

# Obstacles to inferring mechanistic similarity using Representational Similarity Analysis

Marin Dujmović<sup>1\*</sup>, Jeffrey S Bowers<sup>1</sup>, Federico Adolfi<sup>1,2</sup>, and Gaurav  
Malhotra<sup>1</sup>

<sup>1</sup>*School of Psychological Science, University of Bristol, Bristol, UK*

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

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

## Abstract

Representational Similarity Analysis (RSA) is an innovative approach used to compare neural representations across individuals, species and computational models. Despite its popularity within neuroscience, psychology and artificial intelligence, this approach has led to difficult-to-reconcile and contradictory findings, particularly when comparing primate visual representations with deep neural networks (DNNs). Here, we demonstrate how such contradictory findings could arise due to incorrect inferences about mechanism when comparing complex systems processing high-dimensional stimuli. In a series of studies comparing computational models, primate cortex and human cortex we find two problematic phenomena: 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, 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 is a multi-voxel pattern analysis method specifically designed to compare representations between different systems. 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 Fig 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* [4]. 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 25  
of the two activity patterns. One could compare, for example, fMRI signals with single 26  
cell recordings, EEG traces with behavioural data, or vectors in a computer algorithm 27  
with spiking activity of neurons [5]. RSA is now ubiquitous in computational psychology 28  
and neuroscience and has been applied to compare object representations in humans and 29  
primates [6], representations of visual scenes by different individuals [7,8], representations 30  
of visual scenes in different parts of the brain [9], to study specific processes such as cog- 31  
nitive control [10] or the dynamics of object processing [11], and most recently, to relate 32  
neuronal activations in human (and primate) visual cortex with activations of units in 33  
Deep Neural Networks [12–16]. 34

However, this flexibility in the application of RSA comes at the price of increased 35  
ambiguity in the inferences one can draw from this analysis. Since RSA is a second- 36  
order statistic (it looks at the similarity of similarities), it remains ambiguous which 37  
stimulus features drive the observed representational geometry in each system [17]. That 38  
is, two systems that operate on completely different stimulus features can nevertheless 39  
have highly correlated representational geometries. This makes inferences about mech- 40  
anism based on conducting RSA highly problematic. While researchers have recently 41  
highlighted a similar conceptual issue of confounds driving performance for multivariate 42  
decoding methods [18–20], it is less well appreciated for RSA. This is likely because there 43  
is a lack of understanding of how confounds can lead to misleading RSA scores, whether 44  
it is plausible that such confounds exist in datasets used for RSA and whether existing 45  
methods of dealing with confounds can address problematic inferences. 46

In our view, one particularly problematic area is research comparing biological systems 47  
and Deep Neural Networks (DNNs). There are many examples in this domain where 48  
researchers have recently used RSA for making inferences about psychological and neural 49

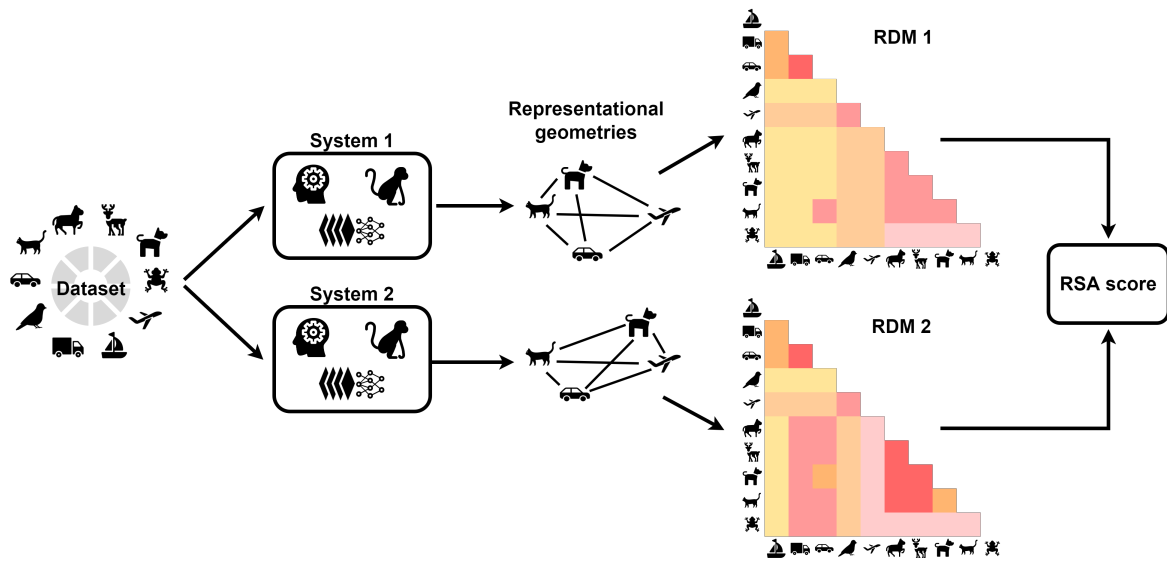


Figure 1: **RSA calculation.** Stimuli from a set of categories (or conditions) are used as inputs to two different systems (for example, a human brain and a primate brain). Activity from regions of interest is recorded for each stimulus. Pair-wise distances in activity patterns are calculated to get 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  $\rho$ , Kendall's  $\tau$ , etc.).

mechanisms. For example, Cichy et al. [15] observed a correspondence in the RDMs of DNNs performing object categorization and neural responses in human visual cortex recorded using MEG and fMRI. Based on this correspondence, the authors concluded that:

...hierarchical systems of visual representations emerge in both the human ventral and dorsal visual stream *as the result of* task constraints of object categorization posed in everyday life, and provide strong evidence for object representations in the dorsal stream independent of attention or motor intention. [pg. 5, emphasis added]

Thus, the correspondence in RDMs is used to infer the mechanism of emergence of visual representations. Based on a similar comparison, Kriegeskorte [21] concluded that <sup>1</sup>:

Deep convolutional feedforward networks for object recognition are not biologically detailed and rely on nonlinearities and learning algorithms that may differ from those of biological brains. Nevertheless they learn internal representations that are highly similar to representations in human and nonhuman primate IT cortex. [pg. 441]

In this paper, we will show – through a series of simulations – that such inferences about similarity in neural representations or mechanism based on RSA are unwarranted. We will show that when target RDMs are obtained using a highly complex model and high-dimensional stimuli (both of which are true for comparisons between biological systems

---

<sup>1</sup>It is not our intention to pick on or criticise these authors here. We have inserted these quotes to point out the broad understanding of this method in the field. Many authors are indeed careful in stating that the term ‘similarity in representations’ is used as a shorthand for a ‘similarity in representational geometries’. Nevertheless, readers are also invited to accept that different systems show similar representational geometries because it is likely that they also use similar mechanisms to transform sensory information into latent representations, or they use similar (downstream) mechanisms to decode these latent representations e.g. [22]. But how safe are these assumptions?

and DNNs) unknown confounds could drive similarity in representational geometries. 70  
Under these conditions, methods developed for dealing with confounds in multivariate 71  
decoding, such as cross-validation [23] and confound regression [24] may be insufficient 72  
for preventing false inferences. Furthermore, we will show that these confounds are not 73  
just possible, but also plausible given the nature of stimuli and structure of datasets 74  
frequently used for testing similarities between DNNs and humans. We will also argue 75  
that representational geometry should *not* be understood as the representation of a system 76  
as it conflicts with how most psychologists and neuroscientists view representations – the 77  
relationship between cognitive states and entities in the external world (see section S1 of 78  
the Supplementary Information a brief history of RSA and its philosophical origins). 79

The structure of the paper is as follows. In Study 1, in a bare-bones setup, we show 80  
that it is possible for two systems to transform input stimuli through known functions 81  
that are vastly different but end up with similar representational geometries. In particu- 82  
lar, the study shows that 1) the presence of second-order confounds in the training data 83  
can lead systems to mimic each other’s representational geometry even in the absence of 84  
mechanistic similarity, and 2) the intrinsic structure of datasets rather than mechanistic 85  
alignment can lead to artifactual modulation of RSA scores. Then in Studies 2 and 3 we 86  
use real neural data collected in previous experiments to show these problems extend to 87  
more complex datasets directly relevant to artificial intelligence and computational neu- 88  
rosience. Finally, in Study 4, we show that not only are misleadingly high RSA scores 89  
possible in practice but they are also highly plausible given the hierarchical structure of 90  
categories in datasets that are routinely used. Since our results have considerable general- 91  
ity with respect to current practices across multiple fields, we discuss the implications for 92  
published results, including a discussion of two alternative philosophical perspectives on 93  
the nature of mental representations that our findings speak to. We conclude by provid- 94

ing some general recommendations regarding how to best compare representations across 95  
systems going forward. 96

## Results 97

### Proof of concept 98

It may be tempting to infer that two systems which have similar representational geome- 99  
tries for a set of concepts do so because they encode similar properties of sensory data and 100  
transform sensory data through a similar set of functions. In this section, we show that 101  
it is possible, at least in principle, for qualitatively different systems to end up with very 102  
similar representational geometries even though they (i) transform their inputs through 103  
very different functions, and (ii) select different features of inputs. 104

**Study 1: Demonstrably different transformations of inputs can lead to low 105  
or high RSA-scores** We start by considering a simple two-dimensional dataset and 106  
two systems where we know the closed-form functions that project this data into two 107  
representational spaces. This simple setup helps us gain a theoretical understanding of 108  
the circumstances under which it is possible for qualitatively different projections to show 109  
similar representational geometries. 110

