# Are we done with object recognition? The iCub robot's perspective.

**Giulia Pasquale · Carlo Ciliberto · Francesca Odone · Lorenzo Rosasco · Lorenzo Natale**

**Abstract** We report on an extensive study of the benefits and limitations of current deep learning approaches to object recognition in robot vision scenarios, introducing a novel dataset used for our investigation. To avoid the biases in currently available datasets, we consider a natural human-robot interaction setting to design a data-acquisition protocol for visual object recognition on the iCub humanoid robot. Analyzing the performance of off-the-shelf models trained off-line on large-scale image retrieval datasets, we show the necessity for knowledge transfer. We evaluate different ways in which this last step can be done, and identify the major bottlenecks affecting robotic scenarios. By studying both object categorization and identification problems, we highlight key differences between object recognition in robotics applications and in image retrieval tasks, for which the considered deep learning approaches have been originally designed. In a nutshell, our results confirm the remarkable improvements yield by deep learning in this setting, while pointing to specific open challenges that need be addressed for seamless deployment in robotics.

G. Pasquale · L. Natale
Humanoid Sensing and Perception, Istituto Italiano di Tecnologia, Via Morego 30, Genova, IT
E-mail: giulia.pasquale@iit.it

G. Pasquale · F. Odone · L. Rosasco
Dipartimento di Informatica, Bioingegneria, Robotica e Ingegneria dei Sistemi, Universitá degli Studi di Genova, Via Dodecaneso 35, Genova, IT

G. Pasquale · L. Rosasco
Laboratory for Computational and Statistical Learning (IIT@MIT), Massachusetts Institute of Technology, Bldg. 46-5155, 43 Vassar Street, Cambridge, MA

L. Rosasco
Center for Brains, Minds and Machines, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, MA

## 1 Introduction

Artificial intelligence has recently progressed dramatically, largely thanks to the advance in deep learning. Computational vision, specifically object recognition, is perhaps the most obvious example where deep learning has achieved so stunning results to raise the question of whether this problem is actually solved (Krizhevsky et al 2012; Simonyan et al 2014; Szegedy et al 2015; He et al 2015, 2016). Should this be the case, robotics would be a main field where the benefits could have far reaching effect. Indeed, the lack of reliable visual skills is largely considered a major bottleneck for the successful deployment of robotic agents in everyday life (Kemp et al 2007).

With this perspective in mind, we have recently started an effort to isolate and quantify the benefits and limitations, if any, of deep learning approaches to visual object recognition in robotic applications (Pasquale et al 2015, 2016a). The remarkable performance of deep learning methods for object recognition has in fact been primarily reported on computer vision benchmarks such as (Griffin et al 2007; Everingham et al 2010, 2015; Russakovsky et al 2015), which are essentially designed for large-scale image retrieval tasks and are hardly representative of a robotic application setting (a motivation common to other recent works such as (Pinto et al 2011; Leitner et al 2015; Oberlin et al 2015; Borji et al 2016)).

Clearly, visual perception is only one of the possible sensory modalities equipping modern robots, that can be involved in the object recognition process (see for example (Luo et al 2017; Higy et al 2016)). In addition it has been shown that the physical interaction with the

environment can be used to aid perception (Pinto et al 2016), demonstrating that there is more than "just" vision to object recognition. Nonetheless, visual recognition of objects is evidently the very first and critical step for autonomous agents to act in the real world. Current deep learning-based artificial systems perform so well in this task, that it seems natural to ask how far they can go, before further perceptual cues are needed. To this end, in this work we focus on the problem of object recognition in robotics using only visual cues.

We consider a prototypical robotic scenario, where a humanoid robot is taught to recognize different objects by a human through natural interaction. We started with the design of a dataset tailored to reflect the visual "experience" of the robot in this scenario. This dataset is rich and easy to expand to include more data and complex perceptual scenarios. It includes several object categories with many instances per category, hence allowing to test both object categorization and identification tasks. Notably, the dataset is segmented into multiple sets of image sequences per object, representing specific viewpoint transformations like scaling, in- and out-of-plane rotations, and so forth. It provides a unique combination that, to our knowledge, was missing in the literature, allowing for articulated analyses of the robustness and invariance properties of recognition systems within a realistic robotic scenario. Since we used the iCub robot (Metta et al 2010), we called it iCubWorld Transformations (iCWT for short).

We performed extensive empirical investigation using different state-of-the-art Convolutional Neural Network architectures, demonstrating that off-the-shelf models do not perform accurately enough and pre-trained networks need to be adapted (fine-tuned) to the data at hand to obtain substantial improvements. However, these methods did not quite provide the close to perfect accuracy one would wish for. We hence proceeded taking a closer look at the results, starting from the question of whether the missing gap could be imputed to lack of data. Investigating this latter question highlighted a remarkable distance between iCWT and other datasets such as ImageNet (Russakovsky et al 2015). We identified clear differences between the object recognition task in robotics with respect to scenarios typically considered in learning and vision. In fact, as we show in the paper, our empirical observations on iCWT are extended to other robotic datasets like, e.g., Washington RGB-D (Lai et al 2011).

Along the way, our analysis allowed also to test the invariance properties of the considered deep learning networks and quantify their merits not only for categorization but also for identification.

The description and discussion of our empirical findings is concluded with a critical review of some of the main venues of improvements, from a pure machine learning perspective but also taking extensive advantage of the robotic platform. Indeed, bridging the gap in performance appears to be an exciting avenue for future multidisciplinary research.

In the next Section we discuss several related works, while the rest of the paper is organized as follows: Sec. 2 introduces the iCWT dataset and its acquisition setting. In Sec. 3 we review the deep learning methods considered for our empirical analysis, which is reported in Sec. 4 for the categorization task and in Sec. 5 for object identification. In Sec. 6 we finally show how the same observations drawn from experiments on iCWT hold also for other robotic datasets, specifically we consider Washington RGB-D (Lai et al 2011). Sec. 7 concludes our study with the review of possible directions of improvement for visual recognition in robotics.

## 1.1 Deep Learning for Robotics

Deep Learning methods are receiving growing attention in robotics, and are being adopted for a variety of problems such as object recognition (Schwarz et al 2015; Eitel et al 2015; Pinto et al 2016; Held et al 2016; Pasquale et al 2016a), place recognition and mapping (Sünderhauf et al 2015, 2016), object affordances (Nguyen et al 2016), grasping (Redmon and Angelova 2015; Pinto and Gupta 2016; Levine et al 2016) and tactile perception (Baishya and Bäuml 2016). We limit our discussion to the work on object recognition, which is more relevant to this paper.

In (Schwarz et al 2015; Eitel et al 2015) the authors demonstrate transfer learning from pre-trained deep Convolutional Neural Networks (CNNs) and propose a way to include depth information from an RGB-D camera by encoding depth data into RGB with colorization schemes. The main idea of (Schwarz et al 2015) is to extract a feature vector from the CNN and train a cascade of Support Vector Machines (SVMs) to discriminate the object's class, identity and position, while (Eitel et al 2015) proposes a pipeline to fine-tune the CNNs weights on RGB-D input pairs. Both works focus on the Washington RGB-D benchmark (Lai et al 2011) and improve its state-of-the-art performance. In this paper, we consider a less constrained setting, in that CNNs are trained with data acquired by the robot during natural interaction (which undergo, therefore, more challenging viewpoint transformations). We adopt similar techniques for transfer learning, but we assess a wide range of architectures and fine-tuning approaches.
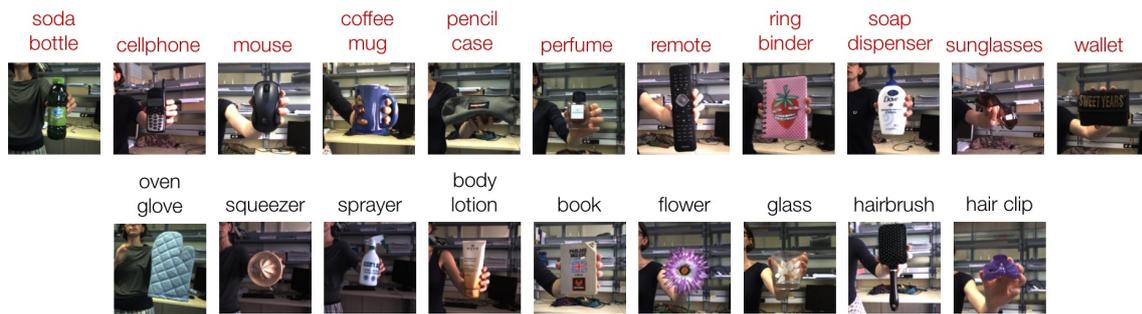
**Fig. 1 iCWT categories**. For each category in iCWT, we report one example image for one instance. (Red) categories appear also in the ILSVRC 2012 dataset. (Black) categories appear in ImageNet but not in ILSVRC (see supplementary material for more details).

The work in (Pinto et al 2016) shows how the robot can use self-generated explorative actions (like pushing and poking objects) to autonomously extract example data and train a CNN. In contrast to (Pinto et al 2016), the use of a human teacher (see Sec. 2 for details on the acquisition setup) gives us more control on the object transformations. On the other hand, our work could be extended by introducing self-supervision using explorative actions similar to the ones in (Pinto et al 2016).

The work in (Held et al 2016) and our work in (Pasquale et al 2016a) are closely related to the work presented in this paper in that they investigate invariance properties of CNNs and how to improve them in order to learn from few examples. They focus, however, on instance recognition, whereas in this paper we include in the analysis – thus significantly extending (Pasquale et al 2016a) – the problem of object categorization. In addition, we perform a detailed investigation of various transfer learning approaches and present a systematic evaluation of the recognition performance for specific viewpoint transformations.

### 1.2 Datasets for Visual Recognition in Robotics

In the literature, several datasets have been used to benchmark visual object recognition in robotics: COIL (Nene et al 1996), ALOI (Geusebroek et al 2005), Washington RGB-D (Lai et al 2011), KIT (Kasper et al 2012), SHORT-100 (Rivera-Rubio et al 2014), BigBIRD (Singh et al 2014), Rutgers Amazon Picking Challenge RGB-D Dataset (Rennie et al 2016) are only some examples. One of the main characteristic of these datasets is to capture images of an object while it undergoes viewpoint transformations. However, these datasets are usually acquired in strictly controlled settings (i.e. using a turntable), because they are aimed to provide also pose annotations. As a consequence, image sequences differ from the actual data that can be gathered by a robot in its operation: they do not show substantial variations of objects' appearance and often background or light changes are missing or under-represented. See (Leitner et al 2015) for a review of the major limitations of current datasets.

The NORB dataset (LeCun et al 2004), albeit acquired in a similar turntable setting, is one of the first benchmarks released in support of the investigation of invariance properties of recognition methods. A similar study focusing on deep learning methods is described in (Borji et al 2016) using the iLab-20M dataset. This aims to be a comprehensive benchmark for visual object recognition that, while representing a high number of object instances, provides also varied images of each object. While presenting a remarkably higher variability, however, also iLab-20M is acquired in a turntable setting (to collect pose annotations), hence suffering from similar limitations.

Our iCWT dataset separates from previous work in that objects are captured during "natural" transformations. Acquisition is performed in a "semi-controlled" setting intended to reproduce typical uncertainties faced by the visual recognition system of a robot during a real-world task. Very few works in the literature consider "real" (i.e. non-synthetic) transformations when evaluating the invariance properties of visual representations (see, e.g., (Goodfellow et al 2009)). To our knowledge, iCWT is the first dataset to address invariance-related questions in robotics. Moreover, it accounts for a much wider range of visual transformations with respect to previous datasets. In the following, we discuss the data collection and the acquisition setting in detail. Note that, while we used an initial subset of iCWT in (Pasquale et al 2016a), in this paper we present the dataset for the first time in its entirety.

### 2 The iCubWorld Transformations Dataset

The iCubWorld Transformations (iCWT) is a novel benchmark for visual object recognition, which we use for

**Fig. 2 Semantic variability.** Sample images for the different object instances in the *mug* category to provide a qualitative intuition of the semantic variability in iCubWorld. See Fig. 14 in the supplementary material for more examples.

the empirical analysis in this work. iCWT is the latest release within the iCubWorld[1] project, whose goal is to benchmark and improve artificial visual systems for robotics. iCubWorld datasets (Fanello et al 2013; Pasquale et al 2015) are designed to record a prototypical visual "experience" of a robot while it is performing vision-based tasks. To this end, we devised and implemented a simple human-robot interaction application, during which we acquire images for the dataset from the robot's cameras.

There is a remarkable advantage in collecting iCubWorld directly from the robot platform. The resulting data collection offers a natural testbed, as close as possible to the real application. This ensures that the performance measured *off-line* on iCubWorld can be expected to generalize well when the system is deployed on the actual robot. Note that this aspect of iCubWorld is extremely relevant since visual biases make it typically difficult to generalize performances across different datasets and applications, as already well known from previous work (Pinto et al 2008; Torralba and Efros 2011; Khosla et al 2012; Hoffman et al 2013; Rodner et al 2013; Model and Shamir 2015; Stamos et al 2015; Tommasi et al 2015) and also shown empirically in Sec. 4.1 of this paper.

Currently, to acquire iCubWorld releases we did not make extensive use of the robot's physical capabilities (e.g., manipulation, exploration, etc.). This was done because latest deep learning methods already achieve remarkable performance by relying solely on visual cues and our goal was to evaluate their accuracy in *isolation* in a robotic setting. While exploiting the robot body could provide further advantages to modeling and recognition (see for instance (Moldovan et al 2012), (Leitner et al 2014) and (Dansereau et al 2016)), it would also prevent us to assess the contribution of the visual component alone.

## 2.1 Acquisition Setup

During data acquisition a human "teacher" shows an object to the robot and pronounces the associated label. The robot exploits bottom-up visual cues to track and record images of the object while the human moves
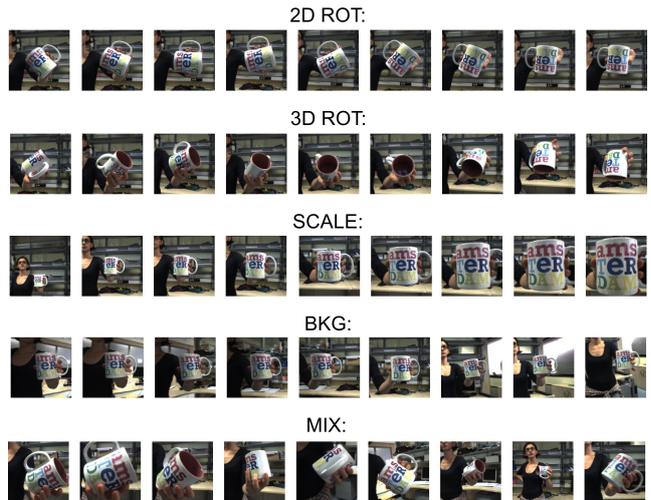
**Fig. 3 Visual transformations.** Excerpts from the sequences acquired for one mug, representing the object while it undergoes specific visual transformations.

**Table 1** Summary of the *iCubWorld Transformations* dataset

| # Categories | # Obj. per Category | # Days | Transformations | # Frames per Session |
|---|---|---|---|---|
| 20 | 10 | 2 | 2D ROT, 3D ROT SCALE, BKG | 150 |
| | | | MIX | 300 |

and shows it from different poses. Annotations are automatically recorded for each image in terms of the object's label (provided by the human) and bounding box (provided by the tracker).

Differently from all previous releases of iCubWorld, for the acquisition of iCWT we decided to rely on a tracker based on depth segmentation (Pasquale et al 2016b) instead of one based on the detection of object motion, because this greatly simplifies the acquisition while providing more stable and precise bounding boxes.

We developed an application to scale the acquisition seamlessly to hundreds of objects, which were collected during multiple interactive sessions. iCWT is available on-line and we plan to make also this application publicly available in order for other laboratories to use the same protocol to collect their own (or possibly contribute to) iCubWorld. The proposed acquisition approach allows to build large-scale data collections, fully annotated through natural interaction with the robot and with minimal manual intervention. Moreover, the application can be directly deployed on other iCub robots and, being relatively simple, can also be adapted to other platforms.

