### 000 TOWARDS FLEXIBLE PERCEPTION WITH 001 002 VISUAL MEMORY

Anonymous authors

Paper under double-blind review

## ABSTRACT

Training a neural network is a monolithic endeavor, akin to carving knowledge into stone: once the process is completed, editing the knowledge in a network is nearly impossible, since all information is distributed across the network's weights. We here explore a simple, compelling alternative by marrying the representational power of deep neural networks with the flexibility of a database. Decomposing the task of image classification into image similarity (from a pretrained embedding) and search (via fast nearest neighbor retrieval from a knowledge database), we build a simple and flexible visual memory that has the following key capabilities: (1.) The ability to flexibly add data across scales: from individual samples all the way to entire classes and billion-scale data; (2.) The ability to remove data through unlearning and memory pruning; (3.) An interpretable decision-mechanism on which we can intervene to control its behavior. Taken together, these capabilities comprehensively demonstrate the benefits of an explicit visual memory. We hope that it might contribute to a conversation on how knowledge should be represented in deep vision models—beyond carving it in "stone" weights.

025 026 027

028

033

003 004

010 011

012

013

014

015

016

017

018

019

021

#### INTRODUCTION 1

029 In the pretty diagrams on "Intro to Machine Learning" slides, an ideal ML workflow looks like this: Data collection, preprocessing, choosing a model, training, evaluation, deployment. Happy 031 ending-the model is deployed, the users love it, and one can finally go on that well-deserved vacation and catch up on the latest AGI memes.

034 Until, of course, the enemy of any ideal world sets in: reality. The real world constantly keeps changing, and so do data requirements. New data and datasets become available, and existing ones 035 become deprecated for a variety of reasons, including concerns around fairness, biases or unsafe 036 content. Knowledge changes, and concepts drift (Tsymbal, 2004; Lu et al., 2018): Phones and cars 037 look different today than they did a few years ago, and different from how they will look in the future. When it comes to data, the only constant is change (Cao & Yang, 2015; Bourtoule et al., 2021; Nguyen et al., 2022; Zhang et al., 2023). Consequently, from a modeling perspective, in 040 order to keep up with this change one would ideally want to constantly re-train or fine-tune models, 041 which is of course not feasible. In short, as anyone who has ever deployed a model has experienced 042 firsthand, one is constantly battling the symptoms of a single underlying cause: the fact that deep 043 learning models have a static knowledge representation entangled in millions or billions of model 044 parameters. We, among many others working on memory (e.g. Weston et al., 2014; Chen et al., 2018; Wu et al., 2021; Iscen et al., 2022; Nakata et al., 2022; Iscen et al., 2023; Prabhu et al., 2023; 045 Gui et al., 2024; Shao et al., 2024; Silva et al., 2024), believe that this is not a great way to represent 046 visual knowledge for deep learning. Instead, we argue that we should build models that cleanly 047 separate representation (how things are represented, e.g. through feature embeddings) from visual 048 memory (what is known). In short, deep learning models need a flexible visual memory: a way to 049 explicitly utilize and edit knowledge. 050

051 In this work, we build a simple visual memory for classification and show that it has seven desirable capabilities, including the ability to flexibly add data across scales (from individual samples 052 to classes and even billion-scale data), the ability to remove data from our model's classification process through machine unlearning and memory pruning, and a simple, interpretable decisionmechanism on which we can intervene to control its behavior. Our main goal is to provide a compelling idea of how beneficial a flexible visual memory for deep learning can be from a variety of perspectives and capabilities. From a technical standpoint, we aim for simplicity: retrieving knearest neighbors (in an embedding feature space) along with their labels to classify a query image. This approach allows us to investigate where a simple visual memory mechanism helps, where its limitations may be, and where there might be opportunities for improvement through a more complex system. We hope that by demonstrating clear benefits from a simple visual memory, this article might contribute to a conversation on how knowledge ought to be represented in deep vision models.

- Here are some highlights of this article:
  - 1. Improved aggregation of retrieved samples: we propose using *RankVoting*, a power-law weighting that surpasses previous SOTA (SoftmaxVoting) for a deep learning based memory.
  - 2. Re-ranking samples using a vision-language model achieves 88.5% top-1 ImageNet validation accuracy, improving over both DinoV2 ViT-L14 kNN and linear probing.
    - 3. Flexible perception: the visual memory achieves perfect unlearning, scales to billion-scale data without additional training, and enables controlling sample influence via *memory pruning*.

We argue that the way current deep learning models represent knowledge (static knowledge representation, hard to update, hard to unlearn, hard to understand how a decision is made) is problematic. As an alternative, we built a working proof-of-concept: By building on the long history of nearest neighbor methods, and "marrying" them with a powerful deep learning representation (such as SSL features from DinoV2) and a billion-scale visual memory.

Related work. The concept of a visual memory has a long history in ML, neuroscience and psy-076 chology. In psychology, exemplar theory posits that humans recognize objects by comparing them to 077 existing examples in visual memory (Medin & Schaffer, 1978; Nosofsky, 1986; Dopkins & Gleason, 078 1997; Jäkel et al., 2008; Nosofsky, 2011), like the ALCOVE model (Kruschke, 2020). In ML, prior 079 to deep learning, instance-based learning (also known as memory-based learning) was a popular alternative to model-based learning (Aha et al., 1991; Quinlan, 1993). For instance, Turk & Pent-081 land (1991) used nearest neighbor methods to classify faces, and Sivic & Zisserman (2003) build 082 a visual memory inspired by text retrieval for object retrieval from videos. In recent years, hybrid 083 approaches have started to combine the benefits of both approaches. Deep neural network variants 084 (model-based since they learn generalized abstractions of data) of k-nearest neighbor algorithms 085 (instance-based since they compare new data to existing exemplars in memory) have been proposed with various motivations, including few-shot learning (Wang et al., 2019b; Yang et al., 2020; Bari 086 et al., 2021), improving adversarial robustness (Sitawarin & Wagner, 2019; Papernot & McDaniel, 087 2018; Rajani et al., 2020), medical image classification (Zhuang et al., 2020), confidence calibration 088 (Papernot & McDaniel, 2018), interpretability (Papernot & McDaniel, 2018; Wallace et al., 2018; 089 Lee et al., 2020; Rajani et al., 2020), image denoising (Plötz & Roth, 2018), retrieval-augmented 090 learning (Khandelwal et al., 2019; Drozdov et al., 2022), anomaly and out-of-distribution detection 091 (Bergman et al., 2020; Sun et al., 2022). Recently, Nakata et al. (2022) tested a kNN-based vi-092 sual memory up to ImageNet-scale (1.28M images), and Khandelwal et al. (2019); Wu et al. (2021) applied kNN-based approaches to neural language models.

094 095 096

064

065

066

067 068

069

## 2 BUILDING A RETRIEVAL-BASED VISUAL MEMORY FOR CLASSIFICATION

Given a dataset  $\mathcal{D}_{\text{test}} := \{(\tilde{x}_1, y_1), \dots, \tilde{x}_n, y_n\}$ , we want to classify each image  $\tilde{x}_i \in \mathcal{D}_{\text{test}}$ . Our classification approach consists of two steps: (i) building a visual memory, and (ii) fast nearest neighbor based inference using the visual memory.

100 101

102

2.1 BUILDING A VISUAL MEMORY

Our visual memory retrieves (image, label) pairs from an image dataset when a query is made by
directly retrieving those images that are considered similar to a test image according to a model.
The model is a fixed pre-trained image encoder, meaning that no training takes place when adding
information to visual memory. No copies of the dataset are stored in the visual memory. Instead,
feature maps are extracted from the model based on a set of images related to the downstream
classification task at hand, such as a standard training set. For our experiments, our visual memory



Figure 1: **Reliability of retrieved memory samples.** This plot visualizes the ImageNet (**left**) and iNaturalist (**right**) top-1 validation accuracy of a single retrieved neighbor depending on the index of the neighbor (index 0: nearest neighbor). In both datasets and across models, the decrease in accuracy with increasing neighbor index follows smooth trajectories and can be approximated by a two-parameter logarithmic fit (black lines).

127 comprises of features extracted from a dataset like the ImageNet-1K (Russakovsky et al., 2015) 128 training set using different encoders like DinoV2 (Oquab et al., 2023) and CLIP (Radford et al., 129 2021). Thus, given a pretrained image encoder,  $\Phi$ , and a dataset of (image, label) pairs  $\mathcal{D}_{\text{train}}$  := 130  $(x_1, y_1), (x_2, y_2), \cdots, (x_N, y_N)$ , we obtain features  $z_i := \Phi(x_i), \forall x_i \in \mathcal{D}_{\text{train}}$ . Subsequently, the 131 feature maps and corresponding label pairs are put in a database thereby creating VisualMemory := 132  $\{(z_1, y_1), (z_2, y_2), \cdots, (z_N, y_N)\}$  for classification. For both DinoV2 and CLIP, we use the last 133 image embedding layer as a feature space.

### 2.2 RETRIEVAL-BASED CLASSIFICATION USING VISUAL MEMORY

Given a query image  $\tilde{x} \in \mathcal{D}_{\text{test}}$ , we extract its feature map,  $\tilde{z} = \Phi(\tilde{x})$ . We then query VisualMemory to extract k feature vectors, Neighbors $(\tilde{x}) := \{(z_{[1]}, y_{[1]}), (z_{[2]}, y_{[2]}((, \cdots, (z_{[k]}, y_{[k]}))\}, \text{ that are$  $closest to the query features <math>\tilde{z}$  using the cosine distance, which is the default retrieval similarity measure for SSL models like DinoV2. Neighbors $(\tilde{x})$ , are ordered by distance i.e.

 $\mathsf{dist}(\tilde{\boldsymbol{z}}, \boldsymbol{z}_{[i]}) \leq \mathsf{dist}(\tilde{\boldsymbol{z}}, \boldsymbol{z}_{[i]}), \ \forall i \leq j.$ 

121

122

123

124

125 126

134

135

140

141

We then assign a weight,  $w_i$ , to each neighbour  $(z_{[i]}, y_{[i]})$  and aggregate the scores for each neighbour with the same label. Finally, we assign that label to the query image with the highest aggregate score. We implemented retrieval based classification using one of the following two approaches:

**1. Fast inference using matrix multiplication on GPUs/TPUs:** For smaller datasets like ImageNet, we saved VisualMemory as a matrix of size num\_images  $\times$  num\_dims. During inference, for an encoded query image of size  $1 \times$  num\_dims, we computed the dot product of this encoded image with every entry in VisualMemory getting a matrix of size num\_images  $\times 1$ . We then computed the *k* nearest neighbors using the arg max operation.

