## How to sample the world for understanding the visual system

Johannes Roth (jroth@cbs.mpg.de), Martin N. Hebart (hebart@cbs.mpg.de)

Max Planck Institute for Human Cognitive and Brain Sciences, 04103 Leipzig, Germany Department of Medicine, Justus Liebig University Giessen, 35390 Giessen, Germany Center for Mind, Brain and Behavior (CMBB), Universities of Marburg, Giessen, and Darmstadt

48

78

79

80

#### Abstract

Understanding vision requires capturing the vast diver- 49 2 sity of the visual world we experience. How can we sam-50 3 ple this diversity in a manner that supports robust, gen-<sup>51</sup> 4 eralizable inferences? While widely-used, massive neu-52 5 roimaging datasets have strongly contributed to our un-53 6 derstanding of brain function, their ability to comprehen-54 7 sively capture the diversity of visual and semantic experi-55 8 ences has remained largely untested. More broadly, the 56 9 factors required for diverse and generalizable datasets 57 10 have remained unknown. To address these gaps, we in-58 11 troduce LAION-natural, a curated subset of 120 million 59 12 natural photographs filtered from LAION-2B, and use it 60 13 as a proxy of the breadth of our visual experience in as-61 14 sessing visual-semantic coverage. Our analysis of CLIP 62 15 embeddings of these images suggests significant repre-63 16 sentational gaps in existing datasets, demonstrating that <sup>64</sup> 17 they cover only a restricted subset of the space spanned 65 18 by LAION-natural. Simulations and analyses of functional 66 19 MRI data further demonstrate that these gaps are associ-67 20 ated with impaired out-of-distribution generalization. Im-68 21 portantly, our results reveal that even moderately sized 69 22 stimulus sets can achieve strong generalization if they 23 are sampled from a diverse stimulus pool, and that this 70 24 diversity is more important than the specific sampling 71 25 strategy employed. These findings not only highlight lim-72 26 itations of existing datasets in generalizability and model  $_{_{73}}$ 27 comparison, but also provide guidance for future stud-74 28 ies to support the development of stronger computational 75 29 models of the visual system and generalizable inferences.  $_{76}$ 30 77

Keywords: stimulus selection; naturalistic fMRI; adaptive sampling; dataset; generalization

33

### Introduction

Humans encounter an incredibly diverse range of visual stim- 81 34 uli, and capturing this breadth is essential for understanding 82 35 how the brain represents visual information. This has moti-83 36 vated the development of ever larger datasets of brain activity 84 37 in response to naturalistic visual stimuli (Allen et al. (2022); 85 38 Chang et al. (2019); Hebart et al. (2023)). While broad sam-86 39 pling suffers from diminishing returns (Allen et al. (2022)), 87 40 large-scale, condition-rich datasets have laid the foundation 88 41 for extensive computational modeling of the visual system, al- 89 42 lowing for a detailed and fine-grained analysis of its function 90 43 (Contier et al. (2024); Takagi & Nishimoto (2023)). 44 While these datasets have been invaluable in advancing 92 45

our understanding, particularly within specific domains like ob- 93
 ject and scene processing, their ability to capture the visual- 94

semantic diversity of our world remains largely untested. This has an important consequence: If a dataset lacks diversity, then inferences may not generalize. This is particularly relevant in cognitive computational neuroscience, where recent insights from model comparison, encoding, and decoding often rely on such datasets (Doerig et al. (2023)). While previous research has highlighted a lack of semantic diversity in some datasets (Shirakawa et al. (2024)), much less is known about how the visual-semantic coverage of real-world concepts affects generalization performance.

To address these gaps, we used a three-pronged approach. First, we assessed how well existing stimulus sets cover the breadth of visual experience by embedding them within a much larger representational image space defined by LAIONnatural, a newly curated subset of 120 million naturalistic photographs (images depicting real-world scenes or objects, excluding heavily edited or synthetic content) filtered from LAION-2B (Schuhmann et al. (2022)). Second, we evaluated whether a sufficiently diverse stimulus set can enable generalization at a scale practical for vision research. Finally, in simulations and analyses of functional MRI data, we determined the effect of sampling strategy on generalization performance.

#### **Results**

# Significant gaps in visual-semantic coverage in existing stimulus sets

Determining how well existing stimulus sets represent our natural visual experience is not possible due to the lack of a detailed understanding of its contents and the relative frequency of different concepts or "classes" of experience. To approximate broader coverage, we used the large and highly diverse LAION-2B dataset (~2 billion image-text pairs). To focus on high-quality natural photographs, which more closely reflect our visual experience, we (1) filtered LAION-2B to select only unique, high-resolution images and excluded NSFW content and (2) applied a classifier trained on 25,000 hand-labeled images to restrict images to natural photographs (Fig. 1A, see Fig. S1 for an overview of natural/non-natural images). We term this new image set "LAION-natural" (~120 million photographs). To verify the diversity of LAION-natural, we tested if all concepts found in common stimulus sets, including ImageNet (Deng et al. (2009)), MS COCO (Lin et al. (2014)), and ecoset (Mehrer et al. (2021)), could be found in the dataset. Even though this approximation of natural vision is likely an incomplete characterization and will contain specific biases inherent in the dataset, we can treat coverage and generalization abilities on LAION-natural as a proxy for coverage of the visual world to understand limitations in existing datasets



Figure 1: A) Filtering procedure to generate LAION-natural from LAION-2B. B) t-SNE projection of 5,500 cluster centroids from LAION-natural, with colors reflecting presence of existing datasets (cutoff: 2 or more images). Inspection of these clusters revealed various concepts not covered by THINGS or NSD, including landscapes, natural events, crowds, or non-Western public figures (see Fig. S3). C) Percentage of clusters covered by THINGS, NSD, or both (cutoff: 2 or more images). D) Percentage of LAION-natural that is in-distribution, based on Principal Component Analysis (PCA)-based outlier detection. While neither THINGS nor NSD was able to explain the variance in LAION-natural, a random subset of only 6,000 LAION-natural samples still captured 93.51%.

