CAPTURED BY CAPTIONS: ON MEMORIZATION AND ITS MITIGATION IN CLIP MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multi-modal models, such as CLIP, have demonstrated strong performance in aligning visual and textual representations, excelling in tasks like image retrieval and zero-shot classification. Despite this success, the mechanisms by which these models utilize training data, particularly the role of memorization, remain unclear. In uni-modal models, both supervised and self-supervised, memorization has been shown to be essential for generalization. However, it is not well understood how these findings would apply to CLIP, which incorporates elements from both supervised learning via captions that provide a supervisory signal similar to labels, and from self-supervised learning via the contrastive objective. To bridge this gap in understanding, we propose a formal definition of memorization in CLIP (CLIPMem) and use it to quantify memorization in CLIP-style models. Our results indicate that CLIP's memorization behavior falls between the supervised and selfsupervised paradigms, with "mis-captioned" samples exhibiting highest levels of memorization. Additionally, we find that the text encoder contributes more to memorization than the image encoder, suggesting that mitigation strategies should focus on the text domain. Building on these insights, we propose multiple strategies to reduce memorization while at the same time improving utility—something that had not been shown before for traditional learning paradigms where reducing memorization typically results in utility decrease.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

Multi-modal models, such as CLIP (Radford et al., 2021), have demonstrated strong performance in representation learning. By aligning visual and textual representations, these models achieve state-of-the-art results in tasks like image retrieval (Baldrati et al., 2022a;b), visual question answering (Pan et al., 2023; Song et al., 2022), and zero-shot classification (Radford et al., 2021; Ali & Khan, 2023; Wang et al., 2023; Zhang et al., 2022). Despite these successes, the mechanisms by which multi-modal models leverage their training data to achieve good generalization remain underexplored.

In uni-modal setups, both supervised (Feldman, 2020; Feldman & Zhang, 2020) and self-supervised (Wang et al., 2024b), models have shown that their ability to *memorize* their training data is essential for generalization. In supervised learning, memorization typically occurs for mislabeled samples, outliers (Bartlett et al., 2020; Feldman, 2020; Feldman & Zhang, 2020), or data points seen towards the end of training (Jagielski et al., 2022), while in self-supervised learning, high memorization is experienced particularly for atypical data points (Wang et al., 2024b). However, it is unclear how these findings extend to models like CLIP which entail elements from both supervised learning (via captions as supervisory signals) and self-supervised learning (via contrastive loss functions).

Existing definitions of memorization offer limited applicability to CLIP and cannot fully address this
gap. The standard definition from supervised learning (Feldman, 2020) relies on one-dimensional
labels and the model's ability to produce confidence scores for these labels, whereas CLIP outputs
high-dimensional representations. While the SSLMem metric (Wang et al., 2024b), developed for selfsupervised vision models, could, in principle, be applied to CLIP's vision encoder outputs, it neglects
the text modality, which is a critical component of CLIP. Additionally, measuring memorization in
only one modality, or treating the modalities separately, risks diluting the signal and under-reporting
memorization. Our experimental results, as shown in Section 4.3, confirm this concern. Therefore,
new definitions of memorization tailored to CLIP's multi-modal nature are necessary.



Figure 1: **Examples of data with different levels of memorization.** Higher memorization scores indicate stronger memorization. We observe that atypical or distorted images, as well as those with incorrect or imprecise captions, experience higher memorization compared to standard samples and easy-to-label images with accurate captions. Results are obtained on OpenCLIP (Ilharco et al., 2021), with encoders based on the ViT-Base architecture trained on the COCO dataset.

- The only existing empirical work on quantifying memorization in CLIP models (Jayaraman et al., 2024) focuses on Déjà Vu memorization (Meehan et al., 2023), a specific type of memorization. The success of their method relies on the accuracy of the integrated object detection method and on the availability of an additional public dataset from the same distribution as CLIP's training data, limiting practical applicability. To overcome this limitation, we propose *CLIPMem* that measures memorization directly on CLIP's output representations. Specifically, it compares the alignment—*i.e.*, the similarity between representations—of a given image-text pair in a CLIP model trained with the pair, to the alignment in a CLIP model trained on the same data but without the pair.
 - In our empirical study of memorization in CLIP-like models using CLIPMem, we uncover several 081 key findings. First, examples with incorrect or imprecise captions ("mis-captioned" examples) exhibit the highest levels of memorization, followed by atypical examples, as illustrated in Figure 1. Second, 083 removing these samples from training yields significant improvements in CLIP's generalization 084 abilities. These findings are particularly noteworthy, given that state-of-the-art CLIP models are 085 usually trained on large, uncurated datasets sourced from the internet with no guarantees regarding the correctness of the text-image pairs. Our results highlight that this practice not only exposes 087 imprecise or incorrect data pairs to more memorization, often recognized as a cause for increased 088 privacy leakage (Carlini et al., 2019; 2021; 2022; Song et al., 2017; Liu et al., 2021), but that it also negatively affects model performance. Furthermore, by disentangling CLIP's two modalities, we are 090 able to dissect how memorization manifests within each. Surprisingly, we find that memorization 091 does not affect both modalities alike, with memorization occurring more in the text modality than in the vision modality. Building on these insights, we propose several strategies to reduce memorization 092 while simultaneously improving generalization—a result that has not been observed in traditional 093 supervised or self-supervised learning, where any reduction of memorization causes decreases in 094 performance. Finally, at a deeper level, our analysis of the model internals, following Wang et al. (2024a), shows that CLIP's memorization behavior sits between that of supervised and self-supervised 096 learning. Specifically, neurons in early layers are responsible for groups of data points (e.g., classes), similar to models trained using supervised learning, while neurons in later layers memorize individual 098 data points, as seen in self-supervised learning. 099
 - 100 In summary, we make the following contributions:

066

067

068

069

071

101

102

103

104

- We propose CLIPMem, a metric to measure memorization in multi-modal vision language models.
- Through extensive evaluation, we identify that "mis-captioned" and "atypical" data points experience the highest memorization, and that the text encoder is more responsible for memorization than the image encoder.
- Based on our insights, we propose and evaluate multiple strategies to mitigate memorization in CLIP. We show that in CLIP, contrary to traditional supervised and self-supervised learning, a reduction of memorization does not need to imply a decrease in performance.

2 BACKGROUND AND RELATED WORK

110 CLIP. Contrastive Language-Image Pretraining (CLIP) (Radford et al., 2021) trains an image encoder 111 f_{img} and a text encoder f_{txt} to map image-text pairs into a shared latent space. It trains these encoders 112 by maximizing similarity for matching image-text pairs while minimizing it for non-matching pairs, 113 using a contrastive loss function \mathcal{L} :

115 116

114

108

109

117

 $\mathcal{L} = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(\sin(f_{\text{img}}(x_i), f_{\text{txt}}(y_i))/\tau)}{\sum_{j=1}^{N} \exp(\sin(f_{\text{img}}(x_i), f_{\text{txt}}(y_j))/\tau)},$

118 where sim(\cdot, \cdot) is cosine similarity, τ is a temperature parameter, and N is the batch size. The popular Language augmented CLIP (LaCLIP) (Fan et al., 2023) extends this approach by introducing text 119 augmentions alongside the original image augmentations (crops) during training to reduce overfitting. 120 We study the impact of this practice on memorization and find it to be a suitable mitigation method. 121

122 Memorization. Memorization refers to a model's tendency to store detailed information from 123 training examples, rather than learning general patterns. (Zhang et al., 2016; Arpit et al., 2017; Chatterjee, 2018; Feldman, 2020). This can lead to privacy risks when sensitive data is memorized 124 (Carlini et al., 2019; 2021; 2022; Song et al., 2017). To date, memorization has been studied mainly 125 within single modalities. In supervised learning, models tend to memorize mislabeled (Feldman, 126 2020), difficult, or atypical examples (Arpit et al., 2017; Sadrtdinov et al., 2021), which can improve 127 generalization on long-tailed data (Feldman, 2020; Feldman & Zhang, 2020). Similarly, in self-128 supervised learning (SSL) in the vision domain (Wang et al., 2024b), atypical samples experience 129 higher memorization, and reducing memorization in SSL encoders leads to decreased performance in 130 downstream tasks. A similar connection between memorization and generalization has been observed 131 in the language domain (Antoniades et al., 2024; Tirumala et al., 2022). However, these papers 132 consider single-modality models. How those insights transfer to multi-modal models remains unclear. 133

Memorization in self-supervised learning. Our CLIPMem builds on the SSLMem metric introduced 134 by Wang et al. (2024b). This metric measures how much an SSL encoder memorizes a data point 135 x by comparing the alignment of representations from its augmented views. Let $f: \mathbb{R}^n \to \mathbb{R}^d$ be 136 an SSL encoder trained on an unlabeled dataset $S = \{x_i\}_{i=1}^m$ using an SSL algorithm \mathcal{A} . The data 137 augmentations are represented as $Aug(x) = \{a(x) | a \in Aug\}$, where a is a transformation function 138 applied to the data point x, mapping from $\mathbb{R}^n \to \mathbb{R}^n$. The encoder's output representation for a given 139 data point x is denoted as f(x). For a trained SSL encoder f, the alignment loss for a data point x is 140 defined as

