# CS 676A Project

Group -G26

L.S. Vishnu Sai Rao (12376) and Saurabh Kataria (12637)

# Chosen Paper

**Topic**
- Weakly Supervised Object Detection and localisation

**Paper title and authors**
- Is object localization for free?-weakly-supervised learning with convolutional neural networks
- Oquab, Maxime, Léon Bottou, Ivan Laptev, and Josef Sivic

**Published in**
- Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015

# What is Weakly supervised object detection?

- Training data-set contains images labelled only with lists of objects they contain and not their locations

- Weakly supervised object detection is important because
  - Annotating locations in an image is an expensive process.
  - 'label-only' annotations are often readily available in large amounts, e.g. in the form of text tags or full sentences even geographical meta-data

- The paper employs a CNN architecture to achieve this task.

# Contributions of the paper

Accurate image level labels

Approximate locations of objects

comparable to its fully-supervised counterparts

Is Object localisation free ?

# Example location predictions for images from the Microsoft COCO dataset



Image credit: Oquab, Maxime, et al. "Is object localization for free?-weakly-supervised learning with convolutional neural networks." CVPR 2015

# Procedure: Earlier architecture



Image credit: Oquab, Maxime, et al. "Learning and transferring mid-level image representations using convolutional neural networks." CVPR 2014

# Description of previous architecture

- The procedure is divided into two steps: pre-training and training.

- In the pre-training step, the earlier convolutional layers are trained using the ImageNet database, which consists of tightly cropped images of single objects. This step enables the architecture to recognise individual objects.

- In the training step, two fully connected adaptation layers are added at the end of architecture, which adapts the new architecture to recognise individual objects in a cluttered image with multiple objects in it. A sliding window method with fixed patch size is used to look at different sections of the image.

# Procedure: Modifications



- Treated the fully connected layers as **convolutions** which helps to deal with nearly arbitrary-sized images as input.

- Explicitly searched for the highest scoring object position in the image by adding a single global **max-pooling layer** at the output.

