003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

LOCAL VS DISTRIBUTED REPRESENTATIONS: WHAT IS THE RIGHT BASIS FOR INTERPRETABILITY?

Anonymous authors

Paper under double-blind review

ABSTRACT

Much of the research on the interpretability of deep neural networks has focused on studying the visual features that maximally activate individual neurons. However, recent work has cast doubts on the usefulness of such local representations for understanding the behavior of deep neural networks because individual neurons tend to respond to multiple unrelated visual patterns, a phenomenon referred to as "superposition". A promising alternative to disentangle these complex patterns is learning sparsely distributed vector representations from entire network layers, as the resulting basis vectors seemingly encode single identifiable visual patterns consistently. Thus, one would expect the resulting code to align better with humanperceivable visual patterns, but supporting evidence remains, at best, anecdotal. To fill this gap, we conducted three large-scale psychophysics experiments collected from a pool of 560 participants. Our findings provide (i) strong evidence that features obtained from sparse distributed representations are easier to interpret by human observers and (ii) that this effect is more pronounced in the deepest layers of a neural network. Complementary analyses also reveal that (*iii*) features derived from sparse distributed representations contribute more to the model's decision. Overall, our results highlight that distributed representations constitute a superior basis for interpretability, underscoring a need for the field to move beyond the interpretation of local neural codes in favor of sparsely distributed ones.

027 028 029

031

1 INTRODUCTION

One of the goals of explainable AI (XAI) in computer vision is to identify the visual features and characterize the representations used by deep neural networks (DNNs) to categorize images (Ribeiro et al., 2016; Sundararajan et al., 2017; Smilkov et al., 2017; Petsiuk et al., 2018; Selvaraju et al., 2017; Linsley et al., 2019; Fel et al., 2021; 2023c;a; Novello et al., 2022; Zhou et al., 2016; Bau et al., 2017; Cammarata et al., 2020b; Kim et al., 2018; Ghorbani et al., 2019). In general, identifying these features requires uncovering the visual patterns that drive the activation of units within a network.

038 This goal is shared with the study of biological vision, where there is extensive research over the last several decades focused on identifying the "preferred stimulus" of individual neurons in the visual 040 cortex (Hubel & Wiesel, 1959; Lettvin et al., 1959; Tsunoda et al., 2001; Wang et al., 1996; Pasupathy 041 & Connor, 2001; Quiroga, 2005). This approach to visual neuroscience reflected the dominant theory 042 at the time, known as the "grandmother (cell)" theory, which postulates that information in the 043 visual system is stored **locally** – at the level of single neurons – and that the visual system contains 044 specific neurons that respond to particular objects or people (including one's own grandmother). Early XAI advances inspired by neuroscience research similarly focused on understanding local representations (Zhou et al., 2016; Bau et al., 2017). This led to the development of more sophisticated 046 optimization methods to synthesize maximally activating images for individual neurons (Erhan et al., 047 2009; Zeiler & Fergus, 2014a; Yosinski et al., 2015; Olah et al., 2017; 2020; Nguyen et al., 2016a;b; 048 Cammarata et al., 2020b). 049

This parallelism between XAI and neuroscience extends beyond vision as recent work has found neurons in multi-modal systems that respond to very high-level concepts beyond simple image appearance, including hand-drawing and text (Goh et al., 2021). Interestingly, the authors identified a "Halle Berry" neuron in CLIP, reminiscent of the neuroscience finding reported two decades ago in the human brain (Quiroga, 2005).



Figure 1: (a) • Local (neuron) versus • Distributed (sparsely distributed vector) visual representations. The activation of individual neurons may be driven by multiple unrelated visual elements (depicted in the images at the bottom) whereas distributed representations, obtained via dictionary learning methods, break down complex patterns into simpler ones corresponding to single visual features. (b) In practice, dictionary learning methods "disentangle" local activations to yield a new vector basis whose activation is driven by single features. The hope for interpretability is that those features align better with the set of features that humans can interpret $S = \{f_1, f_2, ..., f_n\}$.

075 At the same time, a paradigm shift is taking place in neuroscience, where the study of neural 076 populations is quickly superseding the study of single neurons (for a review see Ebitz & Hayden 077 (2021)) because the neural code is believed to be sparse and **distributed** rather than local (Haxby et al., 078 2001; Quiroga et al., 2008; 2013). Interestingly, a similar shift is emerging in XAI, because local 079 representations are known to suffer from the "superposition" problem (Arora et al., 2018; Cheung et al., 2019; Olah et al., 2020; Elhage et al., 2022; Fel et al., 2023b): the number of features captured by DNNs might be larger than the number of neurons. Therefore, the neurons' activations might 081 be driven by multiple unrelated features. To address this challenge, the XAI community has started using dictionary learning methods (Fel et al., 2023c;b; Bricken et al., 2023; Templeton et al., 2024) to 083 project the activations of DNNs onto a new basis of vectors, each activated by a single distinct feature. 084 From an interpretability perspective, representations driven by single features are more desirable 085 because they are expected to be easier to understand by a human, *i.e.*, a unit responding to a single visual pattern –compared to multiple patterns– is inherently easier to interpret. An implicit hypothesis 087 is that applying dictionary learning methods to network activations helps break down complex visual 880 patterns into simpler ones corresponding to single features (Fig. 1.a), which, in turn, might be 089 easier for humans to interpret (Fig. 1.b). Scaling up standard interpretability evaluations to study representations poses significant challenges (Colin et al., 2022). A practical alternative is to assess 090 the ambiguity—or perplexity—of the visual features derived from these representations (Borowski 091 et al., 2021) (see 3.2.1 for a more thorough elaboration on this point). 092

Despite the growing consensus that distributed representations constitute a stronger basis for interpretability compared to local ones (single neurons), empirical evidence remains scarce. This paper aims to fill this gap. Specifically, we contend that a representation can be considered superior if the features derived from it are more intelligible, *i.e.*, easier for humans to make sense of while being demonstrably used by the model in its decision-making process. By means of computational and psychophysics experiments, we set out to find which of the local vs. distributed representations better meets these two conditions.

¹⁰⁰ In sum, the main contributions of this paper are as follows:

055

056

060 061 062

063

064

065

066 067

068

069

071

072

- We conduct three large-scale psychophysics experiments for a total of 15,720 responses from a pool of 560 participants, to evaluate the visual ambiguity of the features derived from local vs distributed representations (see Fig 1). In the process, we identify a potential semantic bias in the experimental protocols commonly used in the field (Borowski et al., 2021; Zimmermann et al., 2023), and provide an approach to at least partially mitigate it.
- Our findings provide strong evidence that (*i*) features derived from distributed representations are significantly easier for humans to interpret than features derived from local

representations, and (*ii*) this effect is even more pronounced in the deepest layers of a neural network.

• Additionally, we observe that (*iii*) models rely significantly more on features derived from distributed representations compared to those derived from local representations. Overall, our results suggest that distributed representations provide a substantially better foundation for the interpretability of models than local representations.

