

INTRODUCTION

Convolutional Neural Networks (CNNs) have achieved an impressive success in solving many problems in several fields including computer vision and image processing, attracting a huge interest from the industrial sector. However, running these deep neural networks in embedded systems with limited hardware represent a challenging task for several engineering applications.

REDUCED SSD-TYPE OBJECT DETECTOR

Assumption: SSD-type architecture, composed of:

- a base net (a CNN);
- some additional convolutional layers;
- two siblings predictors, one for localization prediction and one for class prediction.

1. Network Splitting

Let $Obj_Det : \mathbb{R}^{in} \rightarrow \mathbb{R}^{n_{class} \times 4}$ be an object detector. It can be described as the composition of L functions f_j , representing the different layers of the net:

$$Obj_Det = f_{L+1} \circ f_L \circ \dots \circ f_1 \quad (1)$$

Denoting with ℓ the cut-off index, we can define the pre-model and the post-model as:

$$\begin{aligned} basenet_{pre}^{\ell} &= f_1 \circ f_2 \circ \dots \circ f_{\ell}, \\ basenet_{post}^{\ell} &= f_{\ell+1} \circ f_{\ell+2} \circ \dots \circ f_{L+1}. \end{aligned} \quad (2)$$

NOTE: ℓ has to be chosen carefully! The output of the pre-model $\mathbf{x}^{(\ell)}$ lies in a

OBJECTIVES

Goals of the project:

1. Reduction of the memory storage required for an object detector;
2. Application in an embedded system;
3. Accurate performances;
4. Real-time predictions.

high-dimensional space \rightarrow project into a low-dimensional one.

2. Dimensionality Reduction

Let $\mathbf{S} = [\mathbf{x}^{(l),1}, \dots, \mathbf{x}^{(l),N_{train}}]$ be the snapshot matrix. We can compute the SVD of it:

$$\mathbf{S} = \Psi \Sigma \Theta^T, \quad (3)$$

- the columns of Ψ are the left-singular vectors, also called POD modes;
- Σ contains the corresponding eigenvalues.

Fixed r to be the reduced, we define the projection matrix Ψ_r by keeping the first r modes.

$$\mathbf{z}^i = \Psi_r^T \mathbf{x}^{(l),i}, \quad \text{for } i = 1, \dots, N_{train}. \quad (4)$$

