

---

# Probabilistic Active Few-Shot Learning in Vision-Language Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Pre-trained vision-language models (VLMs) have shown to be an useful model  
2 class for zero- and few-shot learning tasks. In this work, we investigate probabilistic  
3 active few-shot learning in VLMs by leveraging post-hoc uncertainty estimation  
4 and targeted support set selection. To equip VLMs with a notion of uncertainty  
5 on the target task, we utilize a Laplace approximation to the posterior of the VLM  
6 and derive a Gaussian approximation to the distribution over the cosine similarities.  
7 Further, we propose a simple adaptive target region selection based on  $k$ -nearest  
8 neighbour search and evaluate on a series of selection strategies from the Bayesian  
9 experimental design literature. Our experiments on standard benchmarks show that  
10 leveraging epistemic uncertainties leads to improved performance and that further  
11 improvements can be obtained by targeting the selection towards the query region.

## 12 1 Introduction

13 The rise of foundation models [4, 6, 9, 30] has led to their increasing adoption in downstream tasks  
14 where data is scarce [16, 42]. Moreover, in many real-world settings it is imperative that predictions  
15 are reliable and that sources of uncertainties are captured and incorporated to avoid failure modes. The  
16 paradigm of *active few-shot learning* (or *active fine-tuning*) [1, 17, 40] aims to tackle the challenge of  
17 actively selecting a support set (training set for adaption) that is most informative for the downstream  
18 task. However, classical approaches, *e.g.*, from the coreset literature [36] or information theory [14],  
19 typically do not incorporate all sources of uncertainties into their metric of informativeness. Recent  
20 works in Bayesian active learning [15] aim to address this issue by performing selection of support  
21 set candidates based on their effect on the epistemic uncertainty of the model [11] or the predictive  
22 distribution [3]. Moreover, progress in Bayesian deep learning [29] has resulted in methods that can  
23 efficiently estimate epistemic uncertainties in a post-hoc manner [23, 8], making them particularly  
24 attractive for active few-shot learning of large scale models.

25 In this work, we investigate probabilistic active few-shot learning for vision-language models (VLMs)  
26 and show benefits of incorporating uncertainties in the support set selection process as well as  
27 targeting the selection towards the query region. For this, we propose an uncertainty estimation-based  
28 approach by leveraging a Laplace approximation [23] to the posterior of a pre-trained CLIP [30]  
29 model. We derive a Gaussian approximation to the distribution over cosine similarities between  
30 the image and text embeddings, and investigate different scoring mechanisms for the support set  
31 candidate selection. In addition, we propose a simple adaptive target region selection based on  
32  $k$ -nearest neighbour ( $k$ -NN) search. In our experiments, we evaluate two few-shot classification  
33 settings (*i*) support set selection from a large cross-domain training data source and (*ii*) selection from  
34 the training set. We find improved performance over naïve selection for uncertainty-based selection  
35 methods and further improvements when the selection is based on an adaptive target region.

36 **Fig. 1** illustrates the setting we are considering in  
 37 this work: Given a pre-trained VLM, we aim to  
 38 predict labels for a query set of images of a novel  
 39 downstream task. The VLM agent  $\mathcal{M}_0$  is asked  
 40 to first estimate its uncertainty over the predic-  
 41 tions on the query set, where the difficulty of the  
 42 prediction is proportional to the predictive uncer-  
 43 tainty. To avoid failure modes, the agent can select  
 44 a small number of labelled support set candidates  
 45  $\mathcal{S}$  from a large data source and use them to update  
 46 its internal state. Finally, the updated model  $\mathcal{M}_1$   
 47 is used to predict the labels for the query set.

48 Our main contributions are the following: (i) We  
 49 propose a post-hoc method for obtaining a distribu-  
 50 tion over the cosine similarities from a pre-trained  
 51 VLM without needing architecture changes or fur-  
 52 ther training. (ii) We apply our method in active  
 53 learning and assess various scoring mechanisms for support set selection. (iii) We show on benchmark  
 54 data sets that accounting for epistemic uncertainties improves performance and that targeted candidate  
 55 selection results in further improvements.

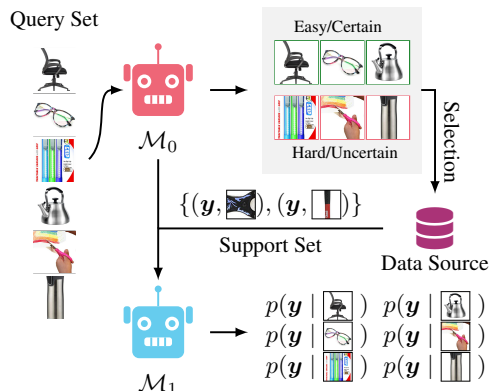


Figure 1: Illustration of the setting.

## 56 2 Methods

57 We denote vectors by bold lower-case letters (*e.g.*,  $\mathbf{x}$ ,  $\mathbf{a}$ ) and use bold upper-case letters for matrices  
 58 (*e.g.*,  $\mathbf{X}$ ,  $\mathbf{P}$ ). Further, sets are denoted in upper-case calligraphic letters (*e.g.*,  $\mathcal{D}$ ,  $\mathcal{I}$ ) and model  
 59 parameters or hyper-parameters are denoted using Greek letters (*e.g.*,  $\alpha$ ,  $\theta$ ). In particular, let  $\mathbf{x}_i \in$   
 60  $\mathbb{R}^{p_{\text{IMG}}}$  and  $\mathbf{y}_j \in \mathbb{R}^{p_{\text{TXT}}}$  denote the  $i^{\text{th}}$  image and  $j^{\text{th}}$  text description, respectively. Further, we use  
 61  $\phi : \mathbb{R}^{p_{\text{IMG}}} \rightarrow \mathbb{R}^{d_{\text{IMG}}}$  and  $\psi : \mathbb{R}^{p_{\text{TXT}}} \rightarrow \mathbb{R}^{d_{\text{TXT}}}$  to denote the image and text encoders of the VLM, where  
 62  $p_{\text{IMG}}$  and  $p_{\text{TXT}}$  denote the respective input dimensionality and  $d_{\text{IMG}}$ ,  $d_{\text{TXT}}$  is the dimensionality of the  
 63 respective feature space. The embeddings are projected into a joint space, given as  $\mathbf{g} = \mathbf{P}\phi(\mathbf{x})$  and  
 64  $\mathbf{h} = \mathbf{Q}\psi(\mathbf{y})$ , using linear projections denoted by  $\mathbf{P} \in \mathbb{R}^{d \times d_{\text{IMG}}}$  and  $\mathbf{Q} \in \mathbb{R}^{d \times d_{\text{TXT}}}$ , respectively.

65 VLMs (*e.g.*, [30]) are typically trained by minimizing the InfoNCE loss [28], which is the sum of two  
 66 cross-entropy terms, one for each relational direction—image to text (IMG  $\rightarrow$  TXT) or text to image  
 67 (IMG  $\leftarrow$  TXT). The loss is given as  $\mathcal{L}(\mathbf{X}, \mathbf{Y}) = 1/2\mathcal{L}_{\text{CE}}^{\text{IMG} \rightarrow \text{TXT}}(\mathbf{X}, \mathbf{Y}) + 1/2\mathcal{L}_{\text{CE}}^{\text{IMG} \leftarrow \text{TXT}}(\mathbf{X}, \mathbf{Y})$  with  
 68 cross-entropy loss terms defined over the cosine similarities between the embeddings, *i.e.*,

$$\mathcal{L}_{\text{CE}}^{\text{IMG} \rightarrow \text{TXT}}(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n -\log \left( \frac{\exp(\hat{\mathbf{g}}_i^\top \hat{\mathbf{h}}_i)}{\sum_{j=1}^n \exp(\hat{\mathbf{g}}_i^\top \hat{\mathbf{h}}_j)} \right), \quad (1)$$

69 where  $\hat{\mathbf{g}}$  and  $\hat{\mathbf{h}}$  are the unit-length normalized embeddings. For further details see [App. A.2](#).

70 In this work, we utilize post-hoc uncertainty estimation based on the Laplace approximation [23] to  
 71 estimate uncertainties over the model parameters. This approach has found increasing application in  
 72 contemporary deep learning (*e.g.*, [8, 20, 25]) and uses a Gaussian approximation to the posterior  
 73 distribution. Utilising a Laplace approximation allows us to induce uncertainty over the feature  
 74 embeddings of both encoders and results in a distribution over cosine similarities, which in turn  
 75 enables quantifying model uncertainties in a principled manner. **Fig. 2** illustrates the propagation of  
 76 uncertainties in our setup by estimating uncertainties over the projection matrices.

77 **Laplace approximation** One of the main computational challenges associated with the Laplace  
 78 approximation is related to the estimation of the Hessian matrix of the log joint w.r.t. the model  
 79 parameters. Since a naïve approach is computationally impractical in the case of VLMs, we chose  
 80 to estimate the Kronecker-factored Generalized Gauss–Newton (GGN) approximation [33, 24].  
 81 Moreover, we apply the Laplace approximation only for the projection matrices  $\mathbf{P}$  and  $\mathbf{Q}$  of the  
 82 image and text encoders. Hence, resulting in GGN approximations  $\text{GGN}_{\text{IMG}}$  and  $\text{GGN}_{\text{TXT}}$  given in  
 83 form of their Kronecker factors, see [App. C.1](#) for details.

