

---

# Supplementary material for “Improving neural network representations using human similarity judgments”

---

Anonymous Author(s)

Affiliation

Address

email

## 1 A Experimental details

### 2 A.1 Model features

3 We extract penultimate layer features of four different ImageNet-models — AlexNet [7], VGG-16  
4 [21], ResNet-18, and ResNet-50 [2] — and image encoder features of four different image/text models  
5 — CLIP RN50 and CLIP ViT-L/14 trained on WIT [16]; CLIP ViT-L/14 trained on Laion-400M  
6 [19] and Laion-2B [20] respectively. For extracting the model features, we use the Python library  
7 `thingsvision` [8].

### 8 A.2 gLocal probing

9 To optimize the gLocal transforms, we use standard SGD with momentum and perform cross-  
10 validation according to the procedure proposed in Muttenthaler et al. [10]. For finding the optimal  
11 gLocal transform, we perform an extensive grid search over four different hyperparameter values —  
12 the learning rate,  $\eta$ , the strength of the regularization term  $\lambda$ , the global-local trade-off parameter  
13  $\alpha$ , and the temperature parameter,  $\tau$ , used in the softmax expression for the local contrastive loss  
14 term (see Eq. 5). Specifically, we perform an extensive grid search over the Cartesian product of the  
15 following sets of hyperparameters:

- 16 •  $\eta \in \{0.0001, 0.001, 0.01, 0.1\}$ ,
- 17 •  $\lambda \in \{0.01, 0.1, 1.0, 10.0\}$ ,
- 18 •  $\alpha \in \{0.05, 0.1, 0.25, 0.5, 1.0\}$ ,
- 19 •  $\tau \in \{0.1, 0.25, 0.5, 1.0\}$ .

20 We use the same  $\eta$  and  $\lambda$  grids for global probing. We use PyTorch [11] for implementing the  
21 probes and PyTorch lightning to accelerate training. We choose the gLocal transform that achieves  
22 the lowest alignment loss (see alignment term in Eq. 6). Among the values in the above grid, we  
23 find that a combination of ( $\alpha = 0.1, \lambda = 0.1, \eta = 0.001$ ) yields the lowest alignment loss/highest  
24 probing odd-one-out accuracy for both CLIP RN50 and CLIP ViT-L/14 (WIT) (see Figure A.1).  
25 A combination of ( $\alpha = 0.25, \lambda = 0.1, \eta = 0.001$ ) gives the second lowest alignment loss/highest  
26 probing odd-one-out accuracy for CLIP RN50 and CLIP ViT-L/14 (WIT).

27 For each  $(\alpha, \lambda)$  combination we select that combination with the best probing odd-one-out accuracy  
28 on a held-out test among the set of possible learning rate,  $\eta$ , and temperature value,  $\tau$ , combinations  
29 determined by the above grid. We observe that  $\eta = 0.001$  generally gives the best results across  
30 the different  $(\alpha, \lambda)$  combinations, whereas performance is fairly insensitive to the value of  $\tau$ . Since  
31 neither  $(\alpha = 0.1, \lambda = 0.1)$  nor  $(\alpha = 0.25, \lambda = 0.1)$  are values at the edges of the hyperparameter

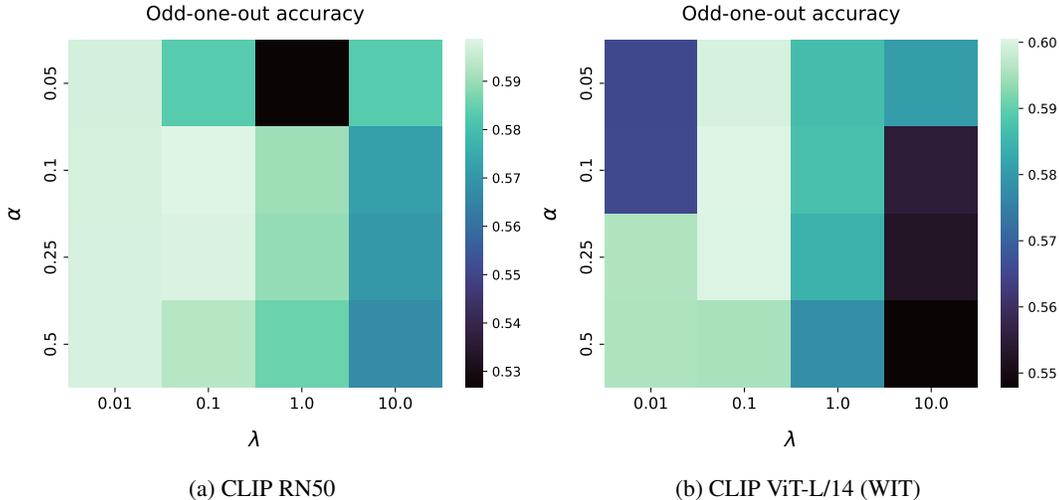


Figure A.1: Among all hyperparameter combinations considered in our grid search, a combination of ( $\alpha = 0.1, \lambda = 0.1, \eta = 0.001$ ) for Eq. 6 in §3 yields the best odd-one-out accuracy on a held-out test set for both CLIP RN50 and CLIP ViT-L/14 (WIT).

32 grid, it is plausible to assume that both the contrastive local loss and the regularization term in Eq. 6  
 33 in §3 are necessary to obtain a transformation that leads to a *best-of-both-worlds* representation.

34 Although our goal has been to find a transform that induces both increased representational alignment  
 35 and improved downstream task performance, we considered  $\alpha = 1.0$  to examine whether downstream  
 36 task performance can potentially be improved by excluding the alignment loss. Note that  $\alpha = 1.0$   
 37 causes the optimization process to ignore the alignment loss. Unsurprisingly we did not find that to be  
 38 the case. We remark that minimizing both the local contrastive loss and the regularization preserves  
 39 the local similarity structure of the original representation space but does not inject any additional  
 40 information into the representations. Moreover, it is non-trivial to choose a transform that works well  
 41 across all downstream tasks without including the alignment loss. Therefore, we exclude  $\alpha = 1.0$  in  
 42 Figure A.1.

43 **Compute.** We used a compute time of approximately 400 hours on a single Nvidia A100 GPU  
 44 with 40GB VRAM for all linear probing experiments — including the hyperparameter sweep. The  
 45 computations were performed on a standard, large-scale academic SLURM cluster.

