

---

# Predictive Inequity in Object Detection

---

Benjamin Wilson<sup>1</sup> Judy Hoffman<sup>1</sup> Jamie Morgenstern<sup>1</sup>

## Abstract

In this work, we investigate whether state-of-the-art object detection systems have *equitable predictive performance* on pedestrians with different skin tones. This work is motivated by many recent examples of ML and vision systems displaying higher error rates for certain demographic groups than others. We annotate an existing large scale dataset which contains pedestrians, BDD100K, with Fitzpatrick skin tones in ranges [1-3] or [4-6]. We then provide an in depth comparative analysis of performance between these two skin tone groupings, finding that neither time of day nor occlusion explain this behavior, suggesting this disparity is not merely the result of pedestrians in the 4-6 range appearing in more difficult scenes for detection. We investigate to what extent time of day, occlusion, and reweighting the supervised loss during training affect this predictive bias.

## 1. Introduction

The methods and models developed by the machine learning community have begun to find homes throughout our daily lives: they shape what news stories and advertisements we see online, the engineering of new products sold in stores, the content of our emails, and, increasingly, the allocation of resources and surveillance. Both private and governmental organizations have increasingly begun to use such statistical methods. Examples of the latter include predictive policing, recidivism prediction, and the allocation of social welfare resources (Rubin, 2010; Maloof, 1999; Kriegler & Berk, 2010; Chouldechova et al., 2018). A particularly driving application domain for ML in the private sector is the design of autonomous vehicles. Autonomous vehicles may greatly reduce transit costs of goods and reduce individuals’ reliance on owning personal vehicles.

Recognizing key objects such as pedestrians and road signs plays a key role in these systems, helping determine when

a car must brake or swerve to avoid fatalities. The few autonomous vehicle systems already on the road have shown an inability to entirely mitigate risks of pedestrian fatalities (Levin & Wong, 2018). A natural question to ask is *which* pedestrians these systems detect with lower fidelity, and why they display this behavior. In this paper, we study the performance of several models used in state-of-the-art (He et al., 2017) object detection, and show *uniformly poorer performance of these systems when detecting pedestrians with Fitzpatrick skin types between 4 and 6*. This behavior suggests that future errors made by autonomous vehicles may not be evenly distributed across different demographic groups.

We then investigate *why* standard object detection might have higher predictive accuracy for pedestrians lower on the Fitzpatrick scale. The training set has roughly 3.5 times as many examples of lower-Fitzpatrick scored pedestrians compared to higher scored pedestrians, which suggests several different sources of predictive disparity between the two groups might be at work. First, one would expect to have lower generalization error on the larger subset of data. Second, many standard loss functions will prioritize accuracy on the larger subset of the data.

These two behaviors, and others, are often conflated and described by both industry and researchers as “biased data”, a shorthand for a milieu of different issues of sampling, measurement, and weighting of different design goals. Sampling issues arise when a dataset does a poor job representing the larger population (e.g., a dataset with mostly men, or no examples of women who successfully repaid mortgages). Issues of measurement arise when the features collected are insufficient to accurately measure and predict the intended outcome variable (such as banking records having insufficient information to predict creditworthiness in “unbanked” communities, where participation in lending circles (White, 2016) and other less centralized systems better predict loan repayment). Designing an objective function (and constraints) for training an ML system tacitly weights different model behaviors over others, such as the weighting of false positives versus false negatives. Relatively little work in the fair ML community has explicitly teased apart these sources for particular examples of inequitable predictive behavior. **We explicitly aim to measure three possible sources of predictive imbalance: whether time**

<sup>1</sup>Georgia Tech. Correspondence to: Benjamin Wilson <benjaminwilson@gatech.edu>.

of day, occlusion, or if the loss function prioritizing accuracy on the larger population heavily impacts this behavior.

## 2. Related Work

Predictive disparities of ML systems have recently been given much attention. These examples appear in numerous domains, a few of which we mention here. Advertising systems show ads based upon numerous demographic features; as a result, a number of findings have shown certain gender or racial groups receive certain ads at much higher rates than others (Zhao et al., 2017; Datta et al., 2018; Sweeney, 2013). Certain predictive policing systems has been shown to have differing predictive performance based on race (Angwin et al., 2016; Selbst, 2017; Lum & Isaac, 2016).

Most closely related to our current work are examples of vision-based systems with differing predictive qualities for women or people of color. Facial recognition systems (and other systems which use facial images) have garnered the lion’s share of the press in this space.

Early warnings that facial recognition might have higher accuracy on white men showed that this problem might be somewhat mitigated by training systems separately for different demographic groups (Klare et al., 2012). Nevertheless, recent, state-of-the art systems designed by many major tech conglomerates have continued to face scrutiny for the behavior of their facial recognition systems. Commercial gender prediction software has been shown to have much worse accuracy on women with Fitzpatrick skin types 4-6 compared to other groups (Buolamwini & Gebru, 2018); this work inspired our use of the Fitzpatrick skin scale to categorize pedestrians. The ACLU found that Amazon’s facial recognition system incorrectly matched a number of darker-skinned members of congress to mugshots from arrests across the country (Snow, 2018).

Our work instigates the measurement of predictive imbalance for a different set of ML-guided vision systems, namely that of driving-centric object detection. This work is particularly timely as several locations have recently allowed autonomous vehicles to operate on public roads, several casualties have resulted. We do not focus on the ethical dilemmas algorithms might ultimately face<sup>1</sup>, but instead on the simpler question of whether several simple building blocks used for research-grade pedestrian detection have similar ability to detect pedestrians with different skin tones.

We study the problem of pedestrian detection in road scenes from the perspective of an autonomous vehicle. Many datasets have been introduced in the computer vision com-

munity for developing methods for recognizing pedestrians at a variety of distances (Geiger et al., 2012; Ess et al., 2008; Dollár et al., 2012) and for recognizing all objects in road scenes relevant to the autonomous driving task (Cordts et al., 2016; Yu et al., 2018). State-of-the-art region-proposal based methods for general object detection such as Faster R-CNN (Ren et al., 2015) or Mask R-CNN (He et al., 2017) often form the backbone for the best performing pedestrian detection models for unoccluded pedestrians close to the vehicle for which more pixels are available for recognition. While the same region proposal based method can be effective for proposing bounding boxes for small (distant) pedestrians (Zhang et al., 2016), often unique representations are needed for simultaneous detection of small (distant), large (nearby), and occluded pedestrians. Recent works have addressed this is by incorporating multiple representations into their model, one for each pedestrian scale (Li et al., 2018), or for different body parts to handle occlusion (Zhang et al., 2018). In this work we compare detection performance of nearby and largely unoccluded pedestrians for which skin-tone is readily identifiable and therefore focus our analysis on the core technology that persists between object and pedestrian detection systems.

## 3. Preliminaries

We begin with an overview of the main concepts used in this work: the problem of pedestrian detection; the classification of people into groups based on skin tone and other characteristics known as Fitzpatrick skin typing; and predictive disparity, or the difference in predictive performance of a learning system on two different groups of datapoints.