142

$$\mathcal{L}_{\text{align}}(f, x) = \mathop{\mathbb{E}}_{x', x'' \sim \operatorname{Aug}(x)} [d\left(f(x'), f(x'')\right)], \tag{1}$$

where x', x'' are augmented views of x and $d(\cdot, \cdot)$ is a distance metric, typically the ℓ_2 distance. 143 144 SSLMem is then defined as

$$SSLMem(x) = \mathop{\mathbb{E}}_{g \sim \mathcal{A}(S \setminus x)} \mathcal{L}_{align}(g, x) - \mathop{\mathbb{E}}_{f \sim \mathcal{A}(S)} \mathcal{L}_{align}(f, x)$$
(2)

147 with f being an SSL encoder trained with data point x, and g, an encoder trained without x but 148 otherwise on the same dataset. While this framework measures memorization using alignment loss 149 for single-modality encoders, this approach is unsuitable to leverage the signal over both modalities 150 from multi-modal encoders like CLIP, as we also highlight empirically in Section 4.3. However, we 151 can build on the main concepts from SSLMem to define a new metric that can evaluate memorization 152 in CLIP, by considering both image and text representations, as we will detail in Section 3.

153 **Memorization in CLIP**. Even though CLIP is a widely used vision-language encoder, there has 154 been limited work on measuring memorization in CLIP. The only existing work (Jayaraman et al., 155 2024) applies the empirical Déjà Vu memorization framework from (Meehan et al., 2023) to CLIP. It 156 measures memorization by computing the overlap between unique objects in potentially memorized 157 images and their nearest neighbors-identified in the CLIP embedding space-from a public dataset. 158 However, the reliance on external public data from the same distribution, along with the required 159 accuracy of the object detection (which may not perform well for all samples, especially atypical ones (Kumar et al., 2023; Dhamija et al., 2020), limits the applicability of this approach. We further 160 expand on this in Appendix A.1. In contrast, our CLIPMem operates directly on CLIP's output 161 representations and returns a joint score over both modalities.

162 3 DEFINING MEMORIZATION OVER MULTI-MODAL ENCODERS

3.1 PROBLEM SETUP

166 Consider a single image-text pair (I, T) from a dataset S and two CLIP models: a model f and a 167 reference model g, trained on dataset S and $S' = S \setminus \{(I, T)\}$, respectively. We aim to quantify the 168 memorization of (I,T) in f, trained on this data point, by leveraging model g not trained on the data point but otherwise on the same data, in a leave-one-out style of defining memorization (Feldman, 169 2020). We denote the image encoder in CLIP as f_{img} : Image $\rightarrow \mathbb{R}^d$ and the text encoder as 170 f_{txt} : Text $\to \mathbb{R}^d$. For the image-text pair (I, T), we denote with $f_{\text{img}}(I)$ the output representation of 171 f's image encoder on image I and with $f_{txt}(T)$ the output representation of f's text encoder on text 172 T. To evaluate the *alignment* between the image and text representations, *i.e.*, to quantify how similar 173 the two representations are, we use cosine similarity $sim(f_{img}(I), f_{txt}(T))$, as defined in the original 174 CLIP paper (Radford et al., 2021). 175

176

164

165

177

3.2 ALIGNMENT WITH CONTRASTIVE OBJECTIVE

178 During training, the contrastive objective in CLIP maximizes the cosine similarity for correct image-179 text pairs while minimizing the cosine similarity for all the other N-1 incorrect pairs in any 180 given training mini-batch with N training samples. This means that for a given image I and text 181 T, the training objective pulls $f_{img}(I)$ and $f_{txt}(T)$ closer together in the latent space, while pushing 182 $f_{\text{img}}(I)$ away from the representations of all other N-1 unrelated texts, and $f_{\text{txt}}(T)$ away from 183 all other images. Hence, the intuition is that the quality of alignment in f, unlike in uni-modal self-supervised learning (Wang et al., 2024b), depends not only on the model's ability to create 185 well-aligned text and image representations for a given text-image pair, but also on its ability to create 186 distant representations for the N-1 other representations.

To formalize this intuition into a metric that quantifies the alignment of f on the image-text pair (I, T), we define $\widehat{T_{test}}$ as a set of N-1 randomly chosen testing samples that were not used in training f or g. Furthermore, when applicable, we denote random augmentations of the training data—e.g., text augmentations in versions like LaCLIP (Fan et al., 2023)—as $T' \sim \operatorname{Aug}(T)$ for texts and $I' \sim \operatorname{Aug}(I)$ for images. Then, we define the alignment score of f on (I, T) as

193

194

195

196 197

10

199 200

201 202

203

208

209

where high scores indicate a better alignment of f on (I, T). In case no text augmentations are applied, as in standard CLIP training, the first term is calculated only over T.

$$\begin{split} \mathcal{A}_{\mathrm{align}}(f,I,T) = & \mathbb{E}_{\substack{(I',T') \sim \mathrm{Aug}(I,T)}} \left[\mathrm{sim}(f_{\mathrm{img}}(I'), f_{\mathrm{txt}}(T')) \right] \\ & - \mathbb{E}_{\substack{(_,t) \in \widehat{T_{test}}}} \left[\mathrm{sim}(f_{\mathrm{img}}(I), f_{\mathrm{txt}}(t)) \right] - \mathbb{E}_{\substack{(i,_) \in \widehat{T_{test}}}} \left[\mathrm{sim}(f_{\mathrm{img}}(i), f_{\mathrm{txt}}(T)) \right], \end{split}$$

3.3 DEFINING MEMORIZATION IN CLIP

Given our definition of alignment scores, we can define our CLIPMem in a similar vein to the definition of memorization in supervised learning (Feldman, 2020), in the leave-one-out style. Given the image-text pair (I, T) from dataset S and two CLIP models, f and g, trained on dataset S and $S' = S \setminus \{(I, T)\}$, respectively, we define CLIPMem as

$$\text{CLIPMem}(I,T) = \mathcal{A}_{\text{align}}(f,I,T) - \mathcal{A}_{\text{align}}(g,I,T).$$
(4)

(3)

If a model f has a significantly higher alignment score than model g on (I, T), this means that fmemorizes this data point. Note that taking the difference between f and g is crucial to get a solid estimate of memorization. This is because without "context", a high or low alignment score of fdoes not express much information. The alignment of f can be high without memorizing (I, T), for example, if (I, T) is a simple (but not memorized) training example. In this case, the reference model g will also have a high score, such that the difference is again small. Thanks to this design of our CLIPMem, it will then correctly report low memorization.

216 4 EMPIRICAL EVALUATION

4.1 EXPERIMENTAL SETUP

218

219

220 **Models and training.** We build our experiments on OpenCLIP (Cherti et al., 2023), an open-221 source Python version of Open-CLIP (Ilharco et al., 2021). The standard architecture used for the 222 experiments builds on ViT-Base, but we also include experiments using ViT-Large. We train the model on the COCO dataset (Lin et al., 2014). Since COCO is much smaller than OpenCLIP's 224 standard training datasets, we reduce the training batch size to 128 and increase the epoch number 225 from 32 to 100 to achieve similar performance. All other settings strictly follow OpenCLIP. For training DINO, as an example of an SSL vision encoder, we follow the default setting of Caron et al. 226 (2021). The supervised model is trained as a multi-label classifier, also based on ViT-Base (with an 227 additional fully connection layer) based on the first-level annotation captions in the COCO dataset. A 228 full specification of our experimental setup is detailed in Appendix A.2. Additional experiments for 229 measuring memorization on the BLIP (Li et al., 2022) model are presented in Appendix A.6. 230

Datasets. We use COCO (Lin et al., 2014), CC3M (Sharma et al., 2018), and the YFCC100M (Thomee et al., 2016a) datasets to pre-train the OpenCLIP models. For CC3M, we randomly sample 75000 examples from 2.91M total. We evaluate the models based on linear probing accuracy on ImageNet (Deng et al., 2009) with an added classification layer trained on top of the output representations. We use YFCC100M to simulate an infinite data regime, *i.e.*, using a single training run where no data point is repeated whereas we train iteratively using CC3M and COCO.

Measuring memorization. We follow Wang et al. (2024b) to approximate our CLIPMem. Since 237 training a separate pair of models for every data point would be computationally intractable, we 238 measure memorization across multiple data points simultaneously. Therefore, we divide the training 239 set into four subsets: (1) S_S , data points that both model f and g were trained on, (2) S_C , data points 240 used only for training f, (3) S_I , data points used only for training g, and (4) S_E , external "test" data 241 points that none of the models was trained on. Note that $|S_C| = |S_I|$, such that f and g have the same 242 number of training data points in total. For our experiments, following a similar approach to Wang 243 et al. (2024b), we want to strike a balance when choosing the size of S_C . If the size is too large, then 244 f and g might differ too much and not yield a strong memorization signal, but if it is too small, we 245 would only have a memorization signal for too few data points. Concretely, for COCO and CC3M, 246 we set $|S_S| = 65000$ and $|S_C| = |S_I| = |S_E| = 5000$. Memorization is reported as an average over 247 all data points in S_C for model f, or per individual data point in S_C .