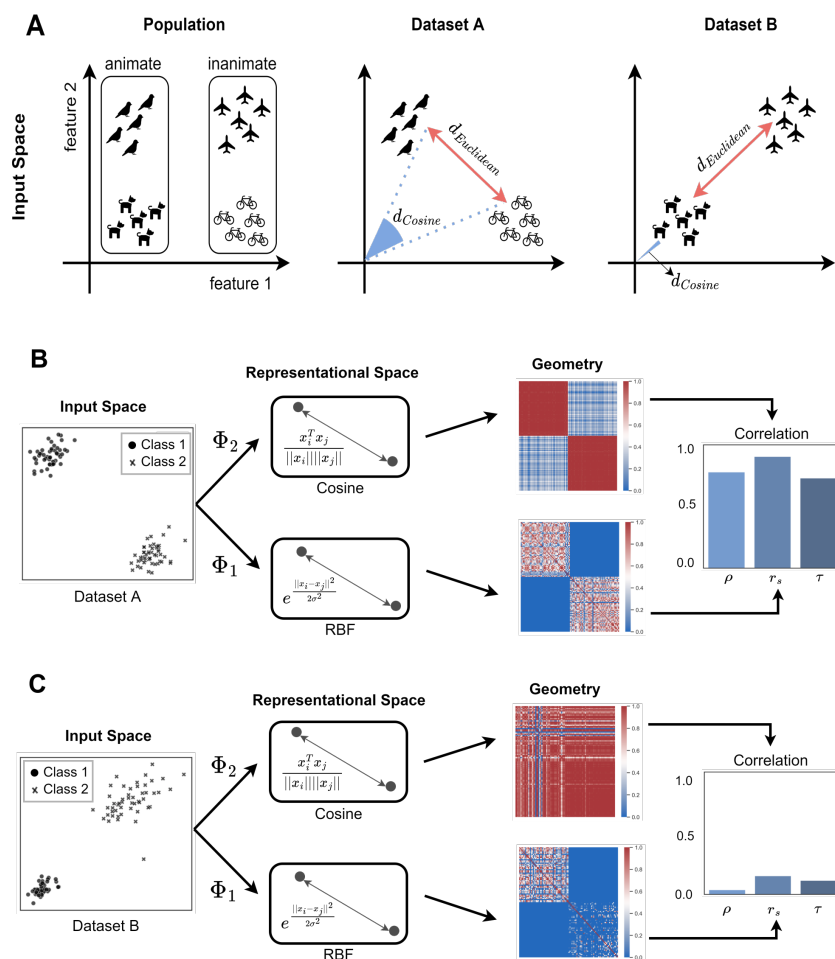
Consider a population of animate and inanimate objects that consist of four categories 111  
of objects – birds, dogs, airplanes and bicycles. Each object in this population will have a 112  
set of stimulus features, using which one can map each exemplar from all four categories 113  
into a feature space. In Fig 2A (left), we show a hypothetical 2D feature space where 114  
exemplars from each category cluster together. Furthermore, we consider two datasets 115  
sampled from this population – Dataset A (Fig 2A, middle) which consists of birds and 116

bicycles and Dataset B (Fig 2A, right) which consists of dogs and airplanes. Both datasets 117  
consist of animate and inanimate objects, but they differ in how items in each category 118  
are represented in the input space. 119

Now, consider two information-processing systems that re-represent Dataset A into 120  
two different latent spaces (Fig 2B). These could be two recognition systems designed to 121  
distinguish animate and inanimate categories. We assume that we can observe the repre- 122  
sentational geometry of the latent representations of each system and we are interested in 123  
understanding whether observing a strong correlation between these geometries implies 124  
whether the two systems have a similar *representational space* – that is, they project in- 125  
puts into the latent space using similar functions. To examine this question, we consider a 126  
setup where we know the functions,  $\Phi_1$  and  $\Phi_2$ , that map the inputs to the latent space in 127  
each system. We will now demonstrate that even when these functions are qualitatively 128  
different from each other, the geometry of latent representations can nevertheless be highly 129  
correlated. We will also show that the difference in representational spaces becomes more 130  
clear when one considers a different dataset (Dataset B), where inputs projected using 131  
the same functions now lead to a low correlation in representational geometries. 132

We can compute the geometry of a set of representations by establishing the pair-wise 133  
distance between all vectors in each representational space  $\Phi$ . There are many different 134  
methods of computing this representational distance between any pair of vectors, all de- 135  
riving from the dot product between vectors (see, for example, Fig 1 in [25]). Previous 136  
research has shown that the choice of the distance metric itself can influence the inferences 137  
one can draw from one’s analysis [25,26]. However, here our focus is not the distance met- 138  
ric itself, but the fundamental nature of RSA. Therefore, we use the same generic distance 139  
metric – the dot product – to compute the pair-wise distance between all vectors in both 140  
representational spaces. In other words, the representational distance  $d[\Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j)]$ , 141





**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 ( $\rho$ ), Spearman ( $r_s$ ) and Kendall's ( $\tau$ ) 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*.

between the projections of any pair of input stimuli,  $\mathbf{x}_i$  and  $\mathbf{x}_j$  into a feature space  $\Phi$ , is 142  
proportional to the inner product between the projections in the feature space: 143

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

And we can obtain the representational geometry of the input stimuli  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  in 144  
any representational space  $\Phi$  by computing the pairwise distances,  $d[\Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j)]$  for all 145  
pairs of data points,  $(i, j)$ . Here, we assume that the projections  $\Phi_1$  and  $\Phi_2$  are such that 146  
these pairwise distances are given by two positive semi-definite kernel functions  $\kappa_1(\mathbf{x}_i, \mathbf{x}_j)$  147  
and  $\kappa_2(\mathbf{x}_i, \mathbf{x}_j)$ , respectively: 148

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

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

Now, let us consider two qualitatively different kernel functions:  $\kappa_1(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$  149  
is a radial-basis kernel (where  $\sigma^2$  is the bandwidth parameter of the kernel), while  $\kappa_2(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$  150  
is a cosine kernel. We have chosen  $\Phi_1$  and  $\Phi_2$  such that they are two fundamen- 151  
tally different projections of the inputs  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  – while  $\Phi_2$  maps a 2D input  $\mathbf{x}_i$  into 152  
a 2D feature space,  $\Phi_1$  maps the same 2D input into an infinite-dimensional space. Nev- 153  
ertheless, since cosine and RBF kernels are Mercer kernels, we can compute the distances 154  
(as measured by the dot product) between each pair of projected vectors using the kernel 155  
trick [27, 28]. That is, we can find the distance between any pair of points in the represen- 156  
tational space by applying the kernel function to those points in the input space. These 157

pairwise distances are shown by the kernel matrices in Fig 2B. 158

Next, we can determine how the geometry of these projections in the two systems 159  
relate to each other by computing the correlation between the kernel matrices, shown on 160  
the right-hand-side of Fig 2B. We can see from these results that the kernel matrices are 161  
highly correlated – i.e., the input stimuli are projected to very similar geometries in the 162  
two representational spaces. 163

If one did not know the input transformations and simply observed the correlation 164  
between kernel matrices, it would be tempting to infer that the two systems  $\Phi_1$  and  $\Phi_2$  165  
transform an unknown input stimulus  $\mathbf{x}$  through a similar set of functions – for example 166  
functions that belong to the same class or project inputs to similar representational spaces. 167  
However, this would be an error. The projections  $\Phi_1(\mathbf{x})$  and  $\Phi_2(\mathbf{x})$  are fundamentally 168  
different –  $\Phi_1$  (radial basis kernel) projects an input vector into an infinite dimensional 169  
space, while  $\Phi_2$  (cosine kernel) projects it onto a unit sphere. The difference between these 170  
functions becomes apparent if one considers how this correlation changes if one considers a 171  
different set of input stimuli. For example, the set of data points from Dataset B (sampled 172  
from the same population) are projected to very different geometries, leading to a low 173  
correlation between the two kernel matrices (Fig 2C). 174

In fact, the reason for highly correlated kernel matrices in Fig 2B is not a similarity 175  
in the transformations  $\Phi_1$  and  $\Phi_2$  but the structure of the dataset. The representational 176  
distance between any two points in the first representational space,  $d[\Phi_1(\mathbf{x}_i), \Phi_1(\mathbf{x}_j)]$ , is 177  
 $e^{\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$ . That is, the representational distance in  $\Phi_1$  is a function of their Euclidean 178  
distance  $\|\mathbf{x}_i - \mathbf{x}_j\|$  in the input space. On the other hand, the representational dis- 179  
tance between any two points in the second representational space,  $d[\Phi_2(\mathbf{x}_i), \Phi_2(\mathbf{x}_j)]$ , is, 180  
 $\frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$ . That is, the representational distance in  $\Phi_2$  is a function of their cosine 181  
distance in the input space. These two stimulus features – Euclidean distance and cosine 182

distance – are *confounds* that lead to the same representational geometries for certain 183  
datasets. In Dataset A, the stimuli is clustered such that the Euclidean distance between 184  
any two stimuli is correlated with their cosine distance (see Fig 2A, middle). However, 185  
for Dataset B, the Euclidean distance is no longer correlated with the angle (see Fig 2A, 186  
right) and the confounds lead to different representational geometries, as can be seen in 187  
Fig 2C. Thus, this example illustrates two effects: (i) a *mimic* effect, where two systems 188  
that transform sensory input through very different functions end up with similar repre- 189  
sentational geometries (Fig 2B), and (ii) a *modulation* effect, where two systems that are 190  
non-identical have similar representational geometries for one set of inputs, but dissimilar 191  
geometries for a second set (compare Fig 2B and 2C). 192

**Study 2: Complex systems encoding different features of inputs can show a** 193  
**high RSA-score** Study 1 made a number of simplifying assumptions – the dataset was 194  
two-dimensional, clustered into two categories and we intentionally chose functions  $\Phi_1$  195  
and  $\Phi_2$  such that the kernel matrices were correlated in one case and not correlated in the 196  
other. It could be argued that, even though the above results hold in principle, they are 197  
unlikely in practice when the transformations and data structure are more complex. For 198  
example, it might be tempting to assume that accidental similarity in representational 199  
geometries becomes less likely as one increases the number of categories (i.e., clusters 200  
or conditions) being considered. However, In Fig 3 we illustrate how complex systems 201  
transforming high-dimensional input from a number of categories may achieve high RSA 202  
scores. Even though one system extracts surface reflectance and the other extracts global 203  
shape, they can end up with very similar representational geometries. This would occur 204  
if objects similar in their reflectance properties were also similar in shape (e.g., glossy 205  
balloons and light bulbs) and if objects dissimilar according to reflectance properties were 206

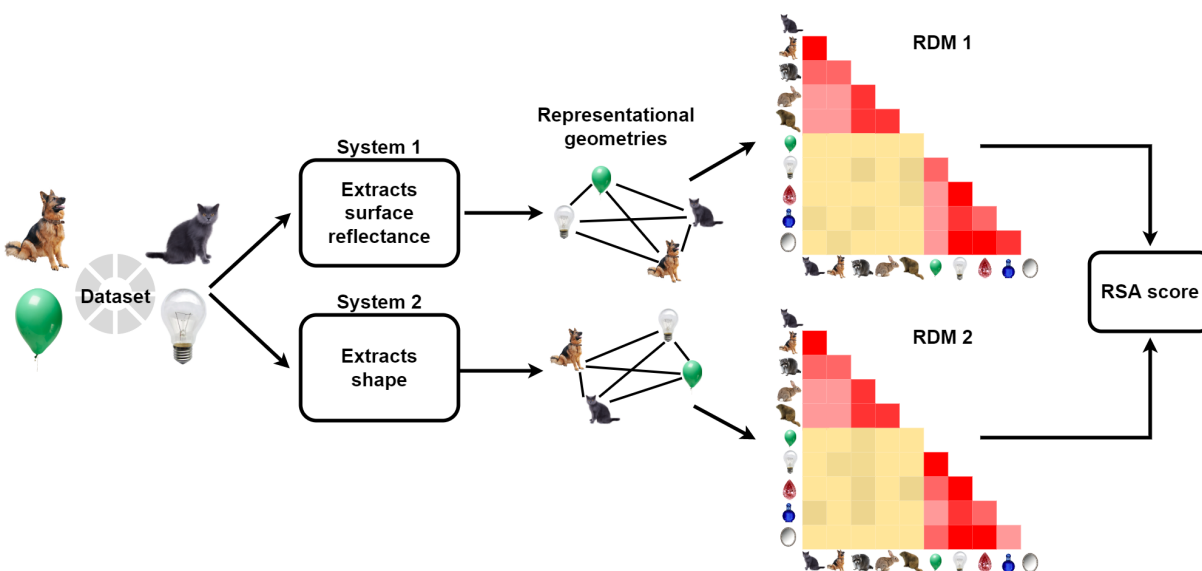


Figure 3: **Example of a second-order confound.** Two systems, one forming representations based on surface reflectance of objects (while ignoring all other features such as colour or texture) and the other based on global shape (while ignoring other features), can have very similar representational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

also dissimilar in shape (e.g., dogs and light bulbs). This is the mimic effect, where  
representational geometries of these two systems end up being similar because reflectance  
and shape are second-order confounds in this dataset. Conducting RSA on this dataset  
will show a high correlation in RDMs, even though the latent representations in these  
systems are related to very different stimulus features.

