Hello! Thank you for checking out my assignment submission!

Please wait while the analysis is being generated.

The animated icon to the top right shows that the system is processing the data.

Plots and data will appear here sequentially as they are generated. :D

We start by generating some statistics for our dataset splits.

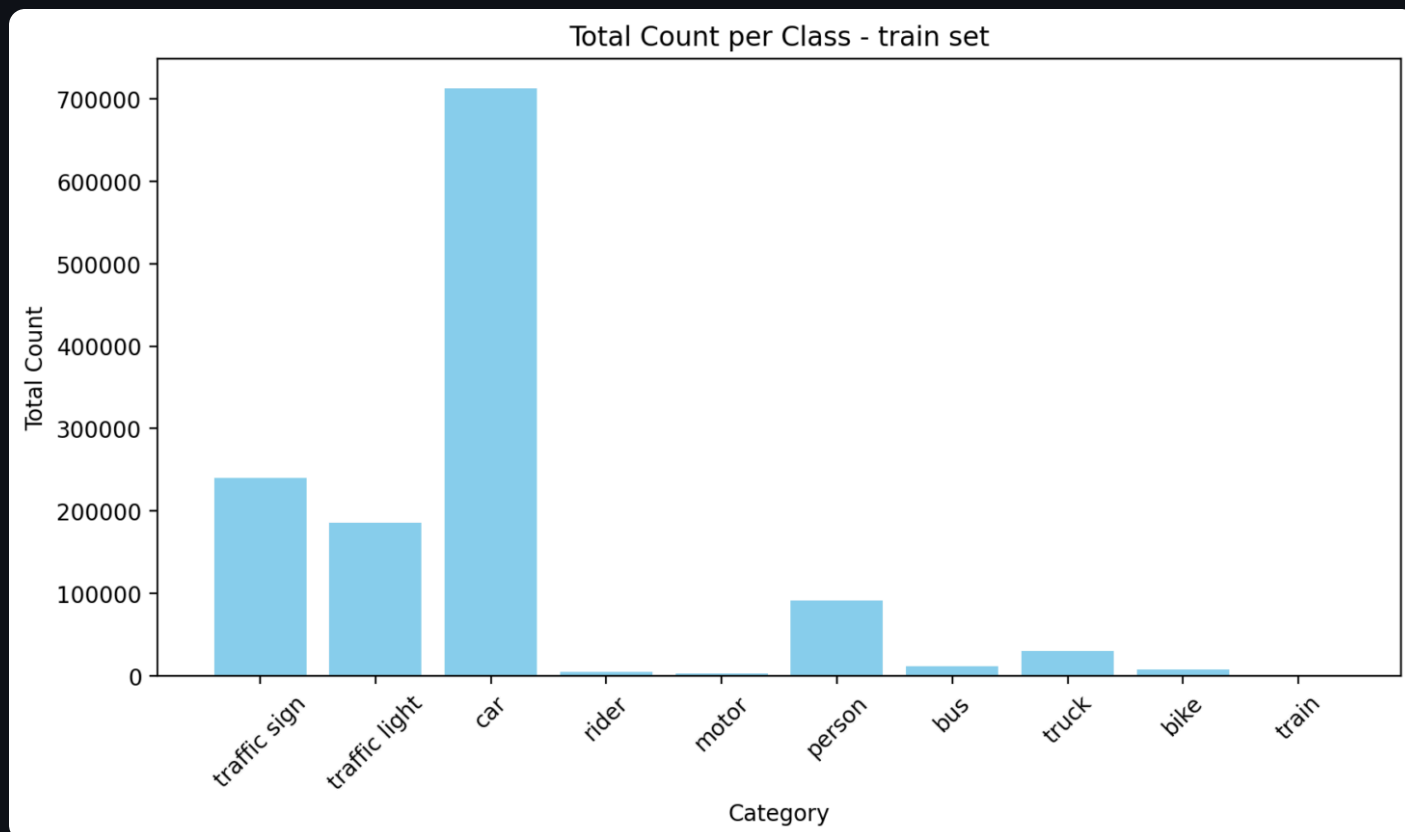This is what the train statistics looks like:

| | traffic sign | traffic light | car | rider | motor | person |
|---|---|---|---|---|---|---|
| total_count | 239686 | 186117 | 713211 | 4517 | 3002 | |
| occluded_count | 26974 | 5911 | 483121 | 4028 | 2297 | |
| truncated_count | 6699 | 4923 | 66529 | 227 | 272 | |
| max_area | 917709.771 | 302654.3096 | 612645.1832 | 257238.437 | 316440.6389 | 3449 |
| min_area | 3.5722 | 0.9366 | 0.8714 | 5.6019 | 30.9835 | |
| sum_area | 287367954.3891 | 94430896.467 | 6720212736.6057 | 28521339.8285 | 22767430.4156 | 2690238 |
| mean_area | 1198.9351 | 507.3738 | 9422.4749 | 6314.2218 | 7584.0874 | 294 |
| max_width | 1279.2694 | 557.525 | 1278.5979 | 553.1279 | 773.633 | 84 |
| min_width | 0.4538 | 0.1066 | 0.4261 | 1.2874 | 4.4931 | |
| sum_width | 7736763.6674 | 2954222.5906 | 53217431.1529 | 196499.8462 | 201506.9276 | 252783 |