2 RELATED WORK

117 From local to distributed representations in XAI. Early research on explainable AI (XAI) 118 in computer vision developed attribution methods (Zeiler et al., 2011; Zeiler & Fergus, 2014a;b; 119 Sundararajan et al., 2017; Smilkov et al., 2017; Fong & Vedaldi, 2017; Ancona et al., 2018; Shrikumar 120 et al., 2017; Chattopadhay et al., 2018; Fel et al., 2021; Novello et al., 2022) to understand specific 121 model predictions. These methods predominantly aimed to identify "where" are the most important 122 pixels of an image, given a specific model prediction. Unfortunately, these methods fell short in 123 explaining the "what" (Kim et al., 2018; 2022; Taesiri et al., 2022; Colin et al., 2022), i.e., the visual features that the models rely on to make their predictions. As a result, new XAI methods, 124 including Feature Visualization approaches (Nguyen et al., 2016b; Olah et al., 2017), were developed 125 to provide insights into the features that neurons or model layers respond to. These efforts revealed 126 that the neurons' activations can be driven by visually distinct features (Nguyen et al., 2016b; 2019; 127 Cammarata et al., 2020a; Bricken et al., 2023). A similar behavior was observed in other fields, and 128 notably in Natural Language Processing (NLP) (Elhage et al., 2022). 129

- The neurons' tendency to respond to multiple unrelated visual elements indicates that single neurons 130 might not align well with the model's internal representations. This phenomenon, known as super-131 position (Arora et al., 2018; Cheung et al., 2019; Olah et al., 2020; Elhage et al., 2022; Fel et al., 132 2023b) or feature collapse (Fel et al., 2023c), suggests that there could be significantly more features 133 than neurons in a model. Consequently, interpreting neurons might be no more meaningful than 134 interpreting arbitrary directions in the feature space. To achieve effective model interpretation, it 135 is, therefore, essential to identify an interpretable basis that facilitates the extraction of meaningful 136 features. This observation has led to increased interest, over the past five years, towards examining 137 deep learning models by considering the distributed nature of latent space representations. More 138 specifically, it spurred the development of concept extraction methods (Ghorbani et al., 2019; Zhang 139 et al., 2021; Fel et al., 2023c; b; Graziani et al., 2023; Vielhaben et al., 2023) that leverage dictionary learning, especially overcomplete dictionaries. In NLP, sparse autoencoders (SAEs) (Elhage et al., 140 2022; Bricken et al., 2023; Cunningham et al., 2023; Tamkin et al., 2023) are extensively studied 141 and regarded as a promising direction for discovering interpretable bases. In this paper, we are 142 interested in comparing the suitability of local vs distributed representations to serve as a basis for 143 the interpretability of deep learning models in computer vision. 144
- 144 145

111

112

113

114 115

116

Human-evaluation of interpretability. Since the ultimate goal of XAI is to make complex models
 understandable to humans, it is essential to incorporate human evaluation in measuring interpretability
 through psychophysics experiments. Human evaluation serves as a critical benchmark for assessing
 whether the explanations provided by XAI systems are comprehensible, useful, and actionable for
 end-users. This human-centric approach ensures that interpretability methods are not just theoretically
 sound but also practically effective in real-world applications.

151 To the best of our knowledge, Borowski et al. (2021) were the first to quantify the interpretability 152 of features through psychophysics experiments. Their approach involved visualizing unit responses 153 by contrasting maximally and minimally activating stimuli. Specifically, Borowski et al. (2021) 154 focused on studying features from local representations (i.e., single unit activations) of an Inception 155 V1 (Szegedy et al., 2015) trained on ImageNet (Deng et al., 2009; Russakovsky et al., 2015). In their 156 experiments, participants were shown maximally and minimally activating images from the ImageNet 157 ILSVRC 2012 validation set (Russakovsky et al., 2015) to illustrate a specific feature. They were then 158 asked to select which of two query images also activated the feature (see Fig. 2 for an example trial). 159 As a proxy for interpretability, they measured the visual coherence of maximally activating stimuli of a given feature. It is maximal when the visual feature is unambiguous. The main conclusion from 160 these experiments is that, from a human-centric perspective, maximally activating natural images 161 are more effective for studying features than synthetically generated feature visualizations (Olah

162 et al., 2017). Zimmermann et al. (2021) proposed a variation of this task to investigate if humans 163 gained causal insights from those visualizations. They tested the participants' ability to predict 164 the effect of an intervention, such as occluding a patch of the image, on the activation of the unit. 165 Both by means of a large-scale crowdsourced psychophysics experiment and measurements with 166 experts, they found that synthetic feature visualizations (Olah et al., 2017) helped humans perform the task successfully. However, these visualizations did not provide a significant advantage over 167 other visualizations, such as exemplary images. Finally, Zimmermann et al. (2023) extended the 168 work of Borowski et al. (2021) to a broader set of deep learning architectures used for computer vision tasks, including a ResNet50 (He et al., 2016). Their main conclusion is that increasing the 170 scale of the model does enhance the interpretability of the features. While both our work and that 171 of Zimmermann et al. (2023) build upon the experimental protocol of Borowski et al. (2021), they 172 focus on scaling Borowski et al. (2021)'s insights, whereas we adapt the protocol to generate insights 173 into a different and novel research question: comparing the suitability of features derived from local 174 vs distributed representations for interpretability. 175

Interestingly, in a similar vein but within the context of NLP, (Bricken et al., 2023) compare the interpretability of 162 features extracted from both local and distributed representations. One of the authors scored each feature using samples drawn uniformly across the spectrum of activation. Their findings indicate that features obtained from distributed representations are substantially more interpretable than those from single neurons. The research presented in this paper differs from this work in the application domain, the experimental design, single participant in their study vs 15,720 responses from a pool of 560 participants in our study.

182 183

184 185

187 188 189

3 Methodology

In this section, we first provide an overview of the technical methods used to create the conditions for the three psychophysics experiments described below.

3.1 TECHNICAL METHODS

Model All the experiments described in this paper were performed on a ResNet50 (He et al., 2016)
 from the Torchvision (Marcel & Rodriguez, 2010) library, pre-trained on ImageNet-1k (Deng et al., 2009).

194 **Sparse Distributed Representations** To compute sparse distributed representations for the dis-195 tributed condition in our psychophysics experiments (see Section 3.2), we employed CRAFT (Fel 196 et al., 2023c), an off-the-shelf dictionary learning method. Specifically, given a model $f: \mathcal{X} \to \mathcal{A}$ 197 that maps from an input space $\mathcal{X} \subseteq \mathbb{R}^d$ to an activation space $\mathcal{A} \subseteq \mathbb{R}^p$ (e.g., any layer of the network), 198 we compute the activations $\mathbf{A} = f(\mathbf{X}) \in \mathbb{R}^{n \times p}$, where $\mathbf{X} = [x_1, \dots, x_n]^\top \in \mathbb{R}^{n \times d}$ represents a 199 set of n input data points. Each row $a_i \ge 0$ of A contains the non-negative activations for a given 200 data point x_i , due to the use of ReLU activations. To obtain the sparse distributed representations, 201 Non-negative Matrix Factorization (NMF) is applied to approximate A as:

$$\left(\mathbf{Z}^{\star}, \mathbf{D}^{\star}\right) = \underset{\mathbf{Z} \geq \mathbf{0}, \ \mathbf{D} \geq \mathbf{0}}{\arg\min} \left\| \mathbf{A} - \mathbf{Z} \mathbf{D}^{\top} \right\|_{F},$$

where $\|\cdot\|_{F}$ denotes the Frobenius norm. Here, $\mathbf{Z} \in \mathbb{R}^{n \times k}$ are the **codes** (or concept coefficients), and $\mathbf{D} \in \mathbb{R}^{p \times k}$ forms the **dictionary** (or concept bank). Both \mathbf{Z} and \mathbf{D} are constrained to have nonnegative entries and tend to be sparse due to the properties of NMF. The dictionary matrix \mathbf{D} provides a new set of basis vectors (concepts) aligned with the activation patterns of the neural network, while \mathbf{Z} contains the coefficients representing the original activations \mathbf{A} in terms of these concepts. CRAFT is particularly well suited for our purposes as it has been shown to extract interpretable features from deep neural networks (Fel et al., 2023c).

