

The pitfalls of measuring representational similarity using representational similarity analysis

Marin Dujmović^{1*}, Jeffrey S Bowers¹, Federico Adolfi^{1,2}, and Gaurav
Malhotra¹

¹*School of Psychological Science, University of Bristol, Bristol, UK*

²*Ernst-Strüngmann Institute for Neuroscience in Cooperation with Max-Planck Society,
Frankfurt, Germany*

**marin.dujmovic@bristol.ac.uk*

Abstract

A core challenge in cognitive and brain sciences is to assess whether different biological systems represent the world in a similar manner. Representational Similarity Analysis (RSA) is an innovative approach to address this problem and has become increasingly popular across disciplines ranging from artificial intelligence to computational neuroscience. Despite these successes, RSA regularly uncovers difficult-to-reconcile and contradictory findings. Here, we demonstrate the pitfalls of using RSA and explain how contradictory findings arise due to false inferences about representational similarity based on RSA-scores. In a series of studies that capture increasingly plausible training and testing scenarios, we compare neural representations in computational models, primate cortex and human cortex. These studies reveal two problematic phenomena that are ubiquitous in current research: a “mimic” effect, where confounds in stimuli can lead to high RSA-scores between provably dissimilar systems, and a “modulation effect”, where RSA-scores become dependent on stimuli used for testing. Since our results bear on a number of influential findings and the inferences drawn by current practitioners in a wide range of disciplines, we provide recommendations to avoid these pitfalls and sketch a way forward to a more solid science of representation in cognitive systems.

Introduction

How do other animals see the world? Do different species represent the world in a similar manner? How do the internal representations of AI systems compare with humans and animals? The traditional scientific method of probing internal representations of humans and animals (popular in both psychology and neuroscience) relates them to properties of the external world. By moving a line across the visual field of a cat, Hubel & Wisel [1] found out that neurons in the visual cortex represent edges moving in specific directions. In another Nobel-prize winning work, O’Keefe, Moser & Moser [2,3] discovered that neurons in the hippocampus and entorhinal cortex represent the location of an animal in the external world. Despite these successes it has proved difficult to relate internal representations to more complex properties of the world. Moreover, relating representations across individuals and species is challenging due to the differences in experience across individuals and differences of neural architectures across species.

These challenges have led to recent excitement around Representation Similarity Analysis (RSA) which appears to overcome many of these obstacles. RSA usually takes patterns of activity from two systems and computes how the distances between activations in one system correlate with the distances between corresponding activations in the second system (see Figure 1). Rather than compare each pattern of activation in the first system directly to the corresponding pattern of activation in the second system, it computes a second-order measure of similarity, comparing the systems based on their *representational geometries*. The advantage of looking at representational geometries is that one no longer needs to match the architecture of two systems, or even the format of the initial activity patterns (see Supplementary Information, Section A for a brief history of RSA and its philosophical origins). One could compare, for example, fMRI signals with single cell

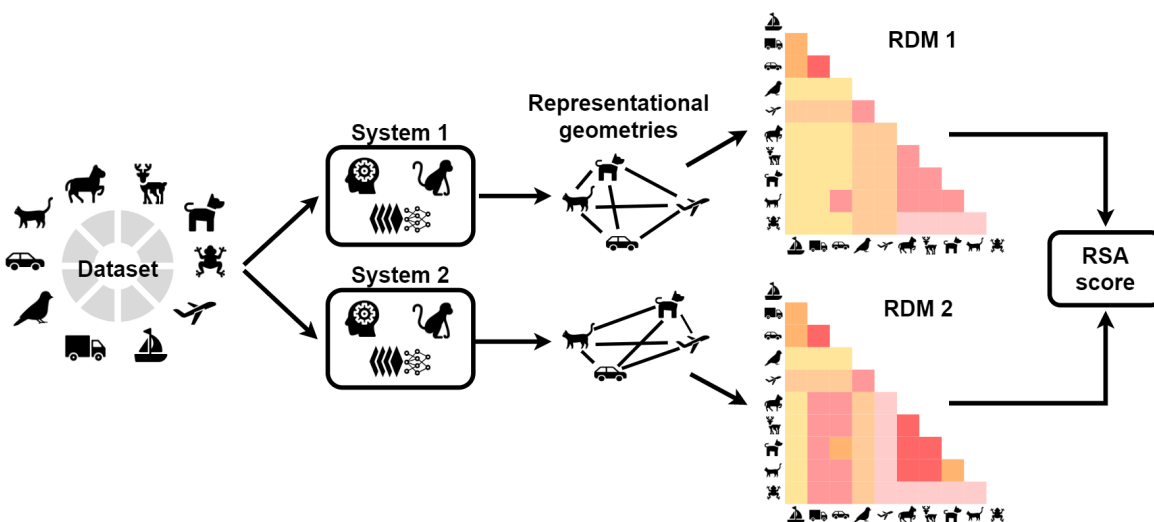


Figure 1: **RSA calculation.** A series of stimuli from a set of categories (or conditions) are used as inputs to two different systems (for example, a human brain and a primate brain). Activity from regions of interest is recorded for each stimulus. Pair-wise distances in activity patterns are calculated to get the representational geometry of each system. This representational geometry is expressed as a representational dissimilarity matrix (RDM) for each system. Finally, an RSA score is determined by computing the correlation between the two RDMs.

recordings, EEG traces with behavioural data, or vectors in a computer algorithm with 25
spiking activity of neurons [4]. RSA is now ubiquitous in computational psychology and 26
neuroscience and has been applied to compare object representations in humans and pri- 27
mates [5], representations of visual scenes by different individuals [6, 7], representations of 28
visual scenes in different parts of the brain [8], to study specific processes such as cognitive 29
control [9] or the dynamics of object processing [10], and most recently, to relate neuronal 30
activations in human (and primate) visual cortex with activations of units in Deep Neural 31
Networks [11–15]. 32

However, some recent research suggests that RSA may be an unreliable measure of how 33
similarly two systems represent the world. For example, many studies [16–20] have shown 34
that Convolutional Neural Networks (CNNs), trained on standard image datasets, such 35
as ImageNet, classify input images based on shortcuts, such as their texture. Activations 36
in these same networks also show a high RSA with activations in the human and primate 37
inferior temporal cortex [11, 12], even though it is well-known that humans primarily 38
represent objects based on their global properties such as shape, rather than shortcuts, 39
such as texture [21–23]. Similarly, some studies using RSA have shown that the hierarchy 40
of representations in the ventral visual stream in humans and primates correlates with 41
the hierarchy of representations in the layers of a CNN – i.e., deeper layer in a CNN 42
have a higher RSA with deeper layer in the visual ventral stream [11]. But Xu & Vaziri- 43
Pashkam [24] have recently shown that this correspondence is dataset-dependent and does 44
not replicate for some naturalistic and artificial stimuli. 45

How is it possible for two systems to have a high RSA score but represent different 46
features of inputs? Through a series of studies that capture increasingly plausible training 47
and testing scenarios, we demonstrate the properties of datasets and procedures that, in 48
practice, lead to high RSA scores between mechanistically dissimilar systems. The ex- 49
periments showcasing these pitfalls span the entire spectrum from artificial intelligence to 50
computational neuroscience, involving comparisons within and between sets of artificial 51
and biological systems. In particular, we shed light on two problematic phenomena that 52
bear on any efforts to compare systems based on RSA: 1) the presence of confounds in the 53
training data which leads systems to mimic each other’s representational geometry even 54
in the absence of mechanistic similarity, 2) the artifactual modulation of RSA scores due 55
to the intrinsic structure of datasets rather than system alignment. Our demonstrations 56
provide an explanation of how these phenomena, which arise ubiquitously, underlie con- 57

tradictory and paradoxical findings in the literature. Since our results have considerable 58
generality with respect to current practices across multiple fields, we discuss the implica- 59
tions for published results and prevailing interpretations, and provide broadly applicable 60
recommendations to move forward. 61

Results 62

Proof of concept 63

It may be tempting to infer that two systems which have similar representational geome- 64
tries for a set of concepts do so because they encode similar properties of sensory data 65
and transform sensory data through similar set of functions. In this section, we show that 66
it is possible, at least in principle, for qualitatively different systems to end up with very 67
similar representational geometries even though they (i) transform their inputs through 68
very different functions, and (ii) select different features of inputs. 69

**Study 1: Demonstrably different transformations of inputs can lead to low 70
or high RSA-scores** We start by considering a simple two-dimensional dataset and 71
two systems where we know the closed-form functions that project this data into two 72
representational spaces. This simple setup helps us gain a theoretical understanding of 73
the circumstances under which it is possible for qualitatively different projections to show 74
similar representational geometries. 75

Consider a set of stimuli, $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ from two classes that form two clusters in the 76
input space as shown in Figure 2A. Let us assume that each stimulus, \mathbf{x}_i contains mul- 77
tiple features that independently predict the class of the stimulus. We will call each of 78
these predictive features *confounds*. For example, shape and texture can be confounds 79

when classifying an image as belonging to DOG or AEROPLANE classes if each feature can
be independently used to predict whether an image belongs to the DOG or AEROPLANE
class. Consider two recognition systems Φ_1 and Φ_2 that map each input stimulus, \mathbf{x}_i ,
to an internal representation using their respective transformation functions, $\Phi_1(\mathbf{x}_i)$ and
 $\Phi_2(\mathbf{x}_i)$. Furthermore, we will assume that Φ_1 and Φ_2 are qualitatively different functions
and act on different features of the input. We are interested in showing that such qualita-
tively different functions acting on different features can nevertheless end up with similar
representational geometries.

