

On the Illumination Influence for Object Learning on Robot Companions

Ingo Keller^{1,*}, Prof. Dr. Katrin S. Lohan^{1,2}

¹*Department of Mathematical and Computer Science, Heriot-Watt University, Edinburgh, UK,*

²*EMS Institute for Development of Mechatronic Systems, NTB University of Applied Sciences in Technology, Buchs, CH*

Correspondence*:

Ingo Keller

I.Keller@hw.ac.uk

2 ABSTRACT

3 Most collaborative tasks require interaction with everyday objects (e.g., utensils while cooking).
4 Thus, robots must perceive everyday objects in an effective and efficient way. This highlights the
5 necessity of understanding environmental factors and their impact on visual perception, such as
6 illumination changes throughout the day on robotic systems in the real world. In object recognition,
7 two of these factors are changes due to illumination of the scene and differences in the sensors
8 capturing it.

9 In this paper, we will present data augmentations for object recognition that **enhance** a deep
10 learning architecture. We will show how simple linear and non-linear illumination models and
11 feature concatenation can be used to **improve** deep learning-based approaches. The aim of
12 this work is to allow for more realistic Human-Robot Interaction scenarios with a small amount
13 of training data **in combination with incremental interactive object learning**. This will benefit the
14 interaction with the robot to maximize object learning for long-term and location-independent
15 learning in unshaped environments.

16 With our model-based analysis, we showed that changes in illumination affect recognition
17 approaches that use Deep Convolutional Neural Network to encode features for object recognition.
18 Using data augmentation, we were able to show that such a system can be modified towards a
19 more robust recognition **without retraining the network**. Additionally, we have shown that using
20 simple brightness change models can help to improve the recognition across all training set sizes.

21 **Keywords:** Object Learning, Visual Perception, Data Augmentation, Human-Robot Interaction, Long-Term Engagement

1 INTRODUCTION

22 **Using robotic companions in unconstrained, domestic environments poses new challenges to the task of**
23 **object recognition and learning. In this work, we focus on such use-cases in which one cannot draw from**
24 **a large set of images for training since the presentation of particular objects might occur infrequently**
25 **and only over short periods. Moreover, consumer-oriented robotic hardware usually does not allow for**
26 **computationally expensive training (e.g., state-of-the-art deep learning networks). The learning needs to be**
27 **fast enough, in the range of seconds rather than hours, to be of value for the user. Additionally, the objects**

28 that will have to be learned are not necessarily known upfront and might only get presented over time
29 rather than all at once, which limits the possibility of pretraining networks.

30 New developments in deep learning, computer vision and robotics research together with the availability
31 of highly integrated, powerful mobile computer systems has made it possible to create robotic systems that
32 can operate in private households. Interactive household robots have started to appear in consumer stores
33 such as the Jibo¹, the Buddy Robot² or the robotic kitchen³. Also, research is increasingly focusing on
34 providing assistance and companion robots, e.g., Meghdari et al. (2018) for socially supporting children in
35 hospitals or the elderly. With an aim to focus on adapt to emotional features, e.g. Churamani et al. (2018).
36 All of these examples will lead to the availability of personal robotic assistants in the near future, the
37 research of long-term engagement with such systems is still in its infancy. A major obstacle for interaction
38 is the real-world environment since it is less controllable than a laboratory and therefore presents new
39 challenges to state-of-the-art approaches. The ability to cope with changes in the environment will be an
40 important factor in the acceptance of interactive robotic systems.

41 In Human-Robot Interaction (HRI), triadic interactions Imai et al. (2003) are one of the most commonly
42 studied problems to create natural interaction between robots and humans. To jointly manipulate objects,
43 Moldovan et al. (2012) first requires the robotic systems to recognize them reliably. Here, we focus
44 on improving object recognition in noisy environments to learn to cope with real-world constraints. In
45 particular, we research the influence of illumination changes and their impact on object recognition in the
46 context of real-time capable recognition systems without prior knowledge.

47 One state-of-the-art visual recognition systems for incremental interactive object learning is provided
48 by the iCub community. The approach utilizes a combination of a Deep Convolutional Neural Network
49 (DCNN) for feature generation with a Multiclass Support Vector Machine (SVM) for classification of
50 objects that were shown to the iCub Pasquale et al. (2015a). An exhaustive evaluation of the performance
51 of the combined networks can be found in Sharif Razavian et al. (2014). The system provides a method for
52 long-term object learning due to its incremental training on images which are acquired by interaction with
53 the robot. Classifiers can be trained near real-time and are usable in real-world scenarios in which novel
54 objects can appear at any time.

55 Training each robot individually on the objects of its particular environment is inefficient and doesn't scale.
56 Therefore, pretrained robots that can expand their knowledge on-the-fly are required for the real world.
57 Hence, we looked into the re-usability of existing datasets for training classification models following the
58 off-the-shelf approach for feature generation of Sharif Razavian et al. (2014). Reusing feature generation
59 models across different robotic platforms would dramatically reduce computational requirements and is
60 preferred over training individual robots. Therefore, we focus on global illumination changes that might
61 occur due to light changes throughout the day, different viewpoints or sensors used. As indicated by our
62 previous work Keller and Lohan (2016), light changes have a negative impact on state-of-the-art object
63 recognition.

64 The contribution of this work is a systematic analysis of the impact of light changes on using two different
65 light models over a wide range of parameters. The benefits of data augmentation are analyzed for the
66 mentioned feature generation and classification methods. Both methods are treated as black-box systems

¹ <https://www.jibo.com>

² <http://www.bluefrogrobotics.com>

³ <http://www.moley.com>

67 to provide a baseline for further research of methods for illumination robustness **in systems with low**
68 **computational resources for extensive retraining such as the mentioned robot companion systems.**

69 In particular, we research recognition capabilities under the assumption of a low number of input images
70 corresponding to short interaction durations. While the datasets used were created throughout multiple
71 training sessions in the lab, in real-world scenarios, this amount of time to acquire the data might not
72 be available. Therefore, we want to increase the recognition performance in circumstances that provide
73 only a small set of training examples. As will be discussed in section 3, we used the ICUBWORLD28
74 dataset to identify the impact of illumination changes on the object recognition pipeline. With the larger
75 ICUBWORLD TRANSFORMATION dataset, we transferred and tested our method to a broader diversity of
76 object manipulations, background changes, and brightness variance.

2 BACKGROUND

77 2.1 Visual Features

78 Two approaches for feature generation on 2D images can be distinguished in state-of-the-art object
79 recognition. Methods, such as SIFT, SURF, and ORB, use local, keypoint-based feature sets that are
80 generated from template images. These features are robust against a variety of transformations, including
81 scaling and rotation. Recognition of 3D objects can then be achieved by creating object feature databases
82 based on different viewpoints Yu et al. (2014) and using matching techniques such as RANSAC Fischler
83 and Bolles (1981) to determine if a given object's feature set can be considered a model for a set of
84 features found in a test image. While these approaches result in robust recognition and are therefore widely
85 used, they suffer from an increase in computational cost in the cases of a high number of objects or high
86 resolution of images.

