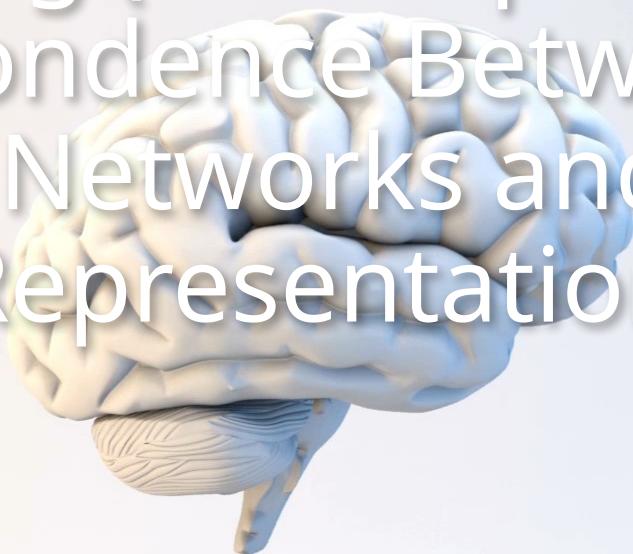


Evaluating (and Improving) the Correspondence Between Deep Neural Networks and Human Representations



**Joshua C. Peterson, Joshua
T. Abbott, Thomas L.
Griffiths (P.A.G)**

(w/ graphics from Ayşe
Aybüke Durmaz, Homa
Priya Tarigopula)

P.A.G Project Overview W



Introduction



Aims



Evaluating the correspondence between representations

Stimuli,
Methods,
Results,
Analysis



Reweighting starts like RSA

- Comparing the representations formed by deep neural networks to those of humans.
- Human representation:
 - Similarity judgments.
 - A similarity function over a set of pairs of data points corresponds to an implicit representation of those points

Evaluating the correspondence between representations



- Obtain the similarity within images using human ratings
- Find the similarity within images of each category by
 - Studying Raw representations of DNNs
 - Studying Transformed Representations of DNNs
- Probe both DNN and human similarities to identify the spatial and taxonomic information encoded
 - MDS, HCA

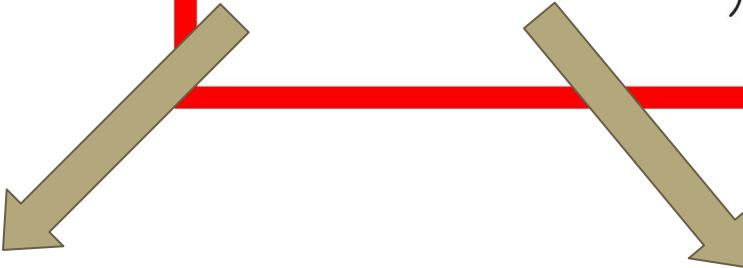
Evaluating the correspondence between representations

- Stimuli
 - 6 categories - 120 images each
- Procedures:
 - pairwise image similarity ratings
 - Amazon Mechanical Turk
 - 0 (“not similar at all”)
 - 10 (“very similar”)
 - six 120x120 similarity matrices



Evaluating the correspondence between representations

$$S = FF^T,$$



**SIMILARITY MATRIX
from humans**

**FEATURE MATRIX
(dot product/ R/
cosine)**

$$A = [x_0, x_1, x_2, \dots, x_n]$$

$$B = [y_0, y_1, y_2, \dots, y_n]$$

$$A \odot B = \sum_{i=0}^n x_i \cdot y_i$$

Evaluating the correspondence between representations

$$\mathbf{S} = \mathbf{F}\mathbf{F}^T,$$

FEATURE MATRIX



Deep neural network representations :

- When deep neural networks are presented with an image, the nodes that comprise the network obtain different activation values.

Multidimensional feature representation

	1	2	...	n
1	a_{11}	a_{12}	...	a_{1n}
2	a_{21}	a_{22}	...	a_{2n}
3	a_{31}	a_{32}	...	a_{3n}
:	:	:	:	:
m	a_{m1}	a_{m2}	...	a_{mn}

120
images

Evaluating the correspondence between representations

$$\mathbf{S} = \mathbf{F}\mathbf{F}^T,$$

FEATURE MATRIX

120
images

Multidimensional feature representation
[ncol :4,096]

$$\begin{matrix} & 1 & 2 & \dots & n \\ 1 & a_{11} & a_{12} & \dots & a_{1n} \\ 2 & a_{21} & a_{22} & \dots & a_{2n} \\ 3 & a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ m & a_{m1} & a_{m2} & \dots & a_{mn} \end{matrix}$$

Evaluating the correspondence between representations

$$\mathbf{S} = \mathbf{F}\mathbf{F}^T,$$

Inner Product of Feature Matrix produces a 120x120 similarity matrix:

FEATURE MATRIX

$$\mathbf{F} = [120 \times 4096]$$

$$\begin{matrix} & 1 & 2 & \dots & n \\ 1 & a_{11} & a_{12} & \dots & a_{1n} \\ 2 & a_{21} & a_{22} & \dots & a_{2n} \\ 3 & a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ m & a_{m1} & a_{m2} & \dots & a_{mn} \end{matrix}$$



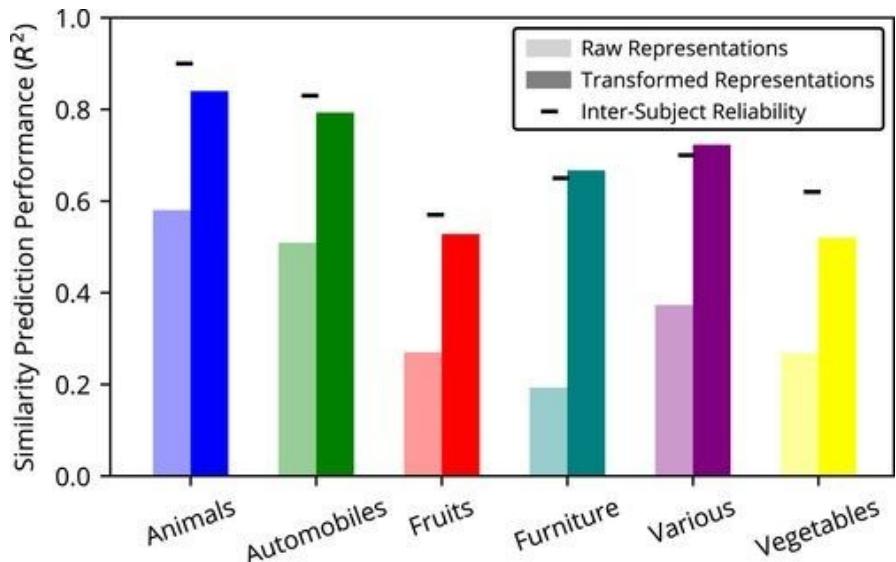
$$\mathbf{F}^T = [4096 \times 120]$$

$$\begin{matrix} & 1 & 2 & \dots & n \\ 1 & a_{11} & a_{12} & \dots & a_{1n} \\ 2 & a_{21} & a_{22} & \dots & a_{2n} \\ 3 & a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ m & a_{m1} & a_{m2} & \dots & a_{mn} \end{matrix}$$

Evaluating the correspondence between representations

Results :

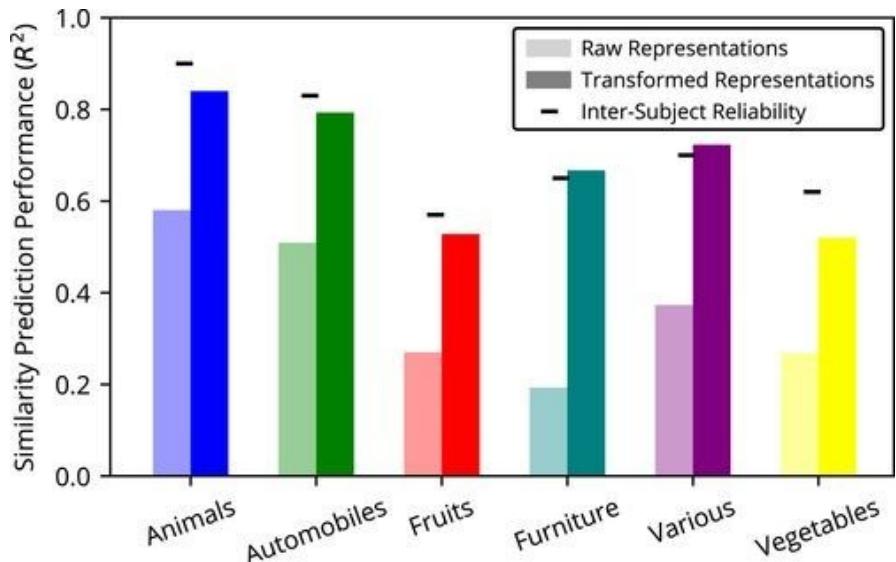
- computing the correlation between the human similarity judgments and the inner products computed in the deep feature spaces. (both 120x120)



Evaluating the correspondence between representations

Results :