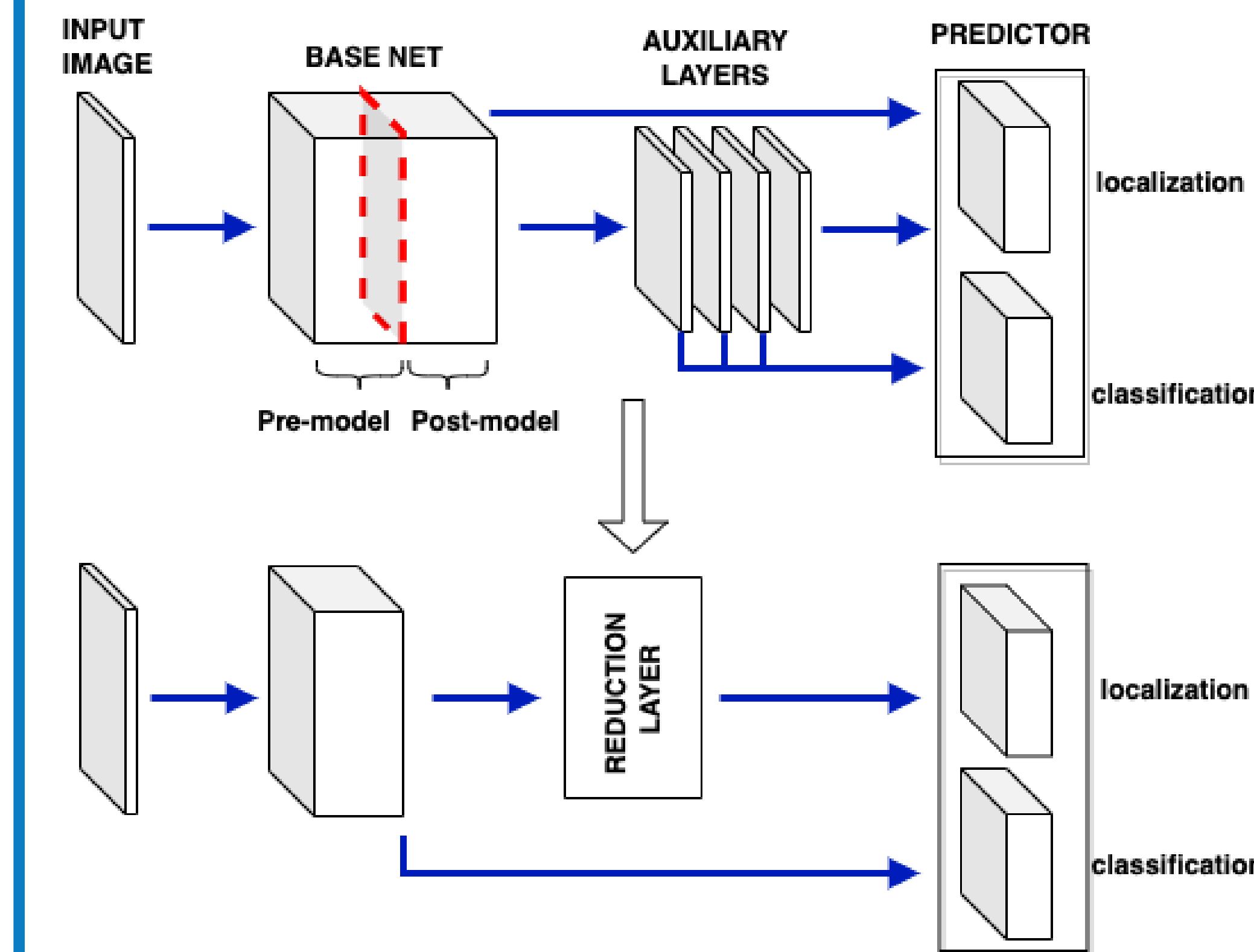
3. Predictor

Same predictor as before:

$$\hat{\mathbf{y}}_{loc}, \hat{\mathbf{y}}_{cls} = predictor(\mathbf{x}^{(l)}, \mathbf{z}) \quad (5)$$

BUT, since the number of inputs is changed, the scale parameter and the number of anchor boxes has to be adjusted.

REDUCED APPROACH- MAIN IDEA



Algorithm 1 Pseudo-code for the construction of the reduced object detector.

Inputs:

- train dataset $\mathcal{D}_{train} = \{\mathbf{x}^{(0),j}, \mathbf{y}_{loc}^j, \mathbf{y}_{cls}^j\}_{j=1}^{N_{train}}$,
- $Obj_Det = [basenet, aux\{layers, predictor\}]$,
- reduced dimension r ,
- index of the cut-off layer ℓ ,
- a test dataset $\mathcal{D}_{test} = \{\mathbf{x}^i, \mathbf{y}_{loc}^i, \mathbf{y}_{cls}^i\}_{i=1}^{N_{test}}$.