212

202 203 204

Feature importance For an image x_i , let v_c be the vector it activates most strongly from an intermediate representation—it can be a sparse distributed vector in **D** or a neuron in \mathcal{A} —, and let y_i be the logit score for x_i . To assess v_c 's importance, we perform an ablation by setting the activation along this dimension to zero for x_i , resulting in modified activations \mathbf{a}'_i . These modified



Figure 2: Illustration of a trial. Example of a trial in our study corresponding to Experiment I, distributed representation condition of a unit located in layer2.0.bn2. Two panels of 9 reference 232 images are located on the left and right-hand side of the display, separated by 2 query images in the center. Participants were asked to select the query image they believed shared the same visual 234 elements as the reference images displayed on the right panel, corresponding to maximally activating 235 stimuli. The less ambiguous this shared visual element is—the more visually coherent the set of images-the more likely participants are to select the correct query. In this case, the correct query is the bottom image depicting a yellow tram. 238

activations are then propagated to the output layer to compute a new logit score, y'_i . One can assess the importance of \mathbf{v}_c by measuring as the difference in logit scores:

$$\Delta \boldsymbol{y}_i = \boldsymbol{y}_i - \boldsymbol{y}_i' \tag{1}$$

The more important the feature driving the activation of \mathbf{v}_c is for the model, the higher the drop in logit score.

3.2 **PSYCHOPHYSICS STUDIES**

231

233

236

237

239 240 241

242

243 244 245

246

247 248

249

251

250 3.2.1 EXPERIMENTAL PROTOCOL

Interpretability is a human-centric attribute. Hence, to shed light on how well local vs distributed 252 representations serve as a basis for the interpretability of deep neural networks, we performed 253 three large-scale online psychophysics experiments. The three experiments focus on evaluating 254 the interpretability of features extracted from both types of representations. In practice, given the challenges of scaling up standard interpretability evaluations to study representations (Colin et al., 256 2022), we adapt the experimental protocol of Borowski et al. (2021) to measure the ambiguity, or 257 perplexity, of the visual feature as a proxy for its interpretability. More precisely, given a feature and 258 a set of images that illustrate it—e.g., maximally activating images—, this protocol measures how 259 visually coherent humans find this set of images. The more visually coherent it is, the more likely 260 they are to correctly identify another member of this set, namely the correct query image. Hence, in 261 our results, we report the visual coherence of a feature as the proportion of participants that correctly identified the query image. This visual coherence should be maximal when the feature represents a 262 single visual pattern that people can interpret. 263

264 Each participant was assigned to one of two conditions: local or distributed representation. After 265 successfully finishing a practice session composed of 9 trials, participants were asked to sequentially 266 perform 40 different trials of the same task, where each trial corresponded to a different network unit. For each selected network unit and trial, participants were shown on the left and right-hand side of 267 the screen two panels of 9 reference images each, separated by 2 images in the center, which were 268 denoted as *queries*, as illustrated in Figure 2. The task across all experiments consisted of selecting 269 the query image that participants believed shared the same visual elements as the 9 reference images

displayed on the right panel. The selection of the reference images shown in each panel differed depending on the experiment, as explained below.

273 3.2.2 NEURON AND ST

274

3.2.2 NEURON AND STIMULI SELECTION

Neuron selection In line with the literature (Zimmermann et al., 2023), we aimed to obtain a 275 representative sample of units within the different layers of the ResNet50. As a starting point, we 276 selected the 80 units reported in (Zimmermann et al., 2023) for a ResNet50 (He et al., 2016) model. 277 Given that CRAFT requires positive activations, we replaced the neurons from the convolutional 278 layers (which could have non-positive activations) to their equivalent in the following batchnorm 279 layer (layer 1.1.conv1 neuron 53 -> layer 1.1.bn1 neuron 53). As a result, we selected 80 neurons 280 distributed across 43 layers that are grouped in four blocks. Thus, the names of the layers used in 281 this paper (layer1 through layer4) are directly obtained from the modules of the ResNet50 Pytorch 282 implementation. 283

284 Stimuli selection for the local representation condition For each of the selected neurons, we 285 identified 2,900 images from the validation set of ImageNet ILSVRC 2012 (Russakovsky et al., 2015) obtained as follows: the 2,500 most strongly activating images and the 400 least strongly 286 activating images. Following the procedure described by Borowski et al. (2021), we illustrated a 287 visual feature through both maximally activating (images that possess the feature, see Fig. 2 right 288 panel) and minimally activating (images that do not possess the feature, see Fig. 2 left panel) stimuli. 289 Again, like in previous work (Borowski et al., 2021), for the former, we selected a random sample of 290 9 images from the top 150 images; for the latter, we uniformly sampled 9 images from the bottom 291 20 images. We created 10 different trials per feature following this procedure to ensure image 292 independence in the results. When participants were assigned to a feature, we randomly sampled one 293 of those 10 trials to illustrate the feature. 294

295 Stimuli selection for the distributed representation condition Our hypothesis is that applying 296 dictionary learning methods to local representation allows us to disentangle the superposition of 297 features that might drive the activation of a neuron. Hence, to test this hypothesis, all the distributed representations studied in this work are obtained from the local representation. In practice, for 298 each neuron in the local condition, we selected the top 300 maximally activating images, *i.e.*, those 299 that most strongly led to the activation of the neuron. We used CRAFT to identify a new basis of 300 distributed vectors through which their activations can be expressed. Following the recommended 301 procedure in (Fel et al., 2023c), we constrained the dimensionality of this new basis to d = 10. From 302 these 10 directions, we selected the one that was the most frequently the most activated across the 303 300 images. Once a direction was chosen, we ranked the 2,900 images according to their activation 304 in that direction to obtain the images that illustrate a visual feature in the distributed representation. 305

306 3.2.3 EXPERIMENTS 307

In this section, we present the details of each of the three psychophysics experiments conducted in this work.

310

Experiment I This experiment is an adaption of the methodology proposed by (Borowski et al., 2021) with two conditions: local vs distributed representations. An illustration of a trial belonging to Experiment I can be found in Figure 2. In this experiment, we performed a total of 1,600 trials (80 neurons × 2 conditions × 10 trials per unit).

314

315 **Experiment II** The main objective of Experiment II was to control for potential semantic con-316 founding variables present in Experiment I. If the reference images that possess/do not possess the 317 feature of interest belong to very different semantic groups (classes), then it could be possible to solve 318 the task through simple semantic grouping. Figure 3 illustrates this phenomenon: in the example 319 depicted in the figure, it is easier to solve the trial by inferring that the feature of interest is not about 320 a monkey than it is to infer the actual visual feature present in the reference images. The goal of 321 Experiment II is to control for this potential semantic confound as follows: given a set of reference images that possess the feature of interest, we extract their semantic labels from ImageNet and aim to 322 find within the 400 minimally activating images a set of 10 images that share the same distribution 323 of semantic labels. We define four levels of semantic similarity (level 0 to level 3). In level 0, the



Figure 3: **Illustration of the role of semantics**. Example of a trial from Experiment I in the local representation condition. In this case, the task can be trivially solved by simply relying on semantics. By observation of the minimally activating stimuli (left panel), it is easy to conclude that the neuron of interest is not a monkey detector, yet, it is hard to articulate what is the visual feature captured by the neuron (images in the right panel).

labels from ImageNet are used to determine the semantic similarity, whereas in *levels* 1 through 3 we obtain the labels by moving up one branch in the WordNet (Fellbaum, 2010) hierarchy. In practice, we performed an iterative search for each trial starting from level0. If there are not 10 images in the 400 minimally activating set that share labels with the reference images at a given level of semantics, the process is restarted, this time using labels from the next and broader level of semantics. We continued this process until 10 images from the minimally activating set shared the distribution of semantic labels with the reference images. In cases where it was impossible to identify 10 images, the feature was excluded from the experiment. Only one feature was fully excluded for both the local and distributed conditions due to this factor.

