



Research Article

Volume 3 Issue 1 - May 2018
DOI: 10.19080/RAEJ.2018.03.555601

Robot Autom Eng J

Copyright © All rights are reserved by Lucas Agudiez Roitman

Real-Time Visual Subject Tracking and Classification by Combining Motion Signal Analysis and Tridimensional - Shape Feature Classifiers with Group-Induction Boosting Algorithms



Lucas Agudiez Roitman*

Department of Computer Science, Stanford University, USA

Submission: February 25, 2018; Published: May 14, 2018

*Corresponding author: Lucas Agudiez Roitman, Stanford Artificial Intelligence Laboratory, Department of Computer Science, Stanford University, Stanford, California, USA, Email: roitman@cs.stanford.edu

Abstract

This paper provides a novel and unprecedented approach for integrating motion features in the detection and classification of moving subjects in a static environment. More specifically, we measure the impact of the use of trajectory history, rotation history, blob orientation, motion frequency in the three axes, motion acceleration, segmentation errors, and flickering scores, and how they can influence classification of moving people, pets, and other objects. We apply our method to data captured by a combined color and depth camera sensor. We find that, while some motion descriptors slightly improve accuracy, the use of them in conjunction outperforms previous approaches in the classification and tracking of real-world moving subjects in real-time.

Keywords: Real-time tracking; Moving subjects; Classification; Motion signal; Motion statistics; Accelerometer; Orientation; Rgb camera; Depth images; Computer vision; Machine learning; Classifiers; Boosting; Artificial intelligence

Introduction

Many home security products that are available on the market promise to detect intruders at home and notify users via text messages. However, these home surveillance platforms often have high rates of false-positives and low tolerance for them. In other words, the user often receives messages because their pet walked in front of the camera or the wind moved the curtains. The user then grows accustomed to these false alarms and therefore ignores any future alarms that could be real threats. Furthermore, when a user wants to play back and watch all indoor moving subject activities, he has to watch all of the false-positive parts of the footage as well, wasting countless hours of time. Other potential uses of home cameras are harmed by the fact that the recognition technology is solely based on naive movement detection.

Our approach uses depth cameras, as well as accelerometers and gyroscopes to easily place the camera on the wall and detect its orientation, create a point cloud or tridimensional representation of the moving subjects, and use statistics and machine learning to more accurately detect and predict the nature of the moving object (Figure 1).



Figure 1:

We also use a group induction¹ method (inferring object type based on close similarity to a labelled object, or geographical proximity to it), which allows us to use smaller amounts of human labelling than similar conventional

¹Because our point clouds are grouped into animations, including multiple frames, all those frames can be labelled as part of the same object, assuming that the object was appropriately tracked. This can help with training the predictors with more data.

Source: Teichman, Alex, and Sebastian Thrun. "Group induction." 2013 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE, 2013.

approaches. Moreover, our results are compatible with semi-supervised learning techniques (meaning we can allow a user to label only a few examples, and we infer the rest from that)². This method enables the user to train the algorithm in a more practical manner, since a customer will often not take the time to label all data, but will agree to label a few of the data points (for example, the false-positives or false-negatives) [1].

Yoctopus accelerometer+gyroscope sensor. We added this accelerometer and gyroscope to the depth camera sensor, in order to programmatically determine its orientation with respect to the horizontal ground plane (Figure 2).

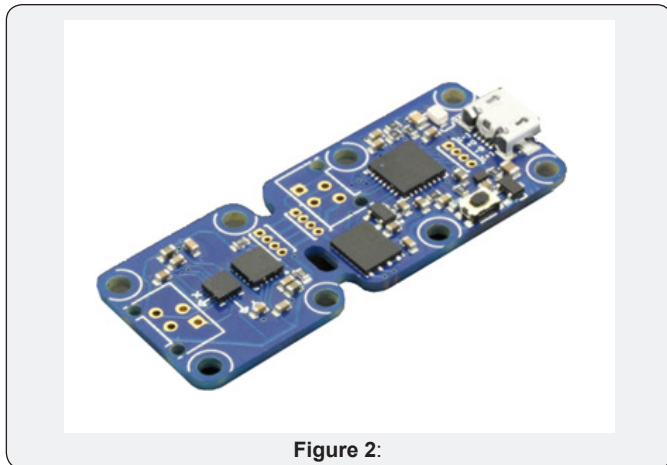


Figure 2:

This is what the depth camera and the IMU sensor look together³. The project also involved a heat sensor (green) but the results of the heat descriptors are reported in a separate paper [2] (Figure 3).

