

3D OBJECT TRACKING SENSOR FUSION PROJECT #3

Jasmine L. Taketa-Tran

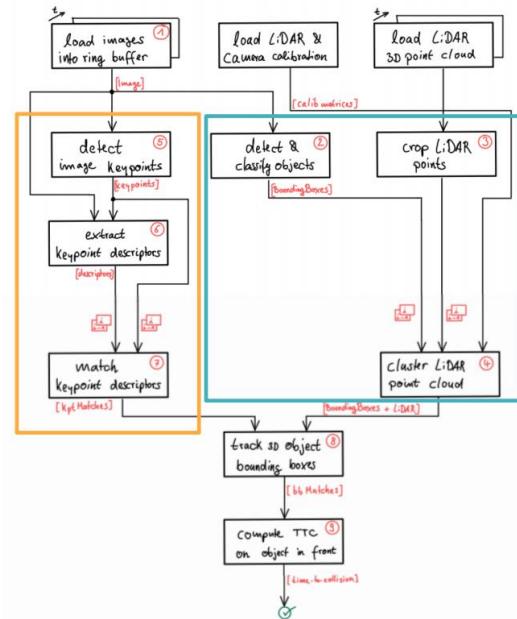
OVERVIEW

The objective of this project is to process 2-D camera and 3-D lidar data to estimate the Time-to-Collision from the ego car to the preceding vehicle in a sequence of KITTI video data frames.

TTC Building Blocks

Course Structure

- **Lesson 3** : Keypoint detection and matching
- **Mid-Term Project** : Develop the matching framework and test several state-of-the-art algorithms.
- **Lesson 4** : Lidar point processing and deep learning for object detection.
- **Final Project** : Track 3D bounding boxes and compute refined TTC



This writeup walks through each rubric and includes a statement, supporting figures, and references to specific places in the code where each step is handled.

FP.1 MATCH BOUNDING BOXES

- i** Implement the method “matchBoundingBoxes”, which takes as input both the previous and current data frames and provides as output the IDs of the matched regions of interest (i.e., the boxID property). Matches must be the ones with the highest number of keypoint correspondences.

In this part of the project, we associated pairs of YOLO-identified bounding boxes across adjacent frames when they contained the highest number of matched keypoints relative to the other bounding box candidates on opposing frames.



Figure 1: YOLO Bounding Boxes on Frames 9 & 10

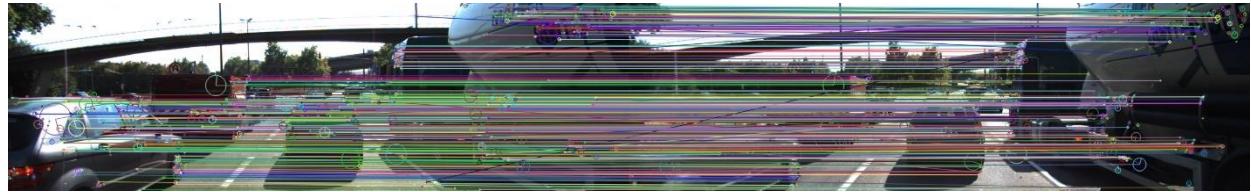


Figure 2: Keypoint Matches Between All Bounding Boxes on Frames 9 & 10 (Overlaps Filtered Out)



Figure 3: Keypoint Matches Between Bounding Boxes on the Preceding Vehicle

The implementation for matching and pairing bounding boxes across frames is shown below. Note that the matching keypoints between the images are stored in the data structure cv::DMatch such that queryIdx is the index of a keypoint in the previous frame, and trainIdx is the index of the corresponding keypoint match in the current frame.

Function Definition for Matching Bounding Boxes in camFusion_student.cpp:

```

void matchBoundingBoxes(std::vector<cv::DMatch> &matches, std::map<int, int> &bbBestMatches, DataFrame &prevFrame, DataFrame &currFrame)
{
    // count the number of keypoint matches that share the same bounding box between the current and previous frames
    cv::Mat count = cv::Mat::zeros(prevFrame.boundingBoxes.size(), currFrame.boundingBoxes.size(), CV_32S);
    for (auto matchpair : matches) // for each matched pair ...
    {
        // find the corresponding keypoint in the current and previous frame
        // note: kptMatches are ordered as follows: kptsPrev uses queryIdx and kptsCurr uses trainIdx
        cv::KeyPoint prevkp1 = prevFrame.keypoints.at(matchpair.queryIdx);
        cv::KeyPoint currkp1 = currFrame.keypoints.at(matchpair.trainIdx);

        for (int prevbb = 0; prevbb < prevFrame.boundingBoxes.size(); prevbb++)
        { // for each bounding box on the previous frame
            if (prevFrame.boundingBoxes[prevbb].roi.contains(prevkp1.pt))
            { // if the bounding box on the previous frame contains this keypoint
                for (int currbb = 0; currbb < currFrame.boundingBoxes.size(); currbb++)
                { // then loop through all bounding boxes on the current frame
                    if (currFrame.boundingBoxes[currbb].roi.contains(currkp1.pt))
                    { // if bounding box on the current frame also contains this keypoint
                        // then increment multimap counter for matches that share the same bounding box
                        count.at<int>(prevbb, currbb) = count.at<int>(prevbb, currbb) + 1;
                    }
                }
            }
        }
    } // end loop over all matched pairs between frames

    // associate bounding boxes that contain the highest number of matched pairs
    for (int prevBBIdx = 0; prevBBIdx < prevFrame.boundingBoxes.size(); prevBBIdx++)
    { // for each bounding box on the previous frame
        int boxID = -1, maxNumMatches = 0;
        for (int currBBIdx = 0; currBBIdx < currFrame.boundingBoxes.size(); currBBIdx++)
        { // for each bounding box on the current frame
            if (count.at<int>(prevBBIdx, currBBIdx) > maxNumMatches)
            { // find which pair of bounding boxes on prev and current frames have the maximum count
                boxID = currBBIdx;
                maxNumMatches = count.at<int>(prevBBIdx, currBBIdx);
            }
        }
        // save the box IDs of all matched pairs in bbBestMatches
        bbBestMatches[prevBBIdx] = boxID;
    }
}

```

Function Call in main():

```

/* +-----+ */
/* | TRACK 3D OBJECT BOUNDING BOXES | */
/* +-----+ */
// associate bounding boxes between current and previous frame using keypoint matches
map<int, int> bbBestMatches;
matchBoundingBoxes(matches, bbBestMatches, *(dataBuffer.end() - 2), *(dataBuffer.end() - 1));
// store best bounding box matches in current frame to the data buffer
(dataBuffer.end() - 1)>bbMatches = bbBestMatches;

```

FP.2 COMPUTE LIDAR-BASED TTC

- i** Compute the time-to-collision for all matched 3D objects using only Lidar measurements from the matched bounding boxes between the current and previous frame. Filter out outlier Lidar points in a statistically robust way to avoid severe estimation errors.

In this stage of the project, we worked to achieve accurate lidar-based measurements of the time-to-collision, starting with some fine-tuning. The following images show the lidar point cloud with the ego car at the origin. Measurements are available from a 360-degree, 3D scanning lidar sensor mounted on the roof of the ego car.

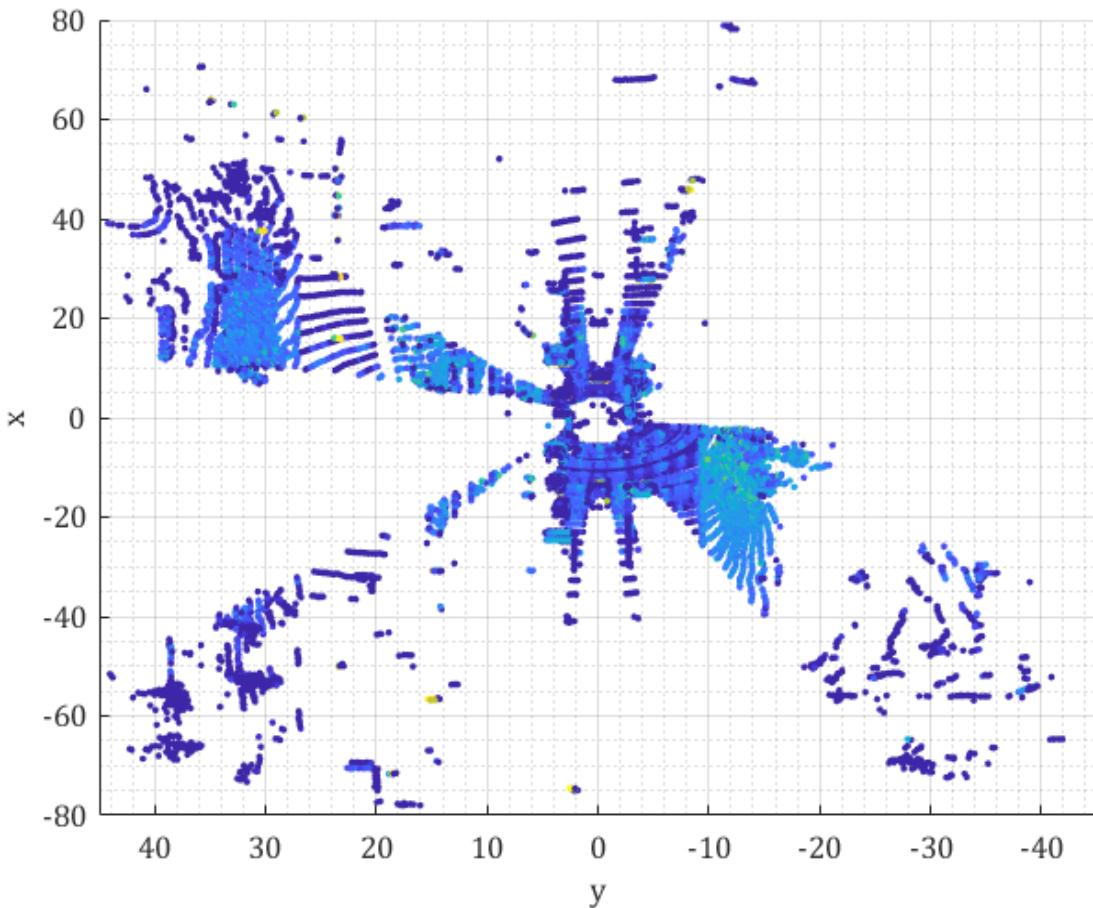


Figure 4: Top-Down View of Lidar Point Cloud

The region of interest across video frames was hand-tuned and cropped using a series of thresholds selected to exclude lidar points that were (1) outside of the ego lane, beyond 2 meters to the left or the right, (2) at or below ground level at -1.5 meters, (3) too far above the roof the car at -0.5 meters, (4) behind the ego car with negative x value, and (5) greater than 20 meters (too far away) in the direction of motion. Cropping the lidar points is crucial for this implementation. The following images show what happens if we do not perform this step:

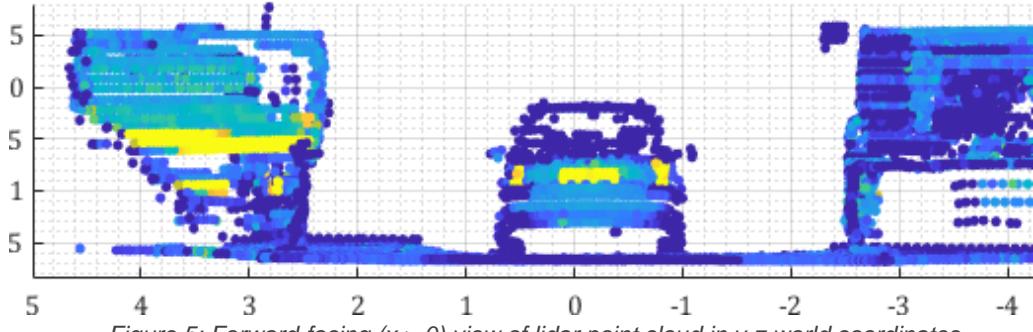


Figure 5: Forward-facing ($x > 0$) view of lidar point cloud in y - z world coordinates

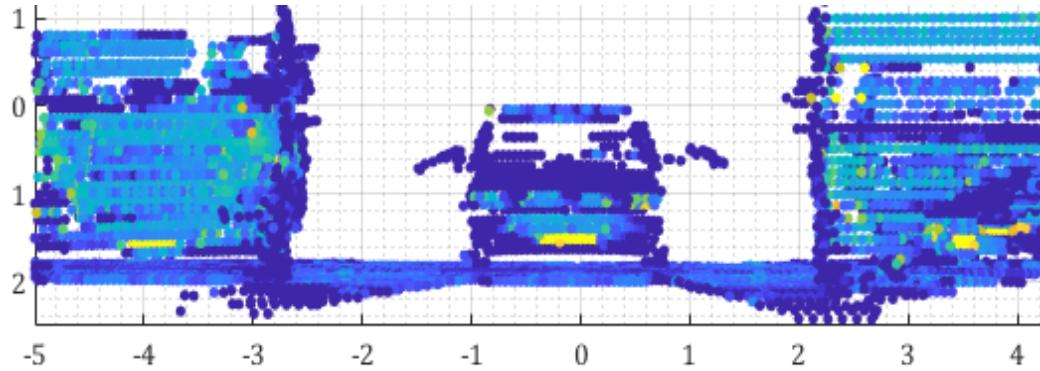


Figure 6: Backward-facing ($x < 0$) view of lidar point cloud in y - z world coordinates

The transformation from 3D world coordinates to 2D camera image results in the lidar points behind the vehicle forming an inverted image stacked above the scene in front of the vehicle:

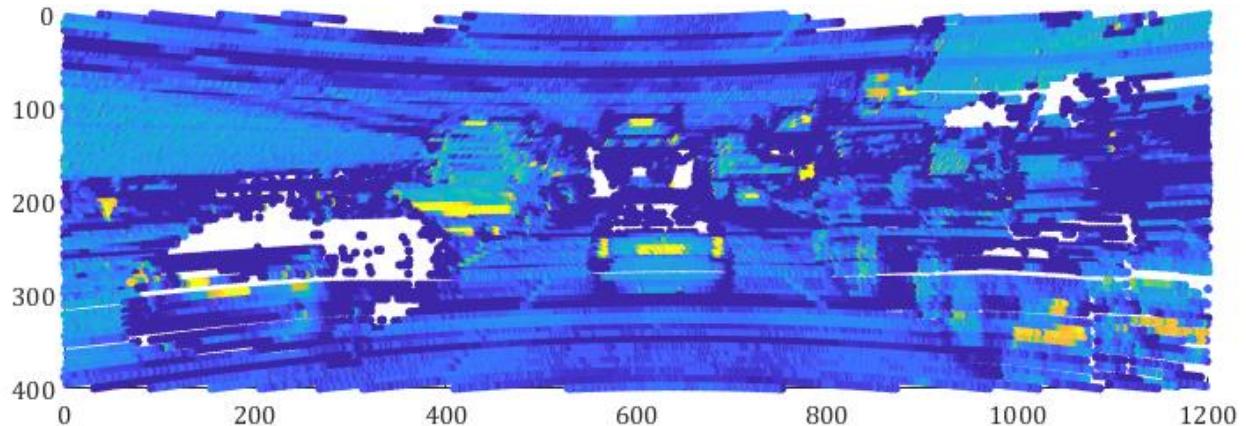


Figure 7: Full lidar point cloud mapped to 2D camera image plane, no cropping

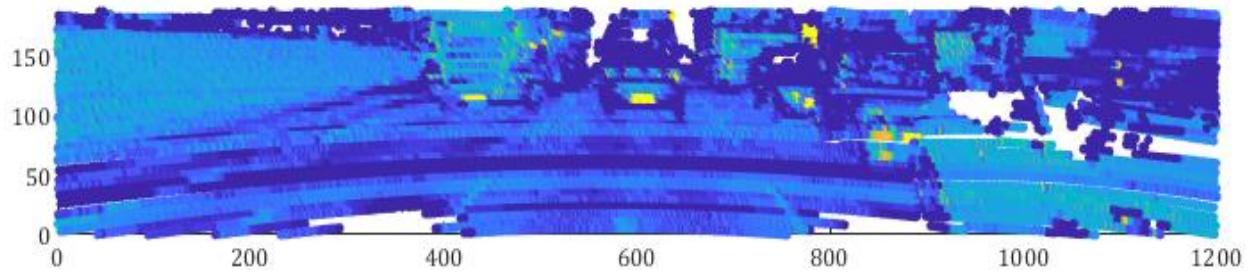


Figure 8: Inverted upper-half of uncropped 2D mapping shows the scene behind the ego car

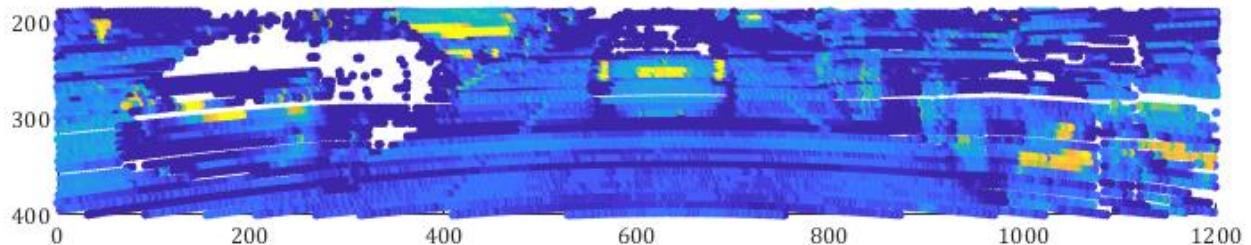


Figure 9: Bottom-half of uncropped 2D camera transformation contains the ego car's forward-facing view

Excluding negative x-values (and nothing else) produces the result shown below. The overlay of 2D-mapped lidar points stops approximately halfway up the image due to physical/hardware constraints in the lidar height and range.



Figure 10: Removed all points behind the ego car and mapped the result to 2D camera image frame

Next, we removed points from the road surface. A top-down view before and after performing this step:

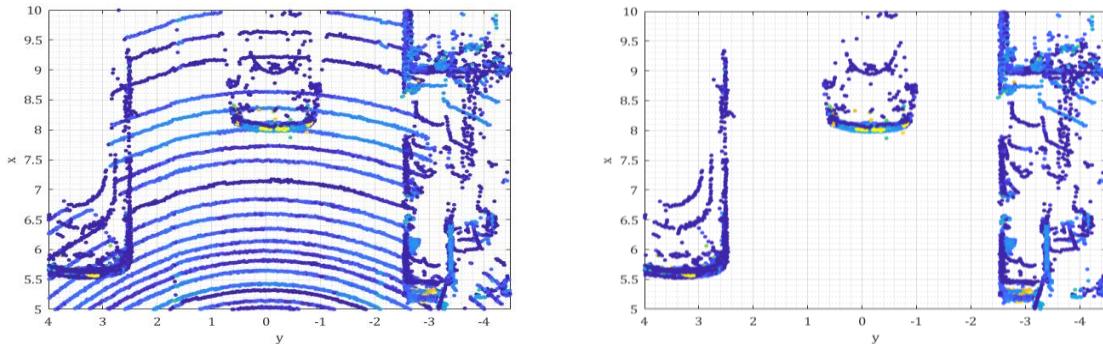


Figure 11: Removed measurements originating from the road surface

Finally, the image below shows the 2D overlay of lidar data after cropping out points according to the specifications described earlier: “not too high”, “not too low”, “not to far in front”, “not behind”, “not outside of the ego lane”.



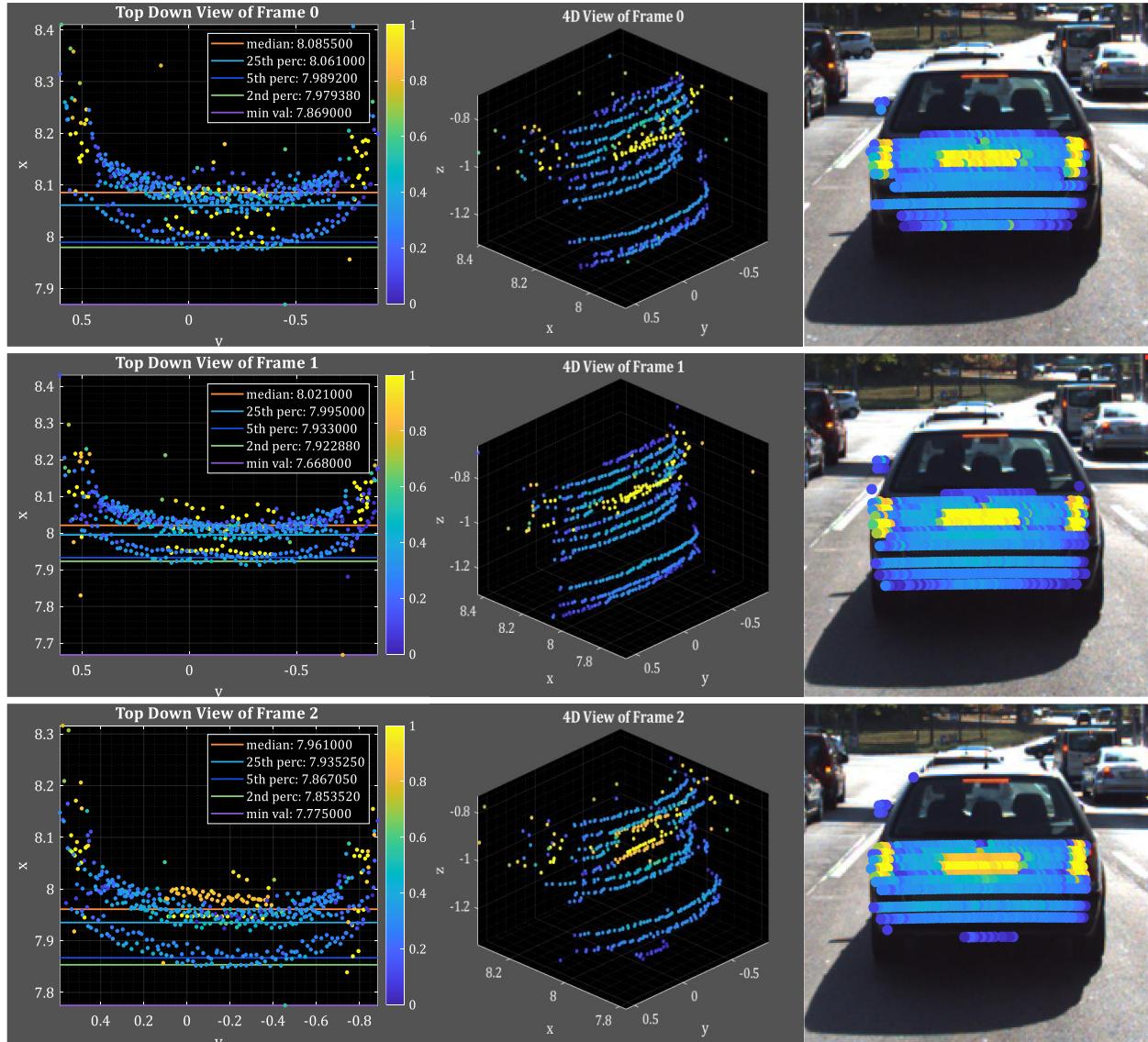
Figure 12: Cropped Lidar Point Cloud Mapped to 2D Camera Image Frame

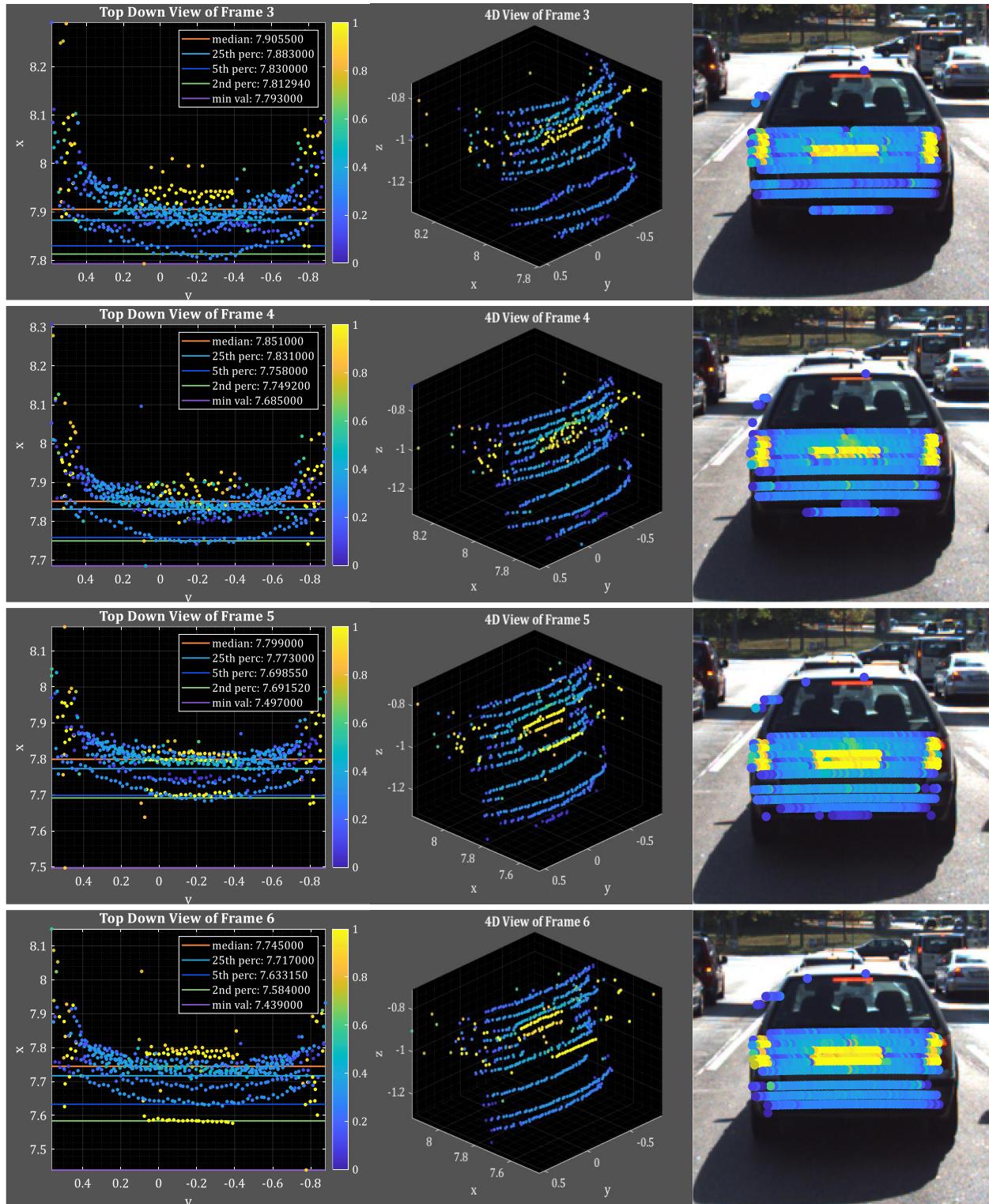
The TTC estimation space was further constrained by applying the (10% reduced) YOLO bounding box to the 2D mapped lidar points. This usually narrowed down the list to one bounding box:

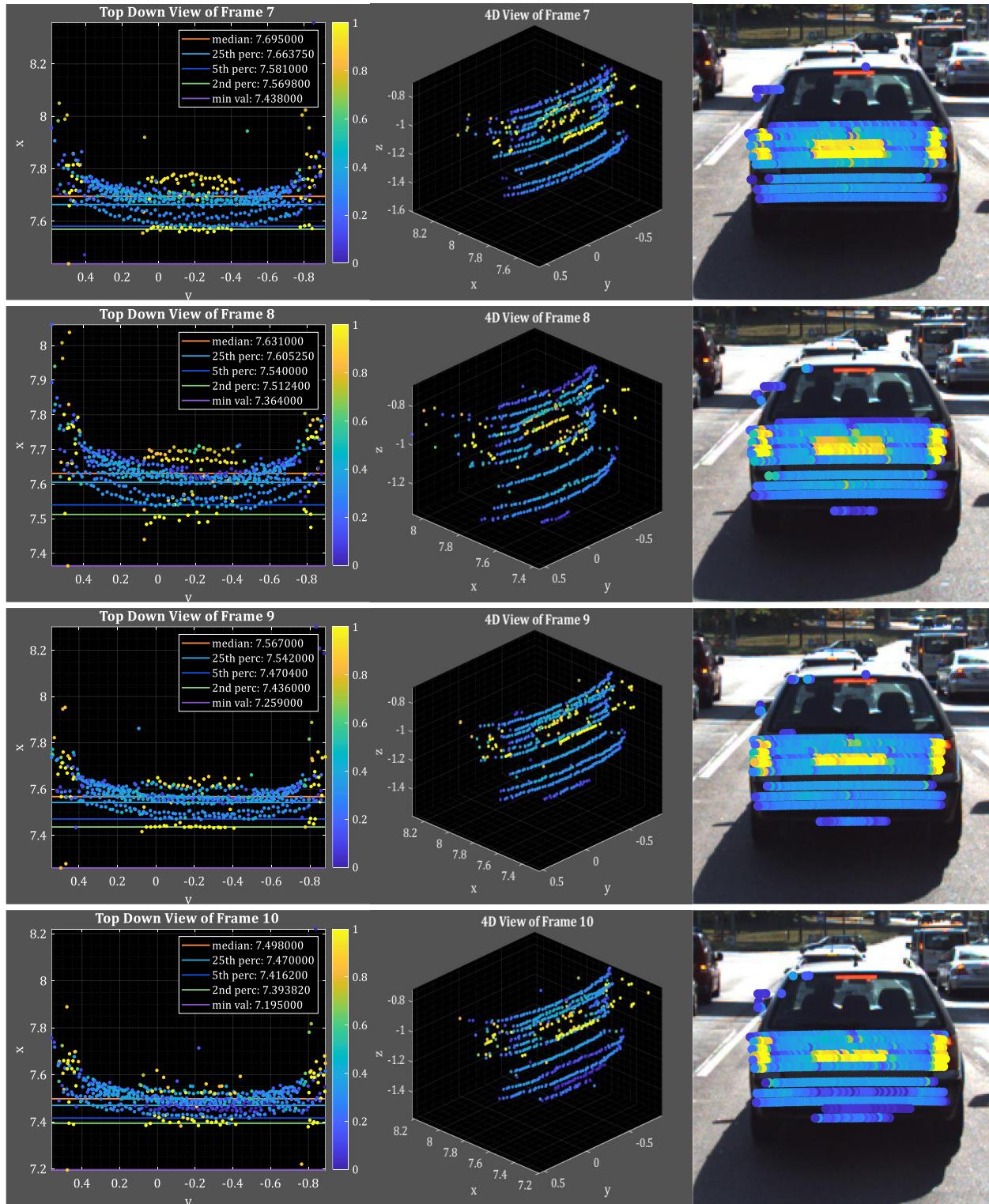


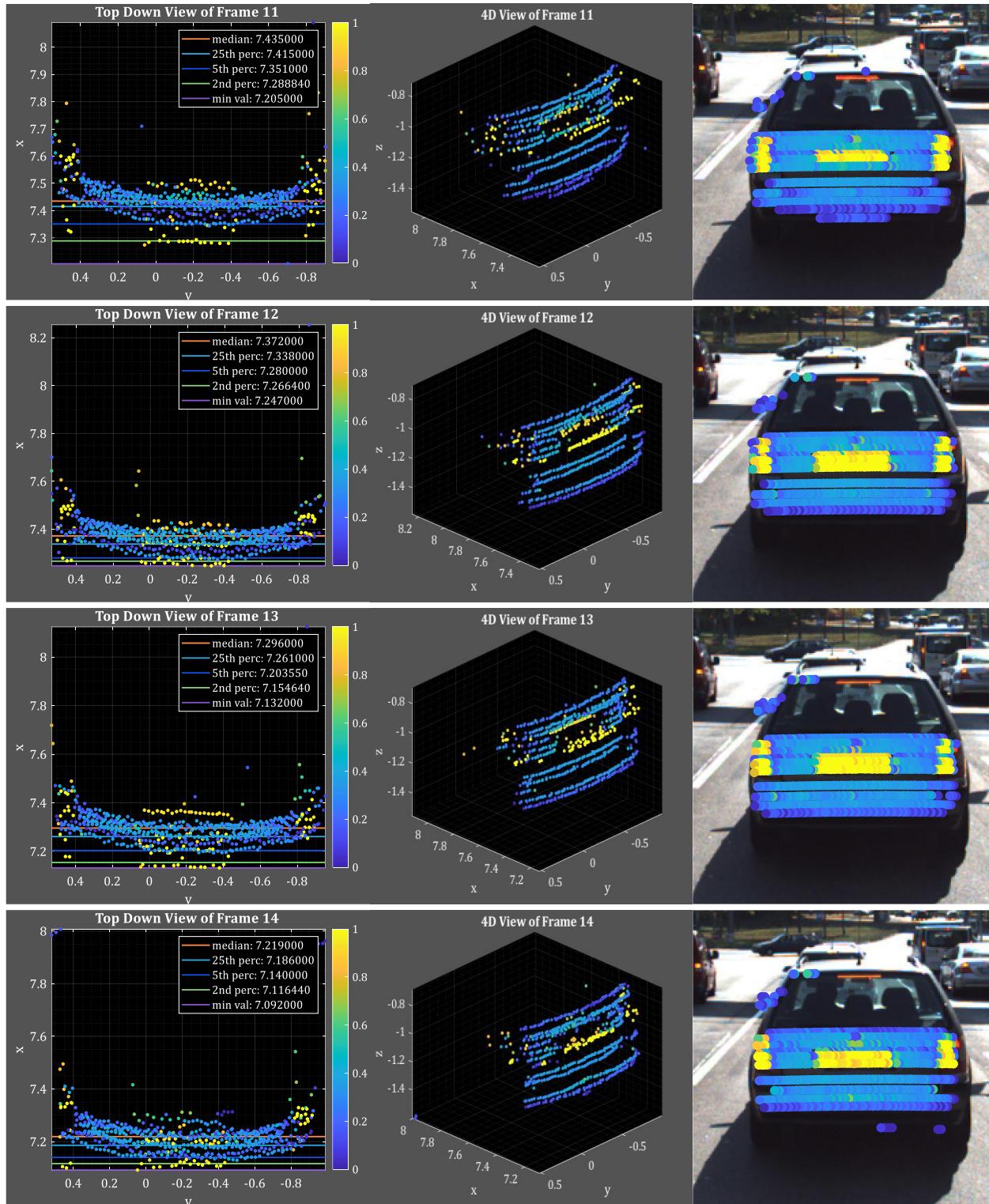
Figure 13: LIDAR Points Contained Within YOLO Bounding Box on Preceding Car (2-frame example)

The distance to the preceding vehicle was estimated based on data within the region of interest. The following top-down views of the lidar measurements bouncing off the preceding vehicle show why the absolute minimum distances could not be used. Visual inspection and review of the percentiles across the distribution of x-distances were used to select “statistically robust” estimates in the presence of noise. In addition to x- and y-positions, the reflectivity of the data points was displayed according to the color scales shown next to each plot. The reflectivity information may show how variations in material properties (car paint, license plate, headlights, etc) and angle of incidence (rear-facing vs edges) contribute to noise in the position estimates.









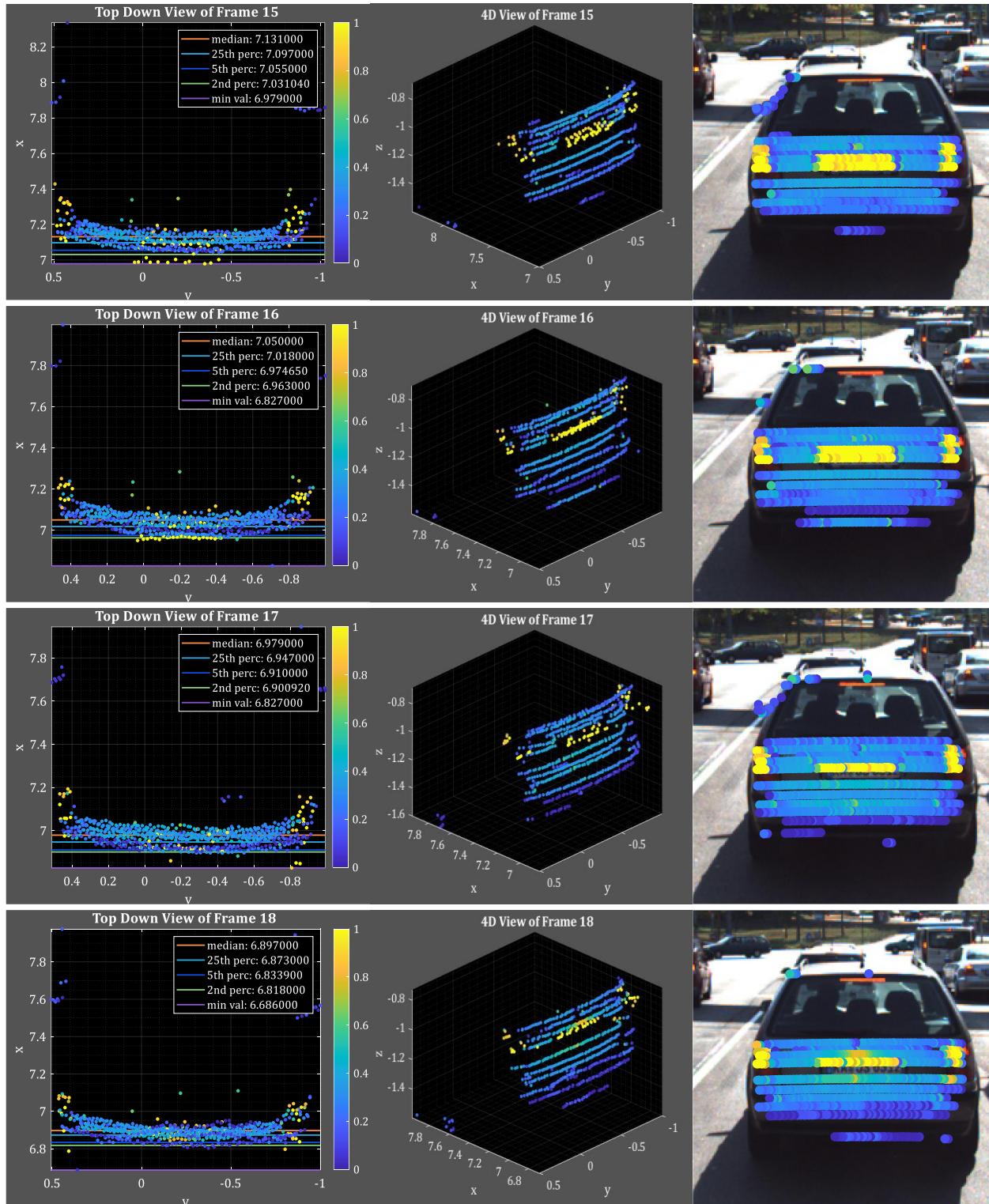


Figure 14: LIDAR point cloud views from different angles on each frame

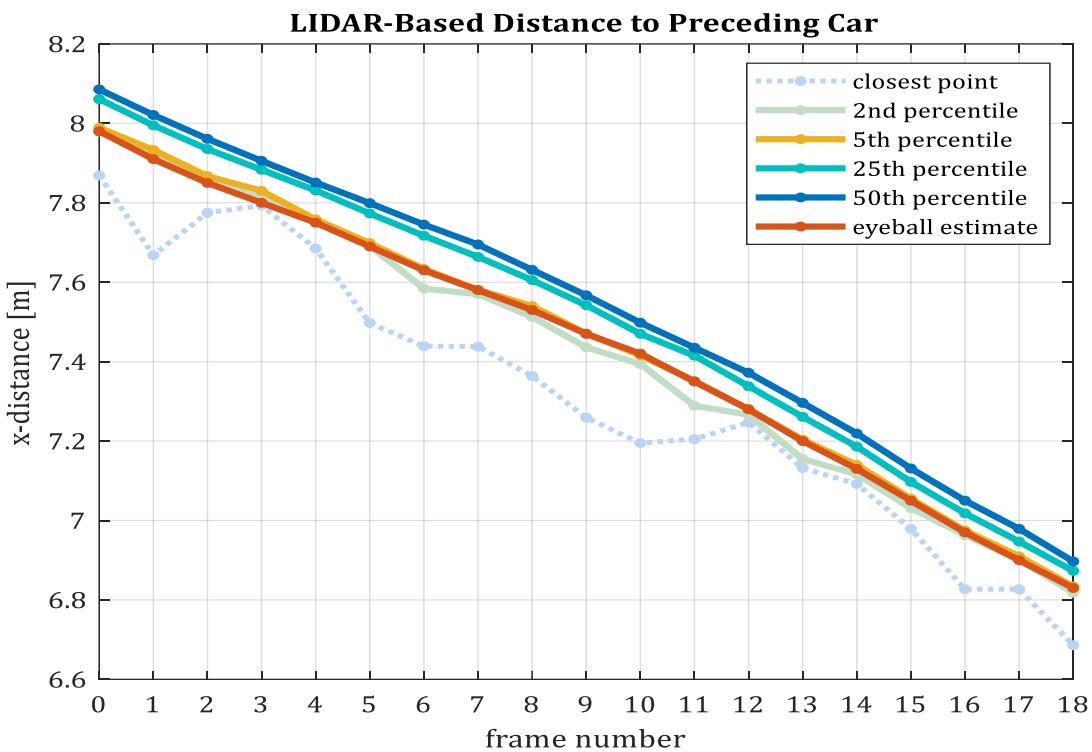


Figure 15: Distribution of lidar x-distances across frames

As discussed in the lecture notes, the ego car is equipped with a lidar sensor that can measure distances to obstacles at sequential time steps. For the purposes of estimating the time-to-collision, we are interested in the distances to the closest 3D point in the path of driving, illustrated here as distance d_0 at time 1, and distance d_1 at time 2:

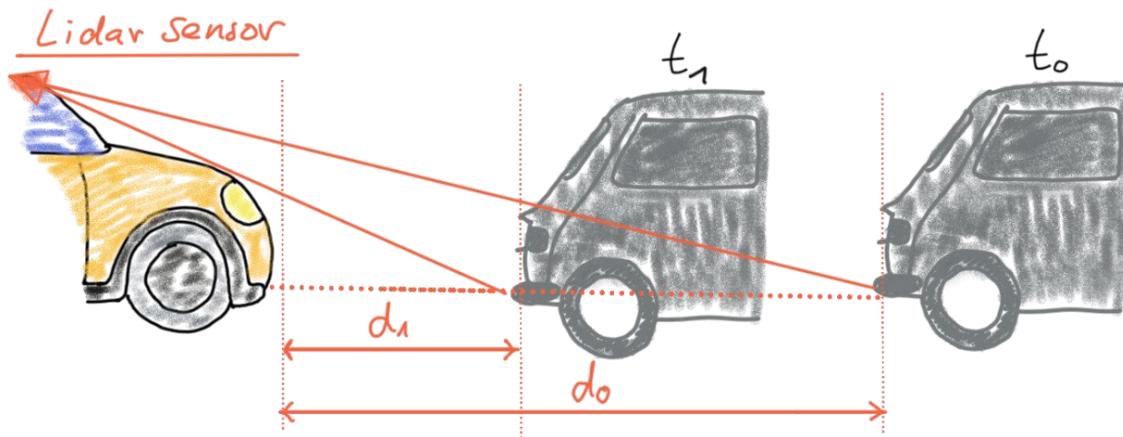


Figure 16: Illustration of the variables needed to compute the lidar-based time-to-collision

Based on the model of constant velocity, we can compute the time-to-collision between two successive lidar measurements as follows:

$$TTC = \frac{d_1 * \Delta t}{d_0 - d_1}$$

Note that the distances d_0 and d_1 , are measured from the front end of the ego car to the back end of the preceding car, not from the roof-mounted lidar sensor to the preceding car. There is an extrinsic calibration from the sensor position to the front-end position of the ego car that needs to be computed. Since we do not have the coefficients for this transformation, I used an “eyeball” estimate of 4.6m for this offset based on a top-down view of when we start to see the road below the front of the vehicle:

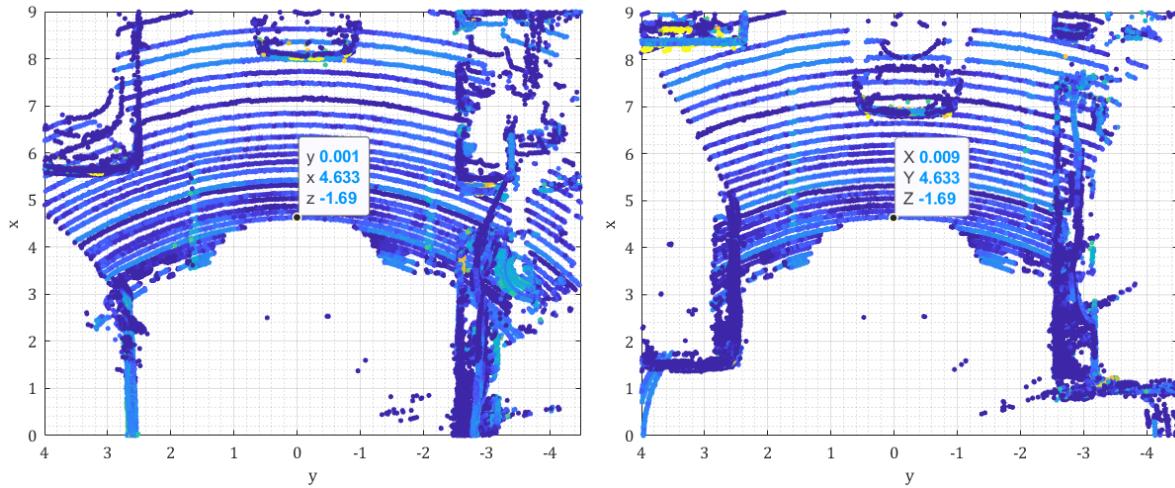


Figure 17: Estimated distance from the roof-mounted lidar sensor to the front-edge of ego car: Frames 0 and 18

Puzzlingly, subtracting an offset of 4.6m from all measured x-distances produced lidar-based TTC estimates that were about 4.6m too small relative to the camera-based TTC estimates (described later in this report). In the absence of truth, I relied on such similarities to cross-validate the implementations. For the purposes of this project, I decided not to subtract the offset from the x-distance measurements. The following graph below shows the TTC estimates *without* the lidar-to-edge distance compensation.

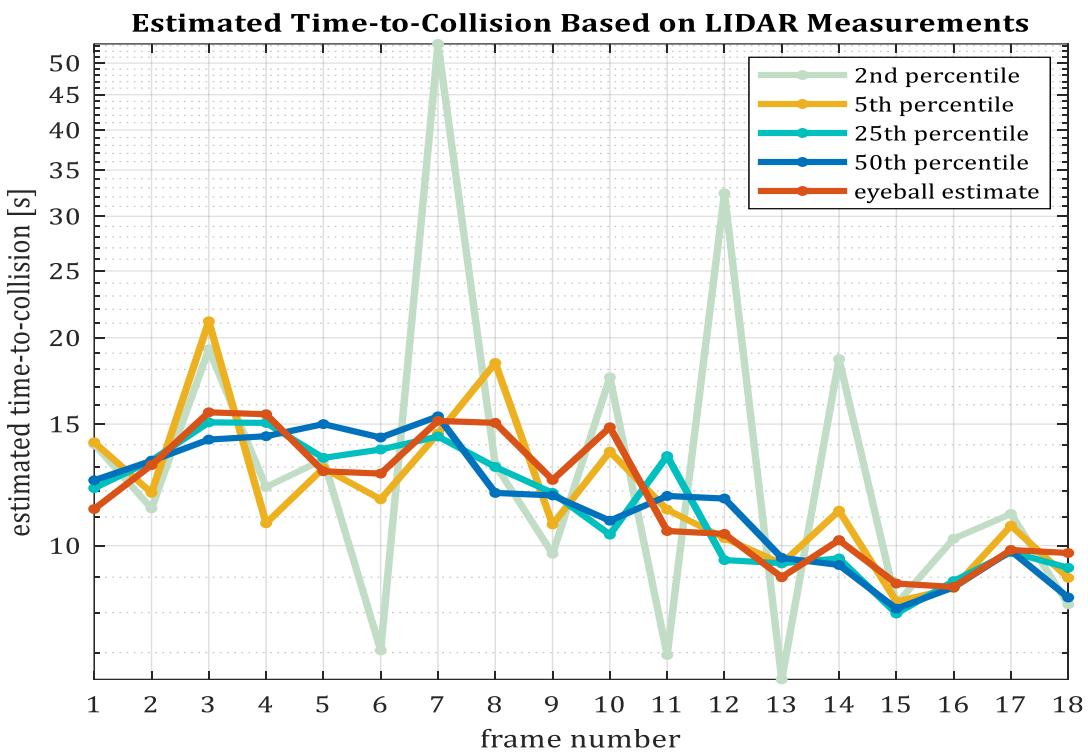


Figure 18: Lidar-based TTC estimates computed at percentiles of the x-distance distribution in the region of interest

Fluctuations in the lidar-based TTC estimates over time demonstrated a high-sensitivity to noise, as shown in the figures above. Based on this graph, it seems that the 25th and 50th raw percentiles provide the most stable frame-to-frame estimates, at the expense of slightly over-estimating the time-to-collision. A temporal smoothing approach could potentially enable us to choose a percentile that runs closer to the back end of the vehicle. For the purposes of this quick-study, I decided to forego implementation of a temporal smoothing filter and selected the x-position corresponding to the 25th percentile of the data set.

Function Definition of LIDAR-based Time-to-Collision in camFusion_student.cpp:

```
// Compute the time-to-collision for all matched 3D objects based on lidar measurements alone
void computeTTCLidar(std::vector<LidarPoint> &lidarPointsPrev,
                      std::vector<LidarPoint> &lidarPointsCurr, double frameRate, double &TTC)
{
    // reformat x (distance) values as a vector of doubles
    std::vector<double> distPrev;
    std::vector<double> distCurr;

    for (auto it = lidarPointsPrev.begin(); it != lidarPointsPrev.end(); ++it)
        distPrev.push_back(it->x);

    for (auto it = lidarPointsCurr.begin(); it != lidarPointsCurr.end(); ++it)
        distCurr.push_back(it->x);

    // take nth smallest point as the distance to use (anything smaller is discarded as noise)
    int nthP, nthC;
    float prctile = 0.25; //0.5; // nth point is based on a fraction of the total number of points
    nthP = floor(prctile * distPrev.size());
    std::nth_element(distPrev.begin(), distPrev.begin() + (nthP - 1), distPrev.end());
    nthC = floor(prctile * distCurr.size());
    std::nth_element(distCurr.begin(), distCurr.begin() + (nthC - 1), distCurr.end());

    // based on the model of constant velocity, the TTC can be computed from two successive lidar measurements as follows:
    TTC = (distCurr[nthC - 1] * (1 / frameRate)) / (distPrev[nthP - 1] - distCurr[nthC - 1]);
}
```

Function Call in main():

```
// if the bounding boxes on the curr and prev frames have lidar points associated with them
// (note: lidar points were previously cropped to capture obstacles only in the ego lane)
if (currBB->lidarPoints.size() > 0 && prevBB->lidarPoints.size() > 0)
{
    double ttclidar = -1;
    // then compute time-to-collision based on LIDAR data for the current match
    computeTTCLidar(prevBB->lidarPoints, currBB->lidarPoints, sensorFrameRate, ttclidar);
```

FP.3 ASSOCIATE KEYPOINT MATCHES WITH BOUNDING BOXES

i Prepare the TTC computation based on camera measurements by associating keypoint correspondences to the bounding boxes that enclose them. All matches that satisfy this condition must be added to a vector representing the respective bounding box. Remove outlier matches based on the Euclidean distance between them in relation to all matches in the bounding box.