87 Another type of object recognition method utilizes Deep Learning to generate features using pretrained
88 networks Sharif Razavian et al. (2014), Fischer et al. (2016). Networks such as the AlexNet Krizhevsky
89 et al. (2012) are trained on large image datasets and, once trained, can be treated as black-box filters that
90 generate features from images. The benefit of this technique is that the resulting feature vectors are of
91 fixed-length, which allows for the use of standard classification methods such as SVMs.

92 We investigate the latter approach **since** it is used in a real object recognition pipeline for incremental
93 learning on an existing humanoid robot platform (iCub) and due to the availability of datasets with many
94 different objects which were captured during interactions with the robot. The next part of this section will
95 give an overview of **the background of illumination variations.**

96 2.2 Illumination Variations

97 **Dealing with illumination variations is a long-standing problem for visual recognition systems, especially**
98 **for the perception of color. Ever since Land (1964) introduced the Mondrian experiment and proposed his**
99 **Retinex model, it has become clear that the human visual system perceives the color of an object not only**
100 **based on its photometric properties. The lighting condition of the surrounding environment is taken into**
101 **account by our brains to tune our sense to perceive a certain color, even if it is not physically present in the**
102 **scene. This problem is known as color constancy and has been addressed by computer vision research in**
103 **many ways. A comprehensive overview can be found in Foster (2011). Furthermore, the perceived color**
104 **of an object can vary between people, and it is thought to depend on people's age, gender, and general**
105 **light exposure habits (early birds vs. night owls) Lafer-Sousa et al. (2015). These findings indicate that the**

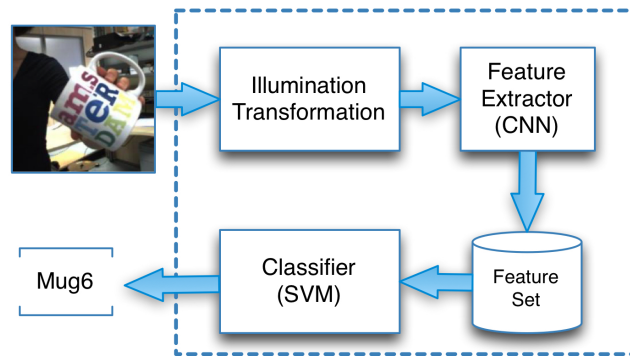


Figure 1. Our adaptation of the iCub’s recognition workflow includes the illumination transformations.

106 human visual system is capable of adapting over time and learning certain illumination scenarios, which it
 107 uses as priors to adjust the perception of color.

108 To capture this ability of the human brain, a wide range of methods have been proposed to tackle
 109 robustness against illumination changes. However, methods that aim to take light variations into account
 110 either suffer from a high computational cost, which renders them infeasible to use in real-time systems with
 111 limited computational resources, or make assumptions that are not met by generic, natural images. The
 112 first case is presented by state-of-the-art methods for object recognition, such as the currently fastest object
 113 detector, YOLOv3 Redmon and Farhadi (2018). Deep Learning methods can achieve a high performance
 114 but require the support of powerful GPUs which are usually not available on the type of robot systems we
 115 are targeting (see Reyes et al. (2018)). The second class of methods can be found in more specific areas.
 116 A prominent example is face recognition, in which most illumination models make assumptions that do
 117 not hold for general objects such as Lambertian surface reflectance, underlying facial models, or more
 118 generally the importance of facial landmarks over general object features (e.g., Le and Kakadiaris (2019)).

119 While state-of-the-art methods have advanced the standard for illumination robustness, they are also
 120 usually computationally expensive. Using pretrained deep learning networks has become feasible using
 121 onboard computing units but retraining them, especially in an incremental and interactive way, is not
 122 yet possible. In the next section, we are introducing a visual pipeline that is capable of handling such
 123 constraints (real-time capability, incremental/interactive object learning, and low amount of training data)
 124 and describe our approach to improve it afterwards.

125 2.3 Visual Recognition System

126 The Interactive Object Learning (IOL) system⁴ is part of the open source software stack for the iCub
 127 robot⁵. Here, we focused on the Feature Extractor and Classifier⁶ (Figure 1).

128 The Feature Extractor is based on the *BLVC Reference CaffeNet*, which is provided by the Caffe library Jia
 129 et al. (2014). This network was trained on the ImageNet dataset Krizhevsky et al. (2012), which contained
 130 more than 1 Million high-resolution images from 1,000 object categories. The Feature Extractor generates a
 131 feature vector that is characteristic for the input image under a given DCNN. The vector corresponds to the
 132 vector representation of the highest convolution layer of the DCNN Pasquale et al. (2015b). A Multiclass

⁴ <https://github.com/robotology/iol>

⁵ <https://github.com/robotology>

⁶ <https://github.com/robotology/himrep>

133 Linear SVM was used for classification of feature vectors. The SVM is trained with a one-vs-all strategy
134 for 1,000 epochs and thus provides a classifier for each object.

135 To use the given implementation without the robot involved, we had to separate the recognition workflow
136 from the rest of the system to use the pipeline in a standalone manner. The separation was necessary to
137 provide an analysis environment that resembles the original system as closely as possible (see Pasquale
138 et al. (2015b)). Since the provided solution is a highly integrated system, we had to rearrange the workflow
139 to fit our needs. This way, we also decreased the processing time for the analysis. Due to performance
140 issues, we replaced the *linearClassifier* module from the pipeline with the *LinearSVC* implementation
141 from the sci-kit package Pedregosa et al. (2011) after ensuring that we achieve comparable results. All
142 feature vectors were generated upfront as we did not require human interaction for our experiments.

143 2.4 Data Augmentation

144 Data augmentation or preprocessing is a way for recognition methods to enhance input signals and to
145 make the recognition more robust against known transformations. It is a standard tool for image recognition.

146 A wide variety of data augmentations have been used to capture different types of invariance, such as
147 translation, rotation, mirroring, distortions, color, and light changes. Bhattacharyya (2011) provided a
148 brief overview of additional color image preprocessing techniques. Ahmad et al. (2017) use detexturized,
149 decolorized, edge enhanced, salient edge map based, and flip/rotate images to improve DCNN-based
150 recognition in visual searches. More specialized versions of preprocessing are available if targeted tasks
151 (e.g., in face recognition) can be narrowed down and underlying information can be modeled more precisely
152 (see Han et al. (2013); Zou et al. (2007)).