2.2 Dataset Overview

iCWT includes 200 objects evenly organized into 20 categories that can be typically found in a domestic environment. Fig. 1 reports a sample image for each category in iCWT: 11 categories (in red in the figure) are also in the ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) 2012 (Russakovsky et al 2015), i.e. we found semantically and visually similar classes among the 1000 of the classification challenge. The remaining 9 categories do not appear ILSVRC but belong (or are similar) to a synset in the larger ImageNet dataset (Deng et al 2009). To provide a qualitative intuition of the semantic variability within a given category, namely the different visual appearance of object instances in the same category, Fig. 2 shows a sample image from the 10 instances in the *mug* category. We refer the reader to the supplementary material (Fig. 14) for example images of all object instances in iCWT.

In iCWT each object is shown in multiple image sequences while it undergoes specific visual transformations (such as rotations, scaling or background changes). For each object instance, we acquired 5 image sequences, while the human supervisor was performing a different visual transformation on the object. Fig. 3 reports excerpts of these sequences, which contain, respectively:

**2D Rotation** The human rotated the object parallel to the camera image plane. The same scale and position of the object were maintained (see Fig. 3, first row).

**3D Rotation** Similarly to 2D rotations, the object was kept at same position and scale. However, this time the human applied a generic rotation to the object (not parallel to the image plane). As a consequence different "faces" of the object where shown to the camera (Fig. 3, second row).

**Scale** The human moved the object towards the cameras and back, thus changing the object's scale in the image. No change in the object orientation (no 2D or 3D rotation) was applied (Fig. 3, third row).

**Background** The human moved the object around the robot, keeping approximately the same distance (scale) and pose of the object with respect to the camera plane. Because the robot tracks the object, the background changes while the object appearance remains approximately the same (Fig. 3, fourth row).

**Mix** The human moved the object freely in front of the robot, as a person would naturally do when showing a new item to a child. In this sequence all nuisances in all combinations can appear (Fig. 3, fifth row).

Each sequence is composed by approximately 150 images acquired at 8 frames per second in the time interval of $20s$, except for the *Mix* sequence that lasted $40s$ and comprises $\sim 300$ images. As anticipated, the acquisition of the 200 objects was split into multiple sessions performed in different days. The acquisition location was always the same (with little uncontrolled changes in the setting across days). The illumination condition was not artificially controlled, since we wanted to investigate its role as a further nuisance: to this end, we acquired objects at different times of the day and in different days, so that lighting conditions are slightly changing across the 200 objects (but not within the five sequences of an object, which were all acquired in the span of few minutes). Moreover, we repeated the acquisition of each object in two different days, so that we ended up with 10 sequences per object, containing 5 visual transformations in 2 different illumination conditions. The adopted iCub's cameras resolution is $640 \times 480$. Both left and right images provided by the iCub's stereo pair were acquired, to allow for offline computation of the disparity map and possibly further improvement of the object's localization and segmentation. We recorded the centroid and bounding box of the object provided by the tracker at each frame. Tab. 1 summarizes the main characteristics of iCub-World Transformations. We refer to the iCub's website [2] for details about the cameras and their setting.

## 3 Methods

### 3.1 Deep Convolutional Neural Networks

Deep Convolutional Neural Networks (CNNs) are hierarchical models iterating multiple processing layers. Their structure aims to map the input (say, an image) into a series of *feature maps* or *representations* that progressively select the visual features which are most relevant to the considered task. The prototypical structure of a CNN (see Fig. 4) alternates blocks of (i) convolution (followed by element-wise non linearities, as sigmoids (Bishop 2006) or ReLUs (He et al 2015)), and (ii) spatial pooling and downsampling. The convolution (plus non linear) layers are such to progressively extract, from the image, maps of selected and more complex features (from edges to local pattern up to object parts), which are relevant to the task

---

[2] http://www.icub.org/
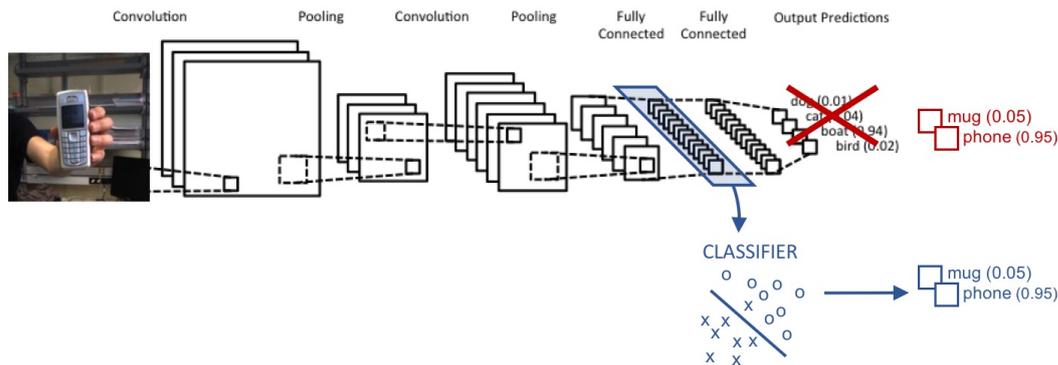[4] https://www.clarifai.com/technology

**Fig. 4** Example of a Convolutional Neural Network (Sec. 3.1) and of the two knowledge transfer approaches considered in this work (Sec. 3.2). **(Blue Pipeline) Feature Extraction**: in this case the response of one of the layers is used as a feature vector for a "shallow" predictor like RLSCs, SVMs (see Sec. 3.2.1), which is trained on the new task. **(Red Pipeline) Fine-tuning**: in this case the network is trained end-to-end to the new task by replacing the final layer and using the original model as a "warm-restart" (see Sec. 3.2.2). Network image from [4].

at hand (Chatfield et al 2014; Yosinski et al 2014; Donahue et al 2014; Zeiler and Fergus 2014). Spatial down-sampling and pooling layers make these features more robust (ideally, invariant) to transformations of the input at increasingly larger scales. A common strategy in image classification is to follow these blocks with one or more fully connected layers (namely, a standard Neural Network). In classification settings, the last layer is a *softmax* function, which maps the output into class-likelihood scores, whose maximum is the predicted class.

The modular structure of CNNs allows training their parameters (namely, convolution filters) simultaneously for all layers (also known as *end-to-end learning*) via back-propagation (LeCun et al 1989). Given the large number of parameters to be optimized (in the order of millions), CNNs typically need large amounts of training data to achieve good performance. Often, the training examples are artificially increased by synthetically modifying the images (*data augmentation*). To further mitigate the risk of overfitting, regularization techniques such as L2 regularization, dropout (Hinton et al 2012; Srivastava et al 2014) or, more recently, batch normalization (Ioffe and Szegedy 2015), have proved helpful.

In this work we investigate the performance of modern CNNs on the robotic setting of iCubWorld. To this end, we selected four architectures achieving the highest accuracy on the ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) (Russakovsky et al 2015). We used their implementation trained on ILSVRC 2012 and publicly available within the CAFFE (Jia et al 2014) framework. Specifically we consider **CaffeNet**[5], a variation of AlexNet (Krizhevsky et al 2012), and **VGG-**

16[6] (Simonyan et al 2014). These two models concatenate a number of convolution layers (5 and 13) with 3 fully connected layers and comprise around 60 and $140M$ parameters respectively. Then we consider **GoogLeNet**[7] (Szegedy et al 2015) and **ResNet-50**[8] (He et al 2016), which slightly diverge from the standard CNN structure described above in that they respectively employ so-called *inception modules* or *residual connections*, and the number of parameters is also reduced ($4M$ for GoogLeNet, 22 layers, and $20M$ for ResNet-50, 50 layers).

## 3.2 Transfer Learning Techniques

The need for large datasets to successfully train deep CNNs could in principle prevent their applicability to problems where training data is scarce. However, recent empirical evidence has shown that the knowledge learned by a CNN on a large-scale problem can be "transferred" to multiple domains. In this section we review two of the most well-established methods which we empirically assess in our experiments, namely *feature extraction* and *fine-tuning*.

### 3.2.1 Feature Extraction

It has been shown that CNNs trained on large, varied image datasets can be used on smaller datasets as generic "feature extractors", leading to remarkable performance (Schwarz et al 2015; Oquab et al 2014;

---

[5] https://github.com/BVLC/caffe/tree/master/models/bvlc_reference_caffenet

[6] http://www.robots.ox.ac.uk/$\sim$vgg/research/very_deep/

[7] https://github.com/BVLC/caffe/tree/master/models/bvlc_googlenet

[8] https://github.com/KaimingHe/deep-residual-networks

**Table 2** Feature extraction layers for the four architectures considered in this work. We used the notation adopted in CAFFE, in which the number identifies the layer number and the label specifies its type (i.e., fully connected or pooling layer).

| Model | Output Layer |
|---|---|
| CaffeNet | fc6 or fc7 |
| GoogLeNet | pool5/7x7_s1 |
| VGG-16 | fc6 or fc7 |
| ResNet-50 | pool5 |

**Table 3** Fine-tuning protocols for *CaffeNet* and *GoogLeNet*. Base LR is the starting learning rate of all layers that are initialized with the original model. The FC layers that are learned from scratch are indicated using their names in CAFFE models ($2^{nd}$ row), specifying the starting learning rate used for each of them. For the other parameters, we refer the reader to CAFFE documentation.

| | **CaffeNet** | | **GoogLeNet** | |
|---|---|---|---|---|
| | *adaptive* | *conservative* | *adaptive* | *conservative* |
| **Base LR** | 1e-3 | 0 | 1e-5 | 0 |
| **Learned FC Layers** | fc8: 1e-2 | fc8: 1e-4<br>fc7: 1e-4<br>fc6: 1e-4 | loss3/classifier: 1e-2<br>loss1(2)/classifier: 1e-3<br>loss1(2)/fc: 1e-3 | |
| **Dropout (%)** | fc7: 50<br>fc6: 50 | | pool5/drop_7x7_s1: 60<br>loss1(2)/drop_fc: 80 | |
| **Solver**<br>**LR Decay Policy** | SGD<br>Polynomial (exp 0.5) | | Adam<br>No decay | |
| **# Epochs** | 6 | 36 | 6 | |
| **Batch Size** | 256 | | 32 | |

Sünderhauf et al 2015; Pasquale et al 2016a). This is typically done by training a classifier (such as a Support Vector Machine (SVM) or a Regularized Least Squares Classifier (RLSC) (Bishop 2006)) on feature vectors obtained as the activations of intermediate network layers. This strategy, depicted in Fig. 4 (Blue pipeline), can also be interpreted as changing the last layer of the CNN and training it on the new dataset, while keeping fixed all other network parameters.

**Implementation Details**. We used this strategy to transfer knowledge from CNNs trained on ImageNet to iCWT. The CNN layers used for our experiments are reported in Tab. 2 for the four architectures considered in the paper. We used RLSC with a Gaussian Kernel as classifier for the CNN-extracted features. In particular we implemented the Nyström sub-sampling approach proposed in (Rudi et al 2015; Rudi and Rosasco 2017), which is computationally appealing for mid and large-scale settings and significantly speeded up our experimental evaluation. We refer to the supplementary material for details on model selection and image preprocessing operations.

*3.2.2 Fine-tuning*

A CNN trained on a large image dataset can also be used to "warm-start" the learning process on other (potentially smaller) datasets. This strategy, known as *fine-tuning* (Chatfield et al 2014; Simonyan et al 2014), consists in performing back-propagation on the new training set by initializing the parameters of the network to those previously learned (see Fig. 4 (Red pipeline)). In this setting it is necessary to adapt the final layer to the new task (e.g., by changing the number of units in order to account for the new number of classes to discriminate).

A potential advantage of fine-tuning is that it allows to adapt the parameters of all layers to the new problem (rather than only those in the final layer). This flexibility however comes at the price of a more involved training process: the choice of the (many) hyper-parameters

available (e.g., which layers to adapt and how strongly) can have a critical impact on performance, especially when dealing with large architectures and smaller training sets. Fine-tuning has been recently used to transfer CNNs learned on ImageNet to robotic tasks (see, e.g., (Eitel et al 2015; Pinto and Gupta 2016; Redmon and Angelova 2015; Pasquale et al 2016a; Nguyen et al 2016)).

**Implementation Details.** We extensively experimented fine-tuning of *CaffeNet* and *GoogLeNet*, for which we report in the paper systematic and statistically robust performance trends on the proposed benchmark. Since *VGG-16* and *ResNet-50* have remarkably longer training times, for these models we performed less systematic experiments, that confirmed analogous trends (which we do not report).

After testing several hyper-parameters settings, we identified two fine-tuning "regimes" as representatives, which we selected for our analysis: one updating only fully-connected layers, while keeping convolution layers fixed, and another one more aggressively adapting all layers to the training set. We refer to these two protocols as *conservative* and *adaptive* and report their corresponding hyper-parameters in Tab. 3 (the two are characterized by different learning rates for each network layer).

We refer to the supplementary material for our analysis of fine-tuning regimes and the model selection protocols implemented in our fine-tuning experiments.
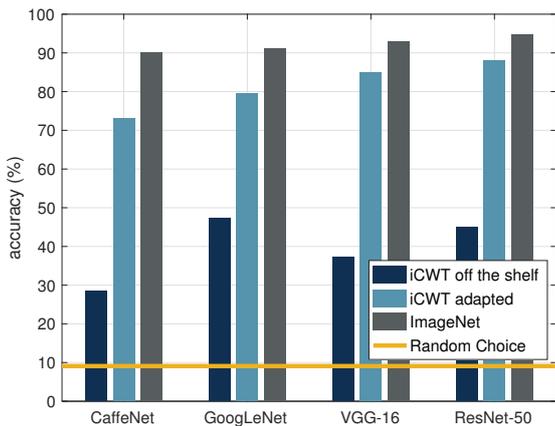
**Fig. 5** Average classification accuracy of off-the-shelf networks (trained on ILSVRC) tested on iCWT (Dark Blue) or on ImageNet itself (Gray). The test sets for the two datasets are restricted to the 11 shared categories (see Sec. 2). (Light Blue) reports the classification accuracy when the same networks are "transferred" to iCubWorld (see Sec. 4.1 for details). The (Orange line) shows the recognition chance of a random classifier.

## 4 Results on Object Categorization

In this Section we present our empirical investigation of deep learning methods in the robotic setting of iCubWorld.

### 4.1 Deep Learning and (the Need for) Knowledge Transfer

Modern datasets for visual recognition, as the ILSVRC, comprise million images depicting objects in a wide range of natural scenes. This extreme variability opens the question of whether datasets such as iCWT, which represent a smaller "reality", could be interpreted as their sub-domains. If this was the case, deep models trained on ImageNet would achieve high performance on iCWT as well.

To address this question, we evaluated four off-the-shelf CNNs (see Sec. 3.1) for the task of image classification on iCWT. For this experiment we restricted the test set to the 11 categories of iCWT that appear also in the ILSVRC (see Sec. 2 and Fig. 1). We compared these results with the accuracy achieved by the same models on the corresponding 11 categories of the ImageNet dataset. The test set for iCWT was composed, for each category, by the images of all 10 object instances, including all 5 transformations for one day and the left camera (unless differently specified, we always used this camera for the experiments), for a total of $\sim 9000$ images per category. For testing on Ima-

geNet, we downloaded the images of the corresponding 11 categories (refer to the supplementary material for the synset IDs), comprising on average $\sim 1300$ images per category.

Fig. 5 reports the average classification accuracy on iCWT (Dark Blue) and ImageNet (Gray). It can be observed that there is a substantial $\sim 50-60\%$ performance drop when testing on iCWT. While a detailed and formal analysis of cross-domain generalization capabilities of deep learning models is outside the scope of this work (see, e.g., (Tommasi et al 2015; Hoffman et al 2013; Rodner et al 2013)), in the supplementary material we qualitatively provide some interesting evidence of this effect. Note that in this case we have restricted the 1000-dimensional output vector provided by the off-the-shelf CNNs, to the considered 11 classes (we report in the supplementary material the same results when considering the entire 1000-dimensional prediction).