84 However, naïvely applying Laplace approximations in VLMs is challenging as the contrastive loss  
 85 entangles  $\mathbf{P}$  and  $\mathbf{Q}$ , which further complicates the estimation of the Hessian. These models are

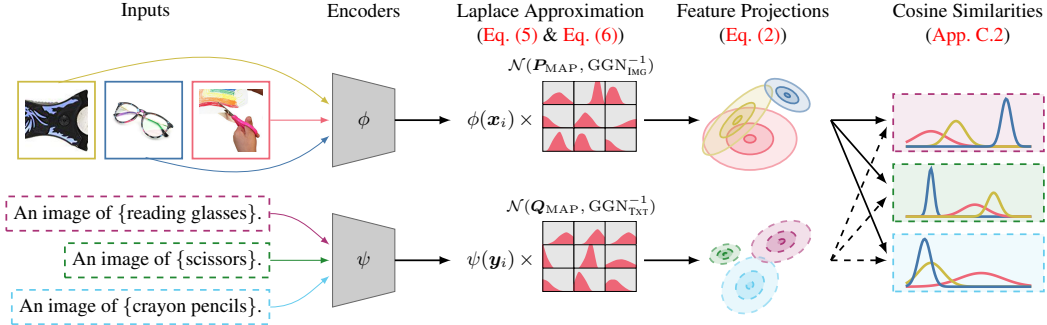


Figure 2: Illustration of uncertainty propagation in VLMs: We estimate uncertainties over the projection matrices of both encoders using a Laplace approximation, which induces distributions over the feature projections. We then approximate the distribution over cosine similarities by a Gaussian.

86 also typically trained with mini-batch sizes of around  $30k$  samples. In order to compute the GGN  
87 approximations in VLMs, we simplify the contrastive loss  $\mathcal{L}$  used for pre-training by assuming  
88 independence between  $\mathbf{P}$  and  $\mathbf{Q}$ . Specifically, we treat each of the two loss terms independent and  
89 consider only  $\mathcal{L}_{\text{CE}}^{\text{IMG} \rightarrow \text{TXT}}$  for the image encoder and  $\mathcal{L}_{\text{CE}}^{\text{IMG} \leftarrow \text{TXT}}$  for the text encoder in the Laplace  
90 approximation. Hence, dropping interactions between the image and text encoders in the Laplace  
91 approximation. Lastly, we use an incremental computation of the Kronecker factors to account for  
92 large mini-batch sizes. Further details and derivations are given in [App. C.1](#).

93 **Distribution over cosine similarities** As the Laplace approximation uses a Gaussian approximation,  
94 the feature embeddings are distributed according to another Gaussian distribution. Specifically, the  
95 distribution over embedding vectors  $\mathbf{g}$  (or  $\mathbf{h}$ ) for a datum  $\mathbf{x}$  (or  $\mathbf{y}$ ) can be expressed as follows due to  
96 linearity, *i.e.*,

$$\mathcal{N}\left(\mathbf{g}, \left(\phi(\mathbf{x})^\top \mathbf{A}_{\text{IMG}}^{-1} \phi(\mathbf{x})\right) \mathbf{B}_{\text{IMG}}^{-1}\right) \quad \text{and} \quad \mathcal{N}\left(\mathbf{h}, \left(\psi(\mathbf{y})^\top \mathbf{A}_{\text{TXT}}^{-1} \psi(\mathbf{y})\right) \mathbf{B}_{\text{TXT}}^{-1}\right), \quad (2)$$

97 where  $\mathbf{A}$  and  $\mathbf{B}$  denote the Kronecker factors of the GNN approximation of the Hessian matrix,  
98 respectively. Unfortunately, the distribution over cosine similarities is in general not Gaussian. How-  
99 ever, by assuming independence between the elements of  $\mathbf{g}$  and  $\mathbf{h}$  and in the limit of  $d \rightarrow \infty$  we can  
100 approximate the distribution over cosine similarities to be Gaussian distributed. We find this approx-  
101 imation to work well in practice, while not accurately capturing the skewness of the distributions.  
102 A detailed derivation and empirical results on the approximation quality are given in [App. C.2](#).

103 **Targeted support set selection** Let  $\mathcal{X}_{\text{test}} = \{\mathbf{x}_i^*\}$  with  $\mathbf{x}_i^* \sim p(\mathbf{x}^*)$  be a set of unseen test data  
104 (query set) with unknown class labels. We aim to find a set  $\{(\mathbf{x}_j, \mathbf{y}_j)\}_j^m$  of support candidates of  
105 cardinality  $m$  with  $\mathbf{x}_j, \mathbf{y}_j \sim p(\mathbf{x}_j, \mathbf{y}_j)$  such that we reduce uncertainty over the class labels of  $\mathcal{X}_{\text{test}}$ .  
106 To approach this problem, we target the selection process towards the predictive distribution of the  
107 query set. In particular, we propose to use a  $k$ -nearest neighbours selection in the joint space to pre-  
108 select support set candidates based on the Wasserstein distance between the distributions over image  
109 embeddings. After pre-selection, we quantify the information gain of the support set candidates either  
110 using the entropy over the predictive distribution, the expected predictive information gain (EPIG,  
111 [3]), or the BALD score [15]. Doing so adaptively targets the candidate search for the the support set  
112 towards the predictive distribution of the query set and reduces the computational complexity of the  
113 selection process. Further details on the selection process and the score functions are given in [App. D](#).

### 114 3 Experiments

115 To evaluate our approach for probabilistic active few-shot learning, we conducted experiments using  
116 pre-trained OpenCLIP models from Hugging Face [18]. We estimated the Laplace approximations  
117 of the OpenCLIP model with ViT-Base backbone and ViT-Huge backbone [10] using a randomly  
118 sampled subset from the Laion-400M data set [35]. Further details are given in [App. E](#).

119 For probabilistic active few-shot learning with VLMs we consider the task of image classification and  
120 present results on the Flowers102 [27], Food101 [5], CIFAR-100 [21], ImageNet-R [13], EuroSAT

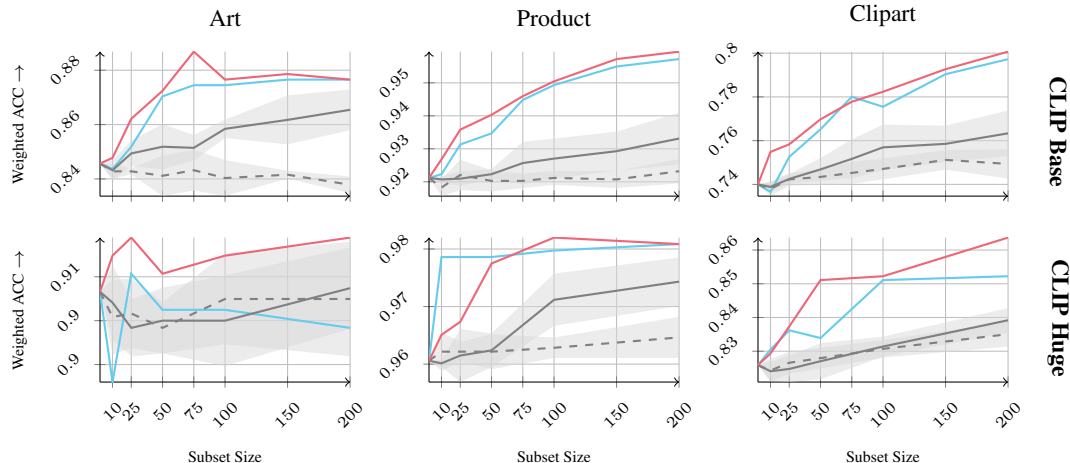


Figure 3: Results on the Office-Home data set with support set selection from all training domains. We observe that incorporating epistemic uncertainties (—) improves over entropy based targeted selection (—) in most of the cases and outperforms naïve random selection (---) and random selection with targeted support set candidates (—). Shaded regions indicate the std over 5 runs.

121 [12], and the Office-Home [39] data sets. To assess the performance of the proposal, we investi-  
 122 gated the following questions: (i) Do approaches that account for epistemic uncertainties improve  
 123 performance? (ii) What is the effect of targeting the support set candidates towards the query region?  
 124 (iii) How does the model capacity affect the performance of the proposed approach?

125 To address these questions, we performed support set selection from all training domains available in  
 126 the Office-Home data set and evaluated on the test set (query set) of each domain independently. In  
 127 Fig. 3 we compare the performance of targeted entropy-based support set selection, random selection,  
 128 random selection with targeted support set region, and the best performing (according to the valida-  
 129 tion loss) acquisition function that incorporates epistemic uncertainties. We find that incorporating  
 130 epistemic uncertainties improves the few-shot learning performance in most cases and generally  
 131 outperforms random selection. Further, we observe that targeted support set selection improves the per-  
 132 formance as indicated by the performance gap between naïve random selection and targeted random  
 133 selection and that the model capacity can have a substantial impact on the performance gains across  
 134 all approaches. A listing of the results using the negative log-predictive density are given in App. E.2.

135 **Single-domain Finetuning** In App. E.2, we show results for single-domain finetuning on standard  
 136 benchmark data sets (e.g. CIFAR-100, Imagenet-R, Flowers102, etc.) using the different support set  
 137 selection methods with the OpenCLIP model. The selection methods using the epistemic uncertainty  
 138 (BALD and EPIG) perform better or on par with the Targeted Maximum Entropy across the different  
 139 subset sizes and data sets, which demonstrates the benefits of using our proposed uncertainty estimates  
 140 for support set selection.

## 141 4 Discussion and Conclusion

142 In this work, we have introduced a probabilistic active few-shot learning approach for VLMs. Our  
 143 approach leverages a Laplace approximation to the posterior of the projection layers of the VLM  
 144 to estimate epistemic uncertainties. We have further introduced an adaptive targeted support set  
 145 candidate selection based on  $k$ -NN selection using the Wasserstein distance between the distributions  
 146 over image embeddings in the joint space. To assess the performance of probabilistic active few-shot  
 147 learning in VLMs, we have conducted two sets of experiments, one in the cross-domain setting on  
 148 the Office-Home data set and one in the single-domain setting on standard benchmark data sets. We  
 149 found that incorporating epistemic uncertainties improves the few-shot learning performance in most  
 150 cases and generally outperforms random selection. Moreover, targeting the selection process towards  
 151 the query region provides further improvements in all cases.