The representational distance, $d[\mathbf{x}_i, \mathbf{x}_j]$, between the projections of any pair of input
stimuli, \mathbf{x}_i and \mathbf{x}_j , is proportional to the inner product between their projection in the
feature space:

$$d[\mathbf{x}_i, \mathbf{x}_j] \propto \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) \quad (1)$$

Thus, we can obtain the representational geometry of the input stimuli, $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, by
computing the pairwise distances, $d[\mathbf{x}_i, \mathbf{x}_j]$ for all pairs of data points, (i, j) . Here, we
assume that the projections Φ_1 and Φ_2 are such that these pairwise distances are given
by two positive semi-definite kernel functions $\kappa_1(\mathbf{x}_i, \mathbf{x}_j)$ and $\kappa_2(\mathbf{x}_i, \mathbf{x}_j)$, respectively:

$$\kappa_1(\mathbf{x}_i, \mathbf{x}_j) = \Phi_1(\mathbf{x}_i) \cdot \Phi_1(\mathbf{x}_j) \quad (2)$$

$$\kappa_2(\mathbf{x}_i, \mathbf{x}_j) = \Phi_2(\mathbf{x}_i) \cdot \Phi_2(\mathbf{x}_j) \quad (3)$$

Now, let us consider two qualitatively different kernel functions: $\kappa_1(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$ is a
radial-basis kernel (where σ^2 is the bandwidth parameter of the kernel), while $\kappa_2(\mathbf{x}_i, \mathbf{x}_j) =$
 $\frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$ is a cosine kernel. Figure 2A shows a dataset of points in a 2D input space that

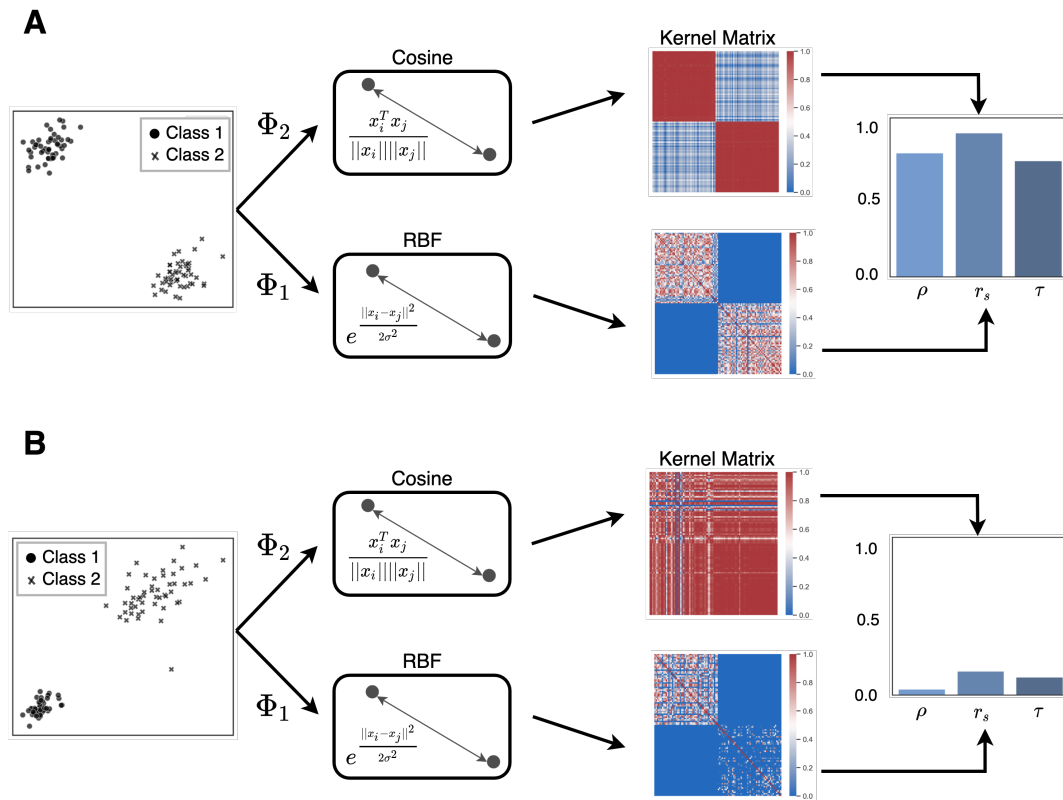


Figure 2: RSA between two systems with known transformations. In each panel a set of 2D stimuli are transformed using two different functions (Φ_1 and Φ_2), which project these stimuli into two different representational spaces. The distance between these projections are given by the RBF and Cosine kernels, respectively (see main text). The geometry of these projections can be visualised using the kernel matrices, which show the pair-wise distances between all stimuli in the representational space. The bar graph on the right-hand-side shows the RSA-score computed as a Pearson correlation (ρ), Spearman's rank correlation (r_s) and Kendall's rank correlation (τ). We can see that the input stimuli in Panel A leads to a high correlation in the representational geometry of the two systems, while the input stimuli in Panel B leads to a low correlation, even though the transformations remain the same.

are projected by two different systems into a cosine and RBF kernel space. Since the 91
cosine and RBF kernels are Mercer kernels [25, 26], each kernel matrix in Figure 2A shows 92
the pairwise distances (as measured by the inner product) between data points projected 93
in the two feature spaces. We can determine how the geometry of these projections in 94
the two systems relate to each other by computing the correlation between the kernel 95
matrices, shown on the right-hand-side of Figure 2A. We can see from these results that 96
the kernel matrices are highly correlated – i.e., the input stimuli are projected to very 97
similar geometries in the two representational spaces. 98

If one did not know the input transformations and simply observed the correlation 99
between kernel matrices, it would be tempting to infer that the two systems Φ_1 and Φ_2 100
transform an unknown input stimulus \mathbf{x} through a similar set of functions – for example 101
functions that belong to the same class or project inputs to similar representational spaces. 102
However, this would be an error. The projections $\Phi_1(\mathbf{x})$ and $\Phi_2(\mathbf{x})$ are fundamentally 103
different – Φ_1 (radial basis kernel) projects an input vector into an infinite dimensional 104
space, while Φ_2 (cosine kernel) projects it onto a unit sphere. The difference between these 105
functions becomes apparent if one considers how this correlation changes if one considers 106
a different set of input stimuli. For example, the set of data points shown at the left of 107
Figure 2B, are projected to very different geometries, leading to a low correlation between 108
the two kernel matrices (right-hand side). 109

In fact, the reason for highly correlated kernel matrices in Figure 2A is not a similarity 110
in the transformations Φ_1 and Φ_2 but the structure of the dataset. The representational 111
distance between any two points \mathbf{x}_i and \mathbf{x}_j in Φ_1 is a function of their Euclidean distance 112
 $\|\mathbf{x}_i - \mathbf{x}_j\|$, while in Φ_2 , it is a function of their cosine distance, $\mathbf{x}_i^T \mathbf{x}_j$. These two features 113
– Euclidean distance and cosine distance – mimic each other for certain datasets. In the 114
dataset in Figure 2A, the stimuli is clustered such that the Euclidean distance between 115

any two stimuli is correlated with their cosine distance. However, for the dataset in 116
Figure 2B, the Euclidean distance is no longer correlated with the angle and the confounds 117
lead to different representational geometries. Thus, this example illustrates how: (i) two 118
systems acting on very different features of inputs can nevertheless end up with similar 119
representational geometries when these features are able to mimic each other, and (ii) 120
when the two systems are non-identical, the correlation in representational geometries 121
will be modulated by the structure of the data – two systems may show a high correlation 122
in their representational geometries on one set but a low correlation on another set. 123

**Study 2: Complex systems encoding different features of inputs can show a 124
high RSA-score** Study 1 made a number of simplifying assumptions – the dataset was 125
two-dimensional, clustered into two categories and we intentionally chose functions Φ_1 126
and Φ_2 such that the kernel matrices were correlated in one case and not correlated in 127
the other. It could be argued that, even though the above results hold in principle, they 128
are unlikely in practice when the transformations and data structure are more complex. 129
Indeed, it is possible that a similarity in representational geometries becomes less likely 130
as one increases the number of categories (i.e., clusters or conditions) being considered. 131

To address this objection, we now consider a more complex setup, where the transfor- 132
mations Φ_1 and Φ_2 are modelled as feedforward deep neural networks (DNNs), trained to 133
classify a high-dimensional dataset into multiple categories. Many studies that use RSA 134
compare systems using naturalistic images as visual inputs [5, 11]. While using naturalis- 135
tic images brings research closer to the real-world, it is also well-known that datasets of 136
naturalistic images frequently contain confounds – independent features that can predict 137
image categories [14]. We will now show how the simplest of such confounds, a single pixel, 138
can lead to a high RSA between two DNNs that encode qualitatively different features of 139

inputs.