**Knowledge Transfer.** The performance in Fig. 5 shows that all models performed much better than chance (Orange line), suggesting that the networks did retain some knowledge about the problem. Therefore, it seems convenient to transfer such knowledge, rather than training an architecture "from scratch", essentially by "adapting" the networks trained on ImageNet to the new setting (see Sec. 3.2).

In Fig. 5 we report the classification performance achieved by models where knowledge transfer has been applied (Light Blue). For these experiments we followed the protocol described in Sec. 3.2.1, where RLSC predictors are trained on features extracted from the deeper layers of the CNNs. We created a training set from iCWT by choosing 9 instances for each category for training and keeping the 10th instance for testing. We repeated the experiment for 10 trials in order to allow each instance of a category to be used in the test set and Fig. 5 reports the average accuracy over these trials. We observe a sharp improvement for all networks, which achieve a remarkable accuracy in the range of $\sim 70-90\%$. While performance on ImageNet is still higher, such gap seems to be reduced for more recent architectures. What are the reasons for this gap? In the following we empirically address this question.

### 4.2 Do we need more data?

While knowledge transfer can remarkably reduce the amount of data needed by deep CNNs in order to learn a new task, the size and richness of the training set remains a critical aspect also when performing transfer learning. A common practice to train deep CNNs in computer vision is to artificially augment the example
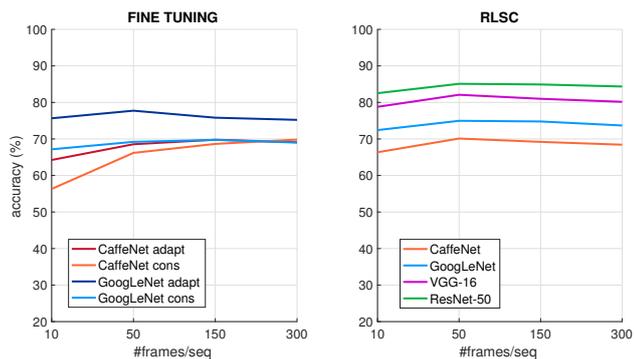
**Fig. 6 Recognition accuracy vs # frames**. (Left) Accuracy of *CaffeNet* and *GoogLeNet* models fine-tuned according to the *conservative* and *adaptive* strategies (see Sec. 3.2.2). (Right) Accuracy of RLSC classifiers trained over features extracted from the 4 architectures considered in this work (see Sec. 3.2.1).

**Fig. 7 Recognition accuracy vs # instances** (number of object instances available during training). (Left) Accuracy of *CaffeNet* and *GoogLeNet* models fine-tuned according to the *conservative* and *adaptive* strategies (see Sec. 3.2.2). (Right) Accuracy of RLSC classifiers trained over features extracted from the 4 architectures considered in this work (see Sec. 3.2.1).

images by applying synthetic transformations like rotations, reflections, crops, illumination changes, etc. From this perspective, in a robotic setting data augmentation can be achieved by simply acquiring more images ("frames") depicting an object, while viewpoint or illumination change naturally. On the other hand, in a typical robotic application it is expensive to gather many different object instances to be shown to the robot.

In this Section, we investigate the impact of these aspects in robot vision, taking iCubWorld as a testbed. Note that in the following we use the term *instance* to refer to a specific object belonging to a given category, while *frame* denotes a single image depicting an object.

### 4.2.1 What do we gain by adding more frames?

To compare the performance of models trained on an increasing number of frames, we consider a 15-class categorization task on iCWT. For each category (see the supplementary material for their list) we used 7 object instances for training, 2 for validation and 1 for testing. We created training sets of increasing size by sampling randomly $N = 10, 50, 150$ and finally 300 frames from each image sequence of an object (we recall that each object in iCWT is represented by 10 sequences containing 5 isolated visual transformations acquired in 2 days). Validation and test sets contained all images available for the corresponding instances. To account for statistical variability, we repeated the experiment for 10 trials, each time leaving out a different instance for testing. For this experiment, we considered only one of the two available days in iCWT. We used only the left camera, apart from when sampling 300 frames, where we drew images also from the right camera. The number of frames per category therefore ranged from 350 to 10500.
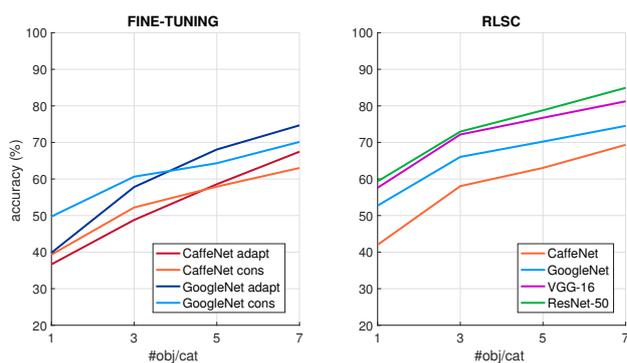
Fig. 6 reports the average classification accuracy of different models as more example frames are provided. Surprisingly, most architectures achieve high accuracy already when trained on the smallest training set and show little or no improvement when new data is available. This finding is in contrast with our expectations, since increasing the dataset size does not seem key to a significant improvement in performance. To further support this finding, in the supplementary material we report results for the same experiment when using less example instances per category. Moreover, in Sec. 6.1 we report similar observations on the Washington RGB-D dataset.

Secondary observations:

- Fine-tuning and RLSC achieve comparable accuracy (both for *CaffeNet* and *GoogLeNet*).
- We confirm the ILSVRC trends, with more recent networks generally outperforming older ones.
- *CaffeNet* performs worse when training data is scarce because of the high number of parameters to be learned in the 3 fully connected layers (see Sec. 3 and the supplementary material).

### 4.2.2 What do we gain by adding more instances?

We evaluated then the impact of experiencing less or more example object instances per category. We consider this process as increasing the *semantic* variability of the training set, in contrast to increasing the *geometric* variability by showing more frames of a given object. To this end, we kept the same 15-class categorization task of Sec. 4.2.1 and created multiple training sets each containing an equal number of 900 frames per category, but sampled from an increasing number of instances per category, namely 1, 3, 5 and 7 (i.e., we took
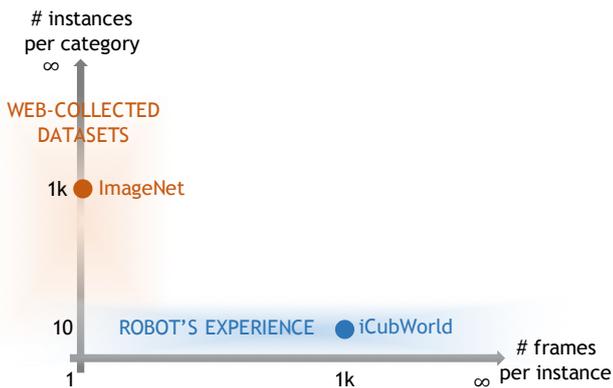
**Fig. 8** Different training regimes in robot vision and image retrieval settings.

all 900 frames from one instance per category, then 300 frames from each of the 3 instances per category, and so forth). Validation and test sets are as in the previous experiment and we repeated again the experiment for 10 trials.

Results in Fig. 7 show that increasing the semantic variability dramatically improves the accuracy of *all* models by a similar margin (more than 20%), and performance does not saturate at 7 example instances per category. To further support this finding, in the supplementary material we report results for the same experiment when discriminating between even less categories (10 and 5) and in Sec. 6.1 we show a similar effect on the Washington RGB-D dataset.

Finally, it is worth noting that, when few example objects per category are available, *adaptive* fine-tuning (Dark Blue and Red) provides worst performance, suggesting that adapting ImageNet features (which are "optimized" for a rich categorization task) is not convenient if the dataset has poor semantic variability, as we also discuss in Sec. 4.2.3.

### 4.2.3 Robot Vision and Image Retrieval

In this Section we considered two major aspects associated to visual object recognition: semantic and geometric variability of objects' appearance. We observed that, while essentially the same problem is addressed both in robot vision and image retrieval, it is cast within two different regimes (see Fig. 8). In both cases it is possible to gather a large amount of image examples. However, typically, in image retrieval, each image depicts a different object instance (as an example, ILSVRC training set comprises $\sim 1000$ images per category, 1 image per instance). On the opposite end of the spectrum, in robot vision it is easy to gather many images of an object (the robot can observe it from multiple viewpoints),

but there is a remarkable limitation in the number of instances that can be experienced.

Our analysis has shown that the limited semantic variability that characterizes our setting dramatically reduces the recognition accuracy, even when relying on pre-trained deep learning models. Moreover, contrarily to our expectations, we observed that feeding more object views to the network (i.e., increasing the geometric variability) does not alleviate this lack: classifiers trained on very few, or many, object views, and fixed semantic variability, achieved identical performance.

This represents a problem, because usually a robot has access to limited instance examples of categories to be learned. In this perspective, adopting "data augmentation" strategies to artificially increase semantic variability, could be a viable solution. In Sec 8 we discuss how this problem may be addressed in the future, while in the following we focus our analysis on the invariance and robustness of deep representations to geometric (viewpoint) transformations.

### 4.3 Invariance to Viewpoint Transformations

Invariance, i.e., robustness to identity-preserving visual transformations, is a desirable property of recognition systems, since it increases the capability of generalizing the visual appearance of an object from a limited number of examples (ideally, just one) (Anselmi et al 2016a,b).

In this Section, we investigate to what extent CNN models are invariant to the viewpoint transformations represented in iCWT.

We considered the same 15-class categorization problem introduced in Sec. 4.2.2, where we use 7 instances per category for training, 2 for validation and 1 for testing. However, in this case we did not mix example images from all the 5 available sequences (*2D Rotation*, *3D Rotation, Scale, Background* and *Mix*). Instead, we performed training (and validation) using only an individual transformation, and then tested the model on the others. We considered two different training set sizes: a larger one, with including all images from each sequence (i.e. $\sim 150 \times 7 = 1050$ images per category), and a smaller one, subsampling images of a factor of 7 (i.e. 150 images per category), in order to reproduce the two extreme sampling conditions as in Fig. 6. As a reference, we considered a 6-th training set (*All*), of same size, obtained by randomly sampling images from the other 5 training sets. This training set is analogous to those used in the experiments in Sec. 4.2.2 and 4.2.1. We consider only one of the two available days in iCWT, since we aim to exclude nuisances due to illumination
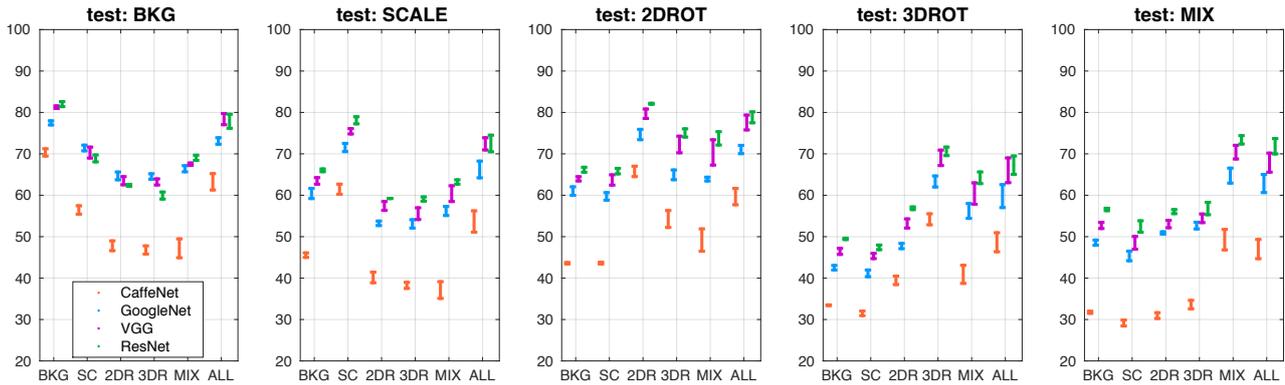
**Fig. 9 Generalization performance across different visual transformations**. Models were trained on one of the 5 transformations in iCWT (horizontal axis), and then tested on all transformations (title of the subplot). Accuracy is reported separately in each plot for each tested transformation. The bars indicate the (small) performance improvement achieved by including all frames of a sequence (higher end of bars) instead of subsampling ∼ 20 of them (lower end of bars).

or setting changes that can happen from one day to another. We repeated the experiment for 10 trials as in previous tests.

Since in this Section we focus on the invariance of representations extracted from networks trained on ImageNet, we report only the accuracy of the approach based on training RLSC on off-the-shelf features (see Sec. 3.2.1). We refer to (Pasquale 2017) for additional results reported by fine-tuning strategies. Fig. 9 shows the generalization capabilities of models trained on a single transformation (indicated on the horizontal axis) and tested on a different one (indicated in title of the subplot). Each bar starts at the accuracy achieved by the small training set and goes up to the accuracy achieved by the large training set, from which we suddenly note that there is no improvement by considering all frames from the sequence of a transformation, even on the transformation itself (i.e., all bars are very short). This explains and confirms the trend observed in Fig. 6.

Overall, the classifiers perform well only on transformations that have been included – even with a few examples – in the training set: best performance is always achieved when training and test set include the same transformation, or when all transformations are included in the training set (*All*). While generalization failure from 2D to 3D rotations of the object is expected, it is quite surprising that generalization also fails between affine transformations, namely *Scale*, *2D Rotation* and *Background*, which the CNN could have learned from ImageNet. We will investigate this aspect in Sec. 5.2, by studying the invariance of CNNs in the context of object *identification*.

It is finally worth noting that training on the *Mix* sequence achieves considerably good performance when

tested on every specific transformation. This suggests that showing an object to the robot in a natural way, with transformations appearing in random combinations (instead of systematically collecting sequences comprising individual transformations), is a good approach to obtain predictors invariant to these transformations.

## 5 Results on Object Identification

Object identification is the task of discriminating between specific instances within a category, and it is clearly of paramount importance for robot to correctly interact with the environment. In this Section we evaluate the performance of deep learning models on this task using iCubWorld.

The problem has been largely addressed with methods based on keypoints and template matching (Lowe 2004; Philbin et al 2008; Collet et al 2009, 2011a,b; Muja et al 2011; Crowley and Zisserman 2014). However, it has recently been observed that approaches that rely on holistic visual representations perform typically better in scenarios characterized by substantial variations of objects' appearance (Ciliberto et al 2013). Following this, we focus the analysis on the methods considered in Sec. 3.

### 5.1 Knowledge Transfer: from Categorization to Identification

We investigated knowledge transfer in the context of object identification, specifically to determine to what extent performance are affected by the number of frames per object. We addressed this question by building an experimental setting similar to the one investigated for
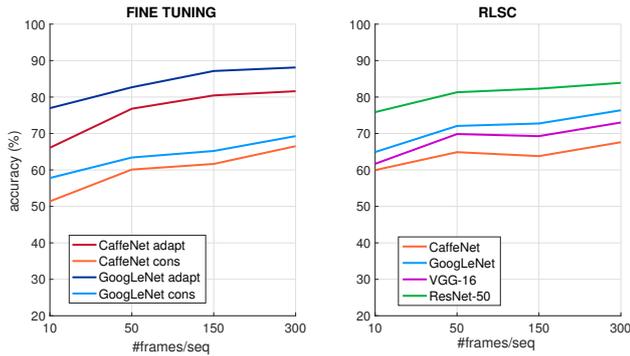
**Fig. 10 Recognition accuracy vs # frames (Identification)**. (Left) Accuracy of *CaffeNet* and *GoogleNet* models fine-tuned according to the *conservative* and *adaptive* strategies (see Sec. 3.2.2). (Right) Accuracy of RLSC classifiers trained over features extracted from the 4 architectures considered in this work (see Sec. 3.2.1).

categorization, i.e., considering tasks on iCWT with increasing number of example images. We selected 50 objects from the dataset by taking all instances from the *book*, *flower*, *glass*, *hairbrush* and *hairclip* categories. Note that these categories do not appear in the ILSVRC. We created four training sets containing respectively $10, 50, 150$ and $300$ images per object, sampled randomly from the 4 transformation sequences *2D Rot, 3D Rot, Scale* and *Bkg* (see Sec. 2). From each training set, 20% images were retained for validation. The images from the *Mix* sequence were used to test the classification accuracy of the methods. As for the categorization experiment, only the images from a single day were used.