### 46 A.3 Few-shot learning

47 Here, we use  $n_s$ -fold cross-validation for finding the optimal  $\ell_2$ -regularization parameter, where  
 48  $n_s$  refers to the number of shots per class. We select the parameter from the following set  
 49 of values,  $\{1e+6, 1, 1e+5, 1e+4, 1e+3, 1, 1e+2, 1e+1, 1, 1e-1, 1e-2, 1e-3, 1e-4\}$ . We use the  
 50 `scikit-learn` [12] implementation of (multinomial) logistic regression and refit the regression after  
 51 selecting the optimal regularization parameter.

52 **Compute.** We used a compute time of approximately 5600 CPU-hours of 2.90GHz Intel Xeon Gold  
 53 6326 CPUs for all few-shot experiments. Computations were performed on a standard, large-scale  
 54 academic SLURM cluster.

### 55 A.4 Anomaly Detection

56 In this section, we outline our anomaly detection experimental setting in more detail. In the anomaly  
 57 detection settings that we consider in our analyses *normal/anomaly* classes are determined via the  
 58 original classes in the data. Here, each of the original classes is once selected as a normal class  
 59 with the remaining classes being anomalous and, vice versa, each class in the data is once selected  
 60 as an anomalous class with the other classes being normal. After embedding the training images  
 61 from either the normal or the anomalous class in a model’s representation space, at inference time a  
 62 model must classify whether a new image belongs to the normal data or whether it deviates from

63 it and is thus considered an anomaly. For each example in the test set, a model yields an anomaly  
64 score where higher scores indicate more probability of an example being anomalous. Using the  
65 binary anomaly labels and the anomaly scores for each of the examples, we can then compute the  
66 *area-under-the-receiver-operating-characteristic-curve* (AUROC) to quantify the performance of the  
67 model.

68 **One-vs-rest.** Given a dataset (e.g., CIFAR-10) with  $C$  classes, one class (e.g., “airplane”) is chosen  
69 to be the normal class and the remaining  $C - 1$  classes of the dataset are considered anomalies. Each  
70 of the  $C$  classes is once selected as a normal class and the AUROC is averaged across the classes.

71 **Leave-one-out (LOO).** In contrast to the “one-vs-rest” setting, in LOO we define one class of the  
72 dataset as an anomaly and the remaining classes as normal. Similarly to the “one-vs-rest” setting,  
73 this results in  $C$  evaluations for a dataset with  $C$  classes.

74 In both “one-vs-rest” and LOO AD settings, we evaluate model representations in the following way:  
75 First, we compute the representations  $\mathbf{X}_{\text{train}}$  of the normal samples in the train set. Then, we compute  
76 the representations of all test set examples  $\mathbf{X}_{\text{test}}$ . For each test set representation, we compute the  
77 cosine similarity to all normal train set representations,  $\mathbf{X}_{\text{train}}$ , and select the  $k$  nearest neighbor  
78 samples that have the highest cosine similarity.

79 The anomaly score of a test set representation is then defined as the average cosine distance to the  $k$   
80 nearest train representations.  $k$  is a hyperparameter that determines the number of nearest neighbors  
81 over which the anomaly score is computed. We choose  $k = 5$  for our experiments but show that  
82 performance is fairly insensitive to the value of  $k$  (see Tab. D.6).

83 **Compute.** For all AD experiments, we used a compute time of approximately 20 hours on a single  
84 Nvidia A100 GPU with 40GB VRAM. Computations were performed on a standard, large-scale  
85 academic SLURM cluster.

## 86 **B What changes in the global structure of the representations after** 87 **alignment?**

88 In this section, we attempt to build some intuitions for how the global structure of the representations  
89 changes after alignment. To do so, we analyze the movements of the representations of items and  
90 superordinate categories in the THINGS dataset. Specifically, we compute cosine differences between  
91 the CLIP-ViT-L/14 representations of each pair of items in THINGS and then compute how these  
92 distances change under the transforms.

93 We show the pairs of items that change the most in distance in Table B.1. Items that are semantically  
94 related, like “curry” and “scrambled egg”, tend to move closer together, and therefore have trans-  
95 formed distances that are smaller than their original distance. By contrast, items like “handcuff” and  
96 “stethoscope”, which are semantically unrelated but perhaps have some slight visual similarity, tend  
97 to move farther apart. The distance changes under the gLocal transforms are correlated with, though  
98 generally less varied than, those under the naively-aligned transform.

99 To more broadly analyze the change in global structure, we then look at how distances between  
100 pairs of items change within and across superordinate categories (the top-down categories from  
101 THINGS). We show the results in Fig. B.1. Under the naively aligned transform, the items within  
102 each superordinate category tend to move slightly closer together — the diagonal is slightly blue —  
103 while the items from different categories tend to move substantially farther apart — the off-diagonal  
104 is mostly red. That is, the representations are broadly moving in a way that reflects the overall human  
105 semantic organization of the categories.

106 There are a few notable standouts: the categories of drink, food, plant, and animal change particularly  
107 much, and in particularly interesting ways. These categories each move much farther relative to all  
108 other categories (such as tool or musical instrument) than those other categories move relative to  
109 each other. This perhaps reflects the particular semantic salience of food, drink, plants, and animals  
110 from a human perspective. Furthermore, food and drink are one of the few pairs of superordinate  
111 categories between which distances actually decrease after the transform, presumably reflecting the  
112 strong semantic ties between these categories. Similarly, animals move less far from plants than from  
113 any other category, perhaps reflecting the fact that the animate/inanimate distinction is one of the  
114 strongest features in human semantic representations [18].

115 Under the gLocal transform, the pattern of changes is strongly correlated with the naively aligned  
 116 transform ( $r = 0.96, p \leq 10^{-16}$ ). However, in keeping with the regularization, the magnitude of the  
 117 changes varies less.

Table B.1: Distances between pairs of individual items from THINGS, ranked by the relative change in cosine distance from before to after naive alignment (normalized by original distance). The top items move much closer together under naive alignment, while the bottom ones move much farther apart. (All results are from CLIP-ViT-L/14.)