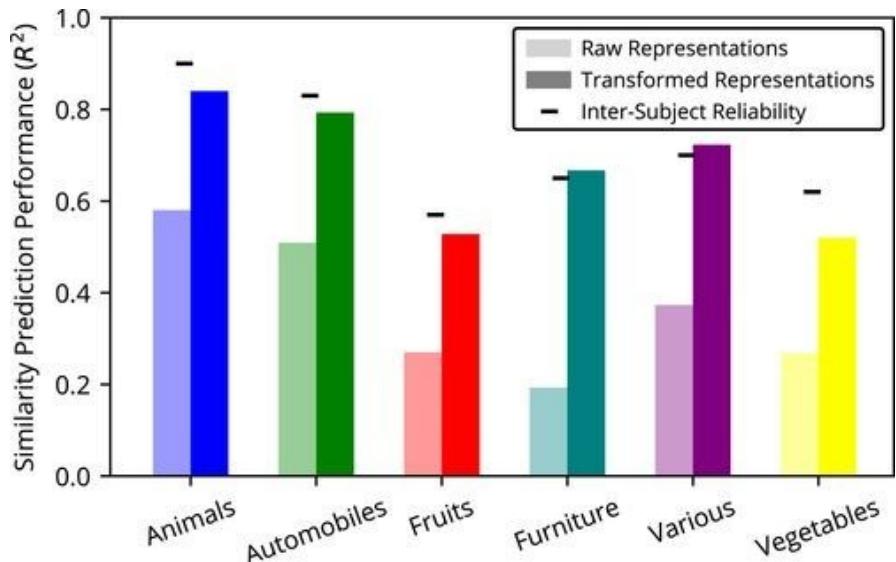
- computing the correlation between the human similarity judgments and the inner products computed in the deep feature spaces. (both 120x120)
- VGG (best performing DNN architecture)



Evaluating the correspondence between representations

Results :

- computing the correlation between the human similarity judgments and the inner products computed in the deep feature spaces. (both 120x120)
- VGG (best performing DNN architecture)
- The raw deep representations provide a reasonable first approximation to human similarity judgments