In this section, we computed the Euclidean distance between all keypoints within the bounding box (usually one, but sometimes two) on the preceding vehicle, which were identified with the help of the highly-constrained lidar space defined in the previous section. “Outlier matches” were identified and removed, and the selected matched keypoints were stored in a dynamic list associated with the respective bounding box structures.

I experimented with 2 different methods for removing outlier matches: (1) RANSAC and (2) an empirically-determined threshold multiplier applied to the mean distance between keypoint matches. Statistically, the approaches produced comparable TTC estimates with reasonable repeatability across frames. However, a cursory visual inspection of both approaches showed that RANSAC did a better job of identifying out “bad” keypoint matches on the road, surrounding objects or shadows.

Earlier implementations of this approach showed significant errors in the estimated TTC due to keypoints on adjacent cars, the road, shadows, or other elements that fell within the bounding box but did not physically belong to the preceding car. To mitigate these effects, I added the following pre-filtering steps:

- 1) Excluding keypoints that did not belong to any bounding box
- 2) Removing keypoints that fell within more than one YOLO bounding box
- 3) Shrinking the bounding box slightly to avoid having too many outlier points around the edges of the vehicle

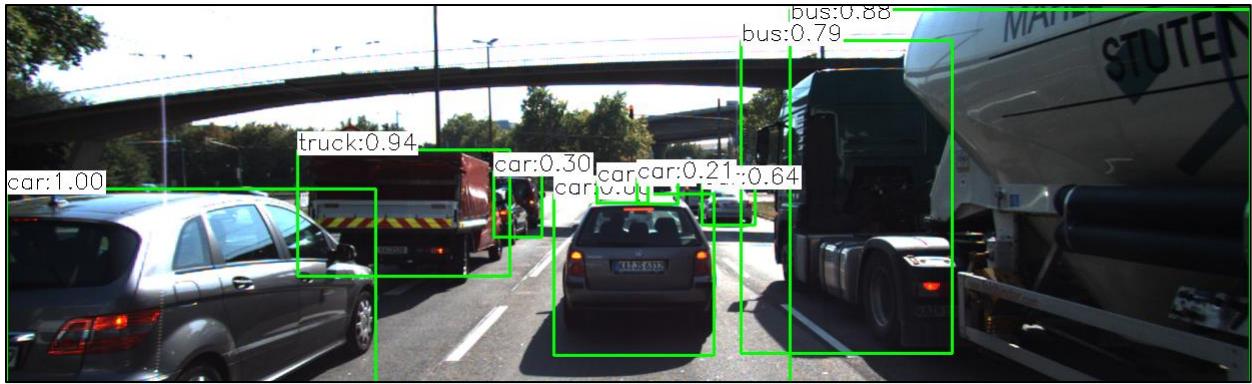


Figure 19: YOLO Bounding Boxes on Frame 7

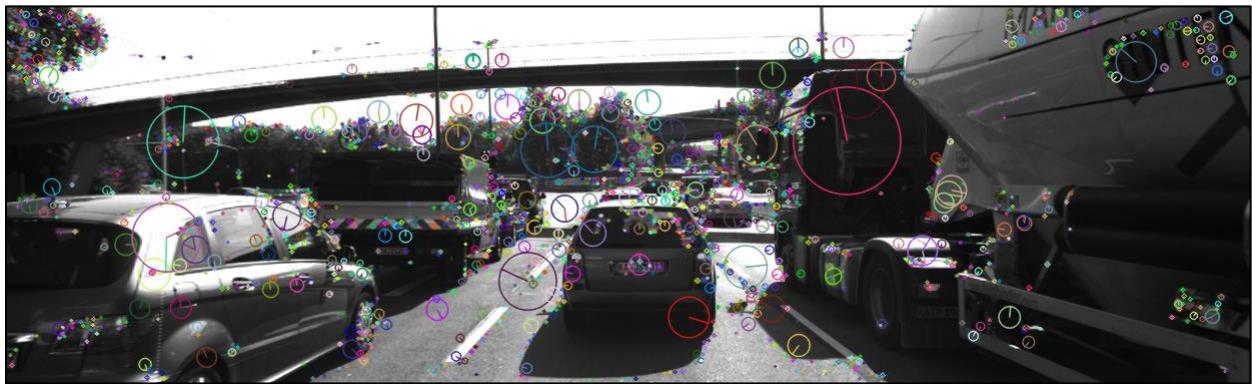


Figure 20: Original SIFT Keypoints on Frame 7

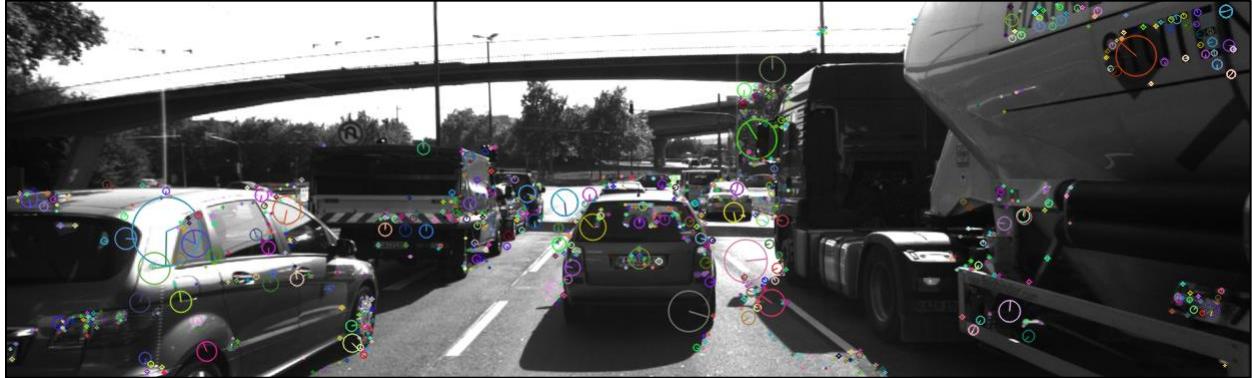


Figure 21: SIFT Keypoints Filtered to Non-Overlapping Bounding Boxes

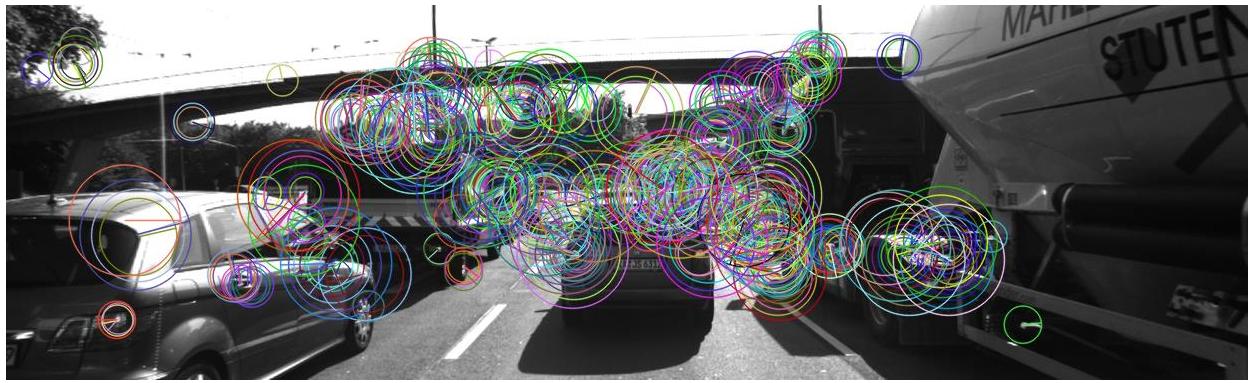


Figure 22: Original ORB Keypoints on Frame 7

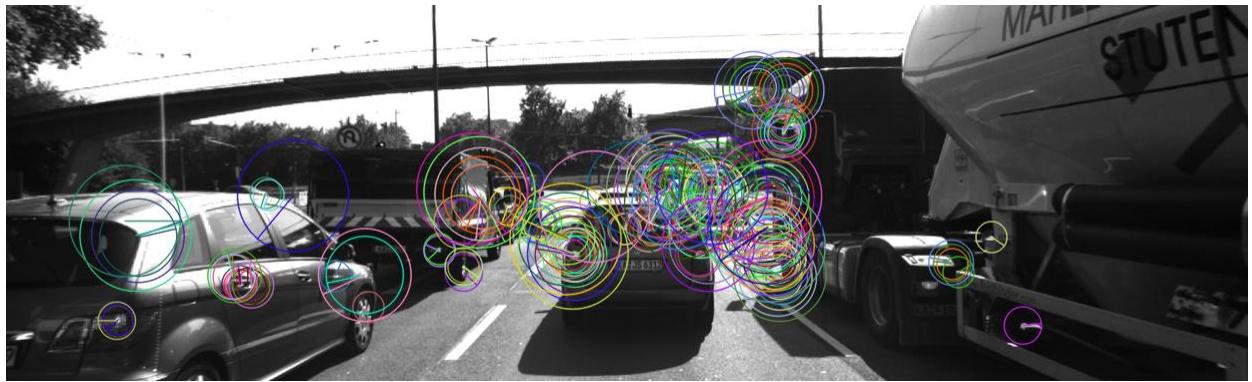


Figure 23: ORB Keypoints Filtered to Non-Overlapping Bounding Boxes



Figure 24: Original Shi-Tomasi Keypoints on Frame 7



Figure 25: Shi-Tomasi Keypoints Filtered to Non-Overlapping Bounding Boxes

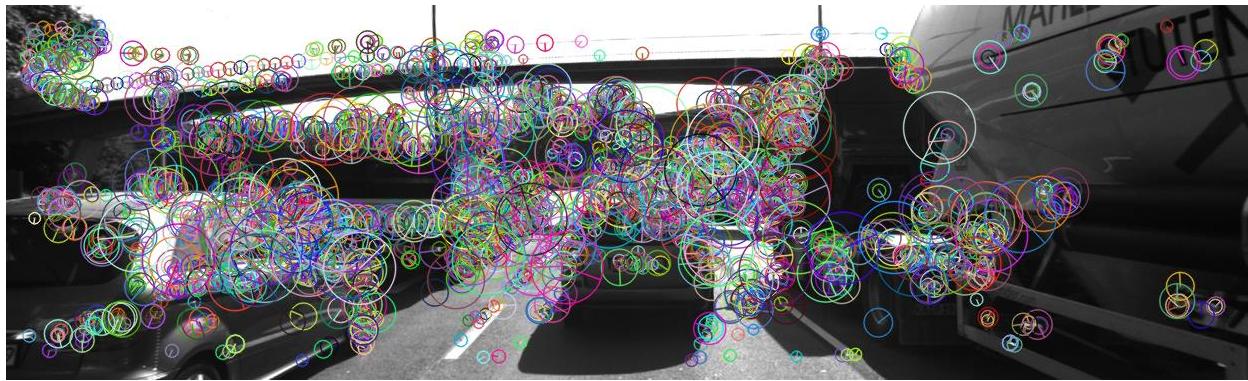


Figure 26: Original BRISK Keypoints on Frame 7

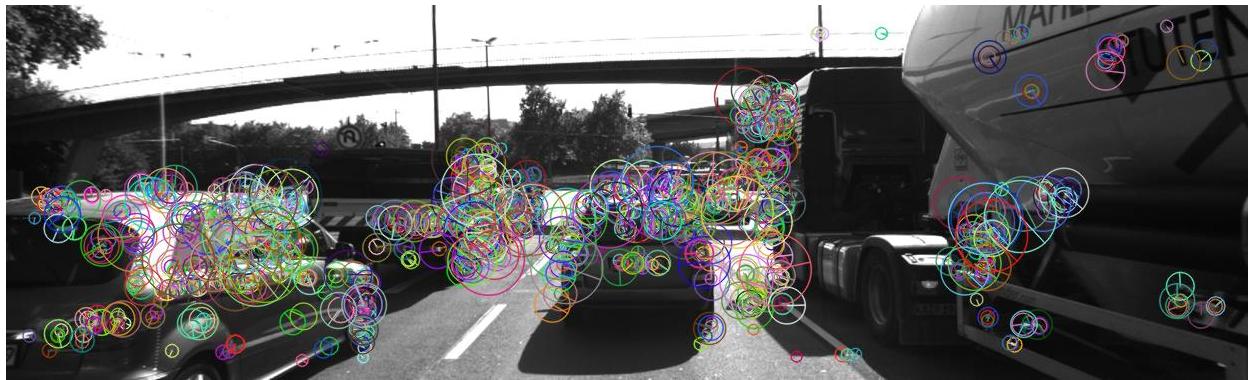


Figure 27: BRISK Keypoints Filtered to Non-Overlapping Bounding Boxes

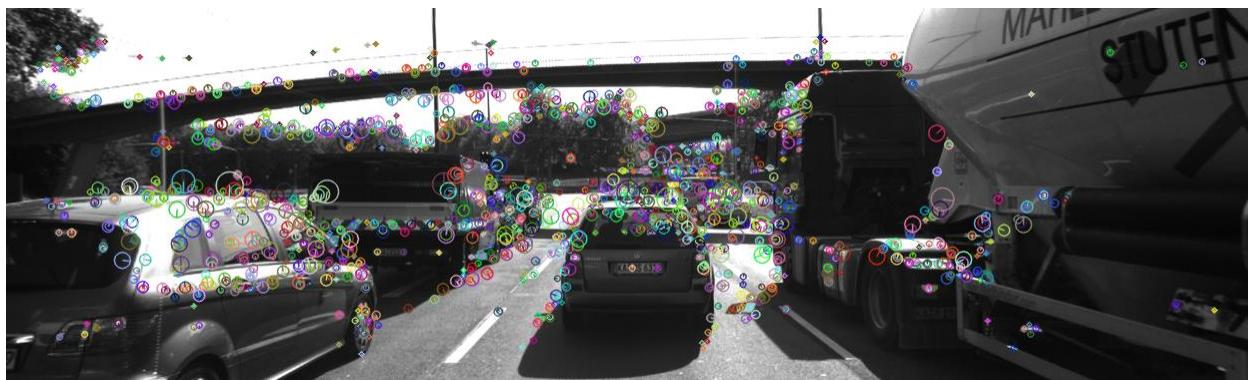


Figure 28: Original AKAZE Keypoints on Frame 7

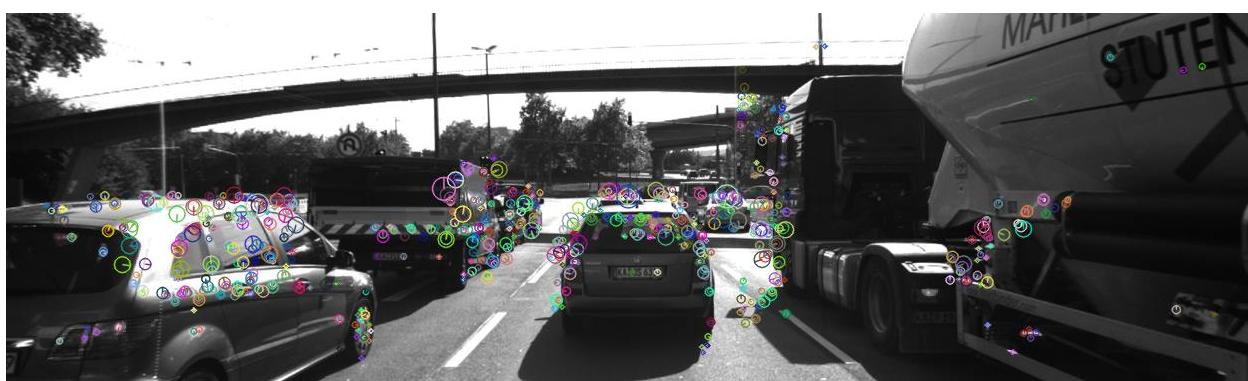


Figure 29: AKAZE Keypoints Filtered to Non-Overlapping Bounding Boxes

Function Definition to Reduce Keypoints to Non-Overlapping Bounding Boxes:

```

std::vector<cv::KeyPoint> eraseKptsInOverlappingBBs(std::vector<BoundingBox> &boundingBoxes,
                                                    std::vector<cv::KeyPoint> &keypoints)
{
    std::vector<cv::KeyPoint> filteredKeypoints;
    // loop over all keypoints and associate them to a 2D bounding box
    for (auto currKpt : keypoints)
    {
        int nBBsContainingKpt = 0;
        for (auto it2 = boundingBoxes.begin(); it2 != boundingBoxes.end(); ++it2)
        {
            // shrink current bounding box slightly to avoid having too many outlier points around the edges
            cv::Rect smallerBox;
            smallerBox.x = it2->roi.x;
            smallerBox.y = it2->roi.y;
            smallerBox.width = it2->roi.width;
            smallerBox.height = it2->roi.height;

            // check whether point is within current bounding box
            if (smallerBox.contains(currKpt.pt))
                nBBsContainingKpt++;

        } // end loop over all bounding boxes

        // in order to avoid inadvertent association of a keypoint on one vehicle with another vehicle,
        // exclude keypoints that are enclosed within multiple bounding boxes from further processing
        if (nBBsContainingKpt == 1) // if keypoint is enclosed by exactly one bounding box
            filteredKeypoints.push_back(currKpt);

    } // end loop over all keypoints
    return filteredKeypoints;
}

```

The following images are examples of keypoint matches associated to the nearest ego lane bounding box between frames, where red markers represent keypoints on the previous frame, and green markers represent the corresponding keypoint matches on the current frame:

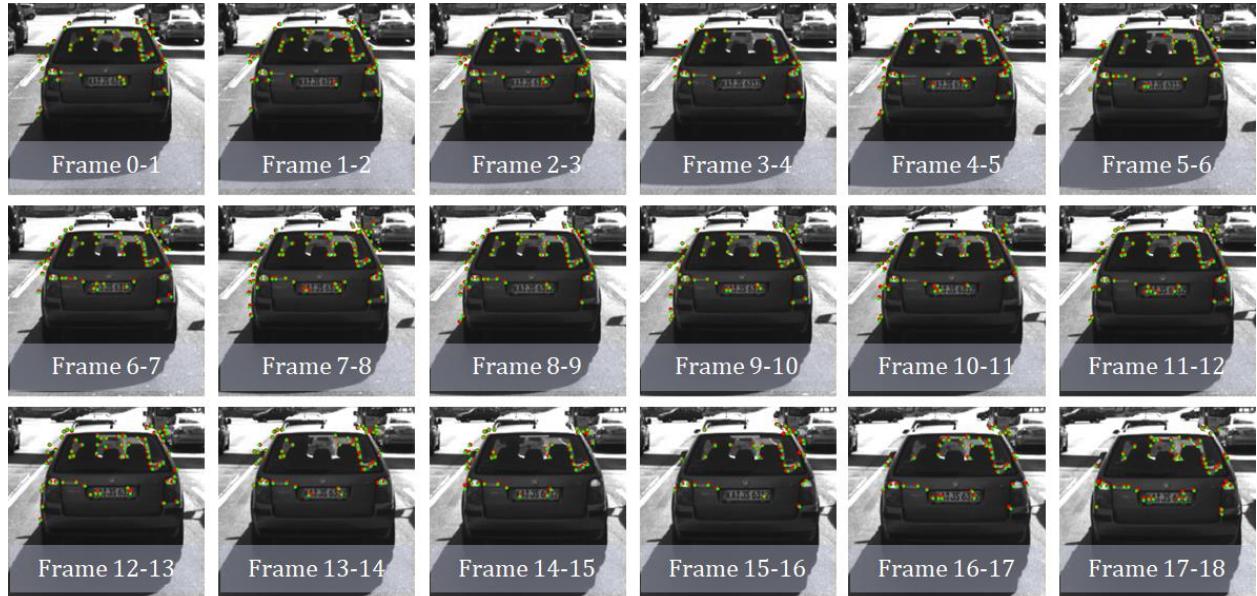


Figure 30: Frame-to-Frame Bounding Box Matches Using Shi-Tomasi Detector - SIFT Descriptor

Function Definition to Associate Keypoint Matches to Bounding Boxes:

```

// associate a given bounding box with the keypoints it contains
void clusterKptMatchesWithROI(BoundingBox &boundingBoxPrev, BoundingBox &boundingBox,
                               std::vector<cv::KeyPoint> &kptsPrev, std::vector<cv::KeyPoint> &kptsCurr,
                               std::vector<cv::DMatch> &kptMatches)
{
    // init variables for mean-based match selection
    double meanDist = 0.0, sumDist = 0.0;
    int numGoodMatches = 0;
    std::map<cv::DMatch, double> bbMatches;
    std::vector<cv::KeyPoint> tmpKpts, tmpKptsPrev;

    // init variables to find good matches using RANSAC
    std::vector<cv::Point2f> pointsPrev, pointsCurr;
    std::vector<cv::DMatch> bbMatches2;

    for (auto match : kptMatches)
    {
        // shrink current bounding box slightly to avoid having too many outlier points around the edges
        cv::Rect smallerBox;
        smallerBox.x = boundingBox.roi.x + 0.05 * boundingBox.roi.width / 2.0;
        smallerBox.y = boundingBox.roi.y + 0.05 * boundingBox.roi.height / 2.0;
        smallerBox.width = boundingBox.roi.width * (1 - 0.01);
        smallerBox.height = boundingBox.roi.height * (1 - 0.2);

        // if the smaller bounding box contains the current keypoint in the current match
        // note: kptMatches are ordered as follows: kptsPrev uses queryIdx and kptsCurr uses trainIdx
        if (smallerBox.contains(kptsCurr[match.trainIdx].pt))
        {
            // convert keypoints to Point2f (for RANSAC)
            pointsPrev.push_back(cv::Point2f(kptsPrev[match.queryIdx].pt.x, kptsPrev[match.queryIdx].pt.y));
            pointsCurr.push_back(cv::Point2f(kptsCurr[match.trainIdx].pt.x, kptsCurr[match.trainIdx].pt.y));

            // save keypoints within the bounding box
            tmpKpts.push_back(kptsCurr[match.trainIdx]);
            tmpKptsPrev.push_back(kptsPrev[match.queryIdx]);

            // compute the distance between the previous and current positions of the keypoint
            double distance = cv::norm(kptsCurr[match.trainIdx].pt - kptsPrev[match.queryIdx].pt);
            bbMatches[match] = distance; // save distance to match structure
            bbMatches2.push_back(match);
            sumDist += distance; // keep a running sum of distances between mapped keypoint positions
            numGoodMatches++; // keep a counter of number of good keypoint matches
        }
    }
}

```

Method 1 uses RANSAC to identify & select “good” keypoint matches:

```

bool useRansac = true;
if (useRansac)
{
    // Compute the homography matrix and inlier (good) matches using RANSAC
    std::vector<uchar> inliers; // cv::Mat inliers;
    cv::Mat H = cv::findHomography(
        cv::Mat(pointsPrev), // matching points in previous frame
        cv::Mat(pointsCurr), // matching points in current frame
        cv::RANSAC, // RANSAC method
        3, // RANSAC reprojection threshold
        inliers); // output mask where 0s indicate outliers

    // select only the inliers
    for (int i = 0; i < inliers.size(); i++)
    {
        if ((int)inliers.at(i))
        {
            boundingBox.kptMatches.push_back(bbMatches2[i]);
            boundingBoxPrev.kptMatches.push_back(bbMatches2[i]);

            boundingBox.keypoints.push_back(tmpKpts[i]);
            boundingBoxPrev.keypoints.push_back(tmpKptsPrev[i]);
        }
    }
}

```

Method 2 applies a threshold multiplier to the mean distance:

```
else
{
    // filter out matches that are not within a threshold of the mean match distance
    int kptNum = 0;
    float outlierThresh = 2.5; // empirically selected threshold
    meanDist = sumDist / numGoodMatches; // compute the average distance
    for (auto matchDist : bbMatches)      // loop over all bounding box matches
    {
        if (matchDist.second < outlierThresh * meanDist) // filter out outliers
        {
            boundingBox.kptMatches.push_back(matchDist.first);
            boundingBoxPrev.kptMatches.push_back(matchDist.first);

            boundingBox.keypoints.push_back(tmpKpts[kptNum]);
            boundingBoxPrev.keypoints.push_back(tmpKptsPrev[kptNum]);
        }
        kptNum++;
    }
}
```

Function Call in main():

```
// assign enclosed keypoint matches to the bounding box
clusterKptMatchesWithROI(*prevBB, *currBB,
    (dataBuffer.end() - 2)->keypoints,
    (dataBuffer.end() - 1)->keypoints,
    (dataBuffer.end() - 1)->kptMatches);
```

FP.4 COMPUTE CAMERA-BASED TTC

- i** Compute the Time-to-Collision for all matched 3D objects using only keypoint correspondences from the matched bounding boxes between the current and previous frames.

A mono camera can estimate the time-to-collision via perspective projection by observing relative changes in relative height (a.k.a. scale change) directly in the image.