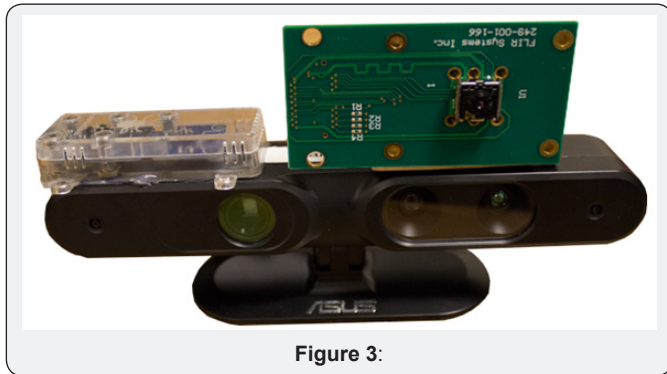


Figure 3:

Background

We used a model built originally for a project dealing with self-driving cars. The model involves a boosted learning classifier based on Adaboost, and the innovative use of group-induction as a method of semi-supervised learning. Although some tests and development were performed with semi-supervised learning, for the most part of the project, and for our research results, we decided to use a fully-supervised model so as to eliminate any possible noise, mislabelling or bias introduced by the use of semi-supervised learning⁴. (Figure 4).

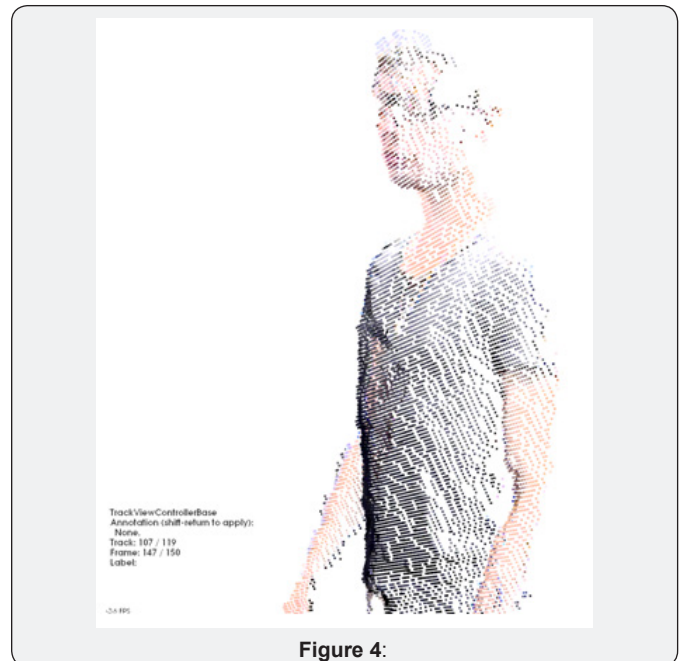


Figure 4:

The data we were working with is originated from a PrimeSense RGBD depth camera (color+depth), to which an accelerometer and gyroscope have been added. The depth image captured from the sensor is used to create a statistical model of the background or environment, which is static. When a part of the image is statistically outside the range of what constitutes the static background, that section is likely to be marked as an object in the foreground. That object is then converted into a point-cloud to more pragmatically represent the data. It is processed throughout our pipeline or saved for later processing⁵.

²Although our research results are compatible with this semi-supervised boosting, we have preferred to use fully supervised learning to generate the results for this paper.

Source: Teichman, Alex, and Sebastian Thrun. "Tracking-based semi-supervised learning." *The International Journal of Robotics Research* 31.7 (2012): 804-818.

³We used an inertial measurement unit, which performs sensor fusion of the data gathered by an accelerometer and gyroscope.

Source: Morrison, Melvin M. "Inertial measurement unit." U.S. Patent No. 4,711,125. 8 Dec. 1987.

⁴We are using the algorithm developed by Alex Teichman in order to classify tracks of all moving objects, instead of tracking a specific class. This method is non-specific to object class.

Source: Teichman, Alex, and Sebastian Thrun. "Practical object recognition in autonomous driving and beyond." *Advanced Robotics and its Social Impacts (ARSO), 2011 IEEE Workshop on. IEEE, 2011.*

⁵We use the depth-image segmentation method created by Alex Teichman and Jake Lussier.

Source: Teichman, Alex, Jake T. Lussier, and Sebastian Thrun. "Learning to Segment and Track in RGBD." *IEEE Transactions on Automation Science and Engineering* 10.4 (2013): 841-852

Machine Learning Algorithm

The classification technique we have decided to use for this experiment was that of the boosting technique⁶. “Boosting refers to a general and provably effective method of producing a very accurate prediction rule by combining rough and moderately inaccurate rules of thumb” in the following manner [3,4]:

$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

Where “in the simplest case, the range of each h_t ” is binary:

$$\epsilon_t = \Pr_{i \sim D_t}[h_t(x) \neq y_i]$$

For “binary alpha t”, we normally set:

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)^7$$

This algorithm is fast and has many tweakable parameter. It is also suitable for the processing of a large list of descriptors. The technique uses an array of weak learners that complement each other to improve overall performance.