2. Fast and scalable nearest neighbor search: We used ScaNN (Guo et al., 2020) for accelerating nearest neighbor search at scale. Specifically, we saved the VisualMemory as a database and used ScaNN for fast lookup of nearest neighbors during inference. This method scales easily to billion-scale memory (cf. Section 3.3). Appendix J details latency and storage; storing features requires only about 1–3% of the space of storing the dataset itself and even with a 1B memory.

We mentioned earlier that we retrieve a set of neighbors, Neighbors( $\tilde{x}$ ) and aggregate information across them to make a classification decision. In order to understand how reliable (i.e., accurate) retrieved memory samples are from the first to the 100th neighbor, we systematically analyze neighbor reliability in Figure 1. As expected, reliability decreases as the neighbor index k increases, but even at large k the neighbors contain above-chance information about the ground truth class. This suggests that aggregating information across different neighbors may be beneficial to decision-making, leading to the question: *What is the best aggregation strategy?* We empirically study this by testing different weighting strategies for aggregation:

176

177

178

179 180 181



Figure 2: Aggregating information across retrieved memory samples. (left) Existing aggregation methods are overconfident in distant neighbors, resulting in the paradox of decaying ImageNet-1K accuracy with more information. The same pattern is also seen for other models and datasets in the Appendix (Figures 7 and 8). (right) This is not the case for RankVoting, a simple power-function based method which reaches higher and stable performance across models and choices of k.

182 **Plurality voting:** Each neighbour in Neighbors( $\tilde{x}$ ) is assigned an equal weight of 1.0. This is the 183 classic, most simple voting method and used for instance by Nakata et al. (2022).

184 185 185 186 187 Distance voting: Each neighbour in Neighbors( $\tilde{x}$ ) is assigned a weight based on its Cosine distance to the query image  $\tilde{x}$  i.e.  $w_i = \exp\left(-\operatorname{dist}(\tilde{z}, \boldsymbol{z}_{[i]})\right)$ . This approach has been used by Khandelwal et al. (2019) for nearest neighbor language models.

**Softmax voting:** Each neighbour is assigned a weight based on the softmax function i.e.  $w_i =$ softmax (dist( $\tilde{z}, z_{[i]}$ ),  $\tau$ ) where  $\tau$  is the temperature. This voting method is considered state-ofthe-art; for example nearest neighbor accuracies of self-supervised models are reported using this method. A temperature of  $\tau = 0.07$  frequently appears in literature (Wu et al., 2018; Caron et al., 2021; Oquab et al., 2023) and is reported as a parameter "which we do not tune" in the Dino paper (Caron et al., 2021, p. 18). We observe that performance is sensitive to this parameter; other temperatures perform worse. We therefore follow the literature in using  $\tau = 0.07$ .

**Rank voting:** We propose using a simple aggregation approach wherein each neighbour is assigned a power-function weight based on its rank in the ordered set Neighbors( $\tilde{x}$ ) i.e.  $w_i = 1/(\alpha + \operatorname{rank}_i)$ where rank<sub>i</sub> is *i* and  $\alpha$  is an offset to avoid division by zero that is set to 2.0. This is similar, though not identical to, Gou et al. (2011) who used power-law weighting in a different context.

In Figure 2a, we compare the top-1 ImageNet validation accuracy of different ranking methods as 199 a function of number of neighbours, with the ImageNet-1K training set as the visual memory using 200 the DinoV2/ViT-L14 model as the featurizer. Paradoxically, existing aggregation methods like plu-201 rality voting, distance-based voting, and softmax voting show *decaying* performance as the provided 202 information (number of nearest neighbors) increases. This suggests that the methods are overconfi-203 dent in distant neighbors, assigning them too much weight. Our simple, parameter-free rank based 204 voting method, however, leads to an increase in performance with more neighbors until a certain k205 after which the performance plateaus, which is the ideal scenario (Figure 2b). Furthermore, rank-206 based voting also outperforms baselines in absolute terms; quantitative comparisons can be found in the Appendix (Tables 4 to 8) where we also study the influence of hyperparameters (Figure 9). 207 This indicates that a simple, power-function based method can reliably integrate information across 208 retrieved memory samples. 209

Gemini re-ranking. Our results above demonstrate that different aggregation strategies have a large impact on downstream performance. How far can we push the upper limit on aggregating information from different neighbors? We perform a controlled experiment using the Gemini 1.5
Flash model (Reid et al., 2024) to test this: We add the 50 nearest neighbors from DinoV2 ViT-L14 for a query image along with their labels into Gemini's context. We then query Gemini to predict the query image's label. This achieves 88.5% ImageNet validation accuracy, a substantial improvement over both DinoV2 ViT-L14 kNN (83.5%) and linear probing (86.3%) performance. Interestingly,

Gemini's performance is mainly driven by the neighbor information through in-context learning
since it only achieves 69.6% accuracy without neighbors (when just the query image is provided to
the model). The performance improvement highlights the potential of using vision-language models
as a visual memory re-ranker. Given that our main goal is to explore a simple visual memory system,
we mostly focus on non-Gemini ranking methods throughout our analysis.

# 3 CAPABILITIES OF A VISUAL MEMORY

Our primary goal is to motivate the concept of a machine *visual memory* from a variety of different perspectives. To this end, we investigate how such a memory can benefit the following capabilities: 3.1 Flexible lifelong learning: adding novel OOD classes; 3.2 Flexibly trading off compute and memory; 3.3 Flexibly adding billion-scale data without training; 3.4 Flexible removal of data: machine unlearning; 3.5 Flexible data selection: memory pruning; 3.6 Flexibly increasing dataset granularity; 3.7 Interpretable & attributable decision-making.

### 3.1 FLEXIBLE LIFELONG LEARNING: ADDING NOVEL OOD CLASSES (DATA AND LABELS)

Standard classifiers, whether trained end-to-end (supervised models) or with a linear classifier (self-supervised models), are not able to handle new information without re-training. For instance, adding new classes or changing labels in an existing model usually involves either re-training or fine-tuning parts of the model. A retrieval-based visual memory, in contrast, is able to process such information in a natural and flexible way, aligning with the requirements of lifelong learning (Parisi et al., 2019). We tested this by adding data for 64 new classes, along with their new labels, to the visual memory of a pre-trained DinoV2 ViT-L14 model (in addition to the ImageNet train set, which is in-distribution for the model). We took the new classes from the NINCO dataset (Bitterwolf et al., 2023), a dedi-cated OOD dataset that is designed to have no overlap with existing ImageNet labels and samples. This requires the model to transfer what it has learned to new, unseen concepts. The new task is therefore harder, as the model has to retrieve images from both in-distribution and OOD classes. The resulting visual memory has 1064 classes (1K from ImageNet and 64 from NINCO). Table 1 shows that with a visual memory it is possible to add new classes such that the in-distribution accu-racy is maintained without catastrophic forgetting (the new classes only change ImageNet validation performance by 0.02–0.04% depending on the aggregation method), while at the same time reaching very high accuracy on the new OOD classes (approx. 87% top-1) without any training. Figure 12 in the appendix confirms that the samples are indeed OOD for the model, as demonstrated by larger distances to nearest neighbors. This highlights that a visual memory is capable of flexibly adding new information—an important capability since the world is not static. Furthermore, the memory is incredibly robust towards label corruption up to 60% random labels, as shown in Appendix D. 

Table 1: Flexible lifelong learning: adding novel OOD classes. A visual memory of DinoV2
ViT-L14 with ImageNet-train (IN-train) as the memory database is able to handle a simple "insert into memory" operation for 64 out-of-distribution classes (data and labels) from the NINCO dataset (Bitterwolf et al., 2023), leading to high performance on the new classes without affecting top-1 ImageNet validation accuracy.

$\begin{array}{c} \text{memory} \rightarrow \\ \text{query} \rightarrow \end{array}$	IN-train	IN-train-	and-NINCO
	IN-val	IN-val	NINCO
no aggregation	81.1	81.1	86.4
Plurality Voting	83.2	83.2	86.9
Distance Voting	83.3	83.3	87.1
Softmax Voting	83.6	83.5	87.5
Rank Voting	83.6	83.6	87.4

# 3.2 FLEXIBLY TRADING OFF COMPUTE AND MEMORY

Next, we turn our attention to studying the scaling behaviour of visual memory with increasing memory model size. We hypothesize that bigger models will be able to attain similar performance as



Figure 3: **Memory scaling: flexibly trading off compute and memory.** ImageNet top-1 validation error decreases systematically as the memory size is increased (i.e., recognition accuracy increases with scale). (**left**) Million-scale memory consisting of ImageNet-train labels. (**right**) Billion-scale memory bank consisting of machine-generated pseudo labels on the JFT dataset (Zhai et al., 2022). Accuracy continues to decrease even with billion-scale data in memory. The roughly constant offset between models of different sizes suggests the possibility of a flexible trade-off: The same error rate can be achieved with a small model and large memory, or a large model and a small memory.

291 smaller models with lesser amount of visual memory. This is because, all else being equal, a bigger 292 model should be a better featurizer that requires fewer examples in memory to represent different 293 concepts. We empirically study the scaling behaviour of visual memory based retrieval systems in Figure 3a using models of different sizes like DinoV2 ViT models of sizes S/14 (21M params), B/14 (86M params), and L/14 (300M params), as well as CLIP ViT models of sizes B/16 and L/14. We 295 plot the top-1 error rate as a function of number of images in visual memory. The plot demonstrates 296 that for each model, the error rate consistently decreases as we increase the visual memory size. 297 Notably, already with a single exemplar per class in memory, ImageNet validation performance is 298 far beyond chance (41% top-1 error for DinoV2 ViT-L14). It also visualizes the possibility of a 299 flexible trade-off between model size and memory size: e.g. for the different DinoV2 models, the 300 S/14, B/14, and L/14 variant achieve similar performance at 1.28M,  $\sim$ 150K, and  $\sim$ 70K memory 301 capacity respectively. In line with Nakata et al. (2022), this indicates that a smaller model with large 302 memory can match the performance of a larger model with smaller memory.

