

Some 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 that addresses this problem by looking for a *second-order isomorphism* in neural activation patterns. This innovation makes it easy to compare latent representations across individuals, species and computational models, and accounts for its popularity across disciplines ranging from artificial intelligence to computational neuroscience. Despite these successes, using RSA has led to difficult-to-reconcile and contradictory findings, particularly when comparing primate visual representations with deep neural networks (DNNs): even though DNNs have been shown to learn and behave in vastly different ways to humans, comparisons based on RSA have shown striking similarities in some studies. Here, we demonstrate some pitfalls of using RSA and explain how contradictory findings can 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, such as comparisons made between human visual representations and those of primates and DNNs, 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 multivariate analyses methods, such as Multi-Voxel Pattern (MVP) Classification, which uses machine learning algorithms to decode neural activity [4]. MVP classification assesses whether a brain region codes for a stimulus feature by examining whether the feature can be easily decoded from

neural response patterns in the region. However, there are at least two issues with using MVP classification for comparing mental representations across individuals. Firstly, just because a stimulus feature can be easily decoded from neural response patterns in a region does not imply that downstream regions in the brain actually decode this information [5]. Different individuals (or species) may use this information in different ways and MVP classification does not provide a way of capturing this. Secondly, there are methodological limitations on mapping brain regions and neural activity patterns between individuals and species. Therefore, even if two individuals represent a visual stimulus in the same manner, a decoder trained on one individual will show a significant performance drop when applied across individuals [6].

A newer addition to multivariate analysis, Representation Similarity Analysis (RSA), is specifically designed to compare representations between different systems and overcomes some 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 representational distance matrices (RDMs), a *second-order* measure of similarity that compares systems based on the relative distances between neural response patterns. This arrangement of neural response patterns in a representational space has been called a system's *representational geometry* [7]. The advantage of looking at representational geometries is that one no longer needs to match the architecture of two systems, or even the feature space of the two 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 recordings, EEG traces with behavioural data, or vectors in a computer algorithm

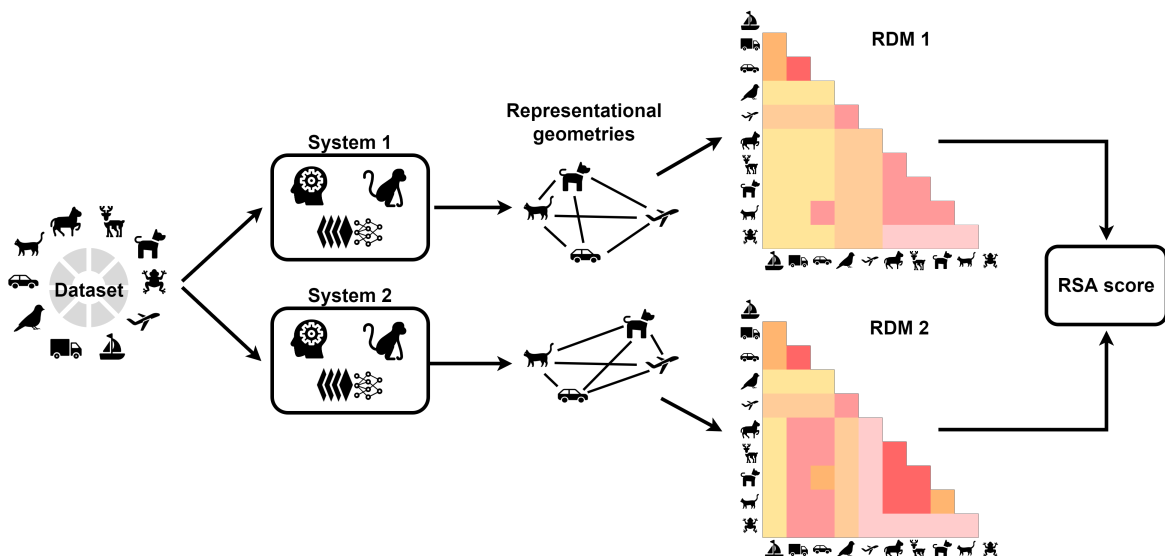


Figure 1: **RSA calculation.** (A) Stimuli from a set of categories (or conditions) are used as inputs to two different systems (for example, a human brain and a primate brain). Activity from regions of interest is recorded for each stimulus. Pair-wise distances in activity patterns are calculated to 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. It is up to the researcher to make a number of choices during this process including the choice of distance measure (e.g., 1-Pearson's r , Euclidean distance etc.) and a measure for comparing RDMs (e.g., Pearson's r , Spearman's ρ , Kendall's τ , etc.).

with spiking activity of neurons [8]. RSA is now ubiquitous in computational psychology 43
and neuroscience and has been applied to compare object representations in humans and 44
primates [9], representations of visual scenes by different individuals [6, 10], representa- 45
tions of visual scenes in different parts of the brain [11], to study specific processes such 46
as cognitive control [12] or the dynamics of object processing [13], and most recently, to 47
relate neuronal activations in human (and primate) visual cortex with activations of units 48
in Deep Neural Networks [14–18]. 49

However, this flexibility in the application of RSA comes at the cost of inferences one 50
can draw from this analysis. If the goal of the neuroscientist, psychologist or AI researcher 51

is to establish whether two systems are similar in mechanism, feature representation or 52
information processing, then RSA may not be the correct analytical method to use. This 53
is because RSA is a second-order measure – it looks at the similarity of similarities – that 54
abstracts over mechanism, feature representations and information processing. This point 55
has been made before. For example, Haxby et al. [4] write that the disadvantage of using 56
RSA is that: 57

...one cannot investigate whether the spaces in different subjects share the 58
same feature tuning functions or how these tuning function codes differ for 59
different brain regions. One cannot predict the response to a new stimulus in 60
a subject on the basis of the responses to that stimulus in other subjects. One 61
cannot predict the tuning function for individual neural features in terms of 62
stimulus features, precluding investigators from predicting the response pat- 63
tern vector for a new stimulus on the basis of its features. [pg. 446] 64

Despite these warnings, RSA continues to be used to infer that different individuals or 65
brain regions or computational models have similar mechanism (that is, they are similar 66
in nested functions and algorithms that transform inputs into neural response vectors). 67
One area where these conclusions are frequently made is the comparison between the 68
hierarchical representations in the visual cortex and Deep Neural Networks (DNNs). For 69
example, Cichy et al. [17] observed a correspondence in the RDMs of DNNs performing 70
object categorization and neural responses in human visual cortex recorded using MEG 71
and fMRI. Based on this correspondence, the authors concluded that: 72

...hierarchical systems of visual representations emerge in both the human 73
ventral and dorsal visual stream as the result of task constraints of object 74
categorization posed in everyday life, and provide strong evidence for object 75

representations in the dorsal stream independent of attention or motor inten- 76
tion. [pg. 5] 77

Thus, the correspondence in RDMs is used to infer the mechanism of emergence of visual 78
representations. Based on a similar comparison, Kriegeskorte [19] concluded that: 79

Deep convolutional feedforward networks for object recognition are not bio- 80
logically detailed and rely on nonlinearities and learning algorithms that may 81
differ from those of biological brains. Nevertheless they learn internal repre- 82
sentations that are highly similar to representations in human and nonhuman 83
primate IT cortex. [pg. 441] 84

While authors are sometimes careful in stating that the term ‘similarity in representations’ 85
is used as a shorthand for a ‘similarity in representational geometries’, they nevertheless 86
also invite the reader to accept that different systems show similar representational ge- 87
ometries because it is likely that they also use similar mechanisms to transform sensory 88
information into latent representations, or they use similar (downstream) mechanisms to 89
decode these latent representations. But how safe are these assumptions? 90

The main goal of our paper is to show that high RSA scores should not be used to infer 91
two systems have similar mechanisms. In Study 1, in a bare-bones setup, we show that 92
it is possible for two systems to transform input stimuli through known functions that 93
are vastly different but end up with similar representational geometries. In particular, 94
the study shows that 1) the presence of second-order confounds in the training data 95
can lead systems to mimic each other’s representational geometry even in the absence of 96
mechanistic similarity, and 2) the intrinsic structure of datasets rather than mechanistic 97
alignment can lead to artifactual modulation of RSA scores. Then in Studies 2 and 3 98
we show these problems extend to more complex datasets directly relevant to artificial 99

intelligence and computational neuroscience by making comparisons within and between 100
sets of artificial and biological systems. Finally, in Study 4, we show that not only are 101
misleadingly high RSA scores possible in practice but they are also highly plausible given 102
the hierarchical structure of categories in datasets that are routinely used. 103

Our demonstrations provide an explanation of how these phenomena, which arise 104
ubiquitously, can lead to incorrect inferences and contradictory or paradoxical findings. 105
For example, it has been recently observed that correlations in representational geometries 106
between human visual cortex and DNNs can vary from being close to the noise ceiling 107
to being uncorrelated based on the visual stimuli used in the experiments [20]. Since our 108
results have considerable generality with respect to current practices across multiple fields, 109
we discuss the implications for published results, including a discussion of two alternative 110
philosophical perspectives on the nature of mental representations that our findings speak 111
to. We conclude by providing some general recommendations regarding how to best use 112
RSA going forward. 113

Results 114

Proof of concept 115

It may be tempting to infer that two systems which have similar representational geome- 116
tries for a set of concepts do so because they encode similar properties of sensory data and 117
transform sensory data through a similar set of functions. In this section, we show that 118
it is possible, at least in principle, for qualitatively different systems to end up with very 119
similar representational geometries even though they (i) transform their inputs through 120
very different functions, and (ii) select different features of inputs. 121

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

Consider a population of animate and inanimate objects that consist of four categories of objects – birds, dogs, airplanes and bicycles. Each object in this population will have a set of stimulus features, using which one can map each exemplar from all four categories into a feature space. In Figure 2A (left), we show a hypothetical 2D feature space where exemplars from each category cluster together. Furthermore, we consider two datasets sampled from this population – Dataset A (Figure 2A, middle) which consists of birds and bicycles and Dataset B (Figure 2A, right) which consists of dogs and airplanes. Both datasets consist of animate and inanimate objects, but they differ in how items in each category are represented in the input space.

