object rotation than the similarity transform. If any improvements are made to this program, fixing the affine transform should be high priority.

The probability threshold was varied too slightly between 0.98 and 0.99 to be used as a tuning variable. It could still function as such, given more research.

Run time is a factor to be considered when implementing the program in an operating room setting.

Can the Scale Invariant Feature Transform algorithm reliably detect objects in an operating room setting? Yes, it can. Conceptually, it could detect instruments with 100% accuracy, provided all instruments are visible in the scene image.

## A Variables and constants

There are many variables and constants I've used in section 2. In this appendix, they are displayed in table 10. Many of these constants are tested or assumptions by Lowe. We do not test most of these, but they are choices that could be adjusted to improve the algorithm.

**Table 10:** Variables, constants and choices made in this application of SIFT.

	Parameter Value	Section
k	$\sqrt{2}$	3.2
$\sigma_0$	1.6	3.2
Extrema threshold	0.03	3.3
Edge ratio	12.1	3.3
Gaussian window $\sigma$	1.5 scale	3.5
Keypoint descriptor area	16x16 pixels	3.5
Histogram vector	4x4x8=128	3.5
Gradient magnitude cap	0.2	3.5
Euclidian distance ratio	0.8	3.6
Min # keypairs sim	2	3.7
Min # keypairs aff	3	3.7
Degrees per bin	10	3.9
Bin size for $d_x, d_y$	10	3.9
Bin size for scale	10	3.9
Bin size for orientation	10	3.9
Bin size for verification	Half bin size	3.9
Probability threshold	0.98	3.11
IoU threshold	0.9	3.12

Sigma multiplier k is chosen as is because the chosen value allows for a downsampling of the image with a factor 2, saving computational time. A greater  $\sigma_0$  improves repeatability of keypoint detection, but increases computational time. At the 1.6, the repeatability is close to optimal.

The Taylor fit allows a value to be placed on the extremum peak. If extrema are lower than the given value, they are removed because these are unstable extrema with low contrast.

To find and remove extrema on the edge, compute the sum of the straight derivatives over the diagonal derivative around the extremum. The ratio between these can be thresholded [7].

In 3.4, 10 degrees per bin was chosen. Fewer degrees per peak lead to more peaks, but lower. This increases precision but also computational time. Fewer degrees per bin would lead to more duplicates, which means more keypoints. On the other hand, more keypoints also mean more false positives.

The Gaussian weighting window  $\sigma$  is chosen to be 1.5 times scale of the keypoint. Larger windows could improve accuracy, as more information is taken into account. However, it could also decrease accuracy, as

at a certain pixel distance useful information carry over is zero.

If peaks are within an 80% threshold, the keypoint is duplicated and given this secondary orientation. This increases the number of keypoints. However, it also increases the number of false positives. At 80%, 15% are duplicated.

The Gaussian weighting window  $\sigma$  in 3.5 is half the descriptor window, which is 8 pixels in our case. Its purpose is to remove is to reduce boundary issues by adding a gradual decline of effectiveness of pixels farther away from the descriptor position. Increasing  $\sigma$  increases the effect of faraway pixels within the window, making the boundary issues more prominent. Decreasing  $\sigma$  also increases boundary effects, as the gradual decline is lessened.

The keypoint descriptor area is chosen to be 16x16 pixels. This can be bigger, but again, it would take into account pixels that wouldn't contribute positive information. It could be smaller, but pixels could be missed that do contribute positive information. The area allows for 4x4 descriptors.

The histograms have 8 pixels, leading to 4x4x8 = 128 element feature vectors. Larger feature vectors can decrease matching by making the histograms too sensitive to distortion. Smaller vectors perform worse as well in terms of matching, though computational time decreases. The descriptor is sensitive to affine change, falling below an 80% matching repeatability for a viewpoint angle difference of 40 degrees or more.

Lowering the cap decreases gradient emphasis. The goal is to make sure gradient contribution isn't overbearing, but also not so low that it's irrelevant. The gradient magnitude cap was determined experimentally. The Euclidian distance ratio of section 3.6 is set to 0.8. This eliminates 90% of false matches whilst discarding less than 5% of true matches.

In subsection 3.9 two keypoint pairs are needed for the similarity transform model and three for the affine transform model.. However, the minimum could be increased to a number greater than 3. This would lead to models being accepted, but also a lower false positive rate.

The bin sizes are chosen for the similarity transform, not the affine transform. This means they are quite broad to catch non-rigid transformations. Decreasing bin size would lead to less models, but a lower false positive rate.

In subsection 3.10, half the bin size of section 3.9 is chosen to test the matches against. This is to remove any outliers that aren't within the scope of the model. A stricter test would remove more potential models, decreasing the false positive rate but increasing the false negative rate. A looser test does the opposite.

In subsection 3.11, the main variable is the probability threshold. This is set at 0.98 by Lowe. Varying this threshold is an excellent way of testing the algorithm. If the threshold increases, and false positives are removed, then the threshold should be higher. The value at which correct identifications are removed is too high. However, if the threshold increases and correct identifications are removed before the false positives, then something else is wrong.

The duplicate removal threshold is set to 0.5. This application of the algorithm, with its many potential instrument identifications per instrument folder, needed a strict duplicate removal threshold. The disadvantage of this is that overlaying instruments have a higher chance of being removed. To counteract this, the instrument photos need to be cropped as close to the instrument as possible, as to decrease the area of the bounding boxes.