153 For general object recognition, computationally inexpensive augmentations that handle light changes are
154 limited and usually involve at least gamma and brightness corrections. For example, Fischer et al. (2016)
155 augment training samples with a gamma adjustment (-0.5, 0.1) and a brightness adjustment (-0.2, 0.2).
156 Dosovitskiy et al. (2015) chose a gamma value between 0.7 and 1.5, incorporated an additive brightness
157 with Gaussian augmentation and a contrast modification. Kim et al. (2019) use brightness and contrast in
158 low-illumination scenarios for video surveillance systems. Howard (2013) uses randomly changed contrast,
159 brightness, color, and random lighting noise to capture light change variance. They base their modifications
160 on Krizhevsky et al. (2012), which also provides the DCNN used in our work and should therefore already
161 capture some variance. However, as we will demonstrate, there is still room for improvement.

162 So far, we discussed augmentations in data space, which are well-established techniques. To a lesser
163 degree, augmentation methods in the feature space are explored. To understand data augmentation for
164 classification, Wong et al. (2016) used warping in the feature space to improve recognition on the MNIST
165 dataset. For our approach, we take the idea of image concatenation (e.g., Saitoh et al. (2017)), apply it as
166 feature concatenation and combine it with gamma and brightness modification.

167 While widely used, to the best of our knowledge, no systematic analysis of the impact of gamma and
168 brightness changes for object recognition have been conducted. The parameters of these modifications used
169 in the literature vary and their separate effects on the recognition are not determinable as they are usually
170 mixed with other types of augmentations and are not reported separately. In our work, we want to provide
171 a starting point for a more systematic analysis. Also, we did not incorporate light modification models,
172 which come with a too high computational cost for online learning such as Gabor filters as suggested by
173 Welke et al. (2006). Additionally, we will demonstrate that even though a DCNN was trained taking light
174 change augmentations into account, it still can benefit if used as a black-box feature generator. Important to

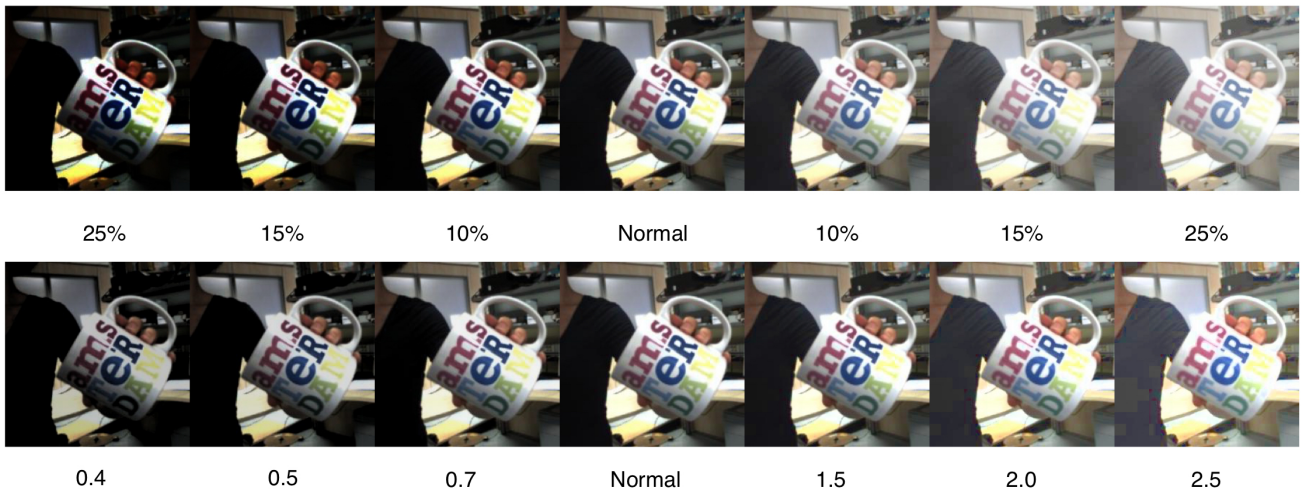


Figure 2. The unaltered sample (middle) and its transformations towards the extreme values: (top) linear model - example value changes in HSV color space with back-transformation into RGB color space, (bottom) non-linear model - example gamma value changes in RGB color space.

175 note is that we are not trying to enhance the feature generator itself as this would involve computationally
 176 expensive training. Instead, we are aiming to generate enhanced feature vectors to help the classification
 177 step under the assumption that creating multiple features from the images can be done in parallel and
 178 therefore do not add much to the computation time for the feature generation.

3 DATASETS AND METHODS

179 In this section, we describe the datasets, the models and the experiments parameters. As mentioned, we
 180 are focusing on global illumination changes as we consider them to be of greater importance than local
 181 ones. While local illumination changes such as reflection or shadow can have a high negative impact on
 182 the recognition as well, in the HRI scenario that we are addressing they are not necessarily persistent
 183 across multiple images during an interaction due to the changing orientation of the objects. Thus, local
 184 illumination changes are not considered in this paper.

185 The experiments are an extension of our earlier work which can be found in Keller and Lohan (2016).
 186 In our previous work, we used a small image dataset in combination with two illumination models to
 187 simulate linear and non-linear brightness changes to understand the impact of changes in light to the
 188 DCNN/SVM-based learning approach. In this paper, we compare the effect on a larger dataset with higher
 189 variability in object presentations as well as present a method to make use of the findings to improve the
 190 recognition process itself.

191 While a practical experiment could have been conducted to show the method's behavior under different
 192 illumination conditions, we chose to start with modified sample images of datasets for repeatability. This
 193 way, we are also able to look separately at linear and non-linear changes while this is much more difficult
 194 to achieve in an experimental setup. Both types of changes might occur at the same time in an experimental
 195 setup in a way that is not trivial to control. For example, a robot might use consumer camera sensors
 196 with a Automatic Gain Control (AGC) that influence the resulting image in a non-linear manner when
 197 the surrounding illumination is changed Fowler (2004). Due to the integration of current image sensors

198 themselves, it can be impossible to deactivate these assistance systems. Linear changes might occur when
199 blinds are used in different positions, limiting the amount of light from the outside.

200 3.1 DS1 - ICUBWORLD28

201 The first dataset is the ICUBWORLD28 dataset from Pasquale et al. (2015a) referred to as *DS1*. It
202 represents the visual perception of the iCub. It was created during a four-day interactive session. It consists
203 of nearly 40,000 images of 28 objects distributed over 7 object **classes** with more than 1,300 images per
204 object.

205 The dataset comes separated by day and is split into a training and test set per day. Since we were
206 interested in the overall performance of the approach, we merged all images per object into one set as this
207 gives the wide range of original illumination changes and allows for analyzing the sample size dependency
208 (see subsection 4.1). From the merged sets a number of training and test sets were randomly selected. First,
209 the 400 images for four test sets were chosen. Afterward, images for the different sized training sets were
210 selected. Thus, the results of the corresponding training and test sets are comparable to each other and are
211 based on balanced sets for all objects.