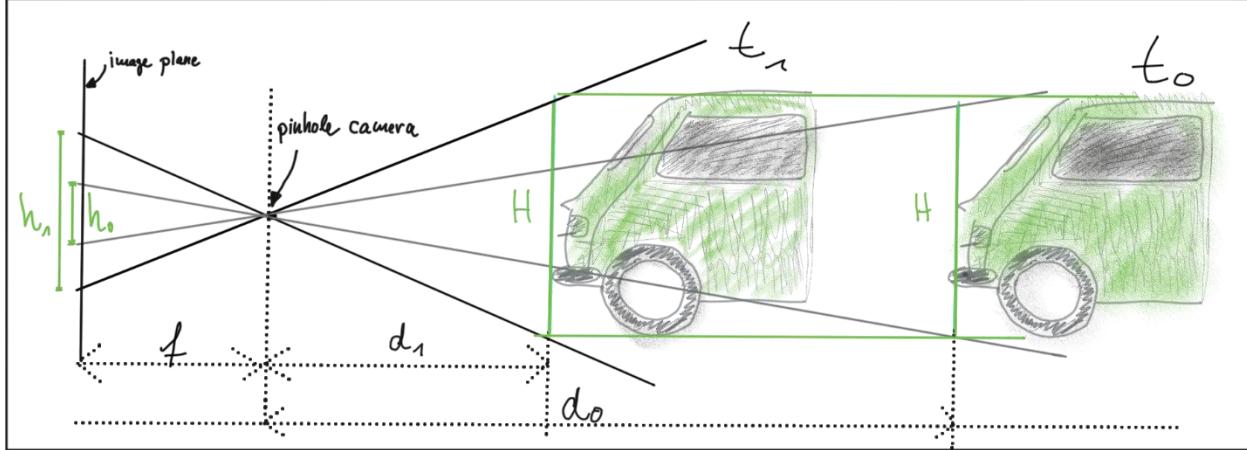


Figure 31: Estimating the Time-to-Collision Using a Single Camera

As was described in the lecture notes, we can substitute the relative distances of the keypoints captured in two temporally separated snapshots of an object moving forward or backward in front of and relative to our ego car to estimate the TTC from a single monocular camera:

$$TTC = \frac{-\Delta t}{\left(1 - \frac{h_1}{h_0}\right)}$$

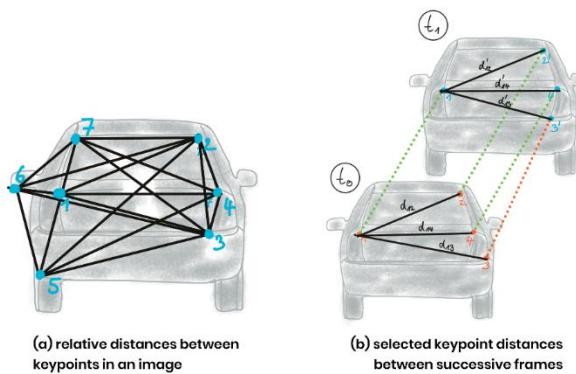


Figure 32: Using Relative Distances Between Texture-Based Keypoints to Estimate the TTC

We bound the problem to a minimum distance of 100 to exclude keypoints that are measuring the same feature in the image and limit the complexity of our first-order estimation process in the presence of multi-modal distance distributions. Reducing the number of keypoints also helps with speed and complexity of the O(n choose 2) operation. We then compute the ratio of the relative distances between keypoint matches, and the resulting TTC estimate between frames in the video sequence.

Some examples are included below (for a complete table of results, see the last section of this report):

Example 1: AKAZE Detector, AKAZE Descriptor

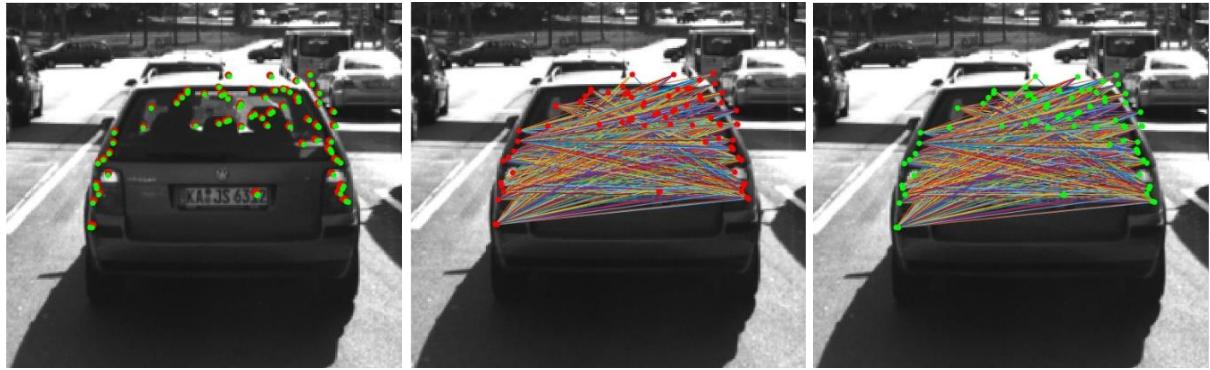


Figure 33: Relative Distances Between AKAZE detectors and AKAZE descriptors on Frames 17 and 18

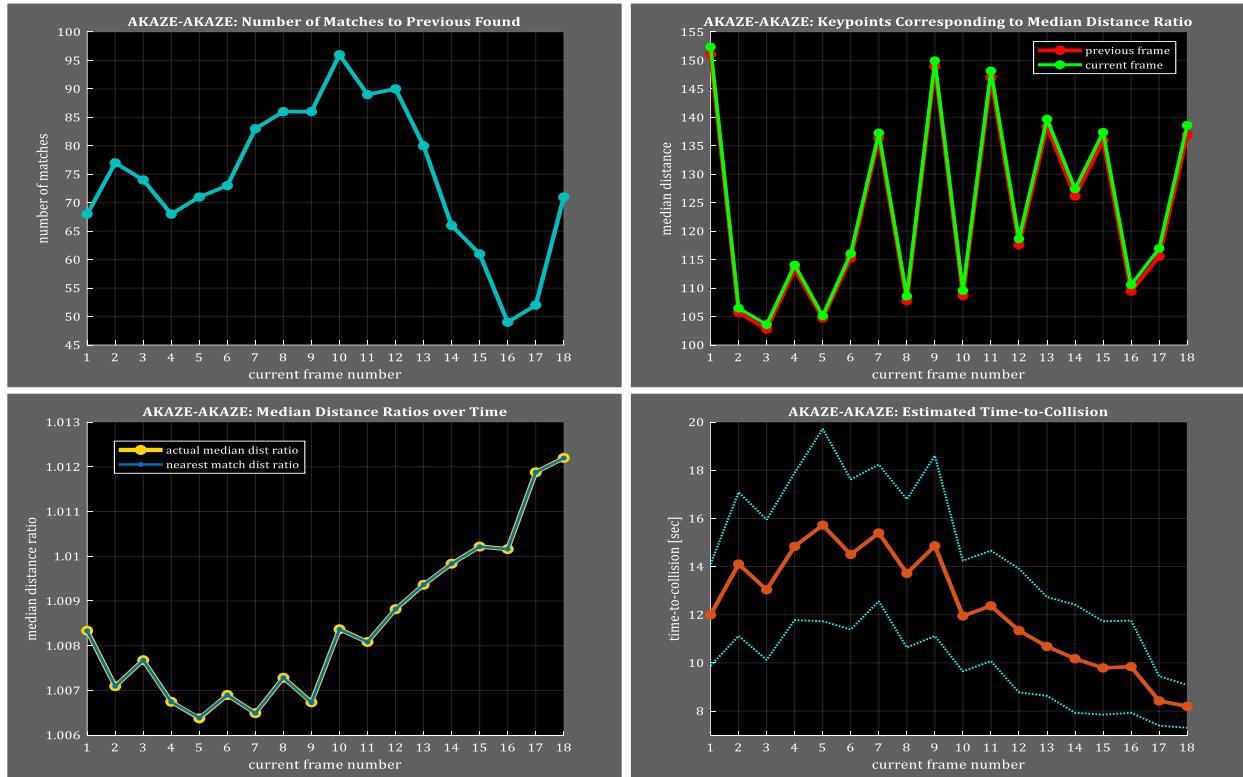


Figure 34: Median Distance Ratio and Time-to-Collision across all frames for AKAZE Detectors & AKAZE Descriptors

Example 2: BRISK Detector, BRIEF Descriptor

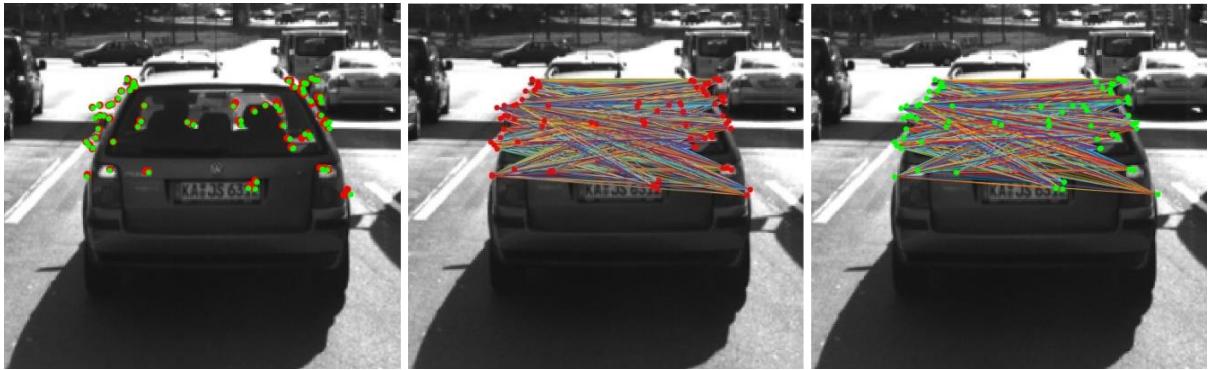


Figure 35: Relative Distances Between BRISK Detectors and BRIEF Descriptors Between Frames 14 and 15

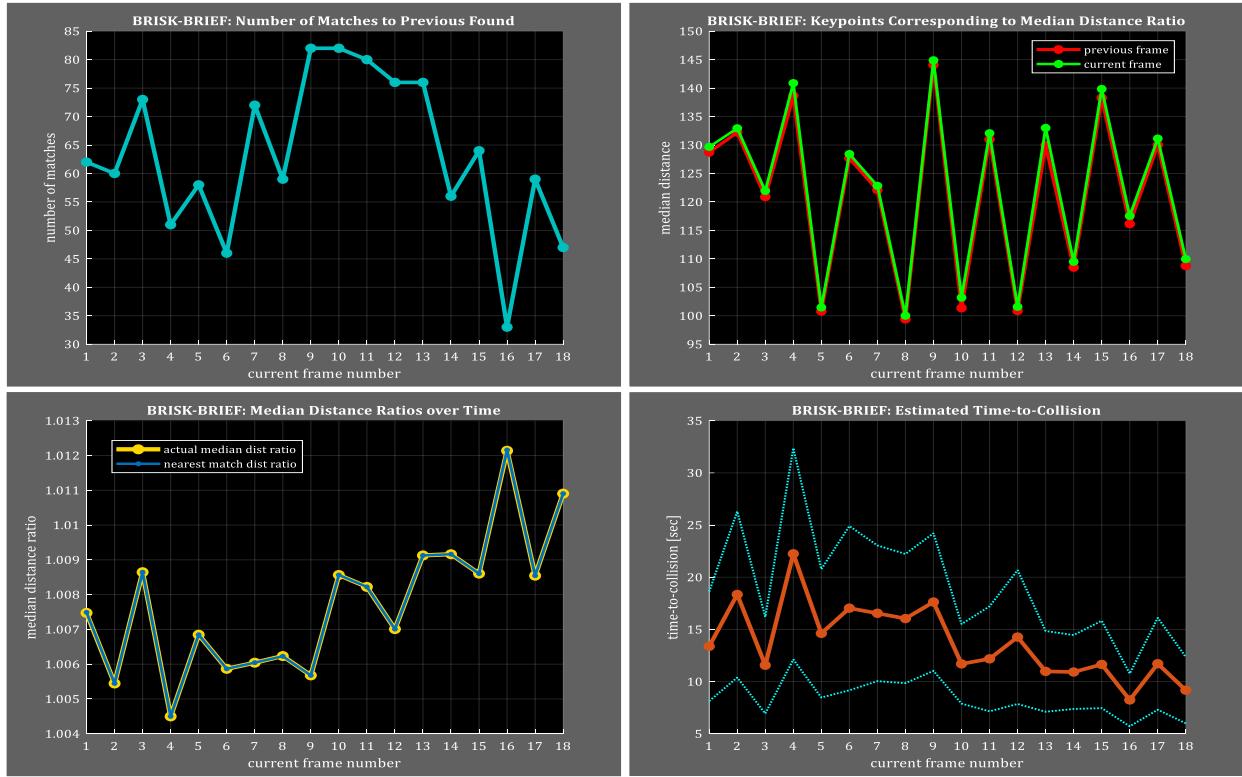


Figure 36: Median Distance Ratio and Time-to-Collision across all frames for BRISK Detectors & BRIEF Descriptors

Example 3: FAST Detector, BRISK Descriptor



Figure 37: Relative Distances Between FAST Detectors and BRISK Descriptors Between Frames 9 and 10

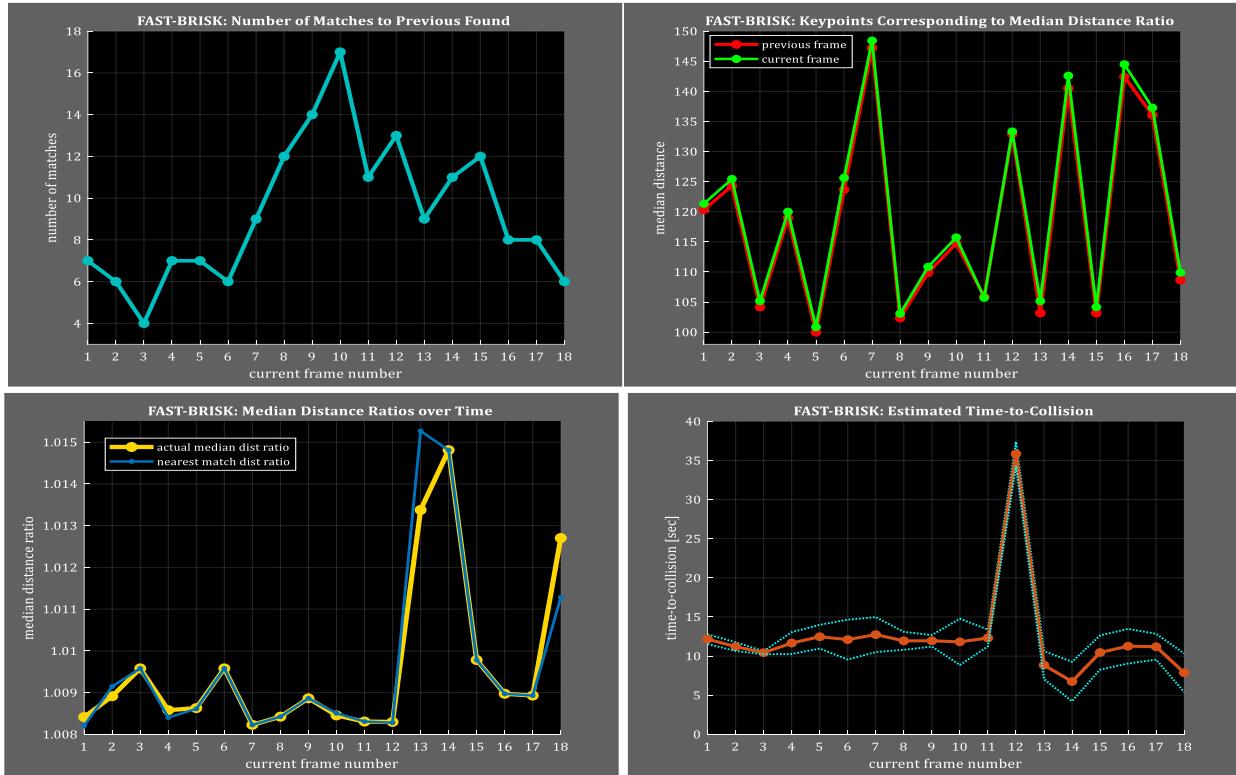


Figure 38: Median Distance Ratio and Time-to-Collision across all frames for FAST Detectors & BRISK Descriptors

Example 4: HARRIS Detector, FREAK Descriptor



Figure 39: Relative Distances Between HARRIS Detectors and FREAK Descriptors Between Frames 6 and 7

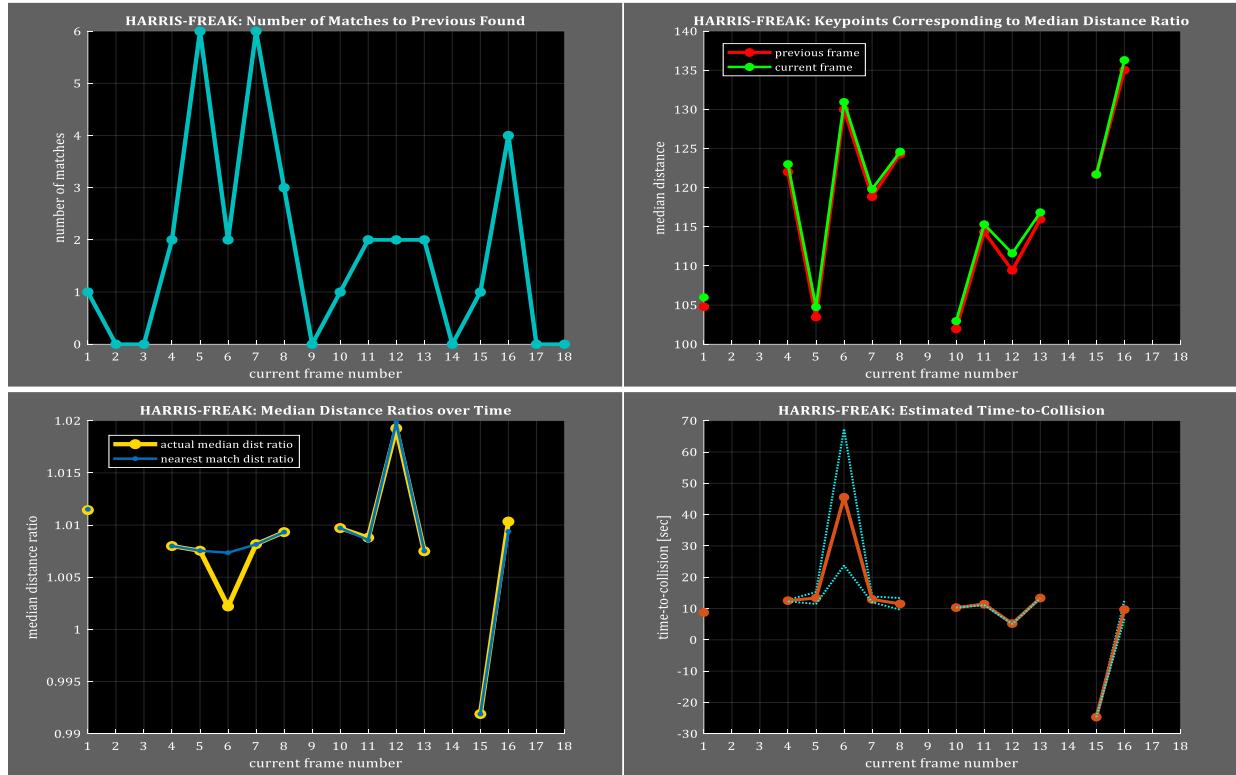


Figure 40: Median Distance Ratio & Time-to-Collision across all frames for HARRIS Detectors & FREAK Descriptors

Example 5: ORB Detector, ORB Descriptor



Figure 41: Relative Distances Between ORB Detectors and ORB Descriptors Between Frames 6 and 7

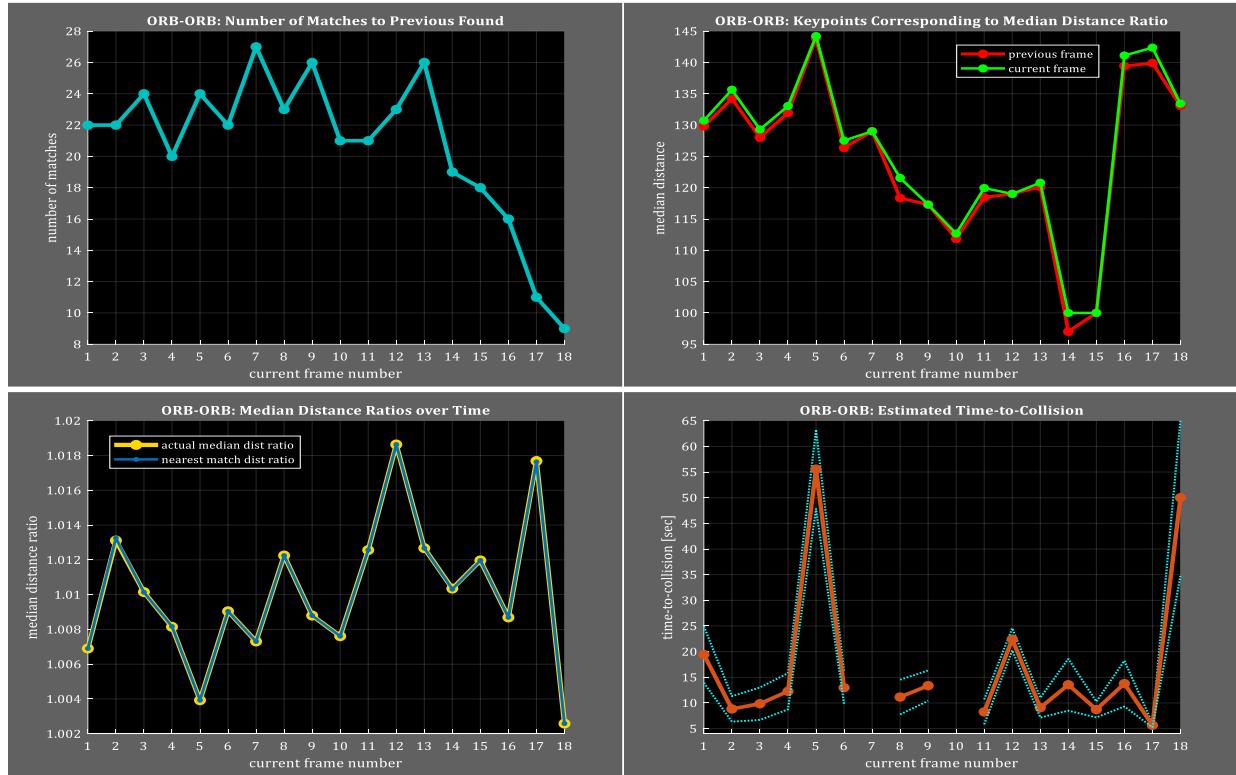


Figure 42: Median Distance Ratio and Time-to-Collision across all frames for ORB Detectors & ORB Descriptors

Example 6: SIFT Detector, SIFT Descriptor



Figure 43: Relative Distances Between SIFT Detectors and SIFT Descriptors Between Frames 14 and 15

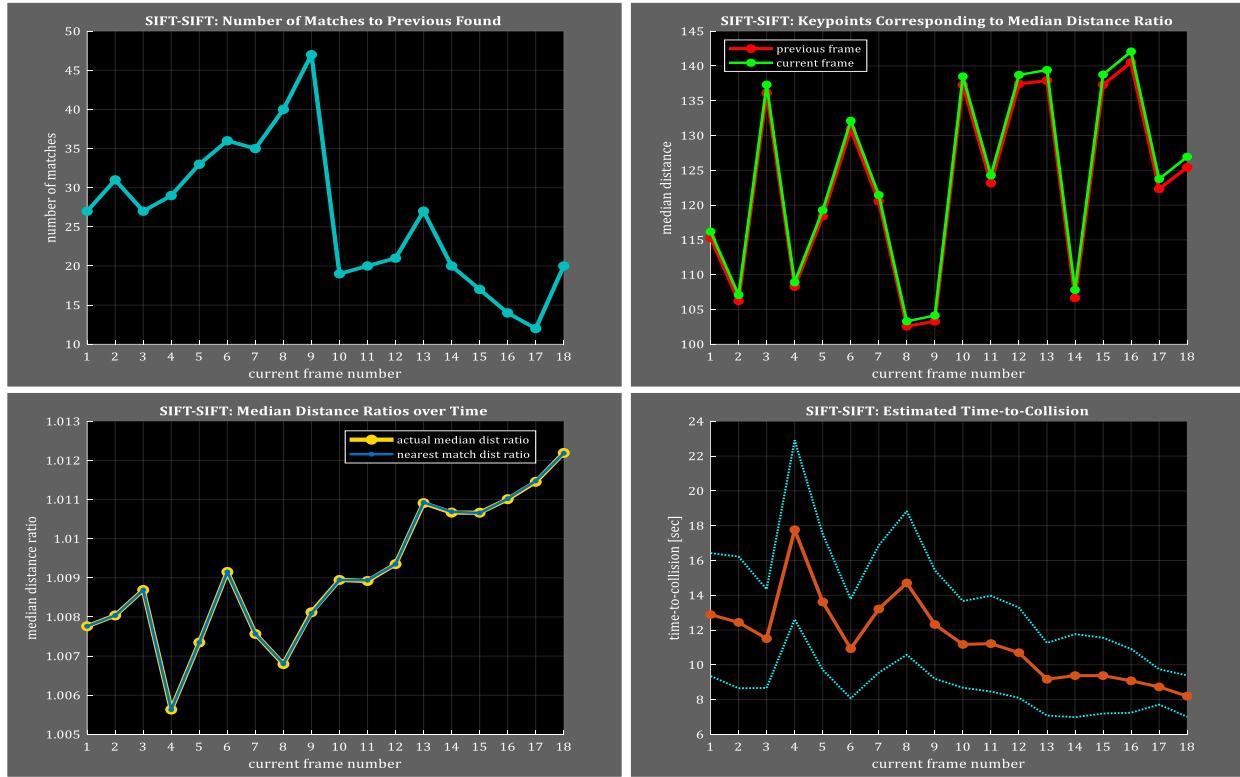


Figure 44: Median Distance Ratio and Time-to-Collision across all frames for SIFT Detectors & SIFT Descriptors

Example 7: Shi-Tomasi Detector, BRIEF Descriptor



Figure 45: Relative Distances Between Shi-Tomasi Detectors and BRIEF Descriptors Between Frames 11 and 12