Now, consider two information-processing systems that re-represent Dataset A into two different latent spaces (Figure 2B). These could be two recognition systems designed to distinguish animate and inanimate categories. We assume that we can observe the representational geometry of the latent representations of each system and we are interested in understanding whether observing a strong correlation between these geometries implies whether the two systems have a similar *representational space* – that is, they project inputs into the latent space using similar functions. To examine this question, we consider a setup where we know the functions, Φ_1 and Φ_2 , that map the inputs to the latent space in each system. We will now demonstrate that even when these functions are qualitatively different from each other, the geometry of latent representations can nev-

ertheless be highly correlated. We will also show that the difference in representational spaces becomes more clear when one considers a different dataset (Dataset B), where inputs projected using the same functions now lead to a low correlation in representational geometries.

We can compute the geometry of a set of representations by establishing the pair-wise distance between all vectors in each representational space Φ . There are many different methods of computing this representational distance between any pair of vectors, all deriving from the dot product between vectors (see, for example, Figure 1 in [21]). Previous research has shown that the choice of the distance metric itself can influence the inferences one can draw from one’s analysis [21, 22]. However, here our focus is not the distance metric itself, but the fundamental nature of RSA. Therefore, we use the same generic distance metric – the dot product – to compute the pair-wise distance between all vectors in both representational spaces. In other words, the representational distance $d[\Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j)]$, between the projections of any pair of input stimuli, \mathbf{x}_i and \mathbf{x}_j into a feature space Φ , is proportional to the inner product between the projections in the feature space:

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

And we can obtain the representational geometry of the input stimuli $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ in any representational space Φ by computing the pairwise distances, $d[\Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j)]$ 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:

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

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

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. In other words, Φ_1 and Φ_2 are two fundamentally different projections of the inputs $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ – while Φ_2 maps a 2D input \mathbf{x}_i into a 2D feature space, Φ_1 maps the same 2D input into an infinite-dimensional space. Nevertheless, since cosine and RBF kernels are Mercer kernels, we can compute the distances (as measured by the dot product) between each pair of projected vectors using the kernel trick [23, 24]. That is, we can find the distance between any pair of points in the representational space by applying the kernel function to those points in the input space. These pairwise distances are shown by the kernel matrices in Figure 2B.

Next, we can determine how the geometry of these projections in the two systems relate to each other by computing the correlation between the kernel matrices, shown on the right-hand-side of Figure 2B. We can see from these results that the kernel matrices are highly correlated – i.e., the input stimuli are projected to very similar geometries in the two representational spaces.

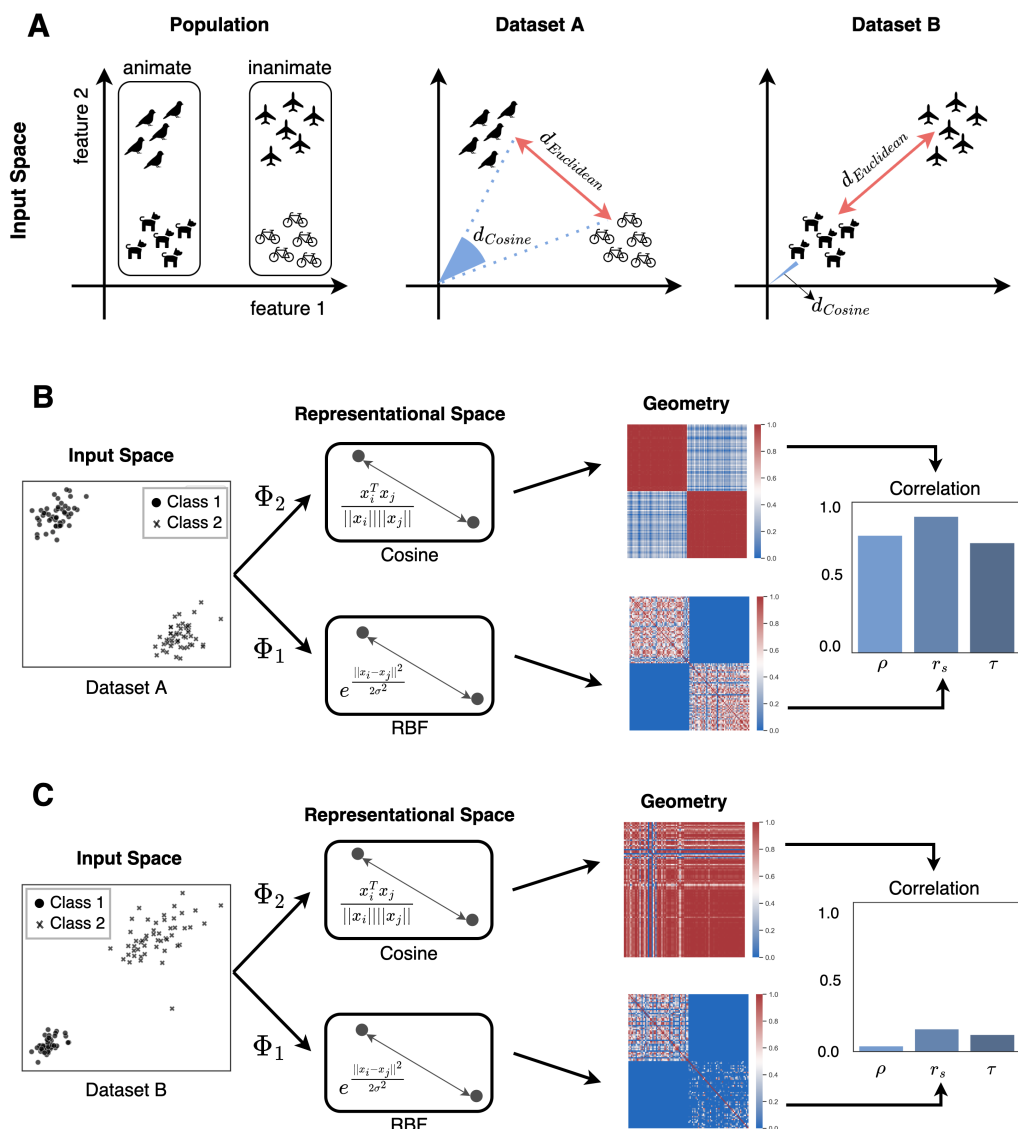


Figure 2: Mimic and modulation effect in representational geometries. (A) An example of a population of animate (birds, dogs) and inanimate (planes, bikes) objects, plotted in a hypothetical 2D stimulus feature space. Two datasets are sampled from this population: In Dataset A (middle), the Euclidean distance (in input space) between categories mirrors the Cosine distance, while in Dataset B (right) it does not. (B) Simulation where two systems transform stimuli in Dataset A into latent representations such that the (dot product) distance between latent vectors is given by RBF and Cosine kernels, respectively. As Euclidean and Cosine distances in the input space mirror each other, the representational geometries (visualised here using kernel matrices) end up being highly correlated (shown using Pearson (ρ), Spearman (r_s) and Kendall's (τ) correlation coefficients on the right). We call this strong correlation in representational geometries despite a difference in input transformation a *mimic effect*. (C) Simulation where objects in Dataset B are projected using same transformations as (B). The (dot product) distance is still given by the same (RBF and Cosine) kernels. However, for this dataset, the Euclidean and Cosine distances in input space do *not* mirror each other and as a consequence, the representational geometries show low correlation. Thus the correlation in representational geometries depends on how the datasets are sampled from the population. We call this change in correlation a *modulation effect*.

If one did not know the input transformations and simply observed the correlation 183
between kernel matrices, it would be tempting to infer that the two systems Φ_1 and Φ_2 184
transform an unknown input stimulus \mathbf{x} through a similar set of functions – for example 185
functions that belong to the same class or project inputs to similar representational spaces. 186
However, this would be an error. The projections $\Phi_1(\mathbf{x})$ and $\Phi_2(\mathbf{x})$ are fundamentally 187
different – Φ_1 (radial basis kernel) projects an input vector into an infinite dimensional 188
space, while Φ_2 (cosine kernel) projects it onto a unit sphere. The difference between these 189
functions becomes apparent if one considers how this correlation changes if one considers a 190
different set of input stimuli. For example, the set of data points from Dataset B (sampled 191
from the same population) are projected to very different geometries, leading to a low 192
correlation between the two kernel matrices (Figure 2C). 193

In fact, the reason for highly correlated kernel matrices in Figure 2B is not a similarity 194
in the transformations Φ_1 and Φ_2 but the structure of the dataset. The representational 195
distance between any two points in the first representational space, $d[\Phi_1(\mathbf{x}_i), \Phi_1(\mathbf{x}_j)]$, is 196
 $e^{\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$. That is, the representational distance in Φ_1 is a function of their Euclidean 197
distance $\|\mathbf{x}_i - \mathbf{x}_j\|$ in the input space. On the other hand, the representational dis- 198
tance between any two points in the second representational space, $d[\Phi_2(\mathbf{x}_i), \Phi_2(\mathbf{x}_j)]$, is, 199
 $\frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$. That is, the representational distance in Φ_2 is a function of their cosine distance 200
in the input space.. These two stimulus features – Euclidean distance and cosine distance 201
– are *confounds* that lead to the same representational geometries for certain datasets. 202
In Dataset A, the stimuli is clustered such that the Euclidean distance between any two 203
stimuli is correlated with their cosine distance (see Figure 2A, middle). However, for 204
Dataset B, the Euclidean distance is no longer correlated with the angle (see Figure 2A, 205
right) and the confounds lead to different representational geometries, as can be seen in 206
Figure 2C. Thus, this example illustrates two effects: (i) a *mimic* effect, where two sys- 207

tems that transform sensory input through very different functions end up with similar 208
representational geometries (Figure 2B) , and (ii) a *modulation* effect, where two systems 209
that are non-identical have similar representational geometries for one set of inputs, but 210
dissimilar geometries for a second set (compare Figures 2B and 2C). 211