Item 1	Item 2	original dist.	naively aligned dist.	gLocal dist.
curry	scrambled egg	0.303120	0.005276	0.401019
otter	warthog	0.305242	0.005530	0.382150
parfait	spaghetti	0.457553	0.009115	0.540346
otter	rhinoceros	0.327497	0.006530	0.456641
		⋮		
stethoscope	wheat	0.263908	1.284535	0.935891
grass	wallet	0.277866	1.347424	1.056572
cat	traffic light	0.285151	1.378671	0.981944
handcuff	sugar cube	0.272936	1.308337	0.904380

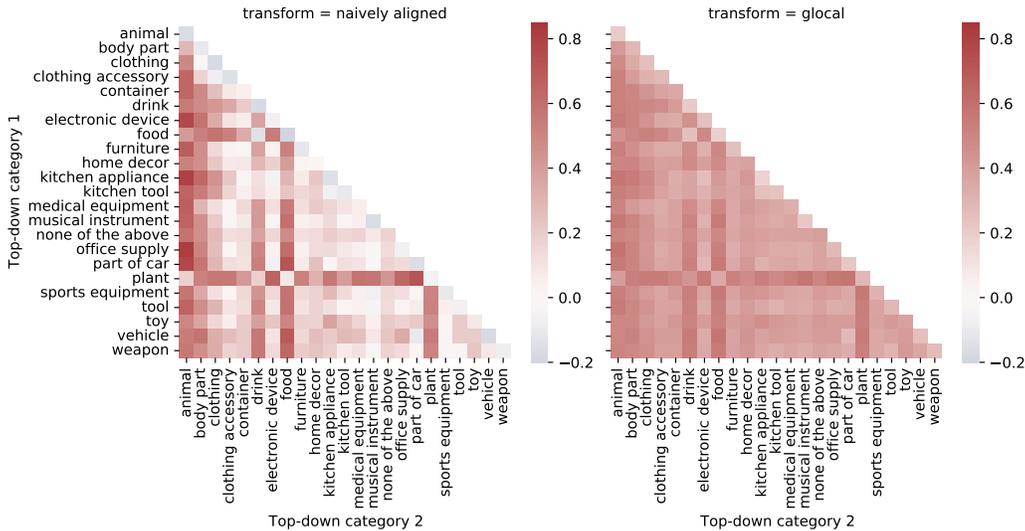


Figure B.1: How does the global structure of the representations change after alignment? Here, we analyze the movements of the representations of pairs of items from different superordinate categories from the THINGS dataset. The squares on the diagonal indicate the change in distance between items within a superordinate category, while the squares off the diagonal indicate changes between pairs of items from the corresponding pair of superordinate categories. A red color indicates the items from the categories move farther apart from each other after alignment, blue indicates moving closer together. Generally, items within a superordinate category move slightly closer together under naive alignment, while those in different categories move farther apart. A similar overall pattern is reflected in both the naively-aligned transform (left) and gLocal (right) ones, though under gLocal alignment there is a greater overall spreading of the representations. (All results are from CLIP-ViT-L/14.)

## 118 C Visualization of neighboring images

119 To provide further insight into the difference between the effects of the naive and gLocal transforms,  
 120 in Figure C.1 we visualize the neighbors of nine anchor images. In order to show a diverse set of  
 121 images, we pick the nearest neighbors in the CLIP ViT-L/14 (WIT) embedding space subject to the  
 122 constraint that each neighbor comes from a different class from the original images and the nearer  
 123 neighbors. In accordance with the results in §4.2, we find that the neighbors in the untransformed and  
 124 gLocal spaces are generally similar, whereas neighbors in the naive representation space are frequently

125 different. The naive transform appears to discard all non-semantic properties of images, whereas the  
126 untransformed and gLocal representation spaces are sensitive to pose, color, and numerosity. In cases  
127 where the closest neighbor differs between the naive and gLocal representations (third and fourth  
128 row), the neighbors in the gLocal representation are arguably better matches to the anchor.



Figure C.1: Comparison of neighbors in the ImageNet validation set for representations with different transforms. We visualize the 10 closest images subject to the constraint that each comes from a unique class. The anchor images are shown in the leftmost column. The three rows corresponding to each anchor image show their nearest neighbors in the untransformed, gLocal transformed, and naively transformed representations.

129 **D Additional results on downstream tasks**

130 In this section, we provide additional few-shot learning and anomaly detection results for all ImageNet  
 131 and image/text models that we considered in our analyses (see §A.1). We start this section by  
 demonstrating a strong relationship between the performances of the different downstream tasks.

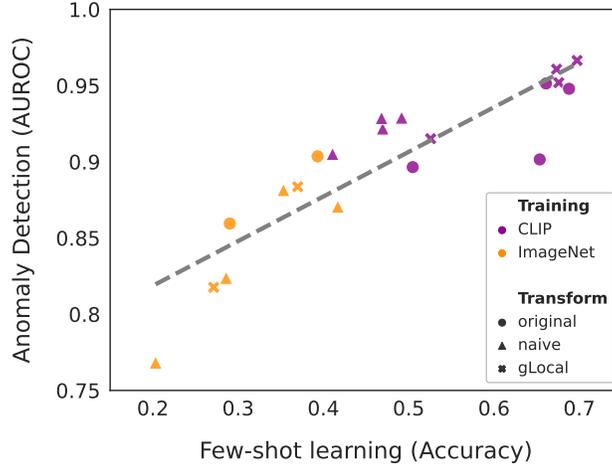


Figure D.1: Here, we show anomaly detection AUROC averaged across all tasks reported in Tab. {2, 3} as a function of the average 4-shot classification performance for all ImageNet and CLIP models (see §A.1), using either the original representations or the representations transformed via the naive or gLocal transformations.

132

133 **Downstream task relationship.** We observe a strong positive relationship ( $r = 0.8872, p \leq 10^{-7}$ )  
 134 between the average few-shot learning and the average anomaly detection performance for all  
 135 ImageNet and image/text models that we considered in our analyses (see Fig. D.1). This observation  
 136 holds for both the original representation space and the representations transformed via the naive or  
 137 gLocal transformations. This indicates that both downstream tasks require similar representations for  
 138 similarly strong performance.

139 **D.1 Few-shot learning**

140 In the following section, we show additional few-shot learning results. Specifically, we report 4-shot  
 141 performance of ImageNet models and show few-shot results as a function of the number of samples  
 142 used during fitting.

143 **Results for ImageNet models.** In Tab. D.1 we report additional 4-shot results for ImageNet models.  
 144 The gLocal transforms improve few-shot accuracy on Entity-{13,30} from BREEDS but the impact on  
 145 few-shot performance is either inconsistent or negative for CIFAR-100 coarse, CIFAR-100, SUN397,  
 146 and DTD of which the latter three are more fine-grained datasets than the other three.

Table D.1: 4-shot FS results using the original or transformed representations.