113

<sup>95</sup> and reveal strategies for broader stimulus sampling.

Having curated a large pool of natural photographs, we<sup>114</sup> 96 next evaluated how much of LAION-natural is covered by115 97 the Natural Scenes Dataset (NSD; Allen et al. (2022)) and<sup>116</sup> 98 THINGS (Hebart et al. (2023)), two of the largest, densely<sup>117</sup> 99 sampled visual neuroimaging datasets. We approximated<sup>118</sup> 100 visual-semantic coverage using CLIP-extracted image fea-119 101 tures, known for their large-scale training datasets and their<sup>120</sup> 102 alignment with both human perceptual ratings (Demircan et121 103 al. (2023); Kaniuth et al. (2024); Muttenthaler et al. (2022))122 104 and cortical activity patterns (Conwell et al. (2024); Wang et123 105 al. (2023)). To evenly distribute the dataset into similarly sized<sup>124</sup> 106 chunks, we divided LAION-natural into 5,500 clusters using125 107 mini-batch k-Means on the image features. 126 108

A 2D t-SNE projection of the cluster centroids revealed low<sup>127</sup>
 overlap between THINGS and NSD (Fig. 1B), likely due to<sup>128</sup>
 their distinct focuses on scenes and objects, respectively. The<sup>129</sup>
 slightly higher coverage of NSD can in part be explained by its<sup>130</sup>

larger dataset size (70,000 vs. 26,107 images). More importantly, the visualization indicates that both datasets exhibited significant gaps in the visual-semantic image space defined by LAION-natural. We quantified this observation by determining the percentage of LAION-natural clusters represented by at least 2 images from either THINGS or NSD, which ensures an inclusive threshold while minimizing the effect of outliers (Fig. 1C). This analysis showed that 49.83% of LAIONnatural clusters were not covered by either dataset. To test whether this finding arises through the direction of comparison, we also clustered both NSD and THINGS (200 clusters, average cluster size 0.5% of total dataset) and assessed how many of them were covered by 10 million random images from LAION-natural. We found that 97.5% of NSD and 100% of THINGS clusters were assigned at least one LAION-natural image, with the few uncovered NSD clusters containing only single images. We additionally evaluated how many of the images in LAION-natural would be considered "in-distribution"



Figure 2: A) OOD accuracy depending on the number of training samples (distributed across 100 clusters). OOD accuracy saturates at 10,000 samples. B) Reducing the number of clusters, while keeping the number of samples constant (6k), reduced OOD accuracy (right side). C) Real-world test for OOD accuracy. LAION-natural is divided into six cluster groups. One of the cluster groups is considered OOD and the remaining groups are used to construct the training set. D) Accuracy for the OOD group from LAION-natural, based on the number of clusters in the training set. Less diverse training sets result in lower OOD accuracy. E & F) OOD accuracy of THINGS and NSD, compared against random and stratified samples from LAION-natural (averaged across groups, controlled for dataset size). Neither of the datasets sufficiently spans the LAION-natural image space to enable OOD generalization.

158

159

160

for THINGS and NSD, respectively, using a PCA-based recon-150
 struction error approach. This revealed that only 62.99% and 151
 60.68% of LAION-natural were in-distribution for THINGS and 152
 NSD, respectively (Fig. 1D).

In contrast, a small random subset of LAION-natural (6,000<sup>154</sup>
 samples) already achieved an in-distribution score of 93.51%.<sup>155</sup>
 Together, these findings demonstrate notable limitations in the<sup>156</sup>
 semantic diversity and coverage of existing stimulus sets.

# When diversified, even moderate-sized stimulus sets can generalize well

Given the limited visual-semantic diversity in existing large-161 141 scale fMRI datasets, we next asked how this affects gener-162 142 alizability of inferences drawn using these data and to what163 143 degree it is possible to collect diverse data at realistic scales164 144 for future studies. Prior work has shown that a lack of diver-165 145 sity can hinder our ability to draw generalizable conclusions166 146 from one of these fMRI datasets (Shirakawa et al. (2024)).167 147 However, it is unclear whether these challenges can even be168 148 mitigated with realistic dataset sizes. While simulations sug-169 149

gest that covering key representational axes can support outof-distribution (OOD) generalization for brain-to-image reconstruction (Shirakawa et al. (2024)), these findings were based on 500,000 samples, an impractical scale for most fMRI studies. Thus, the critical questions of dataset size and dataset diversity on drawing generalizable inferences from fMRI data have remained largely unanswered.