140

Consider the same setup as above, where an input stimulus, \mathbf{x} , is transformed to a representation space by two systems, Φ_1 and Φ_2 . Instead of a two-dimensional input space, \mathbf{x} now exists in a high-dimensional image space and Φ_1 and Φ_2 are two versions of a DNN – VGG-16 – trained to classify input images into different categories. We ensured that Φ_1 and Φ_2 were qualitatively different transformations of input stimuli by making the networks sensitive to different predictive features within the stimuli. The first network was trained on an unperturbed dataset, while the second network was trained on a modified version of the dataset, where each image was modified to contain a confound – a single pixel in a location that was diagnostic of the category (see Figure 3 for the general approach).

141

142

143

144

145

146

147

148

149

The locations of these diagnostic pixels were chosen such that they were correlated to the corresponding representational distances between classes in Φ_1 . Our hypothesis was that if the representational distances in Φ_2 preserve the physical distances of diagnostic pixels in input space, then this confound will end up mimicking the representational geometry of Φ_1 , even though the two systems use qualitatively different features for classification. Furthermore, we trained two more networks, Φ_3 and Φ_4 , which were identical to Φ_2 , except these networks were trained on datasets where the location of the confound was uncorrelated (Φ_3) or negatively correlated (Φ_4) with the representational distances in Φ_1 (see Figure 4 and Methods for details).

150

151

152

153

154

155

156

157

158

Classification accuracy (Figure 5 (left)) revealed that the network Φ_1 , trained on the unperturbed images, learned to classify these images and ignored the diagnostic pixel – that is, it’s performance was identical for the unperturbed and modified images. In contrast, networks Φ_2 (positive), Φ_3 (uncorrelated) and Φ_4 (negative) failed to classify the unperturbed images (performance was statistically at chance) but learned to perfectly classify the modified images, showing that these networks develop qualitatively different

159

160

161

162

163

164

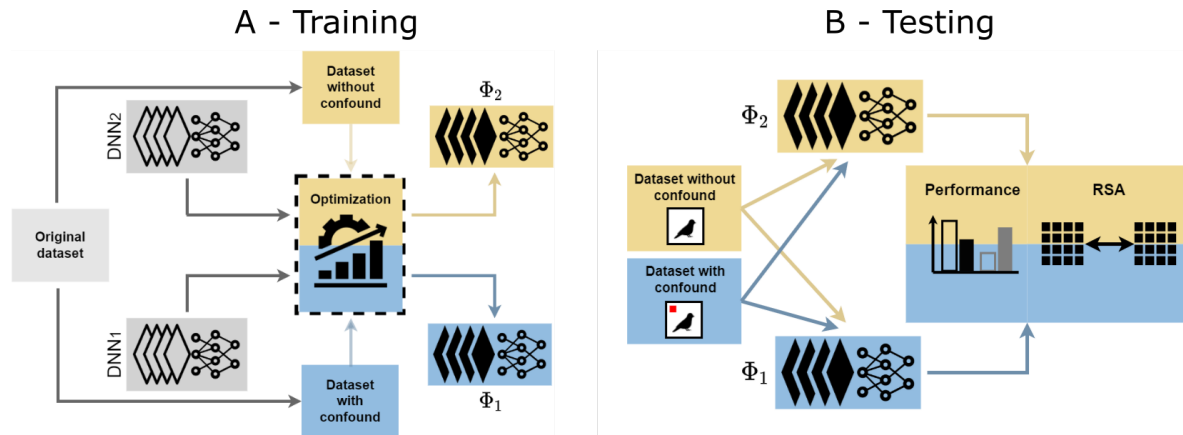


Figure 3: **Training and testing DNNs with different feature encodings.** Panel A shows the training procedure for Studies 2–4, where we created two versions of the original dataset (gray), one containing a confound (blue) and the other left unperturbed (yellow). These two datasets were used to train two networks (gray) on a categorisation task, resulting in two networks that learn to categorise images either based on the confound (projection Φ_2) or based on statistical properties of the unperturbed image (projection Φ_1). Panel B shows the testing procedure where each network was tested on stimuli from each dataset – leading to a 2x2 design. Performance on these datasets was used to infer the features that each network encoded and their internal response patterns were used to calculate RSA-scores between the two networks.

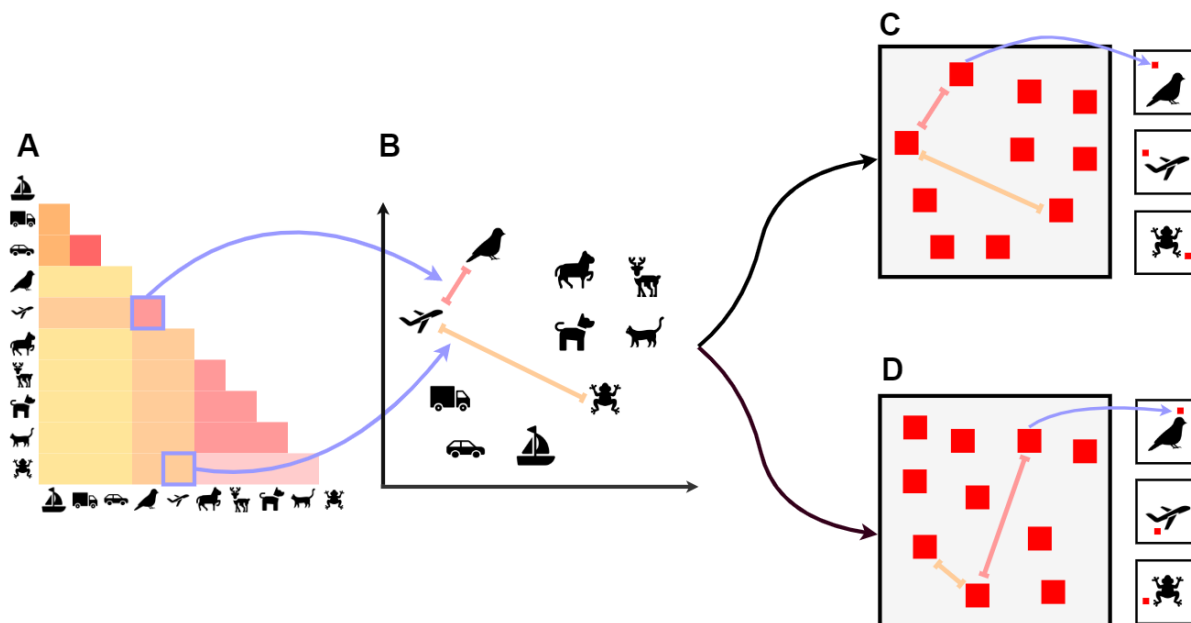


Figure 4: **Study 2 confound placement.** The representational geometry (Panel A and B) from the network trained on the unperturbed CIFAR-10 images is used to determine the location of the single pixel confound (shown as a red patch here) for each category. In the ‘Positive’ condition (Panel C), we determined 10 locations in a 2D plane such that the distances between these locations were positively correlated to the representational geometry – illustrated here as the red patches in Panel C being in similar locations to category locations in Panel B. These 10 locations were then used to insert a single diagnostic – i.e., category-dependent – pixel in each image (Insets in Panel C). A similar procedure was also used to generate datasets where the confound was uncorrelated (Panel D) or negatively correlated (not shown here) with the representational geometry of the network.

