

# Unsupervised Difficulty Estimation with Action Scores

Octavio Arriaga<sup>1</sup> (arriagac@uni-bremen.de) and Matias Valdenegro-Toro<sup>2</sup> (matias.valdenegro@dfki.de)

<sup>1</sup>University of Bremen. <sup>2</sup>German Research Center for Artificial Intelligence. Bremen, Germany



Deutsches  
Forschungszentrum  
für Künstliche  
Intelligenz GmbH



Universität Bremen

## Summary

- **Motivation:** It is well known that for ML models, classifying certain examples might be more difficult, due to variations in object pose, noise, object size, etc. This is important in determining biases in models and datasets, and a method that does not require changes in a model is preferable. We hypothesize that the per-sample loss during training contains information about difficulty.
- **Approach:** We define an action score as the accumulation of per-sample loss values across epochs during training. Our results show that visually difficult examples obtain a higher action scores, while easier examples obtain low scores. This confirms our hypothesis.
- **Contributions:** We propose the action score for difficulty estimation without any kind of supervision. Our approach does not require changes to the model, and allows to extract more information from the training process. We evaluate our method in CIFAR10 classification and SSD object detection on PASCAL VOC, which show that our method can be applied to simple and complex multi-task settings.

## Action Scores

Given a loss function  $\mathcal{L}$  and a model  $m$  with free parameters  $\theta_n$ , we define the action  $\mathcal{A}$  of a sample  $x \in \mathcal{X}$  with labels  $y \in \mathcal{Y}$  as

$$\mathcal{A}(x) = \sum_{n=0}^N \mathcal{L}(y, m(x; \theta_n))$$

where  $n$  represents epochs. Consequently, the action of a sample is the accumulated loss over all epochs. Our method characterizes the action of each sample as a measurement of its difficulty. Therefore, samples with a higher accumulated loss represent samples that are more difficult to learn.

Specifically, we argue that the action is directly proportional to its difficulty i.e.

$$D(x) \propto \mathcal{A}(x)$$

Note that action can be computed only during training, and it is available for training and validation sets.

## Experimental Evaluation and Analysis

- We evaluate on CIFAR10 classification, and SSD object detection on the PASCAL VOC dataset. For SSD, the loss can be decomposed into localization loss (bounding box regression), positive classification loss, and negative classification loss (for background windows). We compute an action score for each loss component separately. We group hard and easy examples by the top and bottom action scores.
- On CIFAR10, the easiest examples all contain objects in a canonical pose that is easy to classify, while the hardest examples have easy to confuse classes, difficult backgrounds, or incorrect labels.
- On PASCAL VOC object detection, we see similar issues, with the localization loss indicating easy examples to ones in easy poses, while the hardest examples contain small objects. Classification losses also show struggles with very small objects or cluttered scenes.

## Conclusion & Future Work

- Our results show that the action score correlates with visual difficulty of training samples, which is useful to debug models and datasets.
- Future work will consider evaluating on more tasks (segmentation, pose estimation, etc), and evaluating the robustness of the action score to variations in the training process.