To address these questions, we simulated synthetic fMRI responses to determine how much data is required for OOD generalization. If a lot of data is required, this indicates that it is not feasible in practice to achieve OOD generalization, while if generalization is possible with fewer samples, this highlights the potential of diverse sampling. We used a teacher-student learning model, inspired by previous simulations (Shirakawa et al. (2024)), and generated stimulus features from 100 clusters via a Gaussian mixture model (GMM). Next, we mapped them to synthetic brain responses using a fixed teacher weight plus Gaussian noise. Finally, we employed Ridge regression for decoding features from synthetic brain data. We evaluated predictions on novel OOD clusters using cluster accuracy, i.e.,



Figure 3: A) From a pool of 500,000 synthetic samples (GMM, C=100), up to 10,000 samples were selected using different sampling strategies (see Methods for details) and evaluated on OOD performance. B) Samples drawn from LAION-natural and evaluated on a held-out cluster group (averaged across groups). C) Effect of sampling strategy on encoding performance in subsets of NSD. Performance is 5-fold cross-validated (20% of NSD as test set), normalized within subjects and averaged across subjects.

how often predicted features correlated most with the source<sup>203</sup>
 cluster centroid. 204

Adjusting the total number of generated samples showed<sup>205</sup> 172 that accuracy started to saturate after 5,000-10,000 samples<sup>206</sup> 173 (distributed across all 100 clusters, Fig. 2A). To determine the207 174 effect of dataset diversity for a dataset with realistic size, we208 175 fixed the number of samples to 6,000 and varied the number209 176 of clusters from which we sampled. This analysis revealed<sup>210</sup> 177 that reducing training set diversity also reduced generalizabil-211 178 ity (Fig. 2B). Together, these findings show that, in principle,212 179 it is possible to identify medium-sized, diverse stimulus sets213 180 that can generalize well in neuroimaging studies, as long as<sup>214</sup> 181 the underlying stimulus pool is diverse. 215 182

However, the simulations assumed evenly distributed clus-216 183 ters, which does not reflect distributions in real-world datasets.217 184 To incorporate dataset realism, we modified our approach by<sup>218</sup> 185 replacing GMM-generated samples with CLIP features from<sup>219</sup> 186 LAION-natural to generate synthetic brain responses. To al.220 187 low for broad yet homogeneous sampling across the entire221 188 dataset, we clustered images into 1,000 clusters. To simulate<sup>222</sup> 189 the effect of uneven distribution and OOD generalization, we223 190 formed 6 OOD groups from these clusters by applying a sec-191 ond layer of clustering. We then used held-out clusters from<sup>224</sup> 192 LAION-natural to evaluate OOD accuracy (Fig. 2C), training<sup>225</sup> 193 the regression model on 6,000 subsamples of the remaining<sub>226</sub> 194 5 groups. Importantly, to test for uneven sampling, we varied 227 195 the number of groups sampled from between 1 and 5. The re-228 196 sults of this simulation are shown in Fig. 2D, revealing a higher 197 accuracy given the noise in the data, but only when incorporat-230 198 ing broad sampling across most groups. The results confirment 199 that, to achieve the highest generalization performance, highese 200 dataset diversity is required. 201 233

202 While our previous simulations examined the effects of non-234

diverse datasets with synthetic subgroups, they did not reflect real, empirical datasets. Thus, we extended this analysis by using THINGS and NSD as training datasets, excluding samples assigned to the OOD group that served as a test set. As a baseline, we sampled from LAION-natural using both random and stratified approaches across clusters. The results revealed that THINGS and NSD strongly underperformed in OOD accuracy compared to stratified sampling across all training clusters in LAION-natural, and also, perhaps surprisingly, when sampled randomly (Fig. 2E/F). Additionally, we assessed how dataset diversity affects OOD performance in real fMRI data, by repeating the clustering analysis on NSD, which mirrored the results from LAION-natural (Fig. S2).

Together, these findings highlight that gaps in visualsemantic coverage reduce OOD accuracy, both in simulations and in common existing datasets. Importantly, assuming results from CLIP embeddings generalize to real fMRI data, our results suggest that constructing a diverse, generalizable stimulus set is feasible within the size constraints of neuroimaging studies.

# Sampling strategy matters less than stimulus pool diversity

Our previous simulation showed that random sampling performed comparably to stratified sampling. This result was unexpected, as targeted sampling could be seen as having a lot of potential for improvements in sampling efficiency and generalizability. Therefore, we addressed the degree to which sampling strategies help to improve dataset efficiency in the presence of a diverse stimulus pool.

To assess the role of sampling strategy in dataset efficiency, we evaluated various established sampling proce-



Figure 4: A) Concept distribution of LAION-natural subset (red) differs from natural language frequency (black). B) Depending on sampling strategy, concept distribution of subset is kept (random, stratified, k-Means, margin) or shifts to natural language (Greedy ED, Core-Set). C) Dimensionality of sampled dataset is highest for Greedy-ED, followed by Core-Set and others.

298

299

dures, including random and stratified sampling across clus-269 235 ters, k-Means clustering-based sampling, Core-Set sampling270 236 (k-Center-Greedy), greedy sampling to maximize effective di-271 237 mensionality (ED), and margin-based uncertainty sampling 272 238 Similar to our initial simulation, we generated 500,000 sam-273 239 ples with a GMM to evaluate the effect of sampling strategy<sub>274</sub> 240 on OOD accuracy across dataset sizes. Our results revealed 275 241 only minor differences between strategies, with strong effects<sub>276</sub> 242 of dataset size, yet small effects of sampling strategy. Only<sub>277</sub> 243 Core-Set and k-Means-based sampling approaches showed 278 244 slight advantages at 2,000 and 3,000 data points (Fig. 3A). 279 245