212 3.2 DS2 - ICUBWORLD TRANSFORMATION

213 The second dataset is the ICUBWORLD TRANSFORMATION Pasquale et al. (2016) referred to as *DS2*.
214 The dataset consists of more than 600,000 images with at least 3,000 images per object. It provides 5
215 different types of visual transformations; 2D, in-plane and 3D free rotations to provide the robot with
216 different viewpoints of the object, scaling transformation in which the human moves the object either
217 closer or further from the robot's position and a transformation in which the human is moving in a circle
218 around the robot to change the background while keeping approximately the same distance from the robot.
219 Additionally, a mixed transformation was included in which the object was presented in a free moving
220 manner. All objects are captured during two different days each. The dataset contains image sequences
221 from 15 object categories with 10 objects in each category giving a total of 150 objects (IROS 2016 subset
222 from Pasquale et al. (2016)).

223 We chose this dataset as it contains more variability in the presentation of the objects (e.g., a wider range
224 of illumination changes) and also contains more objects (150 vs. 28). This way we can test our method on
225 a more challenging task but can compare the results as the acquisition of the second dataset was similar to
226 the first. The preparation of the second dataset follows the method of the first one.

227 3.3 Illumination Change Models

228 To account for linear and non-linear light changes, we generated new images from the dataset. We chose
229 the value modification from the HSV color space as an example of a linear model. HSV stands for *hue*,
230 *saturation*, and *value*, where *value* accounts for brightness. The second model is given by the gamma
231 transformation which serves us as a non-linear modification. The transformation was done using OpenCV
232 Bradski (2000). Figure 2 shows examples for both modifications. While these changes seem to be easy for
233 the human eye they have an impact on recognition systems that operate on pixel values.

234 3.3.1 Linear Model

235 The images were transformed into the HSV color space for modification. This allows for changes to the
236 luminance without interfering with the colors and is one of the common color spaces for classical visual
237 recognition. The model is defined as $V_{out} = V_{in} \pm (V_{max} * V_c)$. After the color space transformation V_c

Table 1. Naming convention for the modified sets

Names	Sets
Linear condition	
0	$V = 0$
-25:1:25	$V \in \{-25\%, 0, 25\%\}$
Non-linear condition	
1	$\gamma = 1$
0.4:1:2.5	$\gamma \in \{0.4, 1, 2.5\}$
Combined condition	
mix	$\gamma \in \{0.4, 2.5\}$ and $V \in \{-25\%, 25\%\}$ and 1 set of unmodified images

238 was changed to 5%, 10%, 15% and 25% for both lighter and darker appearance. For the recognition task,
239 the images are transformed back into RGB color space before being fed to the DCNN.

240 3.3.2 Non-Linear Model

241 The gamma correction is a non-linear transformation that is often used to enhance the visual appearance
242 of images that are under- or overexposed. It is defined as $V_{out} = V_{in}^{1/\gamma}$. For the γ we chose 0.4, 0.5, 0.7 for
243 darker and 1.5, 2.0, 2.5 for brighter images.

244 3.3.3 Parameter Selection

245 The choice of the parameter values is based on selected representatives for non-trivial changes. If the
246 transformation effect gets much stronger differences in color can degrade towards black or white areas
247 destroying the included color information and leading to unrecognizable images. While this is an effect
248 that image recognition approaches have to deal with, it is not the focus of the paper in which we want to
249 improve recognition under different lighting conditions while being in a reasonably well-lit environment.
250 Table 1 gives an overview of the parameter sets we chose and how they are named in our paper.

251 3.3.4 Measurement

252 For our experiments, we report the average accuracy. All reported accuracies are an average over a
253 four-fold experiment run. Minimal and maximal errors can be found in the diagrams but are not reported
254 for a clearer presentation of the findings. Due to their small variability, they did not have an impact on the
255 conclusions. All experiments were supported by test sets with 400 samples per object.

4 EXPERIMENTS AND RESULTS

256 4.1 Sample Size Dependency

257 First, we established a baseline for the recognition task without any modifications to analyze the
258 performance depending on the number of training samples used (Figure 3).

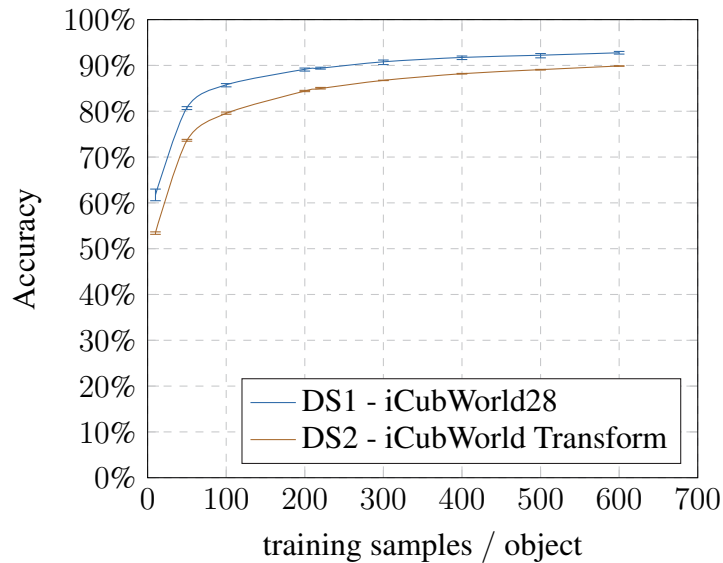


Figure 3. Baseline: Accuracy based on training samples per object.

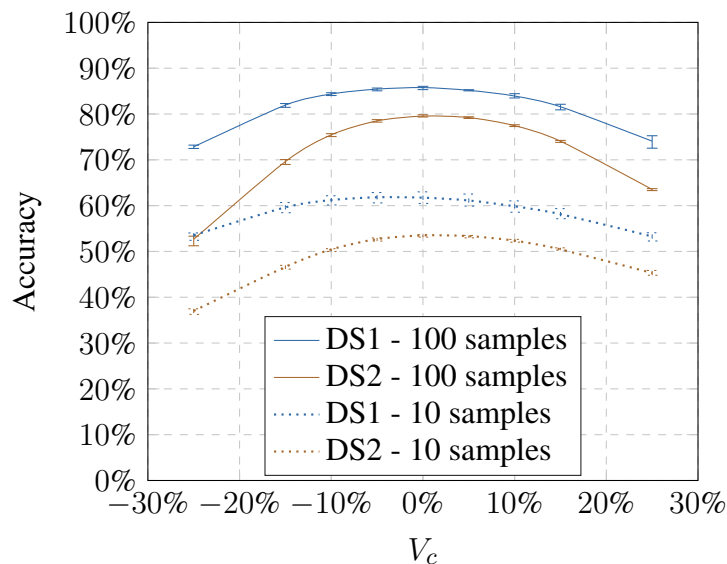


Figure 4. Influence of the linear model (V_c) on the accuracy showing the loss towards stronger transformation values.