303 304

283

284

285

286

287

288

289 290

304 305 3.3

### 3.3 FLEXIBLY ADDING BILLION-SCALE DATA WITHOUT TRAINING

306 Billion-scale dataset with pseudo labels. As demonstrated in Section 3.2, performance systemat-307 ically improves with increased memory size across both small and large models. We here test how 308 far this trend holds beyond relatively small-scale, well-curated settings like ImageNet-1K by scal-309 ing visual memory to the billion-scale unlabeled data regime. We obtain a large-scale dataset from the union of the ImageNet-1K train set and a subset of the JFT-3B dataset (Zhai et al., 2022). To 310 this end, we treat JFT as an unlabeled dataset by ignoring its original labels and instead obtaining 311 pseudo labels by running them through ViT-22B-224px (Dehghani et al., 2023), a highly performant 312 classifier. We excluded images whose labels do not have a correspondence with the ImageNet labels. 313

Scaling. In Figure 3b, we show the downstream ImageNet validation performance of two DinoV2ViTs as a function of memory size. The plot demonstrates that even in the billion-scale data regime,
validation error decreases when increasing memory size without any training. The gain from more
data is most prominent when having fewer samples in memory (e.g., going from 1 to 10 samples per
class). In log-log space, a logarithmic function fits the empirical scaling trend well. In the literature,
simple scaling trends such as the one we observe are powerful predictors of scaling behaviour for
different model and dataset sizes (Hestness et al., 2017; Kaplan et al., 2020; Hoffmann et al., 2022).

Out-of-distribution performance. In order to understand whether the benefits of increased memory
 size transfer to out-of-distribution (OOD) data, we compared DinoV2 ViT-L14 once with ImageNet train in memory and once with JFT pseudo-labels in memory. The models are evaluated on the
 ImageNet-A (Hendrycks et al., 2021), ImageNet-R (Hendrycks et al., 2020), ImageNet-Sketch

324 (Wang et al., 2019a), ImageNet-V2 (Shankar et al., 2020), and ImageNet-ReaL (Beyer et al., 2020) 325 datasets. As an additional well-performing yet "inflexible" baseline, we report linear probing ac-326 curacies from the DinoV2 paper (Oquab et al., 2023). Table 2 shows that visual memory scaled 327 with JFT data improves OOD performance across all datasets compared to an ImageNet-based vi-328 sual memory. Gemini re-ranking again improves leads to performance gains. Overall, the finding that memory scale transfers to OOD improvements is important in the context of continual learning, 329 where a flexible visual memory can easily incorporate newly available data that the model was not 330 trained on and improve performance both in- and out-of-distribution. 331

Table 2: **OOD evaluation.** Out-of-distribution performance improves with larger visual memory size. Across all datasets, a visual memory with JFT memory outperforms ImageNet memory demonstrating advantages of scaling visual memory for OOD performance. Probe details: Appendix I.

Model	Method	IN-A	IN-R	IN-Sketch	IN-V2	IN-ReaL
DinoV2 ViT-L14	linear probe	71.3	74.4	59.3	78.0	89.5
DinoV2 ViT-L14	ImageNet memory + Gemini re-ranking	58.8 68.4	62.8 72.3	61.5 72.5	75.6 81.7	87.1 89.9
DinoV2 ViT-L14	JFT memory + Gemini re-ranking	61.1 69.6	73.7 81.4	68.0 75.0	77.6 82.3	88.2 90.5

### 3.4 FLEXIBLE REMOVAL OF DATA: MACHINE UNLEARNING

347 The world is not static. Thus, in addition to the need to flexibly add novel data, there is often a 348 need to remove the influence of specific training data from a model's decision-making process after 349 it has been trained (Cao & Yang, 2015; Bourtoule et al., 2021; Nguyen et al., 2022; Zhang et al., 2023). A range of intricate methods are being developed to remove or reduce the influence of certain 350 training samples (Gupta et al., 2021; Sekhari et al., 2021; Ullah et al., 2021; Kurmanji et al., 2024; 351 Sepahvand et al., 2024)—a challenging endeavour if knowledge is embedded in millions or billions 352 of model weights. In contrast, for models with an explicit visual memory, machine unlearning 353 becomes as simple as removing the dataset sample from the visual memory. For instance, after 354 adding the NINCO dataset (Bitterwolf et al., 2023) into visual memory, we can remove any NINCO 355 sample with outstanding performance on all three key unlearning metrics reported by Liu (2024): 356 Efficiency: How fast is the algorithm compared to re-training? (Lightning fast.) Model utility: Do 357 we harm performance on the retain data or orthogonal tasks? (Not at all.) Forgetting quality: How 358 much and how well are the 'forget data' actually unlearned? (Completely and entirely.) Can machine 359 unlearning therefore be solved with a visual memory? If the embedding model is trained on data 360 that needs to be unlearned, machine unlearning remains challenging. If, however, the embedding model is trained on a safe, generalist dataset (e.g., a publicly available image dataset) and data that 361 may need to be considered for unlearning later is simply put into the visual memory, then machine 362 unlearning indeed becomes as simple as deleting a datapoint from the visual memory. This can be 363 particularly helpful for tasks that may require private or confidential data—a model can be trained 364 on publicly available datasets to learn general and information features and the private data can be 365 added to a visual memory on local devices for downstream tasks to preserve privacy. 366

367 368

332

333

334

346

### 3.5 FLEXIBLE DATA SELECTION: MEMORY PRUNING

369 The ability to flexibly remove the influence of certain datapoints is not just desirable in the unlearn-370 ing sense, but also advantageous in the context of dataset pruning, an emerging field that analyzes 371 the quality of individual data points. The goal of dataset pruning is to retain only *useful* samples, 372 while removing those that have a neutral or harmful effect on model quality. The key challenge 373 is that in standard black-box models, it is entirely unclear whether any given sample is helpful or 374 harmful. The gold standard is leave-one-out-training (for ImageNet, this would consist of training 375 1.28 million models); current methods seek to approximate this extremely costly approach with various heuristics (Feldman & Zhang, 2020; Chitta et al., 2021; Paul et al., 2021; Sorscher et al., 376 2022; Abbas et al., 2023a). By contrast, the contribution of a data sample to decision-making in 377 a visual memory based system is straightforward. For any given query image  $\tilde{x}$ , the neighbor set



Figure 4: Visualization of memory-based decision-making with and without memory pruning. Given a query image, nearest neighbors are retrieved from memory via Cosine similarity in the embedding space of a model (here: five closest neighbors from the ImageNet train set, embedded via DinoV2 ViT-L14). The model's prediction is based on the weighted aggregation of the neighbor class labels. The rank-based weight decreases with the rank of the neighbor. For soft memory pruning, those weights are adjusted by the reliability of their neighbors. In the specific example here, all five neighbors appear sensible, but they have four different labels. Since the first two neighbors contributed to wrong decisions on the training set, they are downweighted via soft memory pruning, and the prediction changes to the correct class.

398 399

390

391

392

393

394

395

396

397

400 Neighbors( $\tilde{x}$ ) clearly reveals which samples contributed to the decision. Furthermore, this informa-401 tion also highlights whether the samples were helpful (correct label) or harmful (wrong label) for 402 the decision. We, therefore, transfer the concept of dataset pruning to memory, and propose visual 403 memory pruning. To this end, we estimate sample quality by querying the ImageNet training set 404 against a visual memory consisting of the exact same dataset (IN-train, discarding the first neighbor 405 which is identical to the query). This approach requires no more compute than a single forward pass 406 over the training set. We then record the number of times any given neighbor contributed to a wrong decision, resulting in a sample quality estimate. This enables us to exclude low-quality neighbors 407 from the decision-making process by either removing them from the visual memory entirely ("hard 408 memory pruning") or by reducing their weight compared to higher-quality neighbors ("soft memory 409 pruning"). Method details can be found in Appendix H. In Table 3, we show that both memory 410 pruning variants improve ImageNet validation accuracy, with soft pruning leading to larger gains 411 than hard pruning. Figure 4 visualizes the decision-making process for a randomly selected sample 412 where estimating sample reliability improves decision quality. Given that observing the outcome of 413 an intervention is many orders of magnitude faster in visual memory models (as opposed to tradi-414 tional leave-out-training), we are optimistic that the visual memory pruning gains we observed with 415 two simple strategies can be improved further in the future. 416

Table 3: Flexible data selection: memory pruning. ImageNet validation accuracy improves when
removing low-quality samples (hard pruning) or downweighting them (soft pruning). In contrast to
standard black-box models, memory models (here: using DinoV2 ViT-L14) offer a strikingly simple
way to estimate sample quality since their decisions are based on a few retrieved memory samples.

Pruning	PluralityVoting	DistanceVoting	SoftmaxVoting	RankVoting
no pruning (standard)	83.2	83.3	83.6	83.6
hard pruning (ours)	83.3	83.4	83.6	83.7
soft pruning (ours)	83.6	83.6	83.9	84.1

429

### 3.6 FLEXIBLY INCREASING DATASET GRANULARITY

In contrast to static classification, where a model is trained once without updates, a visual memory
 model should be able to flexibly refine its visual understanding as more information becomes avail able. We test this using DinoV2 ViT-L14 embeddings on the iNaturalist21 dataset (iNaturalistTeam,

432 2021), a large-scale imbalanced dataset of animal and plant images containing 10,000 species span-433 ning seven taxonomic levels, from coarse (kingdom) to fine-grained (species). In a leave-one-out 434 fashion, we simulate the discovery of a new species by putting 50 exemplars for each of the 9,999 435 species into memory and then step by step adding more data for the remaining "newly discovered" 436 species-starting from zero exemplars all the way to 50 exemplars (see Algorithm 1 for an algorithmic description). In Section 3.6 we observe the following: (1.) Already before a single example of 437 the new species is added, it can already be placed in the right part of the taxonomic tree well beyond 438 chance (35.2% accuracy at the genus level compared to  $\sim 0\%$  chance). (2.) Accuracy at the species 439 level improves substantially by adding just a handful of images of the target species (e.g., 5–10 im-440 ages); a regime where training a classifier would typically fail due to data scarcity. (3.) Interestingly, 441 adding more samples of the discovered species not only improves species-level accuracy, but also 442 leads to a "rising tide lift" of improvements across all levels of the taxonomic hierarchy. This indi-443 cates that a visual memory is well-suited for hierarchical classification tasks and settings where data 444 for new concepts is initially scarce but becomes more abundant over time-which is often the case 445 in applications like fraud detection, personalized recommender systems, and scientific discovery.