representations compared to normally trained networks. 165

Next we computed pairwise RSA scores between the representations at the last con- 166
volution layer of Φ_1 and each of Φ_2 , Φ_3 and Φ_4 (Figure 5 (right)). When presented un- 167
perturbed test images, the Φ_2 , Φ_3 and Φ_4 networks all showed low RSA scores with the 168
normally trained Φ_1 network. However, when networks were presented with test images 169
that included the predictive pixels, RSA varied depending on the geometry of pixel loca- 170
tions in the input space. When the geometry of pixel locations was positively correlated 171
to the normally trained network, RSA scores approached ceiling (i.e., comparable to RSA 172
scores between two normally trained networks). Networks trained on uncorrelated and 173
negatively correlated pixel placements scored much lower. 174

These results mirror Study 1: we observed that it is possible for two networks (Φ_1 and 175
 Φ_2) to show highly correlated representational geometries even though these networks 176
learn to classify images based on very different features. One may argue that this could 177
be because the two networks could have learned similar representations at the final con- 178
volution layer of the DNN and it is the classifier that sits on top of this representation 179
that leads to the behavioural differences between these networks. But if this was true, it 180
would not explain why RSA scores diminish for the two other comparisons (with Φ_3 and 181
 Φ_4). This modulation of RSA-scores for different datasets suggests that, like in Study 1, 182
the correlation in representational geometry is not because the two systems encode similar 183
features of inputs, but because different features mimic each other in their representational 184
geometries. 185

Re-examining some influential findings 186

In Studies 1 and 2, we showed that it is possible for qualitatively different systems to 187
end up with similar representational geometries. However, it may be argued that while 188

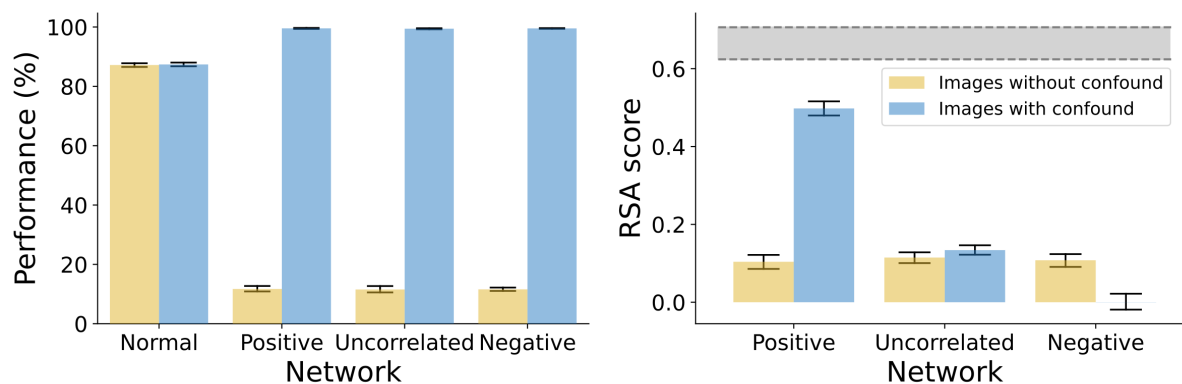


Figure 5: **Study 2 results.** *Left:* Performance of normally trained networks did not depend on whether classification was done on unperturbed CIFAR-10 images or images with a single pixel confound (error bars represent 95% CI). All three networks trained on datasets with confounds could perfectly categorise the test images when they contained the confound (blue bars), but failed to achieve above-chance performance if the predictive pixel was not present (yellow bars). *Right:* The RSA score between the network trained on the unperturbed dataset and each of the networks trained on datasets with confounds. The three networks showed similar scores when tested on images without confounds, but vastly different RSA scores when tested on images with confounds. Networks in the Positive condition showed near ceiling scores (the shaded area represents noise ceiling) while networks in the Uncorrelated and Negative conditions showed much lower RSA.

this is possible in principle, it is unlikely in practice in real-world scenarios. In the fol- 189
lowing two studies, we consider real-world data from some recent influential experiments, 190
recorded from both primate and human cortex. We show how RSA-scores can be driven 191
by confounds in these real-world settings and how properties of training and test data 192
may contribute to observed RSA-scores. 193

Study 3: Neural activations in monkey IT cortex can show a high RSA-score 194
with DNNs despite different encoding of input data In our next study, we con- 195
sider data from experiments comparing representational geometries between computa- 196
tional models and macaque visual cortex [11, 27]. The experimental setup was similar 197
to Study 2, though note that unlike Study 2, where both systems used the same archi- 198
tecture and learning algorithm, this study considered two very different systems – one 199
artificial (DNN) and the other biological (macaque IT cortex). We used the same set of 200
images that were shown to macaques by Majaaj et al. [28] and modified this dataset to 201
superimpose a small diagnostic patch on each image. In the same manner as in Study 2 202
above, we constructed three different datasets, where the locations of these diagnostic 203
patches were either positively correlated, uncorrelated or negatively correlated with the 204
RDM of macaque activations. We then trained four CNNs. The first CNN was pre- 205
trained on ImageNet and then fine-tuned on the unmodified dataset of images shown to 206
the macaques. Previous research has shown that CNNs trained in this manner develop 207
representations that mirror the representational geometry of neurons in primate inferior 208
temporal (IT) cortex [11]. The other three networks were trained on the three modi- 209
fied datasets and learned to entirely rely on the diagnostic patches (accuracy on images 210
without the diagnostic patches was around chance). 211

Figure 6 (right) shows the correlation in representational geometry between the macaque 212

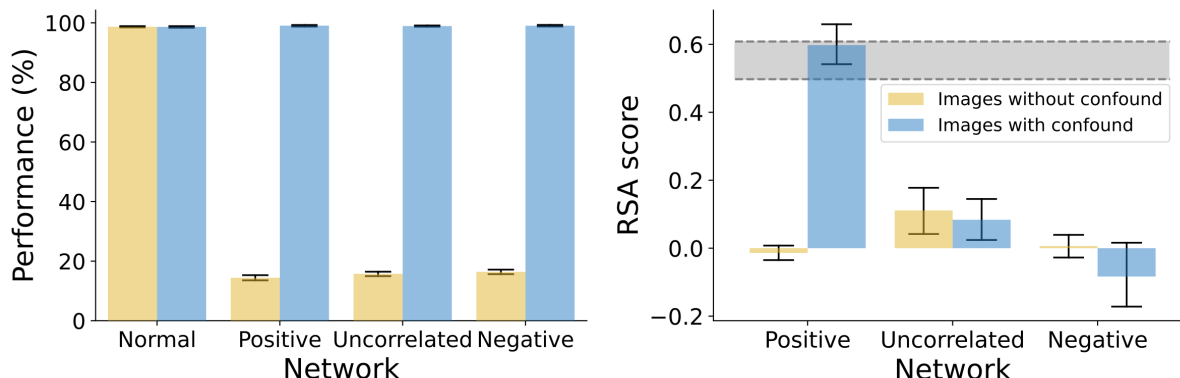


Figure 6: **Study 3 results.** *Left:* Classification Performance of the network trained on unperturbed images (Normal condition) did not depend on the presence or absence of the confound, while performance of networks trained with the confound (Positive, Uncorrelated and Negative conditions) highly depended on whether the confound was present. *Right:* RSA-scores with macaque IT activations were low for all three conditions when images did not contain a confound (yellow bars). When images contained a confound (blue bars), the RSA-scores depended on the condition, matching the RSA-score of the normally trained network (grey band) in the Positive condition, but decreasing significantly in the Uncorrelated and Negative conditions. The grey band represents a 95% CI for the RSA-score between normally trained networks and macaque IT activations.