Experiment III Finally, recognizing that the adopted experimental protocol was originally designed
 to assess whether features are better understood using natural or synthetic visualizations, we devised
 Experiment III to explore the impact of combining natural images with feature visualizations (See
 more details in Section D).

3.2.4 PARTICIPANTS

A total of 560 participants were recruited to take part in Experiments I, II, and III through the online platform Prolific¹. All participants were native English speakers who reported not being visually impaired and completed the study on a laptop or desktop computer (not a mobile phone). They provided informed consent electronically and were compensated \$2.75 for their time ($\sim 10 - 13$ min). The protocol was approved by the University IRB. Based on the power analysis of (Zimmermann et al., 2023), a minimum of 60 participants per condition (120 participants per Experiment) was needed to obtain statistically robust results. Furthermore, participants were required to (1) succeed in at least 5 of the 9 practice or instruction trials and (2) correctly answer at least 4 of the 5 attentiveness tests (catch trials) that were randomly inserted in the experiment. As a result, we analyzed the data corresponding to 138, 133, and 122 participants from Experiments I through III, respectively.

¹www.prolific.com



Figure 4: **Per-layer results for Experiment I (a) and II (b)**. Given a feature and a set of images to illustrate it, we assess how visually coherent participants find this set of images—or how unambiguous the feature is. More precisely, we measure the proportion of participants that are able to identify the query image which is also part of this set of images. In both experiments, a clear trend emerges where features appear significantly less ambiguous in the distributed representation than in the local representation condition, particularly in the deeper layers of the network.

4 Results

396 397 398

399

400

394

In this section, we summarize the results obtained from analyzing the responses from the 138, 133, and 122 participants who successfully completed Experiments I, II, and III, as well as the results from our feature importance analysis.

401 Unless stated otherwise, all of our behavioral data analyses consist of generalized logistic mixed-402 effects regressions (GLMER), with trial accuracy (1 vs. 0) as the dependent variable. The random-403 effects structure included both a by-participant and a by-unit ² random intercept, as well as a by-unit 404 random slope for the condition variable. We used the lme4 (Bates et al., 2015) package in R (R 405 Core Team, 2021) to fit the models and lmerTest to obtain p-values for the fixed effects, with an 406 α -level of 0.05 for statistical significance.

407

Reproduction of previous results. Given that the experimental protocol used in Experiment I, with the local representation condition, is the same as the one described in (Zimmermann et al., 2023), we first evaluate to which degree our results corroborate previous research. For the ResNet50, Zimmermann et al. (Zimmermann et al., 2023) report an average task performance of $83.0\% \pm 2.0^3$. In our experiment, we obtain an average performance of $78.8\% \pm 1.5$. Given the results' similarity and that the selected units are not exactly the same (as described in Section 3), we conclude that Experiment I reproduces previously reported findings regarding local representations. This result also serves as an external validation of our experimental protocol.

415 416

Human performance is superior in the distributed representation condition. Based on our main 417 hypothesis that a distributed representation constitutes a better basis for interpretability than a local 418 representation, we predicted that participants would be better at selecting the maximally activating 419 query image in the distributed condition. In Experiment I, the average performance across participants 420 in the distributed condition was $83.5\% \pm 1.4^4$, when compared to $78.8\% \pm 1.5$ in the local condition. A 421 GLMER with condition as a predictor revealed a statistically significant disadvantage for the local 422 condition: $\beta_{condition} = -0.47, SE = 0.23, z = -2.04, p = 0.04$. This result was corroborated 423 both in Experiments II (see examples of trials in Fig 8, 9, 10) and III, with a mean performance 424 of participants in the distributed condition of $85.1\% \pm 1.4$ vs. $76.2\% \pm 1.6$ in the local condition, 425 $\beta_{condition} = -0.93, SE = 0.24, z = -3.93, p = <.001$ and a mean performance of $80.0\% \pm 1.6$ (distributed) vs. 74.1% \pm 1.7 (local), $\beta_{condition} = -0.42, SE = 0.17, z = -2.47, p = 0.01,$ 426 respectively. 427

428 429

430

²Here, a unit refers to either a neuron in the local representation condition or to a specific direction of the dictionary in the distributed representation condition.

³Values inferred from Figure 3 in their paper.

⁴The values reported correspond to a 95% confidence interval.

432 Semantic control matters, but not significantly. Figures 4a and 4b depict the results obtained 433 without (Experiment I) and with (Experiment II) a semantic control applied to the stimuli, re-434 spectively. At first glance, the results seem consistent in both scenarios, except for the results 435 corresponding to neurons in layer 4 where the performance of participants in the local condition 436 suffers a drastic decrease in Experiment II. This result is coherent with the intuitive observation that the model aggregates class-specific information the closer the layer is to the output layer. Hence, 437 neurons located in deeper layers are more likely to respond highly to features that correspond to 438 certain classes and not at all to features belonging to other classes, *i.e.*, the semantic confounds 439 are expected to be higher for those neurons. To investigate further the role of controlling the se-440 mantics, we pooled the data from Experiments I and II and included both experiment and an 441 interaction term condition:experiment as predictors in a GLMER. We found a significant 442 main effect for condition, $\beta_{neuron} = -0.47, SE = 0.23, z = -2.04, p = 0.04$, but not for 443 experiment. Interestingly, we did not obtain evidence for a significant interaction between 444 experiment and condition, $\beta_{neuron: Exp2} = -0.46, SE = 0.30, z = -1.53, p = 0.13$. As 445 it was not possible to semantically control all trials in Experiment II to the same extent, we also 446 tested for an interaction between condition and semantic control instead of the variable experiment. Semantic control was coded as a categorical variable with 4 levels: no con-447 trol, 1, 2, and 3, corresponding to the previously described levels in Section 3.2.3. We compared 448 two GLMERs: one including only the main effects of condition and semantic control, 449 and another also considering the interaction between these two variables. Based on the Akaike 450 Information Criterion, the model without the interaction was preferred (9116 vs. 9120). In sum, 451 while the nominal disadvantage of the neuron condition is larger in Experiment II (semantic con-452 trol) than in Experiment I (no semantic control), none of the performed statistical tests yielded 453 a significant difference between these two experiments. This result can be partially explained 454 by the unequal distribution of trials across semantic levels (see fig 11) and by the quality of the 455 semantic control achieved once getting to a certain level in the Word-Net hierarchy. We leave 456 to future work a finer-grained study of the role of semantics in the interpretability of features. 457

458

481

The deeper the layer, the more prominent the benefits of the distributed representation. Figures (a, 4b and 7 illustrate the per layer re

461 Figures 4a, 4b and 7 illustrate the per-layer re-462 sults obtained in Experiments I, II and III, re-463 spectively. We find evidence to suggest that 464 the benefits of the distributed representation in-465 crease with the depth of the layer from which 466 we select the unit that the participants were interpreting. Indeed, we identified a significant 467 main effect for unit and also for an interaction 468 term unit: condition, *i.e.*, the advantage of 469 the distributed condition increases as the units 470 belong to deeper layers. This result is consistent 471 across the 3 experiments: $\beta_{depth} = 0.08, SE =$ 472 0.02, z = 3.92, p < .001 and $\beta_{depth:Exp1} =$ 473 -0.06, SE = 0.03, z = -2.22, p = 0.03474 for Experiment I; $\beta_{depth} = 0.07, SE =$ 0.02, z = 3.27, p = 0.001 and $\beta_{depth:Exp2} =$ 475 -0.09, SE = 0.03, z = -3.6, p < .001 for 476 Experiment II; and $\beta_{depth} = 0.05, SE =$ 477 0.02, z = 2.6, p = 0.01 and $\beta_{depth:Exp3} =$ 478 -0.06, SE = 0.03, z = -2.59, p = 0.01 for 479 Experiment III. 480



Figure 5: Feature importance. We measure the importance of a feature as the average drop in logit score Δy for the 300 most activating images when the feature is occluded. Except for *layer*1, we find that the model relies significantly more on features derived from the distributed representation than on features from local representations, z = -5.86, p < .001 (Mann-Whitney U test).