To demonstrate this empirically, we now consider a more complex setup, where the  
transformations  $\Phi_1$  and  $\Phi_2$  are modelled as feedforward deep neural networks (DNNs),  
trained to classify a high-dimensional dataset into multiple categories. Many studies that  
use RSA compare systems using naturalistic images as visual inputs [6, 12]. While using  
naturalistic images brings research closer to the real-world, it is also well-known that  
datasets of naturalistic images frequently contain confounds – independent features that  
can predict image categories [29]. We will now show how the simplest of such confounds,

a single pixel, can lead to a high RSA score between two DNNs that encode qualitatively different features of inputs.

Consider the same setup as above, where an input stimulus,  $\mathbf{x}$ , is transformed to a representation space by two systems,  $\Phi_1$  and  $\Phi_2$ . Instead of a two-dimensional input space,  $\mathbf{x}$  now exists in a high-dimensional image space and  $\Phi_1$  and  $\Phi_2$  are two versions of a DNN – VGG-16 – trained to classify input images into different categories. We ensured that  $\Phi_1$  and  $\Phi_2$  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 Fig 4 for the general approach).

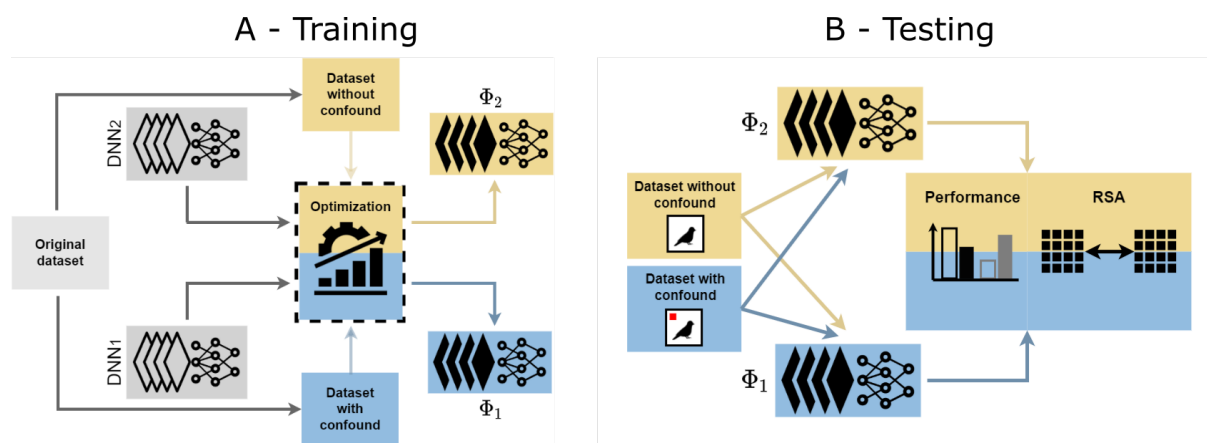


Figure 4: **Training and testing DNNs with different feature encodings.** Panel A shows the training procedure for Studies 2–4, where we created two versions of the original dataset (gray), one containing a confound (blue) and the other left unperturbed (yellow). These two datasets were used to train two networks (gray) on a categorisation task, resulting in two networks that learn to categorise images either based on the confound (projection  $\Phi_2$ ) or based on statistical properties of the unperturbed image (projection  $\Phi_1$ ). 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.

The locations of these diagnostic pixels were chosen such that they were correlated to 230  
the corresponding representational distances between classes in  $\Phi_1$ . Our hypothesis was 231  
that if the representational distances in  $\Phi_2$  preserve the physical distances of diagnos- 232  
tic pixels in input space, then this confound will end up mimicking the representational 233  
geometry of  $\Phi_1$ , even though the two systems use qualitatively different features for clas- 234  
sification. Furthermore, we trained two more networks,  $\Phi_3$  and  $\Phi_4$ , which were identical 235  
to  $\Phi_2$ , except these networks were trained on datasets where the location of the confound 236  
was uncorrelated ( $\Phi_3$ ) or negatively correlated ( $\Phi_4$ ) with the representational distances 237  
in  $\Phi_1$  (see Fig 5 and Methods for details). 238

The locations of these diagnostic pixels were chosen such that they were correlated to 239  
the corresponding representational distances between classes in  $\Phi_1$ . Our hypothesis was 240  
that if the representational distances in  $\Phi_2$  preserve the physical distances of diagnos- 241  
tic pixels in input space, then this confound will end up mimicking the representational 242  
geometry of  $\Phi_1$ , even though the two systems use qualitatively different features for clas- 243  
sification. Furthermore, we trained two more networks,  $\Phi_3$  and  $\Phi_4$ , which were identical 244  
to  $\Phi_2$ , except these networks were trained on datasets where the location of the confound 245  
was uncorrelated ( $\Phi_3$ ) or negatively correlated ( $\Phi_4$ ) with the representational distances 246  
in  $\Phi_1$  (see Fig 5 and Methods for details). 247

Classification accuracy (Fig 6 (left)) revealed that the network  $\Phi_1$ , trained on the 248  
unperturbed images, learned to classify these images and ignored the diagnostic pixel 249  
– that is, it’s performance was identical for the unperturbed and modified images. In 250  
contrast, networks  $\Phi_2$  (positive),  $\Phi_3$  (uncorrelated) and  $\Phi_4$ (negative) failed to classify the 251  
unperturbed images (performance was near chance) but learned to perfectly classify the 252  
modified images, showing that these networks develop qualitatively different representa- 253  
tions compared to normally trained networks. 254

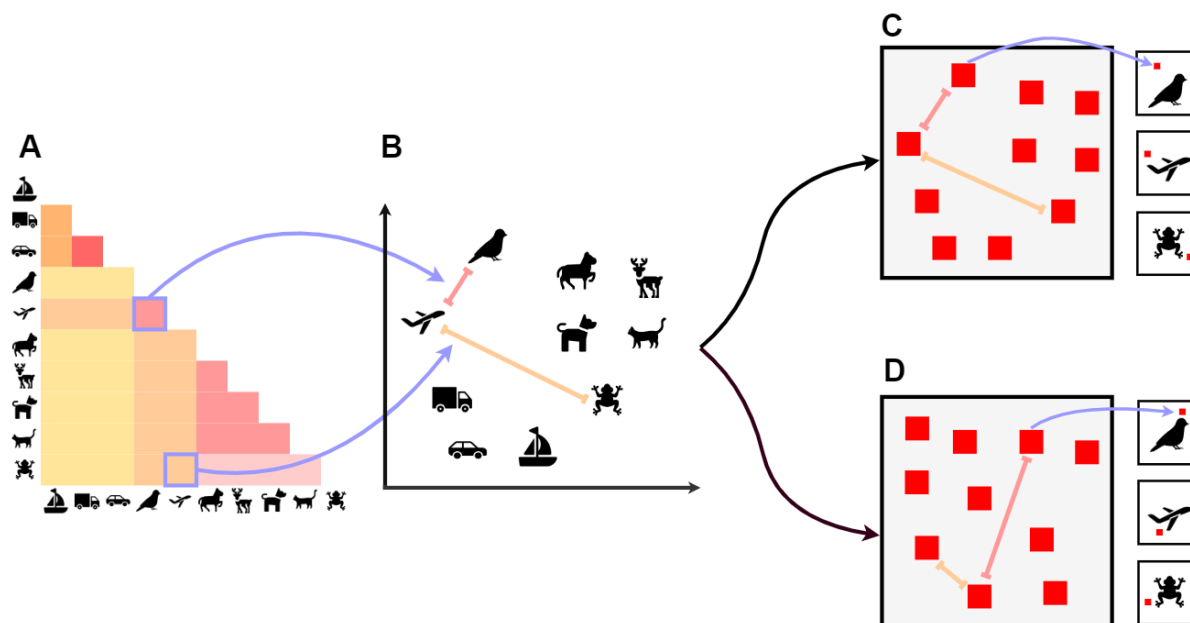


Figure 5: **Study 2 confound placement.** The representational geometry (Panel A and B) from the network trained on the unperturbed CIFAR-10 images is used to determine the location of the single pixel confound (shown as a red patch here) for each category. In the ‘Positive’ condition (Panel C), we determined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

Next we computed pairwise RSA scores between the representations at the last convo- 255  
lution layer of  $\Phi_1$  and each of  $\Phi_2$ ,  $\Phi_3$  and  $\Phi_4$  (Fig 6 (right)). When presented unperturbed 256  
test images, the  $\Phi_2$ ,  $\Phi_3$  and  $\Phi_4$  networks all showed low RSA scores with the normally 257  
trained  $\Phi_1$  network. However, when networks were presented with test images that in- 258  
cluded the predictive pixels, RSA varied depending on the geometry of pixel locations 259  
in the input space. When the geometry of pixel locations was positively correlated to 260  
the normally trained network, RSA scores approached ceiling (i.e., comparable to RSA 261  
scores between two normally trained networks). Networks trained on uncorrelated and 262