Figure 5: Impact of memory bank size on top-1 accuracy across taxonomic levels on iNaturalist. Top-1 accuracy for a target species across different taxonomic levels as the number of exemplars in the memory bank for that species increases from 0 to 50. Each line represents the average accuracy over all 10,000 species in the iNaturalist 2021 dataset, while the number of examples in visual memory is fixed at 50 exemplars for all other species. The black dotted line indicates baseline accuracy from predicting the majority class.

### 3.7 INTERPRETABLE & ATTRIBUTABLE DECISION-MAKING

Unlike a black-box deep learning model, a visual memory offers a natural way to understand a model's specific predictions by attributing them to training data samples (e.g. Papernot & McDaniel, 2018). In Figure 6, we visualize misclassified validation set examples from the ImageNet-A dataset (Hendrycks et al., 2021) using a memory of the ImageNet-1K training set. These randomly selected samples illustrate that many seemingly strange errors (e.g., predicting a type of fence instead of a teddy bear, or a unicycle instead of a bow tie) do in fact appear sensible given the data, raising questions about label quality of ImageNet-A—in a similar vein as label issues identified for ImageNet (Beyer et al., 2020; Shankar et al., 2020; Yun et al., 2021)—rather than about model quality. This issue is quantified in Appendix M, showing that 2 out of 5 model "errors" are instead label errors.

### 4 DISCUSSION

461

462 463

464

465

466

467

468

469

470

471 472

473

474 **Summary.** Typical neural networks are trained end-to-end: perfect for static worlds, yet cumber-475 some to update whenever knowledge changes. This is limiting their potential in real-world settings 476 since the world is constantly evolving. Incorporating a visual memory, in contrast, enables a range of 477 flexible capabilities that embrace change: lifelong learning through incorporating novel knowledge, 478 being able to forget, remove and unlearn obsolete knowledge, flexible data selection through mem-479 ory pruning, and an interpretable decision-making paradigm on which one can intervene to control 480 its behavior. We systematically explored a simple visual memory that decomposes the task of image 481 classification into two primitives, image *similarity* (from a pre-trained embedding representation) 482 and *search* (via fast, scalable nearest neighbor search from a vector database). Our results demon-483 strate that technical improvements like RankVoting improve kNN accuracies for both DinoV2 and CLIP over the widely used SoftmaxVoting method that is sensitive to two hyperparameters (tem-484 perature  $\tau$  and number of neighbors k). Our approach also narrows the accuracy gap between a 485 nearest neighbor memory (best flexibility, perfect unlearning, improved interpretability) and a fixed



Figure 6: **Interpretable decision-making.** A retrieval-based visual memory enables a clear visual understanding of why a model makes a certain prediction. Here, we show four randomly selected misclassified query images from ImageNet-A (Hendrycks et al., 2021) along with five nearest neighbors from DinoV2 ViT-L14 using the ImageNet-1K training set as visual memory. All labels are from the respective datasets (ImageNet-A for query and ImageNet-train for neighbors). While all neighbors visually look reasonable, not all labels do.

517

linear probe (highest accuracy on static image classification). More importantly, we show that visual
 memory enables *flexible* perceptual capabilities.

520 **Limitations and future work.** First, we only considered the task of image classification across a broad range of datasets. It will be interesting to extend the approach to other visual tasks, such 521 as object detection, image segmentation, instance recognition and to image generation where a vi-522 sual memory would be desirable, too (since it is prohibitively expensive to re-train large generative 523 models every time data needs to be removed or added). Secondly, our approach relies on a fixed, 524 pre-trained embedding model; strong distribution shifts may require updating the embedding. Self-525 supervised models are a particularly flexible choice, but it is an open question whether one could 526 train smaller models that excel at their task with the help of a larger memory database. Conceptually, 527 if a model needs to save less information in its weights, it might be possible to reduce the computa-528 tional footprint of such a model. Furthermore, we sometimes observe a trade-off between flexibility 529 and accuracy. Exploring the use of the memory pruning weights as a data selection criterion in the 530 context of dataset pruning to improve over power-law scaling in deep learning (Sorscher et al., 2022) 531 might be an interesting avenue for future work.

Outlook. Deep learning is increasingly becoming a victim of its own success: the more widely it is deployed, the stronger its limitations are felt. While the static nature of end-to-end trained networks can easily be forgotten when focusing on fixed academic benchmarks, the real world is anything but static. Data is constantly evolving, leading to the dreaded "model drift" where once-optimal models gradually become less effective (Bayram et al., 2022). Incorporating an explicit visual memory—however it may be instantiated—appears to be a promising way forward for real-world tasks where flexibility is key. While the specific approach we employ here might well be improved through more complex systems, we hope that the flexible capabilities we demonstrated might inspire and contribute to a conversation on how knowledge ought to be represented in vision models.

540 **Code availability.** Code to replicate experiments from this paper is available via github; for the 541 purpose of the anonymous review period we include it as a supplementary .zip file. 542

### References

543

544

547

548

554

559

565

566

567

576

580

581

- Amro Kamal Mohamed Abbas, Evgenia Rusak, Kushal Tirumala, Wieland Brendel, Kamalika 546 Chaudhuri, and Ari S Morcos. Effective pruning of web-scale datasets based on complexity of concept clusters. In The Twelfth International Conference on Learning Representations, 2023a.
- Amro Kamal Mohamed Abbas, Kushal Tirumala, Daniel Simig, Surya Ganguli, and Ari S Morcos. 549 SemDeDup: Data-efficient learning at web-scale through semantic deduplication. In ICLR 2023 550 Workshop on Mathematical and Empirical Understanding of Foundation Models, 2023b. 551
- 552 David W Aha, Dennis Kibler, and Marc K Albert. Instance-based learning algorithms. Machine 553 learning, 6:37-66, 1991.
- M Saiful Bari, Batool Haider, and Saab Mansour. Nearest neighbour few-shot learning for cross-555 lingual classification. arXiv preprint arXiv:2109.02221, 2021. 556
- Firas Bayram, Bestoun S Ahmed, and Andreas Kassler. From concept drift to model degradation: An overview on performance-aware drift detectors. Knowledge-Based Systems, 245:108632, 2022.
- Liron Bergman, Niv Cohen, and Yedid Hoshen. Deep nearest neighbor anomaly detection. arXiv 560 preprint arXiv:2002.10445, 2020. 561
- 562 Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. Are 563 we done with imagenet? arXiv preprint arXiv:2006.07159, 2020.
  - Julian Bitterwolf, Maximilian Müller, and Matthias Hein. In or out? Fixing ImageNet out-ofdistribution detection evaluation. In International Conference on Machine Learning, pp. 2471– 2506. PMLR, 2023.
- 568 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101-mining discriminative compo-569 nents with random forests. In ECCV 2014, pp. 446-461. Springer, 2014. 570
- Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin 571 Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In 2021 IEEE 572 Symposium on Security and Privacy (SP), pp. 141–159. IEEE, 2021. 573
- 574 Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In 2015 575 IEEE symposium on security and privacy, pp. 463-480. IEEE, 2015.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and 577 Armand Joulin. Emerging properties in self-supervised vision transformers. In Proceedings of 578 the IEEE/CVF international conference on computer vision, pp. 9650–9660, 2021. 579
  - Yanbei Chen, Xiatian Zhu, and Shaogang Gong. Semi-supervised deep learning with memory. In Proceedings of the European conference on computer vision (ECCV), pp. 268–283, 2018.
- Kashyap Chitta, José M Álvarez, Elmar Haussmann, and Clément Farabet. Training data subset 583 search with ensemble active learning. IEEE Transactions on Intelligent Transportation Systems, 584 23(9):14741-14752, 2021. 585
- 586 Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, 587 Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, Rodolphe Jenatton, 588 Lucas Beyer, Michael Tschannen, Anurag Arnab, Xiao Wang, Carlos Riquelme, Matthias Minderer, Joan Puigcerver, Utku Evci, Manoj Kumar, Sjoerd van Steenkiste, Gamaleldin F. Elsayed, 590 Aravindh Mahendran, Fisher Yu, Avital Oliver, Fantine Huot, Jasmijn Bastings, Mark Patrick Collier, Alexey Gritsenko, Vighnesh Birodkar, Cristina Vasconcelos, Yi Tay, Thomas Mensink, Alexander Kolesnikov, Filip Pavetić, Dustin Tran, Thomas Kipf, Mario Lučić, Xiaohua Zhai, 592 Daniel Keysers, Jeremiah Harmsen, and Neil Houlsby. Scaling vision transformers to 22 billion parameters. In International Conference on Machine Learning, 2023.