To verify these findings with a real image pool, we repeated<sup>280</sup> 246 the sampling experiment using LAION-natural. Unlike previ-281 247 ous simulations, where entire groups were left out, we left out282 248 sets of clusters, thus approximating sampling from a broad283 249 dataset. We found that most sampling strategies performed<sub>284</sub> 250 similarly, except for Core-Set sampling, which consistently285 251 outperformed others (Fig. 3B). However, dataset diversity re-286 252 mained the most critical factor, outweighing the choice of sam-287 253 pling strategy. 288 254

To test the degree to which these results would generalize<sup>289</sup> 255 to empirical fMRI data, we applied these sampling strategies 256 to NSD and evaluated generalization performance. Rather290 257 than focusing on decoding, as in the previous analyses, we<sub>pat</sub> 258 focused on encoding, a commonly used approach for NSD,292 259 We trained a Ridge regression model to predict single-trial re-293 260 sponse estimates in the ventral stream using CLIP features<sub>294</sub> 261 (Fig. 3C). Consistent with the previous results, dataset di-295 262 versity remained the key determinant of performance, while 296 263 sampling strategy had a minimal impact. 264 297

#### 265 Sampling strategy can shift concept distribution

The previous results reveal the effect of sampling strategy on generalizability within the stimulus pool, in our case LAION-301 natural. Being derived from the internet, LAION-natural may 302 have biases in concept frequency distributions. While having only minor effects on generalizability, different sampling strategies could affect if and how these biases translate to sampled datasets. To examine this, we used an LLM to generate word labels for a random subset of LAION-natural (100,000 images) and compared the word frequency distribution to that found in natural language use.

We found that LAION-natural overall does not align well with natural language (see Fig. 4A). By calculating the Kullback-Leibler (KL) divergence between the concept distribution of the sampled dataset and LAION-natural or natural language, we found that random, stratified, k-Means-, and marginbased sampling closely mirrored the stimulus pool distribution, whereas Core-Set and ED-based methods were closer to natural language frequencies (see Fig. 4B). Core-Set and ED-based methods also resulted in the highest increase in the effective dimensionality of resulting image sets, measured via their CLIP embeddings (see Fig. 4C). These results underscore the effectiveness of Core-Set sampling for generating new datasets. However, the added benefits in generalization remain minor relative to the role of stimulus diversity.

#### Discussion

Understanding visual representations requires us to capture much of the visual-semantic richness of the visual world. A prominent research paradigm involves collecting extensive data in response to a broad set of stimuli, with the aim of allowing for generalizable inferences (Naselaris et al. (2021)). Our findings demonstrate that, at the level of visual semantics, commonly used stimulus sets fall short of this goal. Their emphasis on a constrained subset of concepts limits the generalizability of insights that can be drawn from them. However, this does not diminish the value of these datasets, which have significantly advanced our understanding of brain function, particularly within the domains of scenes or objects. And while our study focused on visual semantics, many studies us-<sup>357</sup>
 ing these datasets were not carried out at that level, and it is
 possible that existing datasets already have sufficient diversity
 to comprehensively capture purely low-level and mid-level pro cessing. Future studies should explore the degree to which
 generalizable inferences can be drawn in the visual domain
 alone.

However, when inferences are drawn about the entirety of<sub>364</sub> 310 the visual diet or when the aim is to use existing datasets365 311 to build generalizable models of the visual system, including 366 312 deep neural networks (DNNs), our results highlight clear con-367 313 straints in existing datasets and caution against overinterpre-368 314 tation. This insight is particularly relevant given recent findings369 315 suggesting that DNN models of the visual system yield similar<sub>370</sub> 316 performance regardless of architectural differences or training<sub>371</sub> 317 objective (Conwell et al. (2024)), an insight that would strongly372 318 affect the neuroconnectionist research program that requires373 319 a system identification approach for identifying "better" mod-374 320 els of the visual system (Doerig et al. (2023)). In contrast, our375 321 findings highlight that without sufficiently broad datasets, we<sub>376</sub> 322 cannot determine whether models truly converge or if their ap-377 323 parent similarity results from being evaluated on insufficiently<sub>378</sub> 324 diverse datasets. 325

How, then, can we design datasets that are generalizable?380 326 Our simulations indicate that, given a certain "stimulus bud-"81 327 get", visual-semantic breadth of sampling should be priori-  $^{\rm 382}$ 328 tized over depth to ensure maximum possible OOD perfor-383 329 mance. Furthermore, all tested sampling strategies provided<sup>384</sup> 330 sufficient coverage of the stimulus pool, yielding  $\mathsf{comparable}^{^{385}}$ 331 generalization performance. Notably, sampling strategy alone<sup>386</sup> 332 did not compensate for insufficient dataset diversity, empha-387 333 sizing that future studies should prioritize broad stimulus pools<sup>388</sup> 334 389 even when using random sampling. 335 390

While we quantified diversity using CLIP embeddings, even<sub>391</sub> 336 the most diverse existing stimulus pools may omit crucial as 392 337 pects of the visual representational space. Beyond CLIP, the<sub>393</sub> 338 field is in need of a more precise, guantitative definition of di-339 versity to support broader, stratified sampling (Conwell et al.394 340 (2024)). It is also worth noting that many of our findings are  $^{\rm 395}$ 341 based on simulations, where assumptions can affect results.<sup>396</sup> 342 While our results are consistent across simulations and vali-397 343 dated with real fMRI data, future research should further em-398 344 399 pirically evaluate dataset diversity. 345 400