Fig. 10 reports the average accuracy of models trained on a growing number of example images per object. Differently from what observed for categorization (Sec. 4.2.1), in this case adding images of an object greatly improves performance (with a 15-20% margin). We also notice that, differently from categorization, in this case the *adaptive* fine-tuning (Dark Blue and Red) significantly outperforms the other strategies (RLSC and *conservative* fine-tuning). This finding confirms recent empirical evidence (Sharif Razavian et al 2014) that adapting CNNs trained on ImageNet (rather than "from scratch") is beneficial also in identification tasks. It also confirms and extends results from the instance retrieval literature (Babenko et al 2014; Gordo et al 2016), which show that several, different, images of objects are necessary for such adaptation.

## 5.2 Invariance to Viewpoint Transformations

Similarly to Sec. 4.3, we investigated object identification on isolated viewpoint transformations. To this end, we kept the same 50-class identification task of Sec. 5.1 and restricted the training set to contain only images from a single transformation sequence among the 5 available in iCWT. The resulting models were tested separately on the remaining transformations to assess their ability to generalize to unseen viewpoints. As in Sec. 4.3 we considered one day and evaluated only methods based on training RLSC on top of off-the-shelf features.

Fig. 11 reports the accuracy of models trained on a single transformation, using respectively 10 frames per object (lower end of bars) or 150 (upper end of bars), and tested on the others. It can be noticed that the improvement observed in Fig. 10 is confirmed and clarified in this setting. Adding images containing varied transformations (that is, from the *Mix* training set) helps to generalize to the individual transformations. In addition, adding images from a single transformation always improves performance on the others (also the profile of the upper ends of the bars is slightly flatter, i.e., with smaller gaps between one transformation to another).

We conclude that the viewpoint invariance of off-the-shelf representations is small, but adding more views of an object can remarkably boost it. While this is positive, in real world applications, where a robot typically has to learn novel objects on the fly, collecting extensive object views is probably unfeasible. To this end, in the following Section we report on a simple strategy we proposed to increase the invariance of CNNs and hence reduce the number of example frames needed to learn new objects "online".

### 5.2.1 Improving Viewpoint Invariance of CNNs

The performance improvement achievd by the *adaptive* fine-tuning strategy (Fig. 10) may be explained by the fact that adapting the inner layers of CNNs with images representing viewpoint variations improves the invariance of the internal network representation to such transformations.

In this Section we investigate this effect in detail. We consider the same learning setting of Fig. 11, i.e., RLSC trained on few example frames of a single visual transformation. Instead of off-the-shelf features, we compare using features from different *CaffeNet* or *GoogLeNet* models, previously fine-tuned on a selected image set from iCWT. Specifically, we consider the following fine-tuning strategies:

- **iCubWorld identification (iCWT id)**. This set contains all transformation sequences available in iCWT for a number of objects. It is conceived to investigate if the CNN can learn to be invariant
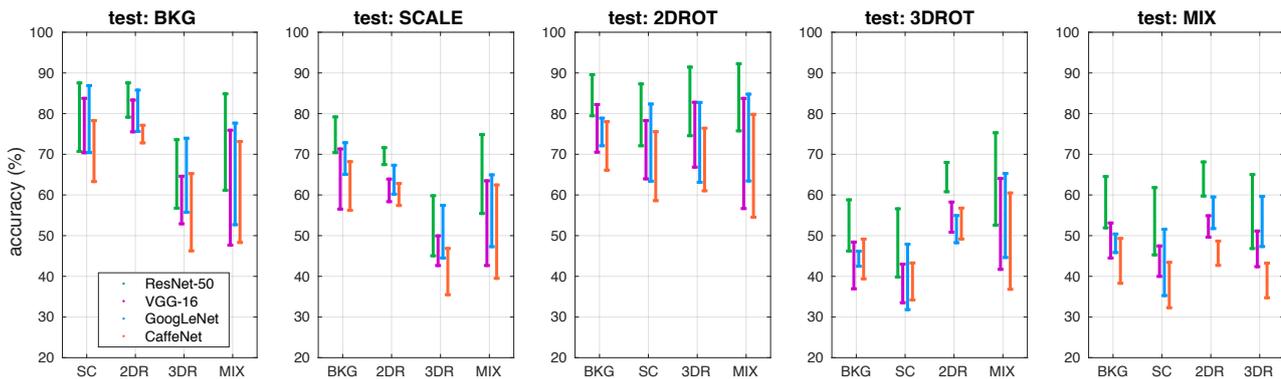
**Fig. 11 Generalization performance across different visual transformations (Identification).** RLSC are trained on one of the 5 transformations in iCWT (horizontal axis), and then tested on the other transformations (title of the subplot). Accuracy is reported separately in each plot for each tested transformation. The bars indicate the performance improvement achieved by including all frames of a sequence (higher end of bars) instead of subsampling ∼ 10 of them (lower end of bars)

.

to experienced viewpoint changes. We used objects belonging to categories not involved later in the considered 50-class identification task (*oven glove*, *squeezer*, *sprayer*, *body lotion* and *soda bottle*). Fine-tuning was performed with the *adaptive* strategy for an identification task. The validation set was obtained following the protocol of Sec. 5.1.

- **iCubWorld categorization (iCWT cat)**. Same set as (iCWT id), but fine-tuning performed on a 5-class categorization task. The validation set was obtained following the protocol of Sec. 4.2.1. This set is conceived to test the same hypothesis of learning viewpoint invariance, but at the category rather than instance level.

- **iCubWorld + ImageNet (iCWT + ImNet)**. Same as (iCWT cat), but images of the 5 categories are sampled also from ImageNet. Note that most iCWT categories do not appear in ILSVRC but are contained in synsets in the larger ImageNet dataset (see supplementary material for the synset list). This set is conceived to investigate if learning category-level viewpoint invariance jointly with semantic variability can be beneficial.

- **ImageNet (ImNet)**. This dataset differs from previous ones and contains the 5 ImageNet synsets corresponding to the 5 categories of the 50 objects on which the CNN would be later used as feature extractor (*book*, *flower*, *glass*, *hairbrush* and *hairclip*, see Sec. 5.1). Fine-tuning was performed on the categorization task. This set is conceived to investigate the effect of focusing the CNN on the categories of the objects to be discriminated later on, by learning category-level features from data available on-line.

Fig. 12 reports the accuracy of RLSC trained on features from *CaffeNet* (top, Orange) or *GooLeNet* (bot-

tom, Blue) fine-tuned on each of these image sets. Training was performed identically to Sec. 5.1, on the smaller set containing 10 examples per instance. It can be noticed that preliminary fine-tuning on sequences available in iCWT for an identification task (iCWT id) is particularly effective, leading to the highest improvement over off-the-shelf features. We also observe that all approaches that involve preliminary fine-tuning over transformation sequences in iCWT do provide an accuracy increase. Interestingly, fine-tuning on corresponding categories in ImageNet (ImNet), degrades performance, possibly since it has increased invariance to intra-class variations between instances, which are in fact relevant for the identification task.

Comparing with Fig. 11, it can be observed that RLSC trained on (iCWT id) features from 10 examples per object now performs better or on par with off-the-shelf features from 150 examples per object. The proposed strategy is a viable approach to build robust feature extractors: training, possibly offline, a CNN to learn invariances by collecting images of objects undergoing various transformations and then using the resulting features online to train new classifiers using only few images.

**Note on Object Categorization.** We repeated the experiment of Fig. 9 using features from these pre-fine-tuned CNNs. However, our results (reported in supplementary material) showed that, in this setting, performance does not improve with respect to off-the-shelf features. This futrther confirms our findings in Sec. 4.2.3, that networks trained on ILSVRC are highly optimized for categorization and adapting them to datasets like iCWT does not provide a clear advantage.
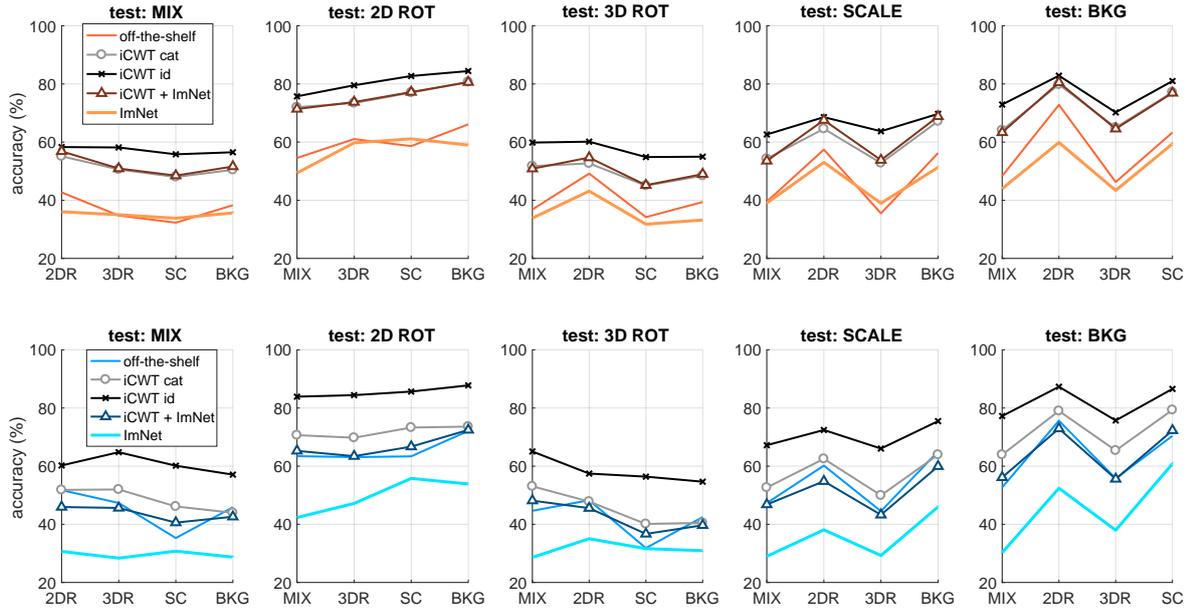
**Fig. 12** Same experiment setting as in Fig. 11, but using different image representations, provided by *CaffeNet* (top, Orange) or *GoogLeNet* (bottom, Blue) network models, previously fine-tuned according to different strategies (see Sec. 5.2.1).

## 6 Results on Washington RGB-D

In this Section we show that the conclusions drawn from our results on iCWT are supported by similar trends on other robotic datasets. Specifically, we consider Washington RGB-D (WRGB-D for short) (Lai et al 2011).

WRGB-D is a turntable dataset comprising 300 objects organized into 51 categories. The number of objects per category ranges from 3 to 14 (6 on average). For each object, 3 RGB-D sequences were acquired by a camera mounted at $\sim$ 1m from the object and $30, 45, 60°$ elevation angles. In each sequence, the camera performed one revolution around the table at constant speed, recording $\sim$ 250 $640 \times 480$ frames. In the literature (and in our evaluation) a reduced dataset version is usually employed, including every fifth frame, cropped to tightly include the object, for a total of 41877 images.

### 6.1 Object Categorization Benchmark

The benchmark categorization task consists in discriminating between the 51 categories, by leaving out one instance per category for testing. We considered this task in our experiments and, as in the literature, we averaged performance on the same 10 trials released by the authors (Lai et al 2011), each time leaving out different 51 instances for testing. We used a random 20% of the training set for validation.

As a first sanity check, we evaluated our methods on this benchmark: Table 4 compares our results with two recent works that achieve the state-of-the-art by employing *CaffeNet* with similar transfer learning approaches. In particular, in (Schwarz et al 2015) SVM classifiers are trained on off-the-shelf features, while in (Eitel et al 2015) the network is fine-tuned. Note that we consider results exploiting only RGB, without depth information. It can be observed that, with *CaffeNet*, our methods achieve comparable results, and the best performance is achieved by *ResNet-50* features.

We then repeated this experiment by increasing either the number of (i) example instances per category or (ii) example images per object. The goal of this evaluation is to reproduce respectively the setting of Fig. 7 and Fig. 6.

When increasing the number of example instances per category, we fixed the overall number of example images such that, similarly to Sec. 4.2.2, we first sampled all example images for a category from one instance (around 150), then 50% images from each of two instances, and so forth, until all instances from all categories were included. Since, differently from iCWT, in WRGB-D each category has a different number of available training instances (from 2 to 13 since one is left for testing), for each category we stopped adding instances when all were included. Specifically, while most categories have at least 3 training instances available, we could train only 25 categories (out of 51) on 5 in-

**Table 4 Categorization benchmark on WRGB-D**: Comparison with state-of-the-art.

| Method | Accuracy (%) |
|---|---|
| Schwarz et al. 2015 | $83.1 \pm 2.0$ |
| Eitel et al. 2015 | $84.1 \pm 2.7$ |
| CaffeNet RLSC | $83.1 \pm 2.8$ |
| CaffeNet FT adapt | $81.9 \pm 2.7$ |
| CaffeNet FT cons | $83.6 \pm 2.4$ |
| GoogLeNet RLSC | $84.6 \pm 2.8$ |
| GoogLeNet FT adapt | $82.4 \pm 2.8$ |
| GoogLeNet FT cons | $83.9 \pm 2.3$ |
| VGG-16 RLSC | $87.5 \pm 2.3$ |
| **ResNet-50 RLSC** | **$89.4 \pm 3.1$** |

**Table 5 Identification benchmark on WRGB-D**: Comparison with state-of-the-art.

| Method | Accuracy (%) |
|---|---|
| Schwarz et al. 2015 | 92.0 |
| Held et al. 2016 | 93.3 |
| CaffeNet RLSC | 94.0 |
| CaffeNet FT adapt | 94.0 |
| CaffeNet FT cons | 92.7 |
| GoogLeNet RLSC | 94.3 |
| GoogLeNet FT adapt | 93.9 |
| GoogLeNet FT cons | 92.5 |
| VGG-16 RLSC | 94.5 |
| **ResNet-50 RLSC** | **96.0** |

stances, and just 10 (5) categories on 7 (or more) instances. Fig. 13(a) reports results of this experiment.

Conversely, Fig. 13(b-e) report results of experiments where we fixed the maximum number of example instances per category to different values (1, 2, 3 and *all*) and, for each value, we repeated the experiment multiple times by progressively increasing the number of sampled frames per instance (starting from 10 to *all* available).

For this evaluation, we considered the three transfer learning techniques based on *CaffeNet*, reported in Fig. 13 in Orange (Dashed, RLSC; Bold, *conservative* fine-tuning) and Red (Bold, *adaptive fine-tuning*). We left out the other methods because we assessed how they compare both on iCWT and on the reference RGB-D categorization task.

These results replicate the findings of iCWT: adding new instances (slope in Fig. 13(a)) leads to mich higher performance improvement than than adding example views (slopes in Fig. 13(b-e)). In this latter case a "jump" in performance is achieved only when a new instance is added to the training set. Note that best results are achieved in Fig. 13(a), when including all training instances, with just $\sim 1/6$ of the frames of those usually used in the reference benchmark (which corresponds to the *all* training set in Fig. 13(e)). We point out that performance saturates at 5 instances per category in Fig. 13(a) because, as explained above, in WRGB-D there are very few categories with more instances to be included.

## 6.2 Object Identification Benchmark

The identification task considered in the literature on WRGB-D (Lai et al 2011), consists in discriminating between the 300 objects in the dataset. For each object, the sequences recorded at 30° and 60° elevation angles are used for training, while the one recorded at 45° for

testing. We considered this task, using a random 20% of the training set for validation.