259 With 10 training samples, the method achieves an avg. accuracy of 61.74% going up to 92.77% for 600
 260 images on the *DS1* (28 classes) and avg. accuracy ranging from 53.50% for 10 images up to 89.88% for
 261 600 images on the *DS2* (150 classes). The avg. accuracy is highly stable over all runs with a variation of
 262 $\pm 1.3\%$ maximum in the worst case (10 training samples / *DS1*). The baseline results show that the gain in
 263 recognition performance using more than 100 training images becomes comparably smaller (Figure 3). As
 264 expected, the recognition pipeline performs similarly on both datasets and benefits from larger training set
 265 sizes in the beginning. The absolute difference between the datasets can be explained due to the second one
 266 containing more challenging presentations as well as more object classes.

267 Next, we tested how the recognition rate changed if the visual pipeline has to recognize images with
 268 model-based altered brightness images. The baseline training was used and presented with altered test

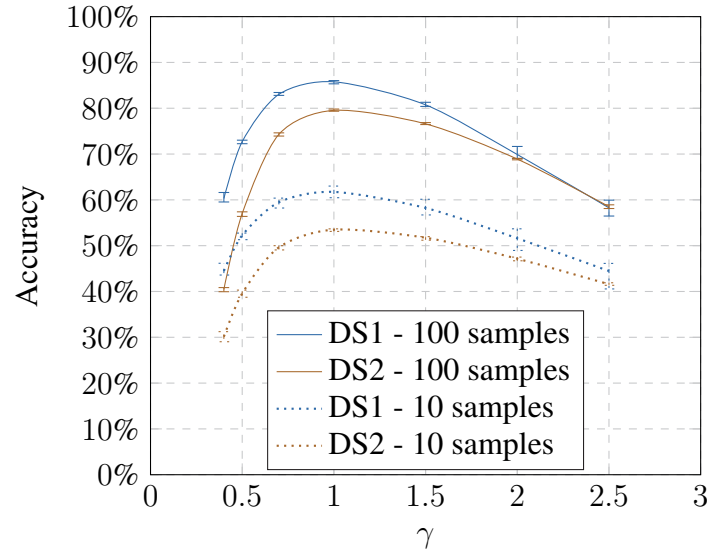


Figure 5. Influence of the non-linear model (γ) on the accuracy showing the loss towards stronger transformation values.

269 images based on the respective models (Figure 4 and Figure 5). The modified test image sets always contain
 270 only one modification to test the behavior on the modification limits. The results show that with an increase
 271 or decrease in brightness the recognition performs worse. The effect is consistent with Keller and Lohan
 272 (2016) although we used a random training and test sample selection this time.

273 Using a larger number of training samples seems preferable due to better recognition and thus compensat-
 274 ing the effects of the illumination influence. However, it defeats the purpose of our approach which is meant
 275 to improve the recognition using small training samples sets. Under real-world constraints, interactions
 276 with the objects might only happen over a very brief period. Additionally, we had to consider that the
 277 training time of a 1-vs-all Multiclass-SVM increases quadratically with the number of classes. The training
 278 with 600 samples per object results in training on 16,800 feature vectors for *DS1* and 90,000 for *DS2* for
 279 each object classifier (28 for *DS1* and 150 for *DS2*). We chose to perform the following experiments on
 280 100 training samples per object since our focus is on small training sets. However, in subsection 4.3, we
 281 will show that the results are generalizable and independent of training set sizes.

282 4.2 Brightness Dependency

283 After identifying the influences of brightness on the recognition process, we modified the training set to
 284 include altered images and tested the trained system against modified test sets (see Figure 6 and Figure 7).
 285 In Keller and Lohan (2016), we showed that it is sufficient to include only the modified images with the
 286 most extreme changes. Adding the modified images, the size of the sets increase from 10 to 30 and from
 287 100 to 300 images respectively. In our previous work, we have shown that this is a sufficient approach for a
 288 comparison (see Keller and Lohan (2016)).

289 For both models, we can show that adding these images increases the avg. accuracy drastically; in the best
 290 case we could achieve an increase from 58.26% up to 81.08% on *DS1* and from 40.43% up to 69.01% on
 291 *DS2* under the non-linear model (see Table 2). Recognition against the baseline training decreases slightly,
 292 suggesting that there is a trade-off between generalization and specificity.

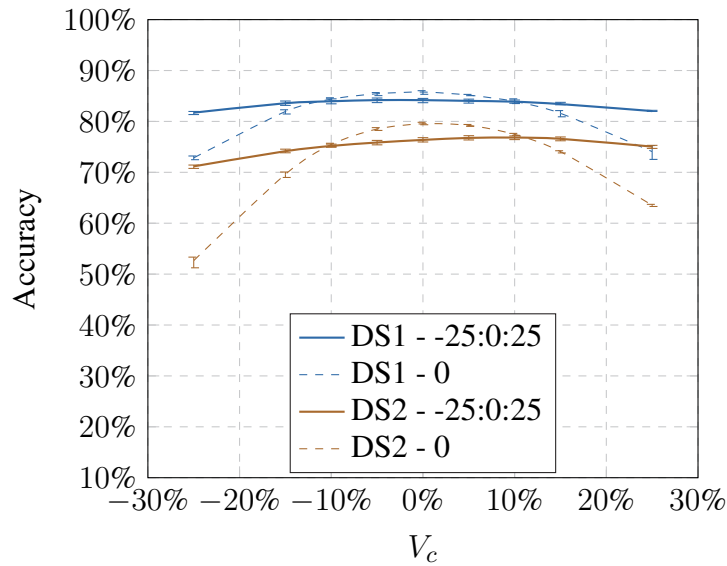


Figure 6. Influence of the brightness change: Using linear modified training samples shows an increase in recognition accuracy towards the extreme values while having a decrease for non-modified images.

Table 2. Selected avg. Accuracy values - modified training sets, 100 samples

<i>DS1</i>			
Linear condition	-25%	not modified	+25%
0	72.84%	85.77%	74.07%
-25:0:25	81.72%	84.18%	82.03%
Non-linear condition	0.4	not modified	2.5
1	60.27%	85.77%	58.26%
0.4:1:2.5	78.64%	82.93%	81.08%
<i>DS2</i>			
Linear condition	-25%	not modified	+25%
0	52.65%	79.54%	63.53%
-25:0:25	71.18%	76.37%	75.05%
Non-linear condition	0.4	not modified	2.5
1	40.43%	79.54%	58.53%
0.4:1:2.5	69.01%	75.70%	74.21%

293 With our first set of experiments, we showed that the investigated visual recognition pipeline is indeed
 294 susceptible to variations in illumination. We also show that by using modified training samples, the adverse
 295 effect of light changes can be circumvented to some degree. While this improvement comes at the cost of a
 296 small drop in recognition in the non-modified condition the overall positive effect justifies this modification
 297 as the recognition shows a much better generalization behavior across model-based changes and hence can
 298 be considered more robust.