- Used a sum of K binary logistic regression based **cost function** that can explicitly model multiple objects present in the image
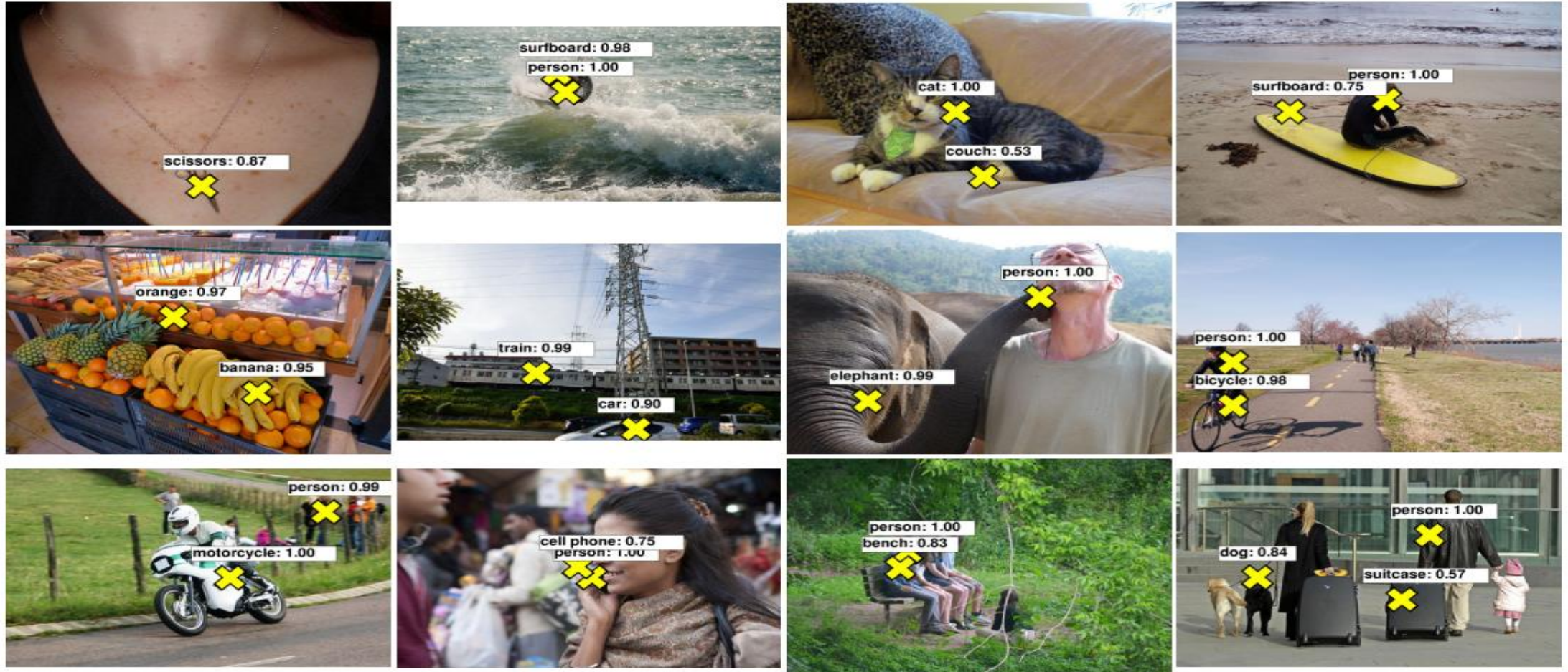
Image credit: Oquab, Maxime, et al. "Is object localization for free?