## References

- 152
- 153 [1] Jihwan Bang, Sumyeong Ahn, and Jae-Gil Lee. Active prompt learning in vision language mod-  
154 els. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
155 pages 27004–27014, 2024. 1
- 156 [2] Shane Barratt. A matrix gaussian distribution, 2018. 10
- 157 [3] Freddie Bickford Smith, Andreas Kirsch, Sebastian Farquhar, Yarin Gal, Adam Foster, and  
158 Tom Rainforth. Prediction-oriented Bayesian active learning. In *International Conference on*  
159 *Artificial Intelligence and Statistics*, 2023. 1, 3, 8, 14
- 160 [4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von  
161 Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the  
162 opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. 1
- 163 [5] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative  
164 components with random forests. In *European Conference on Computer Vision*, 2014. 3, 14, 15
- 165 [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
166 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
167 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1
- 168 [7] Sanghyuk Chun, Seong Joon Oh, Rafael Sampaio De Rezende, Yannis Kalantidis, and Diane  
169 Larlus. Probabilistic embeddings for cross-modal retrieval. In *Proceedings of the IEEE/CVF*  
170 *Conference on Computer Vision and Pattern Recognition*, pages 8415–8424, 2021. 8
- 171 [8] Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer,  
172 and Philipp Hennig. Laplace redux-effortless bayesian deep learning. *Advances in Neural*  
173 *Information Processing Systems*, 34:20089–20103, 2021. 1, 2, 9
- 174 [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of  
175 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,  
176 2018. 1
- 177 [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
178 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
179 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
180 recognition at scale. In *International Conference on Learning Representations*, 2021. 3
- 181 [11] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image  
182 data. In *International conference on machine learning*, pages 1183–1192. PMLR, 2017. 1, 8
- 183 [12] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel  
184 dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal*  
185 *of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.  
186 4, 14, 15
- 187 [13] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,  
188 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin  
189 Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization.  
190 In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, 2021. 3, 14,  
191 15
- 192 [14] Alex Holub, Pietro Perona, and Michael C Burl. Entropy-based active learning for object  
193 recognition. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern*  
194 *Recognition Workshops*, pages 1–8. IEEE, 2008. 1, 8
- 195 [15] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning  
196 for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011. 1, 3, 8, 13
- 197 [16] Shell Xu Hu, Da Li, Jan Stühmer, Minyoung Kim, and Timothy M Hospedales. Pushing the  
198 limits of simple pipelines for few-shot learning: External data and fine-tuning make a difference.  
199 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
200 pages 9068–9077, 2022. 1



- 201 [17] Jonas Hübötter, Bhavya Sukhija, Lenart Treven, Yarden As, and Andreas Krause. Active  
202 few-shot fine-tuning. *arXiv preprint arXiv:2402.15441*, 2024. 1
- 203 [18] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan  
204 Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi,  
205 Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. 3
- 206 [19] Yatai Ji, Junjie Wang, Yuan Gong, Lin Zhang, Yanru Zhu, Hongfa Wang, Jiaying Zhang, Tetsuya  
207 Sakai, and Yujiu Yang. Map: Multimodal uncertainty-aware vision-language pre-training model.  
208 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
209 pages 23262–23271, 2023. 8
- 210 [20] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes  
211 overconfidence in relu networks. In *International conference on machine learning*, pages  
212 5436–5446. PMLR, 2020. 2
- 213 [21] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.  
214 Technical report, University of Toronto, 2009. 3, 14, 15
- 215 [22] Hao Li, Jingkuan Song, Lianli Gao, Pengpeng Zeng, Haonan Zhang, and Gongfu Li. A differ-  
216 entiable semantic metric approximation in probabilistic embedding for cross-modal retrieval.  
217 *Advances in Neural Information Processing Systems*, 35:11934–11946, 2022. 8
- 218 [23] David JC MacKay. Information-based objective functions for active data selection. *Neural*  
219 *computation*, 4(4):590–604, 1992. 1, 2, 9
- 220 [24] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approx-  
221 imate curvature. In *International conference on machine learning*, pages 2408–2417. PMLR,  
222 2015. 2
- 223 [25] Lassi Meronen, Martin Trapp, Andrea Pilzer, Le Yang, and Arno Solin. Fixing overconfidence  
224 in dynamic neural networks. In *IEEE/CVF Winter Conference on Applications of Computer*  
225 *Vision (WACV)*, pages 2680–2690, 2024. 2, 9
- 226 [26] Kimia Nadjahi, Alain Durmus, Pierre E Jacob, Roland Badeau, and Umut Simsekli. Fast  
227 approximation of the sliced-wasserstein distance using concentration of random projections.  
228 *Advances in Neural Information Processing Systems*, 34:12411–12424, 2021. 12
- 229 [27] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large  
230 number of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*,  
231 2008. 3, 14, 15
- 232 [28] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive  
233 predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2
- 234 [29] Theodore Papamarkou, Maria Skoularidou, Konstantina Palla, Laurence Aitchison, Julyan  
235 Arbel, David Dunson, Maurizio Filippone, Vincent Fortuin, Philipp Hennig, José Miguel  
236 Hernández-Lobato, et al. Position: Bayesian deep learning is needed in the age of large-scale ai.  
237 In *International Conference on Machine Learning*, 2024. 1
- 238 [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
239 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
240 models from natural language supervision. In *International conference on machine learning*,  
241 pages 8748–8763. PMLR, 2021. 1, 2, 9, 14
- 242 [31] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine*  
243 *Learning*. The MIT Press, 2006. 11
- 244 [32] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang  
245 Chen, and Xin Wang. A survey of deep active learning. *ACM computing surveys (CSUR)*, 54  
246 (9):1–40, 2021. 8
- 247 [33] Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for  
248 neural networks. In *International conference on learning representations*, 2018. 2, 9

- 249 [34] Subhankar Roy, Martin Trapp, Andrea Pilzer, Juho Kannala, Nicu Sebe, Elisa Ricci, and Arno  
250 Solin. Uncertainty-guided source-free domain adaptation. In *European conference on computer*  
251 *vision*, pages 537–555. Springer, 2022. 9
- 252 [35] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton  
253 Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open  
254 dataset of clip-filtered 400 million image-text pairs, 2021. 3, 14
- 255 [36] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set  
256 approach. In *International Conference on Learning Representations*, 2018. 1, 8
- 257 [37] Vishaal Udandarao, Ankush Gupta, and Samuel Albanie. Sus-x: Training-free name-only  
258 transfer of vision-language models. In *2023 IEEE/CVF International Conference on Computer*  
259 *Vision (ICCV)*. IEEE, 2023. 14
- 260 [38] Uddeshya Upadhyay, Shyangopal Karthik, Massimiliano Mancini, and Zeynep Akata. Problm:  
261 Probabilistic adapter for frozen vision-language models. In *Proceedings of the IEEE/CVF*  
262 *International Conference on Computer Vision*, pages 1899–1910, 2023. 8
- 263 [39] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan.  
264 Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE*  
265 *Conference on Computer Vision and Pattern Recognition*, pages 5018–5027, 2017. 4, 14, 15
- 266 [40] Yichen Xie, Han Lu, Junchi Yan, Xiaokang Yang, Masayoshi Tomizuka, and Wei Zhan. Active  
267 finetuning: Exploiting annotation budget in the pretraining-finetuning paradigm. In *Proceedings*  
268 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23715–23724,  
269 2023. 1
- 270 [41] Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision  
271 tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 9
- 272 [42] Renrui Zhang, Xiangfei Hu, Bohao Li, Siyuan Huang, Hanqiu Deng, Yu Qiao, Peng Gao, and  
273 Hongsheng Li. Prompt, generate, then cache: Cascade of foundation models makes strong  
274 few-shot learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
275 *Recognition*, pages 15211–15222, 2023. 1

---

# Probabilistic Active Few-Shot Learning in Vision-Language Models

## Supplementary Material

---

### 276 A Related Work

#### 277 A.1 Active Learning

278 The active learning setting [32] entails an agent learning a task from an unlabelled dataset, while  
279 simultaneously determining which data points to label for maximal benefit to the target task. The  
280 learner uses an acquisition function to base its sample selection on that should quantify how beneficial  
281 (or informative) this sample will be to learn from for the target task. There exist various acquisition  
282 functions, *e.g.*, (i) entropy-based which aims to minimize the expected entropy after observing  
283 data points [14], and (ii) core-set based methods which are trained to minimize the generalization  
284 error between the unlabelled and labelled sets and use clustering for selection [36]. Uncertainty-  
285 based acquisition functions have been explored to select data points that will mostly reduce the  
286 epistemic uncertainty in the model, *e.g.*, Bayesian Active Learning by Disagreement (BALD) score  
287 [11, 15]. More recently, the expected predictive information gain (EPIG) [3] was proposed to measure  
288 the information gain in the space of predictions rather than parameters. We experiment with the  
289 mentioned uncertainty-based acquisition functions combined with our probabilistic embeddings for  
290 targeted data selection in VLM finetuning.

#### 291 A.2 Probabilistic Vision-Language Models

292 Several works are aiming to extend VLMs to produce predictive uncertainty estimates for various  
293 downstream tasks, *e.g.*, cross-modal retrieval [7, 22] and visual-question answering [19]. These  
294 methods learn probabilistic embeddings on each modality by estimating probability distributions  
295 from the network. However, this approach requires training the networks from scratch, which limits  
296 their applicability to pretrained VLMs (*e.g.* CLIP). To this end, Upadhyay et al. [38] proposed a  
297 post-hoc method called ProbVLM that learns probabilistic embeddings from finetuned adapters on  
298 a frozen VLM backbone. Similar to this work, they also apply their method to the active learning  
299 task and use the uncertainty estimates for selecting informative subsets of training data for finetuning.  
300 However, ProbVLM requires finetuning the probabilistic embeddings on a proxy task, while our  
301 method can be applied directly on the pretrained model.