We first evaluated our methods on this benchmark: Table 5 compares our results with the state-of-the-art achieved in (Schwarz et al 2015) and in (Held et al 2016). This latter is recent paper that aims to improve performance on object identification by exploiting deep architectures (specifically, *CaffeNet*) previously fine-tuned on "multi-view" image sequences. Here we refer to the accuracy that they report by applying a fine-tuning protocol similar to our *conservative* strategy. From Table 5 it can be noticed that our pipeline provides better results, when relying on *CaffeNet*, and best performance is achieved again by *ResNet-50* features.

Similarly to object categorization, we repeated the experiment by progressively increasing the number of example images per instance (from 10 to *all* available). From the results reported in Fig. 13(f), it can be noticed that, as for iCWT, and differently from the categorization setting, the performance gain provided on this task by adding object views to the training set is remarkable. In this case, the accuracy reported in the reference benchmark is achieved only when using all training images.

## 7 Discussion and Future Work

In this work we reported on an extensive experimental validation of the application of latest deep learning methods to visual object recognition in robotics. We challenged these methods on a setting that was specifically designed to represent a prototypical real world scenario. Our results showed that deep learning leads to remarkable performance for both category and instance recognition. We showed also that proper adoption of knowledge transfer strategies – in particular mixing deep learning with "shallow" classifiers – plays a key role, in that they leverage on the visual representation

**Categorization**                                                                   **Identification**
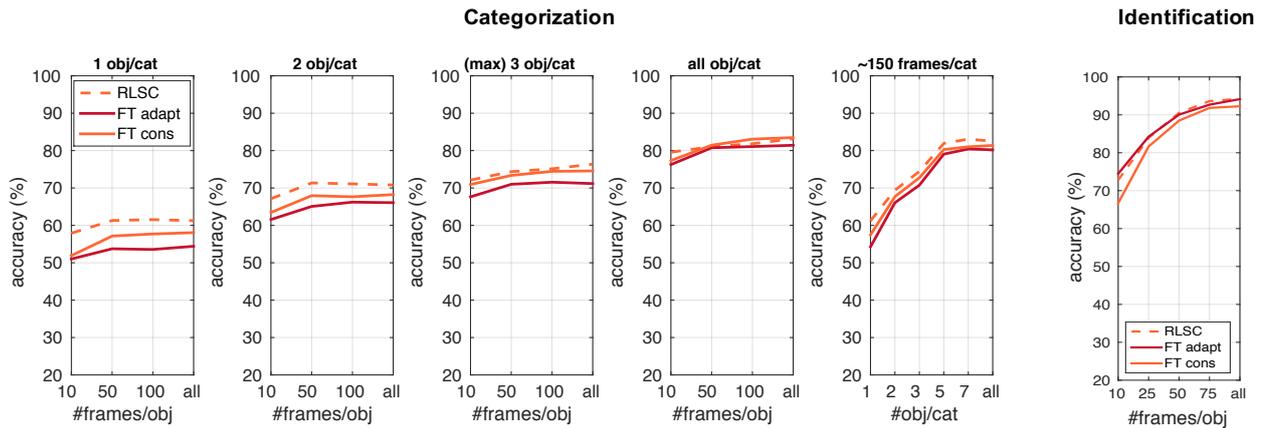


**Fig. 13 Categorization and Identification on WRGB-D varying # frames and # instances.** (a) Object Categorization varying # instances per category (horizontal axis). (b-e) Object Categorization varying # frames per instance (horizontal axis), for different fixed numbers of instances per category (plot title). (f) Object Identification varying the number of #frames per instance (horizontal axis). The three transfer learning methods indicated in the legend are applied to *CaffeNet*.

learned on large-scale datasets to achieve high performance in the robotic application.

However, a substantial gap still exists between the performance that can be obtained in the two domains. Our analysis shows that one reason for this gap is limited semantic variability of data collected in robotics settings (due to the intrinsic cost of acquiring training examples). Moreover, we need to push further these requirements since the error rate of robotic recognition systems will need to be as close as possible to zero in order to be considered for production and deployment in real applications.

In this Section, we consider directions for future research to address these limitations and present how, in our opinion, the visual recognition problem could be addressed in robotics.

**Improving Invariance**. Our experiments in Sec. 4.3 and 5.2 showed that off-the-shelf CNNs are mostly *locally* invariant, i.e., their representation is robust to small viewpoint changes, but these models can learn specific invariances from image sequences that are easy to collect in robotics settings. While local invariance has been of main interest to computer vision in the past (Lowe 2004; Mikolajczyk and Schmid 2004), current research on invariance is focused on learning representations robust to more "global" visual transformations (Anselmi et al 2016a,b). Our results confirm that improving invariance is crucial to perform visual recognition in the real world (Pinto et al 2008, 2011; Borji et al 2016; Poggio and Anselmi 2016). The role of viewpoint invariance in object categorization, which evidenced unexpected results in our setting, must also be further investigated (to this end see also (Zhao et al

2016; Zhao and Itti 2016)).

**"Augmenting" Semantic Variability: Real and Synthetic Data**. Our results pointed out how this aspect is key to object categorization. The simplest method for increasing semantic variability is to share data acquired from different robot platforms. This strategy has been used to learn hand-eye coordination on a manipulator (Levine et al 2018). Similarly, the goal of the Million Object Challenge (Oberlin et al 2015) is to create a sharing platform for data acquired by laboratories owning a Baxter robot[9]. Along a similar direction, we plan to extend iCubWorld with the help of the community of the iCub robot (more than 30 research groups worldwide). This could also allow to extend the analysis presented in this work by introducing much larger variability. Indeed, this is another critical aspect that we started to address in the supplementary material.

A complementary approach follows the idea of data augmentation. More or less sophisticated synthetic image transformations (e.g., 2D rotation, flip, crop/scaling, background and illumination changes) are already standard practice to simulate variations of instance appearance. Visual augmentation to cope for semantic variability is a more challenging problem, although recent work on inverse graphics (Mansinghka et al 2013; Kulkarni et al 2014) is a starting point in this direction. The potentials of this approach are limitless since, by synthetic generation, objects and environment parameters (viewpoint, lighting, texture, etc.) can be endlessly created and tuned to the application requirements. Generalization to real conditions would then clearly depend on the realism of simulated data, but the domain

---

[9] http://www.rethinkrobotics.com/

shift could be eventually tackled with transfer learning approaches. The work of (Handa et al 2016; Zhang et al 2017) for indoor scene recognition and ShapeNet datasets (Chang et al 2015; Wu et al 2015) datasets are examples of this strategy.

**Integrating 3D Information**. The integration of depth and RGB information with deep CNNs has been recently proved to help object recognition in robotics (Schwarz et al 2015; Eitel et al 2015; Redmon and Angelova 2015; Carlucci et al 2016). The spatial structure of the scene can also be exploited as an additional self-supervisory signal, or prior information, to make the prediction more robust. Examples are the work of (Song et al 2015), where memory about the room helps discriminating between old and novel objects. Interestingly, (Pillai and Leonard 2015) implement a SLAM-aware detection system where the object's 3D position is projected back to frames in order to help the recognition. With this approach, object instances could also be autonomously discovered by the robot.

**Exploiting Temporal Coherence**. While it is important to push the limits of visual recognition at the image level, a robot is typically exposed to a continuous stream of frames, where visual information is correlated. Accuracy may be remarkably improved by exploiting this correlation, ranging from simple solutions as temporal averaging of predictions (Pasquale et al 2015) to more complex architectures including recent recurrent networks (Donahue et al 2015). The temporal correlation among consecutive frames can be exploited as self-supervisory signal, too. Recent works exploit this information to learn visual representations in absence of frame-level annotations (Wang and Gupta 2015; Goroshin et al 2015b,a; Jayaraman and Grauman 2015; Agrawal et al 2015).

**Self-supervised Learning**. We opted for the help of a "teacher" for the acquisition of iCWT because we needed a relatively fine control on the object movements to isolate viewpoint transformations. However, making the acquisition self-supervised, i.e., implementing explorative strategies through which the robot autonomously interacts with the environment to collect examples, would allow to extend iCubWorld datasets, while limiting human effort. Training instances in this scenario could be gathered autonomously by detecting invariances in the data which correspond to physical entities (e.g., coherent motion patterns (Wang and Gupta 2015) or bottom-up saliency cues). Strategies specific to the robotic domain could be devised by integrating multiple sensory modalities (Sinapov et al 2014; Higy et al 2016) and a repertoire of explorative actions (Montesano et al 2008; Fitzpatrick et al 2003; Högman et al 2016; Pinto et al 2016)).

**Multi-task Learning**. In this work we adopted a standard classification approach where no similarity or relation is assumed among classes. However, objects and categories have common features (think to object parts like wheels, legs, and so forth). Incorporating information about class similarities can significantly improve performance, especially if training data is limited. In the literature on multi-task learning several approaches have been proposed to enforce these relations on the learning problem, when they are available a-priori (Evgeniou et al 2005; Joachims et al 2009; Fergus et al 2010), or to learn them, when unknown (Argyriou et al 2008; Minh and Sindhwani 2011; Dinuzzo et al 2011; Ciliberto et al 2015).

**Incremental Learning**. A robotic recognition system should guarantee reliability in lifelong scenarios, by learning to adapt to changing conditions. In line with the literature on "learning to learn" and transfer learning (Thrun and Mitchell 1995; Tommasi et al 2010; Kuzborskij et al 2013), strategies to exploit knowledge of previous, well-represented classes, in order to improve prediction accuracy on novel, under-represented ones, ultimately would be a key component of robotic recognition systems to be deployed in real world applications, as observed in recent work (Camoriano et al 2017; Sun and Fox 2016).

## 8 Conclusions

In this paper we presented a systematic experimental study on the application of deep learning methods for object recognition to robot vision settings. For our tests we have devised a prototypical vision task for a humanoid robot in which human-robot interaction is exploited to obtain realistic supervision and train an object recognition system. We presented the iCWT dataset and an in-depth investigation of the performance of state-of-the-art CNNs applied to our scenario. Results confirm deep learning is a remarkable step forward. However, there is still a lot that needs to be done to reach the level of robustness and reliability required by real world applications. We identified specific challenges and possible directions of research to bridge this gap. We are confident that the next few years will be rich of exciting progress in robotics.

## References

Agrawal P, Carreira J, Malik J (2015) Learning to see by moving. In: The IEEE International Conference on Computer Vision (ICCV)

Anselmi F, Leibo JZ, Rosasco L, Mutch J, Tacchetti A, Poggio T (2016a) Unsupervised learning of invariant representations. Theoretical Computer Science 633:112 – 121, DOI https://doi.org/10.1016/j.tcs.2015.06.048

Anselmi F, Rosasco L, Poggio T (2016b) On invariance and selectivity in representation learning. Information and Inference 5(2):134–158

Argyriou A, Evgeniou T, Pontil M (2008) Convex multi-task feature learning. Machine Learning 73

Babenko A, Slesarev A, Chigorin A, Lempitsky V (2014) Neural Codes for Image Retrieval, Springer International Publishing, Cham, pp 584–599. DOI 10.1007/978-3-319-10590-1_38

Baishya S, Bäuml B (2016) Robust material classification with a tactile skin using deep learning. In: IEEE International Conference on Intelligent Robots and Systems

Bishop CM (2006) Pattern Recognition and Machine Learning (Information Science and Statistics). Springer-Verlag New York, Inc., Secaucus, NJ, USA

Borji A, Izadi S, Itti L (2016) iLab-20M: A Large-Scale Controlled Object Dataset to Investigate Deep Learning. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Bottou L (2012) Stochastic Gradient Tricks, vol 7700, Springer, pp 430–445

Camoriano R, Pasquale G, Ciliberto C, Natale L, Rosasco L, Metta G (2017) Incremental robot learning of new objects with fixed update time. In: 2017 IEEE International Conference on Robotics and Automation (ICRA)

Carlucci FM, Russo P, Caputo B (2016) A deep representation for depth images from synthetic data. ArXiv e-prints 1609.09713

Chang AX, Funkhouser T, Guibas L, Hanrahan P, Huang Q, Li Z, Savarese S, Savva M, Song S, Su H, Xiao J, Yi L, Yu F (2015) ShapeNet: An Information-Rich 3D Model Repository. Tech. Rep. arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago

Chatfield K, Simonyan K, Vedaldi A, Zisserman A (2014) Return of the devil in the details: Delving deep into convolutional nets. In: British Machine Vision Conference

Ciliberto C, Fanello SR, Santoro M, Natale L, Metta G, Rosasco L (2013) On the impact of learning hierarchical representations for visual recognition in robotics. In: 2013 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp 3759–3764, DOI 10.1109/IROS.2013.6696893

Ciliberto C, Rosasco L, Villa S (2015) Learning multiple visual tasks while discovering their structure. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 131–139

Collet A, Berenson D, Srinivasa SS, Ferguson D (2009) Object recognition and full pose registration from a single image for robotic manipulation. In: Robotics and Automation, 2009. ICRA '09. IEEE International Conference on, pp 48–55, DOI 10.1109/ROBOT.2009.5152739

Collet A, Martinez M, Srinivasa SS (2011a) The moped framework: Object recognition and pose estimation for manipulation. The International Journal of Robotics Research 30(10):1284–1306, DOI 10.1177/0278364911401765

Collet A, Srinivasay SS, Hebert M (2011b) Structure discovery in multi-modal data: a region-based approach. In: Robotics and Automation (ICRA), 2011 IEEE International Conference on, pp 5695–5702

Crowley E, Zisserman A (2014) The state of the art: Object retrieval in paintings using discriminative regions. In: BMVC

Dansereau DG, Singh SP, Leitner J (2016) Interactive computational imaging for deformable object analysis. In: 2016 International Conference on Robotics and Automation (ICRA)

Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L (2009) ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR09

Dinuzzo F, Ong CS, Gehler P, Pillonetto G (2011) Learning output kernels with block coordinate descent. In: Proceedings of the 28th International Conference on Machine Learning (ICML-11), ACM, New York, NY, USA, pp 49–56

Donahue J, Jia Y, Vinyals O, Hoffman J, Zhang N, Tzeng E, Darrell T (2014) Decaf: A deep convolutional activation feature for generic visual recognition. In: Jebara T, Xing EP (eds) Proceedings of the 31st International Conference on Machine Learning (ICML-14), JMLR Workshop and Conference Proceedings, pp 647–655

Donahue J, Anne Hendricks L, Guadarrama S, Rohrbach M, Venugopalan S, Saenko K, Darrell T (2015) Long-term recurrent convolutional networks for visual recognition and description. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Eitel A, Springenberg JT, Spinello L, Riedmiller M, Burgard W (2015) Multimodal deep learning for robust rgb-d object recognition. In: Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on, pp 681–687, DOI 10.1109/IROS.2015.7353446

Everingham M, Van Gool L, Williams CKI, Winn J, Zisserman A (2010) The pascal visual object classes (voc) challenge. International Journal of Computer Vision 88(2):303–338, DOI 10.1007/s11263-009-0275-4

Everingham M, Eslami SMA, Van Gool L, Williams CKI, Winn J, Zisserman A (2015) The pascal visual object classes challenge: A retrospective. International Journal of Computer Vision 111(1):98–136

Evgeniou T, Micchelli CA, Pontil M (2005) Learning multiple tasks with kernel methods. In: Journal of Machine Learning Research, pp 615–637

Fanello SR, Ciliberto C, Santoro M, Natale L, Metta G, Rosasco L, Odone F (2013) iCub World: Friendly Robots Help Building Good Vision Data-Sets. In: 2013 IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp 700–705, DOI 10.1109/CVPRW.2013.106

Fergus R, Bernal H, Weiss Y, Torralba A (2010) Semantic label sharing for learning with many categories. European Conference on Computer Vision

Fitzpatrick P, Metta G, Natale L, Rao S, Sandini G (2003) Learning about objects through action - initial steps towards artificial cognition. In: 2003 IEEE International Conference on Robotics and Automation, pp 3140–3145 vol.3, DOI 10.1109/ROBOT.2003.1242073

Geusebroek JM, Burghouts GJ, Smeulders AW (2005) The amsterdam library of object images. International Journal of Computer Vision 61(1):103–112, DOI 10.1023/B:VISI. 0000042993.50813.60