Model \ Transform	Entity-13		Entity-30		CIFAR100-Coarse		CIFAR100		SUN397		DTD	
	original	gLocal	original	gLocal	original	gLocal	original	gLocal	original	gLocal	original	gLocal
AlexNet	35.03	<b>39.59</b>	24.78	<b>25.85</b>	<b>30.51</b>	30.17	<b>26.26</b>	21.1	<b>24.19</b>	17.45	<b>33.39</b>	28.66
ResNet-18	56.15	<b>56.47</b>	<b>50.49</b>	50.03	38.3	<b>38.47</b>	<b>35.91</b>	34.42	<b>34.58</b>	33.17	<b>47.01</b>	44.28
ResNet-50	47.44	<b>51.41</b>	47.59	<b>50.22</b>	<b>48.2</b>	47.72	<b>45.29</b>	45.17	<b>44.69</b>	44.62	51.51	<b>51.85</b>
VGG-16	48.34	<b>54.76</b>	42.17	<b>44.04</b>	<b>36.74</b>	33.99	<b>31.77</b>	26.03	<b>34.55</b>	27.71	<b>42.35</b>	35.36

147 **Effect of transforms for different numbers of training samples.** When varying the number of  
 148 training samples for the few-shot experiments described in §4.3 we observe consistent improvements  
 149 of the gLocal transforms across shots. Excluding the high-variance setting of 2-shot learning, we  
 150 either find stable improvements in accuracy for image/text models, or a downward trend for ImageNet  
 151 models on some tasks. This corroborates our findings from §4.3. Results appear to be robust to  
 152 changes in the training set size, in particular for the CLIP models. Yet, we observe the most substantial  
 153 benefits in low data regimes. See Fig. D.2 for more details.

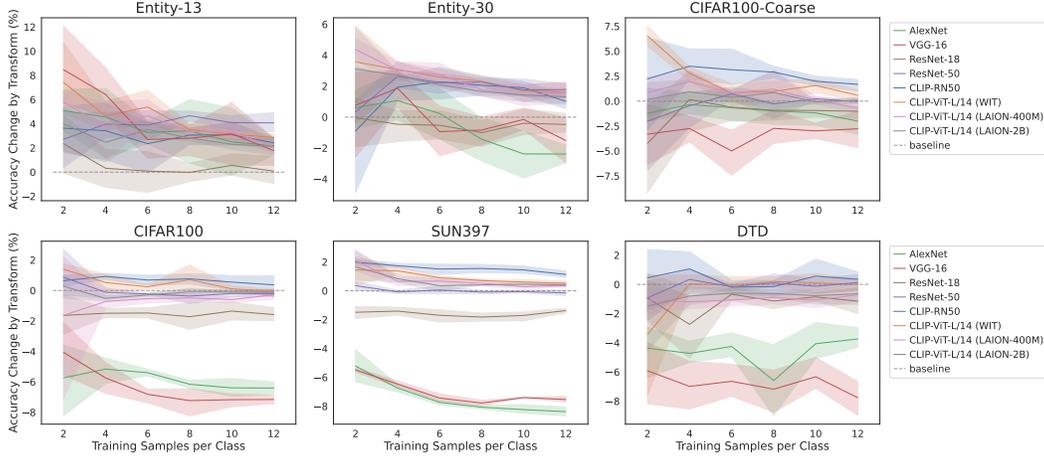


Figure D.2: Change in average accuracy for different numbers of training samples per super-class (top row) or (sub-)class (bottom row) used for few-shot learning. Error bands depict 95% Confidence Intervals (CIs), computed over 5 different runs.

## 154 D.2 Anomaly detection

155 In addition to the results of image/text models for the “one-vs-rest” anomaly detection (AD) setting  
 156 that we presented in §4.4, here we show “one-vs-rest” AD performance of ImageNet models. While  
 157 the gLocal transform considerably improves AD performance over the untransformed representations  
 158 across the different datasets for image/text models (see Tab. {2, 3}, for ImageNet models we do not  
 159 observe any improvements over the original representation space (see Tab. {D.2, D.3}).

160 Furthermore, we present results for the non-standard Leave-one-out (LOO) setting and for “CIFAR10-  
 161 vs-CIFAR100” for all image/text and Imagenet models that we considered. In the “CIFAR10-vs-  
 162 CIFAR100” AD task, all data of CIFAR10 is considered to be the normal class, and each sample from  
 163 the CIFAR100 dataset is considered an anomaly. Similarly to the previously reported AD results, the  
 164 gLocal transform substantially improves AD performance compared to the original representations  
 165 for image/text models across all datasets but does not appear to have a considerable impact on the  
 166 performance of ImageNet models (see Tab. {D.4, D.5}).

Table D.2: One-vs-rest nearest neighbor based AD results; with and without transformation. ImageNet30 results for ImageNet models are omitted due to overlap with train data.

Model \ Transform	CIFAR10		CIFAR100		CIFAR100-Coarse		ImageNet30		DTD	
	original	gLocal	original	gLocal	original	gLocal	original	gLocal	original	gLocal
AlexNet	<b>89.43</b>	85.63	<b>92.34</b>	88.53	<b>87.53</b>	82.75	–	–	<b>86.33</b>	79.51
ResNet-18	<b>92.19</b>	86.70	<b>95.06</b>	90.89	<b>92.16</b>	86.38	–	–	<b>94.38</b>	90.11
ResNet-50	<b>94.74</b>	94.13	<b>96.46</b>	96.18	<b>94.3</b>	94.03	–	–	<b>94.47</b>	94.42
VGG-16	<b>90.33</b>	88.00	<b>93.56</b>	91.97	<b>89.78</b>	88.16	–	–	<b>91.15</b>	85.5

Table D.3: One-vs-rest AD with a class distribution shift between train and test sets; with and without transformation.

Model \ Transform	Entity-13		Entity-30		Living-17		Nonliving-26		Cifar100-shift	
	original	gLocal	original	gLocal	original	gLocal	original	gLocal	original	gLocal
AlexNet	<b>83.84</b>	81.45	<b>85.38</b>	83.71	<b>87.04</b>	79.09	<b>81.45</b>	78.84	<b>80.21</b>	76.37
ResNet-18	<b>91.84</b>	89.45	<b>93.18</b>	91.6	<b>96.82</b>	93.1	<b>90.97</b>	89.87	<b>81.83</b>	77.44
ResNet-50	89.59	<b>91.26</b>	93.51	<b>93.86</b>	<b>98.27</b>	97.98	90.61	<b>91.85</b>	84.73	<b>85.38</b>
VGG-16	<b>89.78</b>	88.87	90.7	<b>91.56</b>	<b>94.72</b>	89.98	<b>89.78</b>	89.32	<b>83.42</b>	81.91

Table D.4: LOO nearest neighbor based AD results and ‘‘CIFAR-10 vs. CIFAR-100’’ AD results; with and without using the gLocal transform.