**Study 2: Complex systems encoding different features of inputs can show a 212
high RSA-score** Study 1 made a number of simplifying assumptions – the dataset was 213
two-dimensional, clustered into two categories and we intentionally chose functions Φ_1 214
and Φ_2 such that the kernel matrices were correlated in one case and not correlated in the 215
other. It could be argued that, even though the above results hold in principle, they are 216
unlikely in practice when the transformations and data structure are more complex. For 217
example, it might be tempting to assume that accidental similarity in representational 218
geometries becomes less likely as one increases the number of categories (i.e., clusters or 219
conditions) being considered. However, In Figure 3 we illustrate how complex systems 220
transforming high-dimensional input from a number of categories may achieve high RSA 221
scores. Even though one system extracts surface reflectance and the other extracts global 222
shape, they can end up with very similar representational geometries. This would occur 223
if objects similar in their reflectance properties were also similar in shape (e.g., glossy 224
balloons and light bulbs) and if objects dissimilar according to reflectance properties were 225
also dissimilar in shape (e.g., dogs and light bulbs). This is the mimic effect, where 226
representational geometries of these two systems end up being similar because reflectance 227
and shape are second-order confounds in this dataset. Conducting RSA on this dataset 228
will show a high correlation in RDMs, even though the latent representations in these 229
systems are related to very different stimulus features. 230

To demonstrate this empirically, we now consider a more complex setup, where the 231

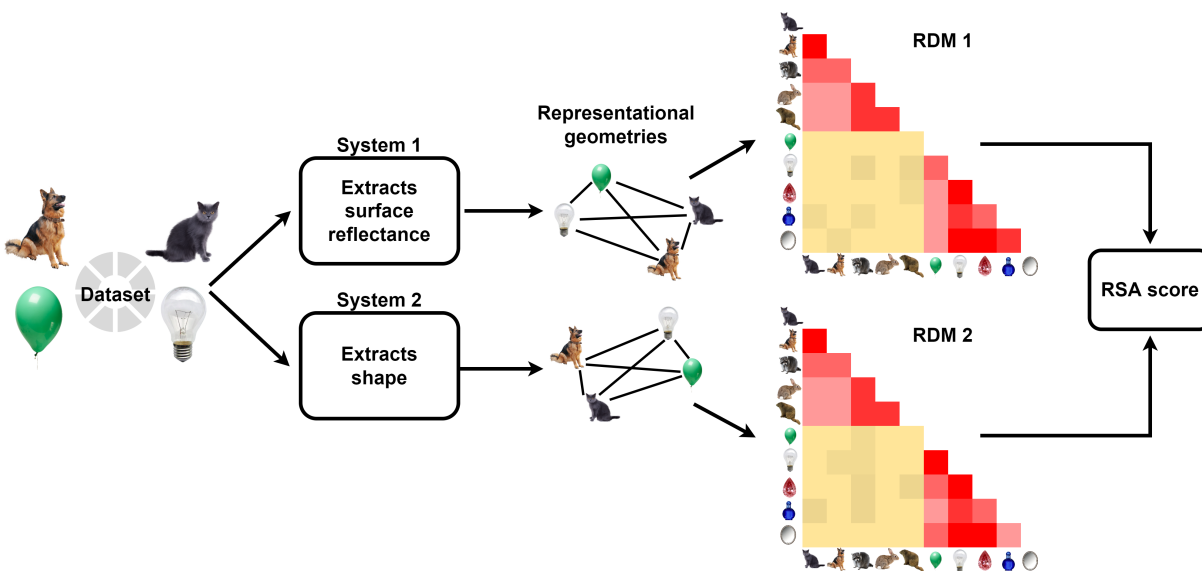


Figure 3: **Example of a second-order confound.** Two systems, one forming representations based on surface reflectance of objects (while ignoring all other features such as colour or texture) and the other based on global shape (while ignoring other features), can have very similar representational geometries. This similarity would lead to a high RSA score but would not justify an inference about the representations being similar.

transformations Φ_1 and Φ_2 are modelled as feedforward deep neural networks (DNNs), 232
 trained to classify a high-dimensional dataset into multiple categories. Many studies that 233
 use RSA compare systems using naturalistic images as visual inputs [9, 14]. While using 234
 naturalistic images brings research closer to the real-world, it is also well-known that 235
 datasets of naturalistic images frequently contain confounds – independent features that 236
 can predict image categories [25]. We will now show how the simplest of such confounds, 237
 a single pixel, can lead to a high RSA score between two DNNs that encode qualitatively 238
 different features of inputs. 239

Consider the same setup as above, where an input stimulus, \mathbf{x} , is transformed to a 240
 representation space by two systems, Φ_1 and Φ_2 . Instead of a two-dimensional input space, 241
 \mathbf{x} now exists in a high-dimensional image space and Φ_1 and Φ_2 are two versions of a DNN – 242
 VGG-16 – trained to classify input images into different categories. We ensured that Φ_1 and 243

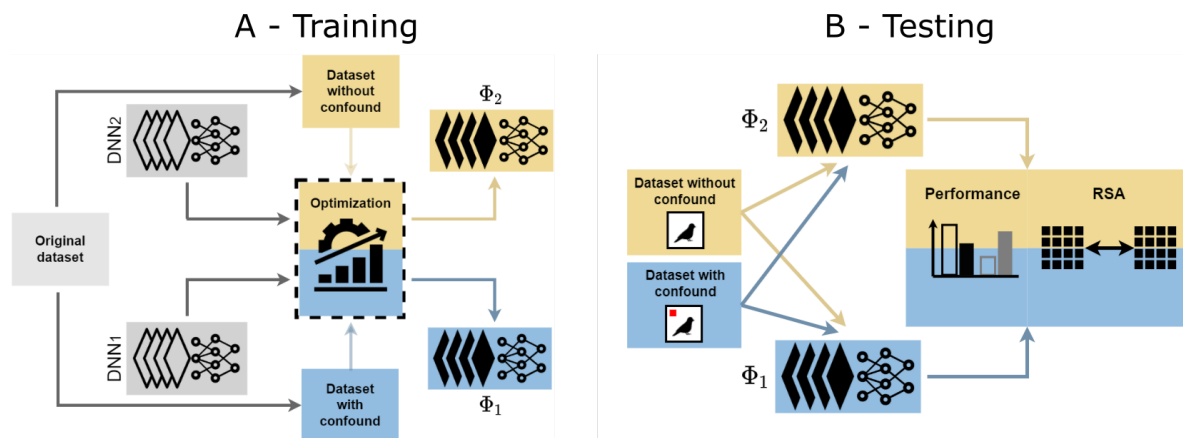


Figure 4: **Training and testing DNNs with different feature encodings.** Panel A shows the training procedure for Studies 2–4, where we created two versions of the original dataset (gray), one containing a confound (blue) and the other left unperturbed (yellow). These two datasets were used to train 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.

Φ_2 were qualitatively different transformations of input stimuli by making the networks 244 sensitive to different predictive features within the stimuli. The first network was trained 245 on an unperturbed dataset, while the second network was trained on a modified version 246 of the dataset, where each image was modified to contain a confound – a single pixel in a 247 location that was diagnostic of the category (see Figure 4 for the general approach). 248

The locations of these diagnostic pixels were chosen such that they were correlated to 249 the corresponding representational distances between classes in Φ_1 . Our hypothesis was 250 that if the representational distances in Φ_2 preserve the physical distances of diagnos- 251 tic pixels in input space, then this confound will end up mimicking the representational 252 geometry of Φ_1 , even though the two systems use qualitatively different features for clas- 253 sification. Furthermore, we trained two more networks, Φ_3 and Φ_4 , which were identical 254

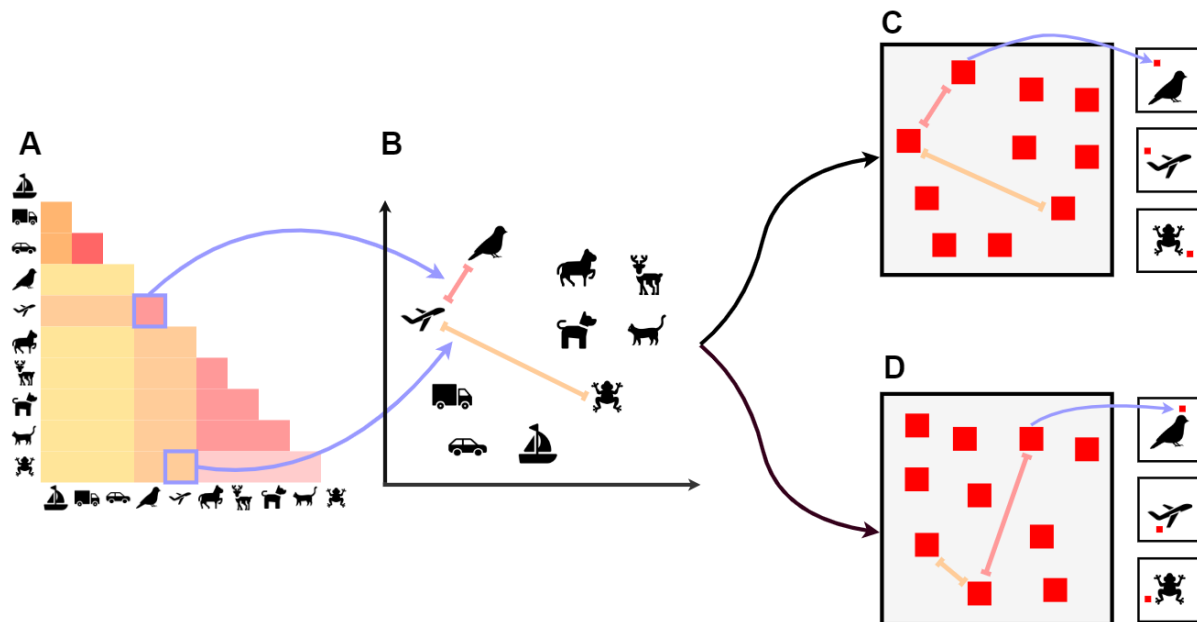


Figure 5: **Study 2 confound placement.** The representational geometry (Panel A and B) from the network trained on the unperturbed CIFAR-10 images is used to determine the location of the single pixel confound (shown as a red patch here) for each category. In the ‘Positive’ condition (Panel C), we determined 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.