Goodfellow I, Lee H, Le QV, Saxe A, Ng AY (2009) Measuring invariances in deep networks. In: Bengio Y, Schuurmans D, Lafferty JD, Williams CKI, Culotta A (eds) Advances in Neural Information Processing Systems 22, Curran Associates, Inc., pp 646–654

Gordo A, Almazán J, Revaud J, Larlus D (2016) Deep Image Retrieval: Learning Global Representations for Image Search, Springer International Publishing, Cham, pp 241–257. DOI 10.1007/978-3-319-46466-4_15

Goroshin R, Bruna J, Tompson J, Eigen D, LeCun Y (2015a) Unsupervised learning of spatiotemporally coherent metrics. In: The IEEE International Conference on Computer Vision (ICCV)

Goroshin R, Mathieu MF, LeCun Y (2015b) Learning to linearize under uncertainty. In: Advances in Neural Information Processing Systems, pp 1234–1242

Griffin G, Holub A, Perona P (2007) Caltech-256 object category dataset. Tech. Rep. 7694, California Institute of Technology

Handa A, Patraucean V, Badrinarayanan V, Stent S, Cipolla R (2016) Understanding real world indoor scenes with synthetic data. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

He K, Zhang X, Ren S, Sun J (2015) Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: The IEEE International Conference on Computer Vision (ICCV)

He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Held D, Thrun S, Savarese S (2016) Robust single-view instance recognition. In: 2016 IEEE International Conference on Robotics and Automation (ICRA), pp 2152–2159, DOI 10.1109/ICRA.2016.7487365

Herranz L, Jiang S, Li X (2016) Scene Recognition With CNNs: Objects, Scales and Dataset Bias. In: Conference on Computer Vision and Pattern Recognition, pp 571–579, DOI 10.1109/CVPR.2016.68

Higy B, Ciliberto C, Rosasco L, Natale L (2016) Combining sensory modalities and exploratory procedures to improve haptic object recognition in robotics. In: IEEE-RAS International Conference on Humanoid Robots

Hinton GE, Srivastava N, Krizhevsky A, Sutskever I, Salakhutdinov R (2012) Improving neural networks by preventing co-adaptation of feature detectors. CoRR abs/1207.0580, 1207.0580

Hoffman J, Tzeng E, Donahue J, Jia Y, Saenko K, Darrell T (2013) One-shot adaptation of supervised deep convolutional models. CoRR abs/1312.6204, 1312.6204

Högman V, Björkman M, Maki A, Kragic D (2016) A Sensorimotor Learning Framework for Object Categorization. IEEE Transactions on Cognitive and Developmental Systems 8(1):15–25, DOI 10.1109/TAMD.2015.2463728

Ioffe S, Szegedy C (2015) Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: Proceedings of the 32nd International Conference on Machine Learning, PMLR, Lille, France, vol 37, pp 448–456

Jayaraman D, Grauman K (2015) Learning image representations tied to ego-motion. In: The IEEE International Conference on Computer Vision (ICCV)

Jia Y, Shelhamer E, Donahue J, Karayev S, Long J, Girshick R, Guadarrama S, Darrell T (2014) Caffe: Convolutional architecture for fast feature embedding. In: Proceedings of the ACM International Conference on Multimedia - MM '14, ACM Press, pp 675–678, DOI 10.1145/2647868. 2654889

Joachims T, Hofmann T, Yue Y, Yu CN (2009) Predicting structured objects with support vector machines. Commun ACM 52(11):97–104, DOI 10.1145/1592761.1592783

Kasper A, Xue Z, Dillmann R (2012) The kit object models database: An object model database for object recognition, localization and manipulation in service robotics. The International Journal of Robotics Research 31(8):927–934

Kemp CC, Edsinger A, Torres-Jara E (2007) Challenges for robot manipulation in human environments [grand challenges of robotics]. IEEE Robotics Automation Magazine 14(1):20–29, DOI 10.1109/MRA.2007.339604

Khosla A, Zhou T, Malisiewicz T, Efros AA, Torralba A (2012) Undoing the Damage of Dataset Bias, Springer Berlin Heidelberg, pp 158–171. DOI 10.1007/ 978-3-642-33718-5_12

Kingma D, Ba J (2015) Adam: A Method for Stochastic Optimization. In: 3rd International Conference for Learning Representations (ICLR)

Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems 25, Curran Associates, Inc., pp 1097–1105

Kulkarni TD, Mansinghka VK, Kohli P, Tenenbaum JB (2014) Inverse graphics with probabilistic cad models. arXiv preprint arXiv:14071339

Kuzborskij I, Orabona F, Caputo B (2013) From n to n+ 1: Multiclass transfer incremental learning. In: Computer Vision and Pattern Recognition (CVPR)

Lai K, Bo L, Ren X, Fox D (2011) A large-scale hierarchical multi-view rgb-d object dataset. In: 2011 IEEE International Conference on Robotics and Automation, pp 1817–1824, DOI 10.1109/ICRA.2011.5980382

LeCun Y, Boser B, Denker JS, Henderson D, Howard RE, Hubbard W, Jackel LD (1989) Backpropagation applied to handwritten zip code recognition. Neural computation 1(4):541–551

LeCun Y, Huang FJ, Bottou L (2004) Learning methods for generic object recognition with invariance to pose and lighting. In: Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on, vol 2, pp II–97–104 Vol.2, DOI 10.1109/CVPR.2004.1315150

Leitner J, Förster A, Schmidhuber J (2014) Improving robot vision models for object detection through interaction. In: 2014 International Joint Conference on Neural Networks (IJCNN), pp 3355–3362

Leitner J, Dansereau DG, Shirazi S, Corke P (2015) The need for more dynamic and active datasets. In: CVPR Workshop on The Future of Datasets in Computer Vision

Levine S, Finn C, Darrell T, Abbeel P (2016) End-to-end training of deep visuomotor policies. Journal of Machine Learning Research 17(39):1–40

Levine S, Pastor P, Krizhevsky A, Ibarz J, Quillen D (2018) Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. The International Journal of Robotics Research 37(4-5):421–436, DOI 10.1177/0278364917710318, URL https://doi.org/ 10.1177/0278364917710318, https://doi.org/10.1177/ 0278364917710318

Lowe DG (2004) Distinctive image features from scale-invariant keypoints. International Journal of Computer Vision 60(2):91–110, DOI 10.1023/B:VISI.0000029664. 99615.94

Luo S, Bimbo J, Dahiya R, Liu H (2017) Robotic tactile perception of object properties: A review. Mechatronics 48:54 – 67, DOI https://doi.org/10.1016/j.mechatronics.2017. 11.002, URL http://www.sciencedirect.com/science/ article/pii/S0957415817301575

Mansinghka V, Kulkarni TD, Perov YN, Tenenbaum J (2013) Approximate bayesian image interpretation using generative probabilistic graphics programs. In: Advances in Neural Information Processing Systems, pp 1520–1528

Metta G, Natale L, Nori F, Sandini G, Vernon D, Fadiga L, von Hofsten C, Rosander K, Lopes M, Santos-Victor J, Bernardino A, Montesano L (2010) The icub humanoid robot: an open-systems platform for research in cognitive development. Neural networks : the official journal of the International Neural Network Society 23(8-9):1125–34, DOI 10.1016/j.neunet.2010.08.010

Mikolajczyk K, Schmid C (2004) Scale & affine invariant interest point detectors. International journal of computer vision 60(1):63–86

Minh HQ, Sindhwani V (2011) Vector-valued manifold regularization. International Conference on Machine Learning

Model I, Shamir L (2015) Comparison of Data Set Bias in Object Recognition Benchmarks. IEEE Access 3:1953–1962, DOI 10.1109/ACCESS.2015.2491921

Moldovan B, Moreno P, van Otterlo M, Santos-Victor J, De Raedt L (2012) Learning relational affordance models for robots in multi-object manipulation tasks. In: Robotics and Automation (ICRA), 2012 IEEE International Conference on, IEEE, pp 4373–4378

Montesano L, Lopes M, Bernardino A, Santos-Victor J (2008) Learning object affordances: From sensory–motor coordination to imitation. IEEE Transactions on Robotics 24(1):15–26

Muja M, Rusu RB, Bradski G, Lowe DG (2011) Rein - a fast, robust, scalable recognition infrastructure. In: Robotics and Automation (ICRA), 2011 IEEE International Conference on, pp 2939–2946, DOI 10.1109/ICRA. 2011.5980153

Nene SA, Nayar SK, Murase H (1996) Columbia object image library (coil-100. Tech. rep., Columbia University

Nguyen A, Kanoulas D, Caldwell DG, Tsagarakis NG (2016) Detecting object affordances with convolutional neural networks. In: 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp 2765–2770, DOI 10.1109/IROS.2016.7759429

Oberlin J, Meier M, Kraska T, Tellex S (2015) Acquiring Object Experiences at Scale. In: AAAI-RSS Special Workshop on the 50th Anniversary of Shakey: The Role of AI to Harmonize Robots and Humans, blue Sky Award.

Oquab M, Bottou L, Laptev I, Sivic J (2014) Learning and transferring mid-level image representations using convolutional neural networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Pasquale G (2017) Deep learning for visual recognition in robotics: Are we ready for real world applications? PhD thesis, Istituto Italiano di Tecnologia, iCub Facility, and University of Genoa, DIBRIS

Pasquale G, Ciliberto C, Odone F, Rosasco L, Natale L (2015) Teaching icub to recognize objects using deep convolutional neural networks. In: Proceedings of The 4th Workshop on Machine Learning for Interactive Systems at ICML 2015, PMLR, Lille, France, vol 43, pp 21–25

Pasquale G, Ciliberto C, Rosasco L, Natale L (2016a) Object identification from few examples by improving the invariance of a deep convolutional neural network. In: 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp 4904–4911, DOI 10. 1109/IROS.2016.7759720

Pasquale G, Mar T, Ciliberto C, Rosasco LA, Natale L (2016b) Enabling depth-driven visual attention on the icub humanoid robot: Instructions for use and new perspectives. Frontiers in Robotics and AI 3(35), DOI 10.3389/frobt.2016.00035

Philbin J, Chum O, Isard M, Sivic J, Zisserman A (2008) Lost in quantization: Improving particular object retrieval in large scale image databases. In: Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, pp 1–8

Pillai S, Leonard J (2015) Monocular slam supported object recognition. In: Proceedings of Robotics: Science and Systems (RSS), Rome, Italy

Pinto L, Gupta A (2016) Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In: 2016 IEEE International Conference on Robotics and Automation (ICRA), pp 3406–3413, DOI 10.1109/ICRA.2016. 7487517

Pinto L, Gandhi D, Han Y, Park YL, Gupta A (2016) The curious robot: Learning visual representations via physical interactions. In: Leibe B, Matas J, Sebe N, Welling M (eds) Computer Vision – ECCV 2016, Springer International Publishing, pp 3–18

Pinto N, Cox DD, DiCarlo JJ (2008) Why is real-world visual object recognition hard? PLoS computational biology 4(1):e27, DOI 10.1371/journal.pcbi.0040027

Pinto N, Barhomi Y, Cox DD, DiCarlo JJ (2011) Comparing state-of-the-art visual features on invariant object recognition tasks. In: 2011 IEEE Workshop on Applications of Computer Vision (WACV), pp 463–470, DOI 10.1109/WACV.2011.5711540

Poggio T, Anselmi F (2016) Visual Cortex and Deep Networks: Learning Invariant Representations. The MIT Press, Cambridge, MA, USA

Redmon J, Angelova A (2015) Real-time grasp detection using convolutional neural networks. In: 2015 IEEE International Conference on Robotics and Automation (ICRA), pp 1316–1322, DOI 10.1109/ICRA.2015.7139361

Rennie C, Shome R, Bekris KE, Ferreira De Souza A (2016) A dataset for improved rgbd-based object detection and pose estimation for warehouse pick-and-place. IEEE Robotics and Automation Letters (RA-L) [Also accepted to appear at the 2016 IEEE International Conference on Robotics and Automation (ICRA)] 1:1179 – 1185

Rivera-Rubio J, Idrees S, Alexiou I, Hadjilucas L, Bharath AA (2014) Small hand-held object recognition test (SHORT). In: 2014 IEEE Winter Conference on Applications of Computer Vision (WACV), pp 524–531, DOI 10.1109/WACV.2014.6836057

Rodner E, Hoffman J, Donahue J, Darrell T, Saenko K (2013) Towards adapting imagenet to reality: Scalable domain adaptation with implicit low-rank transformations. CoRR abs/1308.4200, 1308.4200

Rudi A, Rosasco L (2017) Generalization properties of learning with random features. In: Advances in Neural Information Processing Systems 30, pp 3218–3228

Rudi A, Camoriano R, Rosasco L (2015) Less is more: Nystrom computational regularization. In: Advances in Neural Information Processing Systems, pp 1648–1656

Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg AC, Fei-Fei L (2015) Imagenet large scale visual recognition challenge. International Journal of Computer Vision 115(3):211–252, DOI 10.1007/s11263-015-0816-y

Schwarz M, Schulz H, Behnke S (2015) Rgb-d object recognition and pose estimation based on pre-trained convolutional neural network features. In: 2015 IEEE International Conference on Robotics and Automation (ICRA), pp 1329–1335, DOI 10.1109/ICRA.2015.7139363

Sharif Razavian A, Azizpour H, Sullivan J, Carlsson S (2014) Cnn features off-the-shelf: An astounding baseline for recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops

Simonyan K, Vedaldi A, Zisserman A (2014) Deep inside convolutional networks: Visualising image classification models and saliency maps. In: ICLR Workshop

Sinapov J, Schenck C, Staley K, Sukhoy V, Stoytchev A (2014) Grounding semantic categories in behavioral interactions: Experiments with 100 objects. Robotics and Autonomous Systems 62(5):632–645, DOI 10.1016/j.robot.2012.10.007

Singh A, Sha J, Narayan KS, Achim T, Abbeel P (2014) Bigbird: A large-scale 3d database of object instances. In: 2014 IEEE International Conference on Robotics and Automation (ICRA), pp 509–516, DOI 10.1109/ICRA.2014.6906903

Song S, Lichtenberg SP, Xiao J (2015) Sun rgb-d: A rgb-d scene understanding benchmark suite. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Srivastava N, Hinton GE, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way to prevent neural networks from overfitting. Journal of Machine Learning Research 15(1):1929–1958

Stamos D, Martelli S, Nabi M, McDonald A, Murino V, Pontil M (2015) Learning with dataset bias in latent subcategory models. In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, IEEE, vol 07-12-June, pp 3650–3658, DOI 10.1109/CVPR.2015.7298988

Sun Y, Fox D (2016) NEOL : Toward Never-Ending Object Learning for Robots. 2016 IEEE International Conference on Robotics and Automation (ICRA) pp 1621–1627, DOI 10.1109/ICRA.2016.7487302

Sünderhauf N, Shirazi S, Dayoub F, Upcroft B, Milford M (2015) On the performance of convnet features for place recognition. In: Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on, pp 4297–4304, DOI 10.1109/IROS.2015.7353986

Sünderhauf N, Dayoub F, McMahon S, Talbot B, Schulz R, Corke P, Wyeth G, Upcroft B, Milford M (2016) Place categorization and semantic mapping on a mobile robot. In: 2016 IEEE International Conference on Robotics and Automation (ICRA), pp 5729–5736, DOI 10.1109/ICRA.2016.7487796

Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A (2015) Going deeper with convolutions. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Thrun S, Mitchell TM (1995) Lifelong robot learning. Springer

Tommasi T, Orabona F, Caputo B (2010) Safety in numbers: Learning categories from few examples with multi model knowledge transfer. In: Computer Vision and Pattern Recognition (CVPR)

Tommasi T, Patricia N, Caputo B, Tuytelaars T (2015) A deeper look at dataset bias. In: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), Springer International Publishing, vol 9358, pp 504–516, DOI 10.1007/978-3-319-24947-6_42

Torralba A, Efros AA (2011) Unbiased look at dataset bias. In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp 1521–1528, DOI 10.1109/CVPR.2011.5995347

Wang X, Gupta A (2015) Unsupervised learning of visual representations using videos. In: Proceedings of the IEEE International Conference on Computer Vision, pp 2794–2802

Wu Z, Song S, Khosla A, Yu F, Zhang L, Tang X, Xiao J (2015) 3d shapenets: A deep representation for volumetric shapes. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Yosinski J, Clune J, Fuchs T, Lipson H (2014) Understanding neural networks through deep visualization. In: ICML Workshop on Deep Learning

Zeiler MD, Fergus R (2014) Visualizing and Understanding Convolutional Networks, Springer International Publishing, Cham, pp 818–833. DOI 10.1007/978-3-319-10590-1_53

Zhang Y, Song S, Yumer E, Savva M, Lee JY, Jin H, Funkhouser T (2017) Physically-based rendering for indoor scene understanding using convolutional neural networks. The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)

Zhao J, Itti L (2016) Improved Deep Learning of Object Category using Pose Information. ArXiv e-prints 1607.05836

Zhao J, Chang Ck, Itti L (2016) Learning to Recognize Objects by Retaining other Factors of Variation. ArXiv e-prints 1607.05851

# Are We Done with Object Recognition? The iCub Robot's Perspective. [Supplementary Material]

Giulia Pasquale, Carlo Ciliberto, Francesca Odone, Lorenzo Rosasco and Lorenzo Natale

## A Objects in the iCWT Dataset

Fig. 14 shows one example image for each object in the iCub-World Transformations (iCWT) dataset. We report the associated ImageNet synset in Red for categories among the 1000 classes of the ILSVRC (Russakovsky et al 2015), in Black for others (Deng et al 2009).