This is what the val statistics looks like:

|  | traffic sign | traffic light | car | rider | motor | person |
| --- | --- | --- | --- | --- | --- | --- |
| total_count | 34908 | 26885 | 102506 | 649 | 452 | 132 |
| occluded_count | 4020 | 914 | 69382 | 573 | 337 | 76 |
| truncated_count | 925 | 736 | 9480 | 28 | 48 | 4 |
| max_area | 105944.9079 | 40551.3478 | 455712.645 | 191065.2802 | 185101.9834 | 155240.17 |
| min_area | 3.9836 | 2.3143 | 4.3117 | 24.4497 | 44.8153 | 3.37 |
| sum_area | 41663148.2965 | 13375402.9858 | 962431227.8182 | 3875298.3798 | 3525443.6238 | 38193918.24 |
| mean_area | 1193.5129 | 497.5043 | 9389.0234 | 5971.1839 | 7799.654 | 2879.95 |
| max_width | 952.5307 | 251.0755 | 1096.7758 | 477.2594 | 411.7499 | 365.83 |
| min_width | 0.8161 | 0.1055 | 1.0242 | 1.7482 | 5.9381 | 1.72 |
| sum_width | 1130269.2272 | 426504.5387 | 7612536.1056 | 27800.4585 | 30853.7031 | 364109.76 |

Although these statistics are useful, they are hard to interpret in tabular form.

Let's generate some plots to visualize the data!

It looks like car is the category with the most number of occurences, followed by traffic sign, and then by traffic light.

person class also has a significant presence in the dataset.

train, motor, and rider are have the three lowest counts.