594 595 596	Stephen Dopkins and Theresa Gleason. Comparing exemplar and prototype models of cate- gorization. <i>Canadian Journal of Experimental Psychology/Revue canadienne de psychologie</i> <i>expérimentale</i> , 51(3):212, 1997.
597 598 599 600	Andrew Drozdov, Shufan Wang, Razieh Rahimi, Andrew McCallum, Hamed Zamani, and Mohit Iyyer. You can't pick your neighbors, or can you? when and how to rely on retrieval in the <i>k</i> nn-lm. <i>arXiv preprint arXiv:2210.15859</i> , 2022.
601 602 603	Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. <i>Advances in Neural Information Processing Systems</i> , 33:2881–2891, 2020.
604 605 606	Jianping Gou, Taisong Xiong, Yin Kuang, et al. A novel weighted voting for k-nearest neighbor rule. <i>Journal of Computers</i> , 6(5):833–840, 2011.
607 608 609	Zhongrui Gui, Shuyang Sun, Runjia Li, Jianhao Yuan, Zhaochong An, Karsten Roth, Ameya Prabhu, and Philip Torr. kNN-CLIP: Retrieval enables training-free segmentation on continually expanding large vocabularies. <i>arXiv preprint arXiv:2404.09447</i> , 2024.
610 611 612 613	Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Ku- mar. Accelerating large-scale inference with anisotropic vector quantization. In <i>International</i> <i>Conference on Machine Learning</i> , pp. 3887–3896. PMLR, 2020.
614 615 616	Varun Gupta, Christopher Jung, Seth Neel, Aaron Roth, Saeed Sharifi-Malvajerdi, and Chris Waites. Adaptive machine unlearning. <i>Advances in Neural Information Processing Systems</i> , 34:16319– 16330, 2021.
617 618 619 620 621	Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. <i>arXiv</i> preprint arXiv:2006.16241, 2020.
622 623	Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In <i>CVPR</i> , pp. 15262–15271, 2021.
624 625 626 627	Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. arXiv preprint arXiv:1712.00409, 2017.
628 629 630	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. <i>arXiv preprint arXiv:2203.15556</i> , 2022.
631 632 633	iNaturalistTeam. iNaturalist 2021 competition dataset. https://github.com/visipedia/inat_comp/tree/master/2021, 2021.
634 635	Ahmet Iscen, Thomas Bird, Mathilde Caron, Alireza Fathi, and Cordelia Schmid. A memory trans- former network for incremental learning. <i>arXiv preprint arXiv:2210.04485</i> , 2022.
636 637 638	Ahmet Iscen, Mathilde Caron, Alireza Fathi, and Cordelia Schmid. Retrieval-enhanced contrastive vision-text models. <i>arXiv preprint arXiv:2306.07196</i> , 2023.
639 640 641	Frank Jäkel, Bernhard Schölkopf, and Felix A Wichmann. Generalization and similarity in exemplar models of categorization: Insights from machine learning. <i>Psychonomic Bulletin &amp; Review</i> , 15: 256–271, 2008.
642 643 644	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> , 2020.
646 647	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generaliza- tion through memorization: Nearest neighbor language models. In <i>International Conference on</i> <i>Learning Representations</i> , 2019.

648 649 650	John K Kruschke. ALCOVE: an exemplar-based connectionist model of category learning. In <i>Connectionist psychology</i> , pp. 107–138. Psychology Press, 2020.
651 652	Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded machine unlearning. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
653 654 655	Ritchie Lee, Justin Clarke, Adrian Agogino, and Dimitra Giannakopoulou. Improving trust in deep neural networks with nearest neighbors. In AIAA Scitech 2020 Forum, pp. 2098, 2020.
656 657	Ken Ziyu Liu. Machine unlearning in 2024, Apr 2024. URL https://ai.stanford.edu/ ~kzliu/blog/unlearning.
658 659 660	Jie Lu, Anjin Liu, Fan Dong, Feng Gu, Joao Gama, and Guangquan Zhang. Learning under concept drift: A review. <i>IEEE transactions on knowledge and data engineering</i> , 31(12):2346–2363, 2018.
661 662	Douglas L Medin and Marguerite M Schaffer. Context theory of classification learning. <i>Psycholog-</i> <i>ical review</i> , 85(3):207, 1978.
664 665 666	Kengo Nakata, Youyang Ng, Daisuke Miyashita, Asuka Maki, Yu-Chieh Lin, and Jun Deguchi. Revisiting a knn-based image classification system with high-capacity storage. In <i>European Con-</i> <i>ference on Computer Vision</i> , pp. 457–474. Springer, 2022.
667 668 669	Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. A survey of machine unlearning. <i>arXiv preprint arXiv:2209.02299</i> , 2022.
670 671 672	Robert M Nosofsky. Attention, similarity, and the identification–categorization relationship. <i>Journal of experimental psychology: General</i> , 115(1):39, 1986.
673 674 675	Robert M Nosofsky. The generalized context model: An exemplar model of classification. <i>Formal approaches in categorization</i> , pp. 18–39, 2011.
676 677 678 679	Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khali- dov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. DINOv2: Learning robust visual features without supervision. <i>Transactions on Machine Learning Research</i> , 2023.
680 681	Nicolas Papernot and Patrick McDaniel. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. <i>arXiv preprint arXiv:1803.04765</i> , 2018.
682 683 684	German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. <i>Neural networks</i> , 113:54–71, 2019.
685 686 687	Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Find- ing important examples early in training. <i>Advances in Neural Information Processing Systems</i> , 34:20596–20607, 2021.
688 689 690	Tobias Plötz and Stefan Roth. Neural nearest neighbors networks. Advances in Neural information processing systems, 31, 2018.
691 692 693	Ameya Prabhu, Zhipeng Cai, Puneet Dokania, Philip Torr, Vladlen Koltun, and Ozan Sener. Online continual learning without the storage constraint. <i>arXiv preprint arXiv:2305.09253</i> , 2023.
694 695	J Ross Quinlan. Combining instance-based and model-based learning. In <i>Proceedings of the tenth international conference on machine learning</i> , pp. 236–243, 1993.
696 697 698 699	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>ICML</i> , pp. 8748–8763, 2021.
700 701	Nazneen Fatema Rajani, Ben Krause, Wengpeng Yin, Tong Niu, Richard Socher, and Caiming Xiong. Explaining and improving model behavior with k nearest neighbor representations. <i>arXiv</i> preprint arXiv:2010.09030, 2020.

702 703 704 705	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean- baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem- ini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint</i> <i>arXiv:2403.05530.2024</i>
706	urav.2705.05550, 2024.
707	Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanieev Satheesh, Sean Ma, Zhiheng
708	Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual
709	recognition challenge. International Journal of Computer Vision, 115:211–252, 2015.
710	Arush Salihari Jawaday Asharya Cautan Karath and Ananda Thaartha Surash Damambar what
711	you want to forget: Algorithms for machine unlearning. Advances in Neural Information Pro-
712	cessing Systems, 34:18075–18086, 2021.
713	Nazanin Mohammadi Sepahyand Vincent Dumoulin Eleni Triantafillou and Gintare Karolina Dz-
715	iugaite. Data selection for transfer unlearning. arXiv preprint arXiv:2405.10425, 2024.
716	Vaishaal Shankar Rebecca Roelofs Horia Mania Alex Fang Renjamin Recht and Ludwig
717 718	Schmidt. Evaluating machine accuracy on ImageNet. In <i>International Conference on Machine Learning</i> , pp. 8634–8644. PMLR, 2020.
719	
720	Rulin Shao, Jacqueline He, Akari Asai, Weijia Shi, Tim Dettmers, Sewon Min, Luke Zettlemoyer,
721 722	and Pang Wei Koh. Scaling retrieval-based language models with a trillion-token datastore. <i>arXiv</i> preprint arXiv:2407.12854, 2024.
723	
724	Thalles Silva, Helio Pedrini, and Adín Ramírez Rivera. Learning from memory: Non-parametric
725	memory augmented self-supervised learning of visual features. In Forty-first International Con-
725	ference on Machine Learning, 2024.
720	Chawin Sitawarin and David Wagner. On the robustness of deep k pearest neighbors. In 2010 IEEE
728	Security and Privacy Workshops (SPW), pp. 1–7. IEEE, 2019.
729	Sivic and Zisserman Video Google: a text retrieval approach to object matching in videos. In
730	Proceedings of the IFFF International Conference on Computer Vision pp 1470–1477 IFFF
731	2003
732	2005.
733 734	Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neu- ral scaling laws: beating power law scaling via data pruning Advances in Neural Information
735	Processing Systems, 35:19523–19536, 2022.
736	Yiyou Sun Yifei Ming Xiaojin Zhu and Yixuan Li Out-of-distribution detection with deep nearest
737 738	neighbors. In International Conference on Machine Learning, pp. 20827–20840. PMLR, 2022.
739	Alexey Tsymbal The problem of concept drift: definitions and related work Computer Science
740	Department, Trinity College Dublin, 106(2):58, 2004.
741	Motthew A Turk and Alay D Dantland Ease recognition using signafaces. In Ducase June of the
742	IFEE Computer Society Conference on Computer Vision and Pattern Decognition pp. 596-597
743	IEEE Computer Society Conference on Computer vision and Fattern Recognition, pp. 560–567. IEEE Computer Society 1991
744	ille compater society, 1991.
745	Enayat Ullah, Tung Mai, Anup Rao, Ryan A Rossi, and Raman Arora. Machine unlearning via
746	algorithmic stability. In Conference on Learning Theory, pp. 4126–4142. PMLR, 2021.
740	Eric Wallace, Shi Feng, and Jordan Boyd-Graber. Interpreting neural networks with nearest neigh-
740	bors. EMNLP 2018, pp. 136, 2018.
750	Haphan Wang Songwei Ge, Zachary Linton, and Eric D Ving. Learning robust global represente
751	tions by penalizing local predictive power. In Advances in Neural Information Processing Sys-
752	tems, pp. 10506–10518, 2019a.
753	······, FF. 10000 10010, =0120
754	Yan Wang, Wei-Lun Chao, Kilian Q Weinberger, and Laurens Van Der Maaten. Simpleshot: Re-
755	visiting nearest-neighbor classification for few-shot learning. <i>arXiv preprint arXiv:1911.04623</i> , 2019b.

- Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. arXiv preprint arXiv:1410.3916, 2014.
- Yuhuai Wu, Markus Norman Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing trans formers. In *International Conference on Learning Representations*, 2021.
- Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3733–3742, 2018.
- Xi Yang, Xiaoting Nan, and Bin Song. D2n4: A discriminative deep nearest neighbor neural network
   for few-shot space target recognition. *IEEE Transactions on Geoscience and Remote Sensing*, 58
   (5):3667–3676, 2020.
- Sangdoo Yun, Seong Joon Oh, Byeongho Heo, Dongyoon Han, Junsuk Choe, and Sanghyuk Chun.
   Re-labeling ImageNet: from single to multi-labels, from global to localized labels. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2340–2350, 2021.
- Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12104–12113, 2022.
- Haibo Zhang, Toru Nakamura, Takamasa Isohara, and Kouichi Sakurai. A review on machine
   unlearning. *SN Computer Science*, 4(4):337, 2023.
- Jiaxin Zhuang, Jiabin Cai, Ruixuan Wang, Jianguo Zhang, and Wei-Shi Zheng. Deep kNN for medical image classification. In *Medical Image Computing and Computer Assisted Intervention– MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23*, pp. 127–136. Springer, 2020.