IT activations and activations at the final convolution layer for each of these networks. 213
The correlation with networks trained on the unmodified images is our baseline and shown 214
as the gray band in Figure 6. Our first observation was that a CNN trained to rely on 215
the diagnostic patch can indeed achieve a high RSA score with macaque IT activations. 216
In fact, the networks trained on patch locations that were positively correlated to the 217
macaque RDM matched the RSA score of the CNNs trained on ImageNet and the unmod- 218
ified dataset. This shows how two systems having very different architectures, encoding 219

fundamentally different features of inputs (single patch vs naturalistic features) can show 220
a high correspondence in their representational geometries. We also observed that, like 221
in Study 2, the RSA score depended on the clustering of data in the input space – when 222
patches were placed in other locations (uncorrelated or negatively correlated to macaque 223
RDMs) the RSA score became significantly lower. 224

Study 4: High RSA-scores may be driven by the structure of testing data All 225
the studies so far have used the same method to construct datasets with confounds – we 226
established the representational geometry of one system (Φ_1) and constructed datasets 227
where the clustering of features (pixels) mirrored this geometry. However, it could be 228
argued that confounds which cluster in this manner are unlikely in practice. For example, 229
even if texture and shape exist as confounds in a dataset, the inter-category distances 230
between textures are not necessarily similar to the inter-category distances between shape. 231

However, categories in real-world datasets are usually hierarchically clustered into 232
higher-level and lower-level categories. For example, in the CIFAR-10 dataset, the Dogs 233
and Cats (lower-level categories) are both animate (members of a common higher-level 234
category) and Airplanes and Ships (lower-level categories) are both inanimate (members 235
of a higher-level category). Due to this hierarchical structure, Dog and Cat images are 236
likely to be closer to each other not only in their shape, but also their colour and texture 237
(amongst other features) than they are to Airplane and Ship images. In our next simula- 238
tion, we explore whether this hierarchical structure of categories can lead to a correlation 239
in representational geometries between two systems that learn different feature encodings. 240

For this study, we selected a popular dataset used for comparing representational 241
geometries in humans, macaques and deep learning models [12, 29]. This dataset consists 242
of six categories which can be organised into a hierarchical structure shown in Figure 7. [5] 243

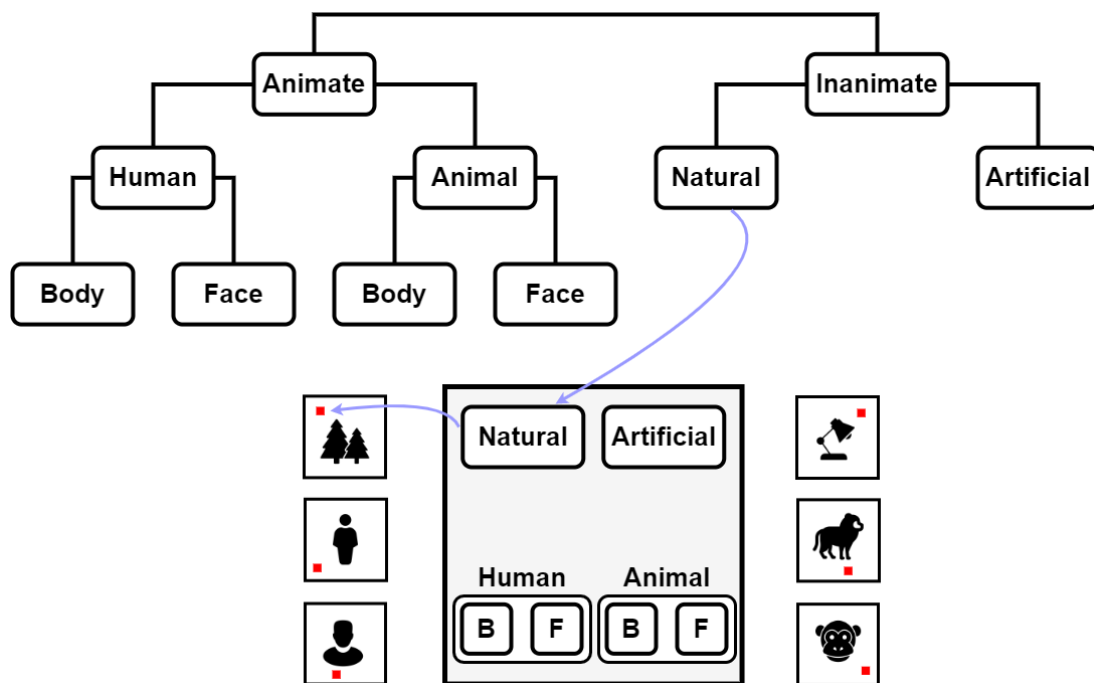


Figure 7: **Exploiting intrinsic dataset hierarchy in order to place confounds.**

The top panel shows the hierarchical structure of categories in the dataset, which was used to place the single pixel confounds. The example at the bottom (middle) shows one such hierarchical placement scheme where the pixels for Inanimate images were closer to the top of the canvas while Animate images were closer to the bottom. Within the Animate images, the pixels for Humans and Animals were placed at the left and right, respectively, and the pixels for bodies (B) and faces (F) were clustered as shown.

showed a striking match in RDMs for response patterns elicited by these stimuli in human 244
and macaque IT. For both humans and macaques, distances in response patterns were 245
larger between the higher-level categories (animate and inanimate) than between the 246
lower-level categories (e.g., between human bodies and human faces). 247

We used a similar experimental paradigm to the above studies, where we trained 248
networks to classify stimuli which included a single predictive pixel. But instead of using 249

an RDM to compute the location of a diagnostic pixel, we used the hierarchical categorical 250
structure. In the first modified version of the dataset, the location of the pixel was based 251
on the hierarchical structure of categories in Figure 7 – predictive pixels for animate 252
kinds were closer to each other than to inanimate kinds, and pixels for faces were closer 253
to each other than to bodies, etc. One such configuration can be seen in Figure 7. In the 254
second version, the predictive pixel was placed at a random location for each category 255
(but, of course, at the same location for all images within each category). We call these 256
conditions ‘Hierarchical’ and ‘Random’. [12] showed that the RDM of average response 257
patterns elicited in the human IT cortex (Φ_1) correlated with the RDM of a DNN trained 258
on naturalistic images (Φ_2). We explored how this compared to the correlation with the 259
RDM of a network trained on the Hierarchical pixel placement (Φ_3) and Random pixel 260
placement (Φ_4). 261

Results for this study are shown in Figure 8. We observed that representational ge- 262
ometry of a network trained on Hierarchically placed pixels (Φ_3) was just as correlated to 263
the representational geometry of human IT responses (Φ_1) as a network trained on natu- 264
ralistic images (Φ_2). However, when the pixel locations for each category were randomly 265
chosen, this correlation decreased significantly. These results suggest that any confound in 266
the dataset (including texture, colour or low-level visual information) that has distances 267
governed by the hierarchical clustering structure of the data could underlie the observed 268
similarity in representational geometries between CNNs and human IT. More generally, 269
these results show how it is plausible that many confounds present in popular datasets 270
may underlie the observed similarity in representational geometries between two systems. 271
The error of inferring a similarity in mechanism based on a high RSA score is not just 272
possible but also probable. 273

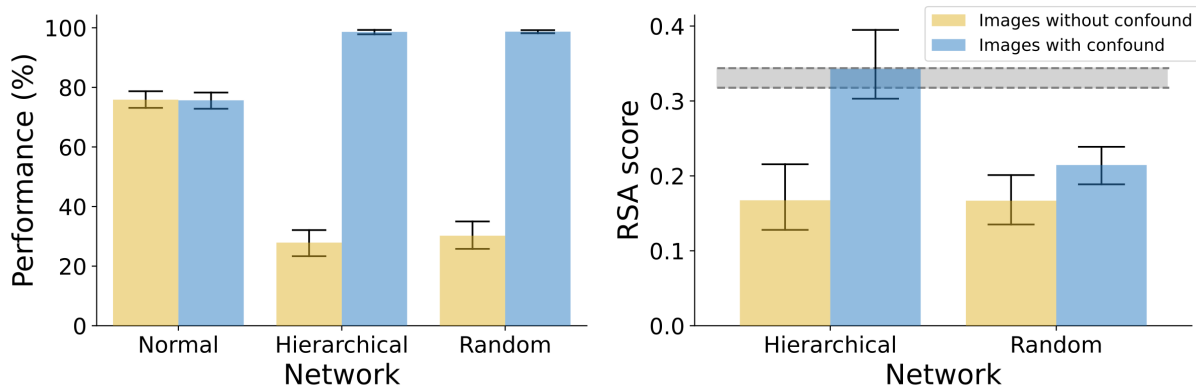


Figure 8: **Study 4 results.** *Left:* Performance of normally trained networks did not depend on whether the confound was present. Networks trained with the confound failed to classify stimuli without the confound (yellow bars) while achieving near perfect classification of stimuli with the confound present (blue bars). *Right:* RSA with human IT activations reveals that, when the confound was present, the RSA-score for networks in the Hierarchical condition matched the RSA-score of normally trained network (gray band), while the RSA-score of the network in the Random condition was significantly lower. The grey band represents 95% CI for the RSA score between normally trained networks and human IT.