I hypothesize from this information, that the model (if trained on this dataset) will perform well if tasked with detecting cars, but will not perform so well when tasked with detecting the three classes with the lowest counts
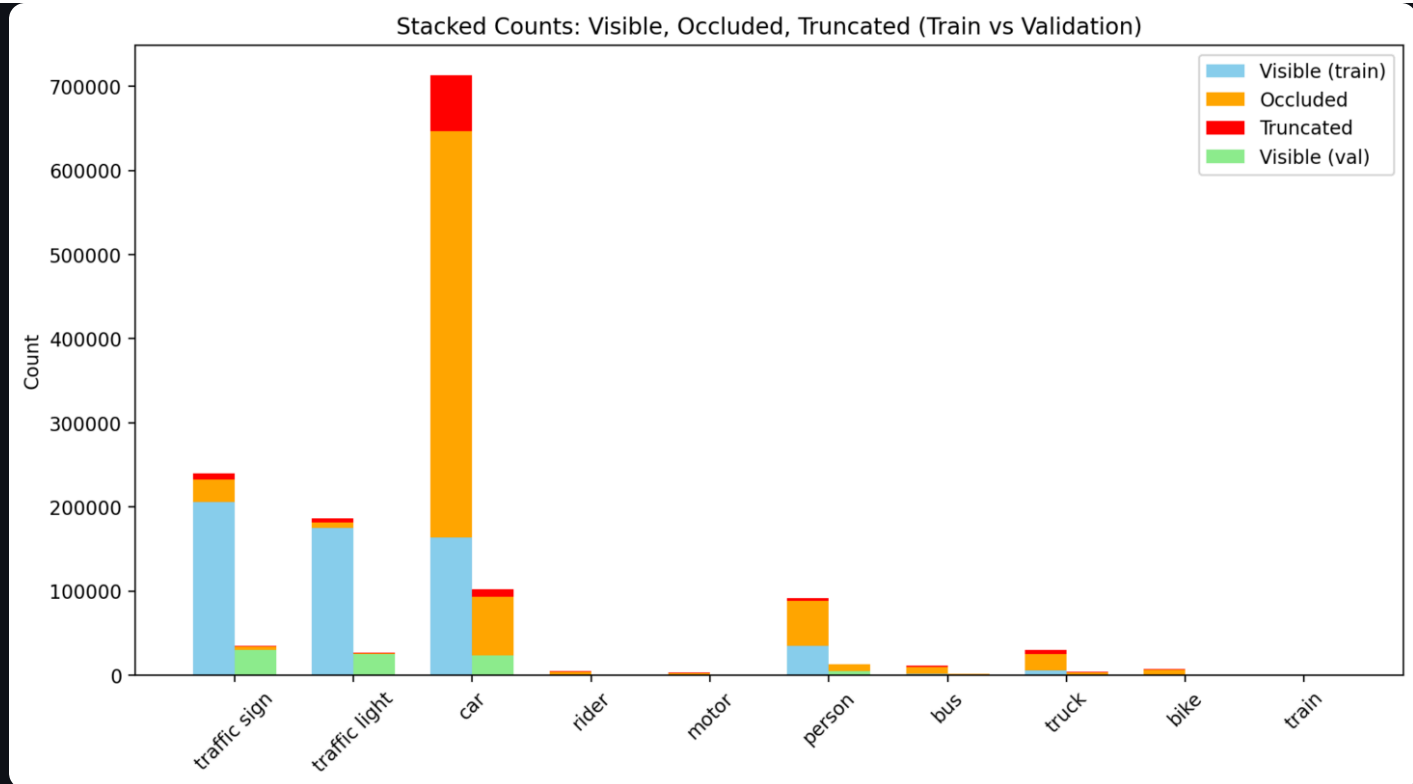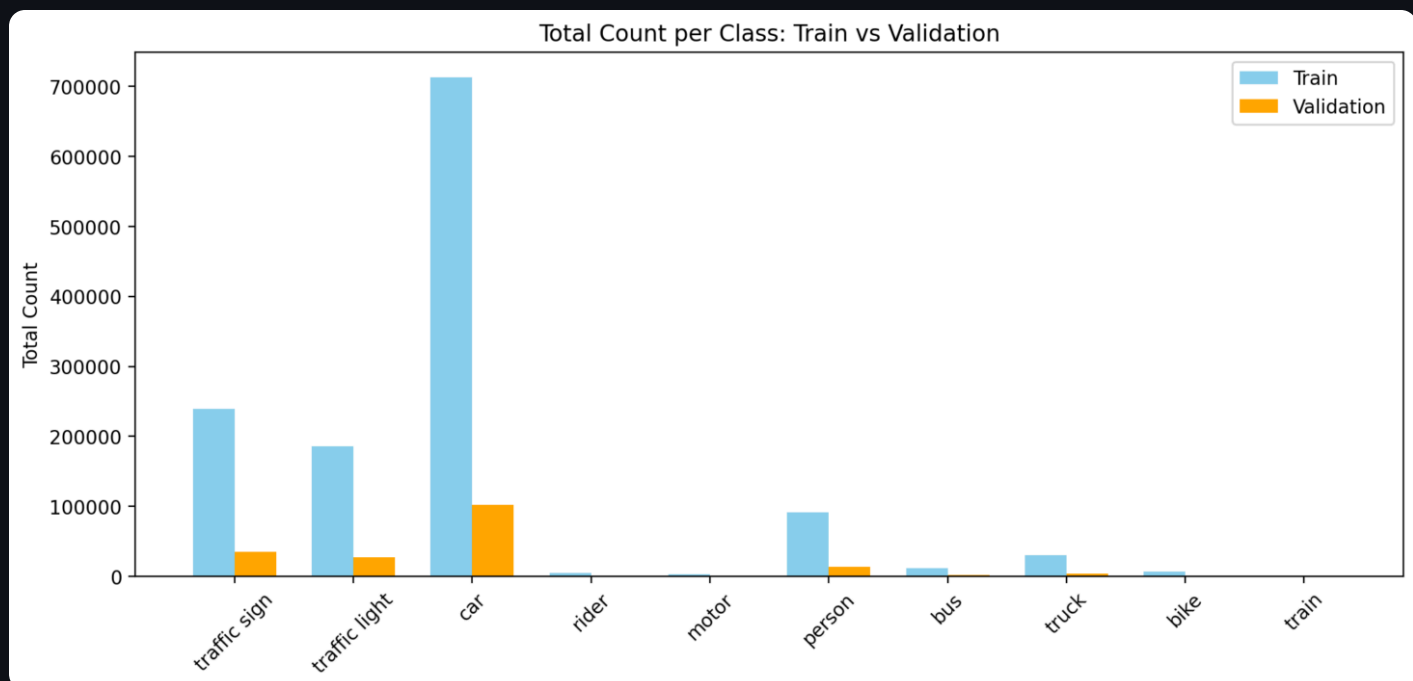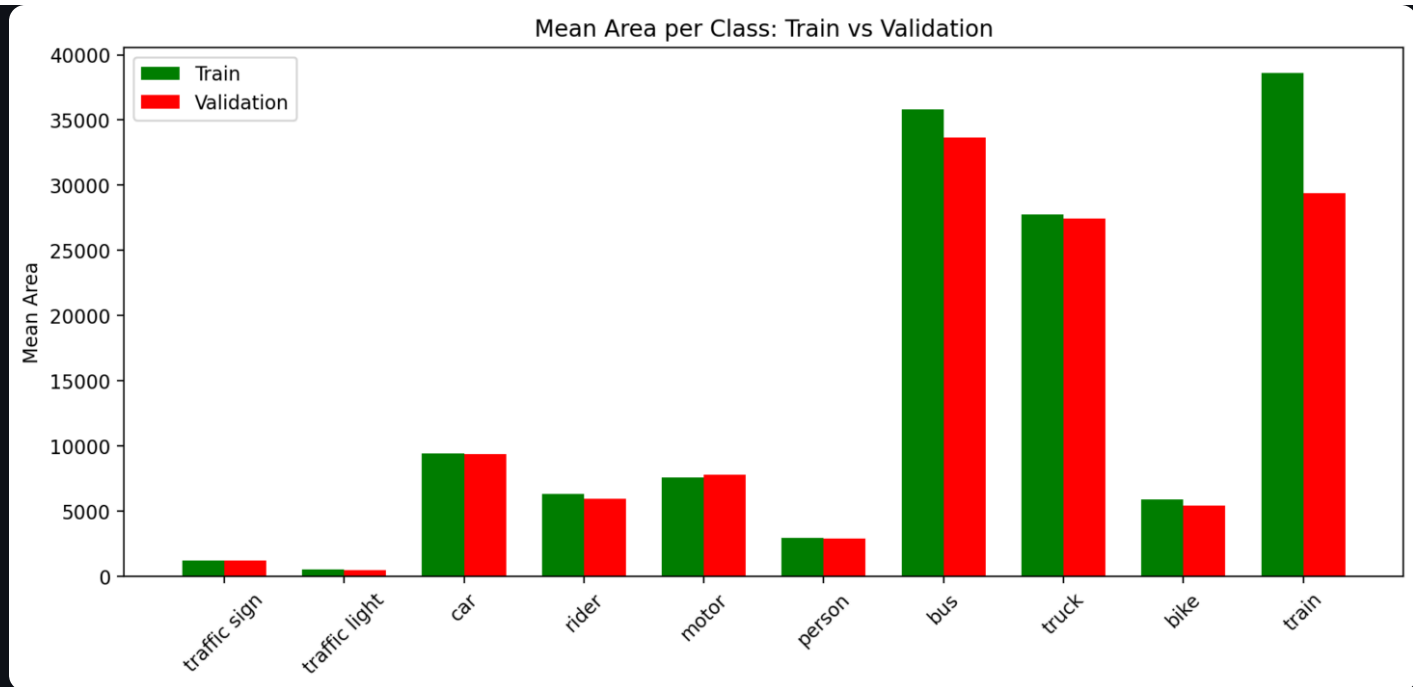