to Φ_2 , except these networks were trained on datasets where the location of the confound 255
was uncorrelated (Φ_3) or negatively correlated (Φ_4) with the representational distances 256
in Φ_1 (see Figure 5 and Methods for details). 257

Classification accuracy (Figure 6 (left)) revealed that the network Φ_1 , trained on the 258
unperturbed images, learned to classify these images and ignored the diagnostic pixel 259
– that is, its performance was identical for the unperturbed and modified images. In 260
contrast, networks Φ_2 (positive), Φ_3 (uncorrelated) and Φ_4 (negative) failed to classify the 261
unperturbed images (performance was near chance) but learned to perfectly classify the 262
modified images, showing that these networks develop qualitatively different representa- 263
tions compared to normally trained networks. 264

Next we computed pairwise RSA scores between the representations at the last con- 265
volution layer of Φ_1 and each of Φ_2 , Φ_3 and Φ_4 (Figure 6 (right)). When presented un- 266
perturbed test images, the Φ_2 , Φ_3 and Φ_4 networks all showed low RSA scores with the 267
normally trained Φ_1 network. However, when networks were presented with test images 268
that included the predictive pixels, RSA varied depending on the geometry of pixel loca- 269
tions in the input space. When the geometry of pixel locations was positively correlated 270
to the normally trained network, RSA scores approached ceiling (i.e., comparable to RSA 271
scores between two normally trained networks). Networks trained on uncorrelated and 272
negatively correlated pixel placements scored much lower. 273

These results mirror Study 1: we observed that it is possible for two networks (Φ_1 and 274
 Φ_2) to show highly correlated representational geometries even though these networks 275
learn to classify images based on very different features. One may argue that this could 276
be because the two networks could have learned similar representations at the final con- 277
volution layer of the DNN and it is the classifier that sits on top of this representation 278
that leads to the behavioural differences between these networks. But if this was true, it 279

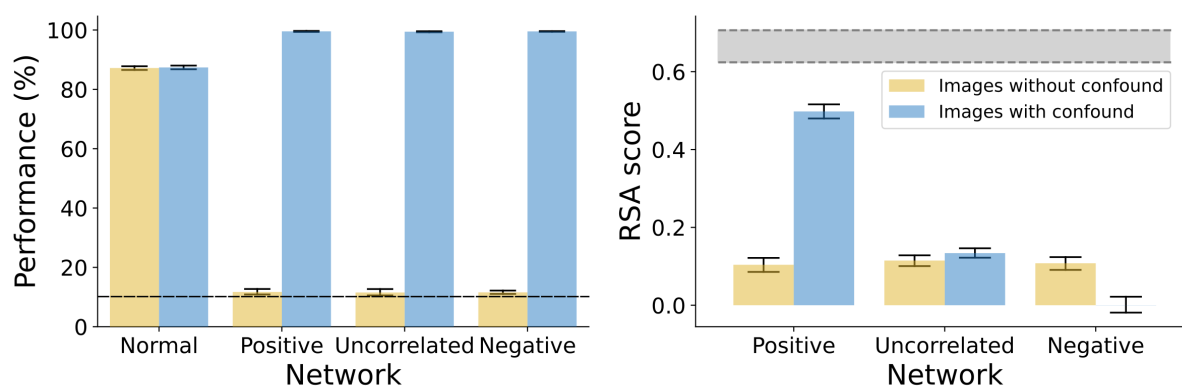


Figure 6: **Study 2 results.** *Left:* Performance of normally trained networks did not depend on whether classification was done on unperturbed CIFAR-10 images or images with a single pixel confound (error bars represent 95% CI, the dashed line represents chance performance). All three networks trained 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.

would not explain why RSA scores diminish for the two other comparisons (with Φ_3 and Φ_4). This modulation of RSA-scores for different datasets suggests that, like in Study 1, the correlation in representational geometry is not because the two systems encode similar features of inputs, but because different features mimic each other in their representational geometries.

Re-examining some influential findings

In Studies 1 and 2, we showed that it is possible for qualitatively different systems to end up with similar representational geometries. However, it may be argued that while this is possible in principle, it is unlikely in practice in real-world scenarios. In the following two studies, we consider real-world data from some recent influential experiments, recorded from both primate and human cortex. We show how RSA-scores can be driven by confounds in these real-world settings and how properties of training and test data may contribute to observed RSA-scores.

Study 3: Neural activations in monkey IT cortex can show a high RSA-score with DNNs despite different encoding of input data

In our next study, we consider data from experiments comparing representational geometries between computational models and macaque visual cortex [14, 26]. The experimental setup was similar to Study 2, though note that unlike Study 2, where both systems used the same architecture and learning algorithm, this study considered two very different systems – one artificial (DNN) and the other biological (macaque IT cortex). We used the same set of images that were shown to macaques by Majaaj et al. [27] and modified this dataset to superimpose a small diagnostic patch on each image. In the same manner as in Study 2 above, we constructed three different datasets, where the locations of these diagnostic

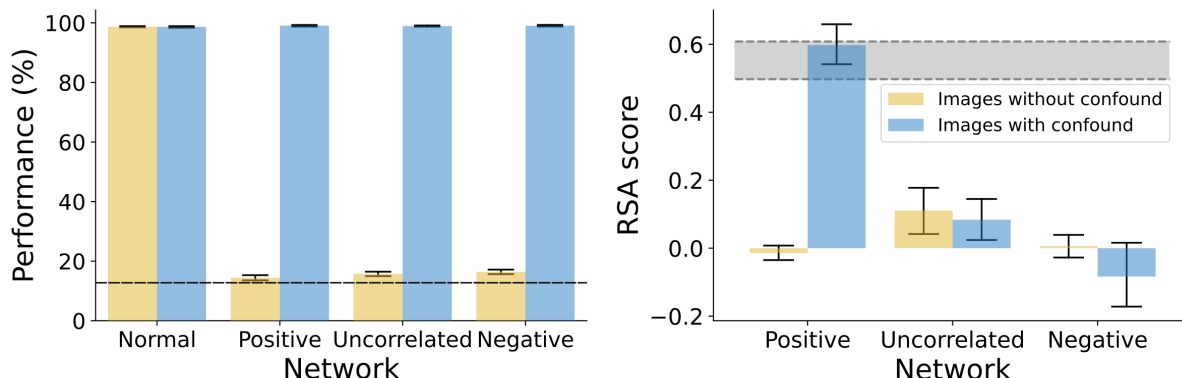


Figure 7: **Study 3 results.** *Left:* Classification Performance of the network trained on unperturbed images (Normal condition) did not depend on the presence or absence of the confound, while performance of networks trained with the confound (Positive, Uncorrelated and Negative conditions) highly depended on whether the confound was present (dashed line represents chance performance). *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.

patches were either positively correlated, uncorrelated or negatively correlated with the 303
RDM of macaque activations. We then trained four CNNs. The first CNN was pre- 304
trained on ImageNet and then fine-tuned on the unmodified dataset of images shown to 305
the macaques. Previous research has shown that CNNs trained in this manner develop 306
representations that mirror the representational geometry of neurons in primate inferior 307
temporal (IT) cortex [14]. The other three networks were trained on the three modi- 308
fied datasets and learned to entirely rely on the diagnostic patches (accuracy on images 309
without the diagnostic patches was around chance). 310

Figure 7 (right) shows the correlation in representational geometry between the macaque 311
IT activations and activations at the final convolution layer for each of these networks. 312
The correlation with networks trained on the unmodified images is our baseline and shown 313
as the gray band in Figure 7. Our first observation was that a CNN trained to rely on 314

the diagnostic patch can indeed achieve a high RSA score with macaque IT activations. 315
In fact, the networks trained on patch locations that were positively correlated to the 316
macaque RDM matched the RSA score of the CNNs trained on ImageNet and the unmod- 317
ified dataset. This shows how two systems having very different architectures, encoding 318
fundamentally different features of inputs (single patch vs naturalistic features) can show 319
a high correspondence in their representational geometries. We also observed that, like 320
in Study 2, the RSA score depended on the clustering of data in the input space – when 321
patches were placed in other locations (uncorrelated or negatively correlated to macaque 322
RDMs) the RSA score became significantly lower. 323

Study 4: High RSA-scores may be driven by the structure of testing data All 324
the studies so far have used the same method to construct datasets with confounds – we 325
established the representational geometry of one system (Φ_1) and constructed datasets 326
where the clustering of features (pixels) mirrored this geometry. However, it could be 327
argued that confounds which cluster in this manner are unlikely in practice. For example, 328
even if texture and shape exist as confounds in a dataset, the inter-category distances 329
between textures are not necessarily similar to the inter-category distances between shape. 330

However, categories in real-world datasets are usually hierarchically clustered into 331
higher-level and lower-level categories. For example, in the CIFAR-10 dataset, the Dogs 332
and Cats (lower-level categories) are both animate (members of a common higher-level 333
category) and Airplanes and Ships (lower-level categories) are both inanimate (members 334
of a higher-level category). Due to this hierarchical structure, Dog and Cat images are 335
likely to be closer to each other not only in their shape, but also their colour and texture 336
(amongst other features) than they are to Airplane and Ship images. In our next simula- 337
tion, we explore whether this hierarchical structure of categories can lead to a correlation 338