## B Affine

$$s = \sqrt{m_1 m_4 - m_3 m_2} \quad (17)$$

$$D = m_1 m_4 - m_3 m_2 \quad (18)$$

$$C = m_2^2 + m_4^2 \quad (19)$$

$$R = \text{sqrts}(s - C \text{ol}) \quad (20)$$

$$\text{Arc}_1 = \left( \frac{R(m_1 m_4 - m_2 m_3 m_4)}{C s D} + m_2^2 m_3 - m_1 m_2 m_4 \right) \quad (21)$$

$$\text{Arc}_2 = \left( \frac{R(m_1 m_2 m_4 - m_2^2 m_3)}{C s D} + m_1 m_4^2 - m_2 m_3 m_4 \right) \quad (22)$$

$$\theta = \arctan \frac{\text{Arc}_1}{\text{Arc}_2} \quad (23)$$

$$c_x = \frac{R}{D} \quad (24)$$

$$c_y = \frac{1}{C} \left( -m_1 m_4 \frac{R}{D} + m_2 m_3 \frac{R}{D} + m_1 m_2 + m_3 m_4 \right) \quad (25)$$

## C Probability

$P(m)$  is the chance that a single keypair is correct, which is the ratio of correctly matched image features to all matched features in a typical image. This is about

0.01. See appendix A. The probability of matched keypairs  $f$  given that  $m$  is not present can be computed with the binomial distribution 26.

$$P(f | -m) = \sum_{j=k}^n \binom{n}{j} p^j (1-p)^{n-j} \quad (26)$$

Here the probability  $p$  is given by the following equation 27.

$$p = d l r s \quad (27)$$

Where  $d$  is the probability of accidentally selection a database match to the current model, which is the matches of the current model divided by all the matches.  $l$  is the probability of accidentally satisfying the location constraints, which is the bin size in both dimensions multiplied, which is  $1/16$ .  $r$  the probability of accidentally satisfying the orientation restrains given our bin sizes, which is  $30/360 = 0.083$  and  $s$  the probability of accidentally satisfying the scale constraints, which is 0.5.

## D Code

```

1 import numpy as np
2 import scipy as sp
3 import scipy.special
4 import shapely.geometry
5 from shapely.ops import cascaded_union
6 from shapely.geometry import Polygon
7 import math
8 from math import pi
9 import matplotlib
10 matplotlib.rcParams['backend'] = 'tkagg'
11 import matplotlib.pyplot as plt
12 import matplotlib.gridspec as gs
13 import cv2
14 from tqdm import tqdm
15
16 import warnings
17 warnings.filterwarnings("error")
18
19 import sys, os
20 from datetime import datetime
21
22 autofill = True
23
24 ## Constants
25 binSizeAngle = 30
26 binSizeScaleFactor = 2
27 binSizeXFactor = 0.25
28 binSizeYFactor = 0.25
29 polygonThreshold = 0.25
30 transformChooser = False
31 # transformChooser = True
32 clusterThreshold = 2

```

```

33
34 probL = binSizeXFactor*binSizeYFactor
35 probR = binSizeAngle/360
36 probS = 1/binSizeScaleFactor
37 probM = .01
38 probThreshold = 0.98
39
40 showBoxPreNMS = True
41 showBoxPostNMS = True
42
43 '''
44 logBase: x is the input of the logarithmic
45          function, base the type of logarithm
46 '''
47 def logBase(x, base=2):
48     return np.log(x)/np.log(base)
49
50 ## Functions
51 '''
52 def betai(a, b, x):
53     '''
54     Calculate incomplete beta function  $I_x(a,b)$ ,
55     taken from paragraph 6.4 of Press, W.H. et
56     al. - Numerical Recipes in C (second
57     edition)
58     '''
59     if x < 0 or x > 1: raise ValueError("betai: x
60     must be 0 <= x <= 1")
61     if x in (0, 1): bt = 0
62     else: bt = np.exp(gammln(a+b)-gammln(a)-gammln
63     (b)+a*np.log(x)+b*np.log(1-x))
64     if x < (a+1)/(a+b+2): return bt*betacf(a, b, x
65     )/a
66     else: return 1-bt*betacf(b, a, 1-x)/b
67
68 def betacf(a, b, x):
69     '''
70     Calculate continued fraction for incomplete
71     beta function, taken from paragraph 6.4 of
72     Press, W.H. et al. - Numerical Recipes in
73     C (second edition)
74     '''
75     maxit = 100
76     eps = 3.e-7
77     fpmin = 1.e-30
78
79     qab = a+b
80     qap = a+1.
81     qam = a-1.
82     c = 1.
83     d = 1.-qab*x/qap
84     if np.abs(d) < fpmin: d = fpmin
85     d = 1/d
86     h = d
87
88     for m in range(1, maxit+1):
89         m2 = 2*m
90         aa = m*(b-m)*x/((qam+m2)*(a+m2))
91
92         d = 1+aa*d
93         if np.abs(d) < fpmin: d = fpmin
94         c = 1+aa/c
95         if np.abs(c) < fpmin: c = fpmin
96         d = 1/d
97         h *= d*c
98         aa = -(a+m)*(qab+m)*x/((a+m2)*(qap+m2))
99         d = 1+aa*d
100        if np.abs(d) < fpmin: d = fpmin
101        c = 1+aa/c
102        if np.abs(c) < fpmin: c = fpmin
103        d = 1/d
104        de = d*c
105        h *= de
106        if np.abs(de-1) < eps: break
107
108        if m > maxit: raise ValueError("betacf: a or b
109        too big, or maxit too small")
110        return h
111
112 def gammln(xx):
113     '''
114     Calculate natural logarithm of gamma function,
115     taken from paragraph 6.1 of Press, W.H.
116     et al. - Numerical Recipes in C (second
117     edition)
118     '''
119     cof = [76.18009172947146, -86.50532032941677,
120     24.01409824083091, -1.231739572450155,
121     0.1208650973866179e-2, -0.5395239384953e
122     -5]
123     y = x = xx
124     tmp = x+5.5
125     tmp -= (x+.5)*np.log(tmp)
126     ser = 1.00000000190015
127     for j in range(6):
128         y += 1
129         ser += cof[j]/y
130     return -tmp+np.log(2.5066282746310005*ser/x)
131
132 '''
133 instrumentCounter: keeps track of which instrument
134                    is being considered at any given iteration of
135                    i
136 photoCount tracks which photo is being looked at
137                    according to the directories, whatIsDetected
138                    gives the name according to the directories, i
139                    is instrument index
140 returns the identifier for the instrument
141 '''
142 def instrumentCounter(photoCount, whatIsDetected, i)
143     :
144     # photoCountSum = 0
145     photoCountSum = 0
146     for instrument in range(len(whatIsDetected)):
147         photoCountSum += int(photoCount[instrument
148         ])
149     if np.floor(i/2) <= photoCountSum:

```

```

126     instrumentIdentifier = whatIsDetected[
127         instrument]
128     return instrumentIdentifier
129 '''
130 keypointDifference: calculates the differences
131     between keypoints.
132     q is the for loop that iterates over the number of
133     matches. featMatches is a list of all the
134     matches between keypoints.keyPointsImg is a
135     list of all the keypoints of the image.
136     keyPointsImgsInstruments is the same, but for all
137     the keypoints of a given instrument of i. i is
138     the index going through all the instrument
139     images
140     Returns the difference in angle, in scale, in x
141     position, in y position and the indices of the
142     image and instrument keypoints.
143 '''
144 def keypointDifference(q,featMatches,keypointsImg,
145     keyPointsImgsInstruments,i):
146     indexImgKP = featMatches[q,0].queryIdx
147     indexInstrumentKP = featMatches[q,0].trainIdx
148     angleImgKP = keypointsImg[indexImgKP].angle*2*
149         pi/360
150     angleInstrumentKP = keypointsImgsInstruments[i
151         ][indexInstrumentKP].angle*2*pi/360
152     scaleImgKP = keypointsImg[indexImgKP].size
153     scaleInstrumentKP = keypointsImgsInstruments[i
154         ][indexInstrumentKP].size
155     xImgKP, yImgKP = keypointsImg[indexImgKP].pt
156     xInstrumentKP, yInstrumentKP =
157         keypointsImgsInstruments[i][
158             indexInstrumentKP].pt
159     dangle = angleInstrumentKP - angleImgKP
160     dangle %= 2*pi
161     dscale = scaleImgKP/scaleInstrumentKP
162     dx = (xInstrumentKP-imgsInstruments[i].shape
163         [1]/2) * dscale
164     dy = (yInstrumentKP-imgsInstruments[i].shape
165         [0]/2) * dscale
166     dx, dy = (np.array([[np.cos(-dangle), -np.sin
167         (-dangle)], [np.sin(-dangle), np.cos(-
168         dangle)]])) @ np.array([[dx], [dy]]))[: ,0]
169     dx, dy = xImgKP-dx, yImgKP-dy
170 '''
171 hashTableFunc stores all keypair contributions to
172     the hough transform in a hashTable.
173     imgGray is the grayed version of the image.
174     imgsGrayInstruments[i] is the same for every
175     instrument image i. binSizeXFactor is the size
176     of the bin defined in variables.
177     Same for binSizeAngle and binSizeScaleFactor.
178     Doesn't return, but updates the hashTable
179 '''
180 def hashTableFunc(imgGray,imgsGrayInstruments,
181     binSizeXFactor,binSizeYFactor,binSizeAngle,
182     binSizeScaleFactor,dx,dy,dscale,dangle,
183     indexImgKP,indexInstrumentKP,hashTable,i):
184     xSizeImg = imgGray.shape[1]
185     ySizeImg = imgGray.shape[0]
186     xSizeInstrument = imgsGrayInstruments[i].shape
187         [1]
188     ySizeInstrument = imgsGrayInstruments[i].shape
189         [0]
190     maxSizeInstrument = max(xSizeInstrument,
191         ySizeInstrument)
192     binSizeX, binSizeY = np.array([binSizeXFactor,
193         binSizeYFactor])*maxSizeInstrument*dscale
194     nrXBins = np.ceil(xSizeImg/binSizeX)
195     nrYBins = np.ceil(ySizeImg/binSizeY)
196     angleHash = (dangle*360/(2*pi))/binSizeAngle
197     scaleHash = logBase(dscale, base=
198         binSizeScaleFactor)
199     xHash = dx*nrXBins/xSizeImg
200     yHash = dy*nrYBins/ySizeImg
201     for iAngle in int(np.floor(angleHash))+np.
202         arange(2):
203         for iScale in int(np.floor(scaleHash))+np.
204             arange(2):
205             for iXHash in int(np.floor(xHash))+np.
206                 arange(2):
207                 for iYHash in int(np.floor(yHash))+
208                     np.arange(2):
209                     key = str(iAngle)+str(iScale)+
210                         str(iXHash)+str(iYHash)
211                     if key in hashTable.keys():
212                         hashTable[key] += [[
213                             indexImgKP,
214                             indexInstrumentKP,
215                             dangle, dscale, dx, dy,
216                             binSizeX, binSizeY]]
217                     else:
218                         hashTable[key] = [[
219                             indexImgKP,
220                             indexInstrumentKP,
221                             dangle, dscale, dx, dy,
222                             binSizeX, binSizeY]]
223     return
224 '''
225 hashTableClusterRemover removes clusters of
226     keypairs if they are less than 3 keypairs or
227     if the model they produce gives an error.
228     hashtableItems is a listed accessible version of
229     the hashTable, q iterates over every cluster