It looks like car is the category with the most number of occurences, followed by traffic sign, and then by traffic light.

person class also has a significant presence in the dataset.

train, motor, and rider are have the three lowest counts.

I hypothesize from this information, that the model (if trained on this dataset) will perform well if tasked with detecting cars, but will not perform so well when tasked with detecting the three classes with the lowest counts

Stacked Counts: Visible, Occluded, Truncated (Train vs Validation)

We see here that a huge number of cars in the train set are either truncated and occluded.

While this is undesirable in a small dataset, the large size of this dataset will help the model identify cars in various settings.
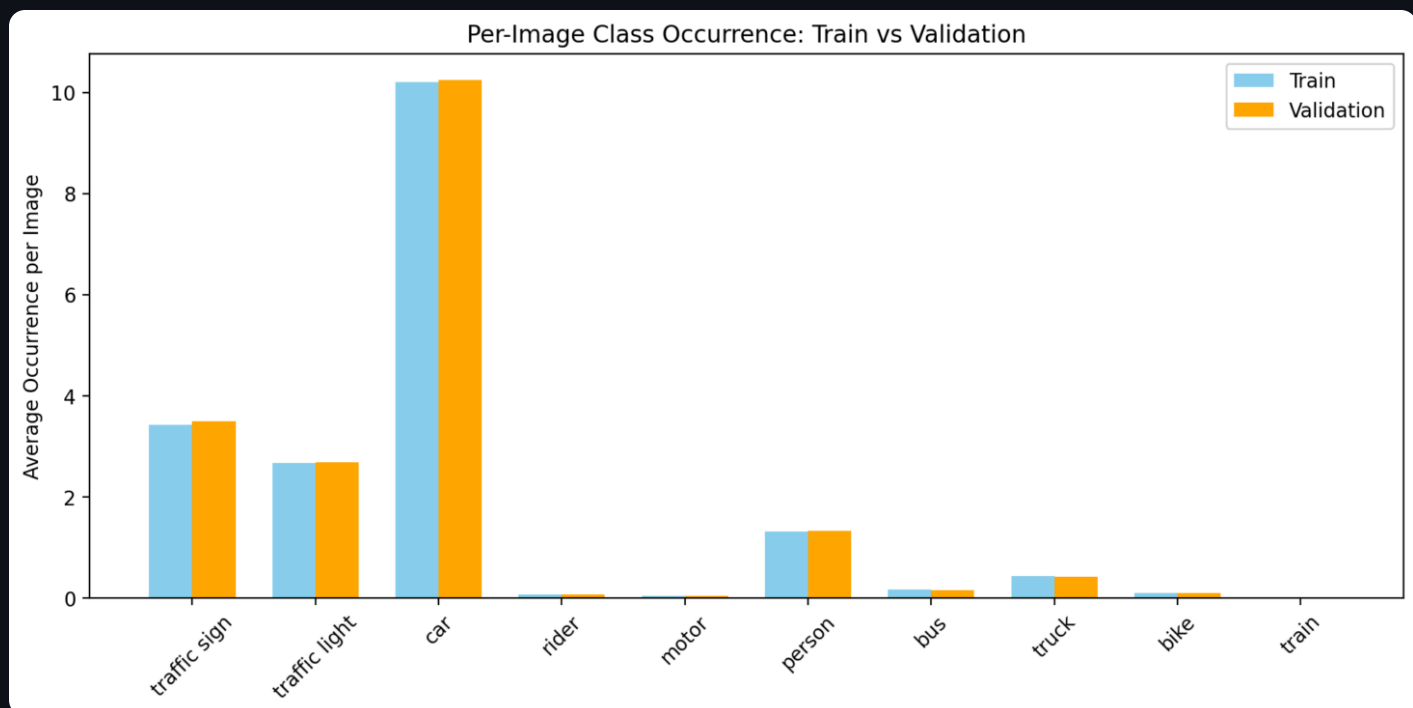