**Pedestrian Detection** Quickly identifying pedestrians has been a long-standing challenge in the Computer Vision community. From security systems to autonomous cars, Pedestrian Detection remains an important and crucial aspect of many Computer Vision models. Common challenges of Pedestrian Detection include: occlusion by other objects or people, changes in clothing, and diverse lighting conditions.

**Fitzpatrick Skin Typing** The Fitzpatrick skin type scale (Fitzpatrick, 1975), introduced to predict a person’s predisposition to burning when exposed to UV light, measures a number of physical attributes of a person including skin, eye, and hair color, as well as a person’s likelihood to freckle, burn, or tan. As a general rule, categories 1-3 correspond to lighter skin tones than 4-6. This categorization aims to design a culture-independent measurement of skin’s predisposition to burn, which correlates with the pigmentation of skin.

**Predictive Inequity** Assuming a fixed partition of the person class (for example, into classes  $LS$  and  $DS$  based on the Fitzpatrick skin scale), we define the *predictive inequity* of a

<sup>1</sup>Should the car crash into the person on the left or right, if those are the only options? Several papers have investigated such questions, e.g. Roff (2018).

model  $f$  for a particular loss function  $\ell$  to be the difference in loss the model  $f$  incurs on members of  $\mathcal{L}\mathcal{S}$  over  $\mathcal{D}\mathcal{S}$ :

$$\mathbb{E}_{\substack{x, x' \sim D \\ x \in \mathcal{L}\mathcal{S}, x' \in \mathcal{D}\mathcal{S}}} [\max\{\ell(f(x)) - \ell(f(x')), 0\}]. \quad (1)$$

This definition measures the *average* additional loss a model would experience for a random member of  $\mathcal{L}\mathcal{S}$  versus  $\mathcal{D}\mathcal{S}$ . While any number of measurements of “fairness” of ML systems have been proposed (e.g., statistical parity, equality of false positives or negatives, calibration, individual fairness (Hardt et al., 2016; Dwork et al., 2012; Kleinberg et al., 2016; Chouldechova, 2017)), many of them involve some per-instance loss function. Our measure does not perfectly capture all aspects of some model  $f$ ’s behavior for all populations, but it does attempt to measure whether a model does a similarly good job minimizing a particular loss function for two different populations simultaneously.

#### 4. Evaluating Predictive Inequity

Our goal is to quantify any disparity in predictive performance of standard recognition models across groups of people with varying skin tones. As no benchmark currently exists for this task, we first describe our methodology for collecting the necessary annotations for performing our evaluations in Section 4.1. Next, we provide evaluation of predictive inequity on our benchmark as well as an in depth analysis of the sources of inequity in Section 4.2. Finally, we propose a simple remedy to reduce this predictive inequity in Section 5.

##### 4.1. Benchmarking Predictive Disparity for Pedestrian Detection

In order to measure predictive disparity of a particular model, we need a partition of a dataset into demographic classes, in our case classes  $\mathcal{L}\mathcal{S}$  and  $\mathcal{D}\mathcal{S}$ . Many instances of bias in ML systems are those where elements of a dataset are not explicitly labeled by their demographic information. In these cases, we cannot necessarily assume the membership of data elements into  $\mathcal{L}\mathcal{S}$  and  $\mathcal{D}\mathcal{S}$  are known, and instead must gather that information as part of our training and evaluation of a system.

For the task of measuring the predictive inequity of object detection on pedestrians of different Fitzpatrick skin types, this corresponds to having a Fitzpatrick skin type label for each ground-truth pedestrian. Standard object detection datasets do not contain this information; given the size of these datasets, we decided to enlist the help of Mechanical Turk workers to categorize each pedestrian in the dataset with the information about their skin tone.

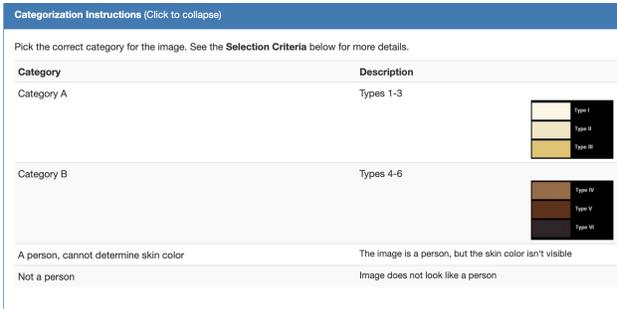


Figure 1. Instructions given to mechanical turk annotators for classifying  $\mathcal{L}\mathcal{S}$  and  $\mathcal{D}\mathcal{S}$  people.



Figure 2. Annotation interface.

##### 4.1.1. A TOOL FOR COLLECTING ANNOTATIONS

The vanilla BDD100K dataset lacks explicit labelings of people by skin color; each image instead is labeled by a set of bounding boxes along with a class label (the finest-grained class of pedestrians is the “person” class). The dataset therefore needed to be augmented with the Fitzpatrick skin type of each pedestrian in order to measure various models’ predictive inequity with respect to this categorization. We outline our approach to gathering these annotations below.

We initially cropped the bounding boxes of individuals that were labeled as the *person* class. We then created tasks from each bounding box, asking for each pedestrian to be classified into one of 4 categories: Fitzpatrick Categories 1-3, Fitzpatrick Categories 4-6, a person whose skin color cannot be determined, and not a person. The last two categories were included for multiple reasons: lighting, size, and occlusion can encumber determining the skin color of an individual with high confidence; moreover, BDD100K contains a small number of mislabeled instances. The instructions presented to Turkers in shown in Figure 1 and an example annotation interface is shown in Figure 2.

Initially, we intended to use the entire set of pedestrians labeled within BDD100K; however, we quickly found issue with this process. We began by manually annotating a small random subset of these pedestrians, but found that even for the same annotator, there were substantial inconsistencies of the labels provided on the same instance when annotations

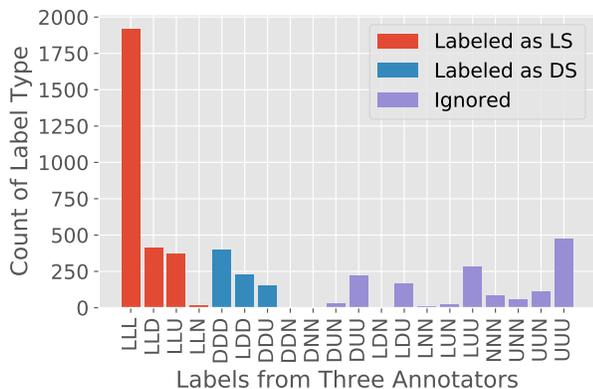


Figure 3. Histogram of the annotator responses. Each of the three annotators was given a choice of labeling as Category A (LS—denoted as L), Category B (DS—denoted as D), unknown (U), or not a person (N). Only instances with a consensus vote for LS or DS were labeled as such.

were collected on separate occasions. This was due to many of the factors stated previously, but we attributed much of the disagreement to the extremely small size of many of the cropped images. We found that this initial experiment had a large amount of disagreement, both amongst Turkers and compared to our labeling. We suspected this was due to the very small size of many of the cropped images.