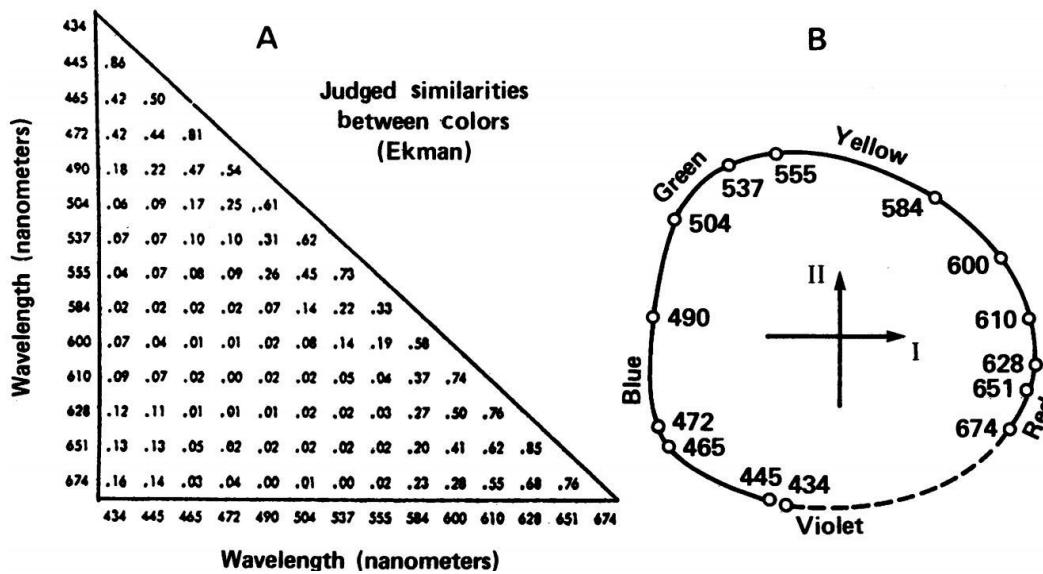


Evaluating the correspondence between representations

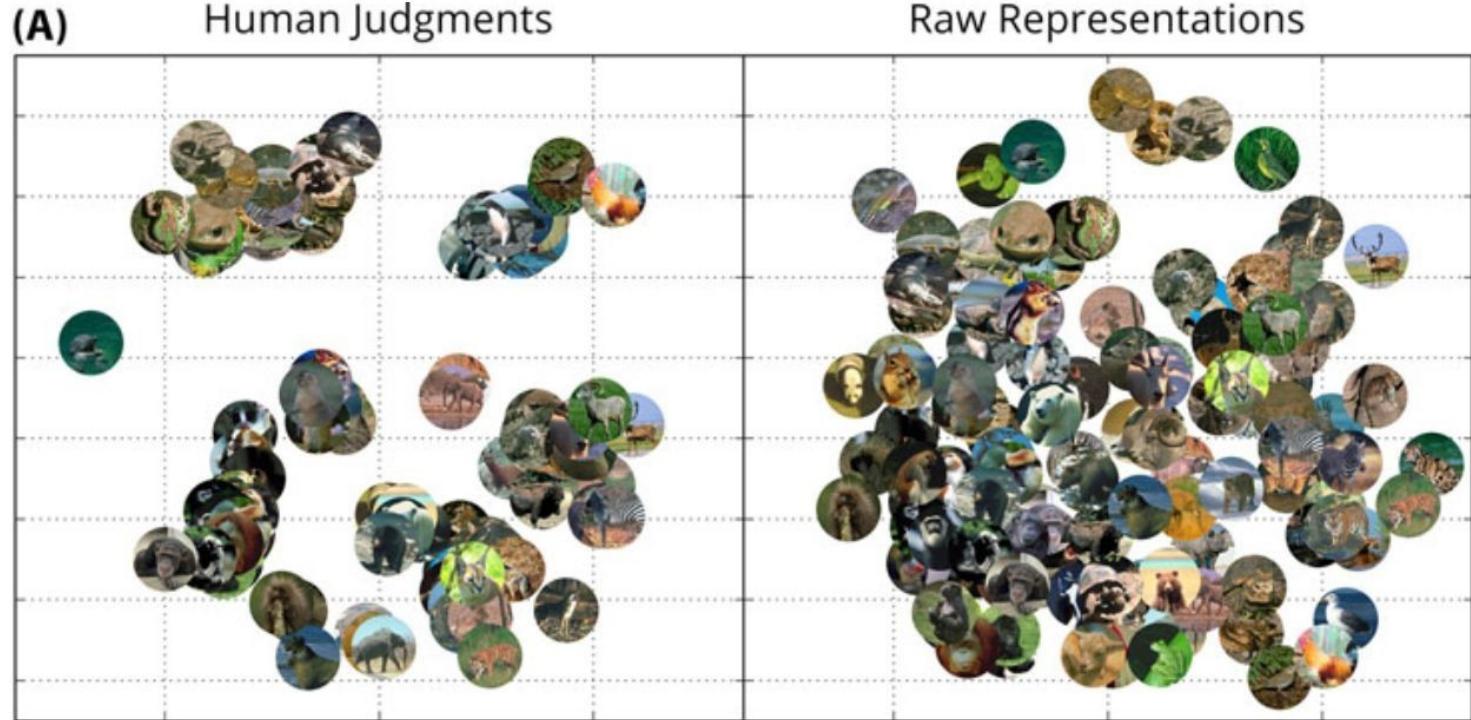
To better understand how DNNs succeed and fail to reproduce the structure of psychological representations, they applied two classic psychological tools:

- non-metric multidimensional scaling
 - converts similarities into a spatial representation
- hierarchical clustering
 - produces a tree structure (dendrogram)

Multidimensional Scaling: from relative distances of n elements to a map of relative distances in 2* dimensions.