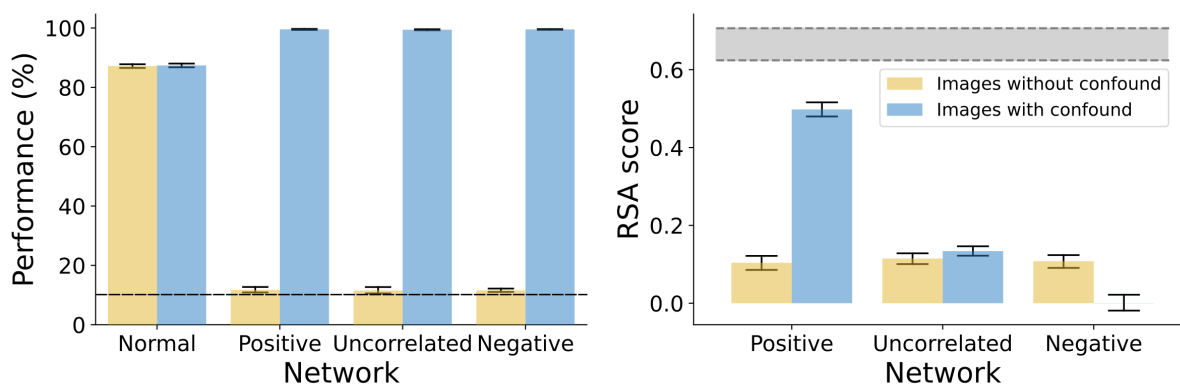


Figure 6: **Study 2 results.** *Left:* Performance of normally trained networks did not depend on whether classification was done on unperturbed CIFAR-10 images or images with a single pixel confound (error bars represent 95% CI, the dashed line represents chance performance). All three networks trained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

negatively correlated pixel placements scored much lower. 263

These results mirror Study 1: we observed that it is possible for two networks ( $\Phi_1$  and 264  
 $\Phi_2$ ) to show highly correlated representational geometries even though these networks 265  
learn to classify images based on very different features. One may argue that this could 266  
be because the two networks could have learned similar representations at the final con- 267  
volution layer of the DNN and it is the classifier that sits on top of this representation 268  
that leads to the behavioural differences between these networks. But if this was true, it 269  
would not explain why RSA scores diminish for the two other comparisons (with  $\Phi_3$  and 270  
 $\Phi_4$ ). This modulation of RSA-scores for different datasets suggests that, like in Study 1, 271  
the correlation in representational geometry is not because the two systems encode similar 272  
features of inputs, but because different features mimic each other in their representational 273  
geometries. 274

## Re-examining some influential findings 275

In Studies 1 and 2, we showed that it is possible for qualitatively different systems to 276  
end up with similar representational geometries. However, it may be argued that while 277  
this is possible in principle, it is unlikely in practice in real-world scenarios. In the fol- 278  
lowing two studies, we consider real-world data from some recent influential experiments, 279  
recorded from both primate and human cortex. We show how RSA-scores can be driven 280  
by confounds in these real-world settings and how properties of training and test data 281  
may contribute to observed RSA-scores. 282

**Study 3: Neural activations in monkey IT cortex can show a high RSA-score 283  
with DNNs despite different encoding of input data** In our next study, we con- 284  
sider data from experiments comparing representational geometries between computa- 285

tional models and macaque visual cortex [12, 30]. The experimental setup was similar 286  
to Study 2, though note that unlike Study 2, where both systems used the same archi- 287  
tecture and learning algorithm, this study considered two very different systems – one 288  
artificial (DNN) and the other biological (macaque IT cortex). We used the same set of 289  
images that were shown to macaques by Majaaj et al. [31] and modified this dataset to 290  
superimpose a small diagnostic patch on each image. In the same manner as in Study 2 291  
above, we constructed three different datasets, where the locations of these diagnostic 292  
patches were either positively correlated, uncorrelated or negatively correlated with the 293  
RDM of macaque activations. We then trained four CNNs. The first CNN was pre- 294  
trained on ImageNet and then fine-tuned on the unmodified dataset of images shown to 295  
the macaques. Previous research has shown that CNNs trained in this manner develop 296  
representations that mirror the representational geometry of neurons in primate inferior 297  
temporal (IT) cortex [12]. The other three networks were trained on the three modi- 298  
fied datasets and learned to entirely rely on the diagnostic patches (accuracy on images 299  
without the diagnostic patches was around chance). 300

Fig 7 (right) shows the correlation in representational geometry between the macaque 301  
IT activations and activations at the final convolution layer for each of these networks. 302  
The correlation with networks trained on the unmodified images is our baseline and shown 303  
as the gray band in Fig 7. Our first observation was that a CNN trained to rely on the di- 304  
agnostic patch can indeed achieve a high RSA score with macaque IT activations. In fact, 305  
the networks trained on patch locations that were positively correlated to the macaque 306  
RDM matched the RSA score of the CNNs trained on ImageNet and the unmodified 307  
dataset. This shows how two systems having very different architectures, encoding fun- 308  
damentally different features of inputs (single patch vs naturalistic features) can show a 309  
high correspondence in their representational geometries. We also observed that, like in 310

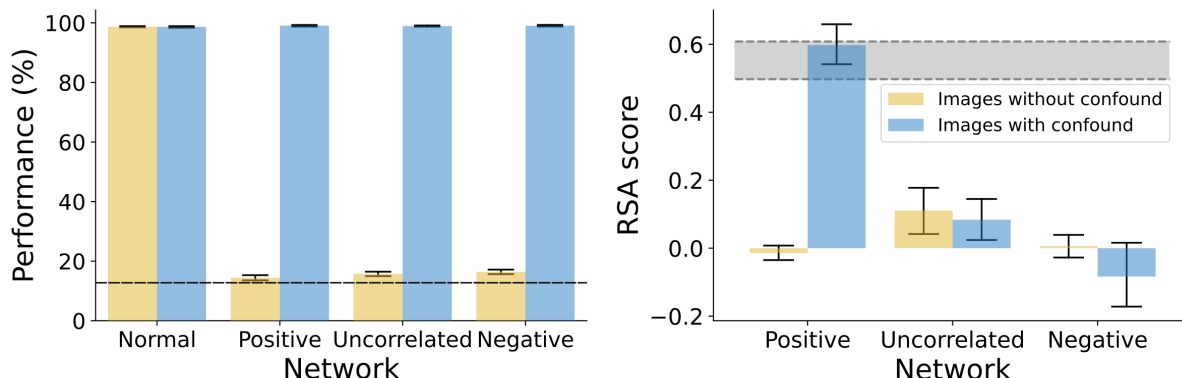


Figure 7: **Study 3 results.** *Left:* Classification Performance of the network trained on unperturbed images (Normal condition) did not depend on the presence or absence of the confound, while performance of networks trained with the confound (Positive, Uncorrelated and Negative conditions) highly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

Study 2, the RSA score depended on the clustering of data in the input space – when 311  
 patches were placed in other locations (uncorrelated or negatively correlated to macaque 312  
 RDMS) the RSA score became significantly lower. 313

**Study 4: High RSA-scores may be driven by the structure of testing data** All 314  
 the studies so far have used the same method to construct datasets with confounds – we 315  
 established the representational geometry of one system ( $\Phi_1$ ) and constructed datasets 316  
 where the clustering of features (pixels) mirrored this geometry. However, it could be 317  
 argued that confounds which cluster in this manner are unlikely in practice. For example, 318  
 even if texture and shape exist as confounds in a dataset, the inter-category distances 319  
 between textures are not necessarily similar to the inter-category distances between shape. 320

However, categories in real-world datasets are usually hierarchically clustered into 321

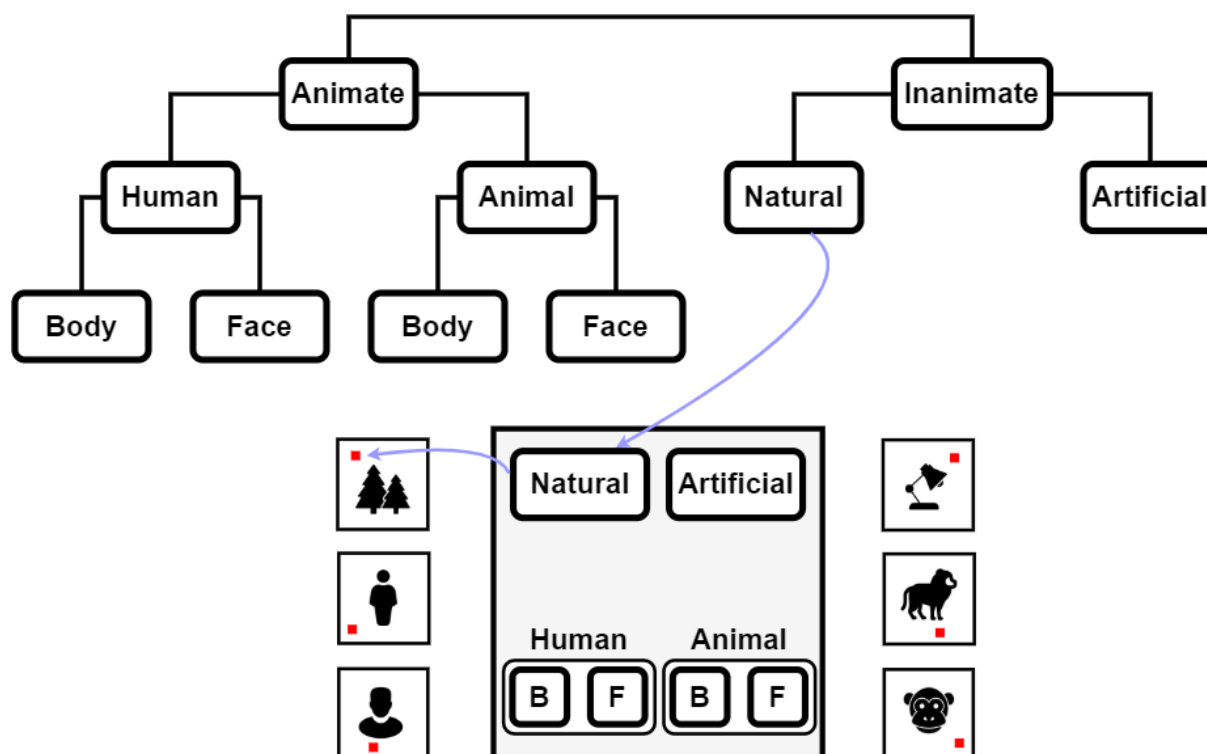


Figure 8: **Exploiting intrinsic dataset hierarchy in order to place confounds.** The top panel shows the hierarchical structure of categories in the dataset, which was used to place the single pixel confounds. The example at the bottom (middle) shows one such hierarchical placement scheme where the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

higher-level and lower-level categories. For example, in the CIFAR-10 dataset, the Dogs 322  
and Cats (lower-level categories) are both animate (members of a common higher-level 323  
category) and Airplanes and Ships (lower-level categories) are both inanimate (members 324  
of a higher-level category). Due to this hierarchical structure, Dog and Cat images are 325  
likely to be closer to each other not only in their shape, but also their colour and texture 326  
(amongst other features) than they are to Airplane and Ship images. In our next simula- 327  
tion, we explore whether this hierarchical structure of categories can lead to a correlation 328  
in representational geometries between two systems that learn different feature encodings. 329

For this study, we selected a popular dataset used for comparing representational 330 geometries in humans, macaques and deep learning models [13, 32]. This dataset consists 331 of six categories which can be organised into a hierarchical structure shown in Fig 8. [6] 332 showed a striking match in RDMs for response patterns elicited by these stimuli in human 333 and macaque IT. For both humans and macaques, distances in response patterns were 334 larger between the higher-level categories (animate and inanimate) than between the 335 lower-level categories (e.g., between human bodies and human faces). 336