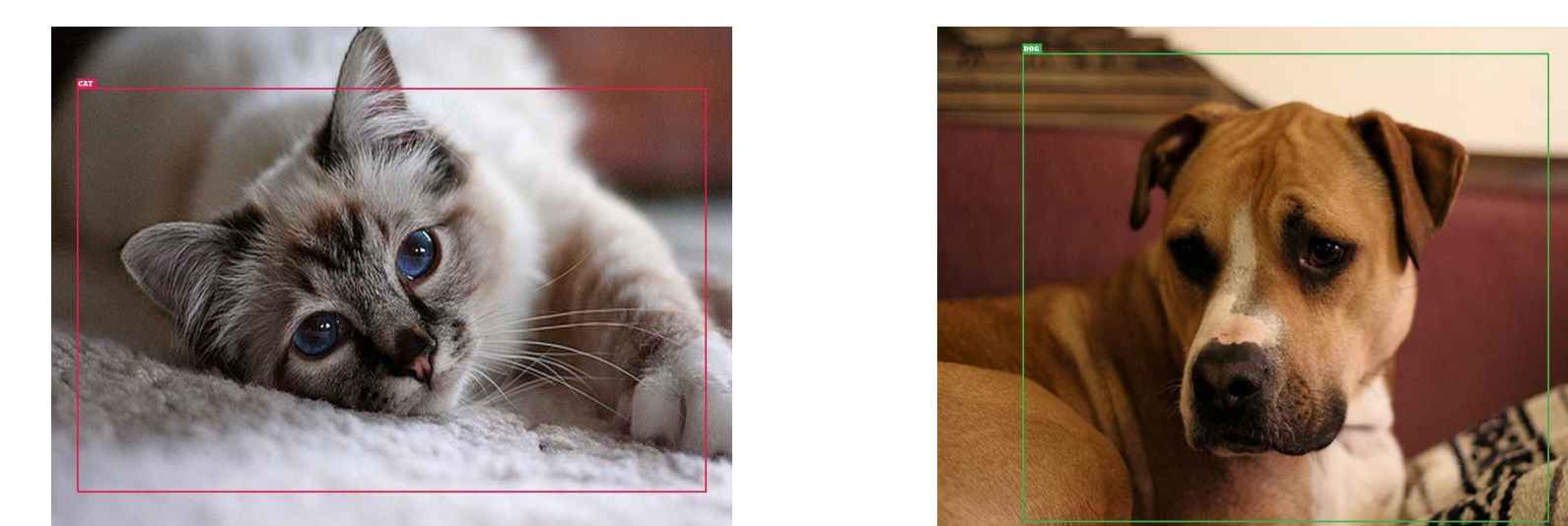
Output: Reduced Object Detector Obj_Det^{red}

- 1: $basenet_{pre}^{\ell}, basenet_{post}^{\ell} = splitting_net(basenet, \ell)$
- 2: $\mathbf{x}^{(\ell)} = basenet_{pre}^{\ell}(\mathbf{x}^{(0)})$
- 3: $\mathbf{z} = reduce(\mathbf{x}^{(\ell)}, r)$
- 4: $\hat{\mathbf{y}}_{loc}, \hat{\mathbf{y}}_{cls} = predictor(\mathbf{x}^{(\ell)}, \mathbf{z})$
- 5: Training of the constructed reduced net using \mathcal{D}_{test} .