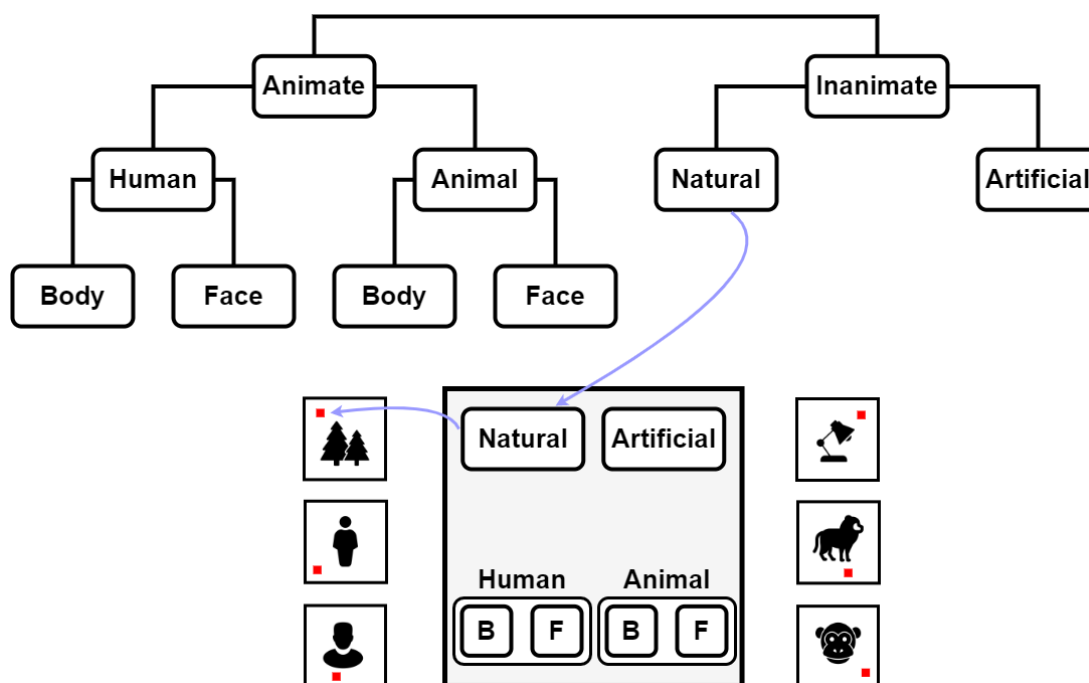


Figure 8: **Exploiting intrinsic dataset hierarchy in order to place confounds.** The top panel shows the hierarchical structure of categories in the dataset, which was used to place the single pixel confounds. The example at the bottom (middle) shows one such hierarchical placement scheme where the 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.

in representational geometries between two systems that learn different feature encodings. 339

For this study, we selected a popular dataset used for comparing representational 340
geometries in humans, macaques and deep learning models [15, 28]. This dataset consists 341
of six categories which can be organised into a hierarchical structure shown in Figure 8. [9] 342
showed a striking match in RDMs for response patterns elicited by these stimuli in human 343
and macaque IT. For both humans and macaques, distances in response patterns were 344
larger between the higher-level categories (animate and inanimate) than between the 345
lower-level categories (e.g., between human bodies and human faces). 346

We used a similar experimental paradigm to the above studies, where we trained 347

networks to classify stimuli which included a single predictive pixel. But instead of using 348
an RDM to compute the location of a diagnostic pixel, we used the hierarchical categorical 349
structure. In the first modified version of the dataset, the location of the pixel was based 350
on the hierarchical structure of categories in Figure 8 – predictive pixels for animate 351
kinds were closer to each other than to inanimate kinds, and pixels for faces were closer 352
to each other than to bodies, etc. One such configuration can be seen in Figure 8. In the 353
second version, the predictive pixel was placed at a random location for each category 354
(but, of course, at the same location for all images within each category). We call these 355
conditions ‘Hierarchical’ and ‘Random’. [15] showed that the RDM of average response 356
patterns elicited in the human IT cortex (Φ_1) correlated with the RDM of a DNN trained 357
on naturalistic images (Φ_2). We explored how this compared to the correlation with the 358
RDM of a network trained on the Hierarchical pixel placement (Φ_3) and Random pixel 359
placement (Φ_4). 360

Results for this study are shown in Figure 9. We observed that representational ge- 361
ometry of a network trained on Hierarchically placed pixels (Φ_3) was just as correlated to 362
the representational geometry of human IT responses (Φ_1) as a network trained on natu- 363
ralistic images (Φ_2). However, when the pixel locations for each category were randomly 364
chosen, this correlation decreased significantly. These results suggest that any confound in 365
the dataset (including texture, colour or low-level visual information) that has distances 366
governed by the hierarchical clustering structure of the data could underlie the observed 367
similarity in representational geometries between CNNs and human IT. More generally, 368
these results show how it is plausible that many confounds present in popular datasets 369
may underlie the observed similarity in representational geometries between two systems. 370
The error of inferring a similarity in mechanism based on a high RSA score is not just 371
possible but also probable. 372

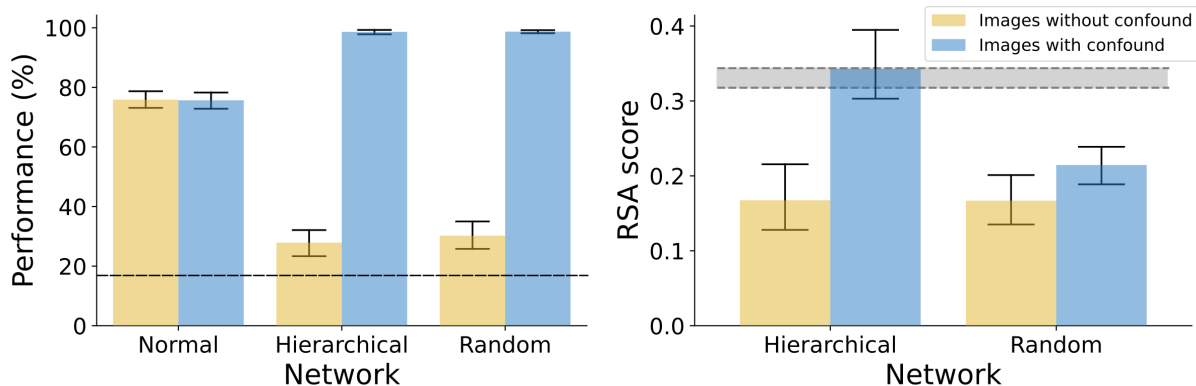


Figure 9: **Study 4 results.** *Left:* Performance of normally trained networks did not depend on whether the confound was present. Networks trained with the confound failed to classify stimuli without the confound (yellow bars) while achieving near perfect classification of stimuli with the confound present (blue bars, dashed line represents chance performance). *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

373

In four studies, we have illustrated a number of conditions under which it can be problem- 374
atic to infer a similarity of representations between two systems based on a correlation in 375
their representational geometries. In particular, we showed that two systems may trans- 376
form their inputs through very different functions and encode very different features of 377
inputs and yet have highly correlated representational geometries. Of course, one may 378
acknowledge that a second-order isomorphism of activity patterns does *not* strictly imply 379
that two systems are similar mechanistically but still assume that it is highly likely to 380
be the case. That is, as a practical matter, a researcher may assume that RSA is a reli- 381
able method to compare systems. However, our findings challenge this assumption. We 382
show how a high RSA score between different systems can not only occur in a bare-bones 383
simulation (Study 1), but also in practice, in high-dimensional systems operating on high- 384

dimensional data (Studies 2–3). Furthermore, we show that the hierarchical structure of 385
datasets frequently used to test similarity of representations lends itself to a high RSA 386
score arising because of second-order confounds present in the dataset (Study 4). There- 387
fore, second-order confounds driving high RSA scores is not only possible but plausible. 388

One limitation of our method is that we manually insert a confound in input stimuli 389
(in Studies 2–4) and train a network based on this confound. Even though our find- 390
ings demonstrate that second-order confounds are plausible, they do not allow us to infer 391
whether such confounds *are* present in existing datasets and driving the observed similar- 392
ity in existing studies. In our view, there are two methods one could use to check whether 393
confounds are driving results of RSA. The best way would be to identify the stimulus 394
features in a dataset that mimic each other in representational space (e.g. shape and re- 395
flectance in Figure 3). This is not straightforward to do in high-dimensional stimuli, such 396
as naturalistic images, which consist of millions of features. However, another approach is 397
more tractable: conduct controlled experiments to establish whether the two systems are 398
representing information in similar ways. We have argued for this approach in relation 399
to making inferences about mechanistic similarity between DNNs and humans [29]. In 400
fact, research relating DNNs to human vision provides a striking case of a disconnect 401
between RSA and behavioural findings from psychology [29–31]. The findings here may 402
explain contradictory RSA scores between DNNs and human visual processing as pointed 403
out by Xu and Vaziri-Pashkam [20]. At the very least, a researcher claiming that two 404
systems are mechanistically similar to one another based on high RSA scores should have 405
an explanation for this discrepancy. 406

A related point has been made by Kriegeskorte and Diedrichson [32] and Kriegeskorte 407
and Wei [33], who point out that two systems may have the same representational geome- 408
try, even if they have a different activity profile over neurons. In this sense, the geometry 409

abstracts away the information about how information was distributed over a set of neu- 410
rons. Kriegeskorte and Diedrichson [32] equate this loss in information to “peeling a layer 411
of an onion” – downstream decoders that are sensitive to the representational geometry 412
rather than activity profiles over neuron populations can focus on difference in information 413
as reflected by a change in geometry and be agnostic to how this information is distributed 414
over a set of neurons. We agree that this invariance over activity profiles is indeed a useful 415
property of representational geometries for downstream decoders. However, we are not 416
aware of any studies that highlight how representational geometries also abstract over be- 417
haviourally relevant stimulus properties . While abstracting over activity profiles may be 418
useful, abstracting over stimulus properties loses an important piece of information when 419
comparing representations across brain regions, individuals, species and between brains 420
and computational models. Our studies show how two systems may appear similar based 421
on their representational geometries in one circumstance (e.g. Figure 2B) but drastically 422
different in another circumstance (Figure 2C). 423