### 302 B Preliminaries

303 This section provides a brief overview of the background concepts relevant to this work.

#### 304 B.1 Vision-Language Models

305 In this work, we consider vision-language models (VLM) learned using the contrastive learning  
306 objective known as InfoNCE. In particular, let  $\mathbf{x}_i \in \mathbb{R}^{p_{\text{IMG}}}$  and  $\mathbf{y}_j \in \mathbb{R}^{p_{\text{TXT}}}$  denote the  $i$ th image and  
307  $j$ th text description, respectively. Further, we use  $\phi : \mathbb{R}^{p_{\text{IMG}}} \rightarrow \mathbb{R}^{d_{\text{IMG}}}$  and  $\psi : \mathbb{R}^{p_{\text{TXT}}} \rightarrow \mathbb{R}^{d_{\text{TXT}}}$  to  
308 denote the image and text encoders of the VLM, where  $p_{\text{IMG}}$  and  $p_{\text{TXT}}$  denote the respective input  
309 dimensionalities and  $d_{\text{IMG}}$ ,  $d_{\text{TXT}}$  is the dimensionality of the respective feature space.



310 To project the embeddings into a joint space, we assume a linear projection layer for both the image  
 311 and the text encoder denoted by  $\mathbf{P} \in \mathbb{R}^{d \times d_{\text{img}}}$  and  $\mathbf{Q} \in \mathbb{R}^{d \times d_{\text{txt}}}$ , respectively. The embeddings in  
 312 the joint space are then given as  $\mathbf{g}_i = \mathbf{P}\phi(\mathbf{x}_i)$  and  $\mathbf{h}_j = \mathbf{Q}\psi(\mathbf{y}_j)$  and we use hat notation to denote  
 313 the unit-length normalized embeddings, e.g.,  $\hat{\mathbf{g}}_i = \frac{\mathbf{P}\phi(\mathbf{x}_i)}{\|\mathbf{P}\phi(\mathbf{x}_i)\|}$ .

314 VLM models (e.g., [30]) are typically trained by minimizing the InfoNCE loss, which is given  
 315 as the sum of two cross-entropy terms, one for each relational direction – image to text (IMG  $\rightarrow$   
 316 TXT) or text to image (IMG  $\leftarrow$  TXT). Specifically, the InfoNCE loss is given as  $\mathcal{L}(\mathbf{X}, \mathbf{Y}) =$   
 317  $1/2\mathcal{L}_{\text{CE}}^{\text{IMG} \rightarrow \text{TXT}}(\mathbf{X}, \mathbf{Y}) + 1/2\mathcal{L}_{\text{CE}}^{\text{IMG} \leftarrow \text{TXT}}(\mathbf{X}, \mathbf{Y})$  with cross-entropy loss terms given as:

$$\mathcal{L}_{\text{CE}}^{\text{IMG} \rightarrow \text{TXT}}(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n -\log \left( \frac{\exp(\hat{\mathbf{g}}_i^\top \hat{\mathbf{h}}_i)}{\sum_{j=1}^n \exp(\hat{\mathbf{g}}_i^\top \hat{\mathbf{h}}_j)} \right) \quad (3)$$

$$\mathcal{L}_{\text{CE}}^{\text{IMG} \leftarrow \text{TXT}}(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n -\log \left( \frac{\exp(\hat{\mathbf{h}}_i^\top \hat{\mathbf{g}}_i)}{\sum_{j=1}^n \exp(\hat{\mathbf{h}}_i^\top \hat{\mathbf{g}}_j)} \right). \quad (4)$$

318 For further details we refer the reader to [30, 41]

## 319 B.2 Bayesian Deep Learning

320 We will briefly review concepts on Bayesian deep learning relevant to this work. Given a dataset  $\mathcal{D} =$   
 321  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  and a probabilistic models with likelihood function  $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$  and prior distribution  
 322  $p(\boldsymbol{\theta})$ , we aim to estimate the posterior distribution  $p(\boldsymbol{\theta}|\mathcal{D})$  of the model parameters  $\boldsymbol{\theta}$  given the  
 323 training data  $\mathcal{D}$ . In the context of deep learning, exact inference of the posterior distribution is at  
 324 least NP-hard in most settings and only becomes tractable if  $p(\boldsymbol{\theta}|\mathcal{D})$  constitute sufficient structure  
 325 [23]. Henceforth, we consider approximate Bayesian inference using the Laplace approximation [23]  
 326 in this work, which has gained increasing popularity in the community (e.g., [33, 8, 25, 34]) as a  
 327 post-hoc techniques to estimate epistemic uncertainties.

328 The Laplace approximation uses a second-order Taylor expansion of the log-joint around the  
 329 maximum-a-posteriori (MAP) estimate  $\boldsymbol{\theta}_{\text{MAP}}$ . The resulting distribution is then approximated  
 330 with a un-normalised Gaussian density, which in turn results in an approximate posterior distribu-  
 331 tion given by a Gaussian distribution located at the MAP estimate, i.e.,  $p(\boldsymbol{\theta}|\mathcal{D}) \approx \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_{\text{MAP}}, \boldsymbol{\Sigma})$ .  
 332 Resulting from the Taylor expansion, the covariance is given by the inverse Hessian at the MAP,  
 333 i.e.,  $\boldsymbol{\Sigma} = (-\nabla^2 \log p(\boldsymbol{\theta}, \mathcal{D})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{\text{MAP}}})^{-1}$ . Predictions are then made based on the posterior predic-  
 334 tive distribution  $p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta}$ , which is typically performed by Monte Carlo  
 335 sampling in case of non-linear likelihoods functions, e.g., classification settings. We refer to ... for a  
 336 detailed review the topic.

## 337 C Derivations

338 This section provides detailed derivations of the equations presented in the main text.

### 339 C.1 Laplace Approximation

340 Then the GGN approximation of the Hessian matrices are given as:

$$\text{GGN}_{\text{IMG}} \approx \underbrace{\left[ \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi(\mathbf{x}_i)\phi(\mathbf{x}_i)^\top \right]}_{=\mathbf{A}_{\text{IMG}}} \otimes \underbrace{\left[ \frac{1}{\sqrt{n}} \sum_{i=1}^n J_{\text{IMG}}(\mathbf{x}_i)^\top \boldsymbol{\Lambda}_{\text{IMG}} J_{\text{IMG}}(\mathbf{x}_i) \right]}_{=\mathbf{B}_{\text{IMG}}} \quad (5)$$

341 where  $J_{\text{IMG}}(\mathbf{x}_i) = \frac{\partial \hat{\mathbf{g}}_i^\top \hat{\mathbf{H}}}{\partial \mathbf{g}_i}$  and

$$\text{GGN}_{\text{TXT}} \approx \underbrace{\left[ \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi(\mathbf{y}_i)\psi(\mathbf{y}_i)^\top \right]}_{=\mathbf{A}_{\text{TXT}}} \otimes \underbrace{\left[ \frac{1}{\sqrt{n}} \sum_{i=1}^n J_{\text{TXT}}(\mathbf{x}_i)^\top \boldsymbol{\Lambda}_{\text{TXT}} J_{\text{TXT}}(\mathbf{x}_i) \right]}_{=\mathbf{B}_{\text{TXT}}}. \quad (6)$$

342 We further incorporate the prior precision  $\lambda$  into the GGN approximation by adding the prior precision  
 343 to the diagonal of the GGN hessian.

$$\text{GGN}_{\text{IMG}} \approx \tau (\mathbf{A}_{\text{IMG}} \otimes \mathbf{B}_{\text{IMG}}) + \lambda \mathbf{I} \quad (7)$$

$$\approx \left( \sqrt{\tau} \mathbf{A}_{\text{IMG}} + \sqrt{\lambda} \mathbf{I} \right) \otimes \left( \sqrt{\tau} \mathbf{B}_{\text{IMG}} + \sqrt{\lambda} \mathbf{I} \right) \quad (8)$$

344 In our experiments, we set  $\tau = 0.75$  for the ViT-Base model and  $\tau = 0.3$  for the ViT-Huge model  
 345 and obtain the prior precision  $\lambda$  through marginal likelihood maximization.

### 346 C.1.1 Obtaining the Posterior Predictive Distribution

347 For conciseness, we denote the posterior precision matrices associated with the image encoder  
 348 as  $\mathbf{A}_{\text{IMG}}$  and  $\mathbf{B}_{\text{IMG}}$ . We have obtained the posterior distribution over the image projection matrix  
 349  $\mathbf{P}$  represented as  $\mathcal{N}(\text{vec}(\mathbf{P}); \text{vec}(\mathbf{P}_{\text{MAP}}), \text{GGN}_{\text{IMG}}^{-1})$ . Given that  $\text{GGN}_{\text{IMG}}^{-1}$  is formulated using the  
 350 Kronecker product of the inverses of these matrices, i.e.,  $\mathbf{A}_{\text{IMG}}^{-1} \otimes \mathbf{B}_{\text{IMG}}^{-1}$ , we proceed to express  
 351 the posterior predictive distribution as a matrix normal distribution  $\mathcal{MN}(\mathbf{P}; \mathbf{P}_{\text{MAP}}, \mathbf{B}_{\text{IMG}}^{-1}, \mathbf{A}_{\text{IMG}}^{-1})$  as  
 352 referenced in [2]:

$$\mathbf{P} \sim \mathcal{MN}(\mathbf{P}_{\text{MAP}}, \mathbf{B}_{\text{IMG}}^{-1}, \mathbf{A}_{\text{IMG}}^{-1}) \quad (9)$$

$$\implies \mathbf{g} = \mathbf{P}\phi(\mathbf{x}) \sim \mathcal{MN}(\mathbf{P}_{\text{MAP}}\phi(\mathbf{x}), \mathbf{B}_{\text{IMG}}^{-1}, \phi(\mathbf{x})^\top \mathbf{A}_{\text{IMG}}^{-1} \phi(\mathbf{x})) \quad (10)$$

$$\implies \mathbf{g} \sim \mathcal{N}(\mathbf{P}_{\text{MAP}}\mathbf{a}, (\phi(\mathbf{x})^\top \mathbf{A}_{\text{IMG}}^{-1} \phi(\mathbf{x})) \mathbf{B}_{\text{IMG}}^{-1}) \quad (11)$$