Total Count per Class: Train vs Validation

Mean Area per Class: Train vs Validation

The train class has a smaller average area in the validation set

But considering the small number of samples, this differnce is expected

Otherwise, all other classes except bus have more or less the same average size in both splits



Per-Image Class Occurrence: Train vs Validation

It can be seen that both splits have a similar average number of class occurences per image
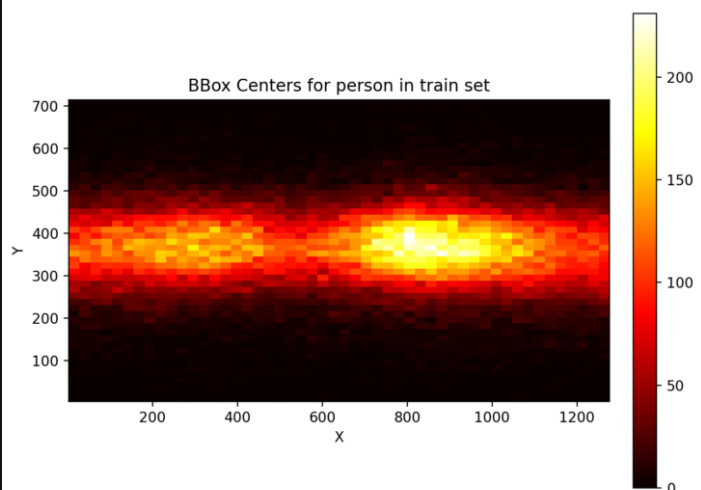
Per-class Boxplot of BBox Area - train
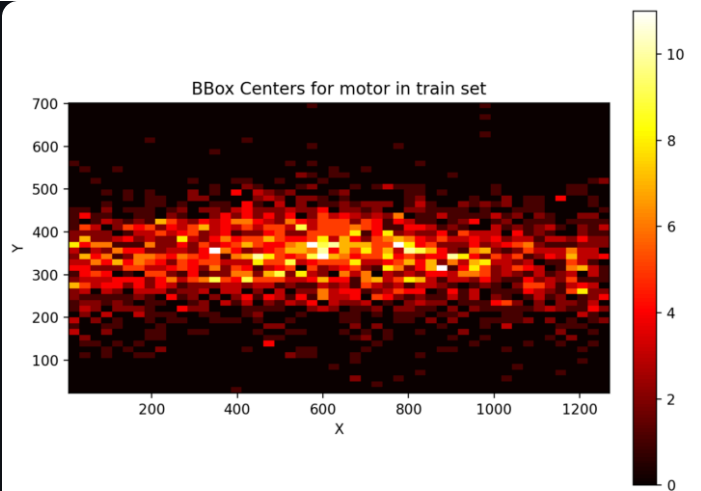
Per-class Boxplot of BBox Area - val