It is important to note how our results differ from previous studies exploring limita- 424
tions of RSA. A number of studies have focused on the importance of how neural data is 425
pre-processed and how the distance between neural patterns is computed. For example, 426
Ramirez [34] found that pre-processing steps, such as centering (de-meaning) activation 427
vectors may lead to incorrect inference about the representational geometry of activations. 428
He demonstrated that subtracting the mean from activations could change the rank or- 429
der of similarity between conditions. In turn, this could lead to clearly distinct RDMs 430
becoming highly correlated and vice-versa. While this is an important methodological 431
point, it is clearly distinct from the point we are making in this study. Indeed, the results 432
here are agnostic of the data pre-processing steps and hold whether or not activations are 433
centered. 434

Some previous studies have also explored how confounds present in data can influence 435
the results of RSA. For example, Henriksson et al. [35] and Cai et al. [36] demonstrated 436
that RDMs measured based on fMRI data can be severely biased because of temporal 437
and spatial correlations in neural activity. These authors have pointed out that if activity 438
patterns from different brain regions are recorded during the same trial, the similarity 439
estimates will be exaggerated due to correlated neural fluctuations in these regions. Sim- 440
ilarly, neural activity is correlated over time, which means estimated similarity based on 441
activity patterns from the same imaging run also introduces a strong bias in RDMs. These 442
sources of bias are important to understand, but they can also be addressed by a more 443
careful task design and analysis [36]. In contrast, the confounds that are highlighted in 444
this study exist in the stimulus itself. Therefore, even if one were to completely mitigate 445
the bias in estimating RDMs, the types of confounds we highlight in our work would still 446
pose problems when drawing inferences from correlation in RDMs. 447

Similarly, previous research has also highlighted the importance of choosing the correct 448
distance metric when using RSA. For example, Ramirez [22] compared Euclidean distance 449
with an angular metric (such as cosine similarity) and showed that the choice of distance 450
metric can reveal different aspects of the same fMRI data. They argued that the Euclidean 451
distance is particularly sensitive to the mean activity over a recorded voxel. Based on this 452
analysis, Ramirez [22] suggested using an angular distance metric, especially when neural 453
signal is aggregated over large number of neurons. Similarly, in another exhaustive study 454
over distance measures, Bobadilla-Suarez et al. [21], evaluated neural similarity using 455
various distance measures, including angle-based measures (cosine, Pearson, Spearman) 456
and magnitude-based measures (Euclidean, Mahalanobis, Minkowski) and found that the 457
choice of metric significantly influenced the measured similarity. They also found that 458
there was no one metric that outperformed all others – rather, the preferred metric varied 459

across different studies, but was consistent across brain regions within a study. The choice 460
of distance metric is again a related but orthogonal issue to the one we highlight in this 461
study. Our results show that representational geometry loses information about stimulus 462
features and different stimulus features (and indeed transformations of input stimulus) 463
can lead to the same geometry. This is fundamental to the nature of representational 464
geometries, rather than a consequence of the distance metric used. 465

Of course, the problem of confounds in stimuli is not unique to RSA and will affect 466
other statistical analyses, including multivariate regression methods such as MVP classi- 467
fication. Indeed, the problem of confounds in stimuli is well appreciated in many different 468
contexts [25, 37, 38], but there has been no consideration of whether these confounds are 469
contributing to RSA findings. Perhaps this is because, unlike for MVP classification, a 470
confound for RSA needs to not only help decode category membership, but also lead to 471
a second-order isomorphism. Nevertheless, as we illustrate in Figure 3, there could be 472
confounds with a second-order similarity structure in many datasets that are the product 473
of unexpected properties of the world or the product of how these datasets are curated or 474
hierarchically organized. This is problematic as we have clearly shown that these second 475
order confounds can drive high RSA scores. 476

A reader could ask why these results matter. Couldn't a researcher take the view 477
that representational geometry *is* representation and therefore, a strong correlation in 478
representational geometries between two systems is sufficient to infer that the systems are 479
representing the world in a similar manner? This question goes to the heart of an existing 480
debate in philosophy, where philosophers distinguish between the *externalist* and *holistic* 481
views on mental representations. According to the first view, the content of representa- 482
tions is determined by their relationship to entities in the external world. This perspective 483
is implicitly taken by most neuroscientists and psychologists, who are interested in com- 484

paring mechanisms underlying cognitive processes – that is, they are interested in the set 485
of nested functions and algorithms responsible for transforming sensory input into a set 486
of activations in the brain. From this perspective, our finding that high RSA scores can 487
be obtained between systems that work in qualitatively different ways poses a challenge 488
to researchers using RSA to compare systems. 489

Alternatively, a researcher may reject an externalist view and adopt the perspective 490
that representations obtain their meaning based on how they are related to each other 491
within each system, rather than based on their relationship to entities in the external 492
world. That is, “representation *is* the representation of similarities” [39]. From this per- 493
spective, as long as the two systems share the same relational distances between internal 494
activations, one can validly infer that the two systems have similar representations. That 495
is, a second-order isomorphism implies a similarity of representations, by definition. This 496
view has been called *holism* in the philosophy of mind [40, 41] and is related to a similar 497
idea of *meaning holism* in language, which is the idea that the meaning of a linguistic 498
expression is determined by its relation to other expressions within a language [42, 43]. 499
For example, Firth [44] (p. 11) writes: “you shall know a word by the company it keeps”. 500
Similarly, Griffiths and Steyvers [45], and Griffiths, Steyvers, and Tenenbaum [46] have 501
adopted meaning holism accounts of semantic representations in neural networks. More 502
recently, Piantadosi and Hill [47] have argued that large language models capture im- 503
portant aspects of meaning and approximate human cognition because they represent 504
relations between concepts and their roles within a representational geometry. Even if a 505
researcher was to adopt this holistic perspective on representations, our results should still 506
be of interest to them as they show that the similarity between representational geometries 507
can vary based on the visual stimulus that is used to compare them (the modulation ef- 508
fect). Additionally, our results show that adopting this view misses the information about 509

differences in mechanistic processes that a psychologist or neuroscientist is frequently in- 510
terested in, for instance, whether the visual system processes surface reflectance or shape 511
(or the location of diagnostic pixels) in order to identify objects. Fodor and Lepore long 512
ago criticized this philosophical stance [41,48], and interestingly, this philosophical debate 513
played an important part in the development of RSA (see Supplementary Information, 514
Section A). Unfortunately, this debate has largely been ignored by researchers who use 515
RSA as a method to compare similarity of systems. 516

In closing, we describe our recommendations for practitioners who find RSA to be 517
useful for their research goals. These will be especially relevant to the large majority 518
of researchers in computational, cognitive, and systems neuroscience, cognitive scientists 519
and AI practitioners, who are interested in mechanistic similarities (i.e., they adopt an 520
externalist position). But they should also be relevant to adopters of the holistic view 521
who are interested in how observed representational geometries depend on the stimulus 522
used to extract them. 523

First, since the intrinsic structure of datasets can artificially modulate RSA scores, 524
researchers should compare systems on a wider variety of datasets and sampling schemes 525
than currently done. Second, given that confounding features can lead to mimicked rep- 526
resentational geometries, researchers should consider running additional controlled ex- 527
periments to rule out this possibility when inferences hinge crucially on it. Third, when 528
studies are conducted to search for evidence of mechanistic similarity between two or more 529
systems, researchers should use a wider range of complementary methods, each addressing 530
the others' blindspots (e.g., RSA combined with neural predictivity [14], MVPC [6, 49], 531
CCA [50], SVCCA [51], CKA [52]). 532

Lastly, perhaps the most important general recommendation we make is that re- 533
searchers should acknowledge, procedurally and in writing, which inferences are afforded 534

by the use of RSA, and what dissimilarities remain possible despite having observed a
given pattern of RSA scores. To this end, we believe that general statements of simi-
larity tend to obfuscate rather than accurately summarize any set of RSA-based results.
Instead, we urge researchers using RSA (1) to justify the use of this method by theoret-
ically motivated interest in representational geometry or otherwise consider other tools
that best fit their goals, and (2) to state in precise terms that RSA scores reflect the
similarity of representational geometries in particular, and generally avoid underspecified
claims of similarity.

Methods

Dataset generation and training

All DNN simulations (Studies 2–4) were carried out using the Pytorch framework [53].
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
from two multivariate Gaussians, $\mathcal{N}(x|\mu, \Sigma)$, where μ is a 2-dimensional vector and Σ is
a 2×2 covariance matrix. In Figure 2A, the two Gaussians have means $\mu_1 = (1, 8)$ and
 $\mu_2 = (8, 1)$ and a covariance matrices $\Sigma_1 = \Sigma_2 = \frac{1}{2}\mathbf{I}$, while in Figure 2B the Gaussians
have means $\mu_1 = (1, 1)$ and $\mu_2 = (8, 8)$ and a covariance matrices $\Sigma_1 = \mathbf{I}$, $\Sigma_2 = 8\mathbf{I}$.
All kernel matrices were computed using the sklearn.metrics.pairwise module of the

scikit-learn Python package.