Discussion

274

In four studies, we have illustrated a number of conditions under which it can be prob- 275
lematic to infer a similarity of representations between two systems based on a correlation 276
in their representational geometries. We showed that two systems may transform their 277
inputs through very different functions and encode very different features of inputs and 278
yet have highly correlated representational geometries. In fact, we showed that this corre- 279
lation can be a product of the structure of the dataset itself. A consequence of this result 280
is that the RSA-score between two systems becomes dataset dependent. For example, 281

one may observe a high RSA-score between a brain region of a primate and human for 282
one dataset (e.g., [5]), but this score may become much lower for another dataset. Thus 283
the observation of a similarity in representational geometry between systems must be 284
interpreted with caution. 285

The significance of these results depends on whether you take an *externalist* or *holistic* 286
view on mental representations. According to the first view, the content of representations 287
is determined by their relationship to entities in the external world. This perspective is 288
implicitly taken by most neuroscientists and psychologists, who are interested in compar- 289
ing mechanisms underlying cognitive processes – that is, they are interested in the set 290
of nested functions and algorithms responsible for transforming sensory input into a set 291
of activations in the brain. From this perspective, our finding that high RSAs can be 292
obtained between systems that work in qualitatively different ways poses a challenge to 293
researchers using RSAs to compare systems. 294

Of course, a researcher with an externalist perspective may acknowledge that a second- 295
order isomorphism of activity patterns does *not* strictly imply that two systems are similar 296
mechanistically but still assume that it is highly likely to be the case. That is, as a practi- 297
cal matter, a researcher may assume that RSAs are a reliable method to compare systems. 298
However, our findings challenge this assumption. We show how a high RSA between dif- 299
ferent systems can not only occur in principle, but also in practice, in high-dimensional 300
systems operating on high-dimensional data. Indeed, we show that the hierarchical struc- 301
ture of datasets frequently used to test similarity of representations lends itself to a high 302
RSA arising because of confounds present in the dataset. Such confounds are commonly 303
found in high-dimensional stimuli such as naturalistic images that are frequently used to 304
measure RSA [11, 29]. Indeed, presence of such confounds may explain why researchers 305
have observed high RSAs between DNNs that classify objects based on texture [16, 17] 306

and the human visual system that classifies by shape [22, 30].

Alternatively, a researcher may reject an externalist view and adopt the perspective that representations obtain their meaning based on how they are related to each other within each system, rather than based on their relationship to entities in the external world. That is, “representation *is* the representation of similarities” [31]. From this perspective, as long as the two systems share the same relational distances between internal activations, one can validly infer that the two systems have similar representations. That is, a second-order isomorphism implies a similarity of representations, by definition. This view has been called *holism* in the philosophy of mind [32, 33] and is related to a similar idea of *meaning holism* in language, which is the idea that the meaning of a linguistic expression is determined by its relation to other expressions within a language [34, 35]. For example, Firth [36] (p. 11) writes: “you shall know a word by the company it keeps”. More recently, Griffiths and Steyvers [37], and Griffiths, Steyvers, and Tenenbaum [38] have adopted meaning holism accounts of semantic representations in neural networks. Our results are *not* problematic for a researcher adopting this holistic perspective. However, our results show that adopting this view misses the information about differences in mechanistic processes that a psychologist or neuroscientist is frequently interested in, for instance, whether the visual system processes shape or texture (or the location of diagnostic pixels) in order to identify objects. Fodor and Lepore long ago criticized this philosophical stance [33, 39], and interestingly, this philosophical debate played an important part in the development of RSA (see Supplementary Information, Section A). Unfortunately, this debate has largely been ignored by researchers who use RSA as a method to compare similarity of systems.

We would also like to make it clear that the results here are not a blanket criticism of the RSA approach as currently practiced. A representational dissimilarity matrix (RDM)

contains important information about the similarity structures of representations. Any 332
mechanistically correct model of an individual or a species must capture this similarity 333
structure. As such, RSA provides a benchmark for *rejecting* possible models. However, 334
the above studies show that RSA may be a misleading benchmark for *selecting* models – 335
two systems may show similar representational geometries and yet work on very different 336
transformations and features of input stimuli (for an in depth discussion about inferring 337
similarity of causal mechanisms from similar outcomes see [40]). 338

A related point has been made by Kriegeskorte and Diedrichson [41] and Kriegesko- 339
rte and Wei [42], who point out that two systems may have the same representational 340
geometry, even if they have a different activity profile over neurons. In this sense, the ge- 341
ometry loses the information about how information was distributed over a set of neurons. 342
Kriegeskorte and Diedrichson [41] equate this loss in information to “peeling a layer of an 343
onion” – downstream decoders that are sensitive to the representational geometry rather 344
than activity profiles over neuron populations can focus on difference in information as 345
reflected by a change in geometry and be agnostic to how this information is distributed 346
over a set of neurons. We agree that this invariance over activity profiles is indeed a 347
useful property of representational geometries for downstream decoders. However, we are 348
not aware of any studies that highlight how representational geometries also abstract over 349
behaviourally relevant stimulus properties (e.g. shape vs texture). While abstracting over 350
activity profiles may be useful, abstracting over stimulus properties loses an important 351
piece of information when comparing representations across brain regions, individuals, 352
species and between brains and computational models. Our studies how two systems 353
may appear similar based on their representational geometries in one circumstance (e.g. 354
Figure 2A) but drastically different in another circumstance (Figure 2B). 355

The key implication of our findings is that researchers should assess RSAs on a wider 356

variety of datasets when comparing systems. Two systems that have the similar representations should show a high RSA irrespective of the stimuli on which they are tested, and testing systems on multiple datasets will reduce the likelihood that confounds or other factors are driving the effects. In practice, observing high RSAs after testing very different datasets, and datasets manipulated to avoid possible confounds, should be required before drawing strong conclusions regarding the similarity of two systems. In this regard, the “controversial stimuli” – images on which different computational models produce distinct responses – developed by Golan & Kriegeskorte [43] is a step in the right direction. By testing on stimuli that produce distinct responses in different models, one can adjudicate between models by comparing their representational geometries to the representational geometry of a target system. Combining RSA results with a range of methods, including experimental studies that stringently test hypotheses about how different systems work, seems the best approach going forward.

Methods

Dataset generation and training

All DNN simulations (Studies 2–4) were carried out using the Pytorch framework [44]. The model implementations were downloaded from the torchvision library. Networks trained on unperturbed datasets in all studies were pre-trained on ImageNet as were networks trained on modified datasets in Study 2. Networks trained on modified datasets in Studies 3 and 4 were randomly initialised. For the pre-trained models, their pre-trained weights were downloaded from torchvision.models subpackage.

Study 1 Each dataset in Study 1 consists of 100 samples (50 in each cluster) drawn 378
from two multivariate Gaussians, $\mathcal{N}(x|\mu, \Sigma)$, where μ is a 2-dimensional vector and Σ is 379
a 2×2 covariance matrix. In Figure 2A, the two Gaussians have means $\mu_1 = (1, 8)$ and 380
 $\mu_2 = (8, 1)$ and a covariance matrices $\Sigma_1 = \Sigma_2 = \frac{1}{2}\mathbf{I}$, while in Figure 2B the Gaussians 381
have means $\mu_1 = (1, 1)$ and $\mu_2 = (8, 8)$ and a covariance matrices $\Sigma_1 = \mathbf{I}$, $\Sigma_2 = 8\mathbf{I}$. 382
All kernel matrices were computed using the `sklearn.metrics.pairwise` module of the 383
`scikit-learn` Python package. 384

Study 2 First, a VGG-16 deep convolutional neural network [45], pre-trained on the 385
ImageNet dataset of naturalistic images, was trained to classify stimuli from the CIFAR-10 386
dataset [46]. The CIFAR-10 dataset includes 10 categories with 5000 training, and 1000 387
test images per category. The network was fine-tuned on CIFAR-10 by replacing the 388
classifier so that the final fully-connected layer reflected the correct number of target 389
classes in CIFAR-10 (10 for CIFAR-10 as opposed to 1000 for ImageNet). Images were 390
rescaled to a size of 224×224 px and then the model learnt to minimise the cross-entropy 391
error using the RMSprop optimizer with a mini-batch size of 64, learning rate of 10^{-5} , 392
and momentum of 0.9. All models were trained for 10 epochs, which were sufficient for 393
convergence across all datasets. 394

Second, 100 random images from the test set for each category were sampled as in- 395
put for the network and activations at the final convolutional layer extracted using the 396
THINGSVision Python toolkit [47]. The same toolkit was used to generate a representa- 397
tional dissimilarity matrix (RDM) from the pattern of activations using 1-Pearson's r 398
as the distance metric. The RDM was then averaged by calculating the median distance 399
between each instance of one category with each instance of the others (e.g., the median 400
distance between `Airplane` and `Ship` was the median of all pair-wise distances between 401

activity patterns for airplane and ship stimuli). This resulted in a 10×10 , category-level, 402
RDM which reflected average between-category distances. 403