### 353 C.1.2 Online Laplace Approximation

354 For the EPIG score, we update our Laplace approximation online after each data point is added to  
 355 the support set. Given the current Laplace approximation of the posterior over the image projection  
 356 matrix  $\mathbf{P}$  we update the posterior distribution as follows:

$$\mathbf{P}_{t+1} = \mathbf{P}_t - \gamma \nabla_{\mathbf{P}} \mathcal{L}_{\text{CE}}^{\text{IMG} \rightarrow \text{TXT}}(\mathbf{x}^*, \mathbf{Y}) \quad (12)$$

$$\mathbf{A}_{\text{IMG}, t+1} = \mathbf{A}_{\text{IMG}, t} + \beta \phi(\mathbf{x}^*) \phi(\mathbf{x}^*)^\top \quad (13)$$

$$\mathbf{B}_{\text{IMG}, t+1} = \mathbf{B}_{\text{IMG}, t} + \beta J_{\text{IMG}}(\mathbf{x}^*)^\top \mathbf{\Lambda}_{\text{IMG}} J_{\text{IMG}}(\mathbf{x}^*) \quad (14)$$

357 From the updated  $\mathbf{A}_{\text{IMG}, t+1}$  and  $\mathbf{B}_{\text{IMG}, t+1}$  we obtain the updated GGN approximation of the Hessian  
 358 matrix:

$$\text{GGN}_{\text{IMG}, t+1} \approx \left( \sqrt{\tau} \mathbf{A}_{\text{IMG}, t+1} + \sqrt{\lambda} \mathbf{I} \right) \otimes \left( \sqrt{\tau} \mathbf{B}_{\text{IMG}, t+1} + \sqrt{\lambda} \mathbf{I} \right) \quad (15)$$

359 After each update, we optimize for the prior precision  $\lambda$  by maximizing the marginal likelihood. For  
 360 our experiments, we set the learning rates  $\gamma = 10^{-3}$  and  $\beta = 1$ .

### 361 C.1.3 Jacobians for the GGN Approximation

362 In the following we derive the Jacobians  $J_{\text{IMG}}(\mathbf{x}_i)$  and  $J_{\text{TXT}}(\mathbf{y}_i)$  used in the Kronecker-factored  
 363 Generalized Gauss-Newton (GGN) approximation of the Hessian matrices.

$$J_{\text{IMG}}(\mathbf{x}_i)^\top = \frac{\partial \hat{\mathbf{H}}^\top \hat{\mathbf{g}}_i}{\partial \mathbf{g}_i} = \hat{\mathbf{H}}^\top \frac{\partial}{\partial \mathbf{g}_i} \frac{\mathbf{g}_i}{\|\mathbf{g}_i\|} = \hat{\mathbf{H}}^\top \frac{\|\mathbf{g}_i\| - \mathbf{g}_i \frac{\partial \|\mathbf{g}_i\|}{\partial \mathbf{g}_i}}{\|\mathbf{g}_i\|^2} = \hat{\mathbf{H}}^\top \frac{\|\mathbf{g}_i\| - \frac{\mathbf{g}_i \mathbf{g}_i^\top}{\|\mathbf{g}_i\|}}{\|\mathbf{g}_i\|^2} \quad (16)$$

$$= \hat{\mathbf{H}}^\top \left( \frac{\mathbf{1}}{\|\mathbf{g}_i\|} - \frac{\mathbf{g}_i \mathbf{g}_i^\top}{\|\mathbf{g}_i\|^3} \right) \quad (17)$$

364 Analogously, we obtain the Jacobian for the text encoder as:

$$J_{\text{TXT}}(\mathbf{y}_i)^\top = \hat{\mathbf{G}}^\top \left( \frac{\mathbf{1}}{\|\mathbf{h}_i\|} - \frac{\mathbf{h}_i \mathbf{h}_i^\top}{\|\mathbf{h}_i\|^3} \right) \quad (18)$$

### 365 C.1.4 Likelihood Hessian for the GGN Approximation

366 The zero-shot classifier induced by CLIP computes unnormalized logits for each class  $c$ , represented  
 367 by  $\hat{\mathbf{g}}_i^\top \hat{\mathbf{h}}_c =: f_c$ . By applying the softmax function, we calculate the probabilities for each class  $c$  as  
 368  $\pi_c = \frac{\exp(f_c)}{\sum_{c'} \exp(f_{c'})}$ . The likelihood Hessian of the cross-entropy loss for this classifier is represented  
 369 by:

$$\Lambda_{\text{IMG}} = \text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi} \boldsymbol{\pi}^\top \quad (19)$$

370 Similarly, the likelihood Hessian for the text encoder follows analogous principles in the text-to-image  
 371 direction. For a more detailed derivation of the likelihood Hessian, we refer to [31]. Rearranging  
 372 terms in the analytical expression for  $J_{\text{IMG}}^\top \Lambda_{\text{IMG}} J_{\text{IMG}}$  facilitates space-efficient computation of the  
 373 GGN approximation.

### 374 C.2 Distribution over Cosine Similarities

375 From the Laplace approximation on the projection layers, we have the distribution over the embed-  
 376 dings  $\hat{\mathbf{g}}$  and  $\hat{\mathbf{h}}_i$ . Specifically, we have  $\mathbb{E}[\hat{\mathbf{g}}]$ ,  $\mathbb{E}[\hat{\mathbf{h}}_i]$ ,  $\text{Var}[\hat{\mathbf{g}}] = \text{Cov}[\hat{\mathbf{g}}, \hat{\mathbf{g}}]$ , and  $\text{Var}[\hat{\mathbf{h}}_i] = \text{Cov}[\hat{\mathbf{h}}_i, \hat{\mathbf{h}}_i]$ .  
 377 Further, we assume that the text- and image embeddings are independent, i.e.,  $\hat{\mathbf{g}} \perp \hat{\mathbf{h}}_i$  for all  $i$ . Then,  
 378 the covariance of the cosine similarities is given by:

$$\text{Cov} \left[ \hat{\mathbf{g}}^\top \hat{\mathbf{h}}_i, \hat{\mathbf{g}}^\top \hat{\mathbf{h}}_j \right] = \text{Cov} \left[ \sum_{m=1}^d (\hat{g})_m (\hat{h}_i)_m, \sum_{n=1}^d (\hat{g})_n (\hat{h}_j)_n \right] \quad (20)$$

$$= \sum_{m=1}^d \sum_{n=1}^d \text{Cov} \left[ (\hat{g})_m (\hat{h}_i)_m, (\hat{g})_n (\hat{h}_j)_n \right] \quad (21)$$

$$= \sum_{m=1}^d \sum_{n=1}^d \mathbb{E} \left[ (\hat{g})_m (\hat{h}_i)_m (\hat{g})_n (\hat{h}_j)_n \right] - \mathbb{E} \left[ (\hat{g})_m (\hat{h}_i)_m \right] \mathbb{E} \left[ (\hat{g})_n (\hat{h}_j)_n \right] \quad (22)$$

$$\stackrel{\hat{\mathbf{g}} \perp \hat{\mathbf{h}}_i}{=} \sum_{m=1}^d \sum_{n=1}^d \mathbb{E} [(\hat{g})_m (\hat{g})_n] \mathbb{E} [(\hat{h}_i)_m (\hat{h}_j)_n] - \mathbb{E} [(\hat{g})_m] \mathbb{E} [(\hat{g})_n] \mathbb{E} [(\hat{h}_i)_m] \mathbb{E} [(\hat{h}_j)_n] \quad (23)$$

$$= \sum_{m=1}^d \sum_{n=1}^d \mathbb{E} [(\hat{g})_m] \mathbb{E} [(\hat{g})_n] \text{Cov} \left[ (\hat{h}_i)_m, (\hat{h}_j)_n \right] + \text{Cov} [(\hat{g})_m, (\hat{g})_n] \mathbb{E} [(\hat{h}_i)_m] \mathbb{E} [(\hat{h}_j)_n] \\ + \text{Cov} [(\hat{g})_m, (\hat{g})_n] \text{Cov} \left[ (\hat{h}_i)_m, (\hat{h}_j)_n \right] \quad (24)$$

$$= \mathbb{E} \left[ \hat{\mathbf{g}}^\top \right] \text{Cov} \left[ \hat{\mathbf{h}}_i, \hat{\mathbf{h}}_j \right] \mathbb{E} [\hat{\mathbf{g}}] + \mathbb{E} \left[ \hat{\mathbf{h}}_i^\top \right] \text{Cov} [\hat{\mathbf{g}}, \hat{\mathbf{g}}] \mathbb{E} [\hat{\mathbf{h}}_j] + \text{tr} \left[ \text{Cov} [\hat{\mathbf{g}}, \hat{\mathbf{g}}] \text{Cov} \left[ \hat{\mathbf{h}}_i, \hat{\mathbf{h}}_j \right] \right] \quad (25)$$

379 Further, assuming  $\hat{\mathbf{h}}_i \perp \hat{\mathbf{h}}_j$  with  $\forall i \neq j$ , we have:

$$\text{Cov} \left[ \hat{\mathbf{g}}^\top \hat{\mathbf{h}}_i, \hat{\mathbf{g}}^\top \hat{\mathbf{h}}_j \right] = \mathbb{E} \left[ \hat{\mathbf{h}}_i^\top \right] \text{Var} [\hat{\mathbf{g}}] \mathbb{E} \left[ \hat{\mathbf{h}}_j \right] + \mathbb{1}_{[i=j]} \mathbb{E} \left[ \hat{\mathbf{g}}^\top \right] \text{Var} \left[ \hat{\mathbf{h}}_i \right] \mathbb{E} [\hat{\mathbf{g}}] \\ + \mathbb{1}_{[i=j]} \text{tr} \left[ \text{Cov} [\hat{\mathbf{g}}, \hat{\mathbf{g}}] \text{Var} \left[ \hat{\mathbf{h}}_i \right] \right], \quad (26)$$

380 which concludes the derivation.