We therefore focus our labeling and subsequent analysis on a filtered set of cropped pedestrian images which contains only those individuals whose bounding box area was greater than or equal to 10,000 pixels, hoping to see an increased agreement among Turkers. With this cutoff, there were only 487 pedestrians in the validation class, so we labeled those ourselves. The training set then had 4979 pedestrians above this cutoff, and for those we created Turker tasks. Due to the ongoing BDD100K challenge, *person* bounding boxes for the testing set were not present; therefore, we did not acquire labels corresponding to skin color for the test dataset. Imposing this minimum size constraint drastically increased the ease with which we were able to hand label the validation set, resulting in 3513 training images with a consensus of LS or DS labels, and 487 in validation.

For each person instance, we collected annotations from three separate Turkers. The set of possible labels we provided them with were: LS, if the person could be clearly identified as from LS; DS if the person could be clearly identified as from DS; U, if the skin color is unknown or too difficult to reliably identify; and N, if the box does not contain a person. A histogram of the received three score annotations received is shown in Figure 3 and shows that the majority of identifiable people fall into LS.

For each bounding box, if at least two of the three Turk-

Subset	LS	DS
<b>Train</b>	2724	789
<b>Validation</b>	387	100

Table 1. Count of labeled instances of people from LS and DS in BDD100K train and validation sets.

ers agreed on a label, we used that label; otherwise, we discarded that image for the purposes of evaluating predictive inequity. The summary of the number of consensus labelings can be found in Figure 3.

#### 4.1.2. AVOIDING ANNOTATOR BIAS

Having humans annotate images of other humans based on skin tone opens up these annotations to scrutiny: will the labels be skewed based on cultural biases? Will they be accurate? How would we tell if the labels were skewed or inaccurate? For this reason and others, we inspected whether the distribution over categories LS versus DS from our hand labeling was similar to that of the distribution from the consensus labeling given by the aggregated Mechanical Turk annotations. We found that the rate of LS to DS was similar for both (22% from DS based on Turker consensus on the training dataset and 20.5% based on our hand labeling of the validation set). This suggests that the consensus labels from MTurk might have similar behavior to our hand labeling, and while this does not rule out bias that both systems share it does suggest some degree of precision between these methods.

In future work, we plan to do further validation between our hand labelings and MTurk labelings, and also consider different aggregation schemes for the Turker labels (Should the labels without universal agreement be labeled by more Turkers? Do our findings hold similarly for images with universal agreement on the skin type as they do on those with just 2 agreements?).

#### 4.1.3. A BENCHMARK FOR STUDYING PREDICTIVE INEQUITY

The Berkeley Driving Dataset (Yu et al., 2018) is one of the most comprehensive driving datasets to date. The dataset is comprised of both bounding box level and segmentation level annotations. The dataset includes 40 different classes common to driving scenes, and the images were taken from 4 different locations: New York, Berkeley, Downtown San Francisco, and the Bay Area near San Jose. Additionally, the dataset includes diverse weather conditions (rain, snow, sunshine), as well as images from various parts of the day (morning, dusk, and nighttime). The split for the bounding box level annotations is 70,000 images in the training set, 10,000 in the validation set, and 20,000 images in the test

set, while the segmentation level annotations include 7,000 images in the training set, 2,000 images in the validation set, and 1,000 images in the test set. BDD100K is intended to be used to train real-time vision models that are included in current state of the art autonomous driving systems. Due to its relatively large size, we decided to use this dataset for experimentation.

## 4.2. Evaluating Predictive Inequity of Standard Object Detection Systems

In this section, we test whether several models displayed higher predictive inequity for pedestrians of Fitzpatrick types 4-6 as compared to those of types 1-3. We begin by defining our evaluation setup as well as the metric used for quantitative evaluation. We evaluate object detection models on the task of recognizing people in the BDD100K validation set using the average precision metric.

**Metric: Average Precision** A main metric used to quantify predictive performance of an object detection model is Average Precision (AP). For a given ground truth box,  $b_{gt}$ , and a predicted box,  $b_{pred}$ , the intersection over union of the predicted box is defined as the area of the intersection divided by the area of the union of the two boxes:

$$\text{IoU} = \frac{b_{gt} \cap b_{pred}}{b_{gt} \cup b_{pred}}$$

A predicted box is considered a true positive if it has IoU greater than a given threshold,  $T$ . Earlier challenges, such as Pascal (Everingham et al., 2010), focused on  $\text{AP}_{50}$ , which corresponds to a threshold,  $T = 0.5$ . Current challenges, such as MS COCO (Lin et al., 2014), evaluate multiple metrics including their overall AP score which consists of averaging across thresholds in the range of  $T \in [0.5, .95]$  by increments of 0.05. Finally, the BDD100K challenge (Yu et al., 2018) from which we derive our benchmark focuses on a stricter localization goal of  $\text{AP}_{75}$ , corresponding to  $T = 0.75$ . In the following sections we report performance for each of the three forms of evaluating average precision mentioned here.

**Data Splits** All of our results are reported on BDD100K validation images. When performance is measured for skin tones from LS and DS, we consider only the subset of the BDD100K validation images for which at least one person from either category is present. All people within an evaluation image which are too small to be reliably annotated into either LS or DS are ignored during per skin-type evaluation. We do this so as not to unnecessarily penalize the person detector for correctly predicting a small person.

**Statistical validity** We take a moment to describe how certain we are that the phenomena described below happen generally over a resampled validation set drawn from

the same distribution as the validation set on which we report our empirical findings. In short, standard holdout-style arguments for a loss function of  $k$  models evaluated on a holdout set of size  $n$  will have  $1 - \delta$  probability of being within  $\pm \sqrt{\frac{\ln \frac{k}{\delta}}{n}}$ . If we view our validation set as a holdout, and consider only those bounding boxes for pedestrians of size at least 10,000 pixels (for which we have Fitzpatrick annotations), we only have holdout set sizes  $n_{LS} = 387$ ,  $n_{DS} = 100$ . This means that our estimates on AP for category B for a single model are only accurate up to 0.17 for  $\delta = .05$ , and up to 0.08 for category A. So, if the true gap between AP for LS versus DS is .05, we would need  $n_{LS} = 12,000$ ,  $n_{DS} = 4,000$  to verify that LS’s AP surpassed that for DS (with probability .95 over the draw of the holdout set).

This suggests that gathering high-confidence comparisons between LS and DS will require much larger datasets (or, at least, many more examples of sufficiently large pedestrians for which we can gather Fitzpatrick annotations). We are not aware of any academic driving dataset with this many instances of large pedestrians, making it difficult without gathering a new dataset to validate our findings.

However, we note that this behavior persists when evaluating on the larger training set and the validation set, in a fairly wide variety of learning settings, for a variety of models, throughout the training process, which should give some degree of confidence that this behavior is not entirely spurious.