Evaluating the correspondence between representations

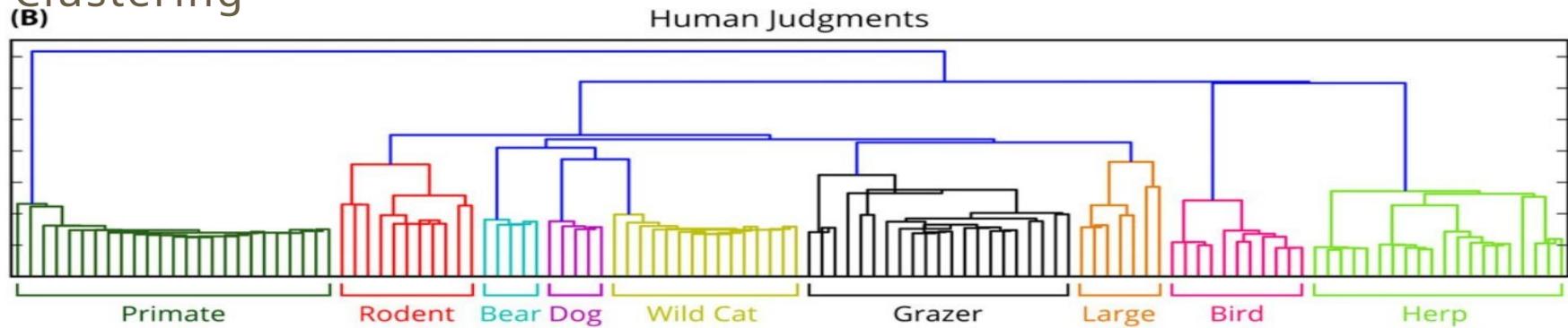


Evaluating the correspondence between representations

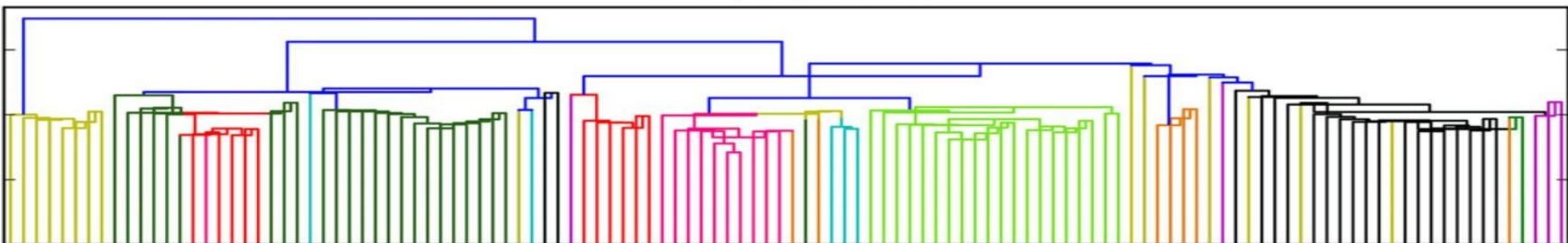
2. Hierarchical Clustering

Clustering

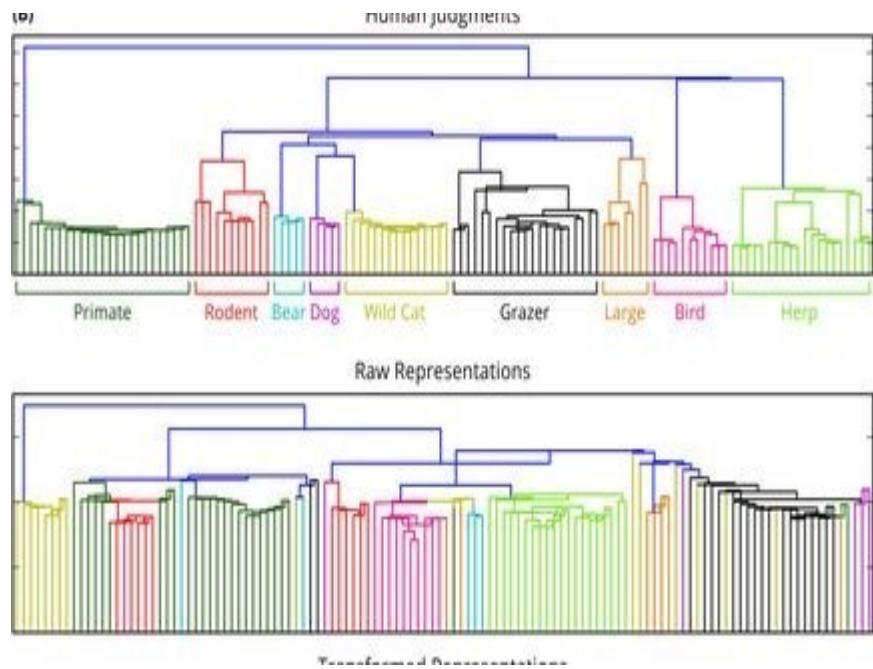
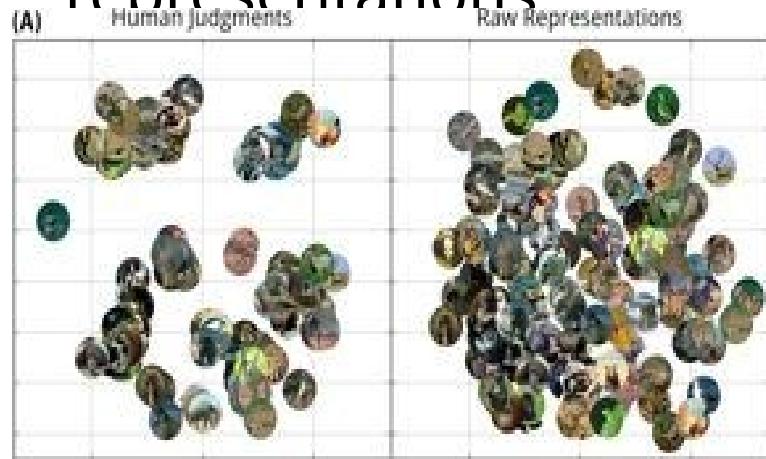
(B)



Raw Representations



Evaluating the correspondence between representations



Human representations exhibit highly distinguished clusters in the spatial projections and intuitive taxonomic structure in the dendograms, neither of which is present in the DNN representations.