Model \ Transform	CIFAR10		CIFAR100		Cifar100-Coarse		Cifar10 vs Cifar100	
	original	gLocal	original	gLocal	original	gLocal	original	gLocal
AlexNet	<b>67.64</b>	62.7	<b>58.94</b>	55.83	<b>63.33</b>	59.37	<b>69.87</b>	68.01
ResNet-18	<b>72.35</b>	65.86	<b>64.86</b>	59.9	<b>71.38</b>	64.52	<b>81.42</b>	75.55
ResNet-50	<b>76.62</b>	75.36	<b>66.91</b>	66.03	<b>74.78</b>	73.96	<b>84.27</b>	84.17
VGG-16	<b>68.45</b>	64.31	<b>59.92</b>	57.89	<b>65.81</b>	64.28	73.55	<b>74.7</b>
CLIP-RN50	70.32	<b>72.46</b>	59.91	<b>61.43</b>	65.63	<b>68.07</b>	72.55	<b>76.78</b>
CLIP-ViT-L/14 (WIT)	84.91	<b>91.33</b>	67.08	<b>72.2</b>	73.48	<b>80.51</b>	85.24	<b>92.42</b>
CLIP-ViT-L/14 (LAION-400M)	<b>93.0</b>	92.37	74.05	<b>74.15</b>	82.13	<b>82.88</b>	94.44	<b>94.68</b>
CLIP-ViT-L/14 (LAION-2B)	93.55	<b>95.23</b>	76.88	<b>77.46</b>	84.67	<b>85.78</b>	93.18	<b>95.26</b>

Table D.5: LOO nearest neighbor based AD results; with and without using the gLocal transform.

Model \ Transform	Entity-13		Entity-30		Living-17		Non-Living-26	
	original	gLocal	original	gLocal	original	gLocal	original	gLocal
AlexNet	<b>62.05</b>	59.73	<b>58.2</b>	56.23	<b>61.02</b>	56.07	<b>56.27</b>	54.93
ResNet-18	<b>74.26</b>	68.88	<b>70.6</b>	64.7	<b>76.48</b>	70.79	<b>66.61</b>	63.82
ResNet-50	72.46	<b>73.67</b>	<b>73.37</b>	72.49	<b>83.5</b>	82.45	68.16	<b>68.29</b>
VGG-16	<b>70.26</b>	68.03	<b>66.38</b>	63.24	<b>73.31</b>	65.99	<b>64.87</b>	62.95
CLIP-RN50	69.99	<b>71.12</b>	63.49	<b>63.72</b>	<b>72.52</b>	70.04	62.55	<b>63.4</b>
CLIP-ViT-L/14 (WIT)	71.68	<b>74.88</b>	68.96	<b>70.58</b>	<b>82.9</b>	82.72	63.94	<b>67.33</b>
CLIP-ViT-L/14 (LAION-400M)	69.98	<b>71.44</b>	65.68	<b>66.26</b>	77.35	<b>77.4</b>	65.19	<b>66.27</b>
CLIP-ViT-L/14 (LAION-2B)	70.77	<b>72.49</b>	66.55	<b>67.68</b>	80.07	<b>80.09</b>	65.43	<b>67.99</b>

167 **The nearest neighbor hyperparameter  $k$ .** From the results reported in Tab. D.6 it can be inferred that  
 168 the nearest neighbor hyperparameter  $k$  does not have a considerable impact on AD task performance  
 169 across the different datasets. Here, we report the impact of  $k$  on the performance of CLIP ViT-L/14  
 170 (WIT) but the observation holds across all image/text models.

Table D.6: Nearest Neighbor AD performance of CLIP ViT-L/14 for different  $k$ .

Dataset \ Transform	$k=2$		$k=5$		$k=10$		$k=20$	
	original	gLocal	original	gLocal	original	gLocal	original	gLocal
CIFAR-10	95.37	<b>98.16</b>	95.14	<b>98.16</b>	94.86	<b>98.11</b>	94.50	<b>98.04</b>
CIFAR-100	91.90	<b>97.08</b>	91.41	<b>97.04</b>	90.93	<b>96.92</b>	90.39	<b>96.75</b>
CIFAR-100-coarse	89.28	<b>95.68</b>	88.50	<b>95.59</b>	87.73	<b>95.4</b>	86.81	<b>95.12</b>
CIFAR-100-shift	74.48	<b>86.18</b>	73.69	<b>86.17</b>	73.00	<b>86.02</b>	72.29	<b>85.82</b>
ImageNet30	98.95	<b>99.78</b>	98.91	<b>99.79</b>	98.85	<b>99.8</b>	98.78	<b>99.8</b>
Entity-13	88.37	<b>92.89</b>	88.54	<b>93.57</b>	88.45	<b>93.94</b>	88.28	<b>94.22</b>
Entity-30	91.26	<b>95.36</b>	91.31	<b>95.77</b>	91.22	<b>95.97</b>	91.03	<b>96.12</b>

### 171 D.3 Global versus gLocal transform

172 Aside from the AD performance of CLIP RN50 and CLIP ViT-L/14 (WIT), the gLocal transform  
 173 leads to more substantial improvements on downstream tasks than the global transform. In Tab D.7,  
 174 we report the average few-shot and anomaly detection performances using the global or gLocal  
 175 transforms. For FS, we average performance over all results reported in Tab. 1, and for AD we  
 176 average performance across all results reported in Tab. {2, 3, D.1}.

## 177 E Representational alignment

### 178 E.1 Human similarity judgments and RSMs

179 **Multi-arrangement task.** Human similarity judgments for King et al. [5] and [1] were obtained  
 180 by using a multi-arrangement task. In a multi-arrangement task, participants are presented with a  
 181 computer screen showing images of a number of different objects. The participants are asked to  
 182 arrange the images into semantically meaningful clusters, given the instruction that images of objects  
 183 that lie close together are considered more similar. From this arrangement, one can infer pairwise  
 184 (dis-)similarities of the objects and average those across all participants to obtain a representative  
 185 (dis-)similarity matrix.

Table D.7: Comparison of the average downstream task performance global and gLocal transforms.