```

```

    in hashTable items, clusterThreshold is 222
    defined in variables and reset is used to 223
    reset the while loop this function is in.
195 returns the hashTableItems and a reset. Also
    updates the hashTable. 224
196 '''
197 def hashTableClusterRemover(hashTableItems,q,
    clusterThreshold,reset): 225
198     nInstrumentKeypoints = len(np.unique([
    hashTableItems[q][1][s][1] for s in range(
    len(hashTableItems[q][1]
199     if nInstrumentKeypoints < clusterThreshold or
    reset == 'remove': 228
200         keyToRemove = hashTableItems[q][0] 229
201         hashTable.pop(keyToRemove) 230
202         del hashTableItems[q] 231
203         reset = True
204     return hashTableItems, reset
205
206 ''' 232
207 similarityTransform gives a model by checking if
    clusters from the hashTable are within the
    bins as defined in variables. This transform
    takes into account rotation and scale. 233
208 All inputs are found in descriptions above.
209 returns modelMat, the model that transforms the
    instrument coordinates to the image
    coordinates. modelParametersMat gives the
    rotation and scale in one matrix. 235
210 modelParametersTrans gives the translation
    separately in one vector. reset, again gives
    the option to reset if one of the keypairs in
    the cluster didn't fall within certain
    parameter bounds 237
211 ''' 241
212 def similarityTransform(i,q,clusterThreshold,
    hashTableItems,keypointsImg, 242
    keypointsImgsInstruments,dangle,dscale,dx,dy,
    reset):
213     affineA = np.zeros((2*(len(hashTableItems[q]
    [1])), 4))
214     affineB = np.zeros((2*(len(hashTableItems[q]
    [1])), 1)) 243
215     keypointIndices = []
216     for r in range(len(hashTableItems[q][1])):
217         [indexImgKP, indexInstrumentKP] =
            hashTableItems[q][1][r][:2] 244
218         xKeypointImg = keypointsImg[indexImgKP].pt
            [0] 245
219         yKeypointImg = keypointsImg[indexImgKP].pt
            [1] 247
220         xKeypointInstrument =
            keypointsImgsInstruments[i][
            indexInstrumentKP].pt[0] -
            imgsInstruments[i].shape[1]/2
221         yKeypointInstrument =
            keypointsImgsInstruments[i][
            indexInstrumentKP].pt[1] -
            imgsInstruments[i].shape[0]/2
    affineA[2*r, :] = np.array([
    xKeypointInstrument, -
    yKeypointInstrument, 1, 0])
    affineA[2*r+1, :] = np.array([
    yKeypointInstrument,
    xKeypointInstrument, 0, 1])
    affineB[2*r] = xKeypointImg
    affineB[2*r+1] = yKeypointImg
    keypointIndices += [indexImgKP,
    indexInstrumentKP]
    modelParameters = np.linalg.lstsq(affineA,
    affineB, rcond=None)[0][:,0]
    modelParametersMat = np.array([[
    modelParameters[0], -modelParameters
    [1]], [modelParameters[1],
    modelParameters[0]]])
    modelParametersTrans = np.array([[
    modelParameters[2]], [modelParameters
    [3]]])
    affineAngle = float(2*pi-np.arctan2(
    modelParameters[1], modelParameters[0])
    )%(2*pi)
    affineScale = float(np.sqrt(modelParameters
    [0]**2 + modelParameters[1]**2))
    affineX = float(modelParameters[2])
    affineY = float(modelParameters[3])
    reset = False
    for t in range(len(hashTableItems[q][1])):
    dangle, dscale, dx, dy, binSizeX, binSizeY
    = hashTableItems[q][1][t][2:]
    binSizeScale = 2**(np.floor(np.log2(dscale)
    )+1)-2**(np.floor(np.log2(dscale)))
    if (np.abs(dangle-affineAngle) <= 0.5*
    binSizeAngle*2*pi/360) and (np.abs(
    dscale-affineScale) <= 0.5*binSizeScale
    ) and (np.abs(dx-affineX) <= 0.5*
    binSizeX) and (np.abs(dy-affineY) <=
    0.5*binSizeY):
    R = np.array([[np.cos(2*pi-affineAngle)
    , -np.sin(2*pi-affineAngle)], [np.
    sin(2*pi-affineAngle), np.cos(2*pi-
    affineAngle)]]])
    modelMat = affineScale*R
    pass
    else:
    keyToChange=hashTableItems[q][0]
    valuesToChange=hashTable.pop(
    keyToChange)
    del valuesToChange[t]
    if len(valuesToChange) != 0:
    hashTable[keyToChange]=
    valuesToChange
    reset = True
    return 0, 0, 0, 0, reset