The boxplots reveal that the train set has several large outliers that cover over 40 percent of the image.
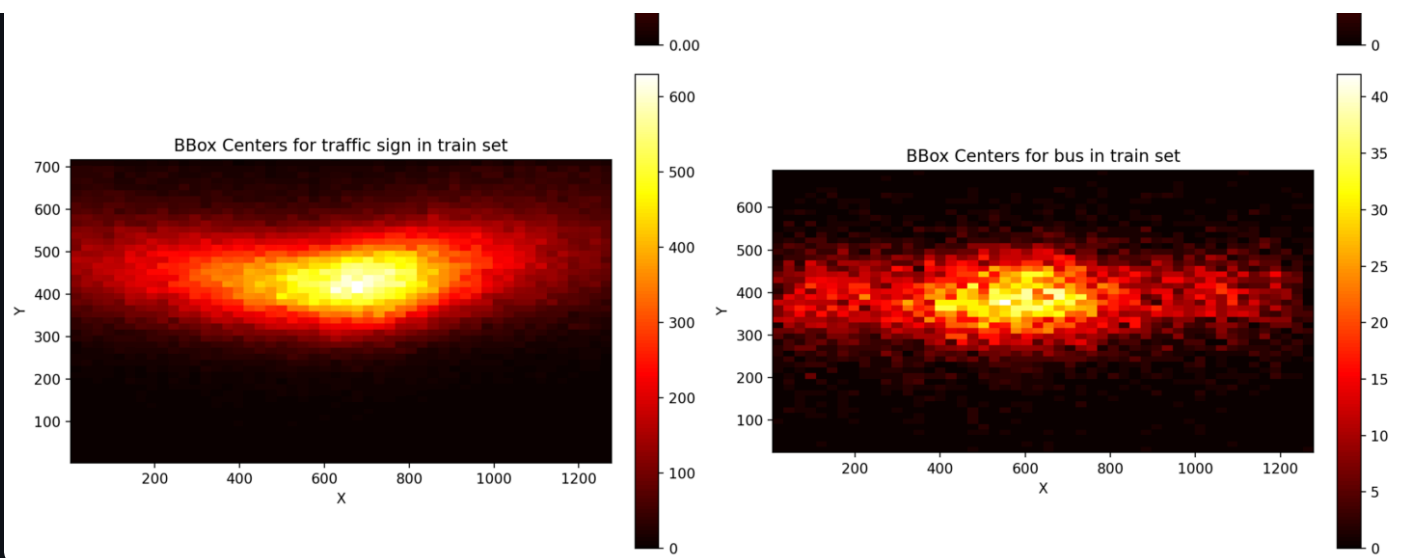
There are two images in the train set with traffic lights that cover almost the entire image.