Moreover, we observed this behavior for a wide range of learning rates (what generally is referred to as a “hyperparameter” of a learning algorithm); this phenomenon was not the result of picking a perfect training schedule. We found consistently that models exhibited between 4 and 10% gaps in these precision metrics (with LS consistently outperforming DS). For all models we train, we run each 10 times and report the mean and standard deviation in the resulting tables. **We emphasize that the standard deviation listed in our experiments is solely a function of randomness in the training procedure and does not provide confidence intervals with respect to sampling error.**

### 4.2.1. TRAINING DATA IMPORTANCE

We first compare performance for a state-of-the-art object detection model, Faster R-CNN (Ren et al., 2015) model with a R-50-FPN (He et al., 2015) backbone, trained using two different sources of annotated data. We consider the standard learning protocol of initializing with ImageNet weights and then further training on a detection dataset.

The MS COCO dataset (Lin et al., 2014) is one of the predominant datasets currently used to evaluate an object detection model’s performance, consisting of 80 different classes,

## Predictive Inequity in Pedestrian Detection

Training Data	AP (%)	AP <sub>50</sub> (%)	AP <sub>75</sub> (%)
MS COCO	24.2	56.9	16.0
BDD100K Train	22.3 ± 2.0	45.4 ± 5.2	18.8 ± 1.0

Table 2. Average precision of the person class on BDD100K validation. Performance here is across all people.

Training Data	AP (%)		AP <sub>50</sub> (%)		AP <sub>75</sub> (%)	
	LS	DS	LS	DS	LS	DS
MS COCO	57.2	53.3	91.3	91.8	63.1	53.6
BDD100K Train	54.6 ± 0.4	51.8 ± 0.9	90.2 ± 0.9	89.9 ± 1.5	59.8 ± 1.0	53.5 ± 1.7

Table 3. Comparing average precision on BDD100K validation set across the larger people labeled as LS or DS for Faster-RCNN R-50-FPN trained either with MS COCO or BDD100K train data (averaged over 10 different trainings).

including a *person* class. Before studying predictive inequity, we begin by verifying the effectiveness of the person detector trained from MS COCO data for recognizing the pedestrians in the BDD100K validation dataset in Table 2. We compare this against the performance of the person detector trained using the BDD100K training set. We find that the MS COCO model outperforms the BDD100K model on the overall AP metric as well as the AP<sub>50</sub> metric, while the BDD100K model outperforms the MS COCO model on the strict localization evaluation of AP<sub>75</sub>. Thus we find that both sources of training data are relevant for assessing person detection on this validation set. Implementation details and hyperparameters are available in the appendix.

This leads us to our main question, do these models perform differently when evaluated on Fitzpatrick skin types [1-3] (LS) vs skin types [4-6] (DS)? To answer this question, we evaluate the same models as above as they perform at recognizing people identified within the BDD100K validation set (described in Section 4.1) annotated with LS or DS labels. Table 3 reports this breakdown. What we find is that consistently, across models trained with either MS COCO or BDD100K train, people from LS are recognized with higher average precision than people from DS. The largest disparity occurs when evaluating using standard BDD100K metric of AP<sub>75</sub>, which requires tight localization. We see a much higher level of predictive inequity using MS COCO weights compared to BDD100k weights on AP<sub>75</sub>, suggesting this problem is not unique to the BDD100k training data. We also verify that this discrepancy between AP on LS vs. DS is not an artifact of a particular point in training and instead persists across training iterations as shown in Figure 4.

**Remark 1** MS COCO contains a broad set of classes, such as “umbrella” and “suitcase”, which are generally not included within a “person” bounding box, while BDD100k’s “person” bounding box often includes these (making the ground truth annotations somewhat different).

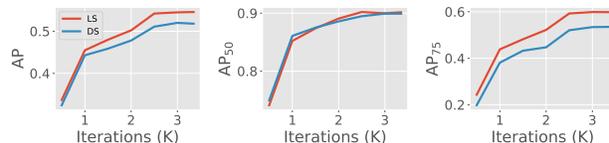


Figure 4. AP performance gap comparing LS and DS individuals using an unweighted model across training iterations on BDD100K. LS consistently has higher AP than DS people.

Model	Backbone	AP (%)		AP <sub>50</sub> (%)		AP <sub>75</sub> (%)	
		LS	DS	LS	DS	LS	DS
Faster R-CNN	R-50-C4	57.3	52.7	90.3	90.0	64.5	53.6
Faster R-CNN	R-50-FPN	57.2	53.3	91.3	91.8	63.1	53.6
Faster R-CNN	R-101-FPN	59.7	56.9	91.6	93.0	70.0	57.3
Faster R-CNN	X-101-32x8d-FPN	60.6	55.4	93.1	90.0	69.8	62.1
Mask R-CNN	R-50-C4	59.3	53.8	91.3	90.5	66.9	51.5
Mask R-CNN	R-50-FPN	58.9	53.2	92.5	92.3	67.6	51.2
Mask R-CNN	R-101-FPN	60.1	54.7	92.4	92.1	65.6	55.2
Mask R-CNN	X-101-32x8d-FPN	60.8	57.1	93.3	91.6	69.2	62.9
Average		59.2	54.6	92.0	91.4	67.1	55.9

Table 4. Performances of the person class on BDD100K validation set for models trained using MS COCO. We note that these models were not trained on BDD100K.

### 4.2.2. IMPORTANCE OF ARCHITECTURE AND MODEL SELECTION

Until now, we have shown evidence of predictive inequity between people of skin tones in the Fitzpatrick range [1-3] (LS) as compared to range [4-6] (DS) for a single object detection model learned using two different sources of data. A reasonable question to ask is whether or not this observation is an artifact of that particular architecture or model. Therefore we study prediction of people across the two skin-type categories across multiple architecture backbones and two state-of-the-art models, Faster R-CNN (Ren et al., 2015) and Mask R-CNN (He et al., 2017). We evaluate the publicly available model weights from training on the MS COCO dataset, released within the Detectron model zoo (Girshick et al., 2018) and report performance in Table 4.

We find that across all models and base architectures studied, performance on LS exceeds that of DS, demonstrating that this phenomenon is not specific to a particular model.

Again, the most striking measure of predictive inequity is observed under the strict localization metric of AP<sub>75</sub>, where the average performance across models studied drops from 67.1% for LS to 55.9% for DS. However, under the weaker localization metric of AP<sub>50</sub>, the gap between the two groupings of people is greatly diminished. In the next section we provide further analysis into the results shown here to explain the discrepancy between the two metrics.