We used a similar experimental paradigm to the above studies, where we trained 337 networks to classify stimuli which included a single predictive pixel. But instead of using 338 an RDM to compute the location of a diagnostic pixel, we used the hierarchical categorical 339 structure. In the first modified version of the dataset, the location of the pixel was based 340 on the hierarchical structure of categories in Fig 8 – predictive pixels for animate kinds 341 were closer to each other than to inanimate kinds, and pixels for faces were closer to 342 each other than to bodies, etc. One such configuration can be seen in Fig 8. In the 343 second version, the predictive pixel was placed at a random location for each category 344 (but, of course, at the same location for all images within each category). We call these 345 conditions ‘Hierarchical’ and ‘Random’. [13] showed that the RDM of average response 346 patterns elicited in the human IT cortex ( $\Phi_1$ ) correlated with the RDM of a DNN trained 347 on naturalistic images ( $\Phi_2$ ). We explored how this compared to the correlation with the 348 RDM of a network trained on the Hierarchical pixel placement ( $\Phi_3$ ) and Random pixel 349 placement ( $\Phi_4$ ). 350

Results for this study are shown in Fig 9. We observed that representational geometry 351 of a network trained on Hierarchically placed pixels ( $\Phi_3$ ) was just as correlated to the rep- 352 resentational geometry of human IT responses ( $\Phi_1$ ) as a network trained on naturalistic 353 images ( $\Phi_2$ ). However, when the pixel locations for each category were randomly cho- 354

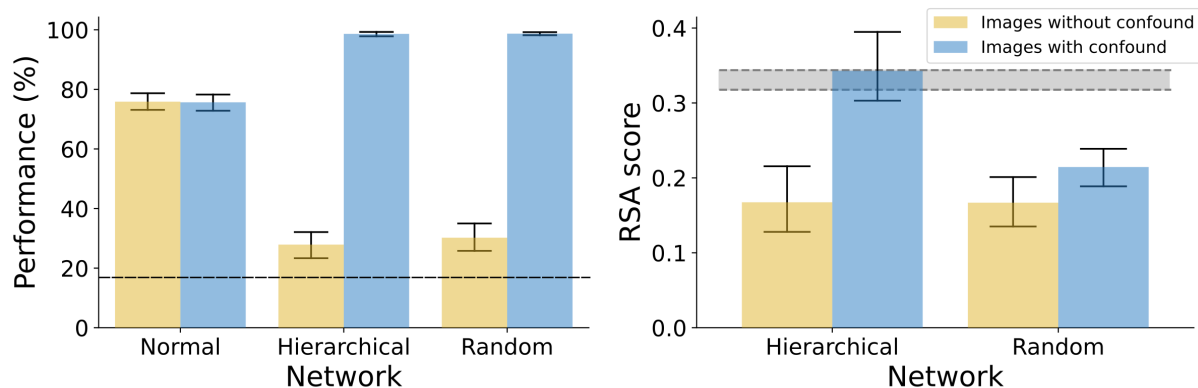


Figure 9: **Study 4 results.** *Left:* Performance of normally trained networks did not depend on whether the confound was present. Networks trained with the confound failed to classify stimuli without the confound (yellow bars) while achieving near perfect classification of stimuli with the confound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

sen, this correlation decreased significantly. These results suggest that any confound in 355  
the dataset (including texture, colour or low-level visual information) that has distances 356  
governed by the hierarchical clustering structure of the data could underlie the observed 357  
similarity in representational geometries between CNNs and human IT. More generally, 358  
these results show how it is plausible that many confounds present in popular datasets 359  
may underlie the observed similarity in representational geometries between two systems. 360  
The error of inferring a similarity in mechanism based on a high RSA score is not just 361  
possible but also probable. 362

## Discussion 363

In four studies, we have illustrated a number of conditions under which it can be problem- 364  
atic to infer a similarity of representations between two systems based on a correlation in 365

their representational geometries. In particular, we showed that two systems may transform their inputs through very different functions and encode very different features of inputs and yet have highly correlated representational geometries. Of course, one may acknowledge that a second-order isomorphism of activity patterns does *not* strictly imply that two systems are similar mechanistically but still assume that it is highly likely to be the case. That is, as a practical matter, a researcher may assume that RSA is a reliable method to compare complex systems representing high-dimensional inputs. However, our findings challenge this assumption. We show how a high RSA score between different systems can not only occur in a bare-bones simulation (Study 1), but also in practice, in high-dimensional systems operating on high-dimensional data (Studies 2–3). Furthermore, we show that the hierarchical structure of datasets frequently used to test similarity of representations lends itself to artificially inflated RSA scores because of second-order confounds present in the dataset (Study 4). Therefore, second-order confounds driving high RSA scores is not only possible but plausible.

## Significance and Implications

These results are critical for any researchers interested in using RSA for comparing mechanistic processes of complex systems. This includes research comparing information processing across species [6], across brain areas [9], between computational models [33] and between artificial intelligence models and brains [12–16]. Our findings are particularly relevant to an ongoing debate about similarity of information processing in Deep Neural Networks and the mammalian visual cortex. Studies comparing Convolutional Neural Networks (CNNs) to the visual cortex present a set of contradictory findings. In some studies, researchers have observed a correlation between the geometry of activation patterns in a CNN trained to classify large datasets of images and the geometry of neural



activation patterns within the macaque or human inferotemporal cortex, when both systems are stimulated using a set of naturalistic images [13]. Based on these observations of similarity (in RDMs), many researchers infer that CNNs provide insight into the mechanisms of information processing in the visual cortex [14, 15]. But other researchers have challenged this claim arguing that there are too many differences – from architectures and algorithms to tasks and environments used for learning – to meaningfully compare these systems, and indeed many behavioural results provide striking contrast between the performance of humans and CNNs [34–36]. Our results show how it is possible (and plausible) for these contradictory results to arise. That is, it is possible that CNNs are extracting very different features of naturalistic images and may even be transforming these images through very different functions compared to biological visual systems and yet end up with correlated representational geometries due to confounds present in datasets and their structural properties. It is important to emphasize that confounds are ubiquitous in datasets [29] leading DNNs to often classify images on the basis of short-cuts [37] and it unclear why confounds would not also drive high RSA scores. Furthermore, our results are consistent with recent results such as Xu and Vaziri-Pashkam [38] who have reported that the correlation between representational geometries (previously reported in many studies) depends on the dataset used for testing. Still, whether this is actually happening will require finding the actual confound driving these effects, a non-trivial task in such high-dimensional data (see below under ‘Limitations’).

We would also like to emphasise that our results are not an indictment of RSA per se. Rather, our critique is aimed at problematic inferences that are frequently being drawn based on this statistical method when comparing complex systems. Like many statistical methods, Representation Similarity Analysis provides insight into data at a certain level of abstraction. There are many fruitful ways to use RSA, particularly by constructing theory-

based representation dissimilarity matrices and comparing observed RDMs with these 415  
theory-driven RDMs. An example of this approach is outlined by Naselaris & Kay [18], 416  
who discuss how RDMs can be constructed based on hypotheses (such as whether the 417  
luminance of an image is driving observed differences between conditions) and compared 418  
with the RDMs of observed data. Used in this way, RSA is a useful tool for model- 419  
comparison, with different target RDMs corresponding to clearly defined hypotheses. The 420  
problem arises when the target RDM is generated by a complex system processing high- 421  
dimensional stimuli. In this setting, our results show that alternative hypotheses about 422  
mechanism are difficult to disambiguate based on RSA. 423

## Philosophical implications 424

But couldn't a researcher take the view that representational geometry *is* representation 425  
and therefore, a strong correlation in representational geometries between two systems is 426  
sufficient to infer that the systems are representing the world in a similar manner? This 427  
question goes to the heart of an existing debate in philosophy, where philosophers distin- 428  
guish between the *externalist* and *holistic* views on mental representations [39]. According 429  
to the first view, the content of representations is determined by their relationship to enti- 430  
ties in the external world. This perspective is implicitly taken by most neuroscientists and 431  
psychologists, who are interested in comparing mechanisms underlying cognitive processes 432  
– that is, they are interested in the set of nested functions and algorithms responsible for 433  
transforming sensory input into a set of activations in the brain. From this perspective, 434  
our finding that high RSA scores can be obtained between systems that work in quali- 435  
tatively different ways poses a challenge to researchers using RSA to compare complex 436  
systems using high-dimensional data where a multitude of feature transformations could 437  
be driving the observed results. 438

Alternatively, a researcher may reject the externalist view and adopt the perspective 439  
that representations obtain their meaning based on how they are related to each other 440  
within each system, rather than based on their relationship to entities in the external 441  
world. That is, “representation *is* the representation of similarities” [40]. From this per- 442  
spective, as long as the two systems share the same relational distances between internal 443  
activations, one can validly infer that the two systems have similar representations. That 444  
is, a second-order isomorphism implies a similarity of representations, by definition. This 445  
view has been called *holism* in the philosophy of mind [39, 41] and is related to a similar 446  
idea of *meaning holism* in language, which is the idea that the meaning of a linguistic 447  
expression is determined by its relation to other expressions within a language [42, 43]. 448  
For example, Firth [44] (p. 11) writes: “you shall know a word by the company it keeps”. 449  
Similarly, Griffiths and Steyvers [45], and Griffiths, Steyvers, and Tenenbaum [46] have 450  
adopted meaning holism accounts of semantic representations in neural networks. More 451  
recently, Piantadosi and Hill [47] have argued that large language models capture impor- 452  
tant aspects of meaning and approximate human cognition providing one assumes meaning 453  
arises from the relations of states rather than an architecture or training. Even if a re- 454  
searcher were to adopt this holistic perspective on representations, our results should still 455  
be of interest to them as they show that the similarity between representational geome- 456  
tries can vary based on the visual stimulus that is used to compare them (the modulation 457  
effect). As all our studies show, representational geometries heavily depend on the dataset 458  
- e.g., a networks trained with a confound in the training set will have vastly different 459  
geometries when presented test data with and without the confound present. Addition- 460  
ally, our results show that adopting this view misses the information about differences 461  
in mechanistic processes that a psychologist or neuroscientist is frequently interested in, 462  
for instance, whether the visual system processes surface reflectance or shape in order to 463

identify objects. Fodor and Lepore long ago criticized this philosophical stance [39, 48], and interestingly, this philosophical debate played an important part in the development of RSA (see section S1 of Supplementary Information). Unfortunately, this debate has largely been ignored by contemporary researchers who use RSA as a method to infer similarity of systems.

## Relation to existing research

A related point has been made by Kriegeskorte and Diedrichsen [22] and Kriegeskorte and Wei [49], who point out that two systems may have the same representational geometry, even if they have a different activity profile over neurons. In this sense, the geometry abstracts away the information about how information was distributed over a set of neurons. Kriegeskorte and Diedrichsen [22] equate this loss in information to “peeling a layer of an onion” – downstream decoders that are sensitive to the representational geometry rather than activity profiles over neuron populations can focus on difference in information as reflected by a change in geometry and be agnostic to how this information is distributed over a set of neurons. We agree that this invariance over activity profiles is indeed a useful property of representational geometries for downstream decoders. However, while abstracting over activity profiles may be useful, abstracting over stimulus properties loses an important piece of information when comparing representations across brain regions, individuals, species and between brains and computational models. Our studies show how two systems may appear similar based on their representational geometries in one circumstance (e.g. Fig 2B) but drastically different in another circumstance (Fig 2C). Furthermore, our results show how such second-order confounds can arise because of properties of transformations (Study 1) and structural properties of datasets used for testing (Study 4).