-weakly-supervised learning with convolutional neural networks." CVPR 2015
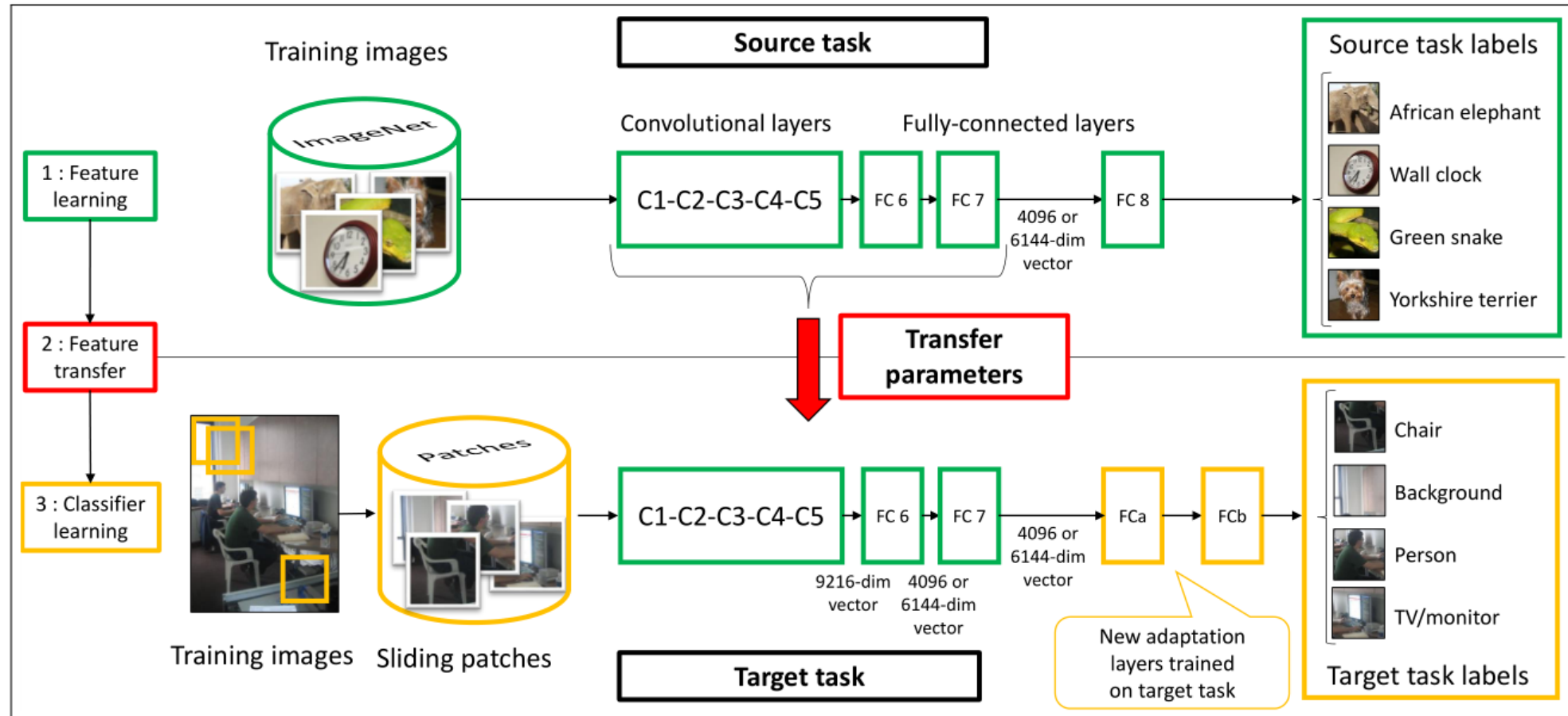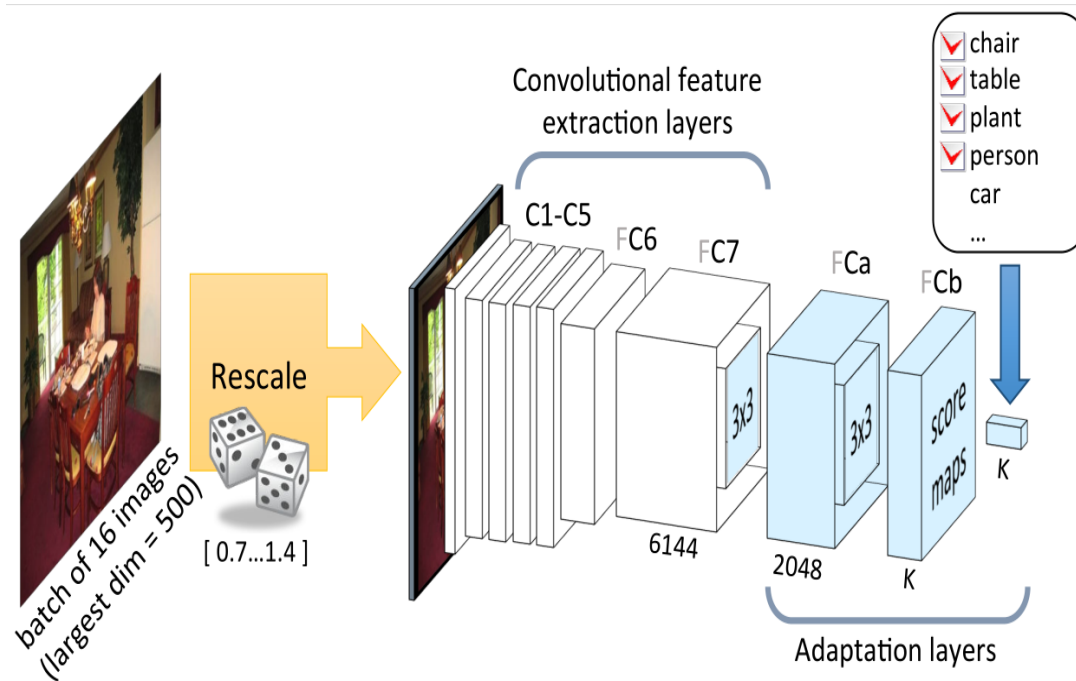
# Descriptions of the modifications