There are three images with buses and four images with trucks that cover over 70 percent of the images in the train set.
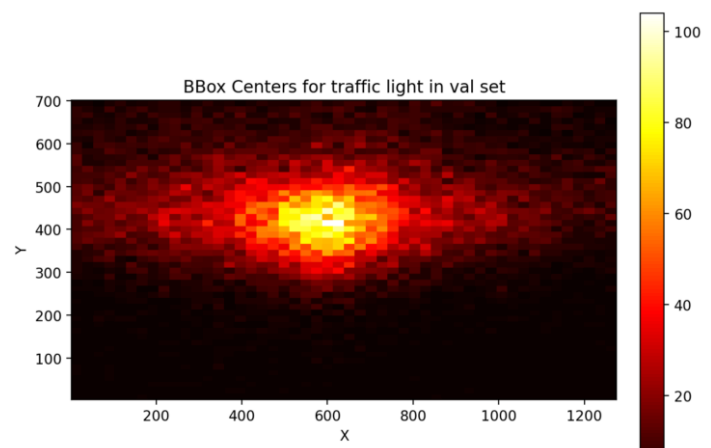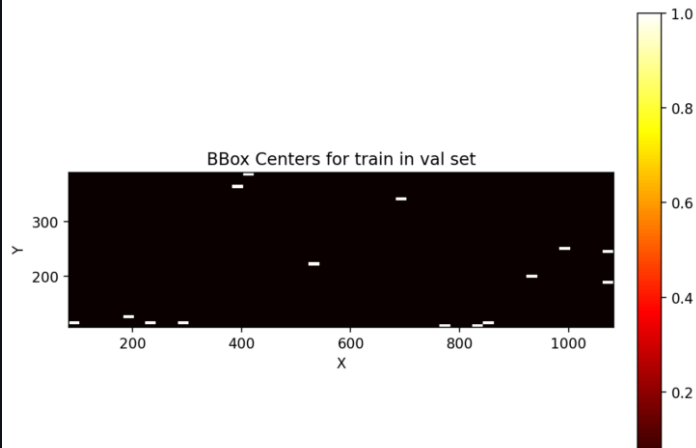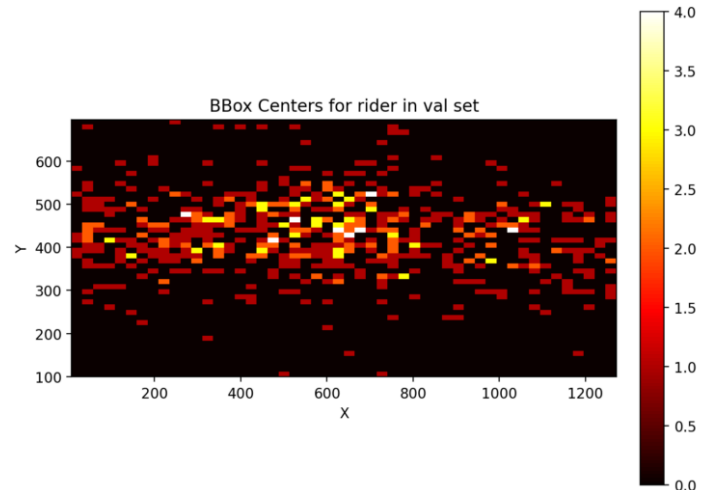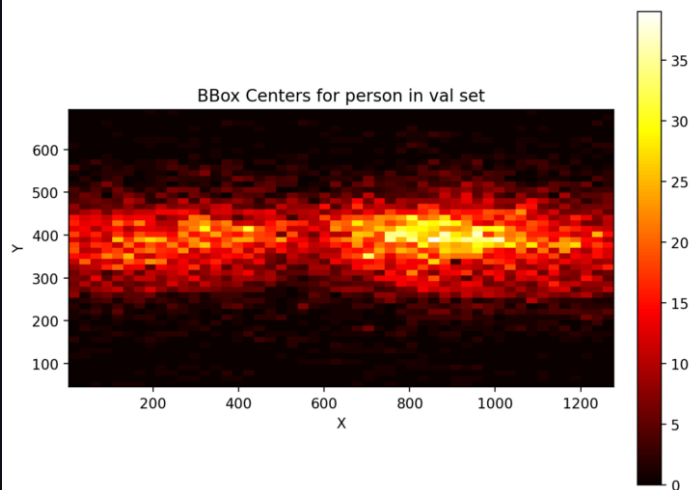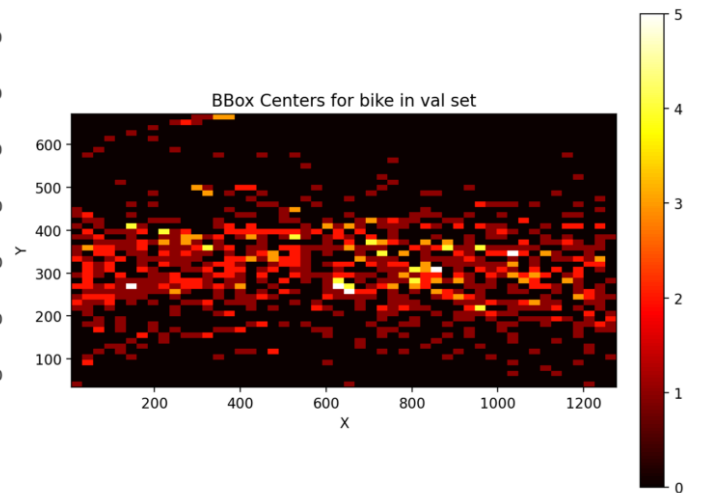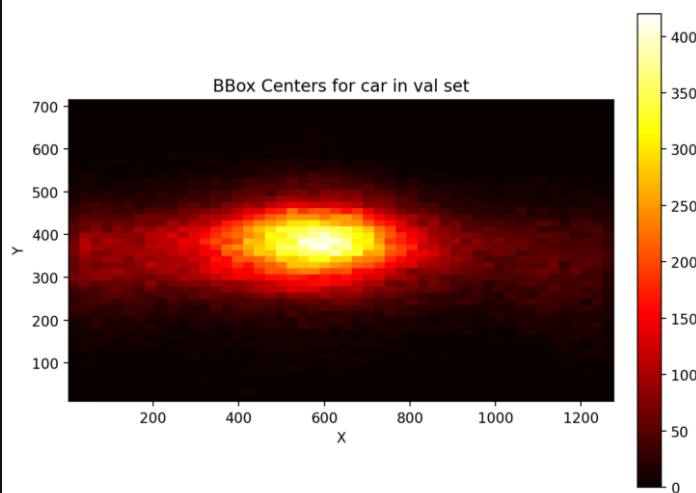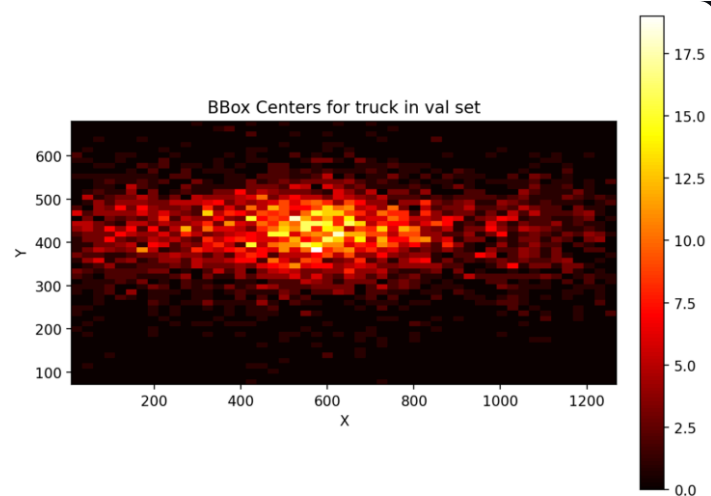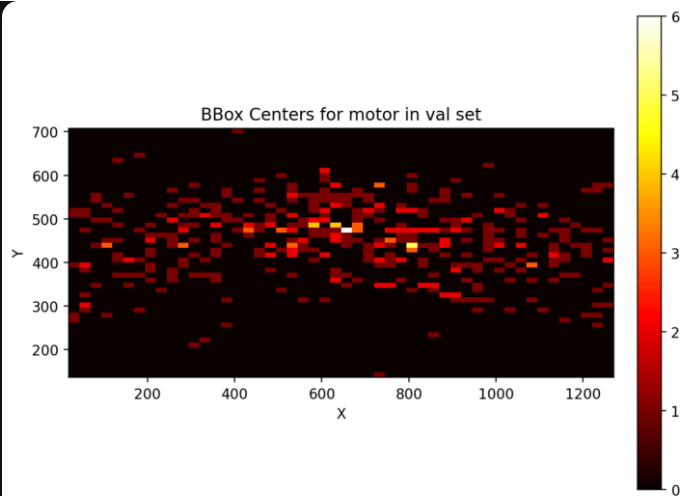
The validation set outliers are much smaller, and there are only about 5 truck images that cover more than 50% of the image.