## B Testing Off-the-shelf CNNs

In Fig. 5 we observed that testing off-the-shelf CNNs to iCWT leads to poor accuracy. In this Section we provide details on this experiment and some qualitative explanations for the observed performance drop.

### B.1 Image Preprocessing

We first evaluated the impact of processing with CNNs coarser or finer regions around the object of interest. Specifically, we compared extracting from the images either (i) a square region of fixed radius centered on the object or (ii) the bounding box provided by depth segmentation (obtained as explained in Sec. 2.1). In both cases, we included more or less background (respectively fixing the radius to 256 or 384 and leaving a margin of 30 or 60 pixels around the bounding box).

Then we reproduced the operations described at the reference page of CAFFE models (see Tab. 6). These basically consist in subtracting the mean image (or the mean pixel) of the training set and running the CNN on a grid of fixed-size crops ($227 \times 227$ for *CaffeNet* and $224 \times 224$ for the other models) at multiple scales. The final prediction is computed by aggregating the predictions of the crops. Since, however, in our case, the considered region is already cropped around the object, we also tried considering only one central crop. Fig. 15 reports the experiment of Fig. 5, when testing off-the-shelf CNNs on iCWT with the described pre-processing options. It can be noticed that (i) a finer localization of the object (Blue and Green) provides better performance and (ii) considering more than the central crop provides a little or no advantage at all.

Based on this finding, for all the experiments reported in the paper we opted for extracting a square $256 \times 256$ region from the image (whose resolution is $640 \times 480$) and considering only the central crop (Light Green in Fig. 15). We opted for a fixed-size region rather than the bounding box from depth segmentation, since this latter has varying shape and should be resized anysotropically to be processed by the CNNs, thus impacting on the viewpoint transformations we aim to study. Note that when applying the networks to ImageNet, we used the multi-crop strategy suggested for each architecture.

**Table 6** Image preprocessing executed before feeding the networks.

| Model | Mean Subtraction | Scaling | Crop Extraction |
|---|---|---|---|
| **CaffeNet ResNet-50** | image | mean image size ($256 \times 256$) | $2 \times 2$ grid + center mirrored |
| **GoogLeNet** | pixel | $256 \times 256$ | $2 \times 2$ grid + center mirrored |
| **VGG-16** | pixel | shorter side to 256, 384, 512 | $5 \times 5$ grid at each scale mirrored |

### B.2 1000-class Categorization Results

In Fig. 5 we compared the accuracy of off-the-shelf CNNs tested on iCWT and ImageNet, selecting the 11 scores of the considered classes from the vector of 1000 scores produced by the CNNs trained on the ILSVRC. In Fig. 16 we report the accuracy of the same experiment, when considering 1000 scores. As it can be noticed, performance drops significantly and proportionally for all models and both test sets (since now chance level for the model is 1/1000). Note that, in this way, the models provide worse-than-chance predictions on iCWT, if chance level for iCWT is considered 1/11).

### B.3 Viewpoint Biases

In this Section we follow-up Sec. 4.1 and show potential biases in ImageNet, that may prevent off-the-shelf CNNs to generalize to iCWT. Keeping our observations qualitative, we analyze frame-level predictions in order to understand which views are actually "harder" to recognize and compare them with prototypical examples in ImageNet.

We report results for *GoogLeNet* (*CaffeNet* behaved similarly) and a subset of the test set considered in Fig. 5: specifically, 2 representative categories and *Scale* and 2*D Rotation* image sequences. As in Fig. 5, the prediction was computed as the maximum among the 11 scores selected from the CNN output. In Fig. 17 we report frame-level predictions on *Scale* (Top) and 2*D Rotation* (Bottom) sequences, separately per category. In each plot, rows represent the sequences of the 10 instances of the category: the frame index is reported in the horizontal axis and each frame is a vertical bar: *White* if the prediction is correct, *Red* if it is wrong, *Black* if the sequence is finished.

It can be noted that, since during the *Scale* acquisition the operator was moving the object back and forth in front of the robot (starting close and moving backward), the CNN fails when the object is far, hence at a smaller scale (red bars mostly concentrated in the right half of rows). This confirms recent studies (Herranz et al 2016) pointing out that the scale bias of object-centric datasets as ImageNet prevents CNNs

soda bottle
n03983396

cellphone
n02992529

mouse
n03793489

coffee mug
n03063599

pencil case
n03908618

perfume
n03916031

remote
n04074963

ring binder
n02840245

soap dispenser
n04254120

sunglasses
n04356056

wallet
n04548362

oven glove
n02885462

squeezer
n04293119

sprayer
n02754103

body lotion
n02862916

book
n02870526

flower
n11669921

glass
n03438257

hairbrush
n03475581

hair clip
n03476684

**Fig. 14** Example images for the 200 objects in iCubWorld Transformations.

**Fig. 15** Average classification accuracy of off-the-shelf networks tested on iCWT segmenting the object according to different strategies.



**Fig. 16** Average classification accuracy of off-the-shelf networks (trained on ILSVRC) tested on iCWT (Dark Blue) or on ImageNet itself (Gray). The test sets for the two datasets are restricted to the 11 categories in common (see Sec. 2).

trained on them to generalize to real world, scene-centric, settings. Another example regards 2D Rotation sequences, where the operator was rotating the object, keeping the same face visible, at a constant speed. In these sequences it can be noticed that the CNN fails at periodic time intervals, corresponding to less common views in ImageNet. In fact, soap dispensers and mugs in ImageNet are mostly placed on tables and, consequently, the CNNs fails when these are rotated from the vertical.

## C Model Selection

In this Section we provide details on the hyper-parameter selection for the methods adopted in the paper (Sec. 3.2).



**Fig. 17** Frame-level predictions of *GoogLeNet* on image sequences containing the *Scale* (Top) and 2D *Rotation* (Bottom) transformations, reported for 2 representative categories. The sequences of the 10 instances belonging to each category are represented as matrix rows (frame index in the horizontal axis). In each row, frame predictions over the sequence are represented as vertical bars: *White* if correct, *Red* if wrong (*Black* if the sequence is ended).

### C.1 Feature Extraction and RLSC

We illustrate some design choices for the pipeline presented in Sec. 3.2.1.

To extract CNN representations, we considered fully-connected layers providing a vector global image representation. These were specified for each architecture in Tab. 2: in particular, for *CaffeNet* and *VGG-16* we tried either *fc6* or *fc7* layers as explained in the following.

We then resorted to (Rudi et al 2015) to implement RLSC with Gaussian Kernel and Nyström subsampling. The major hyper-parameters of this algorithm the number $m$ of training examples sampled to approximate the kernel matrix, the Gaussian's $\sigma$ and the regularization parameter $\lambda$. While the latter two were assigned based on standard cross-validation, that we performed for each experiment, in the following we report how we empirically determined a reasonable range for $m$.

We considered two categorization tasks on iCWT, representative of the smallest and larger tasks that we expected to run in our analysis: respectively $\sim 170$ and $\sim 6300$ examples per class, for 15 classes (similarly to Sec. 4.2.1). Fig. 19 and 18 report the average accuracy respectively for the small and large experiment. We used image representations from the four considered architectures, using either *fc6* layer (Gray) or *fc7* (Pink) for *CaffeNet* and *VGG-16*. We increased logarith-
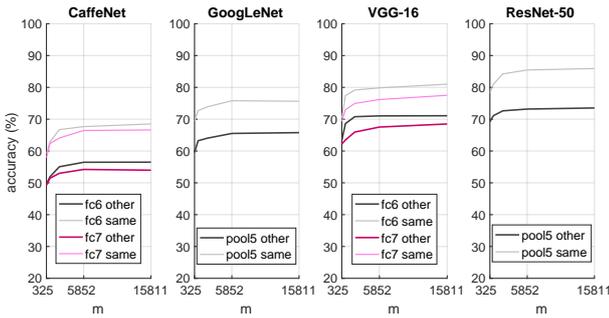
**Fig. 18** Accuracy reported by training RLSC on image representations extracted by the four considered CNNs. A "large" categorization experiment is performed ($\sim N = 95\text{K}$ training examples) and $m$ is varied between $\sqrt{N}$ and 15K (horizontal axis). Performance on the same day of training (Light colors) and on a different one (Dark colors) is reported.

**Fig. 19** Similar to Fig. 18 but on a "small" categorization experiment ($\sim N = 2500$ training examples), varying $m$ between $\sqrt{N}$ and $N$.

mically the value of $m$ (horizontal axis) starting from $\sqrt{N}$, $N$ being the size of the training set. For the small experiment we stopped at $N$, whereas for the large experiment we stopped at values where we observed very little performance improvement. Performance on the same day of training is reported in light Gray/Pink, on a different day in dark Gray/Pink.

From this experiment we observed that *fc6* features consistently outperformed *fc7*, hence we decided to use this layer when extracting representations from off-the-shelf *CaffeNet* or *VGG-16*. Then, we observed that relatively small values of $m$ could provide good accuracies with far smaller training times and therefore we selected $m = \min(15K, N)$ in all experiments.

It is worth noting that we performed a similar experiment (not reported) with CNNs previously fine-tuned on subsets of iCWT and, in that case, *fc7* features were better. Hence, we used this layer in the experiments of Sec. 5.2.1. This is in line with recent work showing that, while lower layers provide more "general" features, more "specialized" features can be extracted from higher layers in models trained on closer domains.

## C.2 Fine-tuning

In this Section we provide details on the fine-tuning procedures and the definition of the two *adaptive* and *conservative* strategies reported in Tab. 3.

### C.2.1 General Protocol

Image preprocessing was similar to the one described in B.1 (except that, as per CAFFE standard, during training a random $227 \times 227$ crop, randomly mirrored, was extracted). Importantly, the training set was shuffled, since we observed that similarity of images within a mini-batch (as consecutive frames would be) negatively affected convergence. We evaluated performance on a validation set every epoch and finally chose the model at the best epoch. We set the mini-batch size to values specified in CAFFE (reported in Tab. 3). The number of epochs was fixed empirically, observing that performance saturated after $\sim 6$ epochs but for the *conservative*
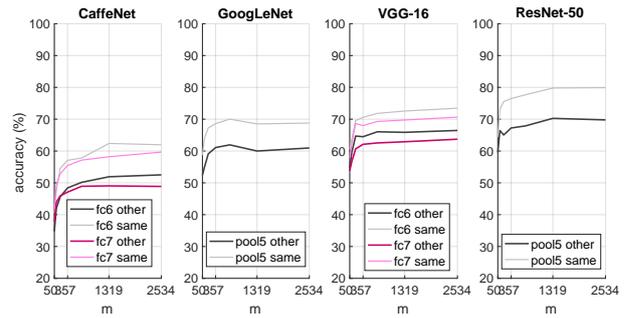
fine-tuning of *CaffeNet*, which involves learning the parameters of the 3 fully-connected layers and takes more epochs to converge (we stopped at 36).

### C.2.2 Hyper-parameters Choice

While model selection is per se an open problem when dealing with deep networks, in our setting this issue is even more complicated by the fact that we do not target a fixed reference task, but plan to span over a wide range of tasks comprising small and large training sets. To this end in this Section we report on the empirical analysis we set up in order to understand the effect and relative importance of the (many) hyper-parameters involved in fine-tuning deep architectures as *CaffeNet* or *GoogLeNet*.

We considered the same two "small" and "large" categorization tasks adopted to perform parameter selection for RLSC as described in C.1 and fine-tuned *CaffeNet* and *GoogLeNet* by varying the values of multiple hyper-parameters. In the following, we report this analysis separately for the two architectures.

**CaffeNet.** We considered the parameters in Tab. 3 and varied them in the following way:

— `Base LR`: the starting learning rate of the layers that are initialized with the parameters of the off-the-shelf model. We tried $10^{-3}$, $5 * 10^{-4}$, $10^{-4}$, $5 * 10^{-5}$, $10^{-5}$, $10^{-6}$, 0.
— `Learned FC Layers`: which fully-connected (FC) layers are learned from scratch with their specific starting LR. We tried to learn (i) only $fc8$ with starting LR set to $10^{-2}$, or (ii) including also $fc7$ and (iii) finally also $fc6$. As an empirical rule, every time we included one more layer to learn from scratch, we decreased the starting LR of these layers of a factor of 10 (hence $10^{-3}$ in (ii) and $10^{-4}$ in (iii)).
— `Dropout`: percentage of dropout in FC layers. We tried 50% (default CAFFE value) or 65%.
— `Solver`: the algorithm used for the stochastic gradient descent. We used the *SGD* solver (Bottou 2012) in CAFFE.
— `LR Decay Policy`: the decay rate of the learning rates. We tried either polynomial decay with exponent 0.5 or $-3$, or step decay decreasing the LR of a factor of 10 every 2 epochs.

We tried all possible combinations of values. For the parameters not mentioned, we kept their value as in CAFFE reference models.
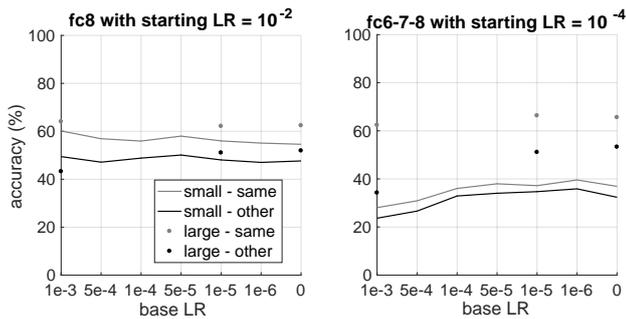
**Fig. 20** Accuracy provided by *CaffeNet* fine-tuned with different strategies: either learning from scratch only $fc8$ (Left) or $fc6$, $fc7$ and $fc8$ (Right). Performance is reported for both the same day of training (Light Gray) and a different one (Dark Gray). We tried multiple values of `base LR` on the "small" training set (Continuous Line), and one small, medium and large values of `base LR` (Dots) on the "large" one.

We observed that the `Dropout` percentage had a small influence and left it to the default value. We also observed that the polynomial `LR Decay Policy` with 0.5 slope consistently provided 5-10% better accuracy. The most critical parameters proved being the `Base LR` and the `Learned FC Layers`. As an example, in Fig. 20 we report the accuracy obtained respectively when learning only $fc8$ (Left) or $fc6$, $fc7$ and $fc8$ (Right). In both cases, we repeated fine-tuning with different `Base LR` for all other layers, varied from $10^{-3}$ to 0 (horizontal axis). Performance is reported, as in Fig. 19 and 18, for both the same day of training (Light Gray) and a different one (Dark Gray). While we tried all values of `Base LR` on the "small" experiment (Continuous Line), for the "large" experiment we limited to one small, medium and large value (Dots).