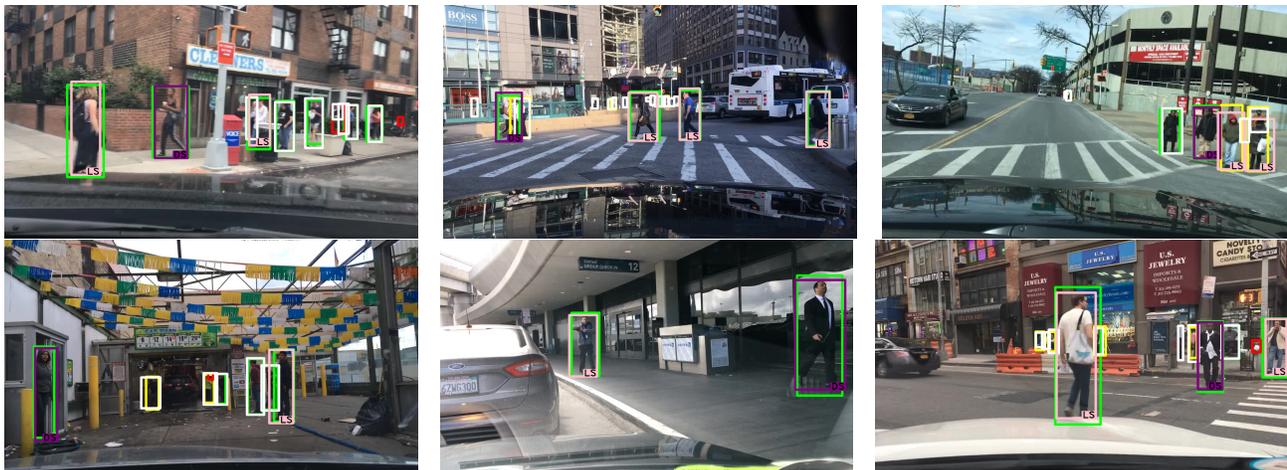


Figure 5. Example detections from Faster R-CNN using the R-50-FPN backbone, trained on BDD100K. For reference, the ground truth annotations for LS and DS are pink and purple respectively. Yellow boxes correspond to true positives under the  $AP_{50}$  metric and false positives under the  $AP_{75}$  metric. Green boxes correspond to true positives under the  $AP_{75}$  metric. All the predictions shown are greater than an 85% confidence threshold.

Model	Backbone	AP (%)		AP <sub>50</sub> (%)		AP <sub>75</sub> (%)	
		LS	DS	LS	DS	LS	DS
Faster R-CNN	R-50-C4	62.0	57.1	93.5	94.1	72.9	60.7
Faster R-CNN	R-50-FPN	61.2	58.0	93.5	95.8	70.9	59.4
Faster R-CNN	R-101-FPN	64.2	60.8	95.7	95.3	77.0	63.2
Faster R-CNN	X-101-32x8d-FPN	64.7	59.4	95.7	92.8	75.9	68.9
Mask R-CNN	R-50-C4	63.4	58.2	94.0	95.0	74.1	58.7
Mask R-CNN	R-50-FPN	63.2	56.3	95.3	94.7	75.6	56.1
Mask R-CNN	R-101-FPN	64.3	58.5	95.3	95.5	73.5	59.1
Mask R-CNN	X-101-32x8d-FPN	64.9	61.4	96.1	95.0	76.6	70.5
Average		63.5	58.7	94.9	94.8	74.6	62.1

Table 5. Average precision on BDD100K validation set with occluded individuals removed for models trained using MS COCO.

### 4.3. Analyzing Sources of Predictive Inequity

In the previous section we demonstrated that using the tight localization metric of  $AP_{75}$  a wide variety of object detection models and training schemes results in predictive inequity between individuals with skin tones from LS and DS. Here we investigate the natural followup question: what causes this discrepancy?

We begin our study by inspecting example output detections in Figure 5. We show here ground truth (white box) with LS or DS indicated where appropriate for people appearing in each image. We then show all boxes with scores greater than a threshold (0.85). We indicate those boxes which are false positives under all metrics in red. Those which have overlap  $\geq 0.5$  with the ground truth, but  $< 0.75$  and are thus false positives under the  $AP_{75}$  metric are indicated in yellow. Finally, those which are true positives under all metrics are shown in green.

By studying such images, we observe potential sources of predictive inequity which we will analyze next in more de-

tail; namely, occlusion, time of day, and loss prioritization.

#### 4.3.1. OCCLUSION AS A POSSIBLE SOURCE OF PREDICTIVE INEQUITY

One observation is that often street scenes contain multiple people on the sidewalk together or crossing the road near each other leading to partially occluded individuals. Recognition under occlusion is known to be more challenging (see Figure 5 top right and if the dataset has a bias by which one type of individual occurs more frequently in crowded scenes with occlusion, it would be unsurprising if that type of individual suffered lower performance in overall recognition.

To remove this confounding factor we study the performance of a reduced ground truth set which contains only the unoccluded people from BDD100K validation for which we have discernible skin type labels. Table 5 reports the three average precision metrics and results in consistent observations with those made in the full validation set which included the occluded individuals. We observe that AP,  $AP_{50}$ , and  $AP_{75}$  all appear to improve for both LS and DS once these occluded examples are removed, but that the gap between LS and DS performance for AP and  $AP_{75}$  remains. Therefore, we conclude that the source of discrepancy between performance on LS and DS is not due to co-occurrence with occluded people.<sup>2</sup>

<sup>2</sup>We remark that considering performance broken down by occlusion was actually motivated by unexpected performance of these MSCOCO weighted models on BDD100K training data; see section 7.1 for further discussion.

### 4.3.2. TIME OF DAY AS A POSSIBLE SOURCE OF PREDICTIVE INEQUITY

Low contrast between a subject and its background may affect the recognition performance of that subject. In outdoor street scenes the natural lighting variations that occur at various times of day will change the contrast of the person itself and between the person and the background, potentially introducing a confounding factor when comparing predictive performance between LS and DS. Therefore, we next measure whether the observed gap in predictive performance is attributable to time of day by reporting performance during daytime or nighttime hours, annotations which are available within BDD100K.

Table 7 reports the recognition performance for individuals from LS and DS who appear in day versus night. We used Faster R-CNN with R-50-FPN backbone trained on BDD100K *train* with  $\alpha_{DS} = 1$  for evaluation. Focusing only on daytime images, for each of AP, AP<sub>50</sub>, and AP<sub>75</sub>, performance of our trained model on the LS category was higher than that on the DS category, and the size of the gap in each case was quite similar to that on the entire validation set. Surprisingly, for nighttime images, the performance of our model on DS pedestrians was higher than on LS pedestrians. However, we note that this may be the result of small nighttime sample size (the validation set has only 144 nighttime pedestrians of sufficiently large size, only 15 of which are designated as DS). Nonetheless, this data suggests an important conclusion that *there is not evidence that time of day is to blame for the predictive inequity observed overall*.

### 4.3.3. LOSS PRIORITIZATION AS A POSSIBLE SOURCE OF PREDICTIVE INEQUITY

We know that there are more than three times as many individuals from LS as from DS in the BDD100K dataset (see Table 1). In the next section we take steps to reduce this higher degree of representation for LS in the training set through supervised loss weightings. This will give some indication as to what fraction of the observed predictive inequity stems from prioritizing loss on pedestrians from LS over those from DS, and whether a simple solution (namely, reweighting the loss function) can effectively treat some fraction of this observed predictive inequity.