557

Study 2 First, a VGG-16 deep convolutional neural network [54], pre-trained on the ImageNet dataset of naturalistic images, was trained to classify stimuli from the CIFAR-10 dataset [55]. The CIFAR-10 dataset includes 10 categories with 5000 training, and 1000 test images per category. The network was fine-tuned on CIFAR-10 by replacing the classifier so that the final fully-connected layer reflected the correct number of target classes in CIFAR-10 (10 for CIFAR-10 as opposed to 1000 for ImageNet). Images were rescaled to a size of 224×224 px and then the model learnt to minimise the cross-entropy error using the RMSprop optimizer with a mini-batch size of 64, learning rate of 10^{-5} , and momentum of 0.9. All models were trained for 10 epochs, which were sufficient for convergence across all datasets.

558

559

560

561

562

563

564

565

566

567

Second, 100 random images from the test set for each category were sampled as input for the network and activations at the final convolutional layer extracted using the THINGSVision Python toolkit [56]. The same toolkit was used to generate a representational dissimilarity matrix (RDM) from the pattern of activations using 1-Pearson's r as the distance metric. The RDM was then averaged by calculating the median distance between each instance of one category with each instance of the others (e.g., the median distance between Airplane and Ship was the median of all pair-wise distances between activity patterns for airplane and ship stimuli). This resulted in a 10×10 , category-level, RDM which reflected median between-category distances.

568

569

570

571

572

573

574

575

576

Third, three modified versions of the CIFAR-10 datasets were created for the 'Positive', 'Uncorrelated' and 'Negative' conditions, respectively. In each dataset, we added one diagnostic pixel to each image, where the location of the pixel depended on the category (See Figure 5). The locations of these pixels were determined using the averaged RDM

577

578

579

580

from the previous step. We call this the target RDM. In the ‘Positive’ condition, we 581
wanted the distances between pixel placements to be positively correlated to the distances 582
between categories in the target RDM. We achieved this by using an iterative algorithm 583
that sampled pixel placements at random, calculated an RDM based on distances between 584
the pixel placements and computed an RSA score (Spearman correlation) with the target 585
RDM. Placements with a score above 0.70 were retained and further optimized (using 586
small perturbations) to achieve an RSA-score over 0.90. The same procedure was also 587
used to determine placements in the Uncorrelated (optimizing for a score close to 0) and 588
Negatively correlated (optimizing for a negative score) conditions. 589

Finally, datasets were created using 10 different placements in each of the three condi- 590
tions. Networks were trained for classification on these modified CIFAR-10 datasets in the 591
same manner as the VGG-16 network trained on the unperturbed version of the dataset 592
(See Figure 4). 593

Study 3 The procedure mirrored Study 2 with the main difference being that the target 594
system was the macaque inferior temporal cortex. Neural data from two macaques, as well 595
as the dataset were obtained from the Brain Score repository [26]. This dataset consists 596
of 3200 images from 8 categories (animals, boats, cars, chairs, faces, fruits, planes, and 597
tables), we computed an 8×8 averaged RDM based on macaque IT response patterns for 598
stimuli in each category. 599

This averaged RDM was then used as the target RDM in the optimization procedure to 600
determine locations of the confound (here, a white predictive patch of size 5×5 pixels) for 601
each category. Using a patch instead of a single pixel was required in this dataset because 602
of the structure and smaller size of the dataset (3200 images, rather than 50,000 images 603
for CIFAR-10). In this smaller dataset, the networks struggle to learn based on a single 604

pixel. However, increasing the size of the patch makes these patches more predictive 605
and the networks are able to again learn entirely based on this confound (see results 606
in Figure 6). In a manner similar to Study 2, this optimisation procedure was used 607
to construct three datasets, where the confound's placement was positively correlated, 608
uncorrelated or negatively correlated with the category distances in the target RDM. 609

Finally, each dataset was split into 75% training (2432 images) and 25% test sets 610
(768 images) before VGG-16 networks were trained on the unperturbed and modified 611
datasets in the same manner as in Study 2. One difference between Studies 2 and 3 612
was that here the networks in the Positive, Uncorrelated and Negative conditions were 613
trained from scratch, i.e., not pre-trained on ImageNet. This was done because we wanted 614
to make sure that the network in the Normal condition (trained on ImageNet) and the 615
networks in the Positive, Uncorrelated and Negative conditions encoded fundamentally 616
different features of their inputs – i.e., there were no ImageNet-related features encoded by 617
representations Φ_2 , Φ_3 and Φ_4 that were responsible for the similarity in representational 618
geometries between these representations and the representations in macaque IT cortex. 619

Study 4 The target system in this study was human IT cortex. The human RDM 620
and dataset were obtained from [9]. Rather than calculating pixel placements based on 621
the human RDM, the hierarchical structure of the dataset was used to place the pixels 622
manually. The dataset consists of 910 images from 6 categories: human bodies, human 623
faces, animal bodies, animal faces, artificial inanimate objects and natural inanimate 624
objects. These low-level categories can be organised into the hierarchical structure shown 625
in Figure 8. Predictive pixels were manually placed so that the distance between pixels 626
for Animate kinds were closer together than they were to Inanimate kinds and that faces 627
were closer together than bodies. This can be done in many different ways, so we created 628

five different datasets, with five possible arrangements of predictive pixels. Results in 629
the Hierarchical condition (Figure 9) are averaged over these five datasets. Placements 630
for the Random condition were done similarly, except that the locations were selected 631
randomly. 632

Networks were then trained on a 6-way classification task (818 training images and 92 633
test images) in a similar manner to the previous studies. As in Study 3, networks trained 634
on the modified datasets (both Hierarchical and Random conditions) were not pre-trained 635
on ImageNet. 636

RDM and RSA computation 637

For Studies 2-4 all image-level RDMs were calculated using $1 - r$ as the distance measure. 638
RSA scores were computed as the Spearman rank correlation between RDMs. 639

In Study 2, a curated set of test images was selected due to the extreme heterogeneity 640
of the CIFAR-10 dataset (low activation pattern similarity between instances of the same 641
category). This was done by selecting 5 images per category which maximally correlated 642
with the averaged activation pattern for the category. Since CIFAR-10 consists of 10 643
categories, the RSA-scores in Study 2 were computed using RDMs of size 50×50 . 644

In Study 3, the dataset consisted of 3200 images belonging to 8 categories. We first 645
calculated a full 3200×3200 RDM using the entire set of stimuli. An averaged, category- 646
level, 8×8 RDM was then calculated using median distances between categories (in 647
a manner similar to that described for Study 2 in the Section ‘Dataset generation and 648
training’). This 8×8 RDM was used to determine the RSA-scores. We also obtained 649
qualitatively similar results using the full 3200×3200 RDMs. These results can be found 650
in the Supplementary Information, Section B. 651

In Study 4, the dataset consisted of 818 training images and 92 test images. Kriegesko- 652

rte et al. [9] used these images to obtain a 92×92 RDM to compare representations between 653 human and macaque IT cortex. Here we computed a similar 92×92 RDM for networks 654 trained in the Normal, Hierarchical and Random training conditions, which were then 655 compared with the 92×92 RDM from human IT cortex to obtain RSA-scores for each 656 condition. 657

Testing 658

In Study 2, we used a 4×2 design to measure classification performance for networks in 659 all four conditions (Normal, Positive, Uncorrelated and Negative) on both unperturbed 660 images and modified images. We computed six RSA-scores: three pairs of networks – 661 Normal-Positive, Normal-Uncorrelated and Normal-Negative – and two types of inputs – 662 unperturbed and modified test images. The noise ceiling (grey band in Figure 6) was de- 663 termined in the standard way as described in [57] and represents the expected range of the 664 highest possible RSA score with the target system (network trained on the unperturbed 665 dataset). 666

In Study 3, performance was estimated in the same manner as in Study 2 (using a 667 4×2 design), but RSA-scores were computed between RDMs from macaque IT activations 668 and the four types of networks – i.e. for the pairs Macaque-Normal, Macaque-Positive, 669 Macaque-Uncorrelated and Macaque-Negative. And like in Study 2, we determined each 670 of these RSA-scores for both unperturbed and modified test images as inputs to the 671 networks. 672

In Study 4, performance and RSA were computed in the same manner as in Study 3, 673 except that the target RDM for RSA computation came from activations in human IT 674 cortex and the networks were trained in one of three conditions: Normal, Hierarchical 675 and Random. 676

Data analysis

677

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

678

679

680

681

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

682

683

684

685

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

686

687

Data Availability

688

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

689

690

691

Acknowledgments

692

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)

693

694

695

References

- [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).
- [2] O'Keefe, J. Place units in the hippocampus of the freely moving rat. *Experimental Neurology* **51**, 78–109 (1976).
- [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).
- [4] Haxby, J. V., Connolly, A. C. & Guntupalli, J. S. Decoding neural representational spaces using multivariate pattern analysis. *Annu Rev Neurosci* **37**, 435–456 (2014).
- [5] Ritchie, J. B., Kaplan, D. M. & Klein, C. Decoding the brain: Neural representation and the limits of multivariate pattern analysis in cognitive neuroscience. *The British Journal for the Philosophy of Science* **70**, 581–607 (2019). URL <https://doi.org/10.1093/bjps/axx023>. <https://doi.org/10.1093/bjps/axx023>.
- [6] Haxby, J. *et al.* A common, high-dimensional model of the representational space in human ventral temporal cortex. *Neuron* **72**, 404–416 (2011). URL <https://www.sciencedirect.com/science/article/pii/S0896627311007811>.
- [7] Kriegeskorte, N. & Kievit, R. A. Representational geometry: integrating cognition, computation, and the brain. *Trends in Cognitive Sciences* **17**, 401–412 (2013). URL <https://www.sciencedirect.com/science/article/pii/S1364661313001277>.
- [8] Kriegeskorte, N., Mur, M. & Bandettini, P. Representational similarity analysis - connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience* **2** (2008).