Model \ Transform	AD		FS	
	global	gLocal	global	gLocal
AlexNet	81.16	<b>81.76</b>	26.65	<b>27.14</b>
ResNet-18	84.62	<b>88.39</b>	40.75	<b>42.80</b>
ResNet-50	93.19	<b>93.23</b>	48.29	<b>48.50</b>
VGG-16	87.32	<b>88.36</b>	35.36	<b>36.98</b>
CLIP-RN50	<b>92.12</b>	91.52	50.02	<b>52.57</b>
CLIP-ViT-L/14 (WIT)	<b>95.49</b>	95.14	65.80	<b>67.44</b>
CLIP-ViT-L/14 (LAION-400M)	95.72	<b>96.08</b>	66.82	<b>67.34</b>
CLIP-ViT-L/14 (LAION-2B)	96.33	<b>96.65</b>	69.73	<b>69.74</b>

186 **Ordinal scale.** In Peterson et al. [13, 14], pairwise similarity judgments were obtained by asking  
 187 human participants to rate the similarity of pairs of objects on an ordinal scale that ranges from 0  
 188 (“not similar at all”) to 10 (“very similar”). The pairwise similarity ratings can be averaged across the  
 189 different participants which in turn yields a matrix of similarities between pairs of objects.

190 **Triplet odd-one-out choices.** The triplet odd-one-out task is a commonly used task in the cognitive  
 191 sciences to infer pairwise object similarity ratings [17, 3, 9]. The triplet odd-one-out task is a  
 192 *three-alternative-forced-choice* task where participants are presented with three objects and have  
 193 to select the one that does not fit. In contrast to the multi-arrangement task or an ordinal scale,  
 194 the triplet odd-one-out task does not naturally yield a similarity matrix. A similarity matrix can  
 195 be obtained, however, by learning representations for the objects being used in the task from the  
 196 human responses. Variational Interpretable Concept Embeddings (VICE) — an approximate Bayesian  
 197 method for inferring mental representations of object concepts from triplet odd-one-out choices — is  
 198 a method that was specifically developed for that purpose. VICE uses variational inference to learn  
 199 representations for the objects in the triplets by fitting the human responses via stochastic gradient  
 200 descent. The method minimizes  $\mathcal{L}_{\text{global}}$  with additional non-negativity and sparsity constraints on the  
 201 representations. More details about the optimization can be found in Muttenthaler et al. [9]. From  
 202 the VICE solution, one can easily compute a representational similarity matrix (RSM). Specifically,  
 203 given learned object representations  $\mathbf{V} \in \mathbb{R}^{n \times d}$ , one first computes the dot-product similarity matrix  
 204  $\mathbf{S}_h := \mathbf{V}\mathbf{V}^\top$  and then exponentiate this matrix elementwise,  $\mathbf{S}'_h := \exp(\mathbf{S}_h)$ . One can then apply  
 205 the softmax function defined in Eq. 1 to every combination of triplets in the exponentiated similarity  
 206 matrix which yields the final RSM for triplet odd-one-out choices from Hebart et al. [4]. The last  
 207 step is performed to guarantee that the pairwise similarities are modeled according to the triplet  
 208 odd-one-out objective function that was used to learn the human object representations  $\mathbf{V}$  (see Eq. 2).

## 209 E.2 Neural network representations and RSMs

210 **Neural network representations.** RSMs for neural network representations are obtained by first  
 211 embedding the same set of images that were presented to the human participants in the  $p$ -dimensional  
 212 latent space of a neural net. The latent space could be any layer of a neural network. Here we use the  
 213 penultimate layer space for ImageNet models and the image encoder space for image/text models.  
 214 We do this because previous work has shown that the penultimate layer space of ImageNet models  
 215 and the image encoder space of image/text models respectively yield the highest similarity to human  
 216 behavior [14, 15, 10]. After embedding the images into the neural net’s latent space, one obtains  
 217 a representation matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$  for the  $n$  images in the data. Instead of simply computing the  
 218 dot-product similarity matrix  $\mathbf{S} := \mathbf{X}\mathbf{X}^\top$ , in RSA one typically uses either a cosine similarity or a  
 219 Pearson correlation kernel to compute the affinity matrix,

$$\cos(\mathbf{x}_i, \mathbf{x}_j) := \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2}; \quad \phi(\mathbf{x}_i, \mathbf{x}_j) := \frac{(\mathbf{x}_i - \bar{\mathbf{x}}_i)^\top (\mathbf{x}_j - \bar{\mathbf{x}}_j)}{\|(\mathbf{x}_i - \bar{\mathbf{x}}_i)\|_2 \|(\mathbf{x}_j - \bar{\mathbf{x}}_j)\|_2},$$

220 where the cosine similarity kernel function  $\cos(\mathbf{x}_i, \mathbf{x}_j)$  or the Pearson correlation kernel function  
 221  $\phi(\mathbf{x}_i, \mathbf{x}_j)$  is applied to every  $(\mathbf{x}_i, \mathbf{x}_j)$  vector pair of the matrix  $\mathbf{X}$  for obtaining the final representa-  
 222 tional similarity matrix  $\mathbf{S}' \in \mathbb{R}^{n \times n}$ . Here, we use the Pearson correlation kernel function  $\phi(\mathbf{x}_i, \mathbf{x}_j)$   
 223 to obtain a neural net’s RSM. Pearson correlation is the centered version of cosine similarity and  
 224 the ranking of the obtained similarities does not differ between the two kernel functions but Pearson  
 225 correlation first centers the vectors to have zero mean and is therefore a more robust measure. For

226 obtaining RSMs with transformed representations, the transforms are first applied to  $\mathbf{X}$  before  
 227 computing  $\mathbf{S}'$ .

### 228 E.3 Representational Similarity Analysis (RSA)

229 **Additional RSMs.** To corroborate our findings from §4.5, here we additionally show RSMs for CLIP  
 230 RN50 and CLIP ViT-L/14 (Laion 2B). In accordance with the different RSMs obtained from the  
 231 representation space of CLIP ViT-L/14 (WIT), there does not appear to be a qualitative difference  
 232 in the global similarity structure between the RSMs obtained from applying either the naive or the  
 233 gLocal transforms to CLIP RN50 or CLIP ViT-L/14 (Laion 2B) (see Fig. E.1). Hence, the gLocal  
 234 transform improves representational alignment while preserving the local similarity structure of the  
 original representation equally well for the different CLIP models, as we show in Tab. 4.