Descriptor Pipeline

Our descriptor pipeline takes the point cloud animations, which consist of a series of frames that represent moving objects. These frames have point cloud stills in them. Together, all point clouds for all frames in the animation constitute *one instance*.

Each frame of the instance (each point cloud) is pushed through the descriptor pipeline individually, even if these individual images correspond to the same instance of the object being observed [5].

As seen on Figure 5, the data travels from Blob Entry Point, the input node, then goes into Blob Projector, which projects the RGBD data from a pixel matrix containing depth and color values into a 3D point cloud, a list of points (X, Y, Z pairs) with color values (R, G, B). Then, the data is sent into two different pods (nodes): HSV Histogram, where the RGB colors are converted into the HSV space⁸, and then a color histogram is computed. The color histogram has H, S, V values for each bin in the histogram [6].

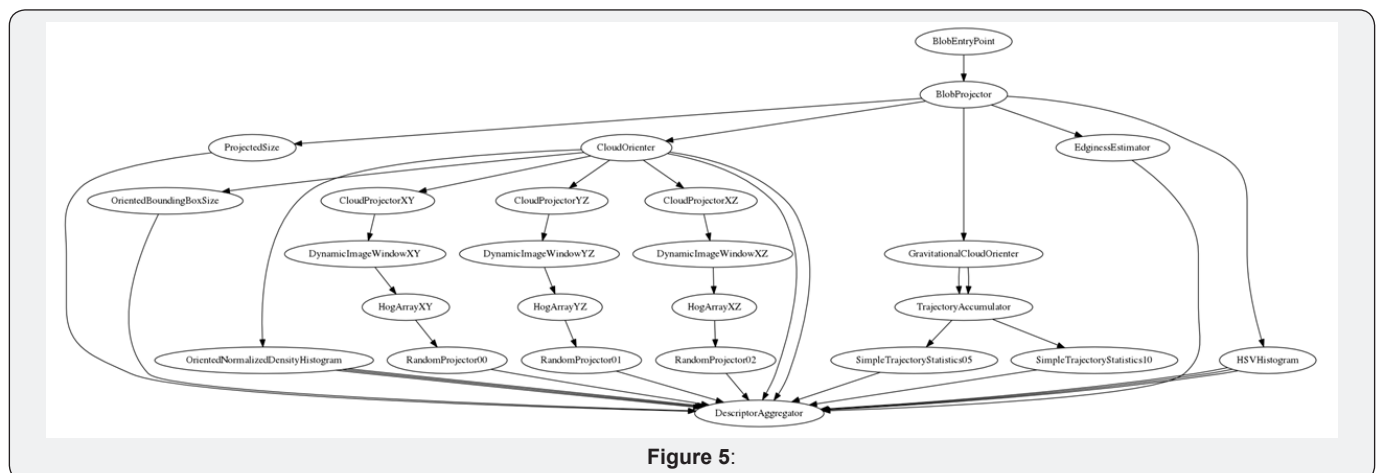


Figure 5:

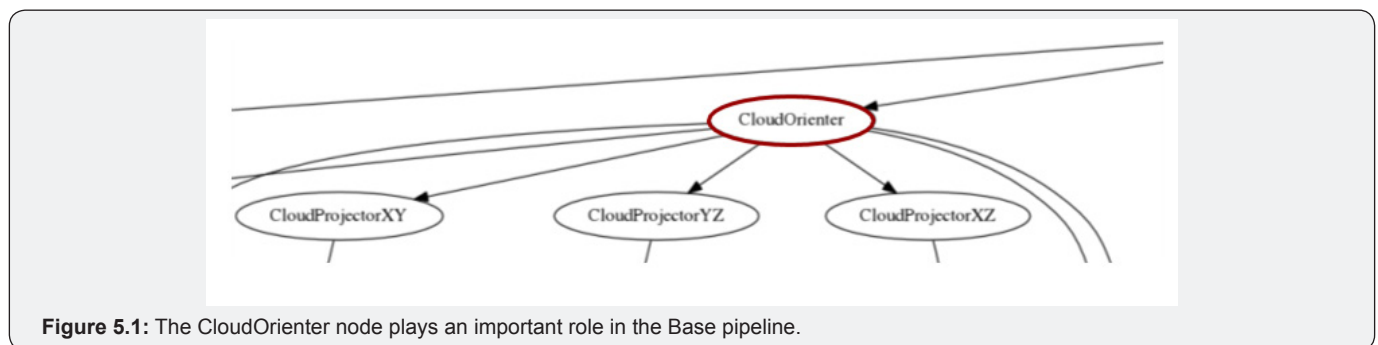


Figure 5.1: The CloudOrienter node plays an important role in the Base pipeline.