- [9] Kriegeskorte, N. *et al.* Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron* **60**, 1126–1141 (2008). 718
719
- [10] 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). 720
721
722
- [11] Michael L. Mack, B. L., Alison R. Preston. Decoding the brain’s algorithm for categorization from its neural implementation. *Current Biology* **23**, 2023–2027 (2013). 723
724
- [12] Freund, M. C., Etzel, J. A. & Braver, T. S. Neural coding of cognitive control: The representational similarity analysis approach. *Trends in Cognitive Sciences* **25**, 622–638 (2021). 725
726
727
- [13] Kaneshiro, B., Perreau Guimaraes, M., Kim, H.-S., Norcia, A. M. & Suppes, P. A representational similarity analysis of the dynamics of object processing using single-trial eeg classification. *PLOS ONE* **10**, 1–27 (2015). 728
729
730
- [14] Yamins, D. L. K. *et al.* Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences* **111**, 8619–8624 (2014). 731
732
733
- [15] Khaligh-Razavi, S.-M. & Kriegeskorte, N. Deep supervised, but not unsupervised, models may explain it cortical representation. *PLOS Computational Biology* **10**, 1–29 (2014). 734
735
736
- [16] Kietzmann, T. C. *et al.* Recurrence is required to capture the representational dynamics of the human visual system. *Proceedings of the National Academy of Sciences* **116**, 21854–21863 (2019). 737
738
739

- [17] Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A. & Oliva, A. Comparison of 740
deep neural networks to spatio-temporal cortical dynamics of human visual object 741
recognition reveals hierarchical correspondence. *Scientific Reports* **6**, 27755 (2016). 742
- [18] Kiat, J. E. *et al.* Linking patterns of infant eye movements to a neural network 743
model of the ventral stream using representational similarity analysis. *Developmental* 744
Science **25**, e13155 (2022). 745
- [19] Kriegeskorte, N. Deep neural networks: A new framework for modeling biological 746
vision and brain information processing. *Annual Review of Vision Science* **1**, 417–446 747
(2015). URL <https://doi.org/10.1146/annurev-vision-082114-035447>. PMID: 748
28532370, <https://doi.org/10.1146/annurev-vision-082114-035447>. 749
- [20] Xu, Y. & Vaziri-Pashkam, M. Limits to visual representational correspondence be- 750
tween convolutional neural networks and the human brain. *Nature Communications* 751
12, 2065 (2021). 752
- [21] Bobadilla-Suarez, S., Ahlheim, C., Mehrotra, A., Panos, A. & Love, B. C. Measures 753
of neural similarity. *Computational Brain & Behavior* **3**, 369–383 (2020). URL 754
<https://doi.org/10.1007/s42113-019-00068-5>. 755
- [22] Ramírez, F. M. Orientation encoding and viewpoint invariance in face recognition: 756
Inferring neural properties from large-scale signals. *The Neuroscientist* **24**, 582– 757
608 (2018). URL <https://doi.org/10.1177/1073858418769554>. PMID: 29855217, 758
<https://doi.org/10.1177/1073858418769554>. 759
- [23] Schölkopf, B. & Smola, F., A. J. and Bach. *Learning with kernels: support vector* 760
machines, regularization, optimization, and beyond (MIT Press, 2002). 761

- [24] 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).
- [25] Torralba, A. & Efros, A. A. Unbiased look at dataset bias. In *CVPR 2011*, 1521–1528 (2011).
- [26] Schrimpf, M. *et al.* Brain-score: Which artificial neural network for object recognition is most brain-like? *bioRxiv preprint: 407007* (2018).
- [27] 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).
- [28] Kriegeskorte, N. Relating population-code representations between man, monkey, and computational models. *Frontiers in Neuroscience* **3**, 363–373 (2009).
- [29] Bowers, J. S. *et al.* Deep problems with neural network models of human vision (2022). URL psyarxiv.com/5zf4s.
- [30] Serre, T. Deep learning: The good, the bad, and the ugly. *Annual Review of Vision Science* **5**, 399–426 (2019). URL <https://doi.org/10.1146/annurev-vision-091718-014951>. PMID: 31394043, <https://doi.org/10.1146/annurev-vision-091718-014951>.
- [31] Dujmović, M., Malhotra, G. & Bowers, J. S. What do adversarial images tell us about human vision? *eLife* **9**, e55978 (2020). URL <https://doi.org/10.7554/eLife.55978>.

- [32] Kriegeskorte, N. & Diedrichsen, J. Peeling the onion of brain representations. *Annual Review of Neuroscience* **42**, 407–432 (2019). 784
785
- [33] Kriegeskorte, N. & Wei, X.-X. Neural tuning and representational geometry. *Nature Reviews Neuroscience* **22**, 703–718 (2021). 786
787
- [34] Ramírez, F. M. Representational confusion: the plausible consequence of demeaning your data. *bioRxiv* (2017). URL <https://www.biorxiv.org/content/early/2017/09/28/195271>. <https://www.biorxiv.org/content/early/2017/09/28/195271.full.pdf>. 788
789
790
791
- [35] Henriksson, L., Khaligh-Razavi, S.-M., Kay, K. & Kriegeskorte, N. Visual representations are dominated by intrinsic fluctuations correlated between areas. *NeuroImage* **114**, 275–286 (2015). URL <https://www.sciencedirect.com/science/article/pii/S105381191500316X>. 792
793
794
795
- [36] Cai, M. B., Schuck, N. W., Pillow, J. W. & Niv, Y. Representational structure or task structure? bias in neural representational similarity analysis and a bayesian method for reducing bias. *PLOS Computational Biology* **15**, 1–30 (2019). URL <https://doi.org/10.1371/journal.pcbi.1006299>. 796
797
798
799
- [37] Mitchell, M. & Krakauer, D. C. The Debate Over Understanding in AI’s Large Language Models (2022). [2210.13966](https://arxiv.org/abs/2210.13966). 800
801
- [38] Geirhos, R. *et al.* Shortcut learning in deep neural networks. *Nature Machine Intelligence* **2**, 665–673 (2020). 802
803
- [39] Edelman, S. Representation is representation of similarities. *Behavioral and Brain Sciences* **21**, 449–467 (1998). 804
805

- [40] Block, N. Advertisement for a semantics for psychology. *Midwest Studies in Philosophy* **10**, 615–678 (1986). 806
807
- [41] Fodor, J. & Lepore, E. *Holism: A Shoppers Guide* (Blackwell, Cambridge, 1992). 808
- [42] Hempel, C. G. Problems and changes in the empiricist criterion of meaning. *Revue Internationale de Philosophie* **4**, 41–63 (1950). 809
810
- [43] Quine, W. V. Main trends in recent philosophy: Two dogmas of empiricism. *The Philosophical Review* **60**, 20–43 (1951). 811
812
- [44] Firth, J. R. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis* 813
(1957). 814
- [45] Griffiths, T. L. & Steyvers, M. A probabilistic approach to semantic representation. 815
In *Proceedings of the Twenty-Fourth Annual Conference of the Cognitive Science Society* (Erlbaum, Hillsdale, NJ, 2002). 816
817
- [46] Griffiths, T. L., Steyvers, M. & Tenenbaum, J. A probabilistic approach to semantic 818
representation. *Psychological Review* **114**, 211–244 (2007). 819
- [47] Piantadosi, S. & Hill, F. Meaning without reference in large language models. In 820
NeurIPS 2022 Workshop on Neuro Causal and Symbolic AI (nCSI) (2022). URL 821
<https://openreview.net/forum?id=nRkJEwmZnM>. 822
- [48] Fodor, J. *Psychosemantics: The Problem of Meaning in the Philosophy of Mind* 823
(MIT Press, Cambridge, 1987). 824
- [49] Haxby, J. V. Multivariate pattern analysis of fMRI: the early beginnings. *Neuroimage* 825
62, 852–855 (2012). 826

- [50] Morcos, A., Raghu, M. & Bengio, S. Insights on representational similar- 827
ity in neural networks with canonical correlation. In Bengio, S. *et al.* (eds.) 828
Advances in Neural Information Processing Systems, vol. 31 (Curran Asso- 829
ciates, Inc., 2018). URL [https://proceedings.neurips.cc/paper/2018/file/](https://proceedings.neurips.cc/paper/2018/file/a7a3d70c6d17a73140918996d03c014f-Paper.pdf) 830
[a7a3d70c6d17a73140918996d03c014f-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/a7a3d70c6d17a73140918996d03c014f-Paper.pdf). 831
- [51] Raghu, M., Gilmer, J., Yosinski, J. & Sohl-Dickstein, J. Svcca: Singular vector 832
canonical correlation analysis for deep learning dynamics and interpretability. In 833
Guyon, I. *et al.* (eds.) *Advances in Neural Information Processing Systems*, vol. 30 834
(Curran Associates, Inc., 2017). URL [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2017/file/dc6a7e655d7e5840e66733e9ee67cc69-Paper.pdf) 835
[2017/file/dc6a7e655d7e5840e66733e9ee67cc69-Paper.pdf](https://proceedings.neurips.cc/paper/2017/file/dc6a7e655d7e5840e66733e9ee67cc69-Paper.pdf). 836
- [52] Kornblith, S., Norouzi, M., Lee, H. & Hinton, G. Similarity of neural network rep- 837
resentations revisited. In *International Conference on Machine Learning*, 3519–3529 838
(PMLR, 2019). 839
- [53] Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning li- 840
brary. In *Advances in Neural Information Processing Systems 32*, 8024–8035 (Curran 841
Associates, Inc., 2019). 842
- [54] Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale 843
image recognition. *arXiv preprint arXiv:1409.1556* (2014). 844
- [55] Krizhevsky, A. Learning multiple layers of features from tiny images. Tech. Rep. 845
(2009). 846
- [56] Muttenthaler, L. & Hebart, M. N. Thingsvision: A python toolbox for streamlining 847
the extraction of activations from deep neural networks. *Frontiers in Neuroinformat-* 848
ics **15**, 679838 (2021). 849

- [57] Nili, H. *et al.* A toolbox for representational similarity analysis. *PLOS Computational Biology* **10**, 1–11 (2014). 850
851