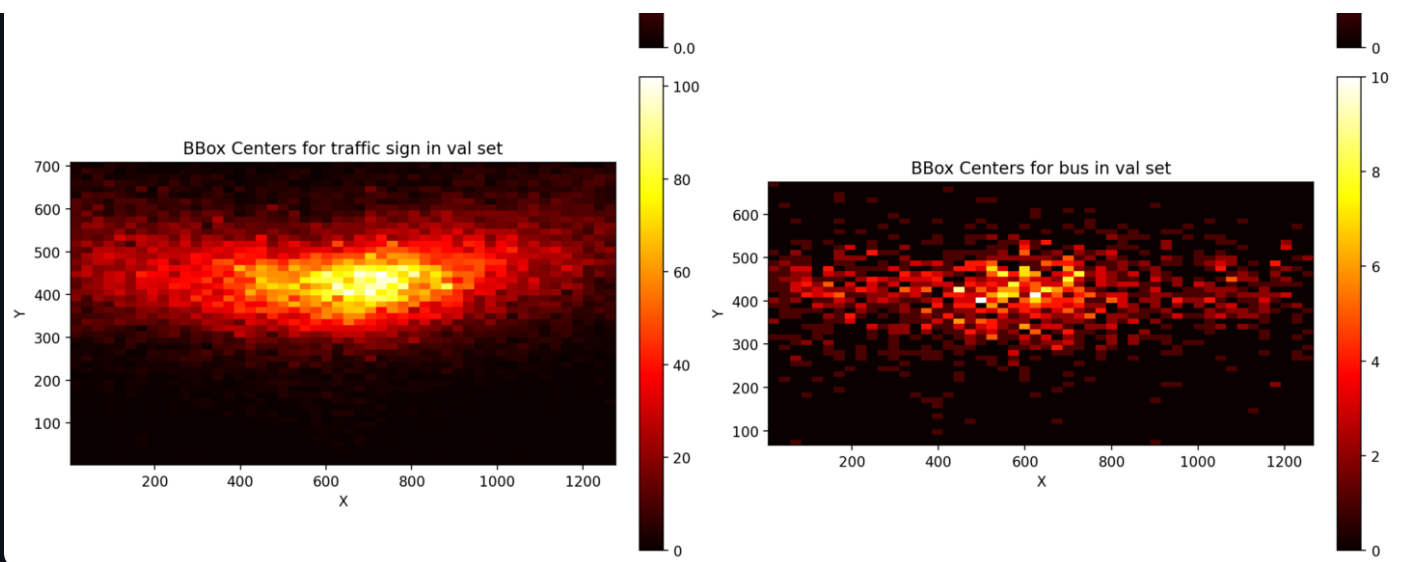Here is what the class-wise heatmaps look like for the train set

Here is what the class-wise heatmaps look like for the val set

BBox Centers for motor in val set

BBox Centers for truck in val set

BBox Centers for car in val set

BBox Centers for bike in val set

BBox Centers for person in val set

BBox Centers for rider in val set

BBox Centers for train in val set

BBox Centers for traffic light in val set

We see from these heatmaps that the distribution of cars, traffic lights, and traffic signs are somewhat uniform

and the distribution of bike, motor, and rider are rather haphazard in both sets, despite their small counts.

The distribution of the person class, which has two modes, is an interesting way to learn that the dataset features several scenes where pedestrians are on either side of the road

That brings us to the end of the analysis.

While I would have liked to do a lot more in this task and visualize individual examples, I have unfortunately run out of time.

Thank you!