```



```

254     return modelMat, modelParametersMat,      282
        modelParametersTrans, keypointIndices,  283
        reset                                  284
255     '''                                     285
256     affineTransform gives a model in the same way as 286
        the similarityTransform, only adding an affine 287
        element.                                     288
257     returns the same only modelParameters also 289
        incorporates the affine element in the matrix. 290
258     The reset can be set to 'remove', which removes an 291
        entire cluster instead of one keypair if the 292
        found model gives an error.                 293
259     '''                                     294
260 def affineTransform(i,q,clusterThreshold,      295
        hashTableItems,keypointsImg,           296
        keypointsImgsInstruments,dangle,dscale,dx,dy, 297
        reset):                                  298
261     affineA = np.zeros((2*(len(hashTableItems[q  299
        ] [1])), 6))                             300
262     affineB = np.zeros((2*(len(hashTableItems[q  301
        ] [1])), 1))                             302
263     keypointIndices = []
264
265     for r in range(len(hashTableItems[q] [1])): 303
266         [indexImgKP, indexInstrumentKP] =
            hashTableItems[q] [1] [r] [:2]
267         xKeypointImg = keypointsImg[indexImgKP].pt 304
            [0]
268         yKeypointImg = keypointsImg[indexImgKP].pt 305
            [1]
269         xKeypointInstrument =
            keypointsImgsInstruments[i] [
            indexInstrumentKP].pt[0] -
            imgsInstruments[i].shape[1]/2
270         yKeypointInstrument =
            keypointsImgsInstruments[i] [
            indexInstrumentKP].pt[1] -
            imgsInstruments[i].shape[0]/2
271
272         affineA[2*r,:] = np.array([
            xKeypointInstrument,
            yKeypointInstrument, 0, 0, 1, 0])
273         affineA[2*r+1,:] = np.array([0, 0,
            xKeypointInstrument,
            yKeypointInstrument, 0, 1])
274         affineB[2*r] = xKeypointImg
275         affineB[2*r+1] = yKeypointImg
276         keypointIndices += [indexImgKP,
            indexInstrumentKP]
277
278     modelParameters = np.linalg.lstsq(affineA, 316
        affineB, rcond=None)[0] [:,0]           317
279     modelParametersMat = np.array([[
        modelParameters[0], modelParameters[1]], [
        modelParameters[2], modelParameters[3]])] 318
280     modelParametersTrans = np.array([[
        modelParameters[4]], [modelParameters
        [5]]])                                     319
281
282     m = modelParameters
283     detM = m[0]*m[3]-m[2]*m[1]
284
285     if detM<=0:
286         reset = "remove"
287         return 0, 0, 0, 0, reset, q
288
289     col2lenM = m[1]**2+m[3]**2
290     rootMtemp = -detM*(detM-col2lenM)
291
292     if abs(rootMtemp) < 0.01:
293         rootMtemp = 0
294
295     if rootMtemp < 0:
296         reset = "remove"
297         return 0, 0, 0, 0, reset, q
298
299     rootM = np.sqrt(rootMtemp)
300
301     affineScale = np.sqrt(detM)
302     arctan2Part1 = 1/(col2lenM*affineScale)*(rootM
        /detM*(m[0]*m[3]**2-m[1]*m[2]*m[3])-m[0]*m
        [1]*m[3]+m[1]**2*m[2])
303     arctan2Part2 = 1/(col2lenM*affineScale)*(rootM
        /detM*(m[0]*m[1]*m[3]-m[1]**2*m[2])+m[0]*m
        [3]**2-m[1]*m[2]*m[3])
304     affineAngle = (2*pi-np.arctan2(arctan2Part1,
        arctan2Part2))%(2*pi)
305     affineAngleCheck = (2*pi-np.arctan2(
        arctan2Part1,-arctan2Part2))%(2*pi)
306     affineSkewX = rootM/detM
307     affineSkewY = 1/col2lenM*(-m[0]*m[3]*rootM/
        detM+m[1]*m[2]*rootM/detM+m[0]*m[1]+m[2]*m
        [3])
308     affineX = modelParameters[4]
309     affineY = modelParameters[5]
310
311     for t in range(len(hashTableItems[q] [1])):
312         dangle, dscale, dx, dy, binSizeX, binSizeY
            = hashTableItems[q] [1] [t] [2:]
313         binSizeScale = 2**((np.floor(np.log2(dscale)
            )+1)-2**((np.floor(np.log2(dscale))))
314
315         if (np.abs(dangle-affineAngle) <= 0.5*
            binSizeAngle*2*pi/360) and (np.abs(
            dscale-affineScale) <= 0.5*binSizeScale
            ) and (np.abs(dx-affineX) <= 0.5*
            binSizeX) and (np.abs(dy-affineY) <=
            0.5*binSizeY):
316             # Construct model matrix
317             Cx = np.array([[1, affineSkewX], [0,
                1]])
318             Cy = np.array([[1, 0], [affineSkewY,
                1]])
319             C = Cx@Cy
320             R = np.array([[np.cos(2*pi-affineAngle)
                , -np.sin(2*pi-affineAngle)], [np.
                sin(2*pi-affineAngle), np.cos(2*pi-
                affineAngle)]]])

```

```

321     modelMat = affineScale*R*OC
322
323     return modelMat, modelParametersMat,
        modelParametersTrans,
        keypointIndices, reset, q
324 else:
325     keyToChange=hashTableItems[q][0]
326     valuesToChange=hashTable.pop(
        keyToChange)
327     del valuesToChange[t]
328     if len(valuesToChange) != 0:
329         hashTable[keyToChange]=
            valuesToChange
330     return 0, modelParametersMat,
        modelParametersTrans,
        keypointIndices, reset, q
331
332 '''
333 probability checks whether a found model is likely
        to be a correct model, and removes those not
        able to make the threshold.
334 keypointsImgCoords is the location of any given
        keypoint in the image. probL,R,S are the
        probabilities that
335 returns the output that can be used to calculate
        bounding boxes
336 '''
337 def probability(modelMat,modelParametersMat,
        keypointsImgCoords,modelParametersTrans,
        imgsInstruments,i,probL,probR,probS,
        hashTableItems,probThreshold,
        instrumentIdentifier,output):
338     keypointsImgModelCoords = np.linalg.inv(
        modelParametersMat)@(keypointsImgCoords-
        modelParametersTrans)
339     keypointsImgModelCoordsInds = np.where((
        keypointsImgModelCoords[0,:] >= -
        imgsInstruments[i].shape[1]/2) & (
        keypointsImgModelCoords[0,:] <
        imgsInstruments[i].shape[1]/2) & (
        keypointsImgModelCoords[1,:] >= -
        imgsInstruments[i].shape[0]/2) & (
        keypointsImgModelCoords[1,:] <
        imgsInstruments[i].shape[0]/2))[0]
340     keypointsImgModelCoords =
        keypointsImgModelCoords[:,
        keypointsImgModelCoordsInds]
341     n = keypointsImgModelCoords.shape[1]
342
343     # 2) Calculate probability of accidentally
        matching a single image feature to the
        current model pose
344     probD = len(keypointsImgsInstruments[i])/np.
        sum([len(keypointsImgInstrument) for
        keypointsImgInstrument in
        keypointsImgsInstruments])
345     p = probD*probL*probR*probS
346
347     # 3) Calculate probability that feature
        matches were accidental, given that the
        model is not present (binomial probability
        CMF, equation 6.4.12 from Press, W.H. et
        al. - Numerical Recipes in C (second
        edition))
348     probMatchesNoModel = betai(len(hashTableItems[
        q][1]), n-len(hashTableItems[q][1])+1, p)
349
350     # 4) Calculate probability that the model is
        present, given the feature matches
351     probModelGivenMatches = probM/(probM +
        probMatchesNoModel)
352     # 5) Accept or reject model
353     if probModelGivenMatches >= probThreshold:
354         output[i].append([modelMat,
            modelParametersTrans, keypointIndices,
            instrumentIdentifier,
            probModelGivenMatches])
355     return
356
357 '''
358 boundingCalculator finds the bounding boxes in the
        scene image according to the instrument
        images. output is an collection of the
        relevant modelParameters,
359 the keypointIndices, the probability and the type
        of instrument.
360 returns a list of bounding boxes, and plots them
        too if desired.
361 '''
362 def boundingCalculator(imgsGrayInstruments,
        imgsInstruments,output):
363     colours = matplotlib.rcParams['axes.prop_cycle
        '].by_key()['color']
364     nInstrumentCols = int(max(1, np.floor(np.sqrt(
        len(imgsGrayInstruments))))))
365     nInstrumentRows = int(np.ceil(len(
        imgsGrayInstruments)/nInstrumentCols))
366     listPointsModel = []
367     for indImg in range(len(imgsGrayInstruments)):
368         y = int(np.floor(round(indImg/
            nInstrumentCols, 1)))
369         x = indImg%nInstrumentCols
370
371         # Show instrument
372         imgInstrument = imgsInstruments[indImg].
            copy()
373         imgInstrument = cv2.cvtColor(imgInstrument,
            cv2.COLOR_BGR2RGB)
374
375         ax = fig.add_subplot(nInstrumentRows,
            nInstrumentCols*2, nInstrumentCols*(y
            +1) + x+1 + y*nInstrumentCols)
376         ax.imshow(imgInstrument)
377         ax.axis('off')
378         rectPoints = np.array(np.meshgrid([0,
            imgInstrument.shape[1]], [0,