⁶We use the boosting technique for our machine learning purposes.

Source: Schapire, Robert E. “The boosting approach to machine learning: An overview.” *Nonlinear estimation and classification*. Springer New York, 2003. 149-171.

⁷From this paper:

Source: Schapire, Robert E. “The boosting approach to machine learning: An overview.” *Nonlinear estimation and classification*. Springer New York, 2003. 149-171.

⁸Data in the <red, green, blue> vector space are converted into <hue, saturation and value> space.

Source: Smith, Alvy Ray. “Color gamut transform pairs.” *ACM Siggraph Computer Graphics* 12.3 (1978): 12-19.

All these individual scalar values are then sent as a list to the descriptor aggregator, which appends them to a list of values for other descriptors.

Projected Size, a pod in which the 3D size is measured and sent to the Descriptor Aggregator pod (at the bottom) in order to append it to the list of descriptors to be computed by the machine learning algorithm⁹. This existing pipeline allows us to add more descriptors and experiment with how they affect classification results.

Other pods/nodes are found as well. We can see another node is CloudOrienter, which transforms the point-cloud and orients it based on its longest axis using the PCA algorithm.¹⁰

Once oriented into its principal components (X is the longest axis, Y the second longest, and Z the last one), the point cloud can be better aligned with previous images of that object. For example, if the camera is observing a pen, it will be aligned on its longest axis, so that it can be compared with other images of a pen no matter their perceived orientation [7].

Then, the oriented pointcloud is projected onto multiple 2D images (different planes: XY, YZ, and XZ) so that we can run a HOG algorithm on each of them.^{11,12}

Then, the results are condensed into a lower-density vector so that the values can be aggregated and sent to the machine learning algorithm [8] (Figure 6).

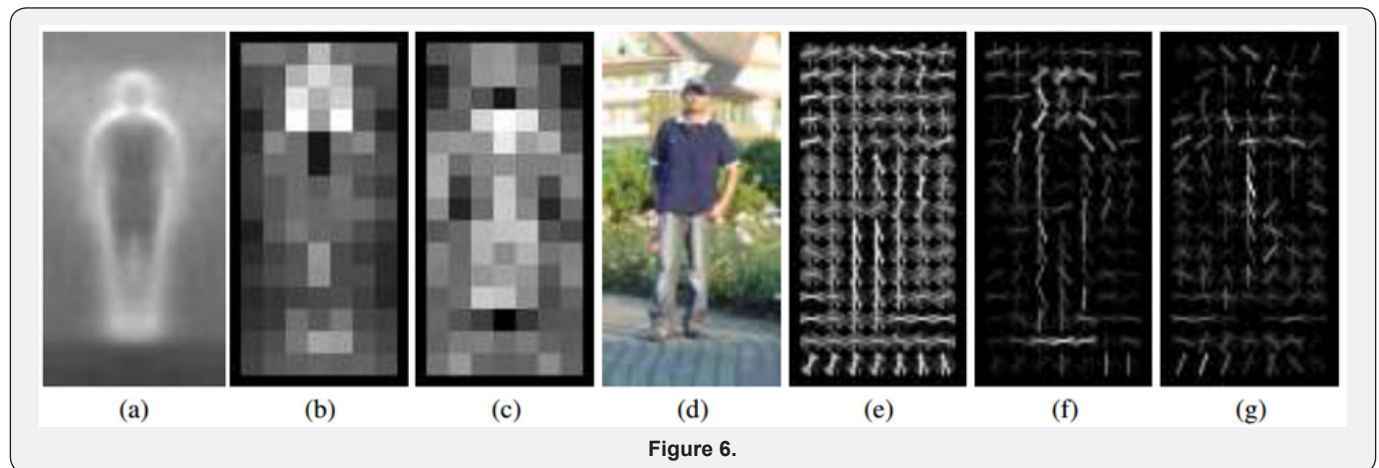


Figure 6.

Method

We placed a Prime Sense RGBD camera near the roof of a room and recorded dozens of hours of footage of activity in the room. We attached an accelerometer and a gyroscope sensor to the camera so that we could later re-orient the point clouds based on the gravity vector measured by the sensors. People, pets and other objects moved throughout the room and were recorded and extracted from the background. These blobs were then separately saved into an SD card. The data was then manually labeled (per instance) using a C++ tool. The data was then passed through the descriptor pipeline [9].