Generating captions and images. To generate additional captions for the training images, we use GPT-3.5-turbo. For each input image, we provide the representation produced by our trained OpenCLIP model and ask GPT to generate five new captions. Generated sample captions are presented in Figure 18. To generate additional images for the COCO dataset, we use Stable Diffusion v1.5 to generate five new images, one corresponding to each of the five per-image captions in the COCO dataset. Sample generated images are presented in Figure 17.

254 255

256

4.2 STUDYING MEMORIZATION USING CLIPMEM

257 We first analyze the general memorization in CLIP in order to identify which data points are 258 memorized. To do this, we quantify CLIPMem over the different training subsets. Our results are 259 presented in Figure 2a. In particular, we observe that CLIPMem for S_C , the data points only used 260 to train model f, is significantly higher than for S_S , the data points shared between the two models. Memorization for S_S is comparable to that for S_E , *i.e.*, the external data not seen during training, 261 indicating that f does not memorize these samples. The data in S_I causes negative CLIPMem scores, 262 indicating that this data is memorized by g, not by f. This is the expected behavior according to the 263 definition of our metric. In Appendix A.4, we additionally highlight that memorization increases 264 with model size, *i.e.*, CLIP based on ViT-Large has a higher overall memorization with an average of 265 0.457 while CLIP based on ViT-Base only reaches 0.438 on average. 266

Additionally, we analyze individual data points by their reported CLIPMem. We give examples of
 highly memorized data points in Figure 1 and highly vs. little memorized samples in Figures 13,14,15,
 and 16 in Appendix A.9. Overall, high CLIPMem samples seem to be difficult examples or examples
 with imprecise or incorrect captions, whereas low CLIPMem samples are simpler and more precisely



Figure 3: Measuring memorization on individual modalities is not able to extract a strong
 signal. (a)–(b) We measure SSLMem (Wang et al., 2024b) on the individual encoders of our CLIP
 model trained on COCO. (c) Our CLIPMem extracts a stronger memorization signal by using both
 modalities in CLIP jointly.

304

305

captioned. In Appendix A.5, we show that these findings also hold when we operate in the *infinite data regime*, *i.e.*, when we perform only a single training run where no data point is repeated.

Motivated by this insight and findings from supervised learning, where models memorize random 306 labels (Zhang et al., 2016) and where mislabeled data experiences highest memorization (Feldman, 307 2020), we test if the same effect can also be observed in CLIP. We "poison" CLIP's training data 308 by randomly shuffling the captions among 500 of 5000 candidate data points in S_C , creating "mis-309 captioned" data points. We train a model on this data and see that the mis-captioned examples 310 experience significantly higher memorization (CLIPMem of 0.586) compared to "clean" data points 311 (CLIPMem of 0.440). Despite CLIP's contrastive training objective, memorization of clean data 312 points is not significantly affected by training the model with the mis-captioned examples, as we can 313 see by their CLIPMem that is 0.438 on the clean model and 0.440 on the poisoned model.

- 314
- 315 316

4.3 MEASURING MEMORIZATION IN ONE MODALITY DOES NOT YIELD A STRONG SIGNAL

To assess the importance of considering both modalities in our definition of CLIPMem, we evaluate whether existing uni-modal encoder memorization metrics (Wang et al., 2024b) yield a sufficiently strong memorization signal in CLIP. Therefore, we apply SSLMem independently to CLIP's vision and text encoders. Since SSLMem relies on augmentations of the encoder input, we use image crops, like during CLIP training for the vision encoder, and the 5 COCO captions as augmentations for the text, like in (Fan et al., 2023). Our results in Figure 3 show that SSLMem and its naive adaptation to CLIP fail to yield a strong memorization signal. In particular, there is a high overlap in scores between the non-memorized samples from S_S , and candidate examples for memorization S_C . Additionally, the highest reported memorization scores for S_C go up to around 0.65 (for SSLMem on the vision encoder) and 0.73 (for SSLMem on the text encoder). In contrast, our new CLIPMem is able to get a distinct signal for the candidates S_C with respect to S_S and reports a much higher memorization of 0.91. Thereby, our CLIPMem prevents under-reporting the actual memorization in CLIP.

328

330

4.4 MEMORIZATION BETWEEN MODALITIES

331 Our results in Figure 3 indicate that memorization is higher in CLIP's text encoder than in the image 332 encoder (the average SSLMem on S_C in the text encoder is 0.209 vs. 0.168 in the image encoder). 333 To provide further insights into how memorization behaves between the modalities in CLIP, we 334 first analyze the use of augmentations. We compare five cases: (1) no additional augmentations 335 beyond the baseline (image cropping), (2) generating one image using a diffusion model for a given 336 original caption, (3) generating five variations of each image using a diffusion model and randomly 337 selecting one for each training iteration while keeping the caption fixed, (4) using the original image but randomly selecting one of the five COCO captions for each training iteration, and (5) randomly 338 pairing each of the five generated images with one of the five COCO captions. 339

As shown in Figure 14, there is quite a variability in the COCO captions for the same sample.
Hence, some images might not fit well with the
chosen training caption. This imprecise captioning can cause an increase in memorization. We
observe that the effect is mitigated when using
the 5 images with the 5 captions (5th case, see

Table 1: Impact of augmentations.

Case	CLIPMem	Lin. Prob. Acc. (ImageNet)
1 Image, 1 Caption	0.438	63.11% ± 0.91%
1 Image (generated), 1 Caption	0.428	63.97% ± 0.79%
5 Images (generated), 1 Caption	0.424	$64.60\% \pm 0.82\%$
1 Image, 5 Captions	0.423	$64.88\% \pm 0.83\%$
5 Images (generated), 5 Captions	0.417	64.79% ± 0.99%

Table 1). This phenomenon results most likely from the increased number of possible image-text pairs (25), such that individual incorrect or imprecise pairs are not seen so often during training. For the third case, *i.e.*, row three in Table 1, we generate five images with a diffusion model based on all five captions per image from the COCO dataset. However, as we only use the first caption during training, this would introduce many mis-captioned images which significantly lowers performance and increases memorization. To avoid this problem, we removed 6000 mis-captioned samples.

Our results in Table 1 highlight that augmenting text during training reduces memorization and increases performance more than augmenting images. However, applying augmentations of both text and images strikes the right balance between the reduction in memorization and the increase in performance. In fact, applying both augmentations reduces memorization most significantly. These results indicate that memorization in CLIP's is tightly coupled to the captions assigned to the training images with imprecise captions having a destructive effect on CLIP performance and memorization.

358 359 360

4.5 RELATION TO CLIP MEMORIZATION TO (SELF-)SUPERVISED MEMORIZATION

We further provide insights on whether CLIP's memorization behavior is more alike to the one of supervised learning or SSL. This question is highly interesting since the captions in CLIP can be considered as a form of labels, like in supervised learning, whereas the contrastive training objective on the dataset resembles more SSL. We perform two experiments to gain a better understanding of the memorization behavior of CLIP with respect to supervised learning and SSL.

First, we compare an SSL vision encoder pair f and q with the same architecture as CLIP's vision 367 encoder but trained from scratch on COCO using DINO, *i.e.*, standard SSL training. We train f and q 368 using the same candidates as the pair of CLIP models in our previous experiments. Then, we use the 369 SSLMem metric from Wang et al. (2024b) to quantify memorization in the CLIP vision encoder and 370 the SSL encoder, respectively. The CLIP vision encoder has a significantly lower SSLMem than the 371 SSL encoder (0.209 vs. 0.279). Hence, CLIP vision encoders experience lower SSL memorization 372 than SSL trained encoders. To further investigate the difference, we also report the overlap between 373 the top 10% memorized samples between the two models, measured according to SSLMem. With an 374 overlap of only 47 out of 500 (9.4%) samples, we find that CLIP memorizes significantly different 375 samples than SSL encoders. Wang et al. (2024b) had performed a similar experiment on SSL vs. supervised learning and found that the two paradigms also lead to different samples being memorized. 376 While this is, on the one hand, an effect of the different objective function, the difference between the 377 memorized samples in CLIP and SSL is likely also closely connected to the additional captions that



Figure 5: Mitigating memorization in CLIP improves downstream generalization. We train CLIP models with different "augmentations" in the textual domain. (a) We use multiple captions per image. (b) We directly noise the text embeddings using Gaussian noise with a mean of 0 and different standard deviations (adding $\mathcal{N}(0, 0.15)$ achieves the optimal balance of smallest memorization and highest performance). Both strategies successfully reduce memorization while improving performance.