## 5. Reducing Predictive Inequity

We now describe our investigation into whether the importance of LS’s loss compared to that of DS’s loss during training is a source of the predictive inequity between the classes. Many loss functions (including those regularly used to train the models studied in this work) decompose into terms for each training example, weighting each example uniformly. This tacit weighting implies that a subset of the

$\alpha_{DS}$	AP (%)		AP <sub>50</sub> (%)		AP <sub>75</sub> (%)	
	LS	DS	LS	DS	LS	DS
1	<b>54.6 ± 0.4</b>	51.8 ± 0.9	<b>90.2 ± 0.9</b>	89.9 ± 1.5	<b>59.8 ± 1.0</b>	53.5 ± 1.7
2	<b>55.3 ± 0.7</b>	53.3 ± 0.6	<b>90.9 ± 0.6</b>	90.4 ± 1.1	<b>60.8 ± 1.7</b>	55.7 ± 3.3
3	<b>55.8 ± 0.6</b>	53.9 ± 1.3	<b>91.3 ± 0.7</b>	90.6 ± 1.0	<b>63.4 ± 1.8</b>	58.3 ± 2.9
5	<b>56.4 ± 0.4</b>	53.9 ± 1.1	<b>91.8 ± 0.5</b>	90.8 ± 1.0	<b>63.0 ± 1.9</b>	57.3 ± 2.3
10	<b>55.8 ± 0.8</b>	54.0 ± 1.0	<b>91.7 ± 0.9</b>	91.0 ± 0.8	<b>62.1 ± 1.5</b>	56.3 ± 2.0

Table 6. Performances of Faster-RCNN using R-50-FPN backbone on BDD100k *validation* with different weightings on DS in the classification network loss function.

training data which is 3.5 times larger than another subset (as LS is compared to DS) will have up to 3.5 times as much impact on such a loss function, possibly steering an algorithm towards models which have lower loss on LS compared to DS. We emphasize here that such issues are not information-theoretic: such issues stem from the design choice of loss function rather than having too little data to give generalizable predictions for the smaller subset.

For the purpose of this discussion we will describe our methodology in terms of a generic loss function (see Appendix for the longer and more precise discussion specific to Faster R-CNN),  $\mathcal{L}(x, y)$ , which takes as input an image crop,  $x$ , and the true label,  $y$ , which in our case will be either person, some other class, or background. Let us indicate the set of people instances which are also labeled as LS or DS as  $X_{LS}$  and  $X_{DS}$  respectively, with the number of instances denoted as  $N_{LS}$  and  $N_{DS}$  respectively. Let all other boxes not containing a person from LS or DS be denoted as  $X_O$ . Then we can define our overall loss in terms of a weighted sum over each instance:

$$\begin{aligned} \text{Total Loss} &= \frac{\alpha_{LS}}{N_{LS}} \sum_{\{x_{LS}, y_{LS}\} \in X_{LS}} \mathcal{L}(x_{LS}, y_{LS}) \\ &+ \frac{\alpha_{DS}}{N_{DS}} \sum_{\{x_{DS}, y_{DS}\} \in X_{DS}} \mathcal{L}(x_{DS}, y_{DS}) \\ &+ \frac{\alpha_O}{N_O} \sum_{\{x_o, y_o\} \in X_O} \mathcal{L}(x_o, y_o) \end{aligned} \quad (2)$$

where  $\alpha_{LS}$ ,  $\alpha_{DS}$ ,  $\alpha_O$  denotes the per class weighting on instances from LS, DS, or other, respectively.

To measure whether our loss function emphasizes accuracy for LS compared to DS and in so doing creates some predictive inequity, we consider several reweightings of the standard loss function used in training Faster R-CNN. The (standard) unweighted or equal weighted loss function for Faster R-CNN has two components, one aimed at affecting the region proposal network and the other for the detection and classification inside these proposed regions. We only consider weightings which affect the latter part of the loss function.

The detection and classification component of the Faster R-CNN loss function combines a cross-entropy term  $\mathcal{L}_{cls}$  and a regularized  $\ell_1$  loss  $\mathcal{L}_{reg}$  term for each example. For the purpose of reweighting, we consider the loss to be the joint cross-entropy and regularized terms and consider variants

Time	AP (%)		AP <sub>50</sub> (%)		AP <sub>75</sub> (%)	
	LS	DS	LS	DS	LS	DS
Day	57.2 ± 0.5	54.6 ± 1.2	91.8 ± 0.8	91.0 ± 1.7	63.9 ± 1.6	58.4 ± 2.2
Night	43.3 ± 0.9	50.8 ± 1.6	81.4 ± 2.4	91.7 ± 3.7	41.7 ± 3.0	53.2 ± 5.0

Table 7. Performances at different times of day using Faster R-CNN with R-50-FPN backbone on BDD100K validation using weights trained on BDD100K train with  $\alpha_{DS} = 1$ . Note the number of daytime and nighttime examples are as follows: Day: LS 297 DS 75, Night: LS 69 DS 15.

of values for  $\alpha_{LS}$ ,  $\alpha_{DS}$ , and  $\alpha_O$ . To reduce the number of hyperparameters to optimize we fix both  $\alpha_{DS}$  and  $\alpha_O$  to be 1 for all experiments and consider the effect of placing greater weight on DS instances through raising  $\alpha_{DS}$ .

Average precision as we vary  $\alpha_{DS}$  weight on instances from DS are reported in Table 6. We find that the gap between the total AP value is reduced for larger values of  $\alpha_{DS}$ , but that the LS to DS AP<sub>75</sub> gap is quite similar for each of the weightings, in the range of 4 to 6%. For some of these weightings, in particular for  $\alpha_{DS} = 3$ , we notice that the AP<sub>75</sub> for DS pedestrians is quite close to the AP<sub>75</sub> for LS pedestrians trained on unweighted examples. Moreover, the performance of the model on LS pedestrians is better under this weighting.

This result suggests that careful reweighting may improve performance on DS pedestrians without sacrificing performance on LS pedestrians, or even improve performance on LS pedestrians as a byproduct. Further analysis is needed to fully understand the effect of re-weighting on the stricter criteria used as part of the total AP metric to fully reduce the inequity. Overall, this finding suggests that predictive inequity stemming from the oft-labeled “too little data” source might actually confound two fundamentally different phenomena: that smaller datasets beget less statistical certainty, but also tend to receive lower emphasis during training.

## 6. Conclusion and Discussion

In this work, we propose the concept of predictive inequity in detecting pedestrians of different skin tones in object detection systems. We give evidence that standard models for the task of object detection, trained on standard datasets, appear to exhibit higher precision on lower Fitzpatrick skin types than higher skin types. This behavior appears on large images of pedestrians, and even grows when we remove occluded pedestrians. Both of these cases (small pedestrians and occluded pedestrians) are known difficult cases for object detectors, so even on the relatively “easy” subset of pedestrian examples, we observe this predictive inequity. We have shown that simple changes during learning (namely, reweighting the terms in the loss function) can partially mitigate this disparity. We hope this study provides compelling

evidence of the real problem that may arise if this source of capture bias is not considered before deploying these sort of recognition models.