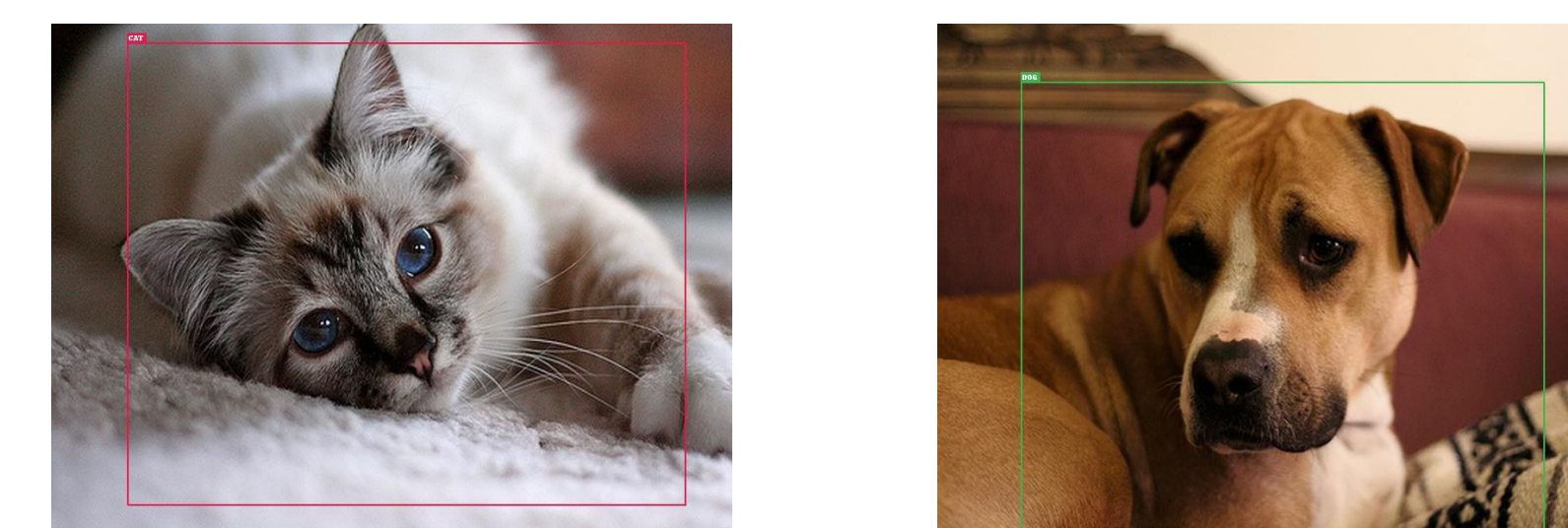
RESULTS1- CATS & DOGS

Network	mAP	Storage (Mb)	Training Time
SSD300	70.2%	91.09	43.5 h
SSD300_red	59%	77.45	26 h