CLIP takes into account. While SSL-trained encoders can memorize atypical images, CLIP encoders can memorize typical images when they have an atypical, imprecise, or incorrect caption.

Additionally, we compare the memorization behavior of CLIP at the neuron level against supervised and SSL-397 trained models. To do so, we train two ViT-Base models on 398 COCO using supervised and SSL (DINO) training. Then, 399 we apply the UnitMem metric (Wang et al., 2024a) to mea-400 sure how much individual neurons memorize individual 401 samples from the training data. High UnitMem indicates 402 neurons memorizing individual data points rather than 403 groups/classes of points. Prior work found that supervised 404 learning results in low UnitMem in lower-layer neurons, 405 indicating group-based learning, while later-layer neurons highly memorize individual data points. In contrast, SSL 406



Figure 4: UnitMem metric: CLIP is between supervised and SSL models.

maintains relatively constant UnitMem, with lower-layer neurons also memorizing individual data 407 points. This difference was attributed to objective functions: supervised learning's cross-entry loss 408 pulls same-class data points together, while SSL's contrastive loss pushes individual data points apart. 409 (Wang et al., 2024a). Our results (Figure 4) reveal that CLIP's memorization behavior falls between 410 supervised learning and SSL. In lower layers, it is much less selective than SSL, *i.e.*, it focuses on 411 groups rather than individual data points, similar to supervised learning. Yet, in later layers, CLIP 412 becomes more selective than SSL, *i.e.*, it memorizes individual data points more in individual neurons, 413 though still less than supervised learning, which exhibits a much higher average per-layer UnitMem. 414 We present additional insights on how memorization evolves over training in Appendix A.7.

415 416

417

4.6 MITIGATING MEMORIZATION WHILE MAINTAINING GENERALIZATION

Experiments in Table 1 suggest that augmentations during training can improve generalization while
also reducing memorization. This is an unexpected synergy, since generalization was shown to decline
when memorization decreases, in both supervised learning (Feldman, 2020) and SSL (Wang et al.,
2024b). To further study how mitigating memorization affects CLIP's downstream generalization, we
explore two orthogonal strategies for "augmenting" text, first in the input space and second in the
embedding space. Additionally, we analyze the effect of removing memorized samples from training.

424 Multiple captions. We vary the number of captions used during training and report memorization 425 and downstream performance in Figure 5a. We find that using more captions during training 426 reduces memorization while improving linear probing accuracy. Our results in Table 3 highlight 427 that using all captions equally often enhances utility without significantly affecting memorization. 428 Since not all dataset have multiple captions, we explore generating them with a language model. Table 7 shows that training CLIP with GPT3.5-generated captions yields extremely similar results 429 in utility and memorization, making this strategy widely applicable. Our findings that modifying 430 text during training can reduce memorization align with Jayaraman et al. (2024), who proposed text 431 randomization, i.e., masking out a fraction of tokens during training as a mitigation for their Déjà Vu



Figure 6: Removing memorized samples according to CLIPMem has a stronger influence on 442 the linear probing accuracy than removing random data points. Removing the mislabeled 443 samples based on CLIPMem improves the performance significantly, followed by a sharper drop 444 when removing atypical samples. 445

memorization. However, unlike our GPT3.5-generated captions, this method reduces performance, 446 likely due to the greater distribution shift introduced by masked tokens. 447

Noising the text embedding during training. To avoid any inherent distribution shifts, we propose 448 applying "augmentations" directly in the embedding space. Specifically, we add small Gaussian noise 449 to text embeddings before computing cosine similarity for contrastive loss. As shown in Figure 5b and 450 Table 8, this strategy is effectively reduces memorization while improving downstream generalization. 451

- Removing memorized samples. Finally, we investigate the effect of removing memorized samples on 452 downstream performance. After training a CLIP model, we identify and remove the most memorized 453 data points, then retrain on the remaining data points. We compare this to two baselines, where we either randomly remove samples or filter out the samples with the lowest CLIP similarity between 455 the training data points' two modalities. We showcase the effect on the downstream linear probing 456 accuracy on ImageNet in Figure 6 with CLIP models trained on COCO and on the CCM3 dataset. 457 For the COCO dataset, when removing up to 100 most memorized data points, we first observe 458 a sharp increase in downstream performance in comparison to removing random samples. Then, 459 downstream performance starts dropping significantly more when removing memorized instead of 460 random samples, until a cutoff point is reached (400-800 removed samples), where further removal 461 based on memorization results in worse performance. For the CC3M dataset, this cutoff occurs later 462 (1,600-3,200 removed samples). While filtering by CLIP similarity also improves performance, it is 463 not as effective as CLIPMem, emphasizing the value of considering memorization as a lens to identify noisy samples. This finding is significantly different than for supervised learning and SSL, where 464 removing highly memorized samples *constantly* harms performance more than the removal of random 465 samples (Feldman, 2020; Wang et al., 2024b). We hypothesize that the effect observed in CLIP might 466 result from the distinction between "mis-captioned" and atypical samples, where the former harms 467 generalization while the latter helps the model learn from smaller sub-populations (Feldman, 2020). 468 We empirically support this hypothesis in Appendix A.3.1. Our finding that CLIP generalization can 469 be improved by identifying inaccurately captioned data points using our CLIPMem and removing 470 them from training is of high practical impact, given that state-of-the-art CLIP models are trained on 471 large, uncurated datasets sourced from the internet with no guarantees regarding the correctness of 472 the text-image pairs. Overall, our results suggest that CLIPMem can help reduce memorization in 473 CLIP while improving downstream generalization.
- 474 475

5 CONCLUSION

476

477 We presented CLIPMem, a formal measure to capture memorization in multi-modal models, such as 478 CLIP. By not only quantifying memorization but also identifying which data points are memorized 479 and why, we provide deeper insights into the underlying mechanisms of CLIP. Our findings highlight 480 that memorization behavior of CLIP models falls between that of supervised and self-supervised 481 models. In particular, CLIP highly memorizes data points with incorrect and imprecise captions, 482 much like supervised models memorize mislabeled samples, but it also memorizes atypical examples. Furthermore, we find that memorization in CLIP happens mainly within the text encoder, which 483 motivates instantiating mitigation strategies there. By doing so, we can not only reduce memorization 484 in CLIP but also *improve* downstream generalization, a result that challenges the typical trade-offs 485 seen in both supervised and self-supervised learning.

486 REFERENCES

- Muhammad Ali and Salman Khan. Clip-decoder: Zeroshot multilabel classification using multimodal
 clip aligned representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4675–4679, 2023.
- Antonis Antoniades, Xinyi Wang, Yanai Elazar, Alfonso Amayuelas, Alon Albalak, Kexun Zhang,
 and William Yang Wang. Generalization vs. memorization: Tracing language models' capabilities
 back to pretraining data. In *ICML 2024 Workshop on Foundation Models in the Wild*, 2024.
- Devansh Arpit, Stanisław Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International conference on machine learning*, pp. 233–242. PMLR, 2017.
- Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Conditioned and composed image retrieval combining and partially fine-tuning clip-based features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4959–4968, 2022a.
- Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Effective conditioned
 and composed image retrieval combining clip-based features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 21466–21474, 2022b.
- Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer:
 Evaluating and testing unintended memorization in neural networks. In 28th USENIX Security
 Symposium (USENIX Security 19), pp. 267–284, 2019.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Nicholas Carlini, Matthew Jagielski, Chiyuan Zhang, Nicolas Papernot, Andreas Terzis, and Florian Tramer. The privacy onion effect: Memorization is relative. *Advances in Neural Information Processing Systems*, 35:13263–13276, 2022.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Satrajit Chatterjee. Learning and memorization. In *International conference on machine learning*, pp. 755–763. PMLR, 2018.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade
 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for
 contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2818–2829, June 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Akshay Dhamija, Manuel Gunther, Jonathan Ventura, and Terrance Boult. The overlooked elephant of object detection: Open set. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1021–1030, 2020.
- Lijie Fan, Dilip Krishnan, Phillip Isola, Dina Katabi, and Yonglong Tian. Improving clip training with
 language rewrites. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.),
 Advances in Neural Information Processing Systems, volume 36, pp. 35544–35575. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/
 2023/file/6fa4d985e7c434002fb6289ab9b2d654-Paper-Conference.pdf.