Overall, this study underscores the necessity of broader<sub>401</sub> 346 coverage in the representational space for making generaliz-402 347 able inferences than provided by existing datasets. By demon-403 348 strating that relatively small yet diverse stimulus sets provide 349 large benefits for out-of-distribution generalization, we provide404 350 a framework for designing stimulus sets that enable large-405 351 scale, condition-rich studies of the visual-semantic system 406 352 Prioritizing diversity and coverage will allow researchers to407 353 construct datasets that better reflect the complexity of naturakos 354 vision, leading to more robust models of how humans perceive409 355 the world. 356 110

#### Methods

#### **Constructing and evaluating LAION-natural**

Filtering LAION-2B to naturalistic images LAION-2B (Schuhmann et al. (2022)) contains  $\sim$ 2 billion images with English captions, but visual inspection revealed a large fraction of unsuitable, non-natural images. We filtered images using metadata, removing NSFW images (often problematic for general participant studies) and those with any dimension smaller than 400px, leaving  $\sim$ 720 million images. Of these,  $\sim$ 440 million were still accessible via URL.

We next removed non-naturalistic images, by establishing three exclusion criteria: watermarks or banners, heavy editing (e.g., strong image filters), and not a real-world scene or object. Based on these, we manually labeled 25,000 images from a pool of 200,000 random LAION-2B samples. Labeling was initialized by first clustering the pool into 400 clusters with mini-batch k-Means on CLIP features from OpenAI's CLIP ViT-32B. These were then manually split into "natural". "not natural" and "mixed", from which 5,000 images were selected from "natural" and "non-natural" clusters. Using entropy-based uncertainty sampling, we iteratively identified the most informative samples to label. Labeling was stopped when accuracy plateaued (5-fold cross-validated). From the resulting dataset, we trained a logistic regression classifier on CLIP features (natural/non-natural), achieving an accuracy of  $\sim$ 82%. For filtering, we used a higher probability threshold to achieve a precision of 90%. Removing images that were labeled nonnaturalistic by the classifier left us with  ${\sim}120$  million images to use for evaluating existing datasets and simulations.

Analyzing random subsets of LAION-natural (6k samples, 10 seeds, as in Fig. 1D), we found semantic coverage of 93.14% of LAION-2B (filtered only for NSFW and resolution, see "Evaluating coverage"), suggesting retained visual-semantic diversity. This was further confirmed by measuring the OOD performance of random NSD-sized subsets of LAION-2B (as in Fig. 2E), which outperformed LAION-natural by only 2.34%.

**Ensuring semantic richness of LAION-natural** To validate that LAION-natural contains all concepts found in ImageNet (Deng et al. (2009)), MS COCO (Lin et al. (2014)), and ecoset (Mehrer et al. (2021)), we extracted text features for each concept and built an approximate nearest-neighbor search index on 5 million randomly sampled LAION-natural images. For each concept, we retrieved the 100 most similar images, based on the cosine similarity of normalized image and text embeddings. Manual inspection confirmed that every concept had at least one corresponding image in LAION-natural.

#### Evaluating coverage

We clustered 10 million random LAION-natural CLIP samples into 5,500 clusters using mini-batch k-Means. Cluster centroids were projected into 2D using t-Distributed Stochastic Neighbor Embedding, t-SNE (Van Der Maaten & Hinton (2008)). We further quantified coverage by looking at the opposite metric, namely measuring how many images in LAION-

natural would be considered outliers with respect to the fea-449 411 ture space of established datasets like THINGS or NSD. To450 412 this end, we used PCA on the CLIP features of the images.451 413 PCA was first trained on the features of the "covering" dataset452 414 (e.g., THINGS or NSD), with the number of principal compo-453 415 nents selected to retain 95% of the variance in that dataset.454 416 This step effectively created a lower-dimensional representa-417 455 tion of the dataset's feature distribution. 418

Next, we calculated the PCA reconstruction error for each456 419 image in both the "covering" dataset and the "covered" dataset<sup>457</sup> 420 (LAION-natural). This error quantified how well an image can458 421 be reconstructed from the principal components learned from<sup>459</sup> 422 the "covering" dataset. A higher error suggests the image's<sup>460</sup> 423 424 features are not well captured by that PCA space. To estab-461 lish a criterion for whether a LAION-natural image was cov-462 425 ered by the THINGS/NSD feature space, we defined an error463 426 threshold by fitting a generalized extreme value (GEV) distri-464 427 bution to the reconstruction errors of the "covering" dataset<sup>465</sup> 428 itself and selecting the 95th percentile from this fitted distri-466 429 bution. A LAION-natural image was considered covered if its 430 reconstruction error, when projected into and reconstructed  $\vec{A_{FR}}$ 431 from the "covering" dataset's PCA space, was below this error 432 threshold. The final coverage percentage was reported as the  $_{470}^{470}$ 433 proportion of "covered" LAION-natural images relative to the  $_{471}$ 434 total number of tested LAION-natural images. 435 472