Transforming Representations

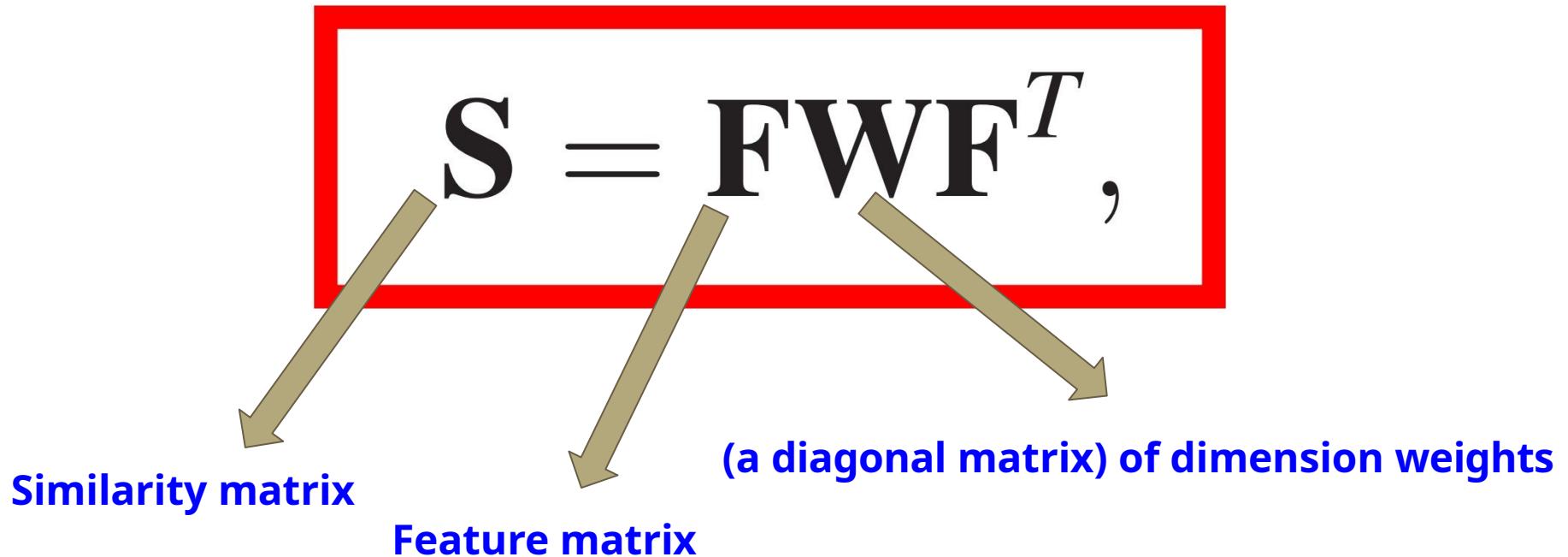
- how can DNN representations be transformed to increase their alignment with psychological representations?
 - with a set of weights on the features used to compute similarity

Transforming Representations

- how can DNN representations be transformed to increase their alignment with psychological representations?
 - with a set of weights on the features used to compute similarity

$$\mathbf{S} = \mathbf{F} \mathbf{W} \mathbf{F}^T ,$$

Transforming Representations



Transforming Representations

- It provides a way to specify the relationship between a feature representation and stimulus similarities.

$$\mathbf{S} = \mathbf{F} \mathbf{W} \mathbf{F}^T ,$$

Transforming Representations

- can be expressed as the **solution to a linear regression problem**
- the predictors for each similarity s_{ij} are the (elementwise) product of the values of each feature for objects i and j
- (i.e., each row of the regression design matrix \mathbf{X} can be written as $\mathbf{F}_i \circ \mathbf{F}_j$)

Dependent Variable → $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

Population Y intercept → β_0

Population Slope Coefficient → β_1

Independent Variable → X_i

Random Error term → ϵ_i

Linear component → $\beta_0 + \beta_1 X_i$

Random Error component → ϵ_i

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} x_{10} & x_{11} & \dots & x_{1, p-1} \\ x_{20} & x_{21} & \dots & x_{2, p-1} \\ \dots & \dots & \dots & \dots \\ x_{n0} & x_{n1} & \dots & x_{n, p-1} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{p-1} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Transforming Representations

- The similarity s_{ij} between objects i and j is therefore modeled as :

$$s_{ij} = \sum_k w_k f_{ik} f_{jk},$$

Weight of that feature

" k " th feature of the image " i "

Similarity between image " i " and " j "

Learning transformations

- Freely identifying the w that best predicts human similarity judgments runs the risk of overfitting, since our DNNs generate thousands of features.

Learning transformations

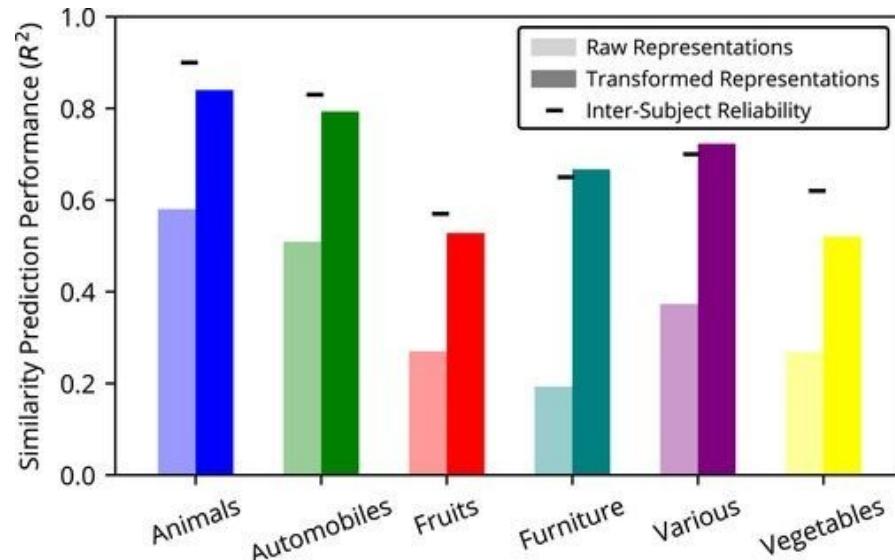
- Freely identifying the \mathbf{w} that best predicts human similarity judgments runs the risk of overfitting, since our DNNs generate thousands of features.
- L2 regularization on \mathbf{w}
 - penalizing models for which the inner product $\mathbf{w}^T \mathbf{w}$ is large

Learning transformations

- Freely identifying the \mathbf{w} that best predicts human similarity judgments runs the risk of overfitting, since our DNNs generate thousands of features.
- L2 regularization on \mathbf{w}
 - penalizing models for which the inner product $\mathbf{w}^T \mathbf{w}$ is large (we want to minimize the square magnitude of weights; produces models with many small weights rather than few large ones.)