- To achieve the same effect of sliding window, the fully connected layers are treated as convolutional layers.

- The global max-pooling layer added at the end of the architecture converts the n x m x K result to 1 x 1 x K. We can use the n x m x K information to make prediction on the location of detected objects. And, the 1 x 1 x K information is used to predict the presence of objects.

- The new cost function enables the architecture to detect multiple objects from one scene collectively, instead of searching for single objects individually.
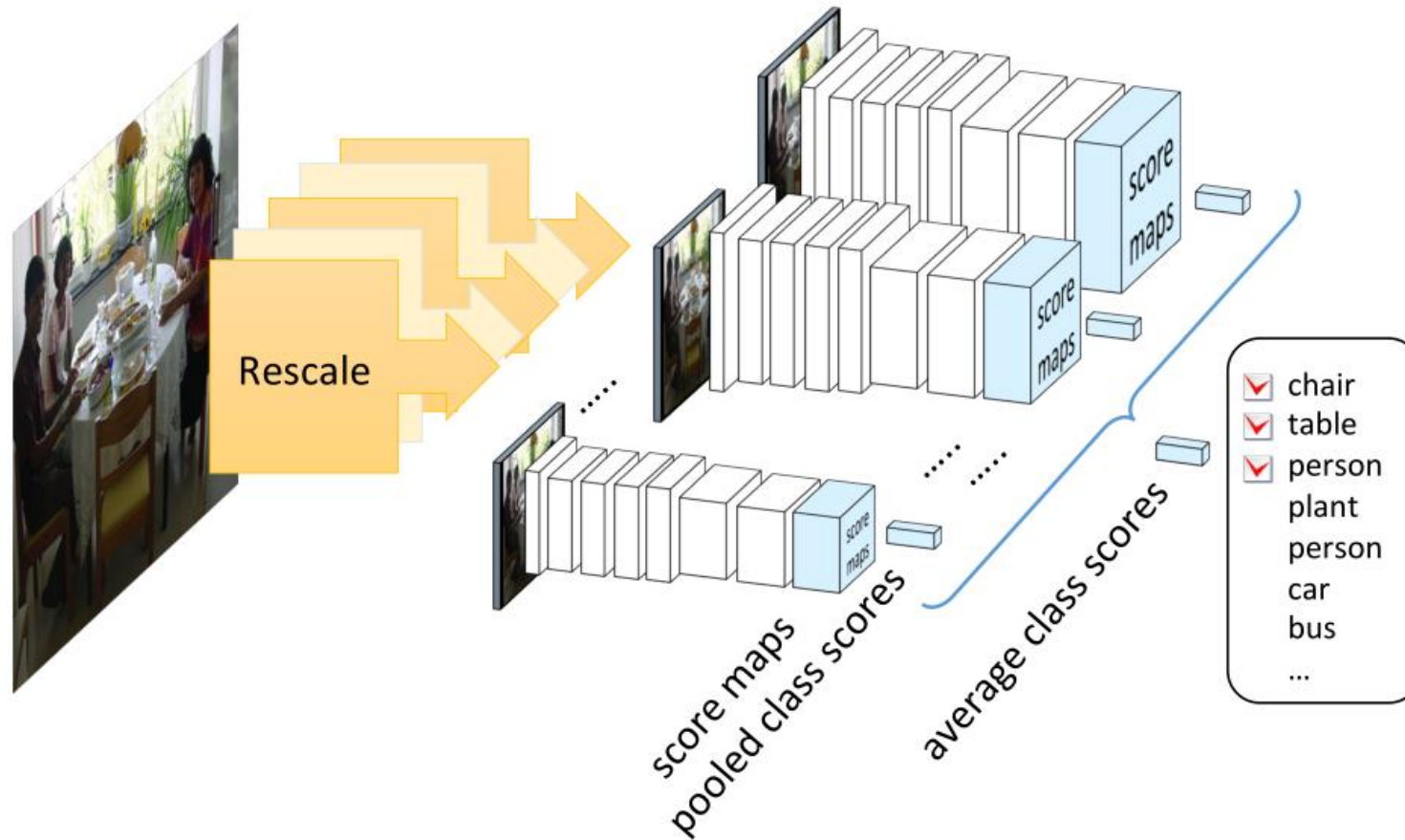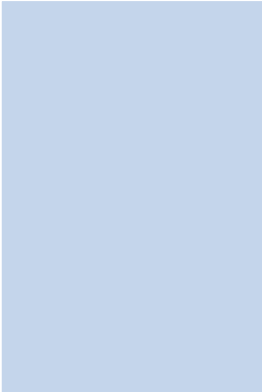
# Procedure: Multi-scaling



Image credit: Oquab, Maxime, et al. "Is object localization for free?-weakly-supervised learning with convolutional neural networks." CVPR 2015
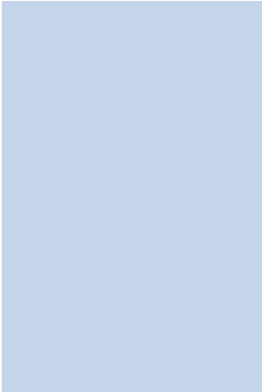
# Role of multi-scaling

- The input image is scaled to various values of s ϵ [0.7, 1.4]. The scaled images are fed to the same network in parallel. The output scores of each network is averaged at the end to give the final score.

- This step enables the architecture to recognise tiny as well as large objects. In other words, this step introduces scale invariance.

# What would we like to contribute?

Predict number of objects in the image as well.