## CIFAR10 Classification

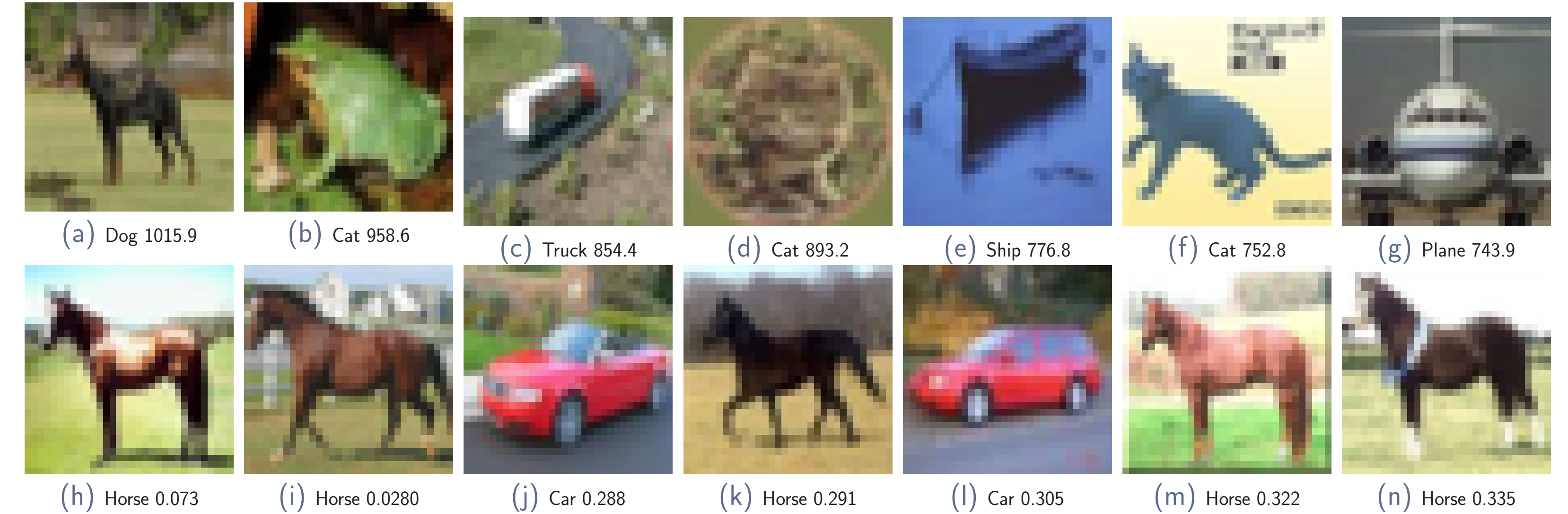


Figure: Most difficult (top-row) and easiest examples (bottom-row) in CIFAR10. Our proposed *action score* is displayed below each image as well as the true label.

## SSD Object Detection on PASCAL VOC - Localization Loss

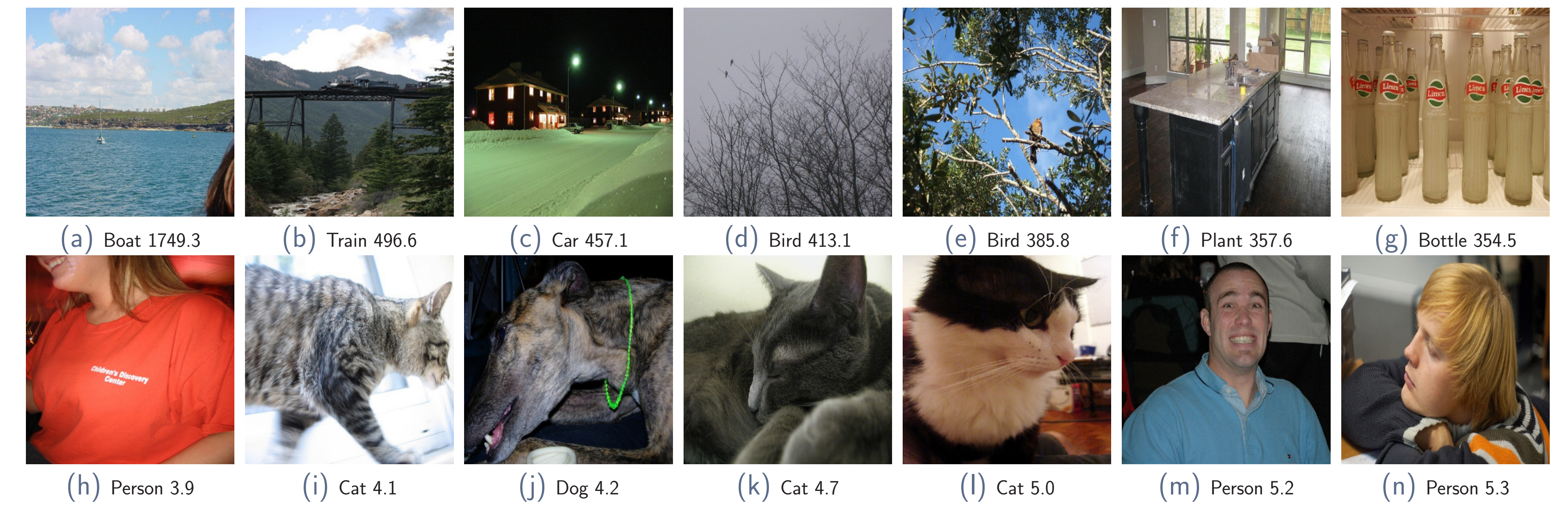


Figure: Most difficult (top-row) and easiest examples (bottom-row) in the VOC 2007-VAL with the SSD localization loss. The *action scores* are displayed below each image as well as the true label.