# Appendix Appendix

812

### 813 We here provide the following supplemental information: 814

- 815 Appendix A Aggregation method comparison on ImageNet-1K
- Appendix B Aggregation method comparison on iNaturalist
- 818 Appendix C Hyperparameter sensitivity analysis
- Appendix D Robustness towards label corruption
   820
- Appendix E Hit rate analysis as an upper bound on aggregation accuracy
- 822 Appendix F Scaling law details
- Appendix G OOD analysis for NINCO dataset
- 825 Appendix H Memory pruning details
- 826 Appendix I Linear probe details
- Appendix J Latency and storage
   828
- 829 Appendix K Algorithm for hierarchical label prediction
- 830 Appendix L Calibration analysis
- Appendix M ImageNet-A error analysis
- 833 Appendix N Compositionality analysis
- 834
- 835 836 837

838

## A AGGREGATION METHOD COMPARISON (IMAGENET-1K)

Table 4: Benchmarking different aggregation variants at different k thresholds, DinoV2 ViT-L14.

Aggregation	@10	@20	@30	@40	@50	@60	@70	@80	@90	@100
PluralityVoting	83.2	82.9	82.6	82.4	82.1	82.0	81.8	81.6	81.5	81.4
DistanceVoting	83.3	83.0	82.7	82.4	82.2	82.1	81.9	81.7	81.6	81.5
Softmax Voting	83.5	83.5	83.4	83.3	83.2	83.1	83.1	83.0	82.9	82.9
RankVoting	83.5	83.6	83.6	83.5	83.5	83.4	83.3	83.3	83.3	83.3

Table 5: Benchmarking different aggregation variants at different k thresholds, DinoV2 ViT-B14.

Aggregation	@10	@20	@30	@40	@50	@60	@70	@80	@90	@100
PluralityVoting	81.8	81.4	81.1	80.9	80.7	80.4	80.2	80.0	79.8	79.6
DistanceVoting	81.9	81.5	81.2	81.0	80.8	80.5	80.3	80.0	79.9	79.7
Softmax Voting	82.0	82.0	81.9	81.8	81.7	81.7	81.6	81.5	81.3	81.3
RankVoting	82.1	82.2	82.1	82.0	82.0	82.0	81.9	81.9	81.9	81.9

Table 6: Benchmarking different aggregation variants at different k thresholds, DinoV2 ViT-S14.

Aggregation	@10	@20	@30	@40	@50	@60	@70	@80	@90	@100
PluralityVoting	78.6	78.2	77.8	77.4	77.1	76.8	76.5	76.3	76.1	75.9
<b>DistanceVoting</b>	78.8	78.4	77.9	77.5	77.2	76.9	76.6	76.4	76.2	76.0
Softmax Voting	78.9	78.9	78.7	78.6	78.5	78.3	78.1	78.0	77.9	77.7
RankVoting	78.9	79.1	79.0	78.9	78.9	78.9	78.9	78.8	78.8	78.8
	Aggregation PluralityVoting DistanceVoting SoftmaxVoting RankVoting	Aggregation@10PluralityVoting78.6DistanceVoting78.8SoftmaxVoting <b>78.9</b> RankVoting <b>78.9</b>	Aggregation         @10         @20           PluralityVoting         78.6         78.2           DistanceVoting         78.8         78.4           SoftmaxVoting <b>78.9</b> 78.9           RankVoting <b>78.9 79.1</b>	Aggregation@10@20@30PluralityVoting78.678.277.8DistanceVoting78.878.477.9SoftmaxVoting <b>78.9</b> 78.7RankVoting <b>78.979.1</b>	Aggregation@10@20@30@40PluralityVoting78.678.277.877.4DistanceVoting78.878.477.977.5SoftmaxVoting <b>78.9</b> 78.978.778.6RankVoting <b>78.979.179.078.9</b>	Aggregation@10@20@30@40@50PluralityVoting78.678.277.877.477.1DistanceVoting78.878.477.977.577.2SoftmaxVoting78.978.978.778.678.5RankVoting78.979.179.078.978.9	Aggregation@10@20@30@40@50@60PluralityVoting78.678.277.877.477.176.8DistanceVoting78.878.477.977.577.276.9SoftmaxVoting <b>78.9</b> 78.978.778.678.578.3RankVoting <b>78.979.179.078.978.978.9</b>	Aggregation@10@20@30@40@50@60@70PluralityVoting78.678.277.877.477.176.876.5DistanceVoting78.878.477.977.577.276.976.6SoftmaxVoting <b>78.9</b> 78.978.778.678.578.378.1RankVoting <b>78.979.179.078.978.978.978.9</b>	Aggregation@10@20@30@40@50@60@70@80PluralityVoting78.678.277.877.477.176.876.576.3DistanceVoting78.878.477.977.577.276.976.676.4SoftmaxVoting <b>78.9</b> 78.978.778.678.578.378.178.0RankVoting <b>78.979.179.078.978.978.978.978.9</b>	Aggregation@10@20@30@40@50@60@70@80@90PluralityVoting78.678.277.877.477.176.876.576.376.1DistanceVoting78.878.477.977.577.276.976.676.476.2SoftmaxVoting <b>78.9</b> 78.978.778.678.578.378.178.077.9RankVoting <b>78.979.179.078.978.978.978.978.878.8</b>

Table 7: Benchmarking different aggregation variants at different k thresholds, CLIP ViT-L14.

@80 @90 @100
77.4 77.2 77.0
77.5 77.3 77.1
78.5 78.4 78.2
79.7 79.7 79.7

Table 8: Benchmarking different aggregation variants at different k thresholds, CLIP ViT-B16.

Aggregation	@10	@20	@30	@40	@50	@60	@70	@80	@90	@100
PluralityVoting	72.8	72.6	72.3	72.0	71.7	71.4	71.2	70.9	70.8	70.5
DistanceVoting	73.1	72.9	72.6	72.3	71.9	71.6	71.4	71.1	70.9	70.6
SoftmaxVoting	73.3	73.3	73.1	72.9	72.7	72.5	72.3	72.1	71.9	71.7
RankVoting	73.0	73.7	73.8	73.8	73.8	73.8	73.7	73.7	73.7	73.7



Figure 7: Aggregation method comparison on the ImageNet-1K validation set (same as Figure 2a but for other models).

Model	Aggegation	IN-val acc (%)
CLIP ViT-L14	CLIP paper (zero-shot)	75.3
CLIP ViT-L14	no aggregation	76.0
CLIP ViT-L14	PluralityVoting	79.2
CLIP ViT-L14	DistanceVoting	79.4
CLIP ViT-L14	SoftmaxVoting	79.6
CLIP ViT-L14	RankVoting	79.9
DinoV2 ViT-L14	DinoV2 paper (kNN Softmax)	83.5
DinoV2 ViT-L14	no aggregation	81.1
DinoV2 ViT-L14	PluralityVoting	83.2
DinoV2 ViT-L14	DistanceVoting	83.3
DinoV2 ViT-L14	SoftmaxVoting	83.6
DinoV2 ViT-L14	RankVoting	83.6

Table 9: Benchmarking different aggregation variants on ImageNet-1K.





Figure 8: Aggregating information across retrieved memory samples on iNaturalist. Same as Figure 2 but for iNaturalist instead of ImageNet. (left) Existing aggregation methods (PluralityVoting, DistanceVoting and SoftmaxVoting) are overconfident in distant neighbors, resulting in the paradox of decaying iNaturalist accuracy with more information. (right) This is not the case for RankVoting which shows strong and stable performance across models and choices of k.

## C HYPERPARAMETER SENSITIVITY ANALYSIS



Figure 9: Sensitivity to hyperparameters for different aggregation methods. Apart from PluralityVoting, all aggregation methods described in Section 2.2 have a hyperparameter ( $\alpha$  for RankVoting,  $\tau$  for SoftmaxVoting). For each model and method, we here plot the maximum performance when aggregating using a certain method, sweeping over the number of neighbors from 1 to 100, as a function of the hyperparameter. This analysis is performed to understand how sensitive the respective method is to the choice of the hyperparameter. Note that the x range is different since for instance the temperature parameter in SoftmaxVoting ranges from [0, 1] while RankVoting for  $\alpha = 0$  is undefined (division by zero). We therefore evaluate a broad range for each method and find that all methods have a regime in which they are relatively stable irrespective of the hyperparameter choice. Since DistanceVoting as implemented by Khandelwal et al. (2019) does not have a hyperparameter, we added a temperature-style parameter  $\xi$  for the purpose of this comparison by setting  $w_i = \exp(-\text{dist}(\tilde{z}, z_{[i]}))^{\xi}$ .

### D ROBUSTNESS TOWARDS LABEL CORRUPTION



Figure 10: Robustness towards label corruption. How robust is a visual memory towards corrupted labels in the memory bank? This plot shows top-1 RankVoting accuracy on the ImageNet validation set as a function of how many labels in the memory (containing ImageNet-1K training set features via DinoV2 ViT-L/14) are corrupted, i.e., assigned to a random class. Intriguingly, performance stays almost unchanged all the way to about 60% (!) corrupted (random) labels in the database.

# 1026 E HIT RATE ANALYSIS



Figure 11: **Hit rate.** This plot shows the probability of the true label being contained in list of labels of the first k retrieved neighbors on ImageNet-1K, for five different models and  $k \in [1, 100]$ . With 100 neighbors, the hit rate approaches 98% for the best model. Conceptually, this is a very high upper bound on the performance that can be achieved by a given featurizer via nearest neighbor retrieval.

### 1049 F SCALING LAW

As we mentioned in Section 3.3, we found that a logarithmic form fits the data well between log<sub>10</sub>(memory size) and log<sub>10</sub>(error rate). Specifically, we found the following functional forms for DinoV2 ViT S14 and DinoV2 ViT L14 respectively via np.polyfit (x, y, dim=1):

**DinoV2 ViT L14:** 
$$y = -0.9434 \cdot \log_{10}(x) + 2.0704$$
  
**DinoV2 ViT S14:**  $y = -1.0942 \cdot \log_{10}(x) + 2.3187$ 

where  $x = \log_{10}$  (memory-size) and  $y = \log_{10}$  (error-rate), where memory-size  $\in [10^3, 10^9]$  and error-rate in [0, 100].