Third, three modified versions of the CIFAR-10 datasets were created for the ‘Positive’, 404
‘Uncorrelated’ and ‘Negative’ conditions, respectively. In each dataset, we added one 405
diagnostic pixel to each image, where the location of the pixel depended on the category 406
(See Figure 4). The locations of these pixels were determined using the averaged RDM 407
from the previous step. We call this the target RDM. In the ‘Positive’ condition, we 408
wanted the distances between pixel placements to be positively correlated to the distances 409
between categories in the target RDM. We achieved this by using an iterative algorithm 410
that sampled pixel placements at random, calculated an RDM based on distances between 411
the pixel placements and computed an RSA-score (Spearman correlation) with the target 412
RDM. Placements with a score above 0.70 were retained and further optimized (using 413
small perturbations) to achieve an RSA-score over 0.90. The same procedure was also 414
used to determine placements in the Uncorrelated (optimizing for a score close to 0) and 415
Negatively correlated (optimizing for a negative score) conditions. 416

Finally, datasets were created using 10 different placements in each of the three condi- 417
tions. Networks were trained for classification on these modified CIFAR-10 datasets in the 418
same manner as the VGG-16 network trained on the unperturbed version of the dataset 419
(See Figure 3). 420

Study 3 The procedure mirrored Study 2 with the main difference being that the target 421
system was the macaque inferior temporal cortex. Neural data from two macaques, as well 422
as the dataset were obtained from the Brain Score repository [27]. This dataset consists 423
of 3200 images from 8 categories (animals, boats, cars, chairs, faces, fruits, planes, and 424
tables), we computed an 8×8 averaged RDM based on macaque IT response patterns for 425

stimuli in each category. 426

This averaged RDM was then used as the target RDM in the optimization procedure to 427
determine locations of the confound (here, a white predictive patch of size 5×5 pixels) for 428
each category. Using a patch instead of a single pixel was required in this dataset because 429
of the structure and smaller size of the dataset (3200 images, rather than 50,000 images 430
for CIFAR-10). In this smaller dataset, the networks struggle to learn based on a single 431
pixel. However, increasing the size of the patch makes these patches more predictive 432
and the networks are able to again learn entirely based on this confound (see results 433
in Figure 5). In a manner similar to Study 2, this optimisation procedure was used 434
to construct three datasets, where the confound's placement was positively correlated, 435
uncorrelated or negatively correlated with the category distances in the target RDM. 436

Finally, each dataset was split into 75% training (2432 images) and 25% test sets 437
(768 images) before VGG-16 networks were trained on the unperturbed and modified 438
datasets in the same manner as in Study 2. One difference between Studies 2 and 3 439
was that here the networks in the Positive, Uncorrelated and Negative conditions were 440
trained from scratch, i.e., not pre-trained on ImageNet. This was done because we wanted 441
to make sure that the network in the Normal condition (trained on ImageNet) and the 442
networks in the Positive, Uncorrelated and Negative conditions encoded fundamentally 443
different features of their inputs – i.e., there were no ImageNet-related features encoded by 444
representations Φ_2 , Φ_3 and Φ_4 that were responsible for the similarity in representational 445
geometries between these representations and the representations in macaque IT cortex. 446

Study 4 The target system in this study was human IT cortex. The human RDM 447
and dataset were obtained from [5]. Rather than calculating pixel placements based on 448
the human RDM, the hierarchical structure of the dataset was used to place the pixels 449

manually. The dataset consists of 910 images from 6 categories: human bodies, human 450
faces, animal bodies, animal faces, artificial inanimate objects and natural inanimate 451
objects. These low-level categories can be organised into the hierarchical structure shown 452
in Figure 7. Predictive pixels were manually placed so that the distance between pixels 453
for Animate kinds were closer together than they were to Inanimate kinds and that faces 454
were closer together than bodies. This can be done in many different ways, so we created 455
five different datasets, with five possible arrangements of predictive pixels. Results in 456
the Hierarchical condition (Figure 8) are averaged over these five datasets. Placements 457
for the Random condition were done similarly, except that the locations were selected 458
randomly. 459

Networks were then trained on a 6-way classification task (818 training images and 92 460
test images) in a similar manner to the previous studies. As in Study 3, networks trained 461
on the modified datasets (both Hierarchical and Random conditions) were not pre-trained 462
on ImageNet. 463

RDM and RSA computation 464

For Studies 2-4 all image-level RDMs were calculated using $1 - r$ as the distance measure. 465
RSA scores were computed as the Spearman rank correlation between RDMs. 466

In Study 2, a curated set of test images was selected due to the extreme heterogeneity 467
of the CIFAR-10 dataset (low activation pattern similarity between instances of the same 468
category). This was done by selecting 5 images per category which maximally correlated 469
with the averaged activation pattern for the category. Since CIFAR-10 consists of 10 470
categories, the RSA-scores in Study 2 were computed using RDMs of size 50×50 . 471

In Study 3, the dataset consisted of 3200 images belonging to 8 categories. We first 472
calculated a full 3200×3200 RDM using the entire set of stimuli. An averaged, category- 473

level, 8×8 RDM was then calculated using median distances between categories (in 474
a manner similar to that described for Study 2 in the Section ‘Dataset generation and 475
training’). This 8×8 RDM was used to determine the RSA-scores. We also obtained 476
qualitatively similar results using the full 3200×3200 RDMs. These results can be found 477
in the Supplementary Information, Section B. 478

In Study 4, the dataset consisted of 818 training images and 92 test images. Kriegesko- 479
rte et al. [5] used these images to obtain a 92×92 RDM to compare representations between 480
human and macaque IT cortex. Here we computed a similar 92×92 RDM for networks 481
trained in the Normal, Hierarchical and Random training conditions, which were then 482
compared with the 92×92 RDM from human IT cortex to obtain RSA-scores for each 483
condition. 484

Testing 485

In Study 2, we used a 4×2 design to measure classification performance for networks in 486
all four conditions (Normal, Positive, Uncorrelated and Negative) on both unperturbed 487
images and modified images. We computed six RSA-scores: three pairs of networks – 488
Normal-Positive, Normal-Uncorrelated and Normal-Negative – and two types of inputs – 489
unperturbed and modified test images. The noise ceiling (grey band in Figure 5) was de- 490
termined in the standard way as described in [48] and represents the expected range of the 491
highest possible RSA score with the target system (network trained on the unperturbed 492
dataset). 493

In Study 3, performance was estimated in the same manner as in Study 2 (using a 494
 4×2 design), but RSA-scores were computed between RDMs from macaque IT activations 495
and the four types of networks – i.e. for the pairs Macaque-Normal, Macaque-Positive, 496
Macaque-Uncorrelated and Macaque-Negative. And like in Study 2, we determined each 497

of these RSA-scores for both unperturbed and modified test images as inputs to the 498
networks. 499

In Study 4, performance and RSA were computed in the same manner as in Study 3, 500
except that the target RDM for RSA computation came from activations in human IT 501
cortex and the networks were trained in one of three conditions: Normal, Hierarchical 502
and Random. 503

Data analysis 504

Performance and RSA scores were compared by running analyses of variance and Tukey 505
HSD post-hoc tests. In Study 2 and 3, performance differences were tested by running a 506
4 (type of training) by 2 (type of dataset) mixed ANOVAs. In, Study 4, the differences 507
were tested by running a 3x2 mixed ANOVA. 508

RSA scores with the target system between networks in various conditions were com- 509
pared by running 3x2 ANOVAs in Studies 2 and 3, and a 2x2 ANOVA in Study 4. We 510
observed that RSA-scores were highly dependent on both the way the networks were 511
trained and also the test images used to elicit response activations. 512

For a detailed overview of the statistical analyses and results, see Supplemental Informa- 513
tion Section C. 514

Data Availability 515

Confound placement coordinates (Studies 2-4), unperturbed datasets (Studies 3 and 4), 516
macaque activation patterns and RDMs (Study 3) and human RDM (Study 4) are avail- 517
able at [OSF](#). 518

Acknowledgments

519

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No 741134)