```

```

    imgInstrument.shape[0]))).T.reshape
    ((-1,2))[0,1,3,2],:]
379 rectPoints = np.append(rectPoints,
    rectPoints[0,None,:], axis=0)
380 ax.plot(rectPoints[:,0], rectPoints[:,1],
    colours[indImg%len(colours)])
381 ax.set_ylim([imgInstrument.shape[0]+1, -2])
382 ax.set_xlim([-2, imgInstrument.shape[1]+1])
383
384 # Show detections
385 currentOutput = output[indImg]
386 modelCounter = 0
387 for o in currentOutput:
388     modelParametersMat = o[0]
389     modelParametersTrans = o[1]
390     rectPointsModel = rectPoints - np.array
391     ([imgInstrument.shape[1]/2,
    imgInstrument.shape[0]/2])
392     rectPointsModel = modelParametersMat @
    rectPointsModel.T +
    modelParametersTrans
393     rectPointsModel = rectPointsModel.T
394     listPointsModel.append([rectPointsModel
    ,o[3],o[4],indImg,modelCounter])
395     modelCounter +=1
396     axImg.plot(rectPointsModel[:,0],
    rectPointsModel[:,1], colours[
    indImg%len(colours)]) # Plot
    bounding rectangle
397 return listPointsModel
398 '''
399 duplicateRemover finds duplicates according to a
    threshold and removes these with a Non-Maximum
    Suppression algorithm.
400 PolygonThreshold is the chosen threshold at which
    duplicates are removed.
401 returns a list of non-duplicate bounding boxes and
    a list of the removed duplicates. The latter
    one is used to remove duplicates from output
    as well.
402 '''
403 def duplicateRemover(listPointsModel,
    polygonThreshold):
404     polygonList = [] #make polygonlist including
    instrument and probability
405     for model in range(len(listPointsModel)):
406         polyCoordinates = ((listPointsModel[model
    ] [0] [0] [0],listPointsModel[model
    ] [0] [0] [1]),(listPointsModel[model
    ] [0] [1] [0],listPointsModel[model
    ] [0] [1] [1]),(listPointsModel[model
    ] [0] [2] [0],listPointsModel[model
    ] [0] [2] [1]),(listPointsModel[model
    ] [0] [3] [0],listPointsModel[model
    ] [0] [3] [1]))
407         currentPolygon = Polygon(polyCoordinates)
408         polygonList.append([currentPolygon,
    listPointsModel[model] [1],
    listPointsModel[model] [2],
    listPointsModel[model] [0],
    listPointsModel[model] [3],
    listPointsModel[model] [4]])
    ## NMS
    polygonProposal = []
    nextBestProposal = 0
    negativeProposal = []
    while len(polygonList) > 0:
        maxProbPolygon = 0
        if len(polygonList) == 1:
            bestPolygon = polygonList.pop(
                maxProbPolygon)
            polygonProposal.append(bestPolygon)
            break
        for polygons in range(len(polygonList)-1):
            if polygonList[maxProbPolygon] [2] >
                polygonList[polygons+1] [2]:
                maxProbPolygon = maxProbPolygon
            else:
                maxProbPolygon = polygons+1
        bestPolygon = polygonList.pop(
            maxProbPolygon)
        polygonProposal.append(bestPolygon)
        toBePopped = []
        for props in range(len(polygonList)):
            polygonU = [polygonProposal[
                nextBestProposal] [0], polygonList[
                props] [0]]
            union = cascaded_union(polygonU)
            intersection = polygonProposal[
                nextBestProposal] [0].intersection(
                polygonList[props] [0])
            IoU = intersection.area/union.area
            if IoU > polygonThreshold:
                toBePopped.append(props)
        toBePopped.reverse()
        for poppers in range(len(toBePopped)):
            thePopped = polygonList.pop(toBePopped[
                poppers])
            negativeProposal.append(thePopped)
        nextBestProposal += 1
    return polygonProposal, negativeProposal
    ## Ask path to image to load
    if autofill:
        print("Autofill_ enabled")
        pathImgSrc = "YourImagePath"
    else:
        pathImgSrc = 0
        while pathImgSrc == 0 or not os.path.isfile(
            pathImgSrc) or not pathImgSrc[-4:].lower()
            in (".jpg", ".png"):
            print("Paste_the_path_to_the_image_to_
            detect_instruments_in_below:" if