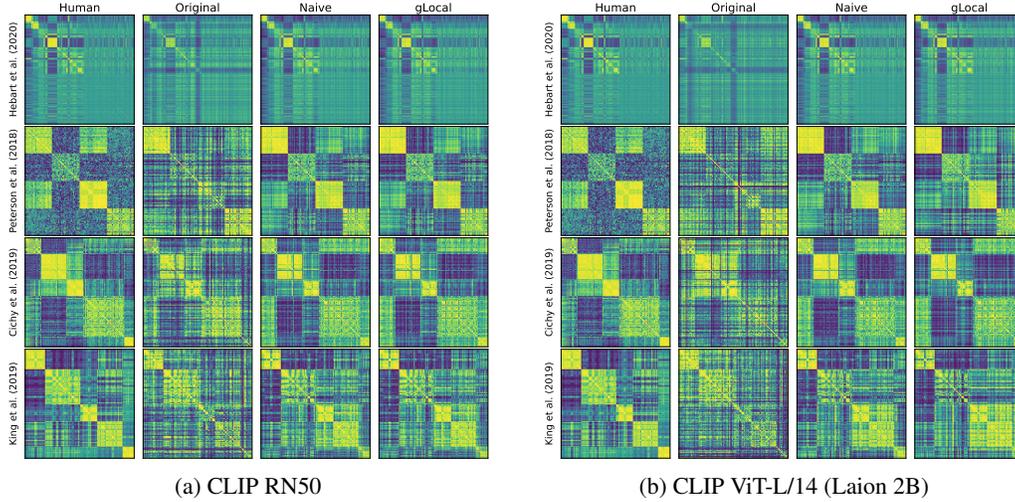


Figure E.1: Here, we show representational similarity matrices (RSMs) for human behavior and CLIP RN50 [WIT; 16] and CLIP ViT-L/14 [Laion 2B; 20] for four different human similarity judgment datasets [3, 14, 1, 5]. We contrast RSMs obtained from the network’s original representation space (second column), the naively transformed representation space [10] (third column), and the representation space obtained by using the gLocal transform (rightmost column) against RSMs directly constructed from human similarity judgments (leftmost column).

235

### 236 F Global transform derivation

237 Here we derive that

$$\min_{\alpha} \|\mathbf{W} - \alpha I\|_F^2 = \|\mathbf{W} - (\sum_{i=1}^p \mathbf{W}_{ii}/p) I\|_F^2.$$

238 First, observe that

$$\begin{aligned} \min_{\alpha} \|\mathbf{W} - \alpha I\|_F^2 &= \min_{\alpha} \sum_{i=1}^p \sum_{j=1}^p (\mathbf{W}_{ij} - \alpha \mathbb{1}_{[i=j]})^2 \\ &= \min_{\alpha} \sum_{i=1}^p \sum_{j=1, j \neq i}^p \mathbf{W}_{ij}^2 + \sum_{k=1}^p (\mathbf{W}_{kk} - \alpha)^2 \\ &= \sum_{i=1}^p \sum_{j=1, j \neq i}^p \mathbf{W}_{ij}^2 + \min_{\alpha} \sum_{k=1}^p (\mathbf{W}_{kk} - \alpha)^2. \end{aligned}$$

239 The minimizer of  $\min_{\alpha} \sum_{k=1}^p (\mathbf{W}_{kk} - \alpha)^2$  is attained with  $\alpha = \sum_{\ell=1}^p \mathbf{W}_{\ell\ell}/p$ . Substituting this  
 240 back into the last equality and reversing the steps from before we have

$$\begin{aligned} \sum_{i=1}^p \sum_{j=1, j \neq i}^p \mathbf{W}_{ij}^2 + \min_{\alpha} \sum_{k=1}^p (\mathbf{W}_{kk} - \alpha)^2 &= \sum_{i=1}^p \sum_{j=1, j \neq i}^p \mathbf{W}_{ij}^2 + \sum_{k=1}^p \left( \mathbf{W}_{kk} - \sum_{\ell=1}^p \mathbf{W}_{\ell\ell}/p \right)^2 \\ &= \sum_{i=1}^p \sum_{j=1}^p \left( \mathbf{W}_{ij} - \left( \sum_{\ell=1}^p \mathbf{W}_{\ell\ell}/p \right) \mathbb{1}_{[i=j]} \right)^2 \\ &= \|\mathbf{W} - (\sum_{\ell=1}^p \mathbf{W}_{\ell\ell}/p) \mathbf{I}\|_{\mathbb{F}}^2, \end{aligned}$$

241 which finishes our derivation.

## 242 G Properties of LCKA

243 Kornblith et al. [6] previously validated linear centered kernel alignment (LCKA) as a way to  
 244 measure similarity between neural network representations. Given representations  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  
 245  $\mathbf{Y} \in \mathbb{R}^{n \times p_2}$  containing embeddings of the same  $n$  images stacked row-wise, LCKA is:

$$\text{LCKA}(\mathbf{X}, \mathbf{Y}) = \frac{\langle \tilde{\mathbf{X}} \tilde{\mathbf{X}}^{\top}, \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^{\top} \rangle_{\mathbb{F}}}{\|\tilde{\mathbf{X}} \tilde{\mathbf{X}}^{\top}\|_{\mathbb{F}} \|\tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^{\top}\|_{\mathbb{F}}} = \frac{\|\tilde{\mathbf{X}}^{\top} \tilde{\mathbf{Y}}\|_{\mathbb{F}}^2}{\|\tilde{\mathbf{X}}^{\top} \tilde{\mathbf{X}}\|_{\mathbb{F}} \|\tilde{\mathbf{Y}}^{\top} \tilde{\mathbf{Y}}\|_{\mathbb{F}}}, \quad (1)$$

246 where  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$  are equal to  $\mathbf{X}$  and  $\mathbf{Y}$  with column means subtracted. (Formally,  $\tilde{\mathbf{X}} = \mathbf{H}\mathbf{X}$  and  
 247  $\tilde{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$  and  $\mathbf{H} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^{\top}$  is the centering matrix, which is a matrix representation of the linear  
 248 operator that subtracts column means.)