## SSD Object Detection on PASCAL VOC - Classification Losses

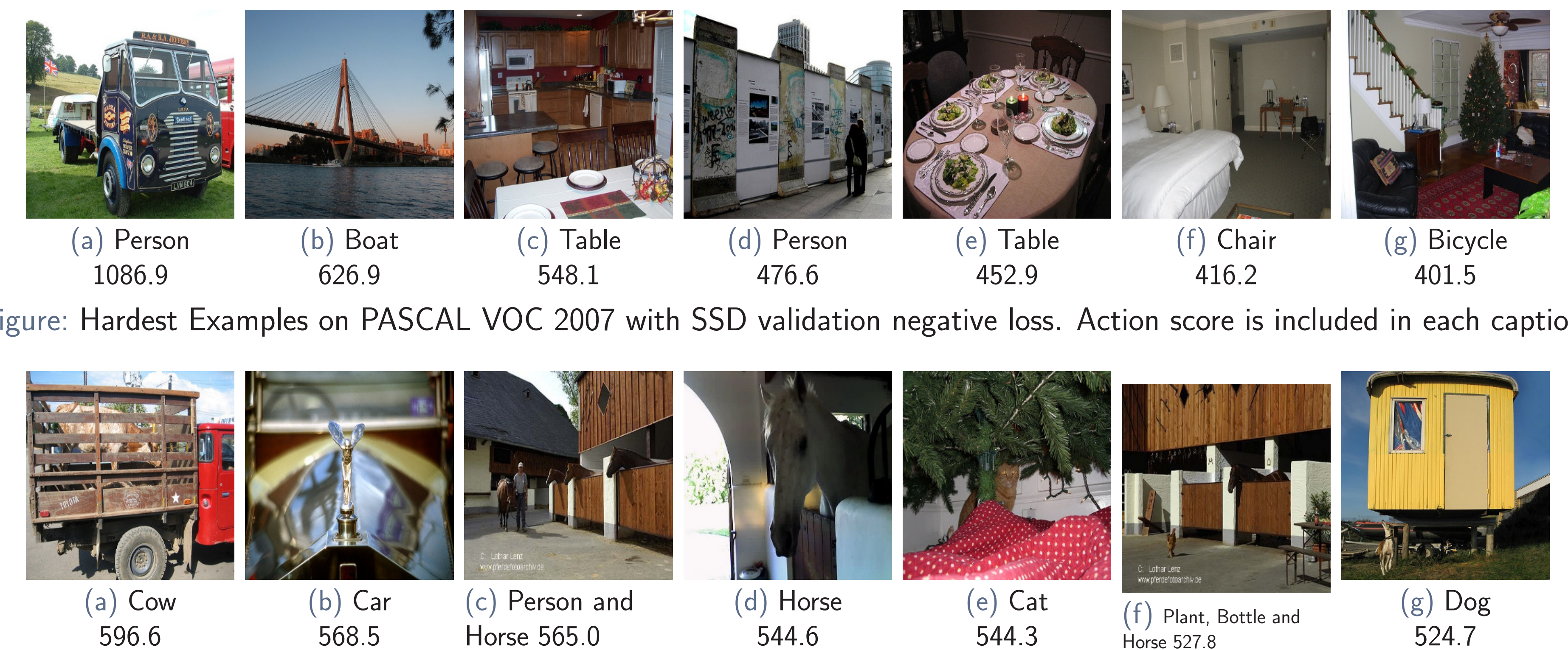


Figure: Hardest Examples on PASCAL VOC 2007 with SSD validation positive loss. Action score is included in each caption.