```

```

    pathImgSrc==0 else "The specified path
    doesn't point to an existing image.
    Please try again:")
453 pathImgSrc = input()
454 if pathImgSrc[0] == "&":
455     pathImgSrc = pathImgSrc[1:]
456     pathImgSrc = pathImgSrc.strip()
457 if (pathImgSrc[0] == "\" and pathImgSrc
    [-1] == "\" or (pathImgSrc[0] == \"'\"
    and pathImgSrc[-1] == \"'\"):
458     pathImgSrc = pathImgSrc[1:-1]
459     print()
460
461
462 ## Ask path to directory of directories with
    instruments to detect
463 if autofill:
464     pathInstrumentsSrc = "YourInstrumentPath"
465
466 else:
467     pathInstrumentsSrc = 0
468     while pathInstrumentsSrc == 0 or not os.path.
        isdir(pathInstrumentsSrc):
469         print("Paste the path to the directory with
            instrument ground-truth images below:")
470         if pathInstrumentsSrc==0 else "The
            specified path doesn't point to an
            existing directory. Please try again:")
471         pathInstrumentsSrc = input()
472         if pathInstrumentsSrc[0] == "&":
473             pathInstrumentsSrc = pathInstrumentsSrc
                [1:]
474             pathInstrumentsSrc = pathInstrumentsSrc
                .strip()
475             if (pathInstrumentsSrc[0] == "\" and
                pathInstrumentsSrc[-1] == "\" or (
                pathInstrumentsSrc[0] == \"'\" and
                pathInstrumentsSrc[-1] == \"'\"):
476                 pathInstrumentsSrc = pathInstrumentsSrc
                    [1:-1]
477                 print()
478                 if pathInstrumentsSrc[-1] != "\\":
479                     pathInstrumentsSrc += "\\
480 # Retrieve filenames for ground-truth images
481 pathsInstrumentsSrc = []
482 photoCount = []
483 whatIsDetected = []
484 pathDirectories = os.listdir(pathInstrumentsSrc)
485 for e in range(len(pathDirectories)):
486     while 1==1:
487         print()
488         print('Detect',pathDirectories[e], '?')
489         isInstrumentPresent = input("Enter [y] or [
            n].")
490         if isInstrumentPresent == "y":
491             print('Scanning for', pathDirectories[
                e])
                break
            elif isInstrumentPresent == "n":
                break
            else:
                print("Please enter [y] or [n].")
        if isInstrumentPresent == "y":
            pathInstrumentDirectory =
                pathInstrumentsSrc + '/' +
                pathDirectories[e]
            pathsInstrumentsSrc.extend([os.path.join(
                pathInstrumentDirectory, f) for f in os
                .listdir(pathInstrumentDirectory) if os
                .path.isfile(os.path.join(
                pathInstrumentDirectory, f)) and f
                [-4:].lower() in (".jpg", ".png")])
            photoCount.append(len(os.listdir(
                pathInstrumentDirectory)))
            whatIsDetected.append(pathDirectories[e])
492
493
494
495
496
497
498
499
500
501
502
503 ## Load image
504 print("Loading image from file..")
505 img = cv2.imread(pathImgSrc)
506 imgGray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
507
508 ## Load instrument images
509 print("Loading ground-truth images from directory
    ..")
510 imgsInstruments = [cv2.imread(f) for f in
    pathsInstrumentsSrc]
511 imgsInstrumentsMirror = [np.flip(imgInstrument,0)
    for imgInstrument in imgsInstruments]
512 imgsInstruments += imgsInstrumentsMirror
513 imgsGrayInstruments = [cv2.cvtColor(imgInstrument,
    cv2.COLOR_BGR2GRAY) for imgInstrument in
    imgsInstruments]
514
515 ## Apply SIFT
516 print("Initiating SIFT algorithm..")
517 sift = cv2.SIFT_create()
518
519 keypointsImg, descriptorsImg = sift.
    detectAndCompute(imgGray, None)
520
521 keypointsImgCoords = np.array([kp.pt for kp in
    keypointsImg]).T
522
523 print("Applying SIFT to ground-truth images..")
524 keypointsImgsInstruments = []
525 descriptorsImgsInstruments = []
526 for imgInstrument in tqdm(imgsGrayInstruments):
527     keypointsAndDescriptorsImgInstrument = sift.
        detectAndCompute(imgInstrument, None)
528     keypointsImgsInstruments += [
        keypointsAndDescriptorsImgInstrument[0],]
529     descriptorsImgsInstruments += [
        keypointsAndDescriptorsImgInstrument[1],]
530
531
532 ## Match features (detect instruments)