Improvement through feature adaptation

- Trained only for the best performing Deep Neural Network, **VGG**
- Variance Explained **doubled** compared to raw representations on all domains



Improvement through feature adaptation

- The stricter case of cross-validation titled ‘CV control’ explained the similarity judgements better than raw representations. In this case no single images occurred in both training folds and test folds of cross-validation
- However, for ‘Transformed model’, the exclusivity was wrt to pair of images.

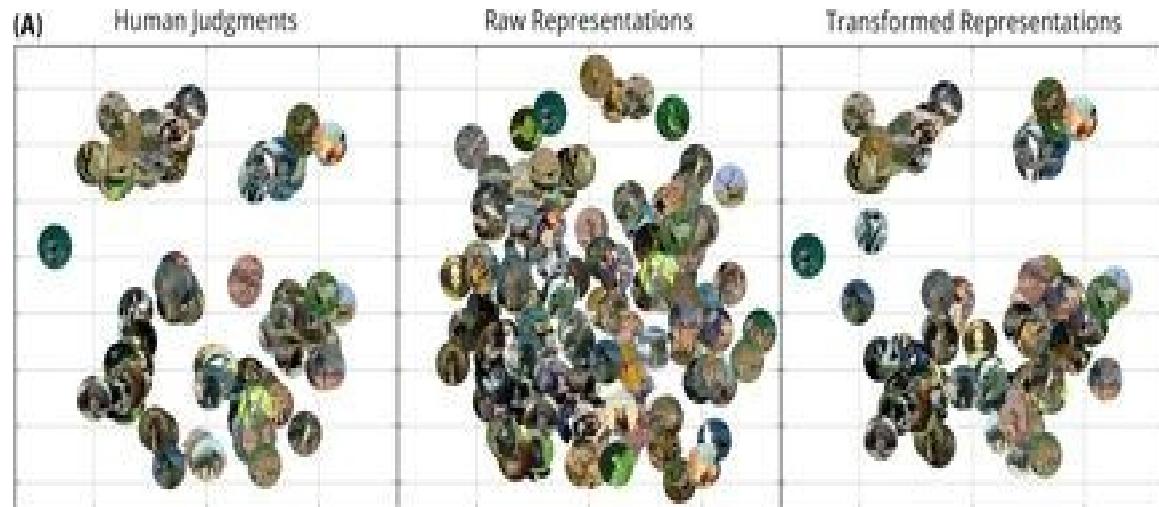
Table 1. Variance explained in human similarity judgments for raw and transformed representations for the best performing network (VGG)

Dataset	Raw R^2	Transformed R^2	CV Control R^2	Human Inter-reliability
Animals	0.58	0.84	0.74	0.90
Automobiles	0.51	0.79	0.58	0.83
Fruits	0.27	0.53	0.36	0.57
Furniture	0.19	0.67	0.35	0.65
Various	0.37	0.72	0.54	0.70
Vegetables	0.27	0.52	0.35	0.62