### G NINCO DATASET



<sup>1073</sup> 

1048

1050

1054 1055

1056

1059 1060

Figure 12: Distance comparison: the NINCO OOD samples are indeed out-of-distribution for the model. In Section 3.1, we described that we can simply plug new out-of-distribution classes into memory and still perform well on both existing data as well as the new classes. This boxplot confirms that the added samples from the NINCO dataset (Bitterwolf et al., 2023) are indeed out-ofdistribution for DinoV2 ViT-L14: The mean (left) and median (right) distances from query to the first 100 neighbors are substantially lower for ImageNet validation images than for OOD samples from NINCO.

Figure 12 confirms that there is a distribution difference between in-distribution data (ImageNet-1081 1K) and OOD data (NINCO). That said, while a distribution shift exists, it is possible that individual NINCO samples were part of the training set for DinoV2. Test-set contamination is generally a concern when working with models trained on large-scale datasets, since test samples may occur as exact, semantic or near-duplicates in large training datasets (e.g. Abbas et al., 2023b). For instance, NINCO contains samples from Food-101 (Bossard et al., 2014) which are also part of LVD-142M dataset used to train DinoV2. That said, the NINCO samples belong to classes which are definitely not part of the ImageNet-train set which serves as a memory bank for our experiments, as ensured by the NINCO dataset collection process (Bitterwolf et al., 2023).

1089

# 1090 H MEMORY PRUNING

1091

For memory pruning from Section 3.5, we implemented two pruning methods: removing unreliable 1093 neighbors from memory entirely ("hard memory pruning"), and reducing their weight ("soft memory 1094 pruning"). We report results on the ImageNet validation set with a (potentially pruned) ImageNettrain set in memory. For hard pruning, we excluded images from memory that contributed to a wrong 1095 decision at least 128 times (this meant excluding 26,257 images for DinoV2 ViT-L14), based on 1096 querying the ImageNet-train set against a memory consisting of the ImageNet-train set and querying 1097 100 neighbors for each sample. In order to obtain a fair comparison, instead of reporting accuracies 1098 for an arbitrary choice of k (the number of neighbors) we instead evaluate accuracy for each k in 1099 [1,100] and report the maximum accuracy obtained in Table 3. This ensures that differences in 1100 observed accuracy can indeed be attributed to memory pruning, as opposed to a choice of k. For soft 1101 pruning, instead of excluding unreliable neighbors entirely as in hard pruning, the neighbor weights 1102 (1.0 for PluralityVoting, or a rank-based weight in case of RankVoting) are instead multiplied by a 1103 reliability factor  $\gamma$  with  $\gamma = \frac{d}{c+v}$  where v is the number of times the image contributed to a wrong 1104 decision on the ImageNet-train set, c = 1 to avoid division by zero, and d = 1.75. This results, 1105 for instance, in  $\gamma = 0.88$  for images that only contribute to a single wrong decision; in  $\gamma = 0.16$ 1106 for images that contribute to ten wrong decisions, and in  $\gamma = 0.02$  for images that contribute to 100 wrong decisions on the training set. Images that never contributed to any wrong decision are 1107 assigned  $\gamma = 1.0$ , i.e. their default weight remains unchanged. 1108

1109

# 1110 I LINEAR PROBE DETAILS

1112 For the linear probe results reported in the paper, we directly used the results that were reported 1113 in the DinoV2 and CLIP papers. For DinoV2, the authors froze the model backbone and trained 1114 the linear layers for 12500 iterations using SGD. Instead of training a single time, they performed 1115 a full grid search sweep over three settings (output layers in 1, 4; pooling token concatenation in 1116 yes, no, and 13 different learning rates), resulting in 52 linear probes. Then, the authors evaluated 1117 the ImageNet validation accuracy for all of those 52 probes and only reported the highest one, as described in Appendix B.3 of the DinoV2 paper. Some may call this test set tuning or double 1118 dipping; the DinoV2 paper describes it as "common practice" (Oquab et al., 2023, p. 31). CLIP 1119 linear probe results are based on a logistic regression classifier learned using scikit-learn's L-BFGS 1120 implementation, and hyperparameter sweeps are performed on a held-out set not used for evaluation, 1121 according to Radford et al. (2021). 1122

1123

1125

# 1124 J LATENCY AND STORAGE

**Latency.** Nearest neighbor retrieval, fortunately, does not need to reinvent the wheel but can, instead, build on top of highly optimized workloads and libraries such as the ScaNN library (Guo et al., 2020). The ScaNN github README shows a latency comparison; with the requirement of perfect recall a million-size memory can handle roughly 500-600 queries per second.

1130

Storage. In addition to latency, storage is another very practical consideration: How much does it take to store features for a large database? To put things into perspective, the ImageNet training dataset requires 154.6 GB of storage, and the ImageNet validation dataset requires 6.0 GB of storage. In comparison, as shown in Table 10, storing DinoV2 or CLIP features for the entire ImageNet

training dataset only requires between 1.9 and 4.9 GB of storage space. Thus compared to storing the training dataset, the model features account for only 1–3% of this size. This means that after constructing the memory, one may decide to keep the dataset which adds 1–3% of storage, or one may decide to delete the dataset only keeping the features which saves 97–99% of storage (compared to the dataset storage requirement). The ratio of features requiring 1–3% of the dataset size doesn't change with dataset scale since it only depends on the embedding model, thus this ratio would hold for very small datasets just as it would for a billion-scale dataset.

1145Table 10: Storage requirements for ImageNet features. Storing features in a memory database1146requires only about 1–3% of the space that is needed to store the dataset (154.6 GB for ImageNet-1147train, 6.0 GB for ImageNet-validation).