The model relies more on features derived from the distributed representation. We evaluate the importance of every feature used in our psychophysics experiment by measuring the average drop in logit score for the 300 most activating images when the feature is occluded (see Eq. 1). Overall, we find that the model relies significantly more on features derived from distributed representations than from local representations: z = -5.86, p < .001 (Mann-Whitney U Test), the only exception being *layer*1 (Fig 5). Interestingly, and in line with the rest of our results, the deeper the layer, the larger the difference between the local vs distributed representations.

5 DISCUSSION

489

490

501

491 **Conclusion** In this work, we have investigated and compared the suitability of local and distributed 492 representations to serve as a basis for the interpretability of deep neural networks. We contend that 493 the best-suited basis possesses a dictionary of features that is (a) more aligned with the set of features 494 that human observers can interpret while (b) being demonstrably more important for decision-making 495 by the model. Through psychophysics experiments, we consistently find that features derived from 496 distributed representations are easier for humans to interpret, particularly when features are derived 497 from deeper layers in the network. Equally important, we find that the model relies significantly 498 more on those features to make its decisions. To the best of our knowledge, our results provide 499 the first empirical evidence of the superiority of distributed representations over local ones for the interpretability of deep neural networks. 500

Limitations Our study is not exempt from limitations. First, the methodology proposed 502 by Borowski et al. (2021) and followed by Zimmermann et al. (2023) utilizes the entire ImageNet 503 validation set with the goal of studying a broad range of stimuli and thereby increasing the likelihood 504 of identifying stimuli that are representative of neurons' selectivity. While such motivation is sound, it 505 poses a significant challenge when performing psychophysics experiments: neurons that are selective 506 to class-specific features (e.g., fish scale) will be maximally activated when presented with stimuli 507 that belong to the corresponding classes (e.g., fish) and minimally activated when provided with 508 stimuli that belong to other classes (e.g., dogs). In those cases, the task can be solved trivially using 509 semantics. The design of Experiment II reflects an initial attempt to mitigate this confounding factor. 510 However, both quantitative results and a manual exploration of the trials by the authors hint at only a 511 partial success of the adopted methodology to address this challenge. We leave to future research a 512 further refinement of the experimental protocol.

513 Second, while the original work (Borowski et al., 2021) this manuscript builds on aimed at measuring 514 the *interpretability* of features, we have been more conservative with the terminology employed 515 in this paper as we believe that the current experimental protocol only measures a loose proxy 516 of interpretability. Interpretability, or understanding, of a complex visual pattern usually requires 517 procedural learning with feedback (Ashby & Maddox, 2005). Given that, in practice, this protocol 518 measures a score based on a set of images presented in a single trial, it seems unreasonable to expect this score to capture interpretability. In contrast, it seems more reasonable to expect this score to 519 characterize something more visual, like the ambiguity of the features, i.e., how complex it might 520 be to interpret the features if one had to. Such ambiguity can be inferred, at least in part, by the 521 visual coherence of the set of images. Zimmermann et al. have indeed shown that responses from 522 participants can be accurately predicted based on visual coherency-or perceptual similarity-of a 523 trial. 524

While we believe that these metrics are useful proxies and that features achieving a high score at them are desirable for interpretability, we want to emphasize that we do not claim that it means that they *are* interpretable. Perceptual similarity differs from interpretability; being able to group images correctly does not necessarily imply a deep understanding of the underlying factors driving those groupings or the ability to explain why the images belong together based on intrinsic features or causal relationships.

Nevertheless, we believe that the adopted experimental protocol serves as a solid foundation for
 developing a scalable and robust paradigm to measure the interpretability of features in deep learning
 models. We leave to future work the development of such a paradigm.

- 534
- 535
- 536
- 53
- 538
- 539

540 ACKNOWLEDGEMENT

541

This work was funded by the ONR grant (N00014-24-1-2026), NSF grant (IIS-2402875 and EAR-1925481) and the ANR-3IA Artificial and Natural Intelligence Toulouse Institute (ANR-19-PI3A-0004). The computing hardware was supported in part by NIH Office of the Director grant #S100D025181 via the Center for Computation and Visualization (CCV) at Brown University.
J.C. and N.O. have been partially supported by Intel Corporation and funding from the Valencian Government (Conselleria d'Innovació, Industria, Comercio y Turismo, Direccion General para el Avance de la Sociedad Digital) by virtue of 2022-2023 grant agreements (convenios singulares 2022, 2023). J.C. has also been partially supported by a grant by Banco Sabadell Foundation.

- 550
- 551 552

561

562

563

567

568

569

570

553 REFERENCES

- Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics*, 6:483–495, 2018.
 - F Gregory Ashby and W Todd Maddox. Human category learning. *Annu. Rev. Psychol.*, 56(1): 149–178, 2005.
- Douglas Bates, Martin Maechler, Ben Bolker, Steven Walker, Rune Haubo Bojesen Christensen,
 Henrik Singmann, Bin Dai, Gabor Grothendieck, Peter Green, and M Ben Bolker. Package 'lme4'.
 convergence, 12(1):2, 2015.
 - David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference* on computer vision and pattern recognition, pp. 6541–6549, 2017.
- Judy Borowski, Roland S Zimmermann, Judith Schepers, Robert Geirhos, Thomas SA Wallis,
 Matthias Bethge, and Wieland Brendel. Exemplary natural images explain cnn activations better
 than state-of-the-art feature visualization. In *International Conference on Learning Representations*,
 2021.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
 Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,
 Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina
 Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and
 Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary
 learning. *Transformer Circuits Thread*, 2023. https://transformer-circuits.pub/2023/monosemanticfeatures/index.html.
- Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah, Michael Petrov, Ludwig Schubert, Chelsea
 Voss, Ben Egan, and Swee Kiat Lim. Thread: Circuits. *Distill*, 2020a.
- Nick Cammarata, Gabriel Goh, Shan Carter, Ludwig Schubert, Michael Petrov, and Chris Olah. Curve detectors. *Distill*, 2020b. doi: 10.23915/distill.00024.003. https://distill.pub/2020/circuits/curve-detectors.
- Aditya Chattopadhay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Gradcam++: Generalized gradient-based visual explanations for deep convolutional networks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2018.
- ⁵⁹³ Brian Cheung, Alexander Terekhov, Yubei Chen, Pulkit Agrawal, and Bruno Olshausen. Superposition of many models into one. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