#### 436 Simulating OOD generalizability

**GMM-based simulation** To evaluate how OOD accuracy<sub>474</sub> changed with dataset diversity and size, we simulated stim<sub>475</sub> ulus features and brain responses using a Gaussian mix<sub>476</sub> ture model (GMM) with 100 clusters in a high-dimensional<sub>477</sub> space (D = 512, same as dimensionality of CLIP embed<sub>478</sub> dings). Cluster centers  $\mu_c$  were drawn from  $N(0, \sigma_{\text{inter}}^2 *_{479} I)$ , with  $\sigma_{\text{inter}}^2 = 100/D$ , and samples were generated from<sub>480</sub>  $N(\mu_c, \sigma_{\text{intra}}^2 * I)$ , with  $\sigma_{\text{intra}}^2 = 10/D$ , ensuring well-separated<sub>481</sub> and evenly spaced clusters. From these synthetic stimu<sub>482</sub> lus features *y*, brain response *x* was generated via a lineat<sub>483</sub> mapping, a common assumption for both encoding and de<sub>484</sub> coding models of brain activity. We specifically used a ran<sub>485</sub> dom teacher weight matrix  $A(N(0, 1/\sqrt{D}))$ , adding Gaussian noise  $\varepsilon$  with variance  $\sigma_{\text{noise}}^2 = 0.25$ :

$$x = A^T * y + \varepsilon$$

To assess generalization, we generated 100 samples from a
new cluster and measured how often the predictions aligned
with the centroid of the OOD cluster. This was repeated 32
times with different OOD clusters for robust estimates.

Simulation from CLIP feature space We extended the pre-487 441 vious analysis to image features from LAION-natural, keep-488 442 ing brain response generation unchanged. We ran mini-489 443 batch k-Means clustering (k = 1,000) on CLIP features de-490 444 rived from LAION-natural and then used agglomerative clus-491 445 tering to group them into six coherent cluster groups. These492 446 groups served to divide the feature space into in-distribution493 447 and out-of-distribution folds. We measured OOD accuracy by494 448

training a model on 6,000 samples from varying parts of the in-distribution space (from combinations of 2, 3, 4, or 5 of the cluster groups, repeated 100 times per condition) and tested on 1,000 OOD samples from the OOD group, by measuring how often predictions aligned with the centroid of their source cluster (cross-validated across all cluster groups).

#### Sampling strategies

473

(Stratified) random sampling Samples were drawn uniformly at random from the entire pool of stimuli, ensuring each selection was unique (sampling without replacement). Alternatively, for stratified random sampling, the aim was to achieve proportional representation from predefined stimulus clusters. An approximately equal number of samples (n\_samples//number of clusters) was drawn from each cluster. If a cluster contained fewer samples, all available samples from that cluster were selected. Both random and stratified sampling were repeated 100 times per dataset size for robust performance estimates.

**k-Means-based sampling** The stimulus pool was clustered using mini-batch k-Means, with the number of clusters set to be the number of stimuli to sample. Data points closest to each resulting cluster centroid (in terms of Euclidean distance) were selected. To increase the speed of finding these nearest neighbors for each centroid, an Annoy index was built using the stimulus features (n\_trees = 50).

**Core-Set sampling** We iteratively selected samples to minimize the maximum distance between any data point and its nearest selected point (Sener & Savarese (2018)), broadly covering the available feature space of the dataset. We used the kCenterGreedy algorithm. This greedy algorithm starts from an empty set and selects the first sample randomly. Subsequently, it iteratively adds the data point from the remaining pool that is farthest from any of the points already selected into the core-set. This ensures that each newly added sample maximally reduces the coverage radius of the selected set, and therefore increases diversity and representation of the overall feature space.

**Sampling to optimize effective dimensionality** Effective dimensionality (ED) measures the number of meaningful axes of variance in a dataset (Del Giudice (2021)). We used the participation ratio of CLIP image features to estimate ED:

$$\mathsf{ED} = \frac{(\sum_{i=1}^{K} \lambda_i)^2}{\sum_{i=1}^{K} \lambda_i^2}$$

where  $\lambda_i$  are the principal components. Intuitively, a low ED suggests an over-representation of semantic concepts - for example, if a dataset contains only mountains and beaches, variance is mostly explained by a single "beach-or-mountain" dimension. A more diverse dataset, also containing meadows, forests, or cities, would require more dimensions.

Given this insight, we also greedily sampled to maximize the ED of the dataset. Initialization of the selected set was performed in one of two ways: either by selecting two sam-

ples uniformly at random, or by using samples closest to clus-549 495 ter centroids derived from a mini-batch k-Means clustering.550 496 Samples were added iteratively. In each step, a candidatessi 497 pool was generated by drawing 10 random samples from each552 498 cluster of a precomputed clustering. To avoid selecting very 499 similar items, candidates that were too close (Euclidean dis-553 500 tance < 0.1) to already selected samples were filtered out.554 501 The ED was then estimated for each remaining candidate, if555 502 it were added to the existing image set. This ED calculation 556 503 was performed with an incremental update formula for the co-557 504 variance matrix for efficiency and was parallelized across 32558 505 CPU cores to speed up selection. The candidate that yielded 559 506 the highest ED for the augmented set was then added to the500 507 selected samples. 508 561

Margin-based, adaptive sampling We also tested if epistemic uncertainty could guide sample selection using a margin-based active learning strategy (Balcan et al. (2007)) with a logistic regression model predicting discretized brain responses from stimulus features. Before model training, the stimulus features (X) were standardized to have zero mean and unit variance.