References

523

- [1] Hubel, D. H. & Wiesel, T. N. Receptive fields of single neurones in the cat’s striate cortex. *The Journal of Physiology* **148**, 574–591 (1959). 524
525
- [2] O’Keefe, J. Place units in the hippocampus of the freely moving rat. *Experimental Neurology* **51**, 78–109 (1976). 526
527
- [3] Hafting, T., Fyhn, M., Molden, S., Moser, M.-B. & Moser, E. I. Microstructure of a spatial map in the entorhinal cortex. *Nature* **436**, 801–806 (2005). 528
529
- [4] Kriegeskorte, N., Mur, M. & Bandettini, P. Representational similarity analysis - connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience* **2** (2008). 530
532
- [5] Kriegeskorte, N. *et al.* Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron* **60**, 1126–1141 (2008). 533
534
- [6] Haxby, J. V., Gobbini, M. I. & Nastase, S. A. Naturalistic stimuli reveal a dominant role for agentic action in visual representation. *NeuroImage* **216**, 116561 (2020). 535
536
- [7] O’Hearn, K., Larsen, B., Fedor, J., Luna, B. & Lynn, A. Representational similarity analysis reveals atypical age-related changes in brain regions supporting face and car recognition in autism. *NeuroImage* **209**, 116322 (2020). 537
538
539

- [8] Michael L. Mack, B. L., Alison R. Preston. Decoding the brain’s algorithm for cate- 540
gorization from its neural implementation. *Current Biology* **23**, 2023–2027 (2013). 541
- [9] Freund, M. C., Etzel, J. A. & Braver, T. S. Neural coding of cognitive control: 542
The representational similarity analysis approach. *Trends in Cognitive Sciences* **25**, 543
622–638 (2021). 544
- [10] Kaneshiro, B., Perreau Guimaraes, M., Kim, H.-S., Norcia, A. M. & Suppes, P. A 545
representational similarity analysis of the dynamics of object processing using single- 546
trial eeg classification. *PLOS ONE* **10**, 1–27 (2015). 547
- [11] Yamins, D. L. K. *et al.* Performance-optimized hierarchical models predict neural 548
responses in higher visual cortex. *Proceedings of the National Academy of Sciences* 549
111, 8619–8624 (2014). 550
- [12] Khaligh-Razavi, S.-M. & Kriegeskorte, N. Deep supervised, but not unsupervised, 551
models may explain it cortical representation. *PLOS Computational Biology* **10**, 1–29 552
(2014). 553
- [13] Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A. & Oliva, A. Comparison of 554
deep neural networks to spatio-temporal cortical dynamics of human visual object 555
recognition reveals hierarchical correspondence. *Scientific Reports* **6**, 27755 (2016). 556
- [14] Kietzmann, T. C. *et al.* Recurrence is required to capture the representational dy- 557
namics of the human visual system. *Proceedings of the National Academy of Sciences* 558
116, 21854–21863 (2019). 559
- [15] Kiat, J. E. *et al.* Linking patterns of infant eye movements to a neural network 560
model of the ventral stream using representational similarity analysis. *Developmental* 561
Science **25**, e13155 (2022). 562

- [16] Geirhos, R. *et al.* Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231* (2018). 563
564
- [17] Geirhos, R. *et al.* Shortcut learning in deep neural networks. *Nature Machine Intelligence* **2**, 665–673 (2020). 565
566
- [18] Hermann, K., Chen, T. & Kornblith, S. The origins and prevalence of texture bias in convolutional neural networks. *Advances in Neural Information Processing Systems* **33** (2020). 567
568
569
- [19] Malhotra, G., Evans, B. D. & Bowers, J. S. Hiding a plane with a pixel: examining shape-bias in cnns and the benefit of building in biological constraints. *Vision Research* **174**, 57–68 (2020). 570
571
572
- [20] Malhotra, G., Dujmovic, M. & Bowers, J. S. Feature blindness: a challenge for understanding and modelling visual object recognition (in press). *PLOS Computational Biology*, preprint *bioRxiv:2021.10.20.465074* (2022). 573
574
575
- [21] Navon, D. Forest before trees: The precedence of global features in visual perception. *Cognitive psychology* **9**, 353–383 (1977). 576
577
- [22] Biederman, I. & Ju, G. Surface versus edge-based determinants of visual recognition. *Cognitive psychology* **20**, 38–64 (1988). 578
579
- [23] Landau, B., Smith, L. B. & Jones, S. S. The importance of shape in early lexical learning. *Cognitive development* **3**, 299–321 (1988). 580
581
- [24] Xu, Y. & Vaziri-Pashkam, M. Limits to visual representational correspondence between convolutional neural networks and the human brain. *Nature Communications* **12**, 2065 (2021). 582
583
584

- [25] Schölkopf, B. & Smola, F., A. J. and Bach. *Learning with kernels: support vector machines, regularization, optimization, and beyond* (MIT Press, 2002). 585
586
- [26] Sahami, M. & Heilman, T. D. A web-based kernel function for measuring the similarity of short text snippets. In *Proceedings of the 15th International Conference on World Wide Web, WWW '06*, 377–386 (Association for Computing Machinery, New York, NY, USA, 2006). 587
588
589
590
- [27] Schrimpf, M. *et al.* Brain-score: Which artificial neural network for object recognition is most brain-like? *bioRxiv preprint: 407007* (2018). 591
592
- [28] Majaj, N. J., Hong, H., Solomon, E. A. & DiCarlo, J. J. Simple learned weighted sums of inferior temporal neuronal firing rates accurately predict human core object recognition performance. *Journal of Neuroscience* **35**, 13402–13418 (2015). 593
594
595
- [29] Kriegeskorte, N. Relating population-code representations between man, monkey, and computational models. *Frontiers in Neuroscience* **3**, 363–373 (2009). 596
597
- [30] Smith, L. B., Jones, S. S., Landau, B., Gershkoff-Stowe, L. & Samuelson, L. Object name learning provides on-the-job training for attention. *Psychological science* **13**, 13–19 (2002). 598
599
600
- [31] Edelman, S. Representation is representation of similarities. *Behavioral and Brain Sciences* **21**, 449–467 (1998). 601
602
- [32] Block, N. Advertisement for a semantics for psychology. *Midwest Studies in Philosophy* **10**, 615–678 (1986). 603
604
- [33] Fodor, J. & Lepore, E. *Holism: A Shoppers Guide* (Blackwell, Cambridge, 1992). 605

- [34] Hempel, C. G. Problems and changes in the empiricist criterion of meaning. *Revue Internationale de Philosophie* **4**, 41–63 (1950). 606
607
- [35] Quine, W. V. Main trends in recent philosophy: Two dogmas of empiricism. *The Philosophical Review* **60**, 20–43 (1951). 608
609
- [36] Firth, J. R. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis* (1957). 610
611
- [37] Griffiths, T. L. & Steyvers, M. A probabilistic approach to semantic representation. 612
In *Proceedings of the Twenty-Fourth Annual Conference of the Cognitive Science Society* (Erlbaum, Hillsdale, NJ, 2002). 613
614
- [38] Griffiths, T. L., Steyvers, M. & Tenenbaum, J. A probabilistic approach to semantic 615
representation. *Psychological Review* **114**, 211–244 (2007). 616
- [39] Fodor, J. *Psychosemantics: The Problem of Meaning in the Philosophy of Mind* 617
(MIT Press, Cambridge, 1987). 618
- [40] Guest, O. & Martin, A. E. On logical inference over brains, behaviour, and artificial 619
neural networks. *PsyArXiv preprint: 10.31234/osf.io/tbmcg* (2021). 620
- [41] Kriegeskorte, N. & Diedrichsen, J. Peeling the onion of brain representations. *Annual 621
Review of Neuroscience* **42**, 407–432 (2019). 622
- [42] Kriegeskorte, N. & Wei, X.-X. Neural tuning and representational geometry. *Nature 623
Reviews Neuroscience* **22**, 703–718 (2021). 624
- [43] Golan, T., Raju, P. C. & Kriegeskorte, N. Controversial stimuli: Pitting neural net- 625
works against each other as models of human cognition. *Proceedings of the National 626
Academy of Sciences* **117**, 29330–29337 (2020). 627

- [44] Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning li- 628
brary. In *Advances in Neural Information Processing Systems 32*, 8024–8035 (Curran 629
Associates, Inc., 2019). 630
- [45] Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale 631
image recognition. *arXiv preprint arXiv:1409.1556* (2014). 632
- [46] Krizhevsky, A. Learning multiple layers of features from tiny images. Tech. Rep. 633
(2009). 634
- [47] Muttenthaler, L. & Hebart, M. N. Thingsvision: A python toolbox for streamlining 635
the extraction of activations from deep neural networks. *Frontiers in Neuroinformat-* 636
ics **15**, 679838 (2021). 637
- [48] Nili, H. *et al.* A toolbox for representational similarity analysis. *PLOS Computational* 638
Biology **10**, 1–11 (2014). 639