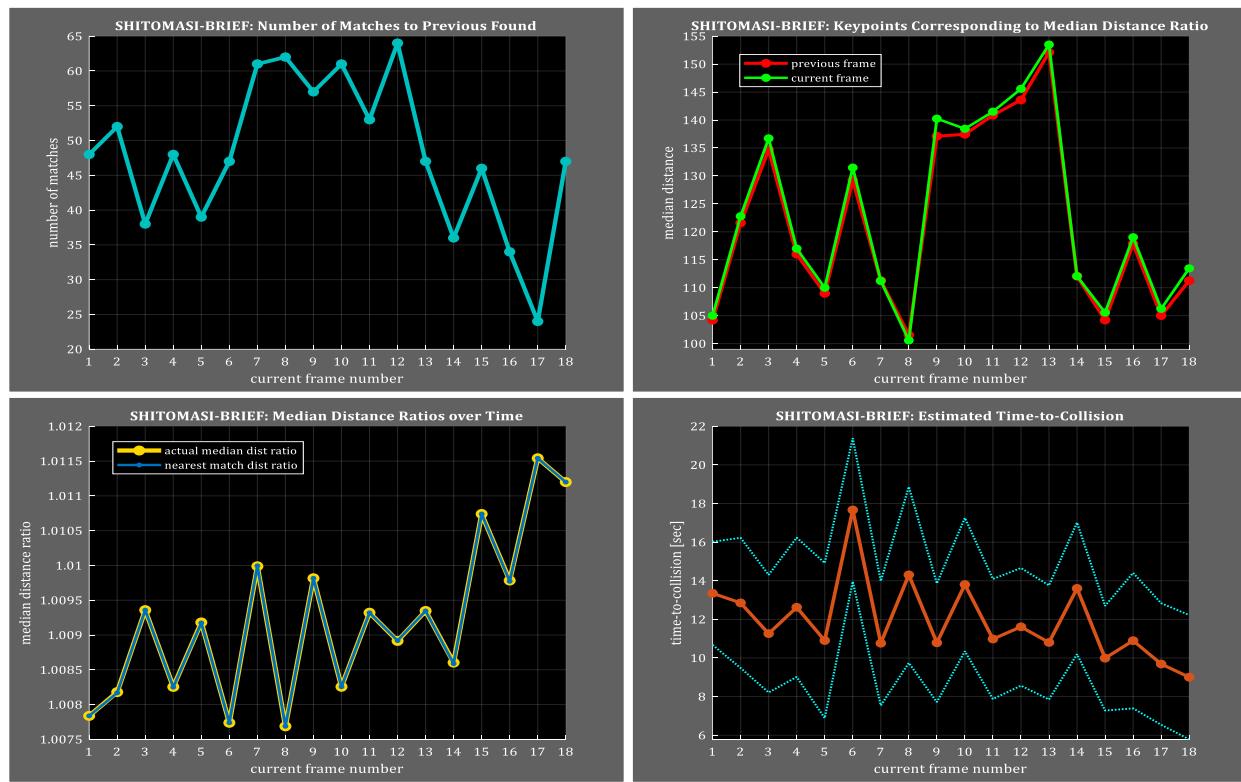


Figure 46: Median Dist Ratio & Time-to-Collision across all frames for Shi-Tomasi Detectors & BRIEF Descriptors

Function Definition for Estimating the Camera-Based TTC in camFusion_Student.cpp:

```

// Compute time-to-collision (TTC) based on keypoint correspondences in successive images
void computeTTCCamera(std::vector<cv::KeyPoint> &kptsPrev, std::vector<cv::KeyPoint> &kptsCurr,
                      std::vector<cv::DMatch> kptMatches, double frameRate, double &TTC, cv::Mat *visImg)
{
    // compute distance ratios for all matched keypoints between the current and previous frame
    vector<double> distRatios;

    for (auto it1 = kptMatches.begin(); it1 != kptMatches.end() - 1; ++it1)
    {
        // get current keypoint and its matched partner in the previous frame
        cv::KeyPoint kpOuterCurr = kptsCurr.at(it1->trainIdx);
        cv::KeyPoint kpOuterPrev = kptsPrev.at(it1->queryIdx);

        for (auto it2 = kptMatches.begin() + 1; it2 != kptMatches.end(); ++it2)
        {
            double minDist = 100.0; // min required distance

            // get next keypoint and its matched partner in the previous frame
            cv::KeyPoint kpInnerCurr = kptsCurr.at(it2->trainIdx);
            cv::KeyPoint kpInnerPrev = kptsPrev.at(it2->queryIdx);

            // compute distances and distance ratios
            double distCurr = cv::norm(kpOuterCurr.pt - kpInnerCurr.pt);
            double distPrev = cv::norm(kpOuterPrev.pt - kpInnerPrev.pt);

            // prevent divide by zero
            if (distPrev > std::numeric_limits<double>::epsilon() && distCurr >= minDist)
            {
                double distRatio = distCurr / distPrev;
                distRatios.push_back(distRatio);
            }
        } // end inner loop over all matched keypoints
    } // end outer loop over all matched keypoints

    // only continue if list of distance ratios is not empty
    if (distRatios.size() == 0)
    {
        TTC = NAN;
        return;
    }

    // compute median distance ratio, removing outlier influence
    std::sort(distRatios.begin(), distRatios.end());

    long medianIdx = floor(distRatios.size() / 2.0);
    double medianDistRatio;
    if (distRatios.size() % 2 != 0) // if the number of distance ratio elements is odd
    {
        // then take the middle value
        medianDistRatio = distRatios[medianIdx];
    }
    else // if the number of distance ratio elements is even
    {
        // then take the average of the two middle values
        medianDistRatio = (distRatios[medianIdx - 1] + distRatios[medianIdx]) / 2;
    }

    double dT = 1 / frameRate;
    TTC = -dT / (1 - medianDistRatio);
}

```

Function Call in main():

```

computeTTCCamera((dataBuffer.end() - 2)->keypoints,
                  (dataBuffer.end() - 1)->keypoints,
                  currBB->kptMatches,
                  sensorFrameRate,
                  ttcCamera,
                  imgIndex,
                  kptTTCfilename);

```

FP.5 PERFORMANCE EVALUATION, PART I

i Find examples where the TTC estimate of the Lidar sensor does not seem plausible. Describe your observations and provide a sound argument as to why you think this happened.

The figure below shows the lidar TTC and relative velocity estimates. The mint green line is the result from the code implementation. The dark orange line is the TTC based my “eyeball” estimate of the distances to the preceding vehicle. The assertion of where and when the TTC is off is based upon comparisons between the implementation results and manual “eyeball” estimates of the distance to the rear of the preceding vehicle from a top-down view of the lidar points.

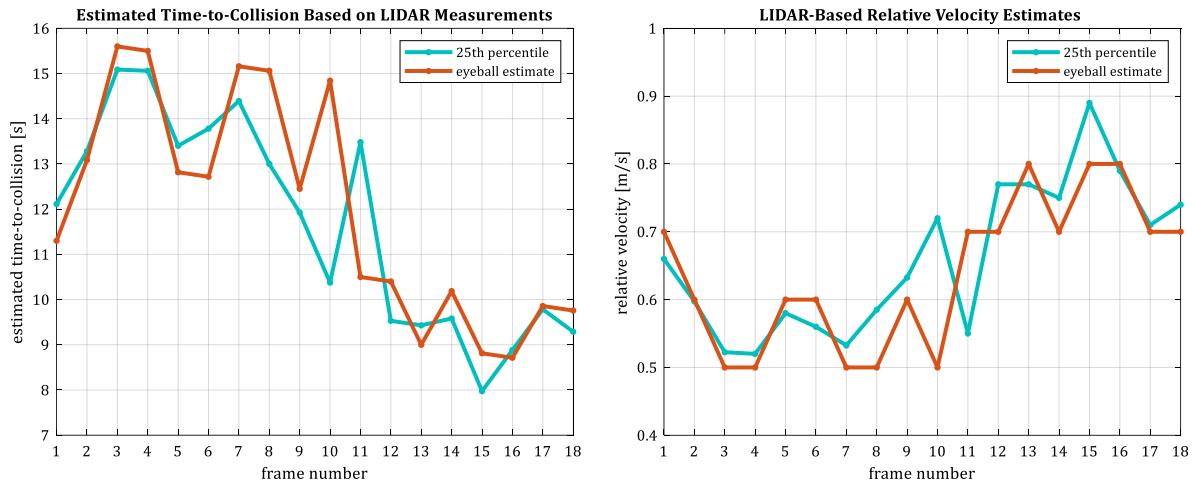
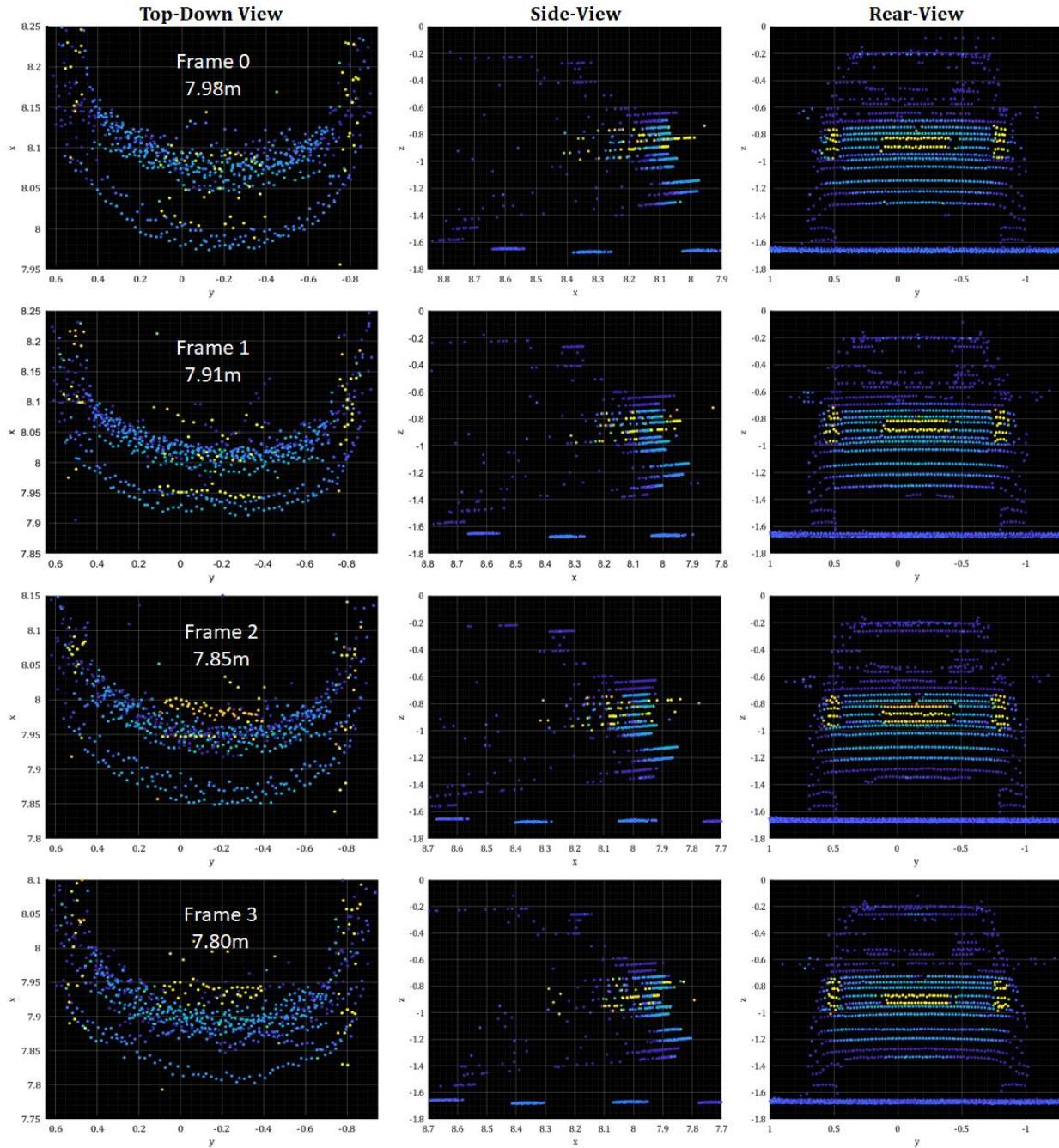
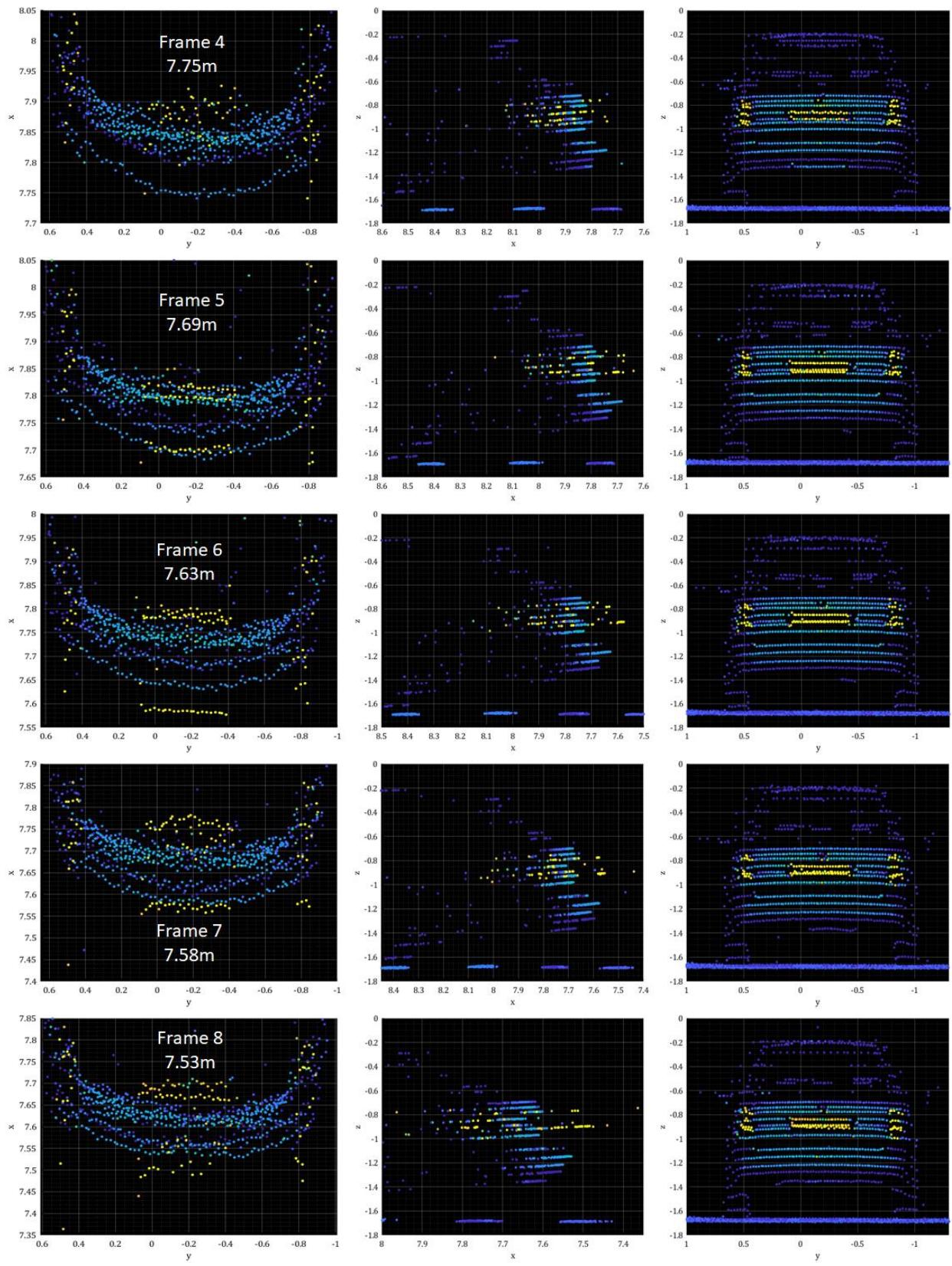


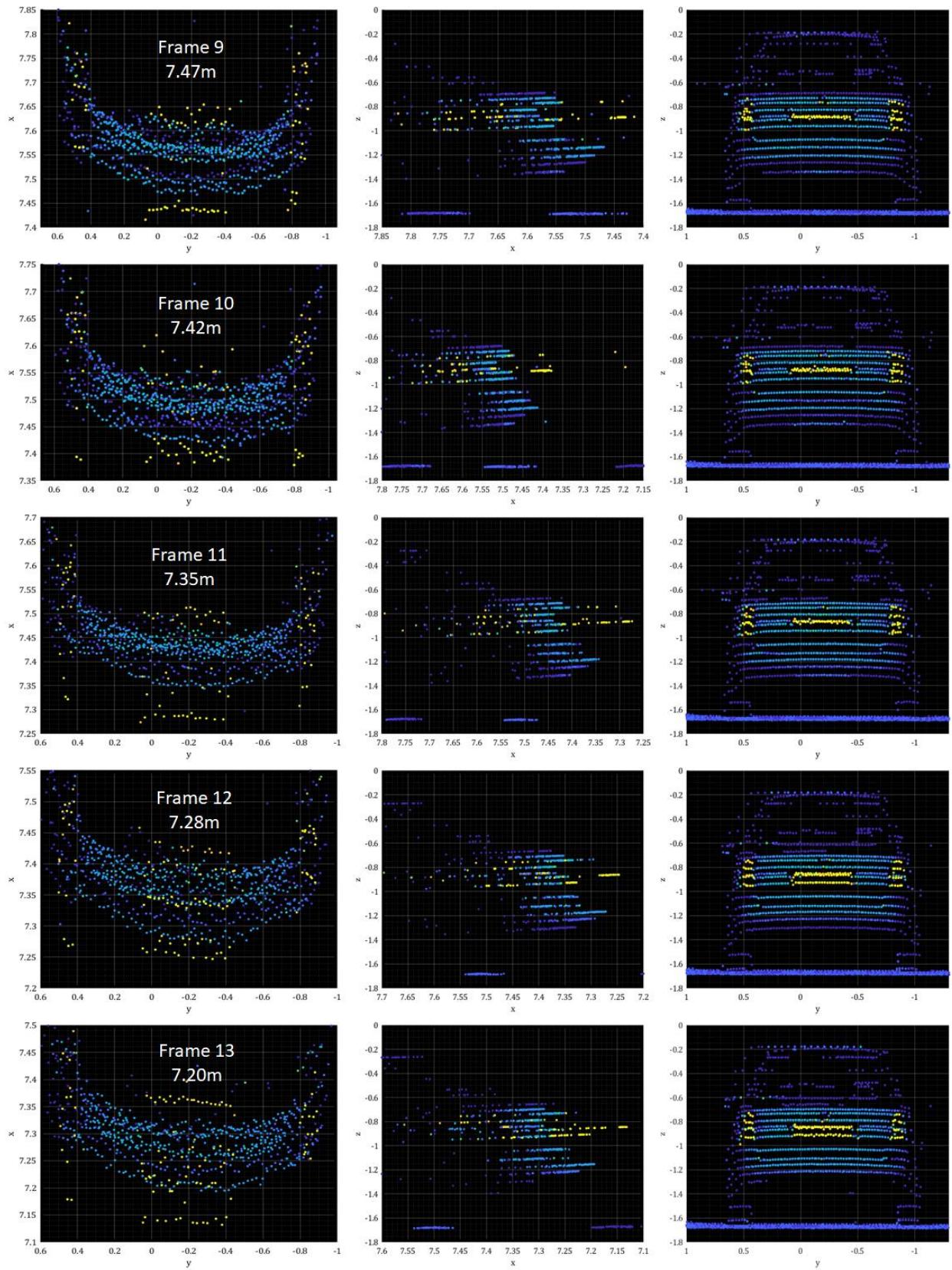
Figure 47: Distribution of LIDAR Time-to-Collision Estimates Based on x-Position Data on the Preceding Car

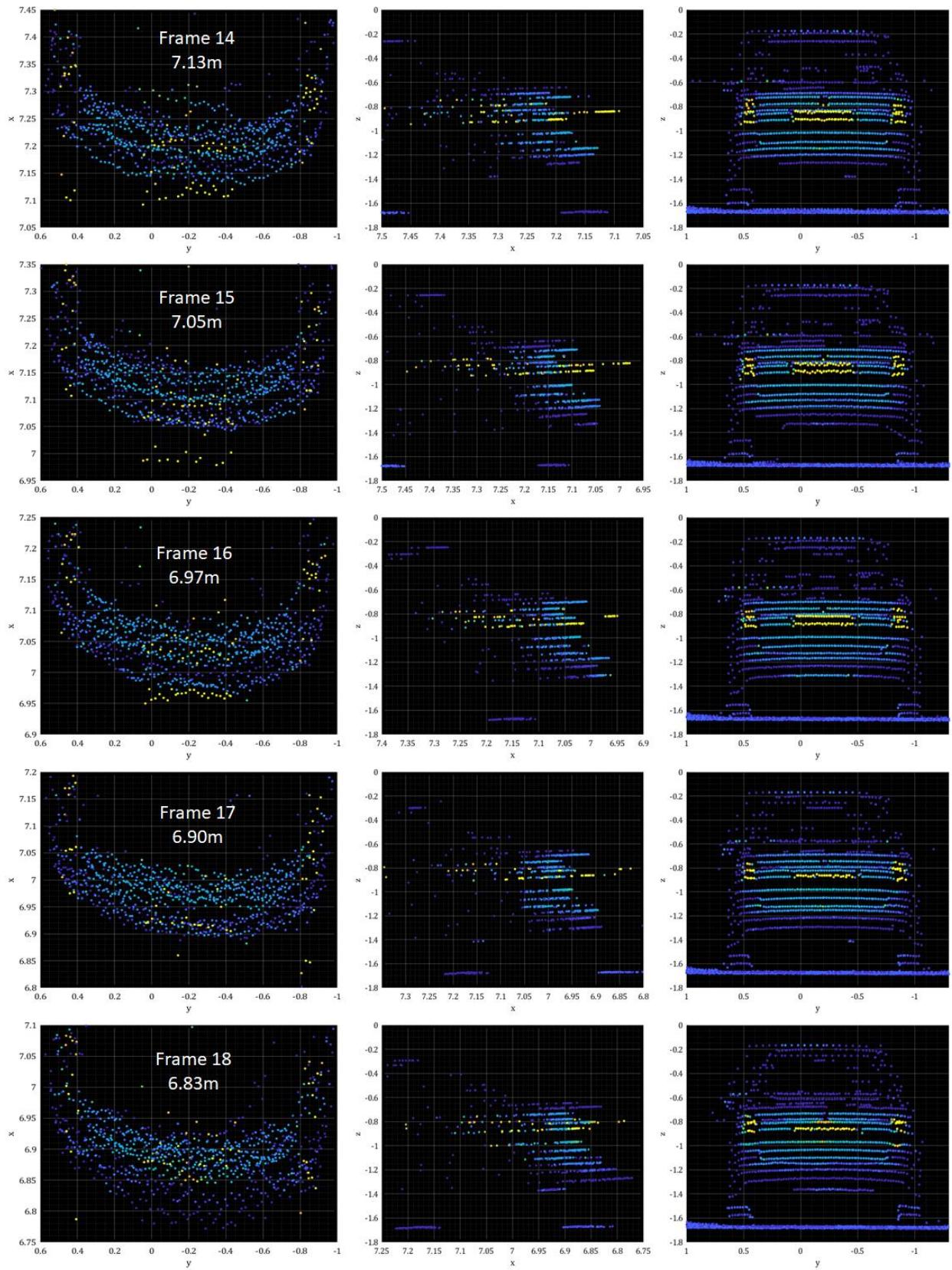
Selection of the lower 25th percentile distance of all lidar points within the region of interest produced relatively stable and plausible estimates of TTC over time. However, the 25th percentile is not the usually the closest point, thus producing estimates of TTC that are biased high relative to truth. This section will explore the errors in the lower percentiles.

I performed the “eyeball” estimates by looking at the top-down views, using some “human judgement” and experience to discern where the closest part of the vehicle might be. This approach was taken due to the absence of a truth estimate to help assess the plausibility of the algorithm results. The x-distance estimates that I selected are shown in the top-down scatter plot for each frame shown below.







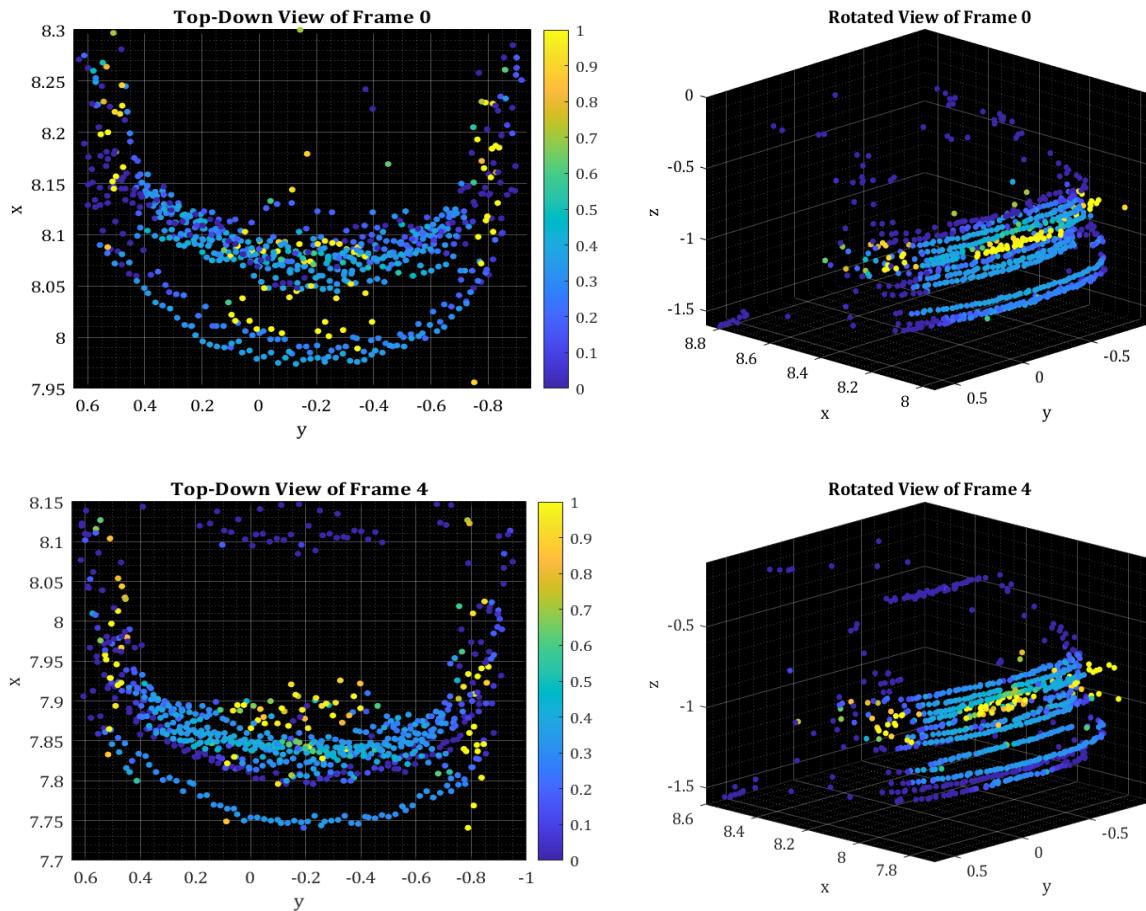


The top-down (y-x) and side (x-z) views of the preceding vehicle show that the lidar estimates in the x-direction tends be very noisy. By contrast, lidar estimates from the rear (y-z) view of the vehicle reveals a high level of detail, precision and accuracy. I used the rear (y-z) views as a color-coded reference to identify the corresponding point position estimates in the x-direction.