594 595 596	Julien Colin, Thomas Fel, Rémi Cadène, and Thomas Serre. What i cannot predict, i do not understand: A human-centered evaluation framework for explainability methods. <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , pp. 2832–2845, 2022.
597 598 599 600	Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen- coders find highly interpretable features in language models. <i>arXiv preprint arXiv:2309.08600</i> , 2023.
601 602 603	Steven de Rooij and Paul Vitányi. Approximating rate-distortion graphs of individual data: Experiments in lossy compression and denoising. <i>arXiv preprint cs/0609121</i> , 2006.
604 605 606	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009.
607 608	R Becket Ebitz and Benjamin Y Hayden. The population doctrine in cognitive neuroscience. <i>Neuron</i> , 109(19):3055–3068, 2021.
610 611 612 613	Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCan- dlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. <i>Transformer Circuits Thread</i> , 2022.
614 615	Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. <i>University of Montreal</i> , 1341(3):1, 2009.
617 618 619	Thomas Fel, Remi Cadene, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, and Thomas Serre. Look at the variance! efficient black-box explanations with sobol-based sensitivity analysis. In Advances in Neural Information Processing Systems (NeurIPS), 2021.
620 621 622 623	Thomas Fel, Thibaut Boissin, Victor Boutin, Agustin Martin Picard, Paul Novello, Julien Colin, Drew Linsley, Tom ROUSSEAU, Remi Cadene, Lore Goetschalckx, et al. Unlocking feature visualization for deep network with magnitude constrained optimization. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023a.
624 625 626 627	Thomas Fel, Victor Boutin, Mazda Moayeri, Rémi Cadène, Louis Bethune, Mathieu Chalvidal, Thomas Serre, et al. A holistic approach to unifying automatic concept extraction and concept importance estimation. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2023b.
628 629 630	Thomas Fel, Agustin Picard, Louis Bethune, Thibaut Boissin, David Vigouroux, Julien Colin, Rémi Cadène, and Thomas Serre. Craft: Concept recursive activation factorization for explainability. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , 2023c.
631 632 633	Christiane Fellbaum. Wordnet. In <i>Theory and applications of ontology: computer applications</i> , pp. 231–243. Springer, 2010.
634 635 636	Ruth C. Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful per- turbation. In <i>Proceedings of the IEEE International Conference on Computer Vision (ICCV)</i> , 2017.
637 638 639	Alex Forsythe, Gerry Mulhern, and Martin Sawey. Confounds in pictorial sets: The role of complexity and familiarity in basic-level picture processing. <i>Behavior research methods</i> , 40(1):116–129, 2008.
640 641 642	Alexandra Forsythe. Visual complexity: Is that all there is? In Engineering Psychology and Cognitive Ergonomics: 8th International Conference, EPCE 2009, Held as Part of HCI International 2009, San Diego, CA, USA, July 19-24, 2009. Proceedings 8, pp. 158–166. Springer, 2009.
643 644 645	Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Texture and art with deep neural networks. <i>Current opinion in neurobiology</i> , 46:178–186, 2017.
646 647	Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. <i>arXiv preprint arXiv:1811.12231</i> , 2018.

648 649 650	Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. In <i>Advances in Neural Information Processing Systems</i> , pp. 9273–9282, 2019.
651 652 653	Gabriel Goh, Nick Cammarata †, Chelsea Voss †, Shan Carter, Michael Petrov, Ludwig Schubert, Alec Radford, and Chris Olah. Multimodal neurons in artificial neural networks. <i>Distill</i> , 2021. doi: 10.23915/distill.00030. https://distill.pub/2021/multimodal-neurons.
654 655 656	Mara Graziani, An-phi Nguyen, Laura O'Mahony, Henning Müller, and Vincent Andrearczyk. Concept discovery and dataset exploration with singular value decomposition. In <i>ICLR 2023</i> <i>Workshop on Pitfalls of limited data and computation for Trustworthy ML</i> , 2023.
657 658 659	James V Haxby, M Ida Gobbini, Maura L Furey, Alumit Ishai, Jennifer L Schouten, and Pietro Pietrini. Distributed and overlapping representations of faces and objects in ventral temporal cortex. <i>Science</i> , 293(5539):2425–2430, 2001.
661 662 663	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> (CVPR), 2016.
664 665 666	David H Hubel and Torsten N Wiesel. Receptive fields of single neurones in the cat's striate cortex. <i>The Journal of physiology</i> , 148(3):574, 1959.
667 668 669 670	Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. In- terpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In <i>International conference on machine learning</i> . Proceedings of the International Conference on Machine Learning (ICML), 2018.
671 672 673	Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, and Olga Russakovsky. HIVE: Evaluating the human interpretability of visual explanations. In <i>Proceedings of the IEEE European Conference on Computer Vision (ECCV)</i> , 2022.
674 675 676	Jerome Y Lettvin, Humberto R Maturana, Warren S McCulloch, and Walter H Pitts. What the frog's eye tells the frog's brain. <i>Proceedings of the IRE</i> , 47(11):1940–1951, 1959.
677 678	Ming Li, Xin Chen, Xin Li, Bin Ma, and Paul MB Vitányi. The similarity metric. <i>IEEE transactions</i> on <i>Information Theory</i> , 50(12):3250–3264, 2004.
679 680 681	Drew Linsley, Dan Shiebler, Sven Eberhardt, and Thomas Serre. Learning what and where to attend. 2019.
682 683	Sébastien Marcel and Yann Rodriguez. Torchvision the machine-vision package of torch. In <i>Proceedings of the 18th ACM international conference on Multimedia</i> , pp. 1485–1488, 2010.
684 685 686 687	Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. <i>Advances in neural information processing systems</i> , 29, 2016a.
688 689 690 691	Anh Nguyen, Jason Yosinski, and Jeff Clune. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. <i>Visualization for Deep Learning workshop, Proceedings of the International Conference on Machine Learning (ICML)</i> , 2016b.
692 693	Anh Nguyen, Jason Yosinski, and Jeff Clune. Understanding neural networks via feature visualization: A survey. <i>arXiv preprint arXiv:1904.08939</i> , 2019.
695 696 697	Paul Novello, Thomas Fel, and David Vigouroux. Making sense of dependence: Efficient black-box explanations using dependence measure. In <i>Advances in Neural Information Processing Systems</i> (<i>NeurIPS</i>), 2022.
698 699	Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. Distill, 2017.
700 701	Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. <i>Distill</i> , 2020. doi: 10.23915/distill.00024.001. https://distill.pub/2020/circuits/zoom-in.

702	Anitha Pasupathy and Charles E Connor. Shape representation in area v4: position-specific tuning
703 704	for boundary conformation. Journal of neurophysiology, 2001.
705 706	Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. In <i>Proceedings of the British Machine Vision Conference (BMVC)</i> , 2018.
707 708	R Quian Quiroga. Invariant visual representation by single neurons in the human brain. <i>Nature</i> , 435 (7045):1102–1107, 2005.
709	
710 711	cell'coding in the medial temporal lobe. <i>Trends in cognitive sciences</i> , 12(3):87–91, 2008.
712 713 714	Rodrigo Quian Quiroga, Itzhak Fried, and Christof Koch. Brain cells for grandmother. <i>Scientific American</i> , 308(2):30–35, 2013.
715 716	R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, 2021. URL https://www.R-project.org/.
717 718 719	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In <i>Knowledge Discovery and Data Mining (KDD)</i> , 2016.
720 721 722	Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. <i>International journal of computer vision</i> , 115:211–252, 2015.
723 724 725 726	Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based lo- calization. In <i>Proceedings of the IEEE International Conference on Computer Vision (ICCV)</i> , 2017.
728 729 730	Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In <i>Proceedings of the International Conference on Machine Learning (ICML)</i> , 2017.
731 732 733	Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. In <i>Workshop on Visualization for Deep Learning, Proceedings of the International Conference on Machine Learning (ICML)</i> , 2017.
734 735 736	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In <i>Proceedings of the International Conference on Machine Learning (ICML)</i> , 2017.
737 738 739	Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Du- mitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1–9, 2015.
740 741 742 743	Mohammad Reza Taesiri, Giang Nguyen, and Anh Nguyen. Visual correspondence-based expla- nations improve ai robustness and human-ai team accuracy. <i>Advances in Neural Information</i> <i>Processing Systems (NeurIPS)</i> , 2022.
744 745	Alex Tamkin, Mohammad Taufeeque, and Noah D Goodman. Codebook features: Sparse and discrete interpretability for neural networks. <i>arXiv preprint arXiv:2310.17230</i> , 2023.
746 747 748 749 750 751 752	Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. <i>Trans-</i> <i>former Circuits Thread</i> , 2024. URL https://transformer-circuits.pub/2024/ scaling-monosemanticity/index.html.
753 754 755	Kazushige Tsunoda, Yukako Yamane, Makoto Nishizaki, and Manabu Tanifuji. Complex objects are represented in macaque inferotemporal cortex by the combination of feature columns. <i>Nature neuroscience</i> , 4(8):832–838, 2001.