It is also important to note how our results differ from previous studies exploring 488  
limitations of RSA. Some of these studies have focused on the importance of how neural 489  
data is pre-processed. For example, Ramirez [50] found that pre-processing steps, such 490  
as centering (de-meaning) activation vectors may lead to incorrect inference about the 491  
representational geometry of activations. They demonstrated that subtracting the mean 492  
from activations could change the rank order of similarity between conditions. In turn, this 493  
could lead to clearly distinct RDMs becoming highly correlated and vice-versa. While this 494  
is an important methodological point, it is clearly distinct from the point we are making 495  
in this study. Indeed, the results here are agnostic of the data pre-processing steps and 496  
hold whether or not activations are centered. 497

Another set of studies have explored how the procedure of data collection can influence 498  
the results of RSA. For example, Henriksson et al. [51] and Cai et al. [52] demonstrated 499  
that RDMs measured based on fMRI data can be severely biased because of temporal 500  
and spatial correlations in neural activity. These authors have pointed out that if activity 501  
patterns from different brain regions are recorded during the same trial, the similarity 502  
estimates will be exaggerated due to correlated neural fluctuations in these regions. Sim- 503  
ilarly, neural activity is correlated over time, which means estimated similarity based on 504  
activity patterns from the same imaging run also introduces a strong bias in RDMs. These 505  
sources of bias are important to understand, but they can also be addressed by a more 506  
careful task design and analysis [52]. In contrast, the confounds that are highlighted in 507  
this study exist in the stimulus itself. Therefore, even if one were to completely mitigate 508  
the bias in estimating RDMs, the types of confounds we highlight in our work would still 509  
pose problems when drawing inferences from correlation in RDMs. 510

A third set of studies have highlighted the importance of choosing the correct distance 511  
metric when using RSA. For example, Ramirez [26] compared Euclidean distance with an 512

angular metric (such as cosine similarity) and showed that the choice of distance metric 513  
can reveal different aspects of the same fMRI data. They argued that the Euclidean 514  
distance is particularly sensitive to the mean activity over a recorded voxel. Based on this 515  
analysis, Ramirez [26] suggested using an angular distance metric, especially when neural 516  
signal is aggregated over large number of neurons. Similarly, in another exhaustive study 517  
over distance measures, Bobadilla-Suarez et al. [25], evaluated neural similarity using 518  
various distance measures, including angle-based measures (cosine, Pearson, Spearman) 519  
and magnitude-based measures (Euclidean, Mahalanobis, Minkowski) and found that the 520  
choice of metric significantly influenced the measured similarity. They also found that 521  
there was no one metric that outperformed all others – rather, the preferred metric varied 522  
across different studies, but was consistent across brain regions within a study. The choice 523  
of distance metric is again a related but orthogonal issue to the one we highlight in this 524  
study. Representational geometry abstracts away information about stimulus features 525  
and how inputs are transformed. Our results demonstrate how different stimulus features 526  
and transformations of input stimulus can lead to the same representational geometry. 527  
This is integral to the nature of representational geometries, rather than a consequence 528  
of the distance metric used. 529

Of course, the problem of confounds in stimuli is not unique to RSA and will affect 530  
other statistical analyses, including multivariate regression methods such as MVP classi- 531  
fication. A number of studies have pointed out this conceptual problem in the context 532  
of multivariate decoding, where authors have argued that successfully decoding a signal 533  
from a neural activation pattern is no guarantee that the signal is encoded by the brain 534  
or decoded by downstream processes [18–20, 53]. In response, researchers have adopted 535  
a variety of methods to deal with such confounds such as cross-validation [23], confound 536  
regression [24], counter-balancing data [54] and commonality analysis [55]. We couldn't 537

agree more with this direction of research and our study highlights two properties of 538  
confounds that makes it especially challenging to compare neural representations with 539  
those in complex models working on high-dimensional data. Firstly, these confounds are 540  
second-order – that is, they are not only category-correlated (as is the case for confounds 541  
highlighted for multivariate decoding), but also mimic the second-order similarity struc- 542  
ture of the variable of interest. Secondly, when using high-dimensional datasets (such as 543  
naturalistic images) and complex target models (such as DNNs) for testing, these con- 544  
founds are unknown to the experimenter and may be present in the entire dataset. This 545  
restricts the utility of existing methods, such as cross-validation and counter-balancing 546  
data, for dealing with these confounds, a point made by researchers employing these 547  
methods [24, 54]. In fact, we are unaware of any statistical methods that can completely 548  
eliminate confounds under these settings. 549

## Limitations and General Recommendations 550

Even though our findings demonstrate that second-order confounds are possible and plau- 551  
sible, they do not allow us to infer whether such confounds *are* present in existing datasets 552  
and driving the observed similarity in existing studies. An important research direction 553  
is to discover these confounds (or lack thereof) and determine the extent to which repre- 554  
sentations in a target model mimic the second-order relationship between neural repre- 555  
sentations in the visual cortex. One way to do this is to systematically eliminate possible 556  
confounds from datasets and check the extent to which this affects previously observed 557  
results. Of course, this is not straightforward to do in high-dimensional stimuli, such 558  
as naturalistic images, which consist of millions of features and combinatorial relation- 559  
ships between these features. Thus identifying confounds in such datasets remains a real 560  
challenge. 561

In closing, we describe our recommendations for practitioners who would like to use RSA for comparing complex systems based on high-dimensional data. First, since the intrinsic structure of datasets can artificially modulate RSA scores, researchers should compare systems on a wider variety of datasets and sampling schemes than currently done. Second, given that confounding features can lead to mimicked representational geometries, researchers should consider running additional controlled experiments that manipulate independent variables designed to test hypotheses to rule out this possibility when inferences hinge crucially on it. This point has recently been made by Bowers et al [36] in relation to testing the similarities of DNN and human vision. Similarly, the ‘controversial stimuli’ designed by Golan et al. [56] should also enable researchers to test representational geometries for stimuli where different models make contrasting predictions. Third, when studies are conducted to search for evidence of mechanistic similarity between two or more systems, researchers should use a wider range of complementary methods in order to increase robustness (e.g., RSA combined with neural predictivity [12], MVPC [7, 57], CCA [58], SVCCA [59], CKA [60]).

Lastly, perhaps the most important general recommendation we make is that researchers should acknowledge, procedurally and in writing, which inferences are afforded 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 similarity 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 theoretically 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

587

### Dataset generation and training

588

All DNN simulations (Studies 2–4) were carried out using the Pytorch framework [61]. 589

The model implementations were downloaded from the torchvision library. Networks 590

trained on unperturbed datasets in all studies were pre-trained on ImageNet as were 591

networks trained on modified datasets in Study 2. Networks trained on modified datasets 592

in Studies 3 and 4 were randomly initialised. For the pre-trained models, their pre-trained 593

weights were downloaded from torchvision.models subpackage. 594

**Study 1** Each dataset in Study 1 consists of 100 samples (50 in each cluster) drawn 595

from two multivariate Gaussians,  $\mathcal{N}(x|\mu, \Sigma)$ , where  $\mu$  is a 2-dimensional vector and  $\Sigma$  596

is a  $2 \times 2$  covariance matrix. In Fig 2A, the two Gaussians have means  $\mu_1 = (1, 8)$  and 597

$\mu_2 = (8, 1)$  and a covariance matrices  $\Sigma_1 = \Sigma_2 = \frac{1}{2}\mathbf{I}$ , while in Fig 2B the Gaussians 598

have means  $\mu_1 = (1, 1)$  and  $\mu_2 = (8, 8)$  and a covariance matrices  $\Sigma_1 = \mathbf{I}$ ,  $\Sigma_2 = 8\mathbf{I}$ . 599

All kernel matrices were computed using the sklearn.metrics.pairwise module of the 600

scikit-learn Python package. 601

**Study 2** First, a VGG-16 deep convolutional neural network [62], pre-trained on the 602

ImageNet dataset of naturalistic images, was trained to classify stimuli from the CIFAR-10 603

dataset [63]. The CIFAR-10 dataset includes 10 categories with 5000 training, and 1000 604

test images per category. The network was fine-tuned on CIFAR-10 by replacing the 605

classifier so that the final fully-connected layer reflected the correct number of target 606

classes in CIFAR-10 (10 for CIFAR-10 as opposed to 1000 for ImageNet). Images were 607

rescaled to a size of  $224 \times 224$ px and then the model learnt to minimise the cross-entropy 608

error using the RMSprop optimizer with a mini-batch size of 64, learning rate of  $10^{-5}$ , 609

and momentum of 0.9. All models were trained for 10 epochs, which were sufficient for  
convergence across all datasets.

Second, 100 random images from the test set for each category were sampled as in-  
put for the network and activations at the final convolutional layer extracted using the  
THINGSVision Python toolkit [64]. The same toolkit was used to generate a representa-  
tional 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  $10 \times 10$ , category-level,  
RDM which reflected median between-category distances.

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 Fig 5). The locations of these pixels were determined using the averaged RDM from  
the previous step. We call this the target RDM. In the 'Positive' condition, we wanted the  
distances between pixel placements to be positively correlated to the distances between  
categories in the target RDM. We achieved this by using an iterative algorithm that  
sampled pixel placements at random, calculated an RDM based on distances between  
the pixel placements and computed an RSA score (Spearman correlation) with the target  
RDM. Placements with a score above 0.70 were retained and further optimized (using  
small perturbations) to achieve an RSA-score over 0.90. The same procedure was also  
used to determine placements in the Uncorrelated (optimizing for a score close to 0) and  
Negatively correlated (optimizing for a negative score) conditions.

Finally, datasets were created using 10 different placements in each of the three condi-

tions. Networks were trained for classification on these modified CIFAR-10 datasets in the same manner as the VGG-16 network trained on the unperturbed version of the dataset (See Fig 4).

**Study 3** The procedure mirrored Study 2 with the main difference being that the target system was the macaque inferior temporal cortex. Neural data from two macaques, as well as the dataset were obtained from the Brain Score repository [30]. This dataset consists of 3200 images from 8 categories (animals, boats, cars, chairs, faces, fruits, planes, and tables), we computed an  $8 \times 8$  averaged RDM based on macaque IT response patterns for stimuli in each category.