- 540
 541
 542
 542
 543
 544
 544
 544
 545
 544
 546
 546
 547
 547
 548
 548
 549
 549
 549
 549
 549
 540
 541
 541
 542
 542
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
 544
- Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long
 tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891,
 2020.
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL https://doi.org/10.5281/ zenodo.5143773.
- Matthew Jagielski, Om Thakkar, Florian Tramer, Daphne Ippolito, Katherine Lee, Nicholas Carlini,
 Eric Wallace, Shuang Song, Abhradeep Guha Thakurta, Nicolas Papernot, et al. Measuring
 forgetting of memorized training examples. In *The Eleventh International Conference on Learning Representations*, 2022.
- Bargav Jayaraman, Chuan Guo, and Kamalika Chaudhuri. Déjà vu memorization in vision-language models, 2024. URL https://arxiv.org/abs/2402.02103.
- ⁵⁵⁸ Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

578

579

580

- Nishant Kumar, Siniša Šegvić, Abouzar Eslami, and Stefan Gumhold. Normalizing flow based
 feature synthesis for outlier-aware object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5156–5165, 2023.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Doll'a r, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. URL http://arxiv.org/abs/1405.0312.
- Hongbin Liu, Jinyuan Jia, Wenjie Qu, and Neil Zhenqiang Gong. Encodermi: Membership inference against pre-trained encoders in contrastive learning. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pp. 2081–2095, 2021.
- 575 Casey Meehan, Florian Bordes, Pascal Vincent, Kamalika Chaudhuri, and Chuan Guo. Do ssl models
 576 have déjà vu? a case of unintended memorization in self-supervised learning. *arXiv e-prints*, pp.
 577 arXiv-2304, 2023.
 - Junting Pan, Ziyi Lin, Yuying Ge, Xiatian Zhu, Renrui Zhang, Yi Wang, Yu Qiao, and Hongsheng Li. Retrieving-to-answer: Zero-shot video question answering with frozen large language models. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 272–283, 2023.
- 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 *International conference on machine learning*, pp.
 8748–8763. PMLR, 2021.
- Ildus Sadrtdinov, Nadezhda Chirkova, and Ekaterina Lobacheva. On the memorization properties of contrastive learning. *arXiv preprint arXiv:2107.10143*, 2021.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*, 2018.
- Congzheng Song, Thomas Ristenpart, and Vitaly Shmatikov. Machine learning models that remember
 too much. In *Proceedings of the 2017 ACM SIGSAC Conference on computer and communications* security, pp. 587–601, 2017.

594 595 596	Haoyu Song, Li Dong, Weinan Zhang, Ting Liu, and Furu Wei. Clip models are few-shot learners: Empirical studies on vqa and visual entailment. In <i>Proceedings of the 60th Annual Meeting of the</i> <i>Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 6088–6100, 2022.
597 598 599 600 601	Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: the new data in multimedia research. <i>Commun. ACM</i> , 59(2):64–73, January 2016a. ISSN 0001-0782. doi: 10.1145/2812802. URL https://doi.org/10.1145/2812802.
602 603 604	Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. <i>Communications of the ACM</i> , 59(2):64–73, 2016b.
605 606 607	Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization without overfitting: Analyzing the training dynamics of large language models. <i>Advances in Neural Information Processing Systems</i> , 35:38274–38290, 2022.
609 610	Wenhao Wang, Adam Dziedzic, Michael Backes, and Franziska Boenisch. Localizing memorization in ssl vision encoders. <i>arXiv preprint arXiv:2409.19069</i> , 2024a.
611 612 613	Wenhao Wang, Muhammad Ahmad Kaleem, Adam Dziedzic, Michael Backes, Nicolas Papernot, and Franziska Boenisch. Memorization in self-supervised learning improves downstream generalization. In <i>The Twelfth International Conference on Learning Representations (ICLR)</i> , 2024b.
615 616 617	Zhengbo Wang, Jian Liang, Ran He, Nan Xu, Zilei Wang, and Tieniu Tan. Improving zero-shot generalization for clip with synthesized prompts. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 3032–3042, 2023.
618 619 620	Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understand- ing deep learning requires rethinking generalization. In <i>International Conference on Learning</i> <i>Representations</i> , 2016.
621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 635 636 637 638 639 640 641 642 643 644 645	Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free adaption of clip for few-shot classification. In <i>European</i> conference on computer vision, pp. 493–510. Springer, 2022.
646 647	

648 APPENDIX А 649

672

673

674

675 676

687

691

650 A.1 EXTENDED BACKGROUND 651

652 Déjà Vu Memorization in CLIP. The Déjà Vu memorization framework (Jayaraman et al., 2024) is the only existing other work that attempts to quantify memorization in vision-language models. 653 It uses the text embedding of a training image caption to retrieve relevant images from a public 654 dataset of images. It then measures the fraction of ground-truth objects from the original image 655 that are present in the retrieved images. If the training pair is memorized, retrieved images have 656 a higher overlap in ground truth objects, beyond the simple correlation. While valuable, several 657 aspects warrant further consideration for broader applicability of the framework. First, its focus 658 on object-level memorization ignores non-object information like spatial relationships or visual 659 patterns that can also influence memorization (Feldman, 2020; Wang et al., 2024b). To perform object 660 retrieval, the framework also relies on object detection and annotation tools, which may introduce 661 variability based on the accuracy and robustness of these tools. Additionally, the assumption that 662 public datasets with similar distributions to the training data are readily available may not always hold, necessitating alternative approaches. Moreover, the framework does not analyze why certain images are memorized limiting detailed analysis. Finally, while Déjà Vu must address the challenge 664 of distinguishing between memorization and spurious correlations, CLIPMem avoids this by directly 665 assessing memorization on the output representations of the model. One notable difference between 666 the results of our approach and Déjà Vu's is that their findings show that their mitigation strategies 667 can reduce memorization, but at the cost of decreased model utility. CLIPMem, in contrast, does not 668 observe trade-offs between memorization and performance. 669

670 A.2 EXTENDED EXPERIMENTAL SETUP 671

General Setup. All the experiments in the paper are done on a server with 4 A100 (80 GB) GPUs and a work station with one RTX 4090 GPU(24 GB). We detail the setup for our model training, both CLIP and SSL (relying on DINO) in Table 2.

	Model Training		Linear Probing			
	CLIP	DINO	Supervised ViT	CLIP	DINO	Supervised ViT
Training Epoch	100	300	100	45	45	45
Warm-up Epoch	5	30	5	5	5	5
Batch Size	128	1024	128	4096	4096	4096
Optimizer	Adam	AdamW	Adam	LARS	LARS	LARS
Learning rate	1.2e-3	2e-3	1e-3	1.6	1.6	1.6
Learning rate Schedule	Cos. Decay	Cos. Decay	Cos. Decay	Cos. Decay	Cos. Decay	Cos. Decay

Table 2: Experimental Setup. We provide details on our setup for encoder training and evaluation.

Experimental Setup for SSLMem. To experimentally evaluate memorization using the SSLMem 686 framework (Wang et al., 2024b), the training dataset S is split into four sets: shared set (S_S) used for training both encoders f and g; candidate set (S_C) used only for training encoder f; independent set 688 (S_I) data used only for training encoder g; and an additional extra set (S_I) from the test set not used 689 for training either f or g. For training encoders, encoder f is trained on $S_S \cup S_C$, while encoder g is 690 trained on $S_S \cup S_I$. The alignment losses $\mathcal{L}_{\text{align}}(f, x)$ and $\mathcal{L}_{\text{align}}(g, x)$ are computed for both encoders, and the memorization score m(x) for each data point is derived as the difference between these 692 alignment losses, normalized to a range between -1 and 1. A score of 0 indicates no memorization, +1 indicates the strongest memorization by f, and -1 indicates the strongest memorization by g. 693

694 Normalization on CLIPMem. For improved interpretability, we normalize our CLIPMem scores to 695 a range of [-1, 1]. A memorization score of 0 indicates no memorization, +1 indicates the strongest 696 memorization on CLIP model f, and -1 indicates the strongest memorization on CLIP model g. We 697 find the normalized CLIPMem score for a dataset using the following process: For each image-text pair (I, T), we first calculate the CLIPMem score as the difference in alignment scores between two 699 CLIP models f and g. Once CLIPMem scores are computed for all data points, we normalize them by dividing each score by the range, which is the difference between the maximum and minimum 700 scores in the dataset. Finally, we report the normalized CLIPMem score for a dataset as the average 701 of these normalized values.

+ Random ClipMem Based 704 705 706 708 709 710 0.63 711 0 50 100 200 400 800 1600 3200 712 # of Sample removed 713 714 Figure 7: **Removing memorized samples.** We show the effect on downstream performance in terms 715 of ImageNet linear probing accuracy and CLIPMem for a CLIP model trained on COCO using 5 text 716 captions instead of 1, like done in Figure 6. We observe the same trend, with the difference that the peak is at roughly 500 removed samples rather than 100. This is likely due to the increase in captions 717 (by factor 5) that causes increase in mis-captioned samples. 718 719 720 0.62 Linear probing Acc. 0.00 0.20 0.20 0.20 0.20 0.20 Random 721 Memorization Based 722 723 724 725 726 727 728 729 200 400 800 1600 3200 Ò 50 100 730 # of Sample removed 731 732 Figure 8: Removing memorized samples in supervised learning. We train a ViT-tiny on CI-733 FAR10 (Krizhevsky et al., 2009) using supervised learning. We use our evaluation setup with S_C , S_S , S_I , and S_E to approximate the memorization metric from Feldman (2020). We use 5000 samples in 734 S_C , but before training, we flip the labels of 200 samples. We calculate memorization over all samples 735 in S_C and test the linear probing accuracy with ImageNet resized to 32*32 on the representations 736 output before the original classification layer. 737 738 739 A.3 ADDITIONAL EXPERIMENTS 740 741 A.3.1 MEMORIZATION VS. GENERALIZATION IN CLIP 742 743 Extending evaluation. In Figure 7, we perform the same experiment as in Figure 6, but on a CLIP

Extending evaluation. In Figure 7, we perform the same experiment as in Figure 6, but on a CLIP
 model trained with 5 captions instead of 1. We observe the same trend, with the difference that the
 peak is at roughly 500 removed samples rather than 100. This is likely due to the increase in captions
 (by factor 5) that causes increase in mis-captioned samples.