756 757 758	Johanna Vielhaben, Stefan Blücher, and Nils Strodthoff. Multi-dimensional concept discovery (mcd): A unifying framework with completeness guarantees. 2023.
759 760	Gang Wang, Keiji Tanaka, and Manabu Tanifuji. Optical imaging of functional organization in the monkey inferotemporal cortex. <i>Science</i> , 272(5268):1665–1668, 1996.
761 762 763	Zhou Wang and Alan C Bovik. Mean squared error: Love it or leave it? a new look at signal fidelity measures. <i>IEEE signal processing magazine</i> , 26(1):98–117, 2009.
764 765 766	Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. <i>IEEE transactions on image processing</i> , 13(4):600–612, 2004.
767 768	Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. <i>arXiv preprint arXiv:1506.06579</i> , 2015.
769 770 771 772	M. D. Zeiler, G. W. Taylor, and R. Fergus. Adaptive deconvolutional networks for mid and high level feature learning. In <i>Proceedings of the IEEE International Conference on Computer Vision (ICCV)</i> , 2011.
773 774	Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In <i>Proceedings of the IEEE European Conference on Computer Vision (ECCV)</i> , 2014a.
775 776 777	Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In <i>Proceedings of the IEEE European Conference on Computer Vision (ECCV)</i> , 2014b.
778 779 780 781	Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A Ehinger, and Benjamin IP Rubinstein. Invertible concept-based explanations for cnn models with non-negative concept activation vectors. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 11682–11690, 2021.
782 783 784 785	Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Object detectors emerge in deep scene cnns. In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> , 2016.
786 787	Roland S Zimmermann, David A Klindt, and Wieland Brendel. Measuring mechanistic interpretability at scale without humans. In <i>ICLR 2024 Workshop on Representational Alignment</i> .
788 789 790 791	Roland S Zimmermann, Judy Borowski, Robert Geirhos, Matthias Bethge, Thomas Wallis, and Wieland Brendel. How well do feature visualizations support causal understanding of cnn activations? <i>Advances in Neural Information Processing Systems</i> , 34:11730–11744, 2021.
792 793 794	Roland S Zimmermann, Thomas Klein, and Wieland Brendel. Scale alone does not improve mechanistic interpretability in vision models. <i>Advances in Neural Information Processing Systems</i> , 36, 2023.
795	
796	
797	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
003	

APPENDIX

812 813

PRACTICE TRIALS А

814 All participants were shown 9 practice trials to help them understand the core idea of the task. The 815 images for these trials did not overlap with the images used in any of the experiments, preventing 816 them from influencing our results. They were selected to represent specific features on the maximum-817 activating grid: checkerboard, veiny, green, round, blue, rough fur, yellow, straight lines, and magenta. 818 Images on the left side of the grid in practice trials had no coherent pattern, and were randomly 819 sampled without replacement from a set of images with features including whiskers, spikes, droplets, 820 and liquid flow.

821 Practice trials were customized for each experiment, to keep all variables controlled. For the 822 feature visualization psychophysics, practice trials were generated using the same methods as the 823 experimental images. They were chosen from neurons in convolutional layers to show the same 824 features that were included in the practice trials for other experiments.

825 826

827

В ATTENTIVENESS TESTS

828 Attentiveness tests were included at random points in the study. Following the methodology adopted 829 in previous work (Borowski et al., 2021), they were based on a simple premise: if the query image 830 is also present in the grid of maximum-activating images, participants paying attention should 831 get near-perfect accuracy. We generated attentiveness tests for each experiment using the same 832 procedure as was used for the experimental trials. To ensure that the attentiveness tests did not overlap 833 with experimental trials, we chose maximum- and minimum-activating images from neurons from 834 convolutional layers. Rather than sampling 10 maximum-activating images, we sampled 9 and put the 835 query in a random position in the maximum-activating grid. The attentiveness tests were generated for all experiments using the same neurons from the model to guarantee consistency. 836

- 837
- 838 839

841

С IMAGE SECTION FOR THE DISTRIBUTED REPRESENTATION CONDITION.

840 We originally tested 2 ways to select the direction from the distributed representation obtained by CRAFT to be kept for the main experiment. The first one is the one explained in the main section 842 of the paper, namely: from the 10 directions, we selected the direction that was the most often the 843 most activating one across the 300 images. However, we considered a second alternative where we 844 selected the direction that was the most often the most activating one across the 2900 images. We ran the psychophysics experiments on both alternatives. The average performance obtained on this 845 alternative was of $81\% \pm 1.4$ and we find no effect of the condition variable between this distributed 846 condition vs the local condition describe in the main paper: $\beta_{condition} = -0.34, SE = 0.25, z =$ 847 -1.39, p = 0.164. Further quantitative and qualitative analysis done on both alternatives have led us 848 to keep the first alternative for the reminder of the experiments.

- 849 850
- 851
- 852
- 853
- 854 855
- 856
- 857
- 858
- 859
- 861 862
- 863

864 D SUPPLEMENTARY INFORMATION ABOUT EXPERIMENT III 865

Experimental protocol Given that the adopted experimental protocol was originally designed to evaluate the informativeness of feature visualizations, we devised Experiment III to shed light on 868 the value that a state-of-the-art feature visualization method (MACO) would add to the task. We subscribe to the idea that feature visualizations are more suitable to complement natural images 870 than to replace them. Thus, Experiment III implements a protocol similar to the *mixed* condition reported by Zimmermann et al. (Zimmermann et al., 2021). In practice, Experiment III was based on Experiment II, but in each trial, 4 of the natural images displayed both on the left and right panels 872 were replaced by 4 feature visualizations, as depicted in Figure 6.

873 874

871

866

867

Method To synthesize feature visualizations for Experiment III (see Section 3.2), we used 875 MACO (Fel et al., 2023a), inspired by the method by Olah et al. (Olah et al., 2017). MACO 876 generates feature visualizations with a natural Fourier amplitude spectrum by fixing the amplitude 877 spectrum to the empirical mean derived from natural images. Specifically, MACO optimizes directly 878 in the Fourier domain $\mathcal{F}(x)$ by modifying the phase φ of the target image while keeping the magni-879 tude r of the spectrum fixed. This constraint on the magnitude helps prevent high-frequency artifacts 880 and ensures that the resulting images remain visually coherent. Formally, let r denote the average amplitude spectrum computed across the ImageNet dataset, and let d be a target concept direction to be maximized. MACO solves the following optimization problem:

883

890

884 885

 $oldsymbol{arphi}^{\star} = rgmax_{oldsymbol{arphi}}\left(oldsymbol{f}\left(\mathcal{F}^{-1}\left(oldsymbol{r}\circ e^{ioldsymbol{arphi}}
ight)
ight)\cdotoldsymbol{d}
ight), \ \ ext{and} \ \ oldsymbol{x}^{\star} = \mathcal{F}^{-1}\left(oldsymbol{r}\circ e^{ioldsymbol{arphi}^{\star}}
ight).$

where x^* is the feature visualization obtained after optimization, \mathcal{F} denotes 2-D Discrete Fourier Transform (DFT) on $\mathcal{X}, \mathcal{F}^{-1}$ its inverse and \circ represents element-wise multiplication. Additionally, 887 MACO introduces an attribution-based transparency mask to highlight spatially important regions in the visualizations, enhancing interpretability. 889