Sampling was initialized with 100 random samples. Tar-516 get brain data (Y) was discretized into three equally popu-517 lated bins per dimension using quantile-based binning (e.g., 518 low, medium, high response categories). A logistic regres-519 sion model was trained for each target dimension (i.e., for 520 each voxel or ROI whose response was being predicted). A 521 candidate pool was drawn by randomly selecting 10 samples 522 per cluster from a precomputed mini-batch k-Means cluster-523 ing (k=1,000). For each candidate, bin probabilities were pre-524 dicted by the trained logistic regression model(s), and uncer-525 tainty was measured as the margin between the top two prob-526 abilities (a lower margin indicates higher uncertainty, as the 527 model is less decisive). Margins were calculated for each di-528 mension independently for a given candidate, and these mar-529 gins were then averaged to get a single uncertainty score for 530 that candidate. In each iteration, the 100 images with the low-531 est average margins (highest uncertainty) were selected and 532 added to the training set, and the model was retrained. This 533 process was repeated until reaching the required dataset size. 534

#### 535 Evaluation of concept distribution

To test to what extent a sample dataset preserved the distri-536 bution of concepts of the stimulus pool, we evaluated how the 537 sampling strategies changed the concept distribution using 538 a subset of LAION-natural (100,000 images). As LAION-2B 539 only provides image captions, and no image-level keywords, 540 we first used Gemini 1.5 Flash (8B), configured with gener-541 ation parameters: temperature=1, top\_p=0.95, top\_k=40, and 542 max\_output\_tokens=8192, to list keywords for each image, us-543 ing the prompt "Describe these images in as many keywords 544 as you like. Return as a list of keywords.". This generated 545 75,535 unique terms, at a cost of  $\sim$ \$2.54. We then filtered 546 these keywords to only include concrete nouns (Concrete-547 ness > 4; Brysbaert et al. (2014)) and availability of natural 548

language frequency (Brysbaert & New (2009)). After filtering, 3,563 keywords remained, which were clustered (HDBSCAN, min\_samples=1) into 231 groups, allowing comparison of cluster occurrence depending on sampling strategy.

#### Implementation details

CLIP features were extracted using the CLIP ViT-32/B model, provided by OpenAI (https://github.com/openai/CLIP). Approximate nearest-neighbor search was implemented with Annoy (https://github.com/spotify/annoy; n\_trees=100). For kCenterGreedy, we used Google's active learning framework. Clustering, t-SNE projection, PCA, Ridge, and logistic classifier fitting were implemented using scikit-learn (Pedregosa et al. (2011)). Active-learning classifier training was implemented with AliPy (Tang et al. (2019)).

## References

616

563

Allen, E. J., St-Yves, G., Wu, Y., Breedlove, J. L., Prince, J. S.,<sup>617</sup>
 Dowdle, L. T., ... Kay, K. (2022). A massive 7t fmri datasef<sup>618</sup>
 to bridge cognitive neuroscience and artificial intelligence.<sup>619</sup>
 *Nature Neuroscience*, *25*(1), 116–126.

Balcan, M.-F., Broder, A., & Zhang, T. (2007). Margin based<sup>621</sup>
 active learning. In *Learning theory* (pp. 35–50). Berlin, Hei-<sup>622</sup>

570 delberg: Springer.

571 Brysbaert, M., & New, B. (2009). Moving beyond kučera<sup>624</sup> 572 and francis: A critical evaluation of current word frequency<sup>625</sup>

norms and the introduction of a new and improved word fre-626

574quency measure for american english. Behavior Research<sup>27</sup>575Methods, 41(4), 977–990.628

Brysbaert, M., Warriner, A. B., & Kuperman, V. (2014). Con-629
 creteness ratings for 40 thousand generally known english<sup>630</sup>
 word lemmas. *Behavior Research Methods*, 46(3), 904–631
 911.

Chang, N., Pyles, J. A., Marcus, A., Gupta, A., Tarr, M. J., &<sup>633</sup>
 Aminoff, E. M. (2019). Bold5000, a public fmri dataset while<sup>634</sup>

viewing 5000 visual images. *Scientific Data*, *6*(1), 49.

583 Contier, O., Baker, C. I., & Hebart, M. N. (2024). Distributed<sup>636</sup>

representations of behaviour-derived object dimensions in<sup>637</sup> the human visual system. *Nature Human Behaviour*, *8*(11),<sup>638</sup> 2179–2193.

587 Conwell, C., Prince, J. S., Kay, K. N., Alvarez, G. A., & Konkle,<sup>640</sup>

T. (2024). A large-scale examination of inductive biases<sub>641</sub> shaping high-level visual representation in brains and ma-642 chines. *Nature Communications*, *15*(1), 9383. 643

chines. Nature Communications, 15(1), 9383.

Del Giudice, M. (2021). Effective dimensionality: A tutorial.<sup>644</sup>
 *Multivariate Behavioral Research*, *56*(3), 527–542.

Demircan, C., Saanum, T., Pettini, L., Binz, M., Baczkowski,<sup>646</sup>
 B. M., Doeller, C. F., ... Schulz, E. (2023). Eval.<sup>847</sup>
 *uating alignment between humans and neural network*<sup>648</sup>
 *representations in image-based learning tasks.* arXiv<sup>649</sup>
 preprint. Retrieved from http://arxiv.org/abs/2306<sub>650</sub>

598 .09377 (arXiv:2306.09377 [cs.LG]) 651

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei,<sup>652</sup>
 L. (2009). Imagenet: A large-scale hierarchical image<sup>653</sup>
 database. In *Proceedings of the 2009 ieee conference or*<sup>654</sup>
 *computer vision and pattern recognition (cvpr)* (pp. 248–655
 255). IEEE.

Doerig, A., Sommers, R. P., Seeliger, K., Richards, B., Ismael,<sup>657</sup>
 J., Lindsay, G. W., ... Kietzmann, T. C. (2023). The neuro-<sup>658</sup>
 connectionist research programme. *Nature Reviews Neu*-<sup>659</sup>
 *roscience*, 24(7), 431–450.