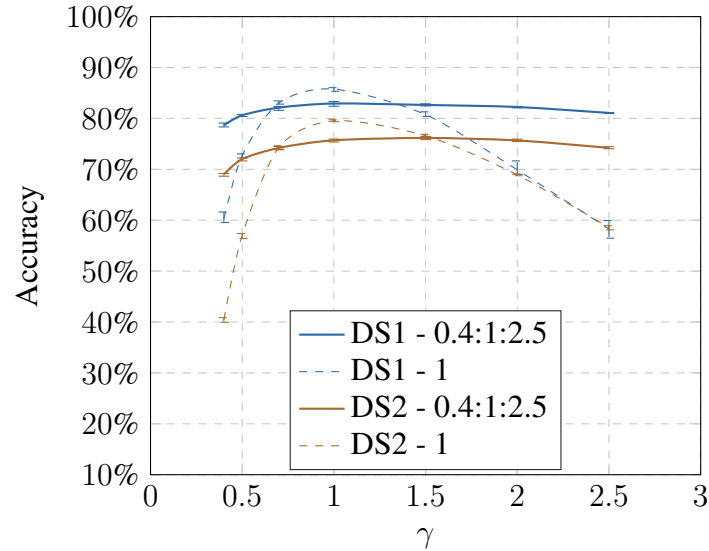


Figure 7. Influence of the gamma change: Using non-linear modified training samples shows an increase in recognition accuracy towards the extreme values while having a decrease for non-modified images.

299 4.3 Feature Fusion

300 In the last set of experiments, we want to make use of these findings to increase the overall performance
 301 of the recognition. So far, our experiments used the approach of training n sample images and testing
 302 individually against independent m test samples and relied on data augmentation in the input space.
 303 However, since we know which modified images belong to each other, we can make use of that additional
 304 knowledge to further improve the recognition. Therefore, we employ data augmentation in the feature
 305 space to bind the original images with their modifications. All feature vectors corresponding to one original
 306 image **and its augmented versions** are concatenated and fed into the SVM; both for training and testing.
 307 Thus, the dimensionality of the search space for the SVM increases and captures different light changes
 308 within one data sample. For example, the *mix* condition contains two feature vectors for the linear model,
 309 two for the non-linear and the non-modified feature vector hence the resulting feature vector is five times
 310 bigger than the ones from our other previous experiments.

311 Table 3 shows the results of the experiment for 10 and 100 training samples. For both training set sizes,
 312 the recognition improved taking the model-based modifications into account. For 10 training samples per
 313 object an increase of the avg. accuracy from 61.74% (baseline vs. baseline) up to 65.28% (mix vs. mix) and
 314 for 100 samples an increase from 85.77% to 90.00% are found for *DS1*. *DS2* shows an improvement from
 315 53.50% up to 57.53% for 10 training samples and from 79.54% up to 85.41% for 100 training samples.

316 In the last step, we compared the baseline results with our most optimal condition (mix vs. mix) (see
 317 Figure 8 and Figure 9). Here we show that the improvement is present in all training set sizes and that its
 318 effect is not dependent on the number of samples per object.

319 The results suggest that the additional information of the altered images are beneficial for the recognition
 320 process increasing the specificity of the method. While the individual models already achieve an improve-
 321 ment, the biggest gain can be seen while combining both models. However, the improvements come at the
 322 cost of a larger feature vector for training which increases the SVM's training time.

Table 3. Fused Feature Recognition avg. accuracy

Training \ Test Set	not modified	linear -25:0:25	non-linear 0.4:1:2.5	mix
<i>DSI- 10 samples</i>				
not modified	61.74%	60.68%	58.72%	59.40%
linear	60.72%	63.79%	-	-
non-linear	59.56%	-	65.09%	-
mix	60.19%	-	-	65.28%
<i>DSI- 100 samples</i>				
not modified	85.77%	84.88%	83.35%	84.10%
linear	83.14%	88.46%	-	-
non-linear	80.67%	-	89.68%	-
mix	81.54%	-	-	90.00%
<i>DS2- 10 samples</i>				
not modified	53.50%	51.28%	50.86%	50.48%
linear	50.59%	56.16%	-	-
non-linear	50.07%	-	57.35%	-
mix	49.69%	-	-	57.53%
<i>DS2- 100 samples</i>				
not modified	79.54%	77.37%	76.59%	76.27%
linear	73.80%	83.77%	-	-
non-linear	71.46%	-	84.74%	-
mix	71.31%	-	-	85.41%

5 DISCUSSION

323 As indicated by our previous work Keller and Lohan (2016), we have shown that illumination changes
 324 have an impact on state-of-the-art object recognition pipelines using DCNNs for feature generation and
 325 Multiclass-SVMs for classification. As expected, our results show that using more samples did improve the
 326 performance. However, our focus was on training with a small number of training samples to allow for
 327 very short interaction periods for data acquisition.

328 By expanding our work to a second dataset, we have shown that the impact of illumination changes from
 329 our first paper is generalizable. Especially, since the ICUBWORLD TRANSFORMATION dataset includes
 330 many more objects and different types of object manipulations in front of the robot together with higher
 331 variability in the background during object presentation.

332 Training with the artificial illumination models results in a slight drop of performance on the unmodified
 333 test sets but results in a major improvement under model-based illumination changes. While these findings
 334 are based on models and thus are not directly translatable to natural light changes, it proves that by using

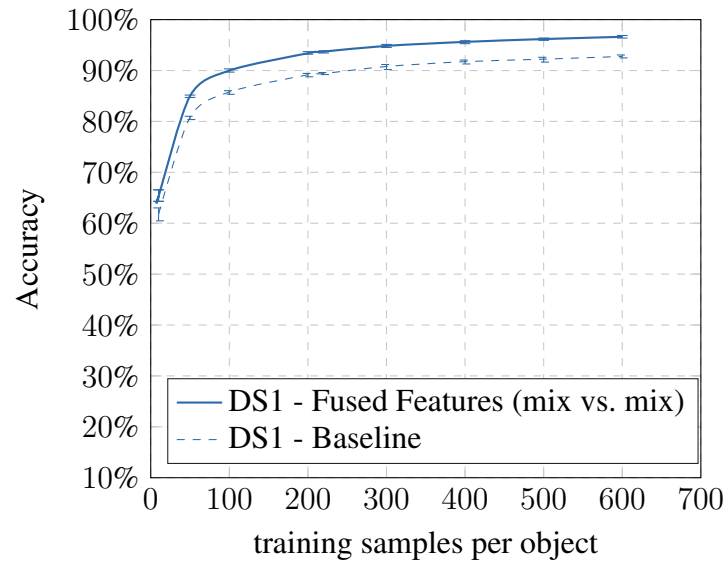


Figure 8. DS1 Baseline vs. Fused Features: Accuracy based on training samples per object.