747 Verifying the hypothesis on memorizing mis-captioned samples through supervised learning. 748 We repeat the same experiment in the supervised learning setup to understand where the increase 749 and then decrease in linear probing accuracy stems from. To test our hypothesis that it stems from 750 "mis-captioned" samples, we "poison" our supervised model by flipping the labels of 200 data points 751 before training. Then, we approximate the memorization metric from Feldman (2020) in our setup 752 and remove highly memorized vs. random data points. In the same vein as in Appendix A.3.1, we 753 first observe an increase in linear probing accuracy when removing memorized samples (instead of random samples). The peak is at roughly 200 data points, *i.e.*, the number of deliberately mislabeled 754 samples. Until the cutoff point at roughly 3200 examples, linear probing accuracy is still higher when 755 removing most memorized rather than random samples, which might suggest that there are other

Table 3: Using different/multiple captions during training. We evaluate CLIPMem how memorization on different data subsets and linear probing accuracy on ImageNet differ when using 1 caption (baseline), 5 COCO captions, one chosen at random at every round (random), and 5 COCO captions, but all chosen equally often, *i.e.*, 20 out of 100 training epochs (balanced). We observe that increasing the number of captions reduces highest memorization. Yet, only when we balance the usage of caption, also model performance increases.

	baseline	random	balanced
Avg. CLIPMem (Top 10 samples)	0.792	0.788	0.790
Avg. CLIPMem (Top 20%)	0.552	0.531	0.540
Linear Probing Acc.	$63.11\% \pm 0.91\%$	$62.44\% \pm 1.18\%$	$64.88\% \pm 0.83\%$

Table 4: The CLIPMem and linear probing accuracy of model trained with original coco captions and captions generated by GPT3.5. For 'Single Caption', only one caption is used during training. For 'Five Caption', all five caption are used equally during training (every caption trained for 20 epoch out of 100). The linear probing accuracy is tested on ImageNet

	COCO		GPT3.5	
	Single Caption	Five Caption	Single Caption	Five Caption
CLIPMem LP Acc	0.438 63 11% + 0 91%	0.423 64 88% + 0 83%	0.430 63 09% + 1 12%	0.411 64 47 + 0 72%

outliers or inherently mislabeled samples whose removal improves model performance. After the cutoff, we observe the behavior as observed in prior work (Wang et al., 2024b; Feldman, 2020) that reducing memorization harms generalization more than reducing random data points from training.

A.3.2 THE EFFECT OF CAPTIONS

In Table 3, we show that using more captions during training reduces memorization and that by using
 each caption at the same frequency over the training epochs, we can additionally improve model
 performance. Additionally, we show that captions generated by GPT3.5 have the same effect as the
 original COCO captions on memorization and linear probing accuracy in Table 4.

A.4 THE EFFECT OF MODEL SIZE

In Table 5, we present how the model size affects the memorization level of CLIP models. Both models are trained using the same dataset and settings. We observe that with more parameters (larger model size), encoders have higher memorization capacity. This aligns with findings from previous research (Wang et al., 2024b; Feldman, 2020; Meehan et al., 2023).

A.5 VERIFICATION OF INFINITE DATA REGIMES

To evaluate CLIPMem over infinite data regimes (*i.e.*, using a single training run where no data point is repeated), we use a subset D (containing 7050000 samples) of YFCC100M dataset (Thomee et al., 2016b) to train another pair of ViT-Base models for only 1 epoch. Following our definition of CLIPMem, we further divide D into S_S with 6950000 samples, S_C with 50000 samples, and S_I with 50000 samples. The reason we use 7M (6950000+50000) samples to train either model f or model q is to make sure the newly trained model has the same number of training samples as the model trained with K-epoch runs (70000 samples/epoch * 100 epoch). The results in Table 6 show that the model trained with infinite data regimes has higher linear probing accuracy on ImageNet as a downstream task and lower memorization scores, as measured by CLIPMem. This aligns with the fact that duplicated data points increase the memorization level and make the model over-fit, hence reducing the generalization (Wang et al., 2024b; Feldman, 2020). The results in Figure 9 show that the most memorized samples according to CLIPMem in the model trained with infinite data regimes are also samples with imprecise or incorrect captions. This aligns with our statements in Section 4.5. 810 Table 5: CLIPMem and linear probing accuracy of models with different sizes. The models are 811 trained using identical settings and the same subset of the COCO dataset. Linear probing accuracy is 812 tested on the ImageNet dataset as the downstream task.

813					
814		Model	CLIPMem L	in. Prob. Acc. (ImageNet)	
815		ViT-base (Baseline in main paper)	0.438	$63.11\% \pm 0.91\%$	
816		ViT-large	0.457	$67.04\% \pm 1.05\%$	
817					
818	M. Anthe	Series taken during heritage surveys in 2007.	1000	Photo By Terri Hodges	
819		captions for these, however happy to		www.craftingthegalaxy.co	m Useage is
820	A CONTRACT	provide some comment if required.		Creative Commons but ci	redit is required
821		Cumbros and Toltos Sossis Bailroad, historiaal			
822	T	narrow gauge railroad between Chama, New		Mellwood Art Center in L	ouisville,
823	2	Mexico, and Antonito, Colorado. These photos		KY (1-27-06)	
824		are from the station and railyard in Chama.			
825					
826		Day Four set here. Entire Japan collection here	. Catal	Lisa Brewster Sent from	n my Palm Pre
827				1	,
828	The second second second second				
829	15	Natural Bridge & Rockshelter Wisconsin State			
830		Natural Area #105 Natural Bridge State Park	The	At Hellfest with After Fore	ever
831	γ	Sauk County	1.0		
832		Mountain Francisco Filord Trin - Drintol TN	-		
833		and environs Photo by Amy C. Evans, SFA	Finance CORE	On the way to visit Weste	rn
834	AL GENI	oral historian June 2009		Kentucky University (9/20	/08)
835	-	www.southerntoodways.org			

Figure 9: Top 10 memorized samples according to CLIPMem in the model trained under infinite data regimes on YFCC100M. The model is trained for one epoch, *i.e.*, seeing each training data point exactly once. Even in this setup, the most memorized samples are still the ones with imprecise or incorrect captions.

To verify the effectiveness of CLIPMem over other similar multi-modal models, we train a BLIP model on COCO dataset following the same settings as the baseline model in the main paper. We present the results for CLIPMem over all 4 data subsets in Figure 10, which is in agreement with the results of the CLIP model in Figure 2a. We also present the UnitMem results for BLIP model in Figure 4, which is also very similar to the result of CLIP models

A.7 MEMORIZATION DISTRIBUTION DURING TRAINING

A.6 EVALUATION ON BLIP

We present the distributions of neurons with highest UnitMem during training in Figure 11. These results highly consistently indicate that in the early stages of training, neuronal memory occurs mainly 852 in the lower layer of the clip model, while in the middle and later stages of training, neuronal memory is more concentrated in the later layer of the model.

857

858

859

836

837

838

839

840 841

842 843

844

845

846

847

848 849

850

851

A.8 HUMAN VS MACHINE GENERATED CAPTIONS

For each image in the COCO dataset, we use GPT 3.5 (specifically, gpt-3.5-turbo) to generate 5 captions (from scratch). We use the following instruction in the OpenAI API:

```
860
      def generate_description_for_image(image_caption, clip_features):
         prompt = f"Here is an image with the caption: '{image_caption}'. "
861
         prompt += f"Based on this caption and the visual features
862
          represented by this embedding '{clip_features}',
         please generate a new detailed description."
```

Table 6: Evaluation of CLIPMem under infinite data regimes, *i.e.*, seeing every data point only
 once during training vs training with 100 epochs. We observe that both setups reach comparable
 downstream accuracy and memorization.





Figure 11: Distribution of top 1%, 3%, and 5% neurons with highest UnitMem during training.
We train a CLIP model on COCO standard image cropping and no text augmentation following the settings of baseline model in main paper. We record the neurons with top 1%, 3%, and 5% of highest UnitMem during training (every epoch). We observe that while during early training stages, memorization focuses on lower layers, during later stages of training, it is mainly in the last layers.