381 **Fig. 4** provides a simulation result illustrating the error induced through the above approximation.  
 382 We find this approximation to work well in practice, but fail to capture potential skewness of the  
 383 distributions. In the future we aim to explore alternative approximations for the distribution over  
 384 cosine similarities.

## 385 D Details on Support Set Selection

386 This section provides further details on the support set selection strategies used in this work.

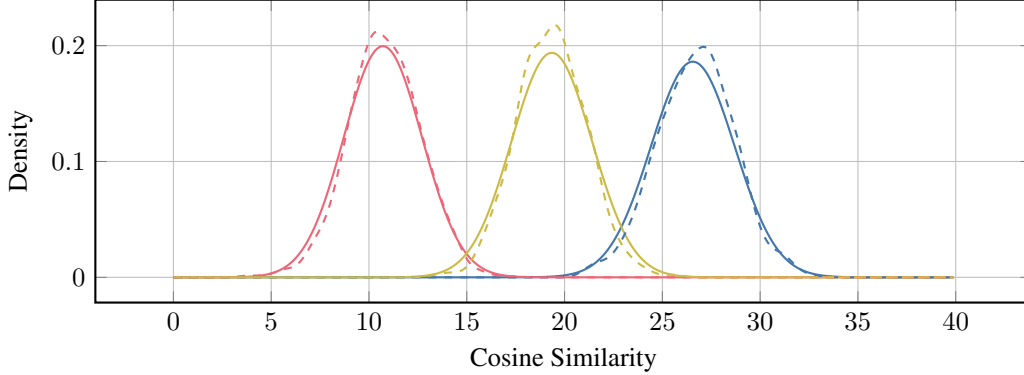


Figure 4: Comparison between Monte-Carlo estimates of the distribution over covariance similarities (---) and our analytic Gaussian approximation (—) for three cases estimated on two random samples from CIFAR-100. Blue lines indicate the distribution over cosine similarities for a **matching pair** (image with corresponding correct label), the distribution over cosine similarities for pairs with **wrong label** and **wrong image** are depicted in red and yellow respectively.

### 387 D.1 k-Nearest Selection

388 Active learning acquisition functions like Maximum Entropy Selection or BALD are often applied  
 389 to the training set, lacking consideration of the target distribution and resulting in unrepresentative  
 390 selections. To address this, we propose the following heuristic: we greedily acquire a maximally  
 391 informative intermediate set  $\mathcal{S}^* \subseteq \mathcal{X}_{\text{test}}$  from the test set, followed by selecting training data points in  
 392 the vicinity of the intermediate set  $\mathcal{S}^*$ . In case of deterministic embeddings one can use the cosine  
 393 similarity or Euclidean distance for this purpose. However, as the embeddings are probabilistic in our  
 394 setting, a point-wise comparison is not possible. Henceforth, we propose to compute the Wasserstein  
 395 distance between the distributions of the embeddings of the test set and the training set, and select  
 396 the training samples with minimal Wasserstein distance to the test set. For multivariate Gaussian  
 397 distributions, the Wasserstein distance can be computed in closed form and is given as:

$$W_2^2(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2 + \text{tr}\left(\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2 - 2(\boldsymbol{\Sigma}_1^{1/2}\boldsymbol{\Sigma}_2\boldsymbol{\Sigma}_1^{1/2})^{1/2}\right) \quad (27)$$

398 where  $\|\cdot\|_2$  denotes the Euclidean norm,  $\text{Tr}(\cdot)$  is the trace operator, and  $\boldsymbol{\Sigma}^{1/2}$  is the matrix square  
 399 root of  $\boldsymbol{\Sigma}$ . As computing the Wasserstein distance exactly is computationally and memory intensive,  
 400 we approximate it by ignoring the correlation terms between the dimensions of the embeddings  
 401 resulting in the Wasserstein distance for univariate Gaussian distributions. We aim to explore more  
 402 sophisticated approximations, e.g., using the sliced Wasserstein distance [26], in future work. Based  
 403 on this distance, we select the training samples closest to the test set in the joint embedding space,  
 404 resulting in:

$$\mathcal{S} = \bigcup_{\mathbf{g}^* \in \mathcal{S}^*} \mathcal{N}_k(\mathbf{g}^*, \mathcal{X}_{\text{train}}), \quad (28)$$

405 with  $\mathcal{N}_k(\mathbf{g}^*, \mathcal{X}_{\text{train}})$  denoting the set of  $k$ -nearest neighbours of  $\mathbf{g}^*$  in the training set  $\mathcal{X}_{\text{train}}$  according  
 406 to the Wasserstein distance over the distributions of the normalized image embeddings. To ensure that  
 407 we select  $k$  distinct training samples for each test sample, we perform an iterative search in which we  
 408 discard the already selected training samples and iteratively increase the search radius until  $k$  distinct  
 409 samples are found. This process is illustrated in Fig. 5.

### 410 D.2 Acquisition Functions

411 **Naive Random** For the *naive random* acquisition function, we randomly sample  $m$  data points from  
 412 the train set  $\mathcal{X}_{\text{train}}$  to form the support set  $\mathcal{S}_{\text{ID}}$ .

413 **Targeted Random** For the *targeted random* acquisition function, we randomly sample  $m$  data  
 414 points from the test set  $\mathcal{X}_{\text{test}}$  to form an intermediate support set  $\mathcal{S}^*$ . According to App. D.1, we then

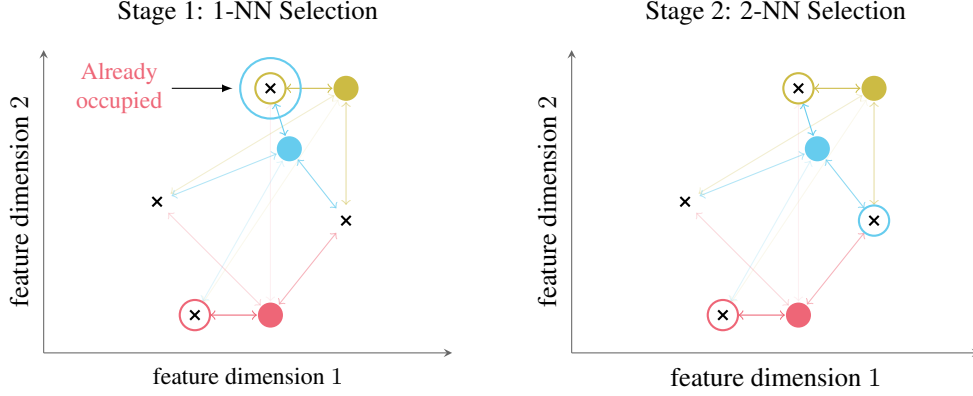


Figure 5: Illustration of the nearest neighbour based support set selection for adaptive targeted selection. The circles  $\bullet$  show test data points with uncertainty scores depicted through their colours: **high**, **medium**, **low**. For each test datum we find the  $k = 1$  nearest neighbour from the support set candidates  $\times$ . If the  $k = 1$  nearest neighbour is already selected, we increase  $k$  for those with occupied neighbours and choose the second nearest neighbour, i.e.,  $k = 2$ . This recursion continues until every test datum has a selected support set candidate. The selected candidates are shown by coloured circles. Note that in case of the **blue** test datum, the closest support set candidate has already been chosen by the **yellow** and hence the second closes candidate is selected in the second stage.

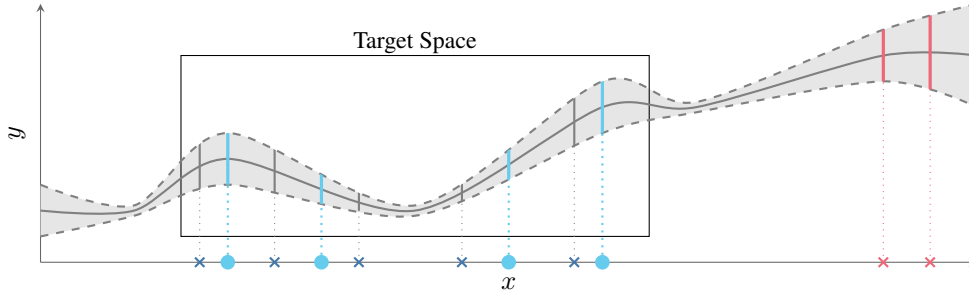


Figure 6: Illustration of targeted support set selection. We aim to select an **informative** support set that reduces the uncertainty over the predictions on the query set  $\bullet$ . Only focusing on the epistemic uncertainties would not lead to a good selection as we would select **uninformative** support set candidates  $\times$  with high epistemic uncertainty. Hence, we target the selection process.

415 select the nearest neighbours to  $\mathcal{S}^*$  from the training set  $\mathcal{X}_{\text{train}}$  based on the cosine similarity of the  
 416 normalized image embeddings to form the support set  $\mathcal{S}_{\text{t-ID}}$ .