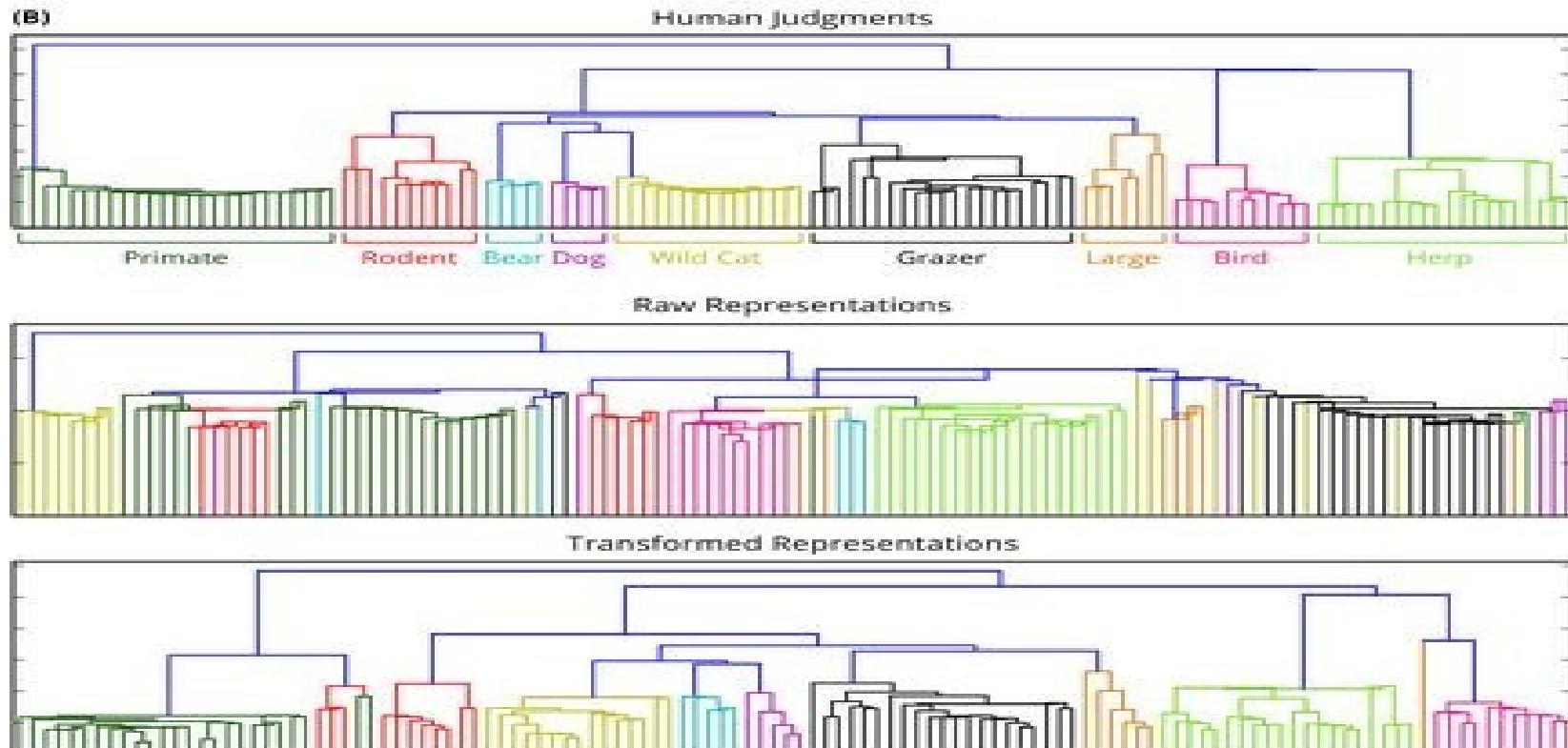
Improvement through feature adaptation

Comparison of the results from MDS :

Stronger resemblance to human judgements



Improvement through feature adaptation

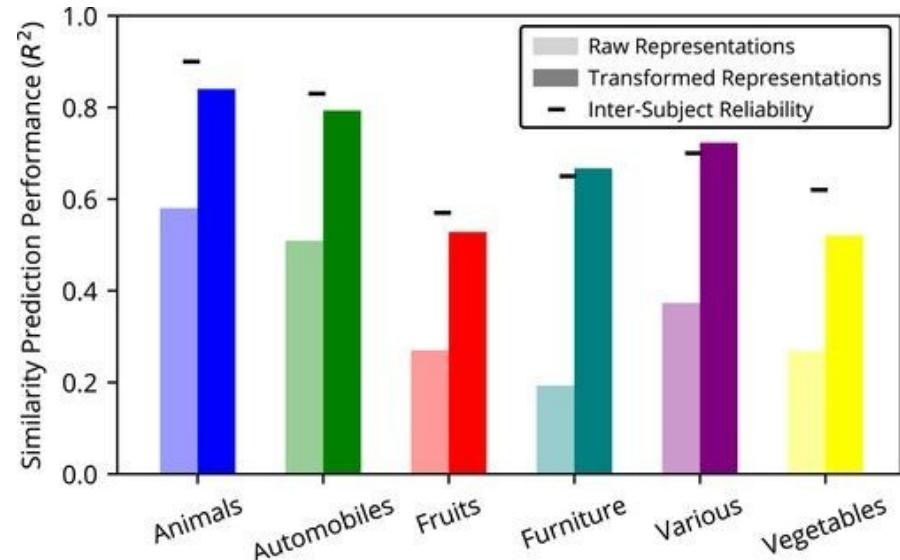


Inter-domain Transfer

- Transformations learnt on one domain do not generalize to others
- Correlations are poor & worse compared to those of

Inter-domain generalization of best performing DNN transformations

Training Set	Test Set	R^2
Animals	Fruits	0.11
Animals	Furniture	0.02
Animals	Vegetables	0.11
Animals	Automobiles	0.17
Animals	Various	0.12



Discussion

- Raw representations from an AI system may approximate well representational space of humans, but MDS and Hierarchy show weakness.
- These representations can be easily transformed to produce better human similarity predictions
- These transformed models can be used to predict similarity judgements for new stimuli

Discussion

- The simple re-weighting of features that the linear transformation performs, can be viewed as an analogue to dimensional attention.
- The ability of transformed representations to generalise for new stimuli empowers studies on cognitive processes relying on such representations in the brain.