249 As Kornblith et al. [6] note, linear CKA can be thought of as measuring the cosine similarity between  
 250 all pairs of principal components (PCs) of  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$ , weighted by the products of the proportions of  
 251 variance these PCs explain in each representation. Formally, let  $\tilde{\mathbf{X}} = \mathbf{U}\Sigma\mathbf{V}^{\top}$  and  $\tilde{\mathbf{Y}} = \mathbf{U}'\Sigma'\mathbf{V}'^{\top}$   
 252 be the singular value decompositions of  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$ . The left-singular vectors  $\mathbf{u}_i = \mathbf{U}_{:,i}$  are the  
 253 (unit-norm) PCs of  $\mathbf{X}$ , and the squared singular values  $\lambda_i = \Sigma_{ii}^2$  are the amount of variance that  
 254 those PCs explain (up to a factor of  $1/n$ ). Given these singular value decompositions, linear CKA is:

$$\text{LCKA}(\mathbf{X}, \mathbf{Y}) = \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \lambda_i \lambda'_j (\mathbf{u}_i^{\top} \mathbf{u}'_j)^2}{\sqrt{\sum_{i=1}^{p_1} \lambda_i^2} \sqrt{\sum_{j=1}^{p_2} \lambda'_j^2}}. \quad (2)$$

## 255 References

- 256 [1] Radoslaw M. Cichy, Nikolaus Kriegeskorte, Kamila M. Jozwik, Jasper J.F. van den Bosch,  
 257 and Ian Charest. The spatiotemporal neural dynamics underlying perceived similarity for  
 258 real-world objects. *NeuroImage*, 194:12–24, 2019. ISSN 1053-8119. doi: [https://doi.org/10.](https://doi.org/10.1016/j.neuroimage.2019.03.031)  
 259 [1016/j.neuroimage.2019.03.031](https://doi.org/10.1016/j.neuroimage.2019.03.031).
- 260 [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for im-  
 261 age recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
 262 *Recognition (CVPR)*, June 2016.
- 263 [3] Martin N. Hebart, Charles Y. Zheng, Francisco Pereira, and Chris I. Baker. Revealing the multi-  
 264 dimensional mental representations of natural objects underlying human similarity judgements.  
 265 *Nature Human Behaviour*, 4(11):1173–1185, October 2020. doi: 10.1038/s41562-020-00951-3.
- 266 [4] Martin N Hebart, Oliver Contier, Lina Teichmann, Adam H Rockter, Charles Y Zheng, Alexis  
 267 Kidder, Anna Corriveau, Maryam Vaziri-Pashkam, and Chris I Baker. Things-data, a multimodal  
 268 collection of large-scale datasets for investigating object representations in human brain and  
 269 behavior. *eLife*, 12:e82580, feb 2023. ISSN 2050-084X. doi: 10.7554/eLife.82580.

- 270 [5] Marcie L. King, Iris I.A. Groen, Adam Steel, Dwight J. Kravitz, and Chris I. Baker. Similarity  
271 judgments and cortical visual responses reflect different properties of object and scene categories  
272 in naturalistic images. *NeuroImage*, 197:368–382, 2019. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2019.04.079>.  
273
- 274 [6] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey E. Hinton. Similarity of  
275 neural network representations revisited. In Kamalika Chaudhuri and Ruslan Salakhutdinov  
276 (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019,*  
277 *9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning*  
278 *Research*, pp. 3519–3529. PMLR, 2019.
- 279 [7] Alex Krizhevsky. One weird trick for parallelizing convolutional neural networks. *arXiv*  
280 *e-prints*, art. arXiv:1404.5997, April 2014.
- 281 [8] Lukas Muttenthaler and Martin N. Hebart. Thingsvision: A python toolbox for streamlining  
282 the extraction of activations from deep neural networks. *Frontiers in Neuroinformatics*, 15:45,  
283 2021. ISSN 1662-5196. doi: 10.3389/fninf.2021.679838.
- 284 [9] Lukas Muttenthaler, Charles Y Zheng, Patrick McClure, Robert A Vandermeulen, Martin N  
285 Hebart, and Francisco Pereira. VICE: Variational Interpretable Concept Embeddings. *Advances*  
286 *in Neural Information Processing Systems*, 35:33661–33675, 2022.
- 287 [10] Lukas Muttenthaler, Jonas Dippel, Lorenz Linhardt, Robert A Vandermeulen, and Simon Korn-  
288 blith. Human alignment of neural network representations. In *11th International Conference*  
289 *on Learning Representations, ICLR 2023, Kigali, Rwanda, Mai 01-05, 2023*. OpenReview.net,  
290 2023.
- 291 [11] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
292 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
293 Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
294 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style,  
295 high-performance deep learning library. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelz-  
296 imer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural*  
297 *Information Processing Systems 32: Annual Conference on Neural Information Processing*  
298 *Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 8024–8035,  
299 2019.
- 300 [12] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,  
301 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-  
302 learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830,  
303 2011.
- 304 [13] Joshua C. Peterson, Joshua T. Abbott, and Thomas L. Griffiths. Adapting deep network features  
305 to capture psychological representations. In Anna Papafragou, Daniel Grodner, Daniel Mirman,  
306 and John C. Trueswell (eds.), *Proceedings of the 38th Annual Meeting of the Cognitive Science*  
307 *Society, Recognizing and Representing Events, CogSci 2016, Philadelphia, PA, USA, August*  
308 *10-13, 2016*. [cognitivesciencesociety.org](http://cognitivesciencesociety.org), 2016.
- 309 [14] Joshua C. Peterson, Joshua T. Abbott, and Thomas L. Griffiths. Evaluating (and improving) the  
310 correspondence between Deep Neural Networks and Human Representations. *Cogn. Sci.*, 42(8):  
311 2648–2669, 2018. doi: 10.1111/cogs.12670.
- 312 [15] Joshua C. Peterson, Ruairidh M. Battleday, Thomas L. Griffiths, and Olga Russakovsky. Human  
313 uncertainty makes classification more robust. In *2019 IEEE/CVF International Conference*  
314 *on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pp.  
315 9616–9625. IEEE, 2019. doi: 10.1109/ICCV.2019.00971.
- 316 [16] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
317 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
318 Sutskever. Learning transferable visual models from natural language supervision. In Marina  
319 Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine*  
320 *Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8748–8763. PMLR,  
321 18–24 Jul 2021.

- 322 [17] Rocco Robilotto and Qasim Zaidi. Limits of lightness identification for real objects under  
323 natural viewing conditions. *Journal of Vision*, 4(9):9–9, 09 2004. ISSN 1534-7362. doi:  
324 10.1167/4.9.9.
- 325 [18] Timothy T Rogers, James L McClelland, et al. *Semantic cognition: A parallel distributed*  
326 *processing approach*. MIT press, 2004.
- 327 [19] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton  
328 Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M:  
329 open dataset of clip-filtered 400 million image-text pairs. *CoRR*, abs/2111.02114, 2021.
- 330 [20] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman,  
331 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick  
332 Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk,  
333 and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text  
334 models. In *NeurIPS*, 2022.
- 335 [21] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
336 image recognition. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on*  
337 *Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track*  
338 *Proceedings*, 2015.