It can be observed that fine-tuning by learning from scratch only the last layer ($fc8$) is more robust to the small training set, for any `Base LR` (Continuos Line in the range 40-60%), achieving best performance with higher `Base LR` ($10^{-3}$). On the other hand, learning from scratch $fc6$-7-8 with small `Base LR` ($10^{-5}$ or 0) achieves best performance on the large training set. This is explained by noting that the three layers $fc6$, $fc7$ and $fc8$ involve a lot of parameters.

Based on these findings, we identified two representative strategies providing best performance respectively in small and large-scale settings: learning $fc8$ with large `Base LR` ($10^{-3}$), that we call *adaptive* strategy, since it quickly adapts all layers to the new training set, and learning $fc6$-7-8 with `Base LR` set to 0, that we call the *conservative* strategy, since it slowly adapts only fully-connected layers.

**GoogLeNet.** As explained in (Szegedy et al 2015), this architecture is composed of a main branch, terminating with one FC layer ($loss3/classifier$), and two identical "auxiliary" branches terminating with two FC layers ($loss1(2)/fc$, $loss1(2)/classifier$). By considering this structure, we explored the following fine-tuning strategies:

- `Base LR`: varied as for *CaffeNet*.
- `Learned FC Layers`: we always learned $loss3/classifier$ from scratch with starting LR set to $10^{-2}$; regarding the auxiliary branches, we tried (i) to cut them out, (ii) to learn also $loss1(2)/classifier$ from scratch with starting LR equal to $10^{-2}$, or, finally, (iii) to learn from scratch both $loss1(2)/fc$ and $loss1(2)/classifier$, with starting LR set to $10^{-3}$.
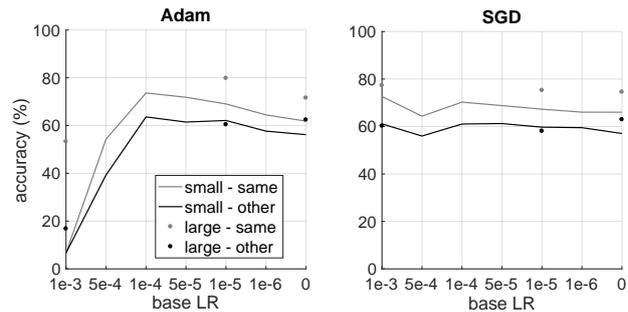
**Fig. 21** Classification accuracy provided by fine-tuning *GoogLeNet* according to different strategies: either using the *Adam* solver (Left) or *SGD* (Right). The rest of the figure is similar to Fig. 20

- `Dropout`: we tried either the default CAFFE values (40% for $loss3/classifier$ and 70% for $loss1(2)/fc$) or 60% for $loss3/classifier$ and 80% for $loss1(2)/fc$.
- `Solver`: we tried either *SGD* or *Adam* (Kingma and Ba 2015) solvers in CAFFE.
- `LR Decay Policy`: when using *SGD*, we used polynomial decay with exponent 0.5 or $-3$; when using *Adam*, we kept the learning rate constant.

As for *CaffeNet*, we tried all combinations and left other parameters to default CAFFE values.

We first observed that for this architecture the impact of the `Learned FC Layers` from scratch was small, with the three strategies behaving similarly. We chose the last one (iii), that was slightly more stable. We also observed a little benefit from using higher `Dropout` percentages.

One critical aspect was instead the choice of the `Solver`. To this end, in Fig. 21 we report the accuracy achieved respectively when using *Adam* (Left) or *SGD* (Right). In the latter case, we applied the polynomial `LR Decay Policy` with 0.5 exponent, since it was consistently better. Performance is reported, as in Fig. 20, for both the same day of training (Light Gray) and a different one (Dark Gray). We fine-tuned again with different `Base LR` from $10^{-3}$ to 0 (horizontal axis), trying all values on the "small" experiment (Continuous Line) and one small, medium and large value for the "large" experiment (Dots).

It can be observed that the *SGD* solver is more robust to variations of the `Base LR`, but we opted for *Adam*, which provides better accuracies for mid-range values of `Base LR`, both for the small and the large setting. Similarly to *CaffeNet*, we identified an *adaptive* fine-tuning strategy with with `Base LR` set to 0 and a more *conservative* strategy with `Base LR` set to $10^{-5}$.

## D Additional Experiments on Categorization

In this Section we report additional results to show that the trends observed in the conditions considered in the paper hold also in other settings.

### D.1 What do we gain by adding more frames?

Fig. 24 reports the experiment of Fig. 6, when performed including only 3 object instances per category in the training
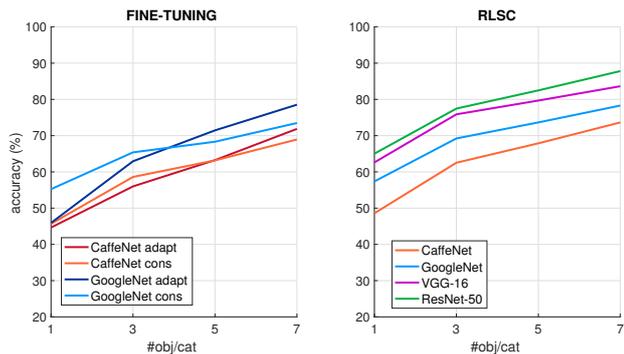
**Fig. 22 Recognition accuracy vs # instances** (number of object instances available during training). Same experiment as Fig. 7 executed for a 10-class categorization problem.
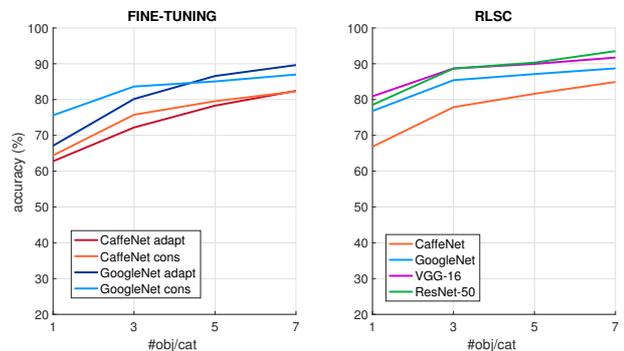


**Fig. 23 Recognition accuracy vs # instances** (number of object instances available during training). Same experiment as Fig. 7 executed for a 5-class categorization problem.
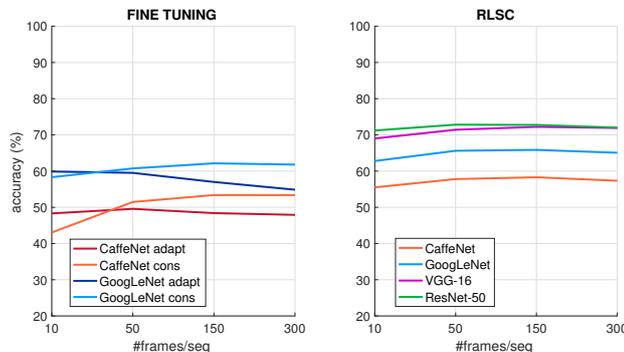


**Fig. 24 Recognition accuracy vs # frames** (number of object views available during training). Same experiment as Fig. 6 but training only on 3 object instances per category.

set. For completeness, here we provide the list of the considered 15 categories: *cellphone, mouse, coffee mug, pencil case, perfume, remote, ring binder, soap dispenser, sunglasses, flower, wallet, glass, hairbrush, hair clip, book*.

This result confirms that adding more object views cannot be used as a viable strategy to improve categorization accuracy, which remains constant and definitely lower than the performance achieved with 7 example instances per category (above $\sim 70\%$, see Fig. 6).

## D.2 What do we gain by adding more instances?

Fig. 22 and 23 report the experiment of Fig. 7, when discriminating between respectively 10 and 5 object categories (rather than 15). As it can be noticed, accuracy increases remarkably as more instances per category are made available, confirming that semantic variability is critical even in settings that involve few categories.

## D.3 Improving Viewpoint Invariance of CNNs

As mentioned at the end of Sec. 5.2.1, here we investigate whether image representations from CNNs fine-tuned on subsets of iCWT can be better than off-the-shelf features also for categorization tasks.

We repeat the 15-class task of Fig. 9, performed by sampling $\sim 20$ frames per sequence, by training RLSC on features from the CNNs fine-tuned as in Sec. 5.2.1. Note that for (**ImNet**) dataset we fine-tuned over the 15 ImageNet synsets corresponding to the 15 categories that involved in the categorization task. From the results reported in Fig. 25 (using the same notation as in Fig. 12), it can be clearly observed that none of these representations is better than off-the-shelf features for the categorization task.

## E Generalization Across Days

In this Section we test robustness to changes of illumination, background, etc., which are neither semantic nor geometric. The purpose of this material is to show how these aspects, albeit not in the scope of the paper, are very important and will be subject of future investigation.

We consider the common situation where the robot is asked to recognize an object that was showed him on a past day, possibly in a slightly different setting. We ask whether even small contextual variations (like a different time of day, or background configuration) can degrade the accuracy of the considered deep learning methods.

To this end, while, in all experiments of the paper, training and testing is always performed on the same day of acquisition, in the following we report the performance drop experienced when testing the same models on a different day (we recall that in the current release of iCWT each sequence acquisition is repeated in two different days).

**Categorization.** Fig. 26 report the performance drop observed respectively for the experiment of Fig. 6 (addition of example frames) and Fig. 7 (addition of object instances). The drop is computed as the difference between the accuracy reported in Fig. 6 (Fig. 7) and the accuracy achieved when testing the models on a different day.

In both experiments, performance degrade (from 5% up to almost 30% for the *adaptive* fine-tuning). A larger drop for fine-tuning (*adaptive* in particular) suggests that this approach is more prone to overfitting the training day. On the contrary, the $\sim 5\%$ drop experienced when using less aggressive strategies as RLSC suggests that features from off-the-shelf networks as *GoogleNet* and *ResNet*-50 can be quite robust.

Note that, when adding more example views, accuracy does not improve on the training day and even degrade on another day (Fig. 26 and Fig. 6). Differently, when adding more example instances, accuracy increases in both cases, but
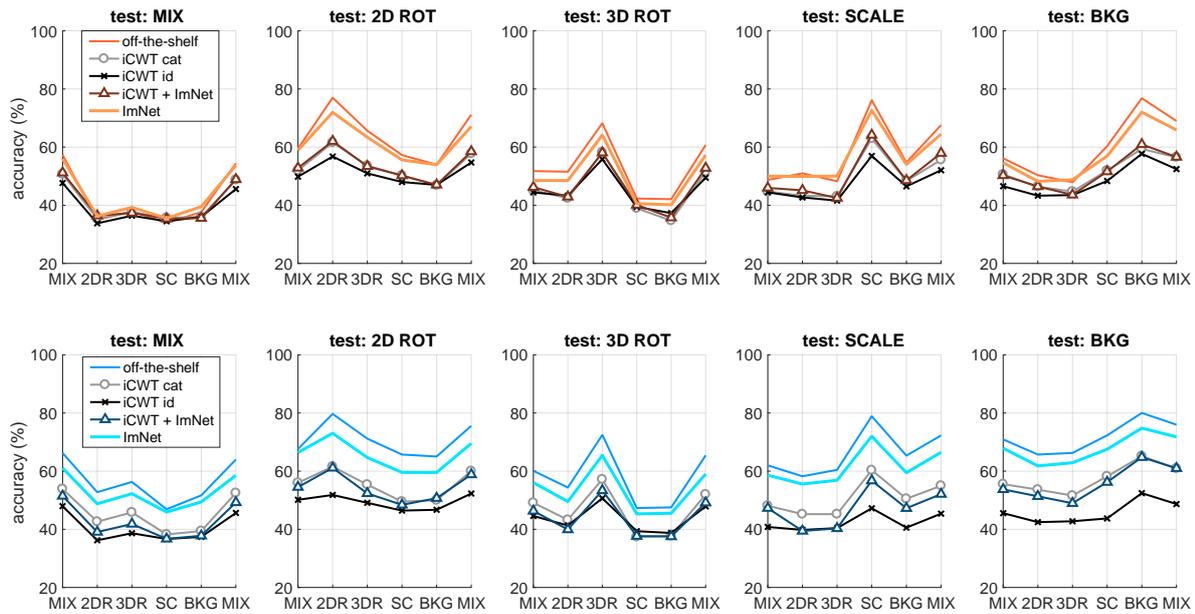
**Fig. 25** Same experimental setting as in Fig. 9: here we compare training RLSC on different image representations, provided by *CaffeNet* (Top, Orange) or *GoogLeNet* (Bottom, Blue) fine-tuned according to different strategies (see Sec. D.3).
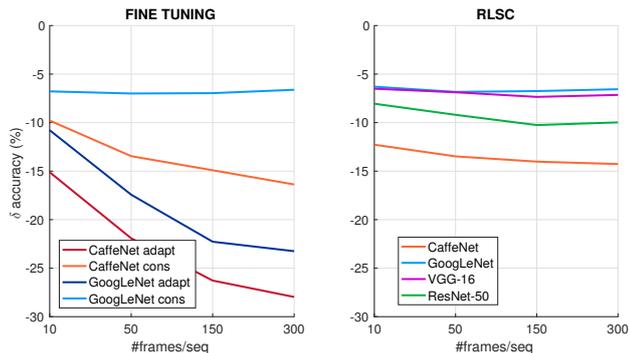


**Fig. 26 Categorization accuracy vs # frames - Generalization across days**. Drop in performance (difference between test accuracy) observed when testing the models trained for the categorization task reported in Sec. 4.2.1 on the same day of training and on a different one.
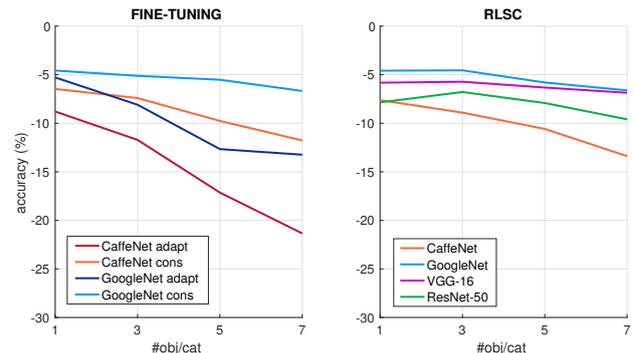
**Fig. 27 Categorization accuracy vs # instances - Generalization across days**. Drop in performance (difference between test accuracy) observed when testing the models trained for the categorization task reported in Sec. 4.2.2 on the same day of training and on a different one.

more on the training day than on another one (Fig. 27 and Fig. 7).

**Identification.** Differently, in this setting *adaptive* fine-tuning does not exhibit a similar dramatic drop: in Fig. 28 the drop with respect to Fig. 10 is small for all methods and does not increase by adding example frames. This positively indicates that adding example views in identification does not overfit the training day, and equally improves performance also on a different day.
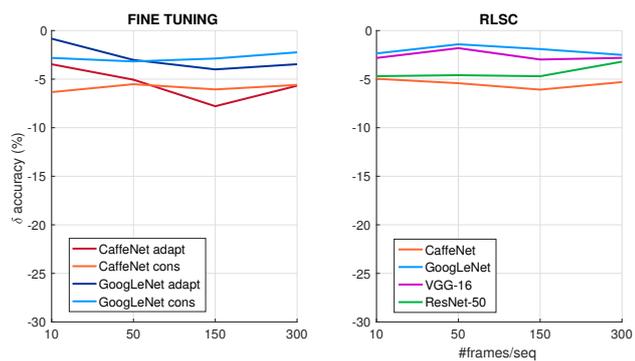
**Fig. 28 Identification accuracy vs # frames - Generalization across days**. Drop in performance (difference between test accuracy) observed when testing the models trained for the identification task reported in Sec. 5.1 on the same day of training and on a different one.