Lidar measurement accuracy is correlated to the amount of light reflected from an object. The lecture notes suggested removing measurements with low reflectivity, explaining that they would be less reliable than points with higher reflectivity. However, I found the opposite to be true ...

The scatter plots showed distinctively higher reflectivity estimates coming off the vehicle license plate and tail-lights. These higher reflectivity points demonstrated some interesting behaviors – sometimes they were noisier than less reflective points (especially on the tail-lights), other times they organized into distinct scan lines (i.e. on the license plate) with conflicting information as to where these surfaces truly existed relative to the rest of the vehicle in the x-direction.

The following images show the lidar measurements on the preceding car (selected bounding box) across different perspectives: (1) a top-down view in world coordinates, and (2) 4D rotated view in world coordinates. In addition to position data, the reflectivity estimates of the measurements have been color-coded according to the scale shown, with 0 being the lowest reflectivity and 1 being the highest reflectivity observed across all lidar measurements at that snapshot in time.



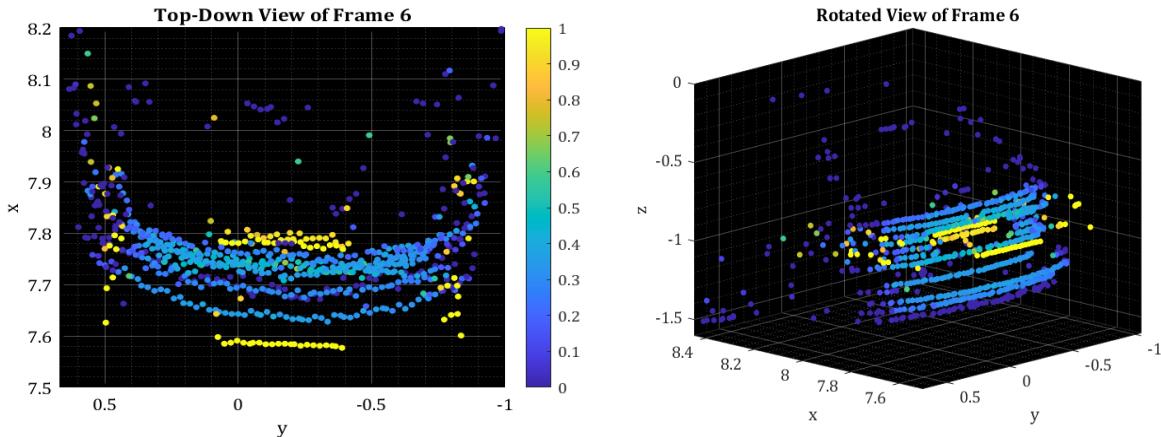


Figure 48: Comparison of Lidar Point Cloud Distributions on High Reflectivity Features

For the tail-lights, the distance measurements in the x-direction are particularly high. This could be due a combination of the surface materials and the shape curvature at the edges of the vehicle, which would affect the directions of scattered light and the intensity of the return signal. Lower reflectivity points along the edge do not show the same magnitude of error, and the noise on the similarly reflective license plate are high on some (but not all) frames.

Relative motion would cause the aspect angle of the lidar sensor to change across frames. It also changes at a much smaller scale within a frame as a function of the elevation of each individual lidar scan. The slight variations in aspect angle and corresponding light scatters at different elevations could explain why some parts of the license plate appear at 7.6m and other parts at 7.8m in Frame 6 shown above. The azimuth off boresight could also explain why the noise is higher on the tail-lights.

The white license plate contains black painted letters and numbers that may contribute to small angle-dependent variations across the surface of the plate. The lidar sensor may also be picking up variations in ambient light, shadows and reflections from passing vehicles. Car exhaust particles and other environmental factors may also contribute to measurement noise as the vehicles gather at the stop light.

Measurement noise in the x-direction is especially problematic for a TTC computation that is highly sensitive to small errors in absolute and relative distance across frames. Recall the equation we used: $TTC = (d_1 * \Delta t) / (d_0 - d_1)$. Smaller denominators produce TTC estimates that approach positive or negative infinity very quickly. Even one centimeter of error in the distance measurement on any frame in our sequence could produce dramatic errors in the estimated TTC.

The TTC equation also makes a simplifying assumption of constant velocity for small Δt . Δt in our case is 1/10th of a second. That sampling rate would be “small” if we were cruising on an uncrowded highway at a higher speed. However, the vehicles in our scenes are approaching a stop light, where the relative velocity between the ego car and the preceding car is not constant -- we can even observe the occasional red flicker of the brake lights as the driver in the preceding vehicle places intermittent pressure on the brake. The figure below shows the relative velocity based on percentiles across the distribution of x-distances in the region of interest. We can see that the magnitude of the fluctuations decrease as the percentiles increase. However, it is still unclear how much of the fluctuations across frames are due to noise versus actual changes in velocity.

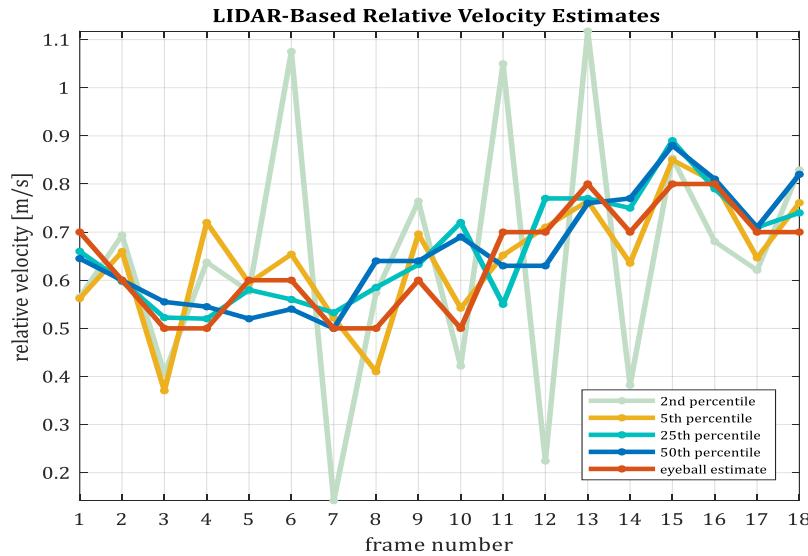


Figure 49: Change in x-distance measurements as a function of time between adjacent frames

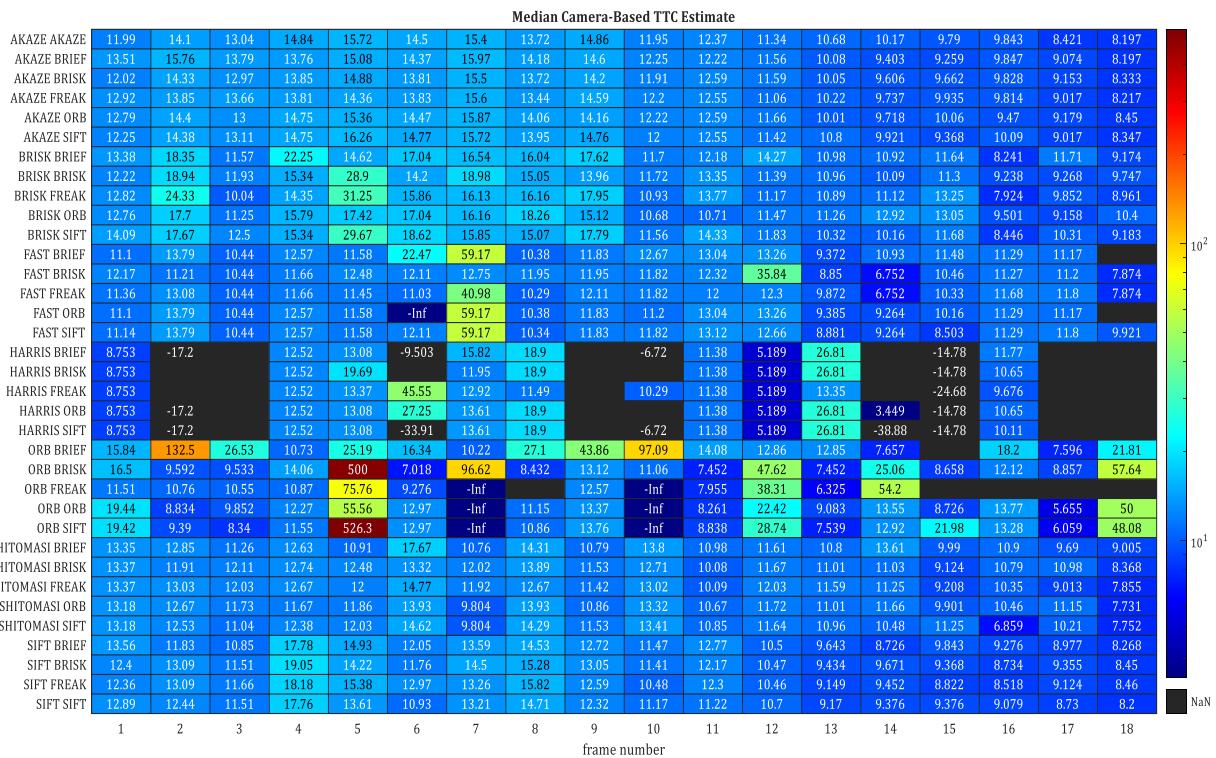
In summary, lidar x-distance measurements can be noisy. Errors in the TTC estimates based on the constant velocity model are expected to be highly-sensitive to measurement noise and fluctuations in velocity that occur below the sampling rate. Furthermore, the constant-velocity assumption does not generalize well to dynamic traffic situations where vehicles are braking hard, or where we don't have a centered reference point to compare to over time.

For real-life applications, it is recommended that we use the constant-acceleration model for computing the lidar-based time-to-collision.

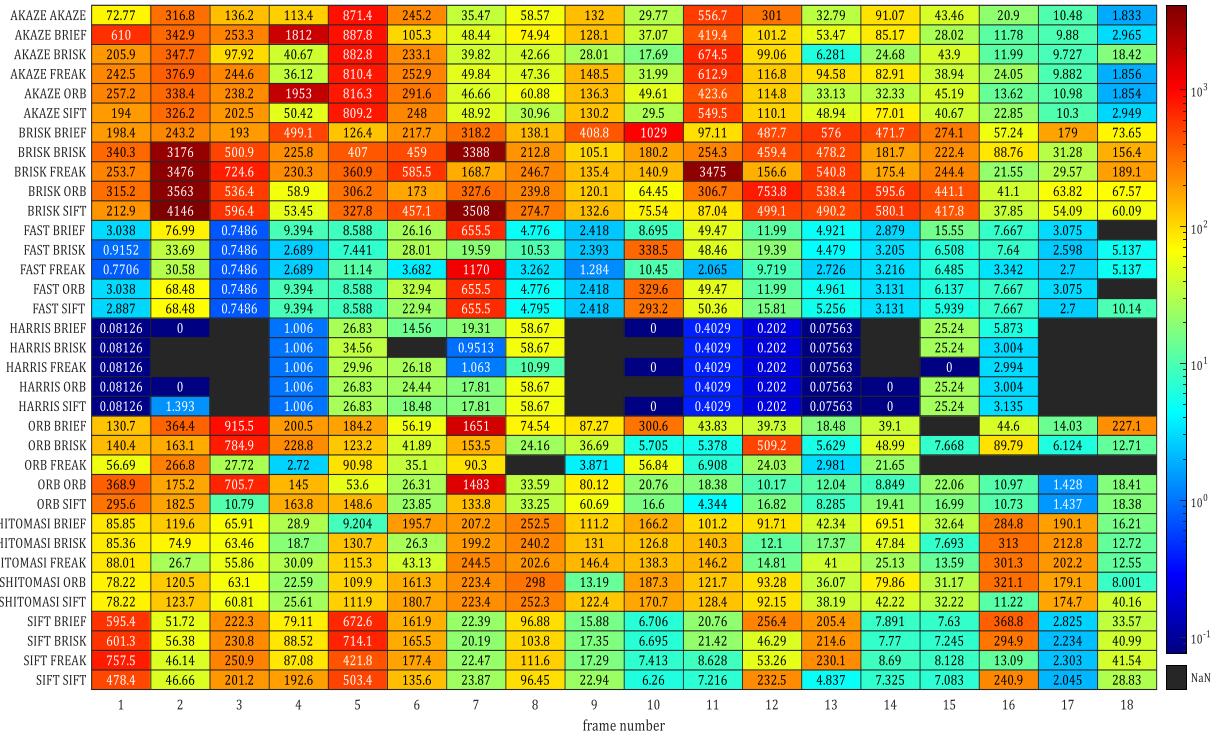
FP.6 PERFORMANCE EVALUATION, PART II

i Run several detector-descriptor combinations and evaluate the differences in TTC estimation. Identify which methods perform the best. Find examples where camera based TTC estimation is way off, describe observations and investigate potential reasons.

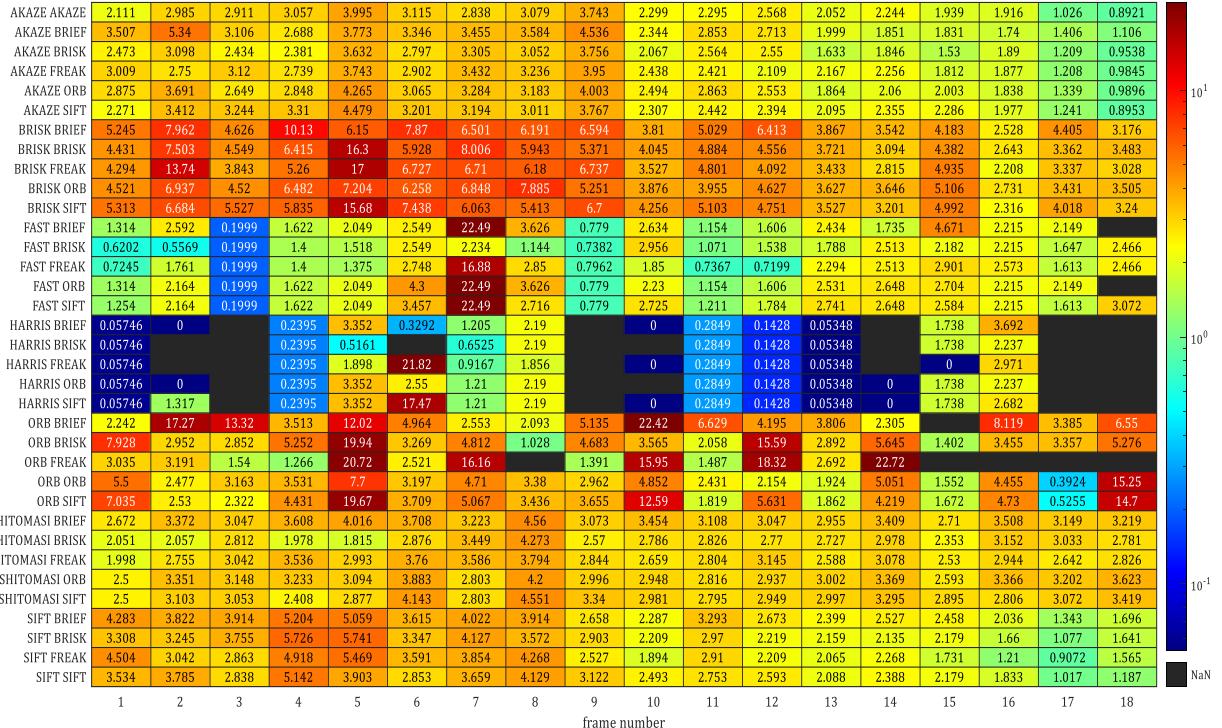
Several detector-descriptor combinations were compared with respect to the TTC estimate on a frame-by-frame basis. The results are shown in the tables below:

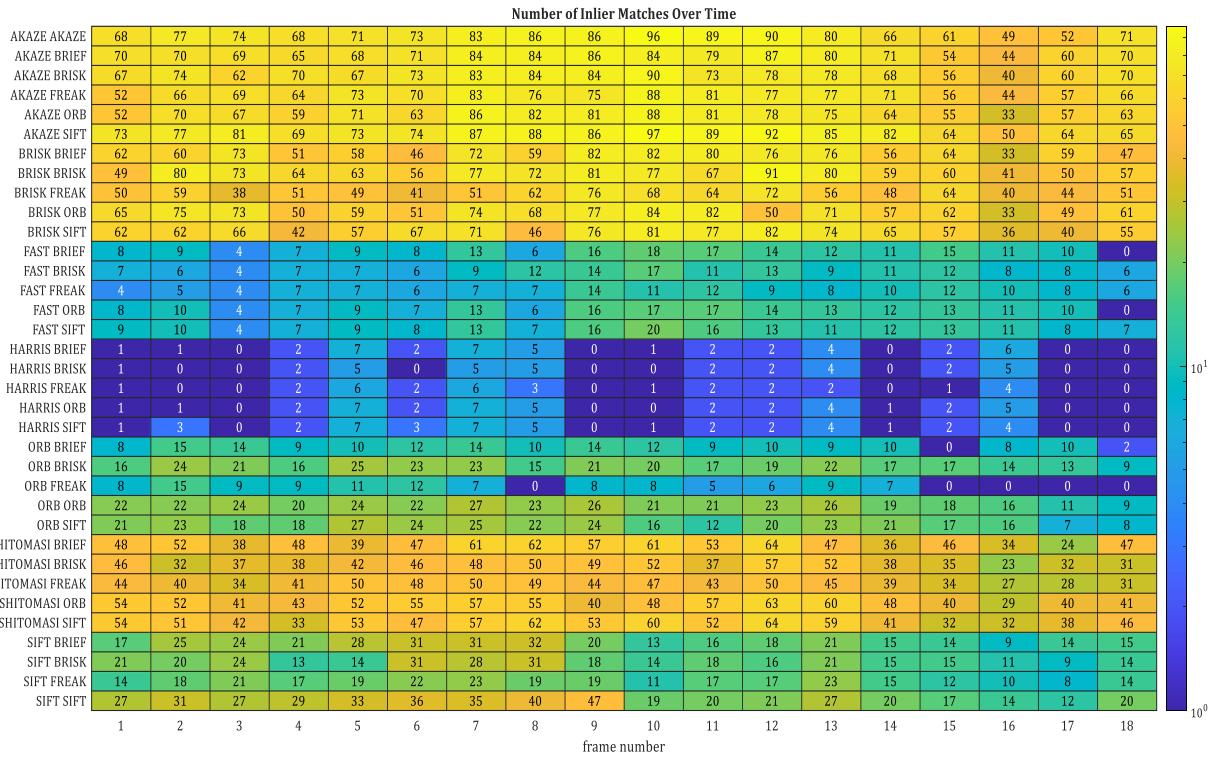


Standard Deviation Across Camera-Based TTC Estimates



Median Absolute Deviation Across Camera-Based TTC Estimates



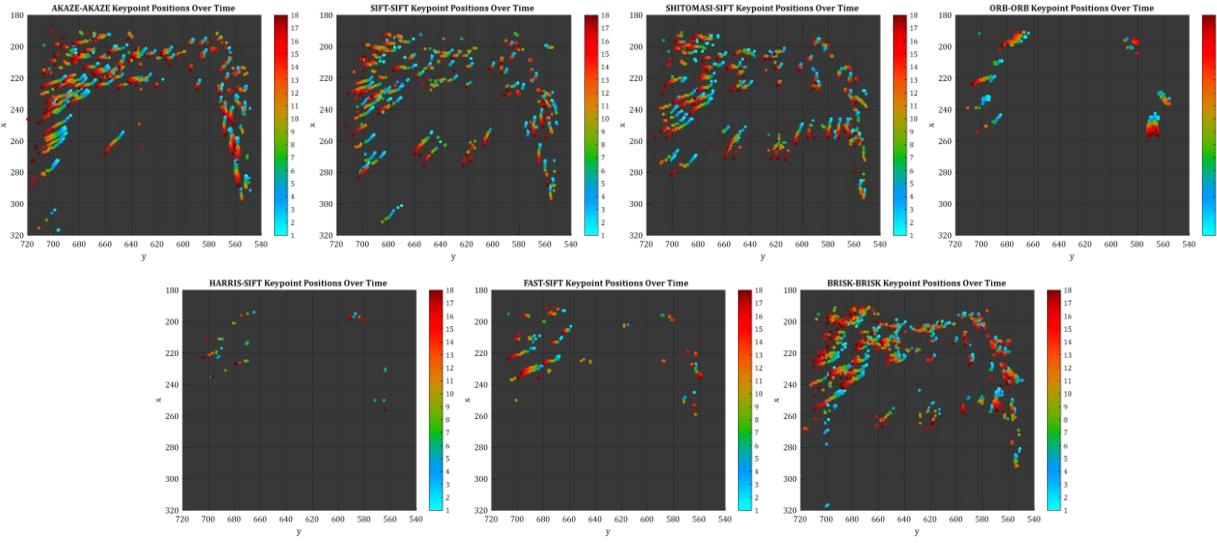


Based on the tables shown above, it appears that TTC estimation may be driven first by the detector type.

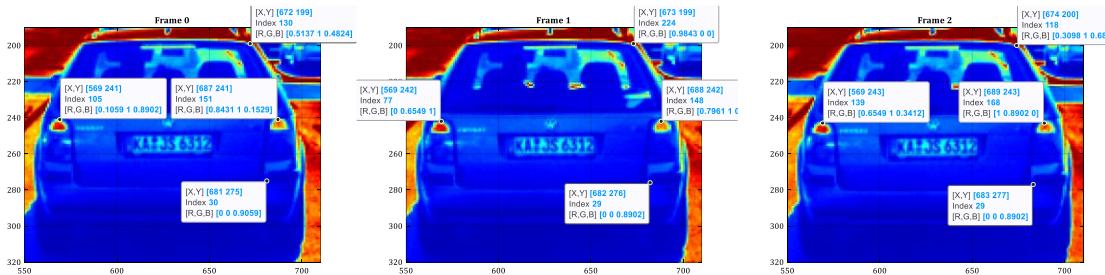
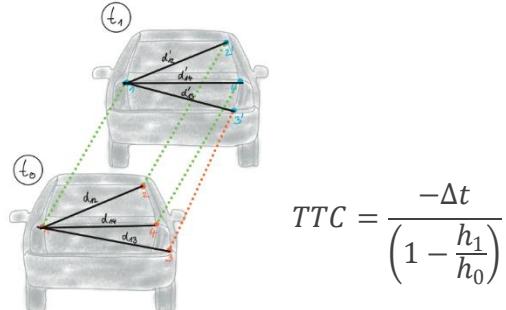


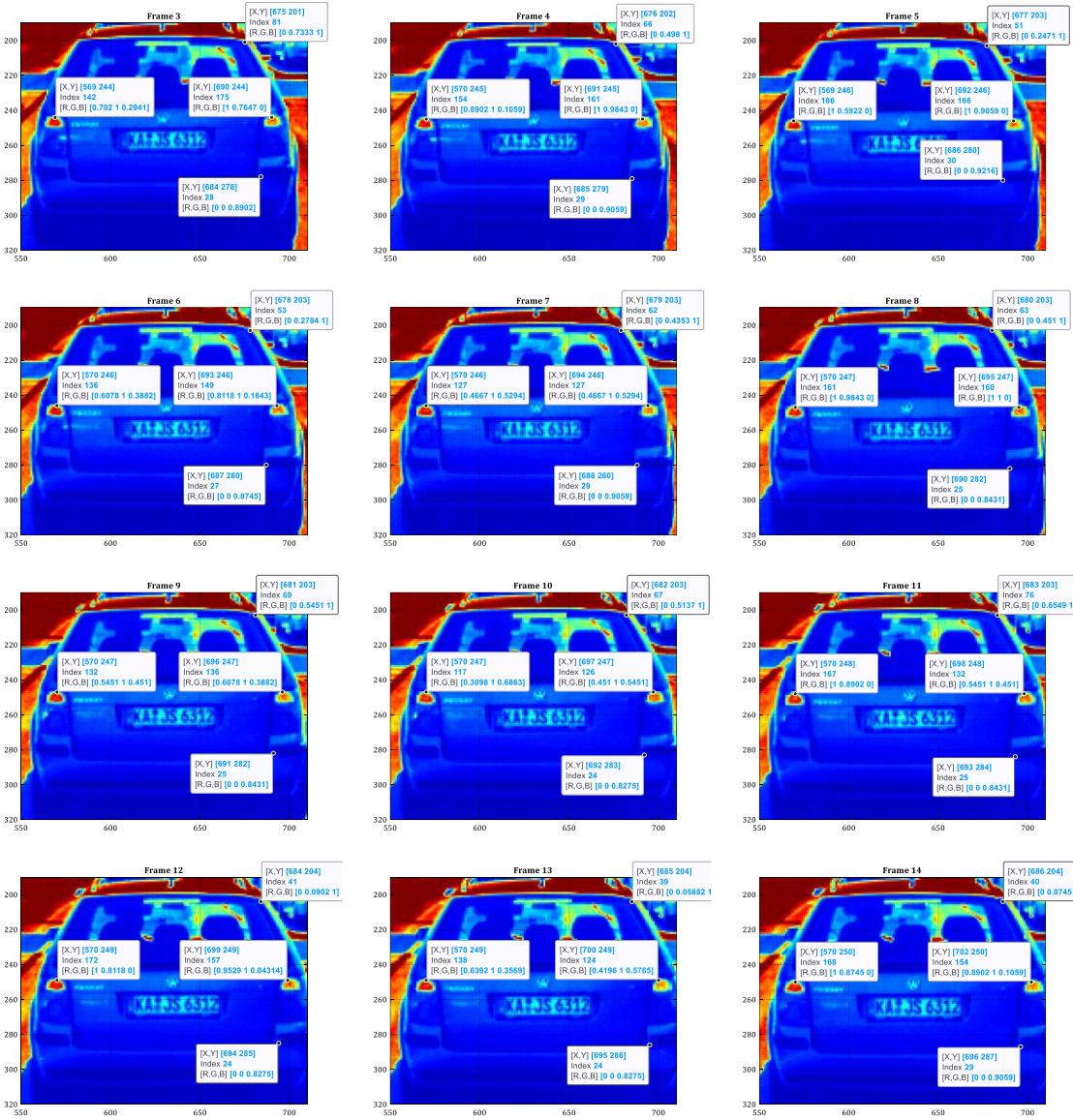
Figure 50: Keypoint Detector Types Evaluated in this Study

Seeking insights into what causes one detector method to perform better than others, I looked at the optical flow of a few combinations. In the following images, the keypoints in the bounding box are plotted over time, where the color gradient varies from Frame 0 (cyan) through Frame 18 (red):



We lack a source of ground truth to evaluate performance, so I performed “eyeball” estimates of keypoints on all frames.