```

```

533 bf = cv2.BFMatcher(normType=cv2.NORM_L2,
                    crossCheck=False)
534
535 # Iterate over objects to detect
536 output = [[] for _ in range(len(imgsInstruments))]
537 allFeatMatches = [[] for _ in range(len(
538     imgsInstruments))]
539 instrumentTable={}
540 for i in tqdm(range(len(imgsInstruments))):
541     # Match features
542     featMatches = np.asarray(bf.knnMatch(
543         descriptorsImg, descriptorsImgsInstruments
544         [i], k=2))
545     imgsInstruments[i]
546     instrumentIdentifier = instrumentCounter(
547         photoCount,whatIsDetected,i)
548
549 # Ratio test
550 featMatches = featMatches[[featMatches[n,0].
551     distance <= .99*featMatches[n,1].distance
552     for n in range(featMatches.shape[0])]]
553 allFeatMatches[i] = featMatches
554 numberOfMatches = featMatches.shape[0]
555 hashTable={}
556
557 for q in range(featMatches.shape[0]):
558     dangle, dscale, dx, dy, indexImgKP,
559     indexInstrumentKP = keypointDifference(
560         q,featMatches,keypointsImg,
561         keypointsImgsInstruments,i)
562     hashTableFunc(imgGray,imgsGrayInstruments,
563         binSizeXFactor,binSizeYFactor,
564         binSizeAngle,binSizeScaleFactor,dx,dy,
565         dscale,dangle,indexImgKP,
566         indexInstrumentKP,hashTable,i)
567
568 # Geometric verification through Least Squares
569     to obtain affine projection
570 hashTableItems = list(hashTable.items())
571 q = 0
572 while q < len(hashTableItems):
573     reset = False
574     hashTableItems,reset =
575         hashTableClusterRemover(hashTableItems,
576         q,clusterThreshold,reset)
577
578     if reset == True:
579         continue
580
581     if transformChooser == False:
582         modelMat, modelParametersMat,
583         modelParametersTrans,
584         keypointIndices, reset =
585         similarityTransform(i,q,
586         clusterThreshold,hashTableItems,
587         keypointsImg,
588         keypointsImgsInstruments,dangle,
589         dscale,dx,dy,reset)
590
591     if reset == True:
592         continue
593
594 elif transformChooser == True:
595     modelMat, modelParametersMat,
596     modelParametersTrans,
597     keypointIndices, reset, q =
598     affineTransform(i,q,
599     clusterThreshold,hashTableItems,
600     keypointsImg,
601     keypointsImgsInstruments,dangle,
602     dscale,dx,dy,reset)
603
604     if reset == 'remove':
605         hashTableItems,reset =
606         hashTableClusterRemover(
607         hashTableItems,q,
608         clusterThreshold,reset)
609
610     if reset == True:
611         continue
612
613     probability(modelMat,modelParametersMat,
614         keypointsImgCoords,modelParametersTrans
615         ,imgsInstruments,i,probl,probR,probS,
616         hashTableItems,probThreshold,
617         instrumentIdentifier,output)
618     q+=1
619
620 ## Show original image and ground-truth images
621     with features drawn
622 ## Show original image with boxes around detected
623     instruments
624
625 if showBoxPreNMS == True:
626     fig = plt.figure()
627     axImg = fig.add_subplot(1, 2, 1)
628     axImg.imshow(img)
629     axImg.axis('off')
630
631 listPointsModel = boundingCalculator(
632     imgsGrayInstruments,imgsInstruments,output)
633 polygonProposal, negativeProposal =
634     duplicateRemover(listPointsModel,
635     polygonThreshold)
636
637 # Count instruments
638 instrumentCount = 0
639 for iden in range(len(whatIsDetected)):
640     for detection in range(len(output)):
641         instrumentCount +=len(output[detection])
642
643 print("Showing",instrumentCount, "detected_
644     instruments")
645
646 #Remove duplicates from output
647 negativeProposal.sort(key=lambda x: (x[4],x[5]),
648     reverse = True)
649
650 for q in range(len(negativeProposal)):
651     del output[negativeProposal[q][4]][
652         negativeProposal[q][5]]
653
654 # Count instruments
655 instrumentCount = 0
656 for iden in range(len(whatIsDetected)):
657     for detection in range(len(output)):

```

```

605     instrumentCount +=len(output[detection])
606
607 print("Showing",instrumentCount, "detected_
        instruments")
608
609 # Show bounding boxes after duplicate remover.
610 if showBoxPostNMS == True:
611     fig = plt.figure()
612     axImg = fig.add_subplot(1, 2, 1)
613     axImg.imshow(img)
614     axImg.axis('off')
615 listPointsModel = boundingCalculator(
        imgsGrayInstruments,imgsInstruments,output)

```

## References

- [1] Bani-Hani, K.E., Gharaibeh, K.A., Yaghan, R.J.: Retained surgical sponges (gossypiboma). *Asian Journal of Surgery* **28(2):109-115** (April 2005). [https://doi.org/10.1016/s1015-9584\(09\)60273-6](https://doi.org/10.1016/s1015-9584(09)60273-6)
- [2] Ebrahim Karami, S.P., Shehata, M.: Image matching using sift, surf, brief and orb: Performance comparison for distorted images. *S2015 Newfoundland Electrical and Computer Engineering Conference* (2015)
- [3] Egorova, N.N., Moskowitz, A., Gelijns, A., Weinberg, A., Curty, J., Rabin-Fastman, B., Kaplan, H., Cooper, M., Fowler, D., Emond, J.C., Greco, G.: Managing the prevention of retained surgical instruments: what is the value of counting? *Annals of Surgery* **247(1):13-8** (January 2008). <https://doi.org/10.1097/SLA.0b013e3180f633be>
- [4] Gawande, A.A., Studdert, D.M., Orav, E.J., Brennan, T.A., Zinner, M.J.: Risk factors for retained instruments and sponges after surgery. *New England Journal of Medicine* **348(3):229-35** (January 2003). <https://doi.org/10.1056/NEJMsa021721>
- [5] Lowe, G, D.: Object recognition from local scale-invariant features. *Proceedings of the Seventh IEEE International Conference on Computer Vision* **pp. 1150-1157 vol.2** (1999). <https://doi.org/10.1109/ICCV.1999.790410>
- [6] Lowe, G, D.: Local feature view clustering for 3d object recognition. *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition* **Vol. 1. IEEE, 2001** (2001)
- [7] Lowe, G, D.: Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision* **60(2), 91-110** (January 2004)
- [8] Lowe, G, D.: Invariant features from interest point groups. *BMVC* **Vol. 4. 2002** (September 2002)
- [9] Organization, W.H., Safety, W.P.: *Implementation manual surgical safety checklist* (first edition) (May 2008)
- [10] Powers, D.M.W.: Evaluation: From precision, recall and f-factor to roc, informedness, markedness correlation. *School of Informatics and Engineering Technical Report SIE-07-001* (2007)
- [11] Rivera, N., Mountain, R., Assumpcao, L., Williams, A.A., Cooper, A., Lewis, D.L., Benson, R.C., Miragliotta, J.A., Marohn, M., Taylor, R.H.: Assist - automated system for surgical instrument and sponge tracking. *IEEE International Conference on RFID* **pp. 297-30** (2008). <https://doi.org/10.1109/RFID.2008.4519358>
- [12] Sambasivarao, K.: Non maximum suppression, a technique to filter the detections of object detectors. *Towards Data Science*, 2019 (October 2019)