Hebart, M. N., Contier, O., Teichmann, L., Rockter, A. H.,<sup>661</sup>
 Zheng, C. Y., Kidder, A., ... Baker, C. I. (2023). THINGS-662
 data, a multimodal collection of large-scale datasets for in-663

vestigating object representations in human brain and be-664

havior. *eLife*, *12*, e82580. Retrieved from https://doi665 .org/10.7554/eLife.82580 doi: 10.7554/eLife.82580

614 Kaniuth, P., Mahner, F. P., Perkuhn, J., & Hebart, M. N. (2024) 667

A high-throughput approach for the efficient prediction of 615

perceived similarity of natural objects. bioRxiv preprint. Retrieved from https://doi.org/10.1101/2024.06.28 .601184 (bioRxiv 2024.06.28.601184) doi: 10.1101/ 2024.06.28.601184

Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., ... Dollár, P. (2014). *Microsoft COCO: Common objects in context.* arXiv preprint. Retrieved from http:// arxiv.org/abs/1405.0312 (arXiv:1405.0312 [cs.CV])

Mehrer, J., Spoerer, C. J., Jones, E. C., Kriegeskorte, N., & Kietzmann, T. C. (2021). An ecologically motivated image dataset for deep learning yields better models of human vision. *Proceedings of the National Academy of Sciences of the United States of America*, *118*(8), e2011417118.

Muttenthaler, L., Dippel, J., Linhardt, L., Vandermeulen, R. A., & Kornblith, S. (2022). *Human alignment of neural network representations.* arXiv preprint. Retrieved from http://arxiv.org/abs/2211.01201 (arXiv:2211.01201 [cs.CV])

Naselaris, T., Allen, E., & Kay, K. (2021). Extensive sampling for complete models of individual brains. *Current Opinion in Behavioral Sciences*, 40, 45–51.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay, É. (2011). Scikitlearn: Machine learning in python. *Journal of Machine Learning Research*, 12(85), 2825–2830.

Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., ... Jitsev, J. (2022). *LAION-5B: An open large-scale dataset for training next generation image text models.* arXiv preprint. Retrieved from http://arxiv .org/abs/2210.08402 (arXiv:2210.08402 [cs.CV])

Sener, O., & Savarese, S. (2018). Active learning for convolutional neural networks: A core-set approach. arXiv preprint. Retrieved from http://arxiv.org/abs/1708 .00489 (arXiv:1708.00489 [cs.CV])

Shirakawa, K., Nagano, Y., Tanaka, M., Aoki, S. C., Majima, K., Muraki, Y., & Kamitani, Y. (2024). Spurious reconstruction from brain activity. arXiv preprint. Retrieved from http://arxiv.org/abs/2405 .10078 (arXiv:2405.10078 [q-bio.NC])

Takagi, Y., & Nishimoto, S. (2023, June). High-resolution image reconstruction with latent diffusion models from human brain activity. In *Proceedings of the 2023 ieee/cvf conference on computer vision and pattern recognition (cvpr)* (pp. 13993–14003). Vancouver, BC, Canada: IEEE/CVF. Retrieved from https://doi.org/10.1109/cvpr52729 .2023.01389 doi: 10.1109/cvpr52729.2023.01389

Tang, Y.-P., Li, G.-X., & Huang, S.-J. (2019). *ALiPy: Active learning in python.* arXiv preprint. Retrieved from http://arxiv.org/abs/1901.03802 (arXiv:1901.03802 [cs.LG])

Van Der Maaten, L., & Hinton, G. (2008, November). Visualizing data using t-sne. *Journal of Machine Learning Research*, 9, 2579–2605. Wang, A. Y., Kay, K., Naselaris, T., Tarr, M. J., & Wehbe, L.
(2023). Better models of human high-level visual cortex
emerge from natural language supervision with a large and
diverse dataset. *Nature Machine Intelligence*, *5*(12), 1415–
1426.

## 674 Code availability

675 All code used for the analyses and generation of figures

 $_{\rm 676}$   $\,$  in this study, including the pre-trained LAION-natural class

sifier mentioned in the text, is publicly available on GitHub:

678 https://github.com/andropar/how-to-sample.

## **Supplementary Material**

### <sup>680</sup> Natural vs. Non-natural images



Figure S1: Examples of images from LAION-2B that were considered either natural (A) or not natural (B). The criteria for natural images were: "no heavy editing (e.g. high saturation / contrast, collages, cropped objects without background) or filter overlaid (e.g. black-and-white)", "no watermarks or text banners" and "must be a real-world object or scene (e.g. no screenshots of websites or video games)".

679

### 681 Evaluation of OOD accuracy using NSD fMRI data



Figure S2: Impact of training set diversity on OOD accuracy in NSD. For each subject, CLIP features of presented images were clustered using mini-batch k-Means. Iteratively using one cluster as the OOD test set, training sets of a fixed size (N=500) were created by stratifying samples from a varying number (k) of the remaining clusters. Thin grey lines represent individual subject data, and the red line shows the mean across subjects, with shaded areas indicating the standard error of the mean. These results suggest that increased visual diversity improves generalization performance, even while keeping the total number of training samples constant.

## 682 Examples of clusters not covered by THINGS or NSD



Figure S3: Examples of distinct clusters that were not covered by THINGS or NSD. These include certain sporting events, architectural styles, landscapes, political figures, images of natural disasters, activities and many more. Clusters were manually selected from the 50 largest clusters not covered by the other datasets, to avoid repetitions in semantic concepts.