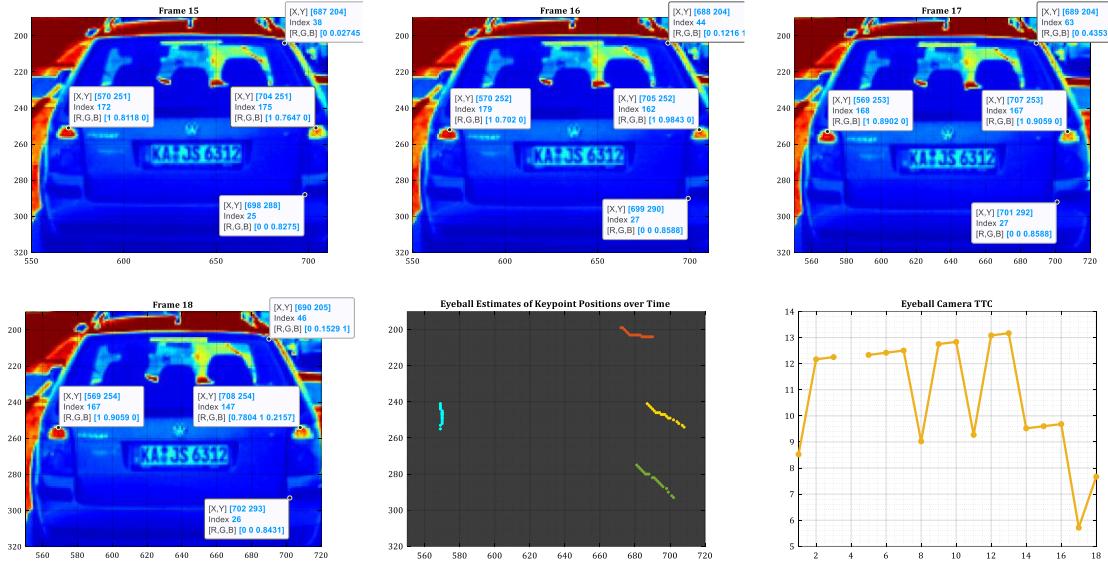


Figure 51: Manual “Eyeball” Estimate of the Camera Keypoint-Based TTC

I also looked for multi-sensor (camera + lidar) consensus / similarities as a surrogate for truth. Given these results, I determined that the group of AKAZE detectors produced temporal TTC trends most similar to the lidar-based TTC estimates.

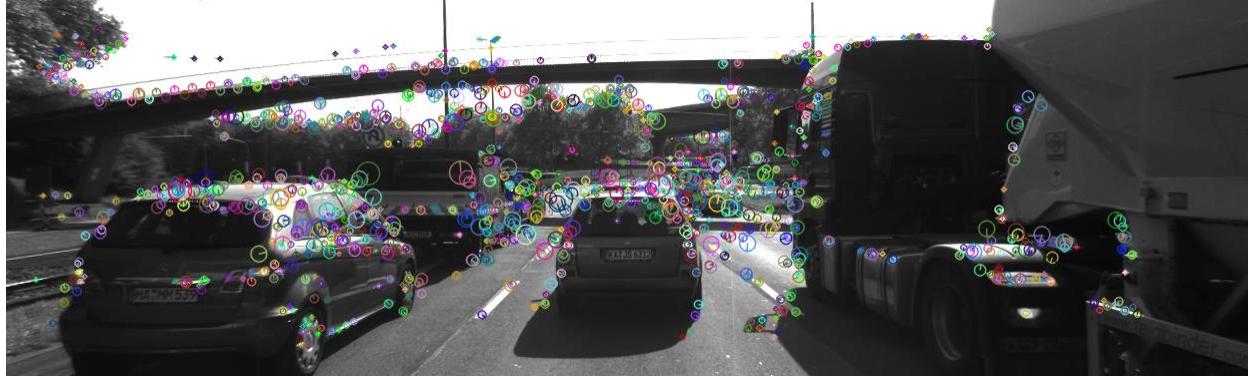


Figure 52: Example of AKAZE Keypoints

In the graph below, I grouped the AKAZE keypoint results by taking the median TTC across all corresponding descriptor methods, at every frame in time.

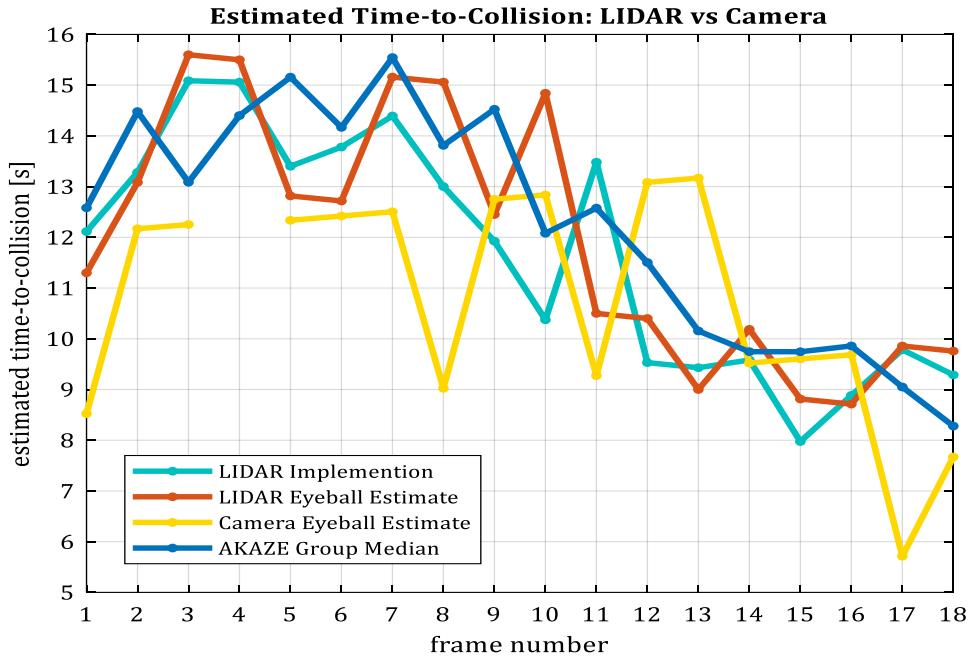
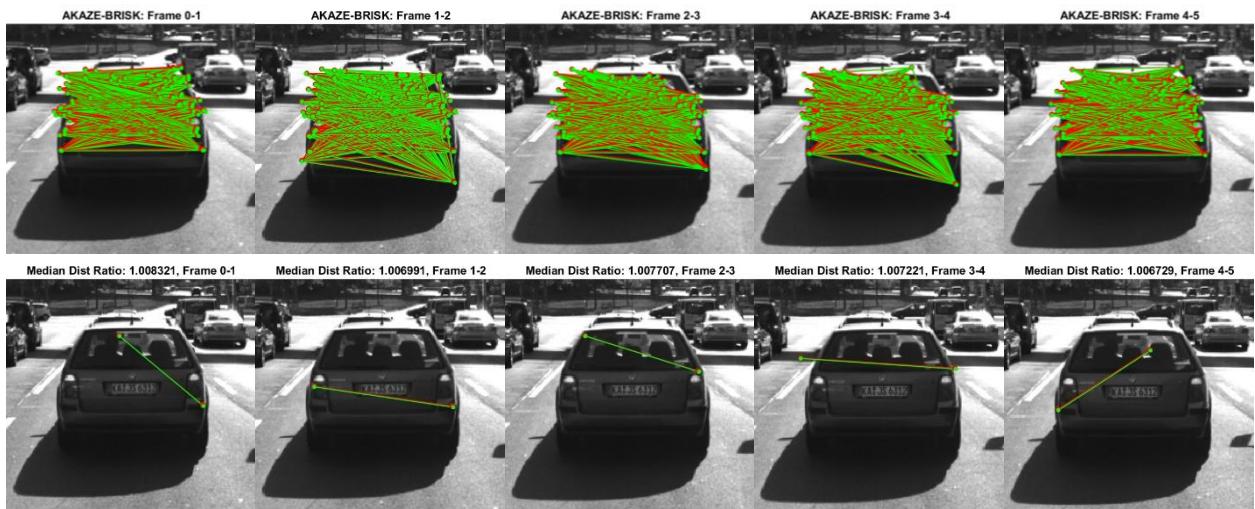


Figure 53: Estimated Time-to-Collision: LIDAR vs Camera

AKAZE-BRISK

I chose AKAZE detectors with BRISK descriptors as a configuration to look at because it produced trends with the same order of magnitude as the LIDAR TTC results. The following images show the keypoints used to estimate TTC on each frame (row 1), followed by the keypoint pairs corresponding to the median distance ratio (row 2). Points and lines in red correspond to the previous frame, while points and lines in green correspond to the current frame:



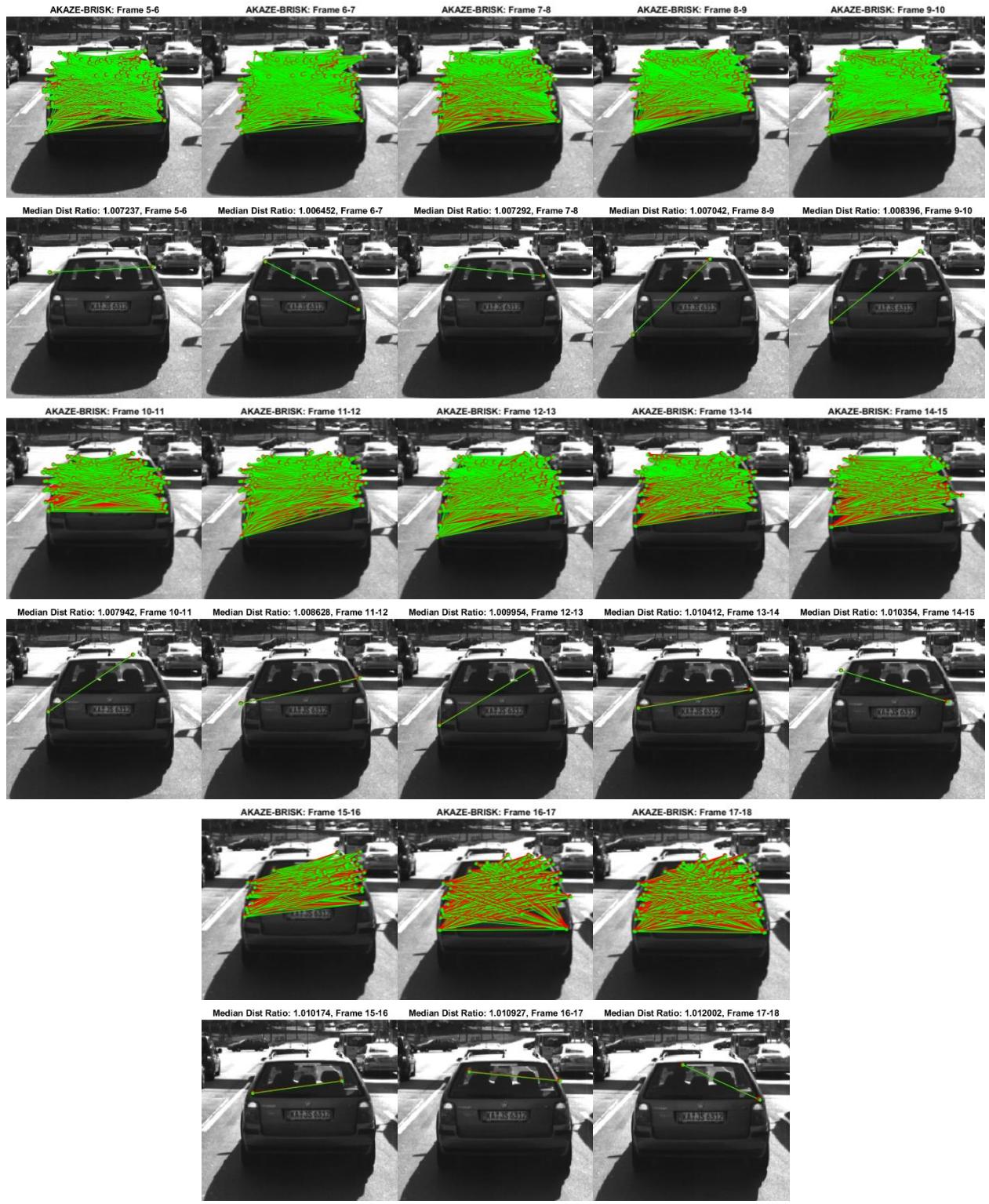


Figure 54: Keypoints Corresponding to the Median Distance Ratio for TTC Estimation

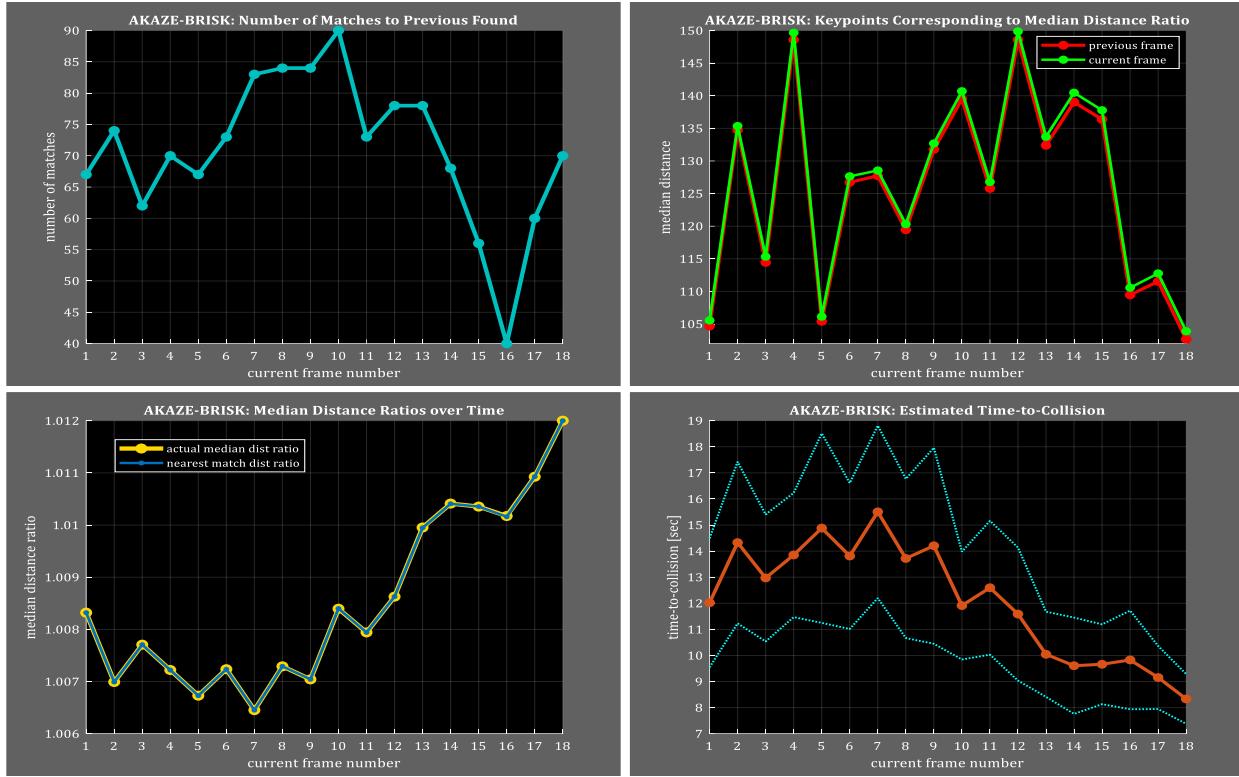


Figure 55: Median Distance Ratio and Estimated TTC based on AKAZE Detector & BRISK Descriptor

On the other end of the spectrum, we looked into why some approaches produced TTC estimates that were way off.

HARRIS-SIFT

The HARRIS detection method produced a very sparse set samples from which we could compute the TTC. Recall the multiple entry criteria that were required for samples to qualify for TTC estimation:

First, Harris Corners were detected in the monochrome image:



Figure 56: Harris Corner Keypoints Detected on Frame 7

Second, when keypoints were found, they had to be matched with keypoints in adjacent frames:



Figure 57: Keypoint Matches Between Frames 6 and 7

Third, keypoints not associated with a YOLO bounding box were excluded. Also, keypoints that fell within more than one bounding box were removed:

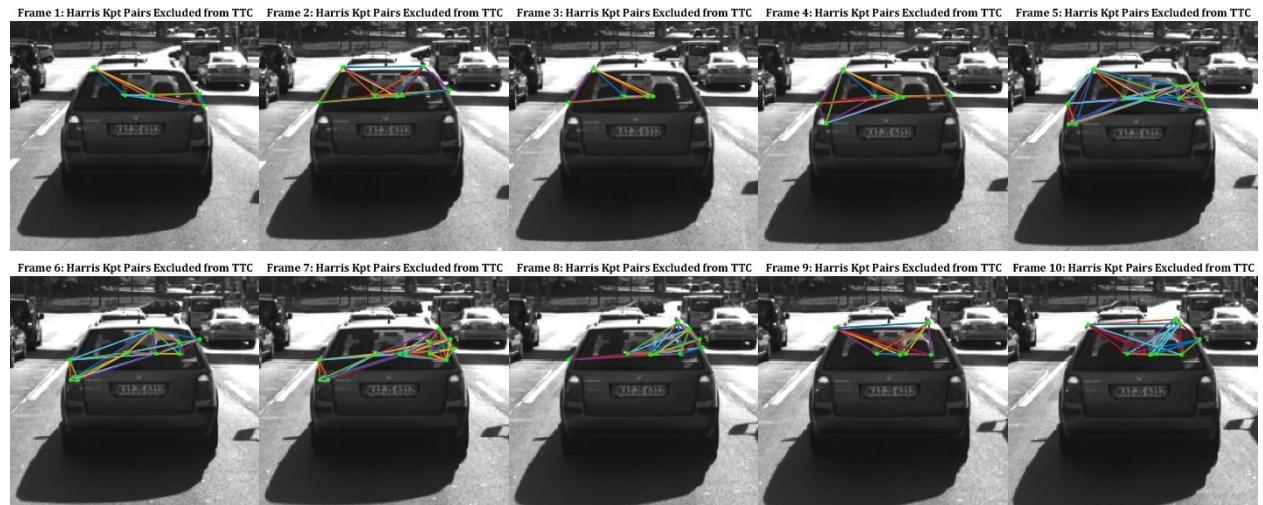


Figure 58: Before (far left) and after (far right) removing keypoints associated with multiple bounding boxes

Fourth, only keypoints within the bounding box of interest* were looked at (*lidar constraints applied earlier usually narrowed it down to one box).

Fifth, intraframe keypoint pairings were required to have a distance greater than or equal to 100cm. The intent was to capture the more distal keypoint pairs stretching across the extent of the vehicle, which would provide better measurements for detecting a change in the projected area than keypoint pairs that were closer together.

The following figures show the HARRIS detection pairs that were excluded from the TTC estimation for having relative distances below 100cm. Note that the filtering criteria has the effect of removing inner keypoints from consideration:



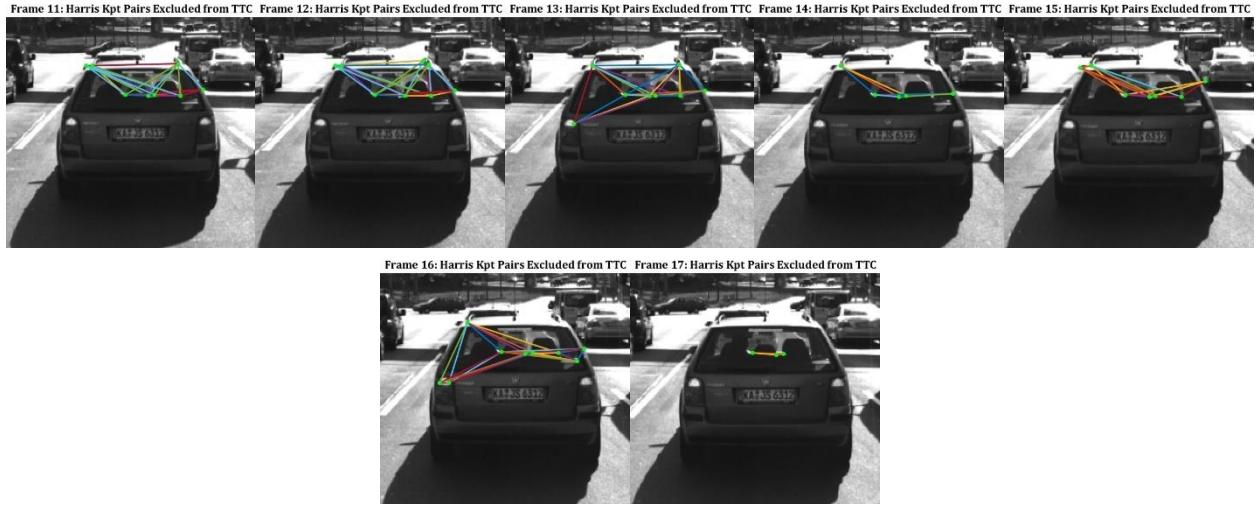


Figure 59: Harris Keypoint Pairs with Distance less than 1m (and therefore excluded from TTC estimation)

Ideally, we would emerge from these 5 stringent entry criteria with a statistically significant number of samples to produce a median distance ratio that captures the frame-to-frame motion.