417 **Targeted Maximum Entropy** For the *entropy* acquisition function, we compute the predictive  
 418 entropy  $\mathcal{H}(y_i^* | \mathbf{x}_i^*)$  for each data point  $\mathbf{x}_i^* \in \mathcal{X}_{\text{test}}$  and select the  $m$  data points with the highest  
 419 entropy. We use the predictive entropy on the MAP estimate of the model parameters to estimate the  
 420 predictive entropy of the model:

$$\mathcal{H}(y | \mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}) = - \sum_{c=1}^C p(y = c | \mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}) \log p(y = c | \mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}) \quad (29)$$

421 According to [App. D.1](#), we then select the most similar data points from  $\mathcal{X}_{\text{train}}$  to form the support set  
 422  $\mathcal{S}_{\text{t-entropy}}$ .

423 **BALD** We compute the BALD score [15] for each data point in  $\mathcal{X}_{\text{train}}$  and select the  $m$  data points  
 424 with the highest score. The score is approximated using nested Monte Carlo sampling as in [15].

$$\text{BALD}(\mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x})} [\mathcal{H}(p(\boldsymbol{\theta})) - \mathcal{H}(p(\boldsymbol{\theta} | \mathbf{x}, y))] \quad (30)$$

$$= \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} [\mathcal{H}(p(y | \mathbf{x}, \boldsymbol{\theta})) - \mathcal{H}(p(y | \mathbf{x}, \mathcal{D}))] \quad (31)$$

Table 1: Data specifications for finetuning data sets with the number of classes  $c$ , training set size  $n_{\text{train}}$ , validation set size  $n_{\text{val}}$ , and test set size  $n_{\text{test}}$ .

Dataset	$c$	$n_{\text{train}}$	$n_{\text{val}}$	$n_{\text{test}}$
Flowers [27]	102	1020	1020	6100
Food-101 [5]	101	75750	15150	25250
CIFAR-10/100 [21]	10/100	50000	10000	10000
ImageNet-R [13]	200	22500	4500	7500
ImageNet1k (subset classes)	200	11168	2792	2298
EuroSAT [12]	10	13500	8100	5400
Office-Home (clipart) [39]	65	2793	699	873
Office-Home (product) [39]	65	2840	711	888
Office-Home (real world) [39]	65	2788	697	872

425 **Targeted BALD** We compute the BALD score (Eq. (31)) for each data point  $\mathbf{x}_i^* \in \mathcal{X}_{\text{test}}$  and select  
 426 the  $m$  data points with the highest score. According to App. D.1, we then select the most similar data  
 427 points from  $\mathcal{X}_{\text{train}}$  to form the support set  $\mathcal{S}_{\text{t-BALD}}$ .

428 **EPIG** The Expected Predictive Information Gain (EPIG) score [3] calculates the expected mutual  
 429 information between the model parameters and the predictive distribution resulting from the acqui-  
 430 sition of a training data point. This method is specifically designed to target relevant information,  
 431 eliminating the need for a k-nearest neighbor search typically used in other acquisition functions.  
 432 The EPIG score is given by

$$\text{EPIG}(\mathbf{x}) = \mathbb{E}_{p_*(\mathbf{x}^*)p_\phi(y|\mathbf{x})} (\mathcal{H}(p_\phi(y^* | \mathbf{x}^*)) - \mathcal{H}(p_\phi(y^* | \mathbf{x}^*, x, y))) \quad (32)$$

$$= \mathbb{E}_{p_*(\mathbf{x}^*)} [\text{D}_{\text{KL}}(p_\phi(y, y^* | \mathbf{x}, \mathbf{x}^*) \| p_\phi(y | \mathbf{x})p_\phi(y^* | \mathbf{x}^*))] \quad (33)$$

$$= \mathbb{E}_{p_*(\mathbf{x}^*)} \left[ \sum_{y \in \mathcal{Y}} \sum_{y^* \in \mathcal{Y}} p_\phi(y, y^* | \mathbf{x}, \mathbf{x}^*) \log \frac{p_\phi(y, y^* | \mathbf{x}, \mathbf{x}^*)}{p_\phi(y | \mathbf{x})p_\phi(y^* | \mathbf{x}^*)} \right] \quad (34)$$

433 and can be approximated using Monte Carlo sampling. For the EPIG selection we perform online  
 434 updates to the model weights using the online Laplace as described in App. C.1.2.

## 435 E Experiments

### 436 E.1 Experimental Details

437 In our experiments we used the a pre-trained CLIP model [30] as the vision-language model with a  
 438 ViT-Base and ViT-Huge backbone. We estimated the Hessians separately for the CLIP image and text  
 439 encoders using the pre-training dataset Laion-400M [35]. For this estimation, we randomly sampled  
 440 a subset of 3 million data points for the CLIP model with a ViT-Base backbone and 0.5 million  
 441 data points for the CLIP model with a ViT-Huge backbone. The pre-training dataset was filtered to  
 442 exclude NSFW content. For the Laplace approximation, we used the GGN approximation of the  
 443 Hessian matrices as described in Sec. 2 and estimated the covariance matrices  $\mathbf{A}$  and  $\mathbf{B}$  for the image  
 444 and text encoders. We use a grid search to find the Hessian scaling  $\tau$  and learned the optimal prior  
 445 precision by maximizing the marginal likelihood of the training data. The grid for the Hessian scale  
 446 was set to  $\tau \in \{0.3, 0.35, 0.4, 0.45, 0.5\}$  for the ViT-Base model and  $\tau \in \{0.6, 0.65, 0.7, 0.75, 0.8\}$   
 447 for the ViT-Huge model.

448 For the *Office-Home* and *Flowers* data sets, we used the pre-defined splits provided by the original  
 449 authors. For *EuroSAT*, we utilized the splits provided by [37]. For *ImageNet-R*, we divided the  
 450 provided training set into a training and validation set with a validation ratio of 0.25 and used the  
 451 provided test set as is. Similarly, for the *Food* and *CIFAR-10/100* data sets, we split the training set  
 452 into a training and validation set with a validation ratio of 0.2 and used the provided test set without  
 453 modifications.

454 In our experiments, we compare the performance of the proposed EPIG acquisition function to  
 455 various baseline acquisition functions: **Naive Random**, **Targeted Random**, **Targeted Maximum**  
 456 **Entropy**, **Targeted BALD**, **EPIG**, and **Targeted EPIG**.



457 **Finetuning Settings** For the finetuning, we trained we create support sets of size  $m \in$   
458  $\{10, 25, 50, 75, 100, 150, 200, 500, 1000\}$  using the cross-entropy loss for 100 epochs. For eval-  
459 uation, we report performance of best checkpoint according to validation loss.

460 **Data sets** We experiment with the following data sets: Flowers102 [27], Food101 [5], CIFAR-10/100  
461 [21], ImageNet-R [13], EuroSAT [12] and Office-Home [39]. Table 1 shows the data split sizes and  
462 number of classes for each dataset.

463 **Metrics** We evaluate each method by measuring the class-weighted accuracy (ACC) on the test  
464 set that weights the accuracy based on the number of samples per class. Moreover, we use the  
465 negative log predictive density (NLPD) to assess the quality of the uncertainty estimates. We report  
466 the performance of each finetuned method at the epoch with the lowest validation loss.

## 467 E.2 Additional Results

468 This section provides additional experimental results and ablations of the proposed method.

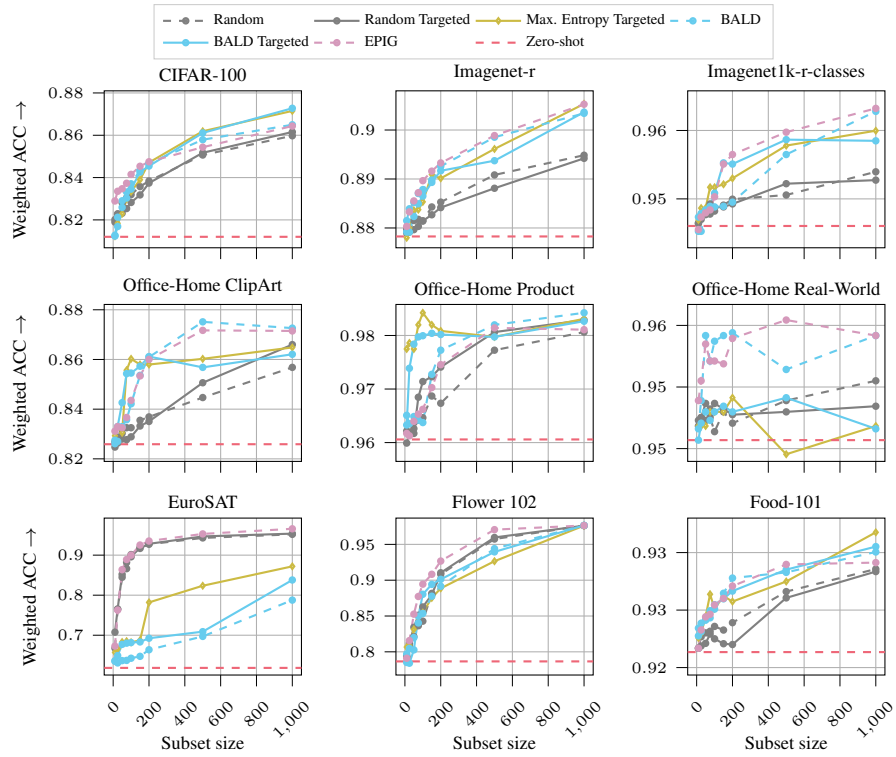
469 **Cross-domain Finetuning Results** Fig. 9 show additional results for the cross-domain setting on  
470 the Office-Home data set for both the base and huge variants of the OpenCLIP model.

471 **Single-domain Finetuning Results** Fig. 7 and Fig. 8 show the results for single-domain finetuning  
472 with support set selection using the huge and base variants of the OpenCLIP model, respectively. We  
473 also show the zero-shot performances from the pretrained CLIP models without any finetuning on the  
474 target task (Zero-shot). Note that we only show the performance for EPIG without targeted support  
475 set selection, as we noticed that EPIG performs competitively against the other selection methods in  
476 this single-domain finetuning setting.