## References

- Angwin, J., Larson, J., Kirchner, L., and Mattu, S. Machine bias. *ProPublica*, May 2016. URL <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- Buolamwini, J. and Gebru, T. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Fairness Accountability and Transparency (FAT)*, 2018.
- Chouldechova, A. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163, 2017.
- Chouldechova, A., Benavides-Prado, D., Fialko, O., and Vaithianathan, R. A case study of algorithm-assisted decision making in child maltreatment hotline screening decisions. In *Conference on Fairness, Accountability and Transparency*, pp. 134–148, 2018.
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., and Schiele, B. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Datta, A., Datta, A., Makagon, J., Mulligan, D. K., and Tschantz, M. C. Discrimination in online advertising: A multidisciplinary inquiry. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, *PMLR*, volume 81, 2018.
- Dollár, P., Wojek, C., Schiele, B., and Perona, P. Pedestrian detection: An evaluation of the state of the art. *PAMI*, 34, 2012.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226. ACM, 2012.
- Ess, A., Leibe, B., Schindler, K., , and van Gool, L. A mobile vision system for robust multi-person tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR’08)*. IEEE Press, June 2008.
- Everingham, M., Van Gool, L., Williams, C. K. I., Winn, J., and Zisserman, A. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88 (2):303–338, June 2010.

- Fitzpatrick, T. B. Soleil et peau. *J Med Esthet*, 2:33–34, 1975.
- Geiger, A., Lenz, P., and Urtasun, R. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- Girshick, R., Radosavovic, I., Gkioxari, G., Dollár, P., and He, K. Detectron. <https://github.com/facebookresearch/detectron>, 2018.
- Hardt, M., Price, E., Srebro, N., et al. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pp. 3315–3323, 2016.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- He, K., Gkioxari, G., Dollár, P., and Girshick, R. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pp. 2980–2988. IEEE, 2017.
- Klare, B. F., Burge, M. J., Klontz, J. C., Bruegge, R. W. V., and Jain, A. K. Face recognition performance: Role of demographic information. *IEEE Transactions on Information Forensics and Security*, 7(6):1789–1801, 2012.
- Kleinberg, J., Mullainathan, S., and Raghavan, M. Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*, 2016.
- Kriegler, B. and Berk, R. Small area estimation of the homeless in los angeles: An application of cost-sensitive stochastic gradient boosting. *The Annals of Applied Statistics*, pp. 1234–1255, 2010.
- Levin, S. and Wong, J. C. Self-driving uber kills arizona woman in first fatal crash involving pedestrian. <https://www.theguardian.com/technology/2018/mar/19/uber-self-driving-car-kills-woman-arizona-tempe>, Mar 2018.
- Li, J., Liang, X., Shen, S., Xu, T., Feng, J., and Yan, S. Scale-aware fast r-cnn for pedestrian detection. *IEEE Transactions on Multimedia*, 20(4):985–996, April 2018. ISSN 1520-9210. doi: 10.1109/TMM.2017.2759508.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.
- Lum, K. and Isaac, W. To predict and serve? *Significance*, 13(5):14–19, 2016.
- Maloof, M. A. A machine learning researchers foray into recidivism prediction. *Technical Report*, 1999.
- Massa, F. and Girshick, R. maskrcnn-benchmark: Fast, modular reference implementation of Instance Segmentation and Object Detection algorithms in PyTorch. <https://github.com/facebookresearch/maskrcnn-benchmark>, 2018. Accessed: [Insert date here].
- Ren, S., He, K., Girshick, R., and Sun, J. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pp. 91–99, 2015.
- Roff, H. M. The folly of trolleys: Ethical challenges and autonomous vehicles, Dec 2018. URL <https://www.brookings.edu/research/the-folly-of-trolleys-ethical-challenges-and-autonomous-vehicles/> amp/?fbclid=IwAR2efjQ8-j6FT0XUN7g\_mZaZJBXpTc5RhrHQAFevHd24hnoMm-QD\_VvjdzI.
- Rubin, J. Stopping crime before it starts. *Los Angeles Times*, 21, 2010.
- Selbst, A. D. Disparate impact in big data policing. *Ga. L. Rev.*, 52:109, 2017.
- Snow, J. Amazon’s face recognition falsely matched 28 members of congress with mugshots, 2018. URL <https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazons-face-recognition-falsely-matched-28>.
- Sweeney, L. Discrimination in online ad delivery. *Queue*, 11(3):10, 2013.
- White, G. B. A newly minted macarthur genius on the financially ‘invisible’, Sep 2016. URL <https://www.theatlantic.com/business/archive/2016/09/quinonez-macarthur-genius-mission-asset-fund/501512/>.
- Yu, F., Xian, W., Chen, Y., Liu, F., Liao, M., Madhavan, V., and Darrell, T. Bdd100k: A diverse driving video database with scalable annotation tooling. *arXiv preprint arXiv:1805.04687*, 2018.
- Zhang, L., Lin, L., Liang, X., and He, K. Is faster r-cnn doing well for pedestrian detection? In Leibe, B., Matas, J., Sebe, N., and Welling, M. (eds.), *Computer Vision – ECCV 2016*, pp. 443–457, Cham, 2016. Springer International Publishing.
- Zhang, S., Yang, J., and Schiele, B. Occluded pedestrian detection through guided attention in cnns. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

Zhao, J., Wang, T., Yatskar, M., Ordonez, V., and Chang, K.-W. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017.

Model	Backbone	AP (%)		AP <sub>50</sub> (%)		AP <sub>75</sub> (%)	
		LS	DS	LS	DS	LS	DS
Faster R-CNN	R-50-C4	54.6	<b>55.9</b>	89.4	<b>91.6</b>	60.2	<b>61.1</b>
Faster R-CNN	R-50-FPN	53.9	<b>54.8</b>	90.0	<b>92.5</b>	<b>58.5</b>	56.9
Faster R-CNN	R-101-FPN	56.7	<b>58.4</b>	90.9	<b>93.5</b>	62.7	<b>64.0</b>
Faster R-CNN	X-101-32x8d-FPN	57.5	<b>59.1</b>	91.5	<b>94.4</b>	<b>63.6</b>	63.1
Mask R-CNN	R-50-C4	56.3	<b>57.5</b>	90.4	<b>92.5</b>	<b>63.0</b>	62.2
Mask R-CNN	R-50-FPN	55.8	<b>56.8</b>	91.0	<b>93.4</b>	<b>61.1</b>	60.7
Mask R-CNN	R-101-FPN	57.0	<b>58.9</b>	91.1	<b>94.0</b>	62.4	<b>64.5</b>
Mask R-CNN	X-101-32x8d-FPN	57.8	<b>59.2</b>	91.6	<b>94.0</b>	64.6	<b>64.9</b>
Average		56.2	<b>57.6</b>	90.7	<b>93.2</b>	62.0	<b>62.2</b>

Table 8. Average precision on BDD100K train set.