Then, we continued in the process to perform k-fold cross-validation (this process that enables us to reduce overfitting, so that our classifier can generalize better - learn to recognize

objects that are not that similar to the training examples). The training data was divided into K chunks and the classifiers were trained on K-1 chunks. Then, the classification and predictions for the remaining chunk were compared against the ground truth that was manually assigned to each instance, and the accuracy values, among others, were computed and saved.

We divided moving objects into multiple classes: cat, person, door, bush, and background, and then we performed 5-way classification. As we can see in the results, there were generally mild improvements as we kept computing new descriptors.

Pipelines Studied

As seen in Figure 6, the Base was the previously existing descriptor pipeline, which orients the point clouds for objects

⁹We use the boosting technique for our machine learning purposes.

Source: Schapire, Robert E. "The boosting approach to machine learning: An overview." *Nonlinear estimation and classification*. Springer New York, 2003. 149-171.

¹⁰Principal component analysis is an algorithm that transforms data points that seem to be correlated into linearly uncorrelated sets of values (principal components). Thus, the variables that are most correlated for the main axis, with other orthogonal axes accommodating the following most correlated sets of variables.

Source: Wold, Svante, Kim Esbensen, and Paul Geladi. "Principal component analysis." *Chemometrics and intelligent laboratory systems 2.1-3* (1987): 37-52.

¹¹Histogram of Oriented Gradients is an algorithm that is often used as a feature descriptor for computer vision tasks in order to detect objects. It was used with a lot of success to detect humans in 2D images.

Source: Dalal, Navneet, and Bill Triggs. "Histograms of oriented gradients for human detection." *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*. Vol. 1. IEEE, 2005.

¹²HOG algorithm steps.

based on their major components and aligns every object on such axes. The Gravity pipeline is the same as the Base descriptor pipeline, but we started aligning objects vertically, based on the camera orientation (measured by the accelerometer and gyroscope). As we can see, performance increases after aligning objects based on the gravity-vector measurement. We can see the Gravity pipeline in Figure 7. Then, Oriented Trajectory is a pipeline to which we also added a trajectory pod, which computes the object velocity and acceleration in X, Y and Z coordinates, aligned with the gravity vector as Z-down. This pipeline can be seen in Figure 8. Finally, since the X and Y orientation is based on the object shape rather than its environment, its use is irrelevant to motion statistics, and thus it cannot be reliably used. Therefore, we added the Plane Trajectory descriptors, which add a vertical speed feature, a horizontal speed feature (on the XY plane), and it also computes and sends the vertical acceleration and horizontal acceleration scalars to the machine learning algorithm. This pipeline can be found in Figure 9. Then, we added mean angular velocity and acceleration, in Figure 10 (Rotation Statistics). Later, we included a change rate estimator that measures the change of the pointcloud between different frames of the moving object

(Figure 11). And finally, we added a Fourier transform node to filter the spectrum of these motion statistics (Figure 12).

Base pipeline (Figure 5)

The base pipeline provides a functioning machine learning architecture that very accurately predicts the class of the moving object. However, in our studies, we will modify and add to this pipeline. As we can see here, the point cloud is sent through the CloudOrienter for most subsequent operations. This will be different in the following modified pipelines. The CloudOrienter rotates a point cloud in order to align its longest component with the X axis, and the following longest orthogonal axis on Y, leaving the remaining orthogonal axis to Z.¹³ (Figure 5.1)

Gravity pipeline (Figure 7)

The Gravity pipeline is an enhanced version of the Base pipeline, and includes a Gravitational CloudOrienter as the main node. Instead of orienting point clouds based on the object's dimensions, this node orients them based on their real-world orientation with respect to the environment's vertical axis. In other words, it uses the measured gravity vector (via accelerometer+gyroscope) to orient the clouds vertically [10].



Figure 7: Gravity Pipeline

The object is oriented vertically based on the IMU sensor that we attached to the depth camera. The Z axis (blue) aligns vertically, with the measured gravity vector.

It is also rotated so that its principal component aligns with the X axis (red) (Figure 7.1 & 7.2).