This averaged RDM was then used as the target RDM in the optimization procedure to determine locations of the confound (here, a white predictive patch of size  $5 \times 5$  pixels) for each category. Using a patch instead of a single pixel was required in this dataset because of the structure and smaller size of the dataset (3200 images, rather than 50,000 images for CIFAR-10). In this smaller dataset, the networks struggle to learn based on a single pixel. However, increasing the size of the patch makes these patches more predictive and the networks are able to again learn entirely based on this confound (see results in Fig 6). In a manner similar to Study 2, this optimisation procedure was used to construct three datasets, where the confound's placement was positively correlated, uncorrelated or negatively correlated with the category distances in the target RDM.

Finally, each dataset was split into 75% training (2432 images) and 25% test sets (768 images) before VGG-16 networks were trained on the unperturbed and modified datasets in the same manner as in Study 2. One difference between Studies 2 and 3 was that here the networks in the Positive, Uncorrelated and Negative conditions were trained from scratch, i.e., not pre-trained on ImageNet. This was done because we wanted

to make sure that the network in the Normal condition (trained on ImageNet) and the 659  
networks in the Positive, Uncorrelated and Negative conditions encoded fundamentally 660  
different features of their inputs – i.e., there were no ImageNet-related features encoded by 661  
representations  $\Phi_2$ ,  $\Phi_3$  and  $\Phi_4$  that were responsible for the similarity in representational 662  
geometries between these representations and the representations in macaque IT cortex. 663

**Study 4** The target system in this study was human IT cortex. The human RDM 664  
and dataset were obtained from [6]. Rather than calculating pixel placements based on 665  
the human RDM, the hierarchical structure of the dataset was used to place the pixels 666  
manually. The dataset consists of 910 images from 6 categories: human bodies, human 667  
faces, animal bodies, animal faces, artificial inanimate objects and natural inanimate 668  
objects. These low-level categories can be organised into the hierarchical structure shown 669  
in Fig 8. Predictive pixels were manually placed so that the distance between pixels for 670  
Animate kinds were closer together than they were to Inanimate kinds and that faces 671  
were closer together than bodies. This can be done in many different ways, so we created 672  
five different datasets, with five possible arrangements of predictive pixels. Results in the 673  
Hierarchical condition (Fig 9) are averaged over these five datasets. Placements for the 674  
Random condition were done similarly, except that the locations were selected randomly. 675

Networks were then trained on a 6-way classification task (818 training images and 92 676  
test images) in a similar manner to the previous studies. As in Study 3, networks trained 677  
on the modified datasets (both Hierarchical and Random conditions) were not pre-trained 678  
on ImageNet. 679

## RDM and RSA computation

680

For Studies 2-4 all image-level RDMs were calculated using  $1 - r$  as the distance measure.

681

RSA scores were computed as the Spearman rank correlation between RDMs.

682

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

683

684

685

686

687

In Study 3, the dataset consisted of 3200 images belonging to 8 categories. We first calculated a full  $3200 \times 3200$  RDM using the entire set of stimuli. An averaged, category-level,  $8 \times 8$  RDM was then calculated using median distances between categories (in a manner similar to that described for Study 2 in the Section ‘Dataset generation and training’). This  $8 \times 8$  RDM was used to determine the RSA-scores. We also obtained qualitatively similar results using the full  $3200 \times 3200$  RDMs. These results can be found in the S2 section of Supplementary Information.

688

689

690

691

692

693

694

In Study 4, the dataset consisted of 818 training images and 92 test images. Kriegeskorte et al. [6] used these images to obtain a  $92 \times 92$  RDM to compare representations between human and macaque IT cortex. Here we computed a similar  $92 \times 92$  RDM for networks trained in the Normal, Hierarchical and Random training conditions, which were then compared with the  $92 \times 92$  RDM from human IT cortex to obtain RSA-scores for each condition.

695

696

697

698

699

700

## Testing

701

In Study 2, we used a  $4 \times 2$  design to measure classification performance for networks in all four conditions (Normal, Positive, Uncorrelated and Negative) on both unperturbed images and modified images. We computed six RSA-scores: three pairs of networks – Normal-Positive, Normal-Uncorrelated and Normal-Negative – and two types of inputs – unperturbed and modified test images. The noise ceiling (grey band in Fig 6) was determined in the standard way as described in [65] and represents the expected range of the highest possible RSA score with the target system (network trained on the unperturbed dataset).

702

703

704

705

706

707

708

709

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

710

711

712

713

714

715

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

716

717

718

719

## Data analysis

720

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

721

722

723

were tested by running a 3x2 mixed ANOVA. 724

RSA scores with the target system between networks in various conditions were com- 725  
pared by running 3x2 ANOVAs in Studies 2 and 3, and a 2x2 ANOVA in Study 4. We 726  
observed that RSA-scores were highly dependent on both the way the networks were 727  
trained and also the test images used to elicit response activations. 728

For a detailed overview of the statistical analyses and results, see section S3 of the Sup- 729  
plementary Information. 730

## Data Availability 731

Confound placement coordinates (Studies 2-4), unperturbed datasets (Studies 3 and 4), 732  
macaque activation patterns and RDMs (Study 3) and human RDM (Study 4) are avail- 733  
able at [OSF](#). 734

## Acknowledgments 735

This project has received funding from the European Research Council (ERC) under the 736  
European Union’s Horizon 2020 research and innovation programme (grant agreement No 737  
741134) 738

## References 739

- [1] Hubel DH, Wiesel TN. Receptive fields of single neurones in the 740  
cat’s striate cortex. The Journal of Physiology. 1959;148(3):574–591. 741  
doi:<https://doi.org/10.1113/jphysiol.1959.sp006308>. 742

- [2] O’Keefe J. Place units in the hippocampus of the freely moving rat. *Experimental Neurology*. 1976;51(1):78–109. doi:[https://doi.org/10.1016/0014-4886\(76\)90055-8](https://doi.org/10.1016/0014-4886(76)90055-8). 743  
744
- [3] Hafting T, Fyhn M, Molden S, Moser MB, Moser EI. Microstructure of a spatial map in the entorhinal cortex. *Nature*. 2005;436(7052):801–806. doi:<https://doi.org/10.1038/nature03721>. 745  
746  
747
- [4] Kriegeskorte N, Kievit RA. Representational geometry: integrating cognition, computation, and the brain. *Trends in Cognitive Sciences*. 2013;17(8):401–412. doi:<https://doi.org/10.1016/j.tics.2013.06.007>. 748  
749  
750
- [5] Kriegeskorte N, Mur M, Bandettini P. Representational similarity analysis - connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*. 2008;2. doi:<https://doi.org/10.3389/neuro.06.004.2008>. 751  
752  
753
- [6] Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, et al. Matching Categorical Object Representations in Inferior Temporal Cortex of Man and Monkey. *Neuron*. 2008;60(6):1126–1141. doi:<https://doi.org/10.1016/j.neuron.2008.10.043>. 754  
755  
756
- [7] Haxby J, Guntupalli J Connolly A, Halchenko Y, Conroy B, Gobbini M et al. A Common, High-Dimensional Model of the Representational Space in Human Ventral Temporal Cortex. *Neuron*. 2011;72(2):404–416. doi:<https://doi.org/10.1016/j.neuron.2011.08.026>. 757  
758  
759  
760
- [8] 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*. 2020;209:116322. doi:<https://doi.org/10.1016/j.neuroimage.2019.116322>. 761  
762  
763  
764



- [9] Michael L Mack BL Alison R Preston. Decoding the Brain’s Algorithm for Catego- 765  
rization from Its Neural Implementation. *Current Biology*. 2013;23(20):2023–2027. 766  
doi:<https://doi.org/10.1016/j.cub.2013.08.035>. 767
- [10] Freund MC, Etzel JA, Braver TS. Neural Coding of Cognitive Control: The 768  
Representational Similarity Analysis Approach. *Trends in Cognitive Sciences*. 769  
2021;25(7):622–638. doi:<https://doi.org/10.1016/j.tics.2021.03.011>. 770
- [11] Kaneshiro B, Perreau Guimaraes M, Kim HS, Norcia AM, Suppes P. A 771  
Representational Similarity Analysis of the Dynamics of Object Process- 772  
ing Using Single-Trial EEG Classification. *PLOS ONE*. 2015;10(8):1–27. 773  
doi:<https://doi.org/10.1371/journal.pone.0135697>. 774
- [12] Yamins DLK, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. Performance- 775  
optimized hierarchical models predict neural responses in higher visual cor- 776  
tex. *Proceedings of the National Academy of Sciences*. 2014;111(23):8619–8624. 777  
doi:<https://doi.org/10.1073/pnas.1403112111>. 778
- [13] Khaligh-Razavi SM, Kriegeskorte N. Deep Supervised, but Not Unsupervised, Models 779  
May Explain IT Cortical Representation. *PLOS Computational Biology*. 2014;10:1– 780  
29. doi:<https://doi.org/10.1371/journal.pcbi.1003915>. 781
- [14] Kietzmann TC, Spoerer CJ, Sörensen LKA, Cichy RM, Hauk O, Kriegeskorte N. 782  
Recurrence is required to capture the representational dynamics of the human visual 783  
system. *Proceedings of the National Academy of Sciences*. 2019;116(43):21854–21863. 784  
doi:<https://doi.org/10.1073/pnas.1905544116>. 785
- [15] Cichy RM, Khosla A, Pantazis D, Torralba A, Oliva A. Comparison of deep 786  
neural networks to spatio-temporal cortical dynamics of human visual object 787

- recognition reveals hierarchical correspondence. *Scientific Reports*. 2016;6:27755. 788  
doi:<https://doi.org/10.1038/srep27755>. 789
- [16] Kiat JE, Luck SJ, Beckner AG, Hayes TR, Pomaranski KI, Henderson JM, et al. Link- 790  
ing patterns of infant eye movements to a neural network model of the ventral stream 791  
using representational similarity analysis. *Developmental Science*. 2022;25(1):e13155. 792  
doi:<https://doi.org/10.1111/desc.13155>. 793
- [17] Haxby JV, Connolly AC, Guntupalli JS. Decoding neural representational 794  
spaces using multivariate pattern analysis. *Annu Rev Neurosci*. 2014;37:435–456. 795  
doi:<https://doi.org/10.1146/annurev-neuro-062012-170325>. 796
- [18] Naselaris T, Kay KN. Resolving Ambiguities of MVPA Using Explicit Mod- 797  
els of Representation. *Trends in Cognitive Sciences*. 2015;19(10):551–554. 798  
doi:<https://doi.org/10.1016/j.tics.2015.07.005>. 799
- [19] Ritchie JB, Kaplan DM, Klein C. Decoding the Brain: Neural Representation and 800  
the Limits of Multivariate Pattern Analysis in Cognitive Neuroscience. *The British* 801  
*Journal for the Philosophy of Science*. 2019;70(2):581–607. doi:10.1093/bjps/axx023. 802
- [20] Weichwald S, Meyer T, Özdenizci O, Schölkopf B, Ball T, Grosse-Wentrup M. Causal 803  
interpretation rules for encoding and decoding models in neuroimaging. *Neuroimage*. 804  
2015;110:48–59. 805
- [21] Kriegeskorte N. Deep Neural Networks: A New Framework for Modeling Biolog- 806  
ical Vision and Brain Information Processing. *Annual Review of Vision Science*. 807  
2015;1(1):417–446. doi:10.1146/annurev-vision-082114-035447. 808