SSD300



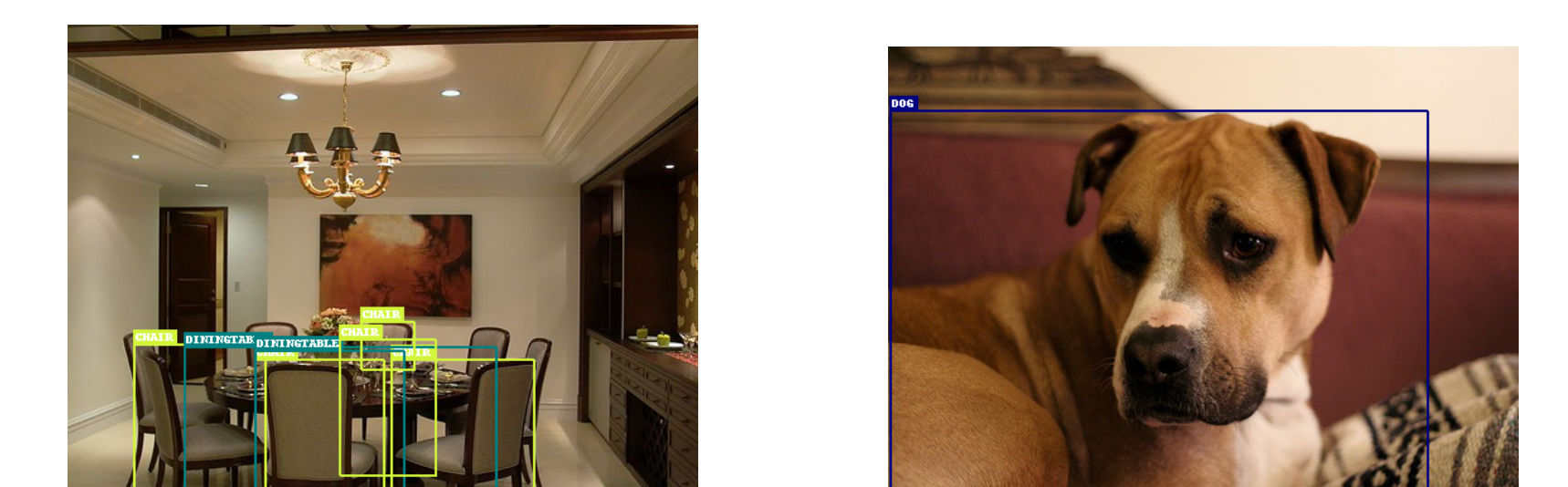
Reduced SSD300



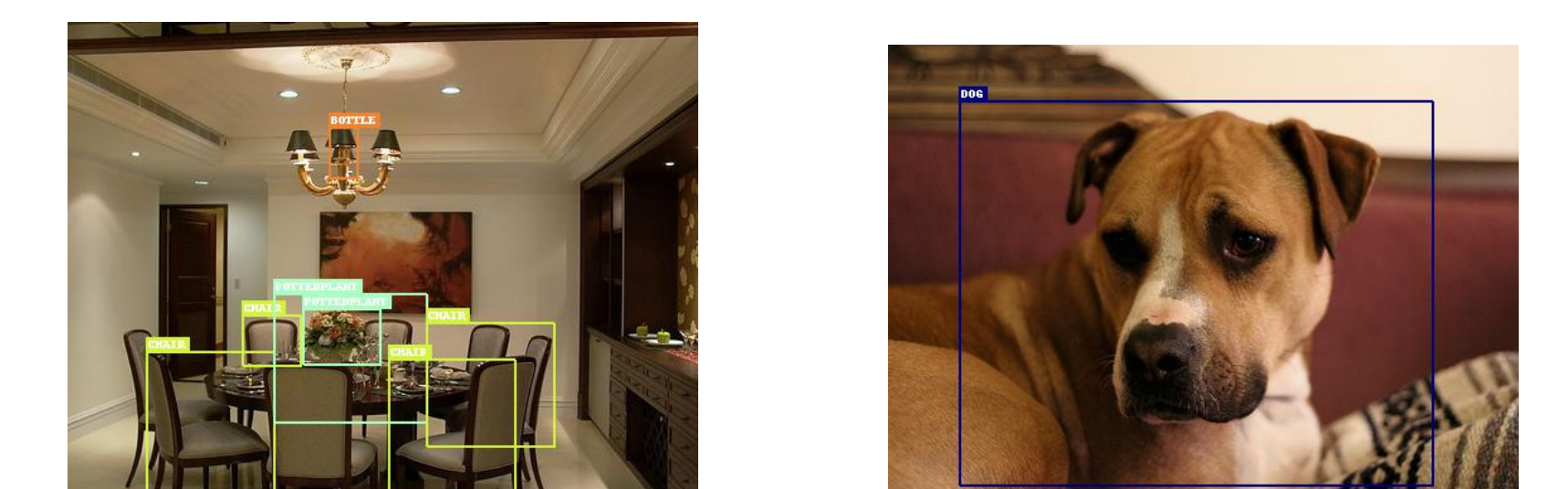
RESULTS-2 PASCAL VOC

Network	mAP	Storage (Mb)	Training Time
SSD300	77.8%	100.23	48 h
SSD300_red	39%	76.23	18 h

SSD300



Reduced SSD300



REFERENCES

- [1] Laura Meneghetti, Nicola Demo, and Gianluigi Rozza. A Dimensionality Reduction Approach for Convolutional Neural Networks. *arXiv preprint arXiv:2110.09163*, 2021.
- [2] Chunfeng Cui, Kaiqi Zhang, Talgat Daulbaev, Julia Gusak, Ivan Oseledets, and Zheng Zhang. Active subspace of neural networks: Structural analysis and universal attacks. *SIAM Journal on Mathematics of Data Science*, 2(4):1096–1122, 2020.
- [3] Jan S. Hesthaven, Gianluigi Rozza, and Benjamin Stamm. *Certified Reduced Basis Methods for Parametrized Partial Differential Equations*. Springer Briefs in Mathematics. Springer, Switzerland, 1 edition, 2015.

FUTURE RESEARCH

- Knowledge Distillation for Object Detection;
- Employment of hyperreduction techniques or POD variants.
- Iterative procedure for the choice of ℓ .

CONTACT INFORMATION

Code <https://github.com/mathLab/Smithers>
Email laura.meneghetti@sissa.it
nicola.demo@sissa.it
gianluigi.rozza@sissa.it