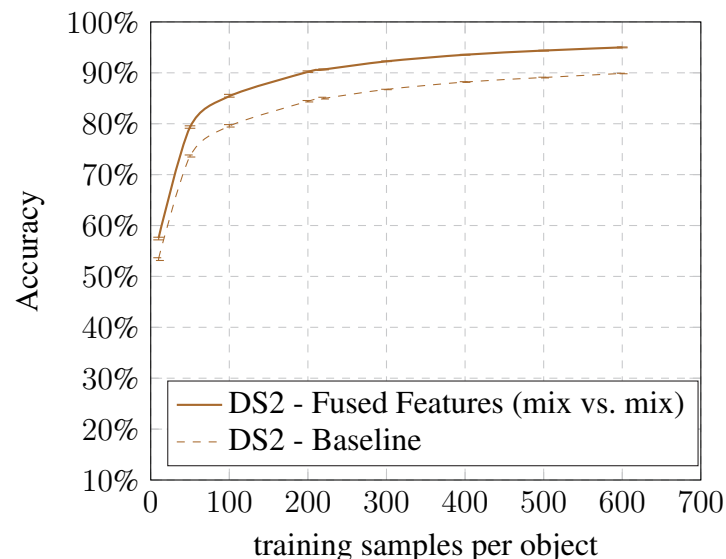


Figure 9. DS2 Baseline vs. Fused Features: Accuracy based on training samples per object.

335 data augmentations in the input space the recognition process can become more robust against light changes,
 336 resulting in a better generalization of the recognition pipeline.

337 By adding knowledge about the data generation for data augmentation in the feature space, namely
 338 concatenating corresponding feature vectors, we have shown how to make use of our findings to improve
 339 the recognition on both datasets which results in a more specialized recognition.

340 The data augmentations increased the accuracy between 3.54% (*DS1*) and 4.03% (*DS2*) for 10 training
 341 samples and between 4.23% (*DS1*) and 5.86% (*DS2*) for 100 training samples. Also, a similar increase
 342 is present in all training sample sizes and datasets. The effect might appear as a small improvement only.
 343 However, the results show in a systematic way which benefit can be expected from the used models and

344 can serve as a baseline to find better ones. It also shows that the impact increases with more diverse light
345 conditions (*DS1* vs. *DS2*).

346 Usually, improvements come at a cost. In our case, it is the additional computation with the *mix* condition
347 of the fused feature approach adding the highest amount. It involves generating four more images for the
348 transformations and the feature generation for them. The added computation for the transformations is
349 small compared to the feature generation. Since the images are independent and thus the feature generation
350 can run in parallel the added time cost for this step is small and easily fits into the online learning pipeline.
351 The limiting process is the second step as the SVM has to process a five-times-larger feature vector. While
352 this might render the data augmentation unattractive for large training sample sizes per object, it is still
353 viable for our focus area of small sample sizes. For example, in the case of 28 classes and 100 samples the
354 computation went up from 0.1s to 0.6s on average still being viable for near real-time purposes.

355 We believe that improving the recognition of objects with our methods during short interactions ($< 1min$)
356 will enhance the reliability of the overall system and therefore enhance the Human-Robot Interaction and
357 hence, the acceptance of the system.

6 CONCLUSION AND FUTURE WORK

358 With our model-based analysis, we showed that changes in illumination affect recognition approaches that
359 use DCNNs to encode features for object recognition. Using data augmentation, we were able to show that
360 a system **using DCNNs and SVMs can be modified towards a more robust recognition even though the**
361 **used DCNN already included an augmentation step towards intensity and color robustness.** Additionally,
362 we have demonstrated that using simple brightness change models can help to improve the recognition
363 across all training set sizes.

364 With our approach, it is easy to adapt existing visual recognition pipelines since only computationally
365 inexpensive data augmentations were used and no modification of either the feature encoder or the
366 classification is needed. **Treating the feature encoder as a black-box system allows to compare different**
367 **networks to the original setup which will be the subject of our future research.**

368 As the next step to improve this approach, we are looking into combining the models and integrating
369 more natural conditions into them. While the artificial models already enhanced the recognition, we believe
370 that with a more realistic representation of natural light changes our approach could be improved. However,
371 the choice of the simple models was due to their low computational cost. This trade-off needs to be taken
372 into account for real-time capable systems like the one we investigation in our paper.

373 To further improve the acquisition of training data for the objects, the lighting conditions could be
374 artificially altered to generate more diversity that could help the recognition process. When inspecting an
375 object for the first time, a flashlight with known spectral properties or RGB LEDs with defined colors could
376 be used. This might overcome the problem to find a model for real-world light changes as the additional
377 knowledge can be used to inform the data augmentation process.

378 Additionally, this approach could be used in conjunction with cloud robotics in which multiple robots
379 with different sensors and in different environments could combine their acquired images to cover more
380 diverse illumination settings.