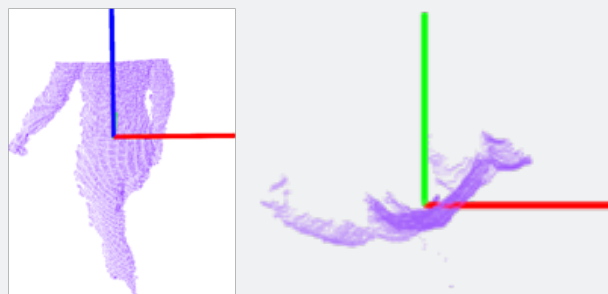


Figure 7.1:

¹³Principal component analysis is an algorithm that transforms data points that seem to be correlated into linearly uncorrelated sets of values (principal components). Thus, the variables that are most correlated for the main axis, with other orthogonal axes accommodating the following most correlated sets of variables.

Source: Wold, Svante, Kim Esbensen, and Paul Geladi. "Principal component analysis." *Chemometrics and intelligent laboratory systems 2.1-3* (1987): 37-52.



Figure 7.2: The Gravitational CloudOrienter node plays an important role this time.

Oriented trajectory pipeline (Figure 8)

In this pipeline, besides the Simple Trajectory Statistics nodes, which accumulated average speed (in any direction),

there is a new node called Oriented Trajectory Statistics, which sends separate X, Y, and Z values for velocity and acceleration in those axes (Figure 8.1).



Figure 8: Oriented Trajectory.



Figure 8.1: The OrientedTrajectory pipeline adds a new statistics node, which is useful for classifying objects with differential vertical and horizontal motion behaviours.

Plane trajectory pipeline (Figure 9)

The trajectory pipeline modifies the Oriented Trajectory Statistics node and adds a couple different computed values.

Besides simply separating into X, Y and Z, the node now computes a horizontal speed (in the 2D horizontal plane), a vertical speed (different from vertical velocity), and a horizontal and vertical acceleration, as well (Figure 9.1).



Figure 9: Plane Trajectory.

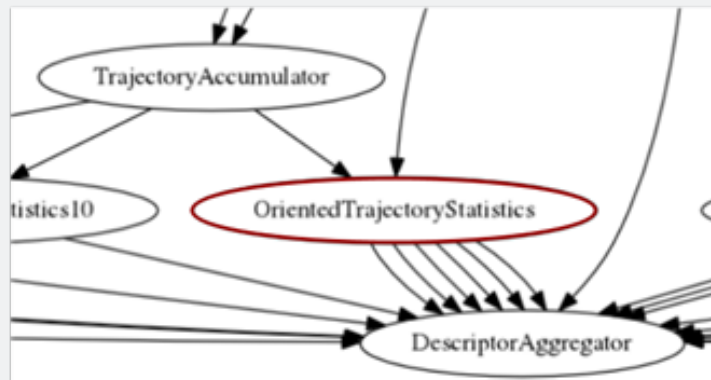


Figure 9.1: The Oriented Trajectory Statistics node is modified and sends more values to the Descriptor Aggregator.

Rotation pipeline (Figure 10)

The rotation pipeline also aggregates data about the point cloud's rotation in each frame, and the object's rotation acceleration as well, using the PCA algorithm to compute

orientations in individual frames¹⁴. This helps the detection technique differentiate between objects that rotate a lot in the 2D horizontal plane at different speeds and acceleration. This happens in the Rotation Statistics node (Figure 10.1).



Figure 10: Rotation.

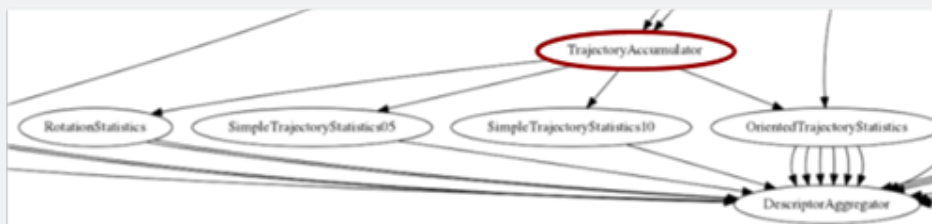


Figure 10.1: The rotation statistics node was added to compute the mean angular velocity and mean angular acceleration values of the moving object.

Change pipeline (Figure 11)

The change pipeline adds a descriptor that computes a score representing how much a point cloud changes among different timestamps or frames in the animation. The object might change in size or position greatly among multiple frames

due to segmentation errors, overlapping objects or simply because the object is moving or changing at high speeds. This descriptor was made to help us detect background noise instances, or objects that are very badly segmented, and easily exclude them from our other classes [11] (Figure 11.1).

¹⁴Principal component analysis is an algorithm that transforms data points that seem to be correlated into linearly uncorrelated sets of values (principal components). Thus, the variables that are most correlated for the main axis, with other orthogonal axes accommodating the following most correlated sets of variables.

Source: Wold, Svante, Kim Esbensen, and Paul Geladi. "Principal component analysis." *Chemometrics and intelligent laboratory systems 2.1-3* (1987): 37-52.