The value of feature visualization remains unclear. Finally, Figure 7 depicts the results obtained 891 from Experiment III, which combined natural images with synthetic images corresponding to feature 892 visualizations using MACO. Similarly to (Zimmermann et al., 2021), we do not find that mixing 893 stimuli from feature visualization with natural images helped participants perform better at the task. 894 In fact, participants performed overall significantly worse in Experiment III than in Experiment II: 895 $\beta_{depth:Exp3} = -0.63, SE = 0.19, z = -3.3, p < 0.001$. Interestingly, the worst performance was 896 observed in the local condition for units located in the deepest layers of the network, with a drop in 897 performance from 80.2% in Experiment I to 68.3% in Experiment III. Furthermore, the difference in the per-layer task performance obtained in layer 4 in each condition is 8.2%, 17.1%, and 18.9%, respectively. We leave to future work the investigation into the specific factors that contributed to 899 this decline in performance. Potential avenues for future research include exploring the cognitive 900 load imposed by mixed stimuli and identifying optimal conditions under which feature visualization 901 might enhance rather than hinder task performance. 902

903

904

905

906

907 908

909

910

911

- 912
- 913
- 914 915

916



Figure 6: **Illustration of a trial from Experiment III**. Example of a trial from *layer*4.2.*bn*3 in the distributed representation condition from Experiment III. Note how the reference images on the left and right panels contain a mix of feature visualizations using MACO and natural images.



Figure 7: **Per-layer results for Experiment III**. Given a feature and a set of images to illustrate it, we assess how visually coherent participants find this set of images—or how unambiguous the feature is. More precisely, we measure the proportion of participants that are able to identify the query image which is also part of this set of images. Similarly to the results from experiment I and II, a clear trend emerges where features appear significantly less ambiguous in the distributed representation than in the local representation condition, particularly in the deeper layers of the network.

Е FURTHER ILLUSTRATION OF LOCAL VS DISTRIBUTED TRIALS



Figure 8: Layer2. This figure illustrates a trial used to assess the features encoded in layer2.0 either by the neuron 52 (a) or at least partially through the neuron 52 (b).



(a) Local representations

(b) Distributed representations

Figure 9: Layer3. This figure illustrates a trial used to assess the features encoded in layer3.1 either by the neuron 957 (a) or at least partially through the neuron 957 (b).



(a) Local representations



(b) Distributed representations

Figure 10: Layer4. This figure illustrates a trial used to assess the features encoded in layer4.2 either by the neuron 259 (a) or at least partially through the neuron 259 (b).





Figure 12: **Per-layer results for Experiment IV (VGG16)**. Given a feature and a set of images to illustrate it, we assess how visually coherent participants find this set of images—or how unambiguous the feature is. More precisely, we measure the proportion of participants that are able to identify the query image which is also part of this set of images. Similarly to the results from experiments I, II, and III, a trend emerges where features appear less ambiguous in the distributed representation than in the local representation condition. Nevertheless, the effect of depth appears more nuanced than in the ResNet50.

1080

1082

1084

1087

1088 1089

1090 1091

1101 1102

1104

1103 F LOCAL VS DISTRIBUTED: RESULTS FROM A VGG16

We extend the experiments to a VGG16 from the Torchvision (Marcel & Rodriguez, 2010) library, pre-trained on ImageNet-1k (Deng et al., 2009). We randomly select 80 neurons across the VGG16, and we follow the same methodology as previously for the selection of stimuli for the local and distributed conditions. For the psychophysics experiment, we follow the experimental protocol of *Experiment II* where we control for semantic confounds. We analyze the data corresponding to 132 participants from Prolific.

1111

1118

Human performance is superior in the distributed representation condition. Based on our main hypothesis that a distributed representation constitutes a better basis for interpretability than a local representation, we predicted that participants would be better at selecting the maximally activating query image in the distributed condition. In Experiment IV, the average performance across participants in the distributed condition was $81.5\% \pm 1.5$, when compared to $77.1\% \pm 1.6$ in the local condition. A GLMER with condition as a predictor revealed a statistically significant disadvantage for the local condition: $\beta_{condition} = -0.45$, SE = 0.22, z = -2.04, p = 0.04.

The benefit of the distributed representation does not increase systematically with depth. Figures 12 illustrates the per-layer results obtained in Experiments IV (VGG16). While the superiority of distributed representation appears again here in the deeper layer, we do not see a systematic trend between depth and the performance of participants in our psychophysics experiments. This observation is reflected in our analysis as we do not find that the advantage of the distributed condition increases as the units belong to deeper layers: $\beta_{depth:Exp4} = -0.06$, SE = 0.05, z = -1.27, p = 0.20.

Takeaway. In general, we expect deeper layers to contain more information, and we expect a model to rely more on superposition to encode information in those layers. Based on those assumptions, we expected participants to systematically benefit more from distributed representation when exposed to features from deeper layers. Yet, while our new results strengthen our claim that there appears to be no downside to studying features using distributed representation, they portray a more complex picture as to when they become necessary.

1132 1133

G DO FEATURES FROM *local* VS. *distributed* REPRESENTATIONS DIFFER IN OBVIOUS WAYS?

1136 1137 1138

1139

1152

1178 1179 Following the useful suggestions of a reviewer we use existing metrics to probe if features from local vs. distributed representations differ in obvious ways.

1140 G.1 Complexity

1142 **Kolmogorov complexity.** The first hypothesis we test is whether features from distributed repre-1143 sentation are easier to interpret than features from local conditions because they are less complex. 1144 While we do not have access to the feature per se, we have access to the images that possess the 1145 features as a proxy. While measuring the complexity of images is still an open research question, 1146 there are works that show that the Kolmogorov complexity of images (Li et al., 2004) correlate well with human-derived ratings for the complexity of natural images (Forsythe et al., 2008; Forsythe, 1147 2009). In practice, we follow previous works (Li et al., 2004; de Rooij & Vitányi, 2006) and use a 1148 standard compression technique (JPEG) to approximate the Kolmogorov complexity of images. The 1149 hypothesis is that the more an image can be compressed, the less complex the features that compose 1150 the images. 1151

MethodologyFor a given target feature, we conducted T trials $T \in [1, 10]$, based on how many
trials can be semantically controlled). In each trial, we had access to 10 images that possess the target
feature. We measured the average Kolmogorov complexity of these 10 images by recording their
compressed file sizes after JPEG compression. We then averaged these complexities across all trials
for each feature, resulting in a measure of Kolmogorov complexity for that feature. This methodology
is applied to every feature.

Results. We did not find a significant correlation between the Kolmogorov complexity of the images and the visual coherency scores from Experiments II (ResNet50) (Fig. 13 and IV (VGG16) (Fig 14). This suggests that while complexity might be a factor, the superiority of features from distributed representations cannot be explained by complexity alone.



Figure 13: Kolmogorov complexity of features from a ResNet50 (Exp II).

1180 G.2 SIMILARITY

Structural similarity Index The second hypothesis we test is if the images used to illustrated distributed features are more similar to the ones used to illustrate local features. The intuition behind this hypothesis is that the more ambiguous the features, the more dissimilar the images used to illustrate it, the worse people will perform in our experiment. To test this hypothesis we used an existing similarity measure, namely the structural similarity index measure (Wang et al., 2004; Wang & Bovik, 2009). This measure quantifies the similarity of 2 images based on luminance, constrast and change in structural information.