	Model	IN-train features (GB)	IN-val features (MB)
		1.0	107
	$\frac{\text{Dinov 2 V11-L/14}}{\text{Dinov 2 V11-L/14}}$	4.9	149/
	$\frac{\text{Dino}\text{V2}\text{V11-B}/\text{14}}{\text{Din}\text{V2}\text{V1T}\text{S}/\text{14}}$	3.7	148
	$\frac{D10V2}{V11-S/14}$	1.9	/5
	CLIP V11-L/14	3.7	148
	CLIP Vi1-B/16	2.5	100
K !	HIERARCHICAL LABEL	PREDICTION ALGOR	ITHM FOR FLEXIBLY
ſ	INCREASING DATASET C	RANULARITY	
Algor	•ithm 1 Hierarchical Label Pro	ediction	
Algor	<b>ithm 1</b> Hierarchical Label Pro	ediction	anda)
Algor Requ	<b>ithm 1</b> Hierarchical Label Pro <b>ire:</b> New example <i>x</i> , Hierarch	ediction hical tree $T$ (with ROOT :	node)
Algor Requ	<b>Tithm 1</b> Hierarchical Label Pro <b>ire:</b> New example $x$ , Hierarch $ur\_node \leftarrow \text{ROOT}$	ediction hical tree $T$ (with ROOT to the second s	node)
Algor Requ 1: ca 2: fo	<b>Tithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarch $ur\_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to bot $aggadidates \downarrow = all abildren$	ediction hical tree $T$ (with ROOT to tom) <b>do</b>	node)
Algor Requ 1: ca 2: fo 3:	<b>Tithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarch $ur\_node \leftarrow \text{ROOT}$ <b>or</b> each level $l$ in $T$ (top to bot $candidates \leftarrow \text{all\_children}$ mage a male a level d	ediction hical tree T (with ROOT in tom) <b>do</b> _of(cur_node)	node)
Algor Requ 1: cri 2: fc 3: 4:	<b>Tithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarch $ur\_node \leftarrow \text{ROOT}$ <b>or</b> each level $l$ in $T$ (top to both $candidates \leftarrow \text{all\_children}$ $max\_p\_value \leftarrow -\infty$ level at level t = 0	ediction hical tree T (with ROOT : tom) <b>do</b> _of(cur_node)	node)
Algor Requ 1: co 2: fo 3: 4: 5: 6	<b>Tithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarch $ur\_node \leftarrow \text{ROOT}$ <b>or</b> each level $l$ in $T$ (top to both $candidates \leftarrow \text{all\_children}$ $max\_p\_value \leftarrow -\infty$ $label\_at\_level \leftarrow \text{NULL}$ <b>for</b> each shild and a simple	ediction hical tree T (with ROOT to tom) <b>do</b> _of(cur_node)	node)
Algor Requ 1: ca 2: fc 3: 4: 5: 6: 7:	<b>Tithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarch $ur_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to bot $candidates \leftarrow \text{all_children}$ $max\_p\_value \leftarrow -\infty$ $label\_at\_level \leftarrow \text{NULL}$ for each child node $c$ in $can$	ediction hical tree T (with ROOT is tom) <b>do</b> _of(cur_node) ndidates <b>do</b> distribution (second second	node)
Algor           Requ           1: ca           2: fc           3:           4:           5:           6:           7:           8:	<b>Fithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarch $ur\_node \leftarrow \text{ROOT}$ <b>or</b> each level $l$ in $T$ (top to both $candidates \leftarrow \text{all\_children}$ $max\_p\_value \leftarrow -\infty$ $label\_at\_level \leftarrow \text{NULL}$ <b>for</b> each child node $c$ in $candle candle candl$	ediction hical tree $T$ (with ROOT is tom) <b>do</b> _of(cur_node) ndidates <b>do</b> _distribution( $x$ , examples(	c))
Algor Requ 1: ca 2: fc 3: 4: 5: 6: 7: 8: 0:	<b>Fithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarch $ur_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to bot $candidates \leftarrow \text{all_children}$ $max\_p\_value \leftarrow -\infty$ $label\_at\_level \leftarrow \text{NULL}$ <b>for</b> each child node $c$ in $can$ $cross\_dist \leftarrow \text{distance}$ $in\_dist \leftarrow \text{distance}$ $maku_a \leftarrow kalara and $	ediction hical tree $T$ (with ROOT is tom) <b>do</b> _of(cur_node) ndidates <b>do</b> distribution( $x$ , examples( $c$ ), examples( $c$ ), examples( $c$ ), examples( $c$ ), examples( $d$ )	c)) amples(c))
Algor           Requ           1: ca           2: fa           3:           4:           5:           6:           7:           8:           9:           10:	<b>Fithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarchical Label Pro- $ur_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to both $candidates \leftarrow \text{all_children}$ $max_p_value \leftarrow -\infty$ $label_at\_level \leftarrow \text{NULL}$ for each child node $c$ in $cancellabel_at\_level \leftarrow \text{MULL}$ for each child node $c$ in $cancellabel_at \leftarrow \text{distancellabel}$ $in\_dist \leftarrow distancellabellabellabellabellabellabellabella$	ediction hical tree $T$ (with ROOT is tom) <b>do</b> _of(cur_node) ndidates <b>do</b> _distribution( $x$ , examples( $c$ ), existing v-Smirnov test(cross_distribution( $cx$ ), existence of the product of the pr	c)) amples(c)) t, in_dist)
Algor           Requ           1: ca           2: fc           3:           4:           5:           6:           7:           8:           9:           10:	<b>Fithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarchical Label Pro- $ur_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to both $candidates \leftarrow \text{all_children}$ $max_p_value \leftarrow -\infty$ $label_at\_level \leftarrow \text{NULL}$ for each child node $c$ in $cancellar$ $cross\_dist \leftarrow \text{distance}\_$ $in\_dist \leftarrow \text{distance}\_$ $p\_value \leftarrow \text{Kolmogorow}$ <b>if</b> $p\_value > max\_p\_value$	ediction hical tree $T$ (with ROOT is tom) <b>do</b> _of(cur_node) ndidates <b>do</b> _distribution( $x$ , examples( $c$ ), existing tribution(examples( $c$ ), existing v-Smirnov test(cross_distribution) number of the product of the produ	c)) amples(c)) t, in_dist)
Algon           Requ           1: c <sup>2</sup> 2: fd           3:           4:           5:           6:           7:           8:           9:           10:           11:	<b>Fithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarchical Label Pro- $ur_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to both $candidates \leftarrow \text{all_children}$ $max_p_value \leftarrow -\infty$ $label_at_level \leftarrow \text{NULL}$ <b>for</b> each child node $c$ in $cancellet$ $cross_dist \leftarrow \text{distance}$ $in_dist \leftarrow \text{distance}$ $p_value \leftarrow \text{Kolmogorow}$ <b>if</b> $p_value > max_p_value \leftarrow p_value \leftarrow$	ediction hical tree T (with ROOT is tom) do _of(cur_node) ndidates do distribution(x, examples(c), exist v-Smirnov test(cross_dist nlue then value	c)) amples(c)) t, in_dist)
Algor Requ 1: <i>c</i> <sup>2</sup> 2: <b>f</b> c 3: 4: 5: 6: 7: 8: 9: 10: 11: 12:	<b>Fithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarchical Label Pro- $ur_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to both $candidates \leftarrow \text{all_children}$ $max_p_value \leftarrow -\infty$ $label_at_level \leftarrow \text{NULL}$ <b>for</b> each child node $c$ in $cancellevel \leftarrow \text{NULL}$ <b>for</b> each child node $c$ in $cancellevel \leftarrow \text{distance_list}$ $p_value \leftarrow \text{Kolmogorow}$ <b>if</b> $p_value > max_p_value \leftarrow p_levelue \leftarrow p_label_at_level \leftarrow c$	ediction hical tree T (with ROOT is tom) do _of(cur_node) ndidates do distribution(x, examples(c), exist v-Smirnov test(cross_dist nlue then value	c)) amples(c)) t, in_dist)
Algor Requ 1: <i>c</i> <sup>2</sup> 2: <b>f</b> c 3: 4: 5: 6: 7: 8: 9: 10: 11: 12: 13:	<b>Fithm 1</b> Hierarchical Label Pro- <b>ire:</b> New example $x$ , Hierarchical Label Pro- $ur_node \leftarrow \text{ROOT}$ or each level $l$ in $T$ (top to both $candidates \leftarrow \text{all_children}$ $max_p_value \leftarrow -\infty$ $label_at_level \leftarrow \text{NULL}$ <b>for</b> each child node $c$ in $cancellevel \leftarrow \text{NULL}$ <b>for</b> each child node $c$ in $cancellevel \leftarrow \text{distance_list}$ $p_value \leftarrow \text{Kolmogorow}$ <b>if</b> $p_value \leftarrow \text{Kolmogorow}$ <b>if</b> $p_value \leftarrow \text{max_p_value}$ $max_p_value \leftarrow p_label_at_level \leftarrow c$ <b>end if</b>	ediction hical tree T (with ROOT is tom) <b>do</b> _of(cur_node) ndidates <b>do</b> .distribution(x, examples(c), exist v-Smirnov test(cross_dist nulue <b>then</b> value	c)) amples(c)) t, in_dist)
Algon           Requ           1: c <sup>2</sup> 2: fd           3:           4:           5:           6:           7:           8:           9:           10:           11:           12:           13:           14:	<b>ithm 1</b> Hierarchical Label Pre- <b>ire:</b> New example x, Hierarch $ur_node \leftarrow ROOT$ or each level l in T (top to bot $candidates \leftarrow$ all_children $max_p_value \leftarrow -\infty$ $label_at_level \leftarrow NULL$ for each child node c in car $cross_dist \leftarrow$ distance_lin_dist $\leftarrow$ distance_lin_dist $\leftarrow$ distance_list $p_value \leftarrow Kolmogoror$ if $p_value > max_p_value \leftarrow p_label_at_level \leftarrow c$ end ifend for	ediction hical tree T (with ROOT is tom) <b>do</b> _of(cur_node) ididates <b>do</b> distribution(x, examples(c), exist v-Smirnov test(cross_dist ilue <b>then</b> value	node) c)) amples(c)) t, in_dist)
Algon           Requ           1: c <sup>2</sup> 2: fd           3:           4:           5:           6:           7:           8:           9:           10:           11:           12:           13:           14:           15:	<b>:ithm 1</b> Hierarchical Label Pre- <b>ire:</b> New example x, Hierarch <i>ur_node</i> $\leftarrow$ ROOTor each level l in T (top to botcandidates $\leftarrow$ all_children $max_p_value \leftarrow -\infty$ label_at_level $\leftarrow$ NULLfor each child node c in car $cross_dist \leftarrow$ distance_list $p_value \leftarrow$ Kolmogororif $p_value \leftarrow$ Kolmogororif $p_value \leftarrow$ Kolmogororif $p_value \leftarrow$ max_p_value $\leftarrow p$ $label_at_level \leftarrow c$ end ifend for $cur_node \leftarrow label_at_level$	ediction hical tree T (with ROOT is tom) <b>do</b> _of(cur_node) adidates <b>do</b> .distribution(x, examples(c),	c)) amples(c)) t, in_dist)



Figure 13: How well are predictions calibrated? Left column: Accuracy vs. confidence from
Softmax of linear classifier for three DinoV2 variants. Right column: Accuracy vs. count of plurality
class among first 100 neighbors for the same three DinoV2 variants. A DinoV2-based kNN classifier
is well calibrated, as is the DinoV2 softmax.

# 1238 M IMAGENET-A ERROR ANALYSIS

1239

As shown in Figure 6, many "errors" on ImageNet-A appear to be perfectly reasonable predictions
 that are caused by dataset label issues as opposed to model mistakes. More randomly selected
 ImageNet-A samples, along with nearest neighbors, are shown in Figure 14. To quantify the issue,

we performed a human experiment on a randomly selected subset of ImageNet-A images (N=100) where the dataset label and the prediction from DinoV2 ViT-L14 with JFT memory disagree. We presented the image alongside the original ImageNet-A label and our model-predicted label to three human observers, asking them to identify which of the labels best describes the image (of course, without telling them which of the labels is the dataset label). The result was that in 39.3% (!) of cases (std:  $\pm 1.25\%$ ), the DinoV2 label was assessed as being better/more suitable than the original dataset label—i.e., roughly 2 out of 5 model "errors" are in fact dataset label errors, quantifying the ImageNet-A label quality issue we alluded to in Figure 6. This percentage can be used to estimate how correcting problematic labels influences performance. Instead of the original model's 61.1%accuracy on ImageNet-A, due to label errors the 'corrected' accuracy is instead 76.4% (a delta of +15.3% in absolute terms or +25.0% in relative terms). 



Figure 14: Interpretable decision-making. A retrieval-based visual memory enables a clear visual understanding of why a model makes a certain prediction. Here, we show four randomly selected misclassified query images from ImageNet-A (Hendrycks et al., 2021) along with five nearest neigh-bors from DinoV2 ViT-L14 using the ImageNet-1K training set as visual memory. All labels are from the respective datasets (ImageNet-A for query and ImageNet-train for neighbors). While all neighbors visually look reasonable, not all labels do.

#### Ν **COMPOSITIONALITY ANALYSIS**

A flexible visual memory also provides a path to analyze representations of various models, particu-larly, how different models represent multiple concepts in an image. We study this for an ImageNettrain visual memory of DinoV2 ViT-L14 and CLIP ViT-L14. We use manually selected query images from outside the ImageNet dataset that have multiple objects from the ImageNet labels. We query the visual memory for nearest neighbors of the query image. Subsequently, we obtain the resid-*ual image* by subtracting the features of the nearest neighbor from the features of the query image. We, then, obtain the nearest neighbors for the residual image from the visual memory. We plot the results in Figure 15 which shows that DinoV2 ViT-L14 and CLIP ViT-L14 represent concepts in their features in a different manner. The nearest neighbors for DinoV2 are mostly images with a single concept (or object) from the query image. The residual image, subsequently, leads to nearest neighbors dominated by another single object in the query image. In contrast, CLIP often results in nearest neighbors that are generally a blend of concepts from the query image. These qualitative explorations are simple demonstrations of the advantages of an interpretable decision-making process provided by a flexible visual memory.



Figure 15: **Compositionality of representations.** The first column indicates a query image; the next three columns are the three nearest neighbors from the training set. The last three columns are the *residual* images, obtained by subtracting the features of the nearest neighbor (2nd column from the left) from the features of the query image (1st column from the left). The nearest neighbors for DinoV2 are mostly images with a single concept (or object) from the query image. The residual image, subsequently, leads to nearest neighbors dominated by another single object in the query image. In contrast, CLIP often finds neighbors that are a blend of concepts from the query image.