Figure 11: Change.



Figure 11.1: The change estimator was added as an independent descriptor.

Fourier pipeline (Figure 12)

The fourier pipeline adds a node that performs the fourier transform, changing the basis of the computed descriptor

data (such as horizontal speed, vertical speed, velocities, accelerations, angular acceleration, change scores, and rotation statistics.¹⁵ (Figure 12.1).

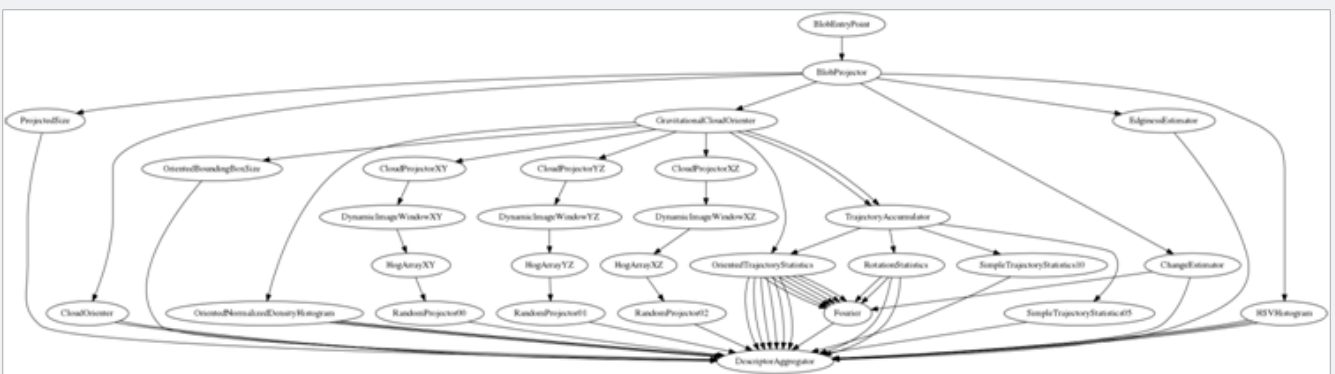


Figure 12: Fourier.

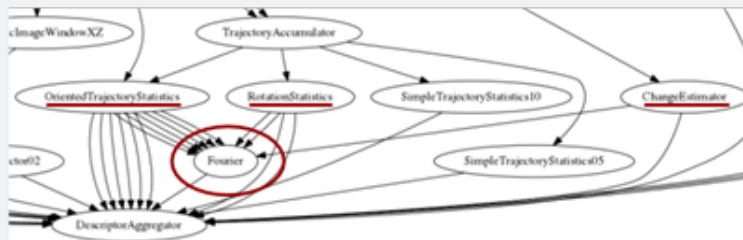


Figure 12.1: The fourier node receives data from other nodes and computes their fourier transform, in order to create descriptors that better adjust to objects that move in different frequencies.

¹⁵We used the fast fourier transform in order to convert the motion signal in time into a spectral representation of each motion frequency and their amplitudes.

Source: Weisstein, Eric W. "Fast fourier transform." (2015).

Results

We first trained our classifier on a dataset recorded indoors, where the only observed objects that moved were cats and people. Some false-positives involved a moving

curtain, background noise, or segmentation errors. Here, we can see that by using the measured gravity vector and aligning detected objects relative to gravity can improve performance (Table 1).

Table 1:

	Base	Gravity
Total test examples (number of objects):	1986	1986
Total accuracy (2-way classification):	0.93001	0.941591
Mean logistic score (max is 0, higher is better):	-0.18789	-0.146692
Mean exponential loss (min is 0, lower is better):	0.478102	0.323899
Test examples:	995	995
Average response:	3.0194	4.14704
True positives:	922	946
True negatives:	925	924
False positives:	66	67
False negatives:	73	49
Accuracy:	0.93001	0.941591
Aggregate precision (tp/(tp+fp)):	0.933198	0.93386
Aggregate recall (tp/(tp+fn)):	0.926633	0.950754
Per-class precision (tp/(tp+fp)):	0.933198	0.93386
Per-class recall (tp/(tp+fn)):	0.926633	0.950754

Following this experiment, the new features and enhanced pipeline were applied to a multi-class problem, with data recorded outdoors, this time (Table 2).