Introduce invariance to rotation and slight distortions
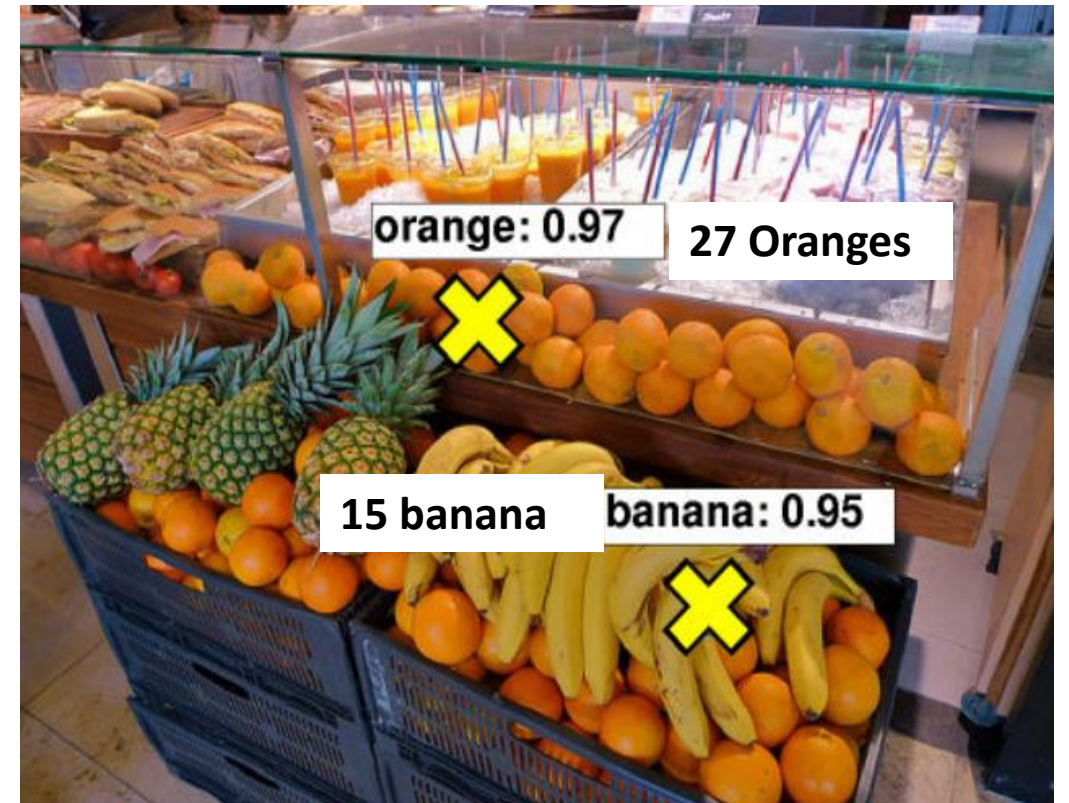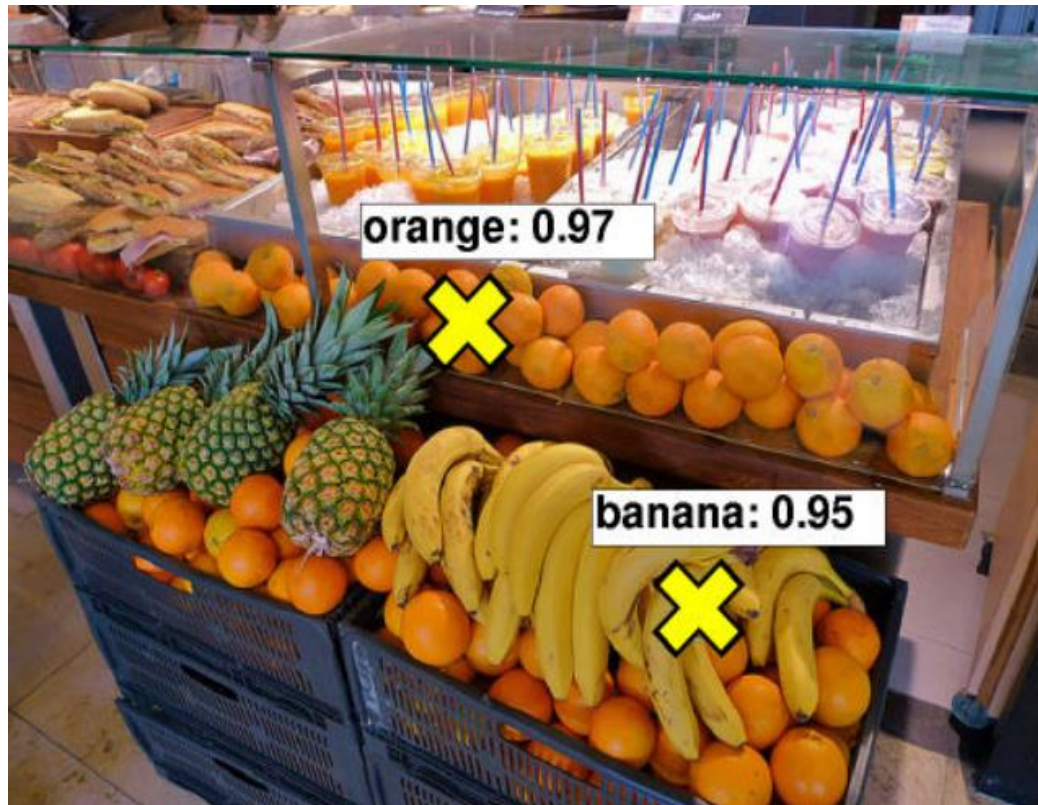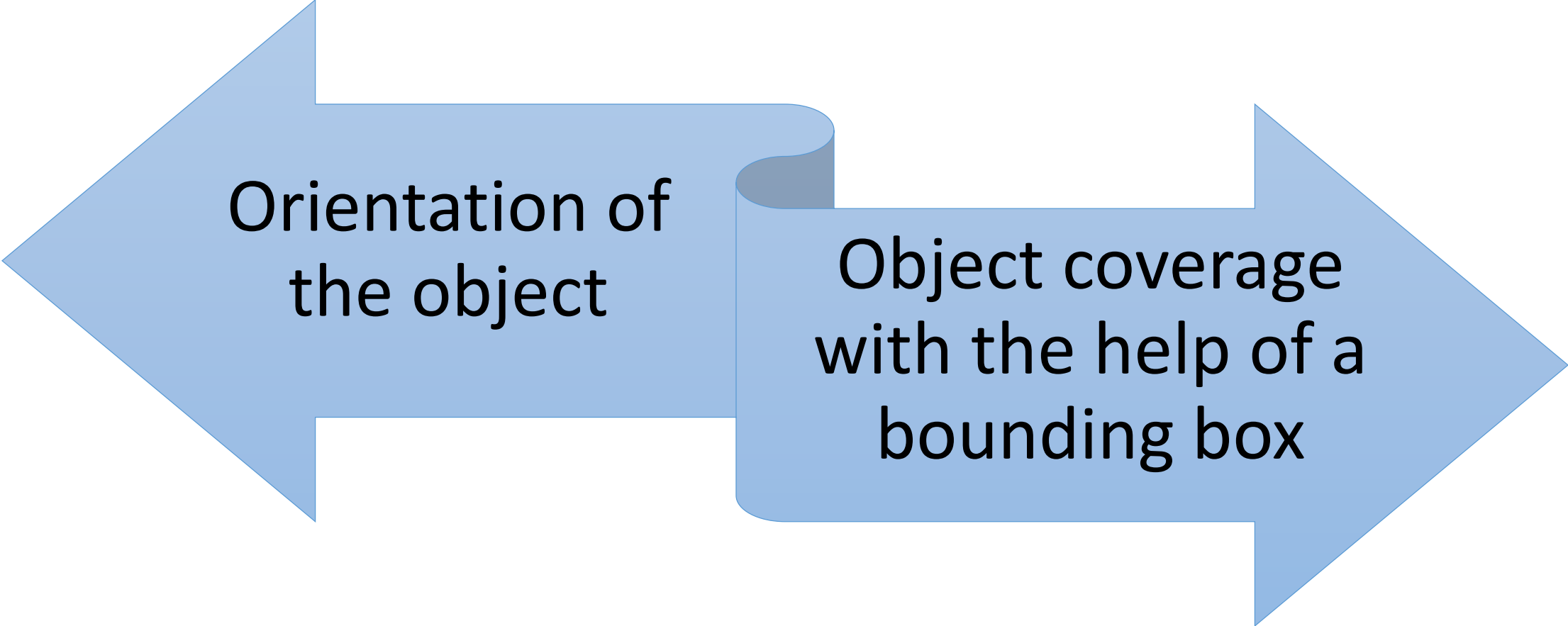
# What would we like to contribute?



Image credit: Oquab, Maxime, et al. "Is object localization for free?-weakly-supervised learning with convolutional neural networks." CVPR 2015

# What else is free?

Orientation of the object

Object coverage with the help of a bounding box