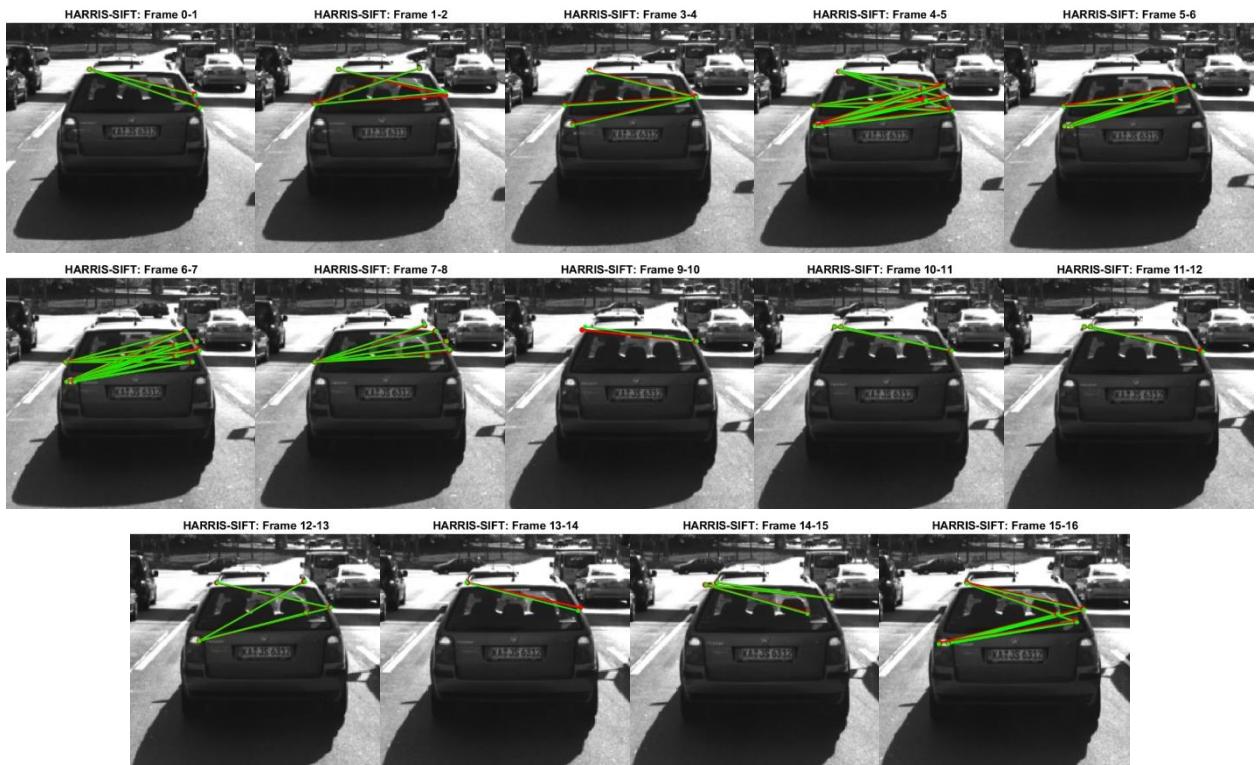


Figure 60: Keypoint pairs with distance greater than 1m: HARRIS Detector - SIFT Descriptor

For all keypoint pairings that manage to pass the stringent list of entry criteria, only the median distance ratio is used to compute the TTC. The more uniform this filtered distribution is, the better the chances of

pulling a good representation. In the case of HARRIS detectors, the small number of surviving samples leaves little room for error in keypoint position accuracy:

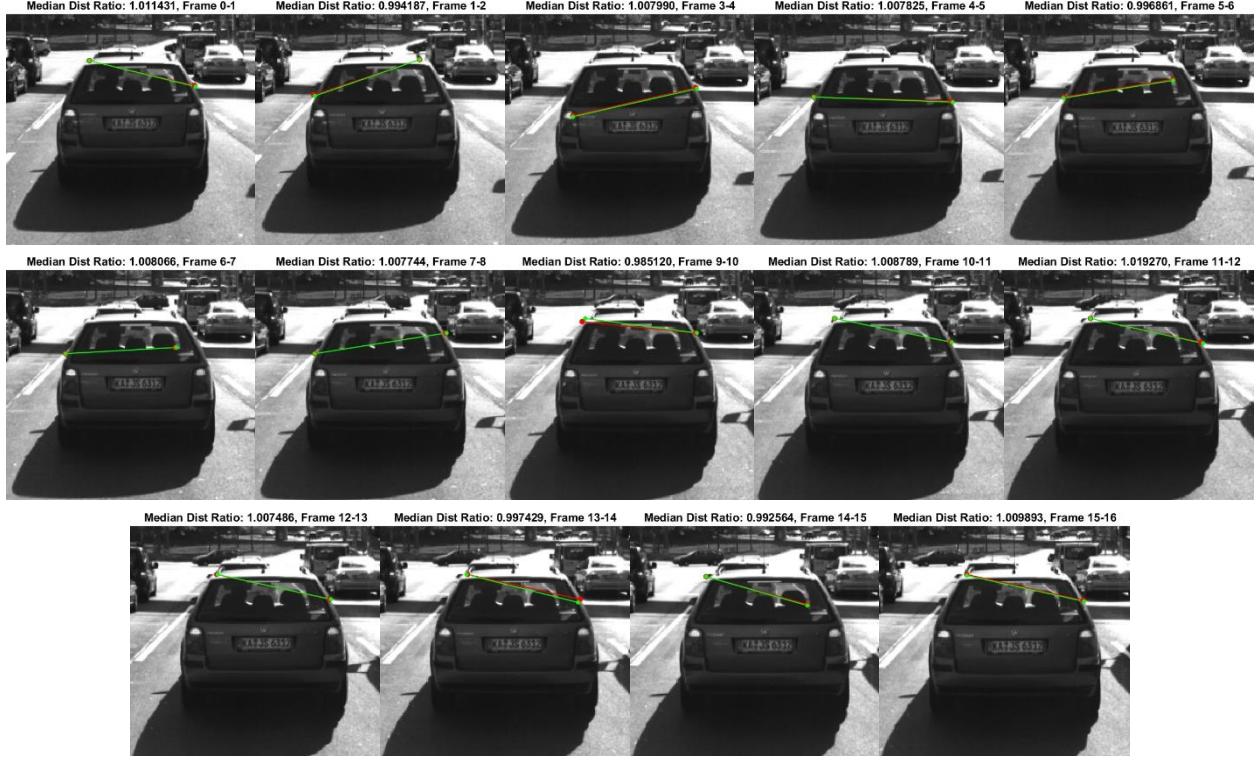


Figure 61: Keypoint Pairs Corresponding to the Median Distance Ratio

If no matching keypoints between a current and previous frame were found within the bounding box, then a TTC could not be computed on the current frame. The figure below shows that we had single-frame TTC dropouts on frames 3 and 9, indicating that no matching keypoints were found in the bounding box between frames 2-3 and 8-9. In addition, both frames 17 and 18 did not produce TTC estimates, which shows that one or both may not have had any keypoints at all.

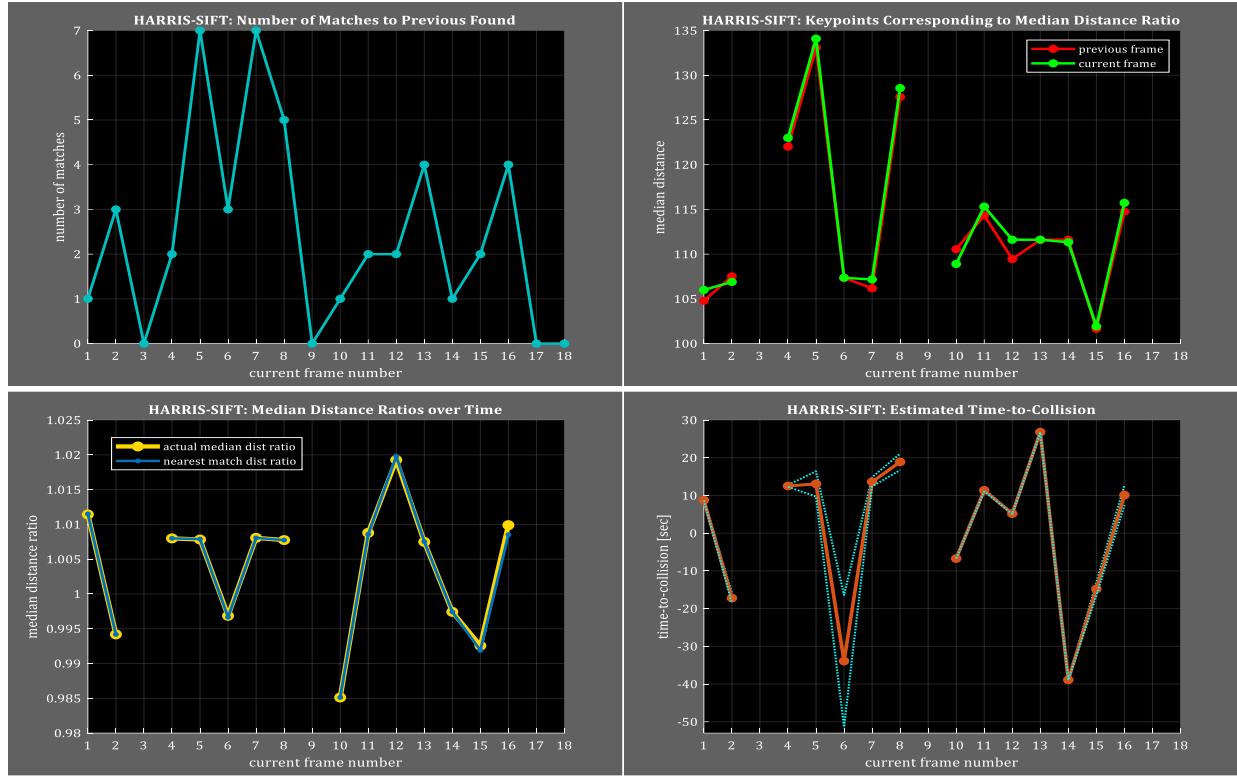


Figure 62: Median Distance Ratio and Estimated TTC based on HARRIS Detector & SIFT Descriptor

Frames 2, 6, 10, 14 and 15 all produced negative TTC estimates, suggesting that the distance between the ego car and preceding vehicle have increased between frames. While that may be the case, these estimates were based on a very small number of sample matches – or 3, 3, 1, 1 and 2 matches, respectively. The keypoints corresponding to the median distance ratio would need to have a very high level of accuracy and consistency in identifying the same point on *both* of the frames being compared. The images below show how small that margin of error is:

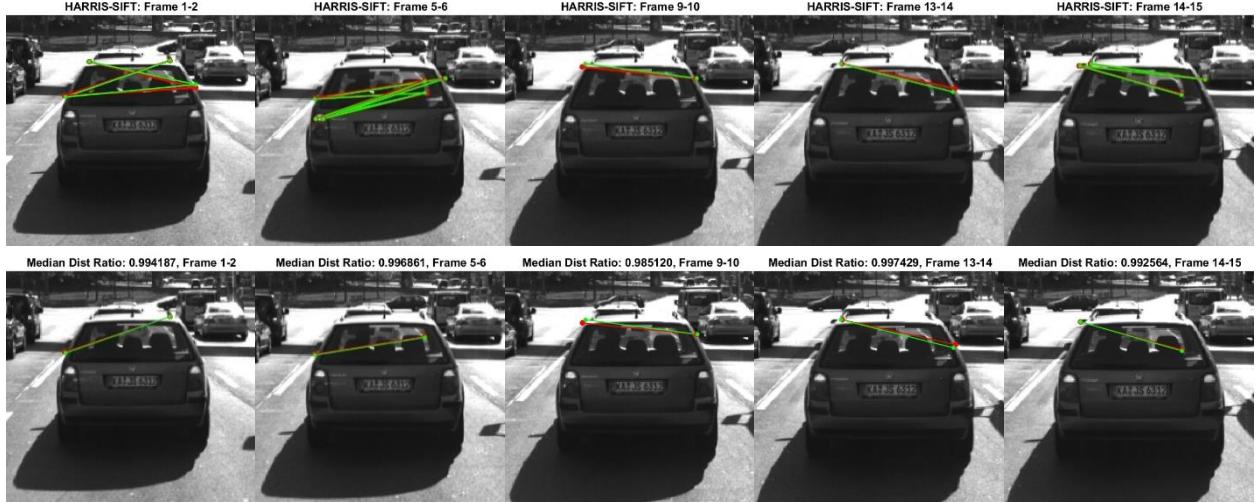


Figure 63: Frames with negative TTC: top row shows all keypoint pairs, bottom row shows median dist ratio pair

ORB-SIFT

ORB detectors produce a large number of keypoints over a relatively small number of distinct locations. Each location consists of a group of concentric circles that represent collection regions of varying size, scale and orientation. These features may prove helpful in keypoint matching, but we still end up with redundant matches within groups.

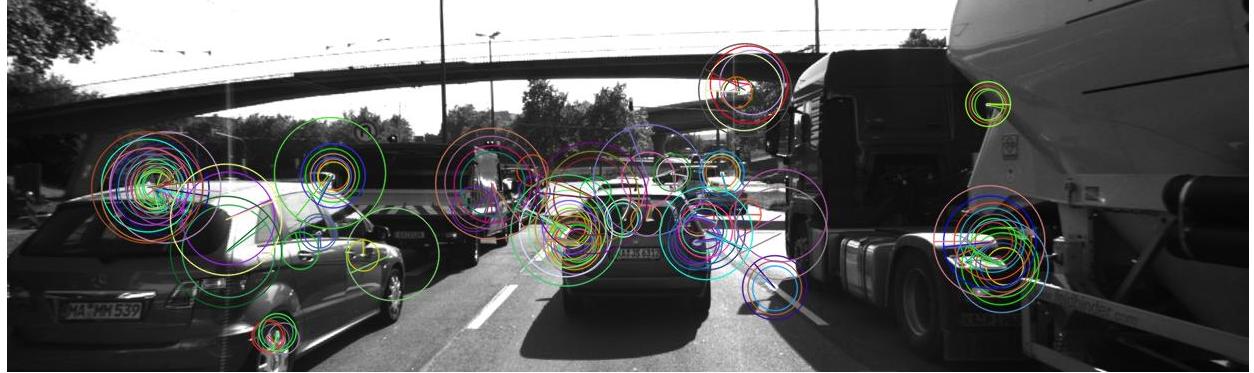


Figure 64: ORB keypoints that were associated with exactly one bounding box on Frame 3 (overlaps removed)

Our camera-based TTC algorithm currently only uses the centroid information from keypoints, effectively collapsing the multi-scale regions into groups of matches with small variations in keypoint positions. This could present a risk factor for our estimation approach, where the smallest errors in keypoint position can be amplified at the TTC estimation level.

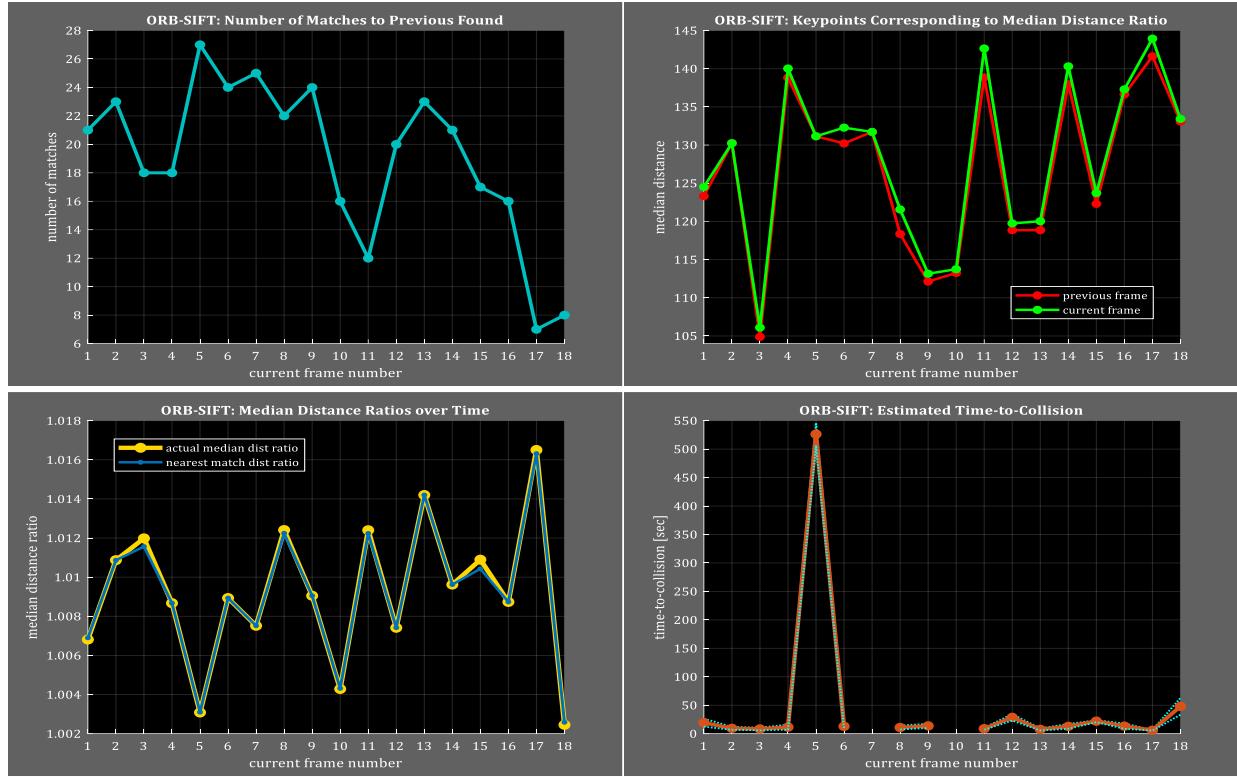
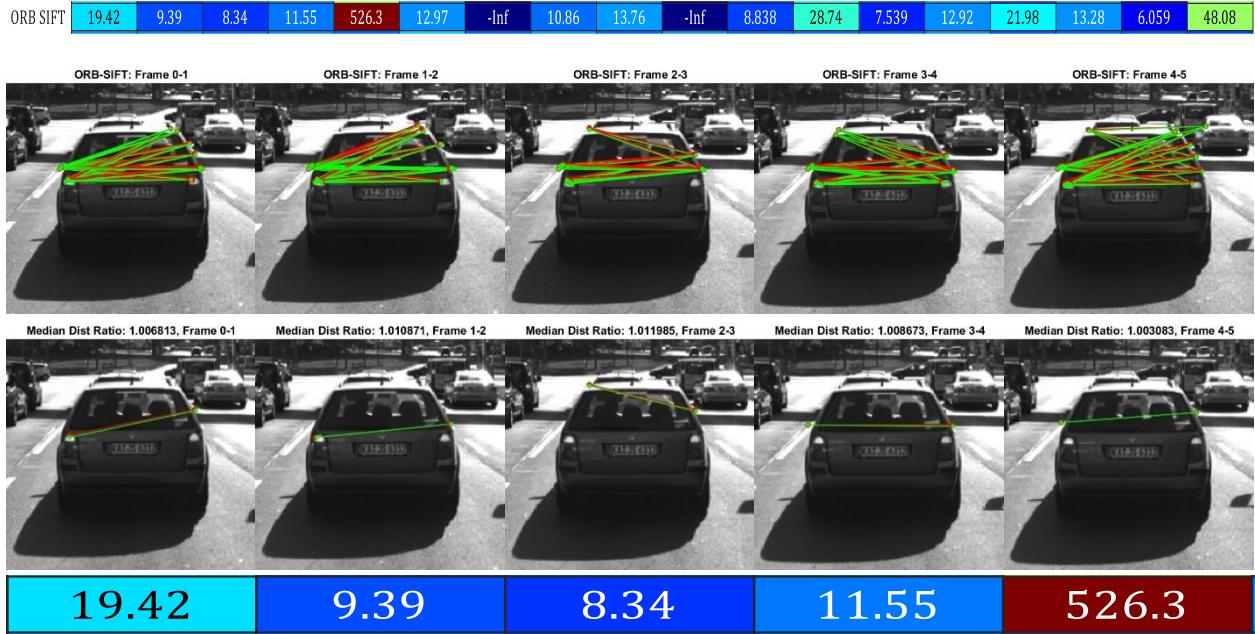


Figure 65: Median Distance Ratio and Estimated TTC based on ORB Detector & SIFT Descriptor



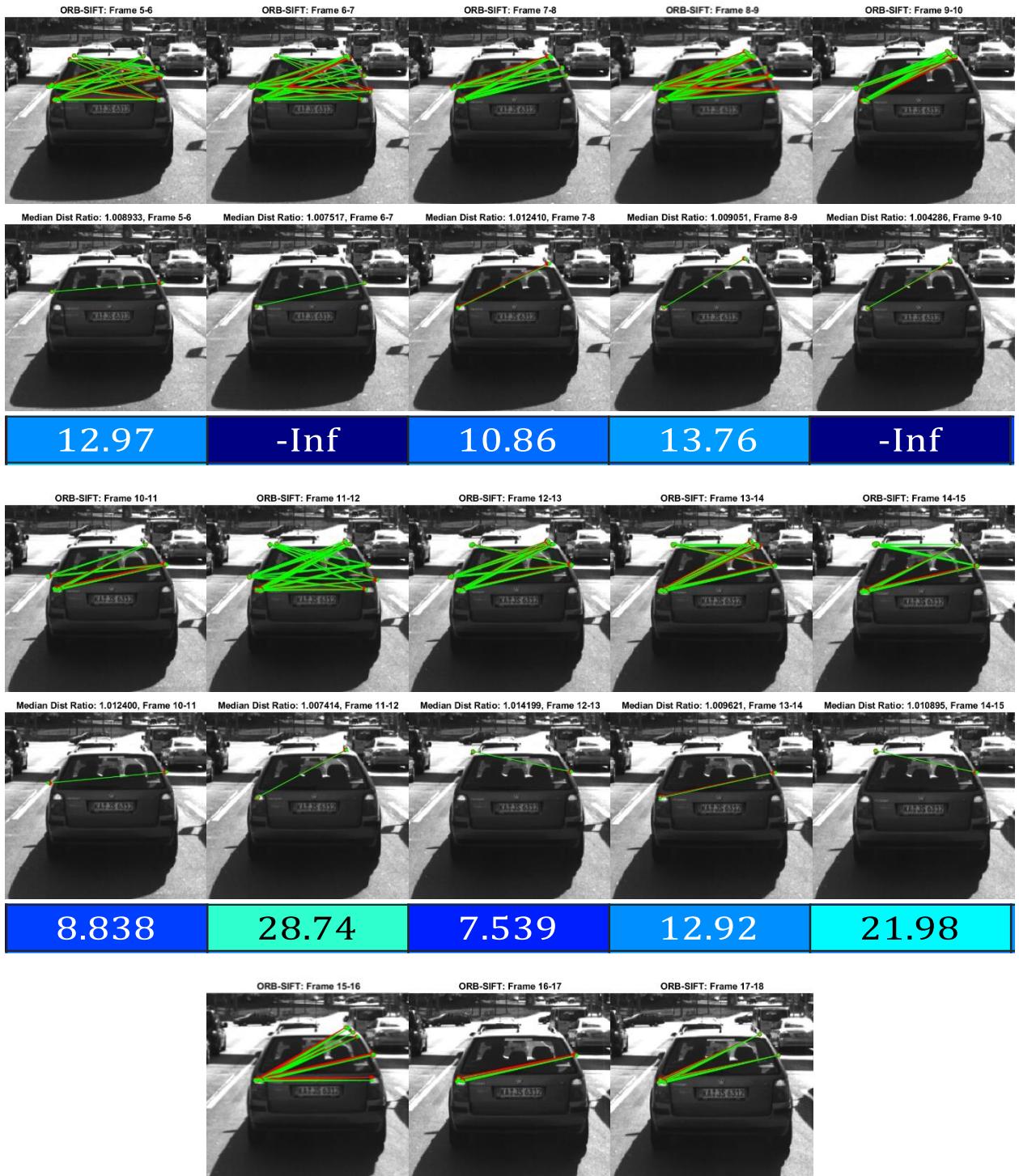




Figure 66: Keypoint Pairs Corresponding to the Median Distance Ratio & TTC

I selected frames 5, 7, 10 and 18 below as examples of cases for which the TTC estimates seemed way off. We observe how ORB detectors amass a large number of keypoints, but those keypoints are clustered around only ~3-8 distinct feature points on the vehicle (and sometimes erroneous points outside of the vehicle). Nonetheless, there are enough observable offsets in the groups of red (previous frame) vs green (current frame) distance pairs in the first row of images below that we can actually visualize the motion from frame-to-frame. The problem then arises when the single median distance ratio selected from these distributions (second row) turns out to be a poor representation of the motion observed in the aggregate sample distribution (first row):

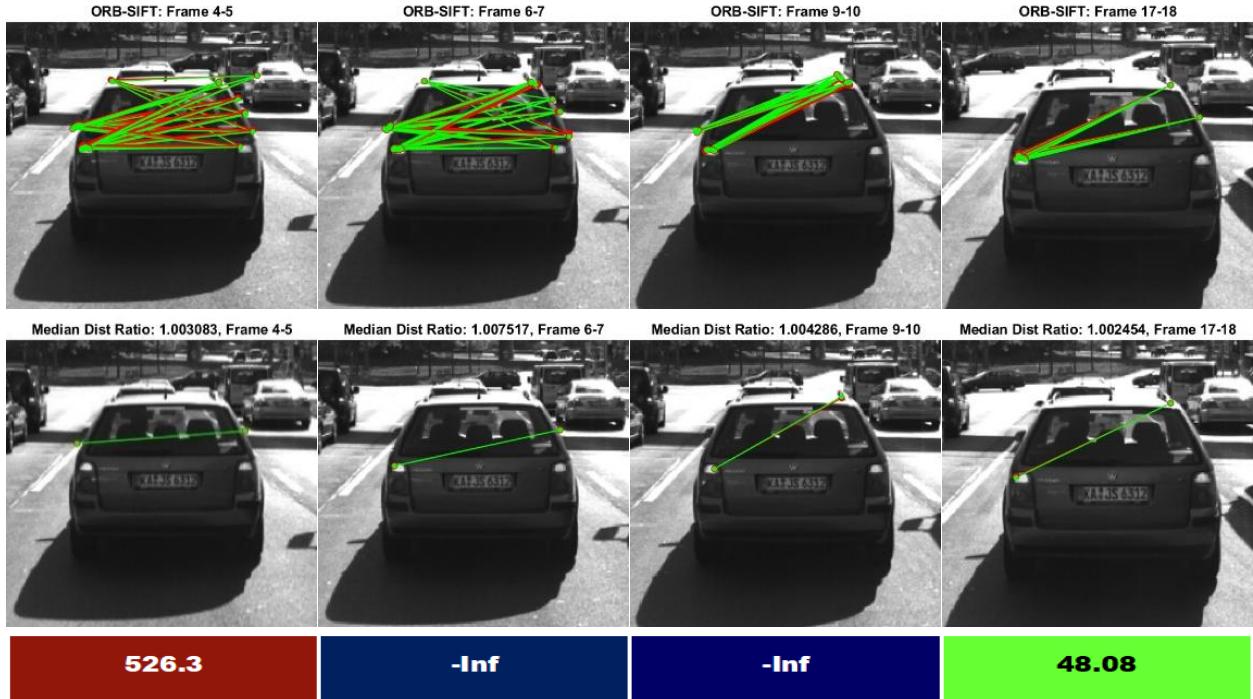


Figure 67: Top Row: Distal Keypoint Pairs, Middle Row: Median Distance Pair, Bottom: Estimated Time-to-Collision

Frames 4-5, 6-7 and 9-10 above seemed to “fail” because almost no motion was detected based on the median distance ratio and corresponding previous and current keypoint positions – and our TTC equation assumes constant velocity during these time intervals. These may have been unlucky draws, because the

offsets can be observed when looking at the full sets of distance ratios. To mitigate these cases, we may want to look at different approaches to estimating/describing the distance ratio, or possibly increasing the temporal gap between successive frames to increase the likelihood of selecting keypoints with observable changes in distance ratio.

Frames 4 & 5 also appear to have selected a keypoint that belongs to a white lane marking on the left-hand side, which would also corrupt our TTC estimate on the preceding vehicle. This issue of keypoints picking up features outside of the vehicle may be happening in other cases as well.

Performance is also dependent upon the detector-descriptor approach used, since inaccuracies at the finer granularities are amplified by the constant-velocity TTC equation. In the case of ORB, many keypoint candidates are presented in the vicinity of salient features, but errors are introduced when we try to collapse the larger collection regions into a single point.