Table 2:

	Base	Gravity	Oriented Trajectory	Plane Trajectory
Total test examples (number of objects):	6468	6468	6468	6468
Total accuracy (5-way classification):	0.91064	0.91141	0.911565	0.912647
Cat accuracy	0.992424	0.992579	0.993352	0.993197
Person accuracy	0.965523	0.964904	0.96475	0.964286
Door accuracy	0.982684	0.982684	0.982993	0.982993
Bush accuracy	0.963358	0.964904	0.963513	0.965213

Following the experiments with the 5 classes, another set of experiments were performed, after removing the “bush” class, since most instances of bushes were marked as background and viceversa. It was hard for the person labelling to discern between a moving bush and random noise in the static background. As we can see, this different labelling scheme increased accuracy on the same dataset. Furthermore, we can tell that improvements were seen when adding more descriptors.

For this set of experiments, we also added a Rotation pod, which computes the object’s angular velocity and angular acceleration. This pipeline can be found in Figure 10. Finally, we added the “ChangeEstimator” pod, which computes a flickering score, a descriptor that evaluates how much the object’s shape changes between multiple frames. If it flickers a lot, it is likely a segmentation error or background noise, and therefore will not count as a moving subject. As we can tell, this descriptor significantly improves classification accuracy. Its pipeline can be seen in Figure 11 (Table 3).

Table 3:

	Oriented Trajectory	Plane Trajectory	Rotation	Change
Total test examples (number of objects):	6468	6468	6468	6468
Total accuracy (4-way classification):	0.94666	0.94682	0.94697	0.94913
Cat accuracy	0.99227	0.992579	0.992733	0.993352
Person accuracy	0.965368	0.96444	0.965213	0.965832
Door accuracy	0.98222	0.982993	0.981911	0.983302

Conclusion

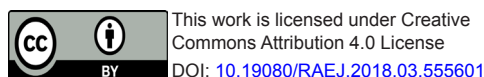
Although some additional descriptors offered mild improvements, in combination, they have delivered consistent improvement in learning accuracy. Additionally, descriptors such as the change-rate descriptor prove very useful to clean up segmentation errors. While aligning objects vertically proves helpful, this should probably be done additionally and not in place of aligning objects on the basis of their principal components (PCA). Future studies should measure the impact of having descriptors using both approaches, combined into the same boosting algorithm's input. The model explained in this paper can be useful for security agencies as well as other applications. It can be used to better identify intruders, differentiate them from other moving objects or subjects, and to create a more robust subject tracking system.

Acknowledgement

We thank Yash Savani for many useful comments, and Alex Teichman for providing the basis software and technology, without which this research would not have been possible. Finally, we thank Sebastian Thrun for providing the lab, funding and advice that has allowed us to perform this study.

References

- Teichman Alex, Jesse Levinson, Sebastian Thrun (2011) Towards 3D object recognition via classification of arbitrary object tracks. Robotics and Automation (ICRA), 2011 IEEE International Conference on. IEEE.
- Smith, Alvy Ray (1978) Color gamut transform pairs. ACM Siggraph Computer Graphics 12(3): 12-19.
- Teichman Alex, Sebastian Thrun (2012) Tracking-based semi-supervised learning. The International Journal of Robotics Research 31(7): 804-818.
- Schapire, Robert E (2003) The boosting approach to machine learning: An overview. Nonlinear estimation and classification. Springer New York, pp. 149-171.
- Teichman Alex, Jake T Lussier, Sebastian Thrun (2013) Learning to Segment and Track in RGBD. IEEE Transactions on Automation Science and Engineering 10(4): 841-852.
- Wold, Svante, Kim Esbensen, Paul Geladi (1987) Principal component analysis. Chemometrics and intelligent laboratory systems 2(1-3): 37-52.
- Teichman Alex, Sebastian Thrun (2011) Practical object recognition in autonomous driving and beyond. Advanced Robotics and its Social Impacts (ARSO), 2011 IEEE Workshop on. IEEE.
- Dalal N, Bill Triggs (2005) Histograms of oriented gradients for human detection. 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05). Vol. 1. IEEE.
- Teichman Alex, Sebastian Thrun (2013) Group induction. 2013 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE.
- Weisstein, Eric W (2015) Fast fourier transform.
- Morrison, Melvin M (1987) Inertial measurement unit. US Patent No. 4,711,125.



This work is licensed under Creative Commons Attribution 4.0 License
DOI: [10.19080/RAEJ.2018.03.555601](https://doi.org/10.19080/RAEJ.2018.03.555601)

Your next submission with Juniper Publishers will reach you the below assets

- Quality Editorial service
- Swift Peer Review
- Reprints availability
- E-prints Service
- Manuscript Podcast for convenient understanding
- Global attainment for your research
- Manuscript accessibility in different formats (Pdf, E-pub, Full Text, Audio)
- Unceasing customer service

Track the below URL for one-step submission

<https://juniperpublishers.com/online-submission.php>