- [22] Kriegeskorte N, Diedrichsen J. Peeling the Onion of Brain Representations. Annual Review of Neuroscience. 2019;42(1):407–432. doi:<https://doi.org/10.1146/annurev-neuro-080317-061906>. 809  
810  
811
- [23] Snoek L, Miletić S, Scholte HS. How to control for confounds in decoding analyses of neuroimaging data. NeuroImage. 2019;184:741–760. doi:<https://doi.org/10.1016/j.neuroimage.2018.09.074>. 812  
813  
814
- [24] Kostro D, Abdulkadir A, Durr A, Roos R, Leavitt BR, Johnson H, et al. Correction of inter-scanner and within-subject variance in structural MRI based automated diagnosing. NeuroImage. 2014;98:405–415. doi:<https://doi.org/10.1016/j.neuroimage.2014.04.057>. 815  
816  
817  
818
- [25] Bobadilla-Suarez S, Ahlheim C, Mehrotra A, Panos A, Love BC. Measures of Neural Similarity. Computational Brain & Behavior. 2020;3(4):369–383. doi:10.1007/s42113-019-00068-5. 819  
820  
821
- [26] Ramírez FM. Orientation Encoding and Viewpoint Invariance in Face Recognition: Inferring Neural Properties from Large-Scale Signals. The Neuroscientist. 2018;24(6):582–608. doi:10.1177/1073858418769554. 822  
823  
824
- [27] Schölkopf B, Smola AJ, Bach F. Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT Press; 2002. 825  
826
- [28] Sahami M, Heilman TD. 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. New York, NY, USA: Association for Computing Machinery; 2006. p. 377–386. 827  
828  
829  
830

- [29] Torralba A, Efros AA. Unbiased look at dataset bias. In: CVPR 2011; 2011. p. 831  
1521–1528. 832
- [30] Schrimpf M, Kubilius J, Hong H, Majaj NJ, Rajalingham R, Issa EB, et al. Brain- 833  
Score: Which Artificial Neural Network for Object Recognition is most Brain-Like? 834  
bioRxiv preprint: 407007. 2018;doi:<https://doi.org/10.1101/407007>. 835
- [31] Majaj NJ, Hong H, Solomon EA, DiCarlo JJ. Simple Learned Weighted Sums 836  
of Inferior Temporal Neuronal Firing Rates Accurately Predict Human Core Ob- 837  
ject Recognition Performance. *Journal of Neuroscience*. 2015;35(39):13402–13418. 838  
doi:<https://doi.org/10.1523/JNEUROSCI.5181-14.2015>. 839
- [32] Kriegeskorte N. Relating Population-Code Representations between Man, Mon- 840  
key, and Computational Models. *Frontiers in Neuroscience*. 2009;3(3):363–373. 841  
doi:<https://doi.org/10.3389/neuro.01.035.2009>. 842
- [33] Raghu M, Unterthiner T, Kornblith S, Zhang C, Dosovitskiy A. Do vision trans- 843  
formers see like convolutional neural networks? *Advances in Neural Information* 844  
*Processing Systems*. 2021;34:12116–12128. 845
- [34] Serre T. Deep Learning: The Good, the Bad, and the Ugly. *Annual Review of Vision* 846  
*Science*. 2019;5(1):399–426. doi:[10.1146/annurev-vision-091718-014951](https://doi.org/10.1146/annurev-vision-091718-014951). 847
- [35] Dujmović M, Malhotra G, Bowers JS. What do adversarial images tell us about 848  
human vision? *eLife*. 2020;9:e55978. doi:[10.7554/eLife.55978](https://doi.org/10.7554/eLife.55978). 849
- [36] Bowers JS, Malhotra G, Dujmović M, Montero ML, Tsvetkov C, Biscione V, et al.. 850  
Deep Problems with Neural Network Models of Human Vision; 2022. Available from: 851  
[psyarxiv.com/5zf4s](https://psyarxiv.com/5zf4s). 852

- [37] Geirhos R, Jacobsen JH, Michaelis C, Zemel R, Brendel W, Bethge M, et al. Shortcut 853  
learning in deep neural networks. *Nature Machine Intelligence*. 2020;2(11):665–673. 854  
doi:<https://doi.org/10.1038/s42256-020-00257-z>. 855
- [38] Xu Y, Vaziri-Pashkam M. Limits to visual representational correspondence be- 856  
tween convolutional neural networks and the human brain. *Nature Communications*. 857  
2021;12:2065. doi:<https://doi.org/10.1038/s41467-021-22244-7>. 858
- [39] Fodor J, Lepore E. *Holism: A Shoppers Guide*. Cambridge: Blackwell; 1992. 859
- [40] Edelman S. Representation is representation of similarities. *Behavioral and Brain* 860  
*Sciences*. 1998;21(4):449–467. doi:<https://10.1017/S0140525X98001253>. 861
- [41] Block N. Advertisement for a Semantics for Psychology. *Midwest Studies in Philos-* 862  
*ophy*. 1986;10(1):615–678. doi:<https://10.1111/j.1475-4975.1987.tb00558.x>. 863
- [42] Hempel CG. Problems and Changes in the Empiricist Criterion of Meaning. *Revue* 864  
*Internationale de Philosophie*. 1950;4(11):41–63. 865
- [43] Quine WV. Main Trends in Recent Philosophy: Two Dogmas of Empiricism. *The* 866  
*Philosophical Review*. 1951;60(1):20–43. doi:<https://doi.org/10.2307/2181906>. 867
- [44] Firth JR. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*. 868  
1957; p. 1–31. 869
- [45] Griffiths TL, Steyvers M. A probabilistic approach to semantic representation. In: 870  
*Proceedings of the Twenty-Fourth Annual Conference of the Cognitive Science So-* 871  
*ciety*. Hillsdale, NJ: Erlbaum; 2002. 872

- [46] Griffiths TL, Steyvers M, Tenenbaum JB. A probabilistic approach 873  
to semantic representation. *Psychological Review*. 2007;114(2):211–244. 874  
doi:<https://psycnet.apa.org/doi/10.1037/0033-295X.114.2.211>. 875
- [47] Piantadosi S, Hill F. Meaning without reference in large language models. In: 876  
*NeurIPS 2022 Workshop on Neuro Causal and Symbolic AI (nCSI)*; 2022. Available 877  
from: <https://openreview.net/forum?id=nRkJEwmZnM>. 878
- [48] Fodor J. *Psychosemantics: The Problem of Meaning in the Philosophy of Mind*. 879  
Cambridge: MIT Press; 1987. 880
- [49] Kriegeskorte N, Wei XX. Neural tuning and representational geometry. *Nature Re-* 881  
*views Neuroscience*. 2021;22(11):703–718. doi:[https://doi.org/10.1038/s41583-021-](https://doi.org/10.1038/s41583-021-00502-3) 882  
00502-3. 883
- [50] Ramírez FM. Representational confusion: the plausible consequence of demeaning 884  
your data. *bioRxiv*. 2017;doi:10.1101/195271. 885
- [51] Henriksson L, Khaligh-Razavi SM, Kay K, Kriegeskorte N. Visual representa- 886  
tions are dominated by intrinsic fluctuations correlated between areas. *NeuroImage*. 887  
2015;114:275–286. doi:<https://doi.org/10.1016/j.neuroimage.2015.04.026>. 888
- [52] Cai MB, Schuck NW, Pillow JW, Niv Y. Representational structure or task 889  
structure? Bias in neural representational similarity analysis and a Bayesian 890  
method for reducing bias. *PLOS Computational Biology*. 2019;15(5):1–30. 891  
doi:10.1371/journal.pcbi.1006299. 892
- [53] Hebart MN, Baker CI. Deconstructing multivariate decoding 893  
for the study of brain function. *NeuroImage*. 2018;180:4–18. 894  
doi:<https://doi.org/10.1016/j.neuroimage.2017.08.005>. 895

- [54] Rao A, Monteiro JM, Mourao-Miranda J. Predictive modelling using neuroimaging data in the presence of confounds. *NeuroImage*. 2017;150:23–49. doi:<https://doi.org/10.1016/j.neuroimage.2017.01.066>.
- [55] Greene MR, Baldassano C, Esteva A, Beck DM, Fei-Fei L. Visual scenes are categorized by function. *Journal of Experimental Psychology: General*. 2016;145(1):82.
- [56] Golan T, Raju PC, Kriegeskorte N. Controversial stimuli: Pitting neural networks against each other as models of human cognition. *Proceedings of the National Academy of Sciences*. 2020;117(47):29330–29337. doi:<https://doi.org/10.1073/pnas.1912334117>.
- [57] Haxby JV. Multivariate pattern analysis of fMRI: the early beginnings. *Neuroimage*. 2012;62(2):852–855. doi:<https://doi.org/10.1016/j.neuroimage.2012.03.016>.
- [58] Morcos A, Raghu M, Bengio S. Insights on representational similarity in neural networks with canonical correlation. In: Bengio S, Wallach H, Larochelle H, Grauman K, Cesa-Bianchi N, Garnett R, editors. *Advances in Neural Information Processing Systems*. vol. 31. Curran Associates, Inc.; 2018. Available from: <https://proceedings.neurips.cc/paper/2018/file/a7a3d70c6d17a73140918996d03c014f-Paper.pdf>.
- [59] Raghu M, Gilmer J, Yosinski J, Sohl-Dickstein J. SVCCA: Singular Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, et al., editors. *Advances in Neural Information Processing Systems*. vol. 30. Curran Associates, Inc.; 2017. Available from: <https://proceedings.neurips.cc/paper/2017/file/dc6a7e655d7e5840e66733e9ee67cc69-Paper.pdf>.

- [60] Kornblith S, Norouzi M, Lee H, Hinton G. Similarity of neural network representa- 918  
tions revisited. In: International Conference on Machine Learning. PMLR; 2019. p. 919  
3519–3529. 920
- [61] Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: An 921  
Imperative Style, High-Performance Deep Learning Library. In: Advances in Neural 922  
Information Processing Systems 32. Curran Associates, Inc.; 2019. p. 8024–8035. 923
- [62] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image 924  
recognition. arXiv preprint arXiv:14091556. 2014;. 925
- [63] Krizhevsky A, Hinton G. Learning multiple layers of features from tiny images. 926  
Toronto, Ontario: University of Toronto; 2009. 927
- [64] Muttenthaler L, Hebart MN. THINGSvision: A Python Toolbox for Streamlining the 928  
Extraction of Activations From Deep Neural Networks. *Frontiers in Neuroinformatics*. 929  
2021;15:679838. doi:<https://doi.org/10.3389/fninf.2021.679838>. 930
- [65] Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N. A Toolbox 931  
for Representational Similarity Analysis. *PLOS Computational Biology*. 2014;10:1– 932  
11. doi:<https://doi.org/10.1371/journal.pcbi.1003553>. 933