Table 7: The machine generated captions provide similar performance to the original humangenerated captions. We report the CLIPMem and linear probing accuracy of model trained with original COCO captions and captions generated by GPT 3.5. For the 'Single Caption', only a single caption is used during training. For 'Five Captions', all five captions are used equally during training (every caption trained for 20 epochs out of 100). The linear probing accuracy is tested on the ImageNet dataset as the downstream task.

	COCO	1	GPT	3.5
	Single Caption	Five Captions	Single Caption	Five Captions
CLIPMem near Probing Accuracy (ImageNet)	$\begin{array}{c} 0.438 \\ 63.11\% \pm 0.91\% \\ \end{array} 6$	$0.423 \\ 4.88\% \pm 0.83\%$	$\begin{array}{c} 0.430 \\ 63.09\% \pm 1.12\% \end{array}$	$\begin{array}{c} 0.411 \\ 64.47 \pm 0.72\% \end{array}$
10000 8000 6000 2000 0.65 0.70 0.75 0.80 0 Cosine Sin (2) COCO (Average	0.85 0.90 0.95 1.00 nilarity	10000 8000 6000 4000 2000 0.65 0.70	0.75 0.80 0.85 0 Cosine Similarit	90 0.95 1.00

Figure 12: Pairwise cosine similarity of 5 captions from COCO and generated by GPT3.5.

Noise	CLIPMem	Lin. Prob. Acc. (ImageNet)
None	0.438	$63.11\% \pm 0.91\%$
N(0.01)	0.435	$63.36\% \pm 0.88\%$
$\mathcal{N}(0.05)$	0.428	$64.02\% \pm 1.12\%$
N(0.10)	0.421	$64.95\% \pm 0.96\%$
$\mathcal{N}(0.15)$	0.417	$65.34\% \pm 0.84\%$
$\mathcal{N}(0.20)$	0.422	$64.83\% \pm 0.92\%$
$\mathcal{N}(0.25)$	0.436	$63.28\% \pm 0.79\%$
N(0.30)	0.447	$61.50\% \pm 0.86\%$
$\mathcal{N}(0.50)$	0.491	$57.04\% \pm 1.11\%$
$\mathcal{N}(0.75)$	0.501	$52.28\% \pm 0.98\%$
$\mathcal{N}(1.00)$	0.504	$51.92\% \pm 1.03\%$

Table 8: Noising text embedding during training. We present the impact of adding noise to the text embedding during training for the ViT-base trained on COCO.





1092 1093 1094

1095

1096

1098 1099

1104

1111

1112 1113 1114

1115

1116 1117 1118

1119

1120

1121 1122

1123

1124

1125 1126 1127

1128 1129 1130



A girl in pink sweater putting a blue umbrella over a yellow fire hydrant.



Closeup of two street signs that read "Airport Pkwy" and "Karmill Ave."



A sign with Oriental writing and the words saying Hyatt on the Bund.



A pink Hello Kitty microwave on a store shelf.



A sign saying "Don't Honk, \$350 Penalty" on a pole.



A group of artistic surfboards are displayed in a tent.

Figure 15: The 10 samples with lowest CLIPMem in the CLIP model trained with 1 caption. We can see that these samples contain clear concepts and precise captions.





two males and a female in a red top holding some flowers



A woman opening the trunk of her car.



a sleeping toddler laying on a womans shoulder



An open point of view of a room with various things all around



A group of people standing on top of a field together.



a pole with some yellow lights in front of a narrow building



a parking lot with a bunch of cars in it



A person is taken in this very picture.



I am unable to see an image above

Figure 16: The 10 samples with highest CLIPMem in the CLIP model trained with 1 caption.

1131 We can see that these samples contain atypical, difficult samples with imprecise or incorrect captions. 1132

1133

A yellow and blue fire hydrant surrounded by leaves

A sign is displayed on a pole

A man in maroon shirt standing

next to a stainless steel refrigerator.

A blue street sign that reads "Thelonius

that says bump.

Monk Circle."