## 7. Appendix

### 7.1. MSCOCO weighted models and occlusion

We also evaluated the MSCOCO weighted models on the BDD100k training set, primarily for the following reason. Because these model were not trained on BDD100K, one can make stronger statistical claims about results of these models evaluated on this larger set (treating it as validation). Interestingly, we found the gap we observed on the validation set did not seem to persist on the larger training set—this presumably meant that either the results on the validation set were the result of sampling error, or there was some large statistical difference between the train/val sets for BDD100k. Upon further inspection, we discovered that occluded pedestrians affected these results significantly. We show the results of these experiments below in Tables 8 and 9.

### 7.2. Loss Function for Faster R-CNN and our Reweighting

Faster R-CNN utilizes a network which suggests regions where an object is likely to be, known as a Region Proposal Network (RPN). These proposed regions serve as an input to a separate detection network, which is then used for both classification and bounding box regression.

Faster R-CNN has two different objectives, which correspond to the RPN and the separate objection detection network; however, we will solely focus on the latter. Faster R-CNN utilizes a set of predefined boxes of different shapes and sizes, defined as *anchors*, to serve as priors for detecting objects. Within each proposed region, each anchor produces both a probability distribution over the set of candidate classes and a set of regression offsets for the bounding boxes.

The full object detection loss can be written as:

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{\text{cls}}} \sum_i L_{\text{cls}}(p_i, p_i^*) + \frac{\lambda}{N_{\text{cls}}} \sum_i p_i^* L_{\text{reg}}(t_i, t_i^*) \quad (3)$$

where  $p_i$  is the predicted probability distribution corresponding to a positive anchor, and  $p_i^*$  corresponds to the one-hot

vector corresponding to the ground truth class.  $t_i$  and  $t_i^*$  represent the set of values that parameterize the predicted bounding box and the ground truth bounding box with respect to the positive anchor as such:

$$\begin{aligned} t_i^x &= \frac{x - x_a}{w_a} & t_i^{*x} &= \frac{x^* - x_a}{w_a} \\ t_i^y &= \frac{y - y_a}{h_a} & t_i^{*y} &= \frac{y^* - y_a}{h_a} \\ t_i^w &= \log\left(\frac{w}{w_a}\right) & t_i^{*w} &= \log\left(\frac{w^*}{w_a}\right) \\ t_i^h &= \log\left(\frac{h}{h_a}\right) & t_i^{*h} &= \log\left(\frac{h^*}{h_a}\right) \end{aligned}$$

$\{x, y\}, \{x^*, y^*\}, \{x_a, y_a\}$  correspond to the centers of the predicted, ground truth, and anchor boxes, and  $\{w, h\}, \{w^*, h^*\}, \{w_a, h_a\}$  correspond to the width and height of these boxes.

$L_{\text{cls}}$  is the Log Loss between the set of  $k$  classes specified, defined as:

$$L_{\text{cls}}(p_i, p_i^*) = - \sum_k p_i^* \log p_i \quad (4)$$

and  $L_{\text{reg}}$  is the smooth  $L_1$  loss between the parameterization of the predicted and ground truth boxes. This is written as:

$$L_{\text{reg}}(t_i, t_i^*) = \sum_{c \in \{x, y, w, h\}} L_{1\text{-smooth}}(t_i - t_i^*) \quad (5)$$

where

$$L_{1\text{-smooth}}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

Both  $N_{\text{cls}}$  and  $N_{\text{reg}}$  are normalization parameters, corresponding to the size of sampled anchors and the number of anchor locations respectively.  $\lambda$  is a balancing parameter for the objective. We use the default parameters in the implementation we used for training.

**Augmented Loss** Given a set of attributes  $\{\text{LS}, \text{DS}, \text{Not a person}, \text{A person (cannot determine skin color)}\}$  for members of the *person* class, we would like to weight the objective based off attribute membership. We can reparameterize our functions, introducing a weight vector,  $\mathcal{W} \in \mathbb{R}^4$ :

$$L(\{p_i\}, \{t_i\}, \{a_i\}) = - \frac{1}{N_{\text{cls}}} \sum_i L_{\text{cls}}(p_i, p_i^*, \mathcal{W}_{a_i}) + \frac{\lambda}{N_{\text{reg}}} \sum_i p_i^* L_{\text{reg}}(t_i, t_i^*, \mathcal{W}_{a_i}) \quad (6)$$

where  $a_i$  represents the index of the corresponding attribute for a given instance. Therefore,  $L_{\text{cls}}$  would become:

$$L_{\text{cls}}(p_i, p_i^*, \mathcal{W}_{a_i}) = - \mathcal{W}_{a_i} \sum_k p_i^* \log(p_i)$$

Predictive Inequity in Pedestrian Detection

Model	Backbone	AP (%)		AP <sub>50</sub> (%)		AP <sub>75</sub> (%)	
		LS	DS	LS	DS	LS	DS
Faster R-CNN	R-50-C4	<b>58.9</b>	58.1	93.4	<b>93.5</b>	<b>66.6</b>	64.7
Faster R-CNN	R-50-FPN	<b>58.0</b>	56.8	93.7	<b>94.2</b>	<b>65.1</b>	60.1
Faster R-CNN	R-101-FPN	<b>61.3</b>	60.2	94.8	<b>95.1</b>	<b>69.4</b>	66.1
Faster R-CNN	X-101-32x8d-FPN	<b>62.0</b>	61.2	95.4	<b>96.1</b>	<b>70.1</b>	66.0
Mask R-CNN	R-50-C4	<b>60.6</b>	59.5	94.0	<b>94.6</b>	<b>69.8</b>	65.3
Mask R-CNN	R-50-FPN	<b>60.1</b>	58.8	94.6	<b>95.1</b>	<b>67.9</b>	63.4
Mask R-CNN	R-101-FPN	<b>61.5</b>	60.6	95.5	<b>95.9</b>	<b>69.0</b>	66.9
Mask R-CNN	X-101-32x8d-FPN	<b>62.4</b>	61.5	95.6	<b>96.3</b>	<b>71.2</b>	68.3
Average		<b>60.6</b>	59.6	<b>94.6</b>	95.1	<b>68.6</b>	65.1

Table 9. Average precision on BDD100K *train* set with occluded individuals removed for models trained using MS COCO. There are 1855 and 570 individuals labeled as LS and DS respectively.

and,

$$L_{\text{reg}}(t_i, t_i^*, \mathcal{W}_{a_i}) = \sum_{c \in \{x, y, w, h\}} \mathcal{W}_{a_i} L_{1, \text{smooth}}(t_i^c - t_i^{*c})$$

### 7.3. Implementation Details

We used a PyTorch implementation of Faster R-CNN (Massa & Girshick, 2018) for all of the experiments listed. All training was done on 1 NVIDIA V100. All default parameters from the PyTorch implementation were used, except the learning rate, batch size, and the step learning schedule were modified to suit the new dataset. We used a learning rate of 0.01 with a mini batch size of 8 images. Additionally, we used a step learning schedule that decays at 0.1 at iterations 2,233 and 2,792. Each of the BDD100K experiments were run for a total of 3,350 iterations.