7 ACKNOWLEDGEMENT

381 This work was supported by ORCA Hub EPSRC under Grant EP/R026173/1 and consortium partners.

REFERENCES

- 382 Ahmad, J., Muhammad, K., and Baik, S. W. (2017). Data augmentation-assisted deep learning of
383 hand-drawn partially colored sketches for visual search. *PLOS ONE* 12. doi:10.1371/journal.pone.
384 0183838
- 385 Bhattacharyya, S. (2011). A Brief Survey of Color Image Preprocessing and Segmentation Techniques.
386 *Journal of Pattern Recognition Research* 6, 120–129. doi:10.13176/11.191
- 387 Bradski, G. (2000). The opencv library. *Dr. Dobb's Journal of Software Tools* 25, 120–126
- 388 Churamani, N., Sutherland, A., and Barros, P. (2018). An Affective Robot Companion for Assisting the
389 Elderly in a Cognitive Game Scenario. *arXiv:1807.09825 [cs]*
- 390 Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., et al. (2015). FlowNet: Learning
391 Optical Flow with Convolutional Networks (Santiago, Chile: IEEE), 2758–2766. doi:10.1109/ICCV.
392 2015.316
- 393 Fischer, L., Hasler, S., Schrom, S., and Wersing, H. (2016). Improving Online Learning of Visual
394 Categories by Deep Features. In *30th Conference on Neural Information Processing Systems* (Barcelona,
395 Spain), 1–5
- 396 Fischler, M. A. and Bolles, R. C. (1981). Random sample consensus: a paradigm for model fitting with
397 applications to image analysis and automated cartography. *Communications of the ACM* 24, 381–395
- 398 Foster, D. H. (2011). Color constancy. *Vision Research* 51, 674–700. doi:10.1016/j.visres.2010.09.006
- 399 Fowler, K. R. (2004). Automatic gain control for image-intensified camera. *IEEE Transactions on*
400 *Instrumentation and Measurement* 53, 1057–1064. doi:10.1109/TIM.2004.831494
- 401 Han, H., Shan, S., Chen, X., and Gao, W. (2013). A comparative study on illumination preprocessing in
402 face recognition. *Pattern Recognition* 46, 1691–1699. doi:10.1016/j.patcog.2012.11.022
- 403 Howard, A. G. (2013). Some Improvements on Deep Convolutional Neural Network Based Image
404 Classification. 1–6
- 405 Imai, M., Ono, T., and Ishiguro, H. (2003). Physical relation and expression: joint attention for human-robot
406 interaction. *IEEE Transactions on Industrial Electronics* 50, 636–643. doi:10.1109/TIE.2003.814769
- 407 Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., et al. (2014). Caffe: Convolutional
408 Architecture for Fast Feature Embedding. In *Proceedings of the 22Nd ACM International Conference on*
409 *Multimedia* (New York, NY, USA), MM '14, 675–678. doi:10.1145/2647868.2654889
- 410 Keller, I. and Lohan, K. S. (2016). Analysis of illumination robustness in long-term object learning. In
411 *2016 25th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*
412 (New York, NY, USA: IEEE), 240–245. doi:10.1109/ROMAN.2016.7745137
- 413 Kim, I. S., Jeong, Y., Kim, S. H., Jang, J. S., and Jung, S. K. (2019). Deep Learning based Effective
414 Surveillance System for Low-Illumination Environments. In *2019 Eleventh International Conference on*
415 *Ubiquitous and Future Networks (ICUFN)*. 141–143. doi:10.1109/ICUFN.2019.8806120
- 416 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional
417 Neural Networks. In *Advances in Neural Information Processing Systems*, eds. F. Pereira, C. J. C.
418 Burges, L. Bottou, and K. Q. Weinberger (Curran Associates, Inc.). 1097–1105
- 419 Lafer-Sousa, R., Hermann, K. L., and Conway, B. R. (2015). Striking individual differences in color
420 perception uncovered by 'the dress' photograph. *Current Biology* 25, R545–R546. doi:10.1016/j.cub.
421 2015.04.053

- 422 Land, E. H. (1964). The Retinex. *American Scientist* 52, 247–253, 255–264
- 423 Le, H. and Kakadiaris, I. (2019). Illumination-Invariant Face Recognition With Deep Relit Face Images. In
424 2019 IEEE Winter Conference on Applications of Computer Vision (WACV). 2146–2155. doi:10.1109/
425 WACV.2019.00232
- 426 Meghdari, A., Shariati, A., Alemi, M., Vossoughi, G. R., Eydi, A., Ahmadi, E., et al. (2018). Arash: A
427 social robot buddy to support children with cancer in a hospital environment. *Proc Inst Mech Eng H*
428 232, 605–618. doi:10.1177/0954411918777520
- 429 Moldovan, B., Moreno, P., Otterlo, M. v., Santos-Victor, J., and Raedt, L. D. (2012). Learning relational
430 affordance models for robots in multi-object manipulation tasks. In *IEEE International Conference on*
431 *Robotics and Automation (ICRA)*. 4373–4378. doi:10.1109/ICRA.2012.6225042
- 432 [Dataset] Pasquale, G., Ciliberto, C., Odone, F., Rosasco, L., and Natale, L. (2015a). Real-world Object
433 Recognition with Off-the-shelf Deep Conv Nets: How Many Objects can iCub Learn?
- 434 Pasquale, G., Ciliberto, C., Odone, F., Rosasco, L., Natale, L., and dei Sistemi, I. (2015b). Teaching iCub
435 to recognize objects using deep Convolutional Neural Networks. In *Proceedings of the 4th Workshop on*
436 *Machine Learning for Interactive Systems*. 21–25
- 437 Pasquale, G., Ciliberto, C., Rosasco, L., and Natale, L. (2016). Object identification from few examples
438 by improving the invariance of a Deep Convolutional Neural Network. In *IEEE/RSJ International*
439 *Conference on Intelligent Robots and Systems (IROS)*. 4904–4911
- 440 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn:
441 Machine Learning in Python. *Journal of Machine Learning Research* 12, 2825–2830
- 442 Redmon, J. and Farhadi, A. (2018). YOLOv3: An Incremental Improvement. *arXiv:1804.02767 [cs]*
443 ArXiv: 1804.02767
- 444 Reyes, E., Gómez, C., Norambuena, E., and Ruiz-del Solar, J. (2018). Near Real-Time Object Recognition
445 for Pepper based on Deep Neural Networks Running on a Backpack. *arXiv:1811.08352 [cs]* ArXiv:
446 1811.08352
- 447 Saitoh, T., Zhou, Z., Zhao, G., and Pietikäinen, M. (2017). Concatenated Frame Image Based CNN for
448 Visual Speech Recognition. In *Computer Vision – ACCV 2016 Workshops* (Cham: Springer International
449 Publishing), vol. 10117. 277–289. doi:10.1007/978-3-319-54427-4_21
- 450 Sharif Razavian, A., Azizpour, H., Sullivan, J., and Carlsson, S. (2014). CNN Features Off-the-Shelf: An
451 Astounding Baseline for Recognition. In *Proceedings of the IEEE Conference on Computer Vision and*
452 *Pattern Recognition Workshops*. 806–813
- 453 Welke, K., Oztop, E., Ude, A., Dillmann, R., and Cheng, G. (2006). Learning feature representations for an
454 object recognition system. In *6th IEEE-RAS International Conference on Humanoid Robots*. 290–295.
455 doi:10.1109/ICHR.2006.321399
- 456 Wong, S. C., Gatt, A., Stamatescu, V., and McDonnell, M. D. (2016). Understanding data augmentation
457 for classification: when to warp? *arXiv:1609.08764 [cs]*
- 458 Yu, J., Zhang, F., and Xiong, J. (2014). An innovative sift-based method for rigid video object recognition.
459 *Mathematical Problems in Engineering*
- 460 Zou, X., Kittler, J., and Messer, K. (2007). Illumination Invariant Face Recognition: A Survey (Crystal
461 City, VA, USA: IEEE), 1–8. doi:10.1109/BTAS.2007.4401921