1101			
1134			
1135			
1136			
1137			
1138	COCO Image COCO Caption	Stabl	e Diffusion Image
1139	A group of people that are touching a animal.		
1140	A group of people walking through and petting a shee A woman is petting a sheep that has been sheered.		
1141	A woman petting a sheep inside a pen at a fair. A sheep in a stall being pet by a woman		
1142			
1143	A man standing next to palm trees next to a bunch of		
1144	A SURFING BOARD STAND WITH A PERSON STANDIN	G State of the second sec	· · · · · · · · · · · · · · · · · · ·
1145	A man standing in front of many surfboards at the bea	sh. 77 10-11 ////	
1146	A paim tree lined beach contains a large surf board rate and people milling about.		
1147	A group of people standing on the sands at the beach		
1148	A tree filled with lots of lemons and green leaves.		
1149	large yellow fresh fruit tree wet from rain Yellow apples hanging from a leafy tree.		
1150	A close up of several oranges in a tree.		
1151	Several round, yenow nuits growing on a tree.	Constant of the second second	
1152			
1153			
1154			
1155			
1156	Figure 17: Samples of images generated by	Stable Diffusion. W	e present the generated images
1157	based on the COCO captions.		
1158			
1159			
1159 1160			
1159 1160 1161			
1159 1160 1161 1162			
1159 1160 1161 1162 1163			
1159 1160 1161 1162 1163 1164			
1159 1160 1161 1162 1163 1164 1165	COCO Captions		GPT3.5 Captions
1159 1160 1161 1162 1163 1164 1165 1166	COCO Captions	each other. three p	GPT3.5 Captions
1159 1160 1161 1162 1163 1164 1165 1166 1167	COCO Captions Three parrots sitting on a wooden branch next to Three brighty-colored parrots are interacting in Three very colorful parrots perioded on different	each other. three t tree. three t	GPT3.5 Captions parrots sitting on a branch Jue and yellow parrots sitting on a branch Jue and yellow parrots sitting on a tree branch
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168	COCO Captions Three parrots sitting on a wooden branch next to Three brightly-colored parrots are interacting in Three very colorful parrots perched on different Three brightly colored birds on top of a cluster o	each other. three t a tree. three t ticks. three t branches. a trio c bair feathers a trio c	GPT3.5 Captions parrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169	COCO Captions Three parrots sitting on a wooden branch next to Three brighty-colored parrots are interacting in Three very colorful parrots perched on different a Three brightly colored birds on top of a cluster of Three parrots sitting on branches and grooming	each other. three r tree. three t ticks. three t 'branches. a trio c heir feathers. a trio c	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a branch room with a cink and mirror in it
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170	COCO Captions Three parrots sitting on a wooden branch next to Three brighty-colored parrots are interacting in Three very colorful parrots perched on different Three brighty colored birds on top of a cluster of Three parrots sitting on branches and grooming A dim bathroom with a light over the sink A bathroom is dimly lit with a single bulb.	each other. three p a tree. three t ticks. three t branches. a trio c heir feathers. a trio c a bath a smal	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a branch room with a sink and mirror in it I bathroom with a sink and mirror
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171	COCO Captions Three parots sitting on a wooden branch next tr Three brighty-colored parots are interacting in Three very colorful parots perched on different : Three brighty colored birds on top of a cluster of Three parots sitting on branches and grooming A dim bathroom with a light over the sink A bathroom with a single bulb. A dimly light bathroom with a single bulb. A dimly light bathroom with a single mult. A bathroom with a single mult.	each other. three p tree. three t ticks. three t branches. a trio c heir feathers. a trio c a bath a smal a bath ht. a smal	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a branch room with a sink and mirror loathroom with a sink and a mirror
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172	COCO Captions Three parrots sitting on a wooden branch next to Three brighty-colored parrots are interacting in Three brighty-colored birds on top of a cluster of Three parrots sitting on branches and grooming Three parrots sitting on branches and grooming A dim bathroom with a light over the sink A bathroom is dimly lit with a single bulb. A dimhy light bathroom with a single bulb. A district on with a single bulb. A district on with a single bulb. A bathroom with a sink, mirror, and over head light bathroom sink with a light hanging above.	each other. three p a tree, three t ticks. three t branches. a trio c heir feathers. a trio c a bath a smal ht. a smal	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch for parrots sitting on a branch room with a sink and mirror l bathroom with a sink and a mirror room with a sink and a mirror room with a sink and a mirror
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173	COCO Captions Image: Construction of the provided in theterpred in the provided in the provided in theterpred	each other. three t tree. three t isks. three t branches. a trio c heir feathers. a trio c a bath ht. a smal bath ht. a smal a bath	GPT3.5 Captions parrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a branch room with a sink and mirror in it I bathroom with a sink and a mirror room with a sink and a mirror it bathroom with a sink and a mirror room with a sink and a mirror room with a sink and a mirror room with a sink and a mirror in it standing next to a car talking on a cell
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174	EXAMPLE EXAMPLE	each other. three t tree. three t branches. a trio o heir feathers. a trio o heir feathers. a trio ot ha bath ht. a smai a bath ht. a smai a bath a bath a bath a bath a bath a a man a man a man	GPT3.5 Captions parrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a branch room with a sink and mirror in it bathroom with a sink and a mirror room with a sink and a mirror it bathroom a sink and a mirror room with a sink and a mirror standing next to a car talking on a cell standing next to a car talking on a cell phone
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175	EXAMPLE EXAMPLE	each other. three p a tree. three t ticks. three t branches. a trio o heir feathers. a trio o a bath a smal a bath a smal a bath a man a ma	GPT3.5 Captions parrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch blathroom with a sink and mirror room with a sink and a mirror is bathroom with a sink and a mirror room with a sink and a mirror soom with a sink and a mirror is standing next to a car talking on a cell standing next to a parked car talking on a cell standing next to a car talking on the phone
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176	EXAMPLE EXAMPLE	each other. three p b tree. three t ticks. three t branches. a trio c heir feathers. a trio c a bath a bath ht. a smal bath a bath ye. a man ge. a man der. a man the the town wall.	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch blathroom with a sink and mirror room with a sink and a mirror lobathroom with a sink and a mirror room with a sink and a mirror soom with a sink and a mirror standing next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on the phone standing next to a car talking on the phone
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	EXAMPLE EXAMPLE	each other. three p s tree. three t biranches. a trio c heir feathers. a trio c a bath a smal a bath ht. a smal a bath a bath stan a man a man a man a man a man a man a trio c a bath a smal a bath stan bath a man a m	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch bathroom with a sink and mirror room with a sink and mirror room with a sink and a mirror bathroom with a sink and a mirror room with a sink a a mirror room with a sink and a mirror room with a sink and mirror room with a sink and mirror room with a sink and mirror room with a sink a a mirror room with a sink and mirror room with
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178	EXAMPLE EXAMPLE EXAMPLE	each other. three p a tree. three t branches. a trio c heir feathers. a trio c heir feathers. a trio c a bath a smal bath ht. a smal a bath a bath a bath a bath a bath a can a man a man der. a man er next to the brown wall. at a fence. two gi A woo	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a branch room with a sink and mirror noom with a sink and mirror l bathroom with a sink and a mirror room with a sink and a mirror room with a sink and a mirror room with a sink and a mirror standing next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on the phone standing next to a car talking on the phone standing next to a car talking on the phone standing next to a car talking on the phone
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179	COCCO Captions Image: State of the state o	each other. three t tree. three t isks. three t branches. a trio c heir feathers. a trio c heir feathers. a trio c heir feathers. a bath ht. a smal a bath ht. a smal a bath ht. a a bath a bath construction of the car. a man a man construction of the brown wall. transform of the car. two gin at a fence. two gin A woo two gin wall. two gin	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a branch room with a sink and mirror room with a sink and a mirror b bathroom with a sink and a mirror room with a sink and a mirror standing next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on the phone standing next to a car talking on the phone raffes standing next to each other raffes standing next to a wooden fence d fence is behind two giraffes raffes standing next to a wooden fence raffes standing next to a moden fence raffes standing in front of a fence
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180	COCCO Captions Image: State of the state o	each other. three r tree. three t 'branches. a trio c heir feathers. a trio c heir feathers. a trio c heir seathers. a trio c a bath ht. a smal a bath ht. a smal a bath maiks away from the car. a man a man der. a man der. a man a man c. a man a man a man c. a man a man a man a man a man a man b the brown wall. two gi t a fence. two gi wall. two gi	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch blathroom with a sink and mirror room with a sink and mirror room with a sink and a mirror lo bathroom with a sink and a mirror room with a sink and a mirror room with a sink and a mirror room with a sink and a mirror it bathroom with a sink and a mirror room with a sink and a mirror atanding next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on a cell standing next to a car talking on the phone raffes standing in front of a wooden fence d fence is behind two giraffes affes standing next to a avooden fence affes standing next to a avooden fence affes standing next to a fence affes standing in front of a fence lice officers on horses lice officers no horses
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181	Example 1 Example 2 Example 2	each other. three r tree. three t icks. three t branches. a trio c heir feathers. a trio c heir feathers. a trio c a bath ht. a smal a bath ht. a smal a bath a bath a man a man der. a man der. a man der. a man a man c man t a fence. two gi x a fence. two gi y wall. two gi	GPT3.5 Captions barrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch blathroom with a sink and mirror room with a sink and mirror room with a sink and a mirror room w
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181	<text><text><text><text><text><text><text><text><text></text></text></text></text></text></text></text></text></text>	each other. three p a tree. three t branches. a trio o heir feathers. a trio o heir feathers. a trio o heir feathers. a trio o heir feathers. a bath ht. a smal a bath ht. a smal a bath ht. a smal a bath a bath ge. a man a man a man a man der. a man der. a man der. a man a man a man a man b a man b a man a man a man a man a man a man a man b a man b a man a man a man a man a man b a man b a man b a man b a man a man a man a man a man a man b a man b a man b a man a man b a man b a man b a man a ma man a m	GPT3.5 Captions parrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch blathroom with a sink and mirror room with a sink and mirror room with a sink and a mirror standing next to a car talking on a cell standing next to a car talking on the phone raffes standing in front of a wooden fence d fence is behind two giraffes raffes standing in front of a fence lice officers on horses lice officers on horses
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182 1183	<text><text><text><text><text><text><text><text><text><text><text><text></text></text></text></text></text></text></text></text></text></text></text></text>	each other. three p a tree. three t ticks. three t branches. a trio o heir feathers. a trio o heir feathers. a trio o a bath ht. a smal a bath walks away from the car. a man a man a. a man der. a man der. a man der. a man trenext to the brown wall. two gi a ta fence. two gi a wall. two gi two gi two gi two gi b of horses. a coup	GPT3.5 Captions parrots sitting on a branch plue and yellow parrots sitting on a branch plue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch to bathroom with a sink and mirror room with a sink and mirror in it bathroom with a sink and a mirror room with a sink and a mirror room with a sink and a mirror room with a sink and a mirror standing next to a car talking on a cell standing next to a car talking on the phone raffes standing in front of a wooden fence affes standing in front of a fence lice officers on horses lice officers on horses lice officers sitting on horses lice officers on horses
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1180 1181 1182 1183 1184	<text><text><text><text><text><text><text><text><text><text><text><text></text></text></text></text></text></text></text></text></text></text></text></text>	each other. three p b tree. three t ticks. three t branches. a trio c heir feathers. a trio c heir feathers. a trio c a bath a bath ht. a smal br. a man der. a man der. a man der. a man der. a man t a fence. two gi t a fence. two gi frowd two po es two po two gi two po two gi two po two po for horses. a coup	GPT3.5 Captions parots sitting on a branch blue and yellow parots sitting on a branch blue and yellow parots sitting on a tree branch of parots sitting on a tree branch of parots sitting on a tree branch of parots sitting on a tree branch blue and yellow parots sitting on a tree branch of parots sitting on a tree branch blue and yellow parots sitting on a tree branch of parots sitting on a tree branch blue and yellow parots sitting on a tree branch or parots sitting on a tree branch blue and yellow parots a branch blue and yellow parots sitting on a cell standing next to a car talking on the phone taffes standing infront of a wooden fence d fence is behind two giraffes raffes standing in front of a fence affes standing in front of a fence affes standing in front of a fence taffes standing in front of a fence taffes standing in front of a fence taffes standing next to a parote branch affes standing in front of a fence affes standing in front of a fence the officers on horses lice officers on horses
1159 1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1180 1181 1182 1183 1184 1185	<text><image/><text><text><text><text><text><text><text><text><text><text><text><text></text></text></text></text></text></text></text></text></text></text></text></text></text>	each other. three p tree. three t branches. a trio c heir feathers. a trio c heir feathers. a trio c heir feathers. a trio c a bath a smal a bath ht. a smal a bath ht. a smal a bath ht. a smal a bath a bath a can a man can er next to the brown wall. two gi a der. a man a man can a man a man a man a man a man a ta fence. two gi wall. two gi two gi two po es two po o of horses. a coup PT3.5. We present	GPT3.5 Captions parrots sitting on a branch blue and yellow parrots sitting on a branch blue and yellow parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch of parrots sitting on a tree branch blathroom with a sink and mirror room with a sink and mirror room with a sink and a mirror bathroom with a sink and a mirror room with a sink and mirror room with a sink and a mirror room with a sink and mirror room with a si



Figure 19: Evaluation of SSLMem and CLIPMem on a CLIP model trained on COCO. Extended version of Figure 3 where we also include SSLMem calculated on encoders trained with 5 captions instead of 1. The trends in both cases are the same. SSLMem for the CLIP Models trained with the 5 captions is slightly higher since SSLMem uses the captions as augmentations for the calculation of the memorization. Overall, our CLIPMem reports the strongest memorization signal for CLIP.

- 1235 1236
- 1237
- 1238
- 1239
- 1240
- 1241