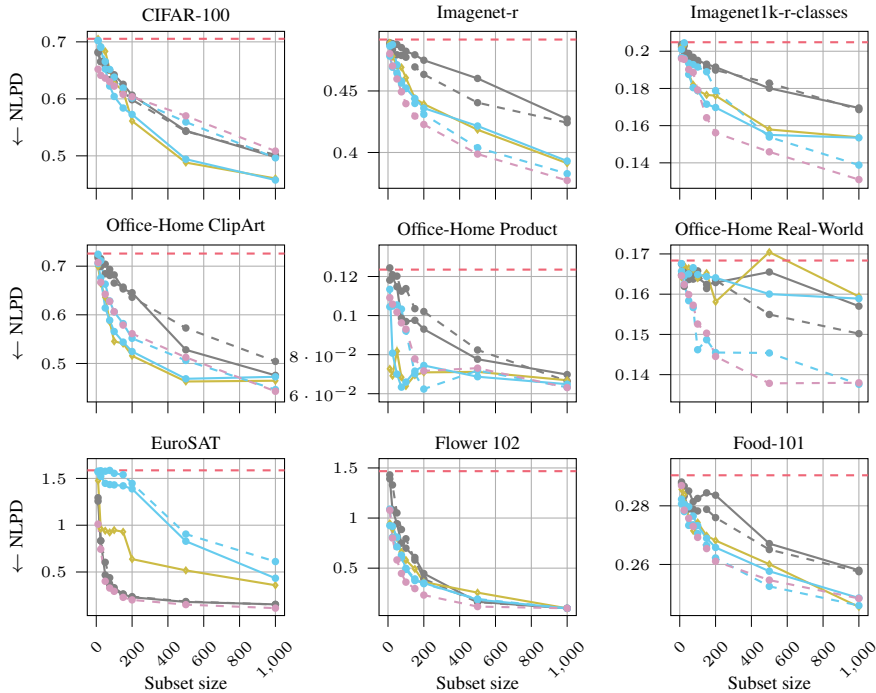
477 We observe that the selection methods using the epistemic uncertainty (BALD and EPIG) perform  
478 better or on par with the Targeted Maximum Entropy across the different subset sizes and data sets.  
479 The accuracy and NLPD become better when increasing the subset sizes, and the huge model variant  
480 (Fig. 7) achieves higher accuracies and lower NLPD on all data sets compared to the base model  
481 variant (Fig. 8) due to its larger model capacity. On EuroSAT, the Random baselines perform on par  
482 with EPIG which possibly is due to that EuroSAT has a small number of classes that can be similar,  
483 *e.g.*, the classes Sea/Lake and River. These results demonstrate the benefits of using our proposed  
484 uncertainty estimates for support set selection.

## 485 E.3 Covariance Analysis

486 In addition to the presented experiments, we performed an ablation on the sensitivity of the covariance  
487 to perturbations in the inputs. As shown in App. E.3, we observe that the covariance over the cosine  
488 similarities encodes meaningful information about the uncertainty of the model predictions under  
489 input perturbations. Further, we observe that the covariance structure captures similarity between  
490 inputs, *e.g.*, semantic similarity between text descriptions, as shown in Fig. 11.

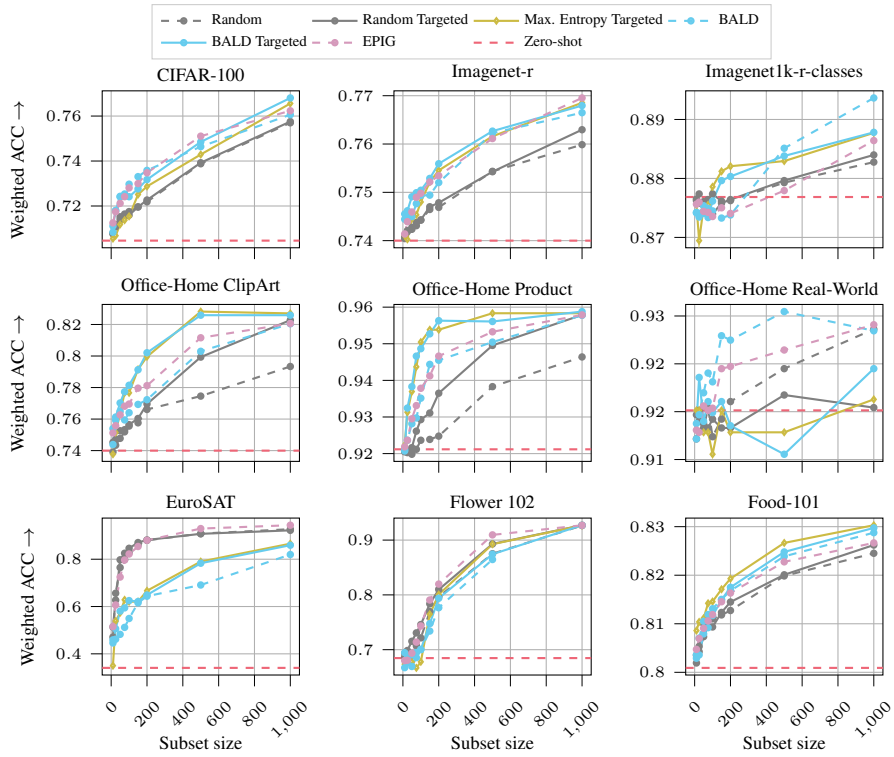


(a) Weighted accuracy (ACC) by the number of samples per class.

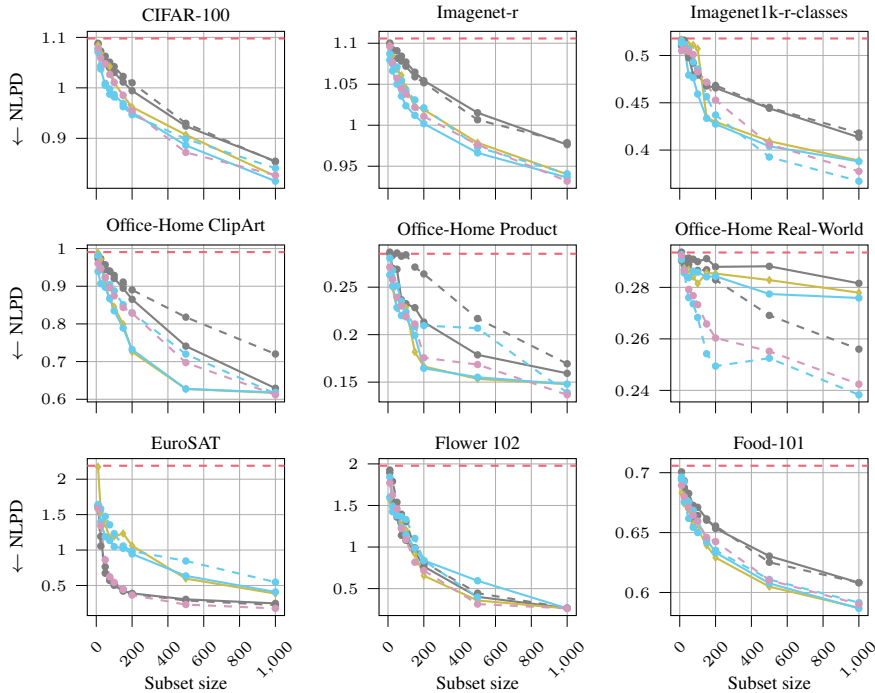


(b) Negative log-probability density (NLPD).

Figure 7: Accuracy and negative log-probability density (NLPD) over subset sizes of the support set across different data sets and subset selection methods using the OpenCLIP huge model variant. Results for random are averaged over 5 seeds.



(a) Weighted accuracy (ACC) by the number of samples per class.



(b) Negative log-probability density (NLPD).

Figure 8: Accuracy and negative log-probability density (NLPD) over subset sizes of the support set across different data sets and subset selection methods using the OpenCLIP base model variant. Results for random are averaged over 5 seeds.

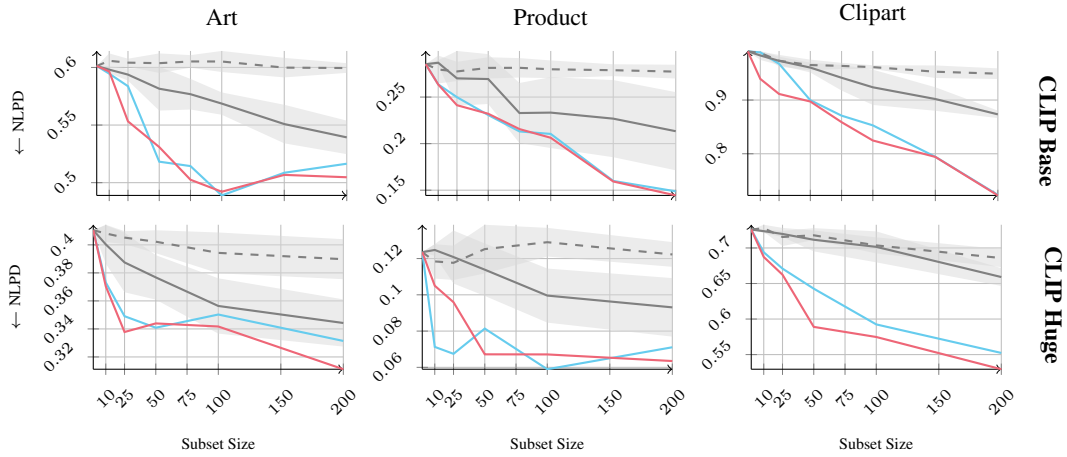


Figure 9: Results on the Office-Home data set with support set selection from all training domains. We depict the performance of the best performing acquisition function incorporating epistemic uncertainties (—), entropy based selection with targeted support set region (—), naïve random selection (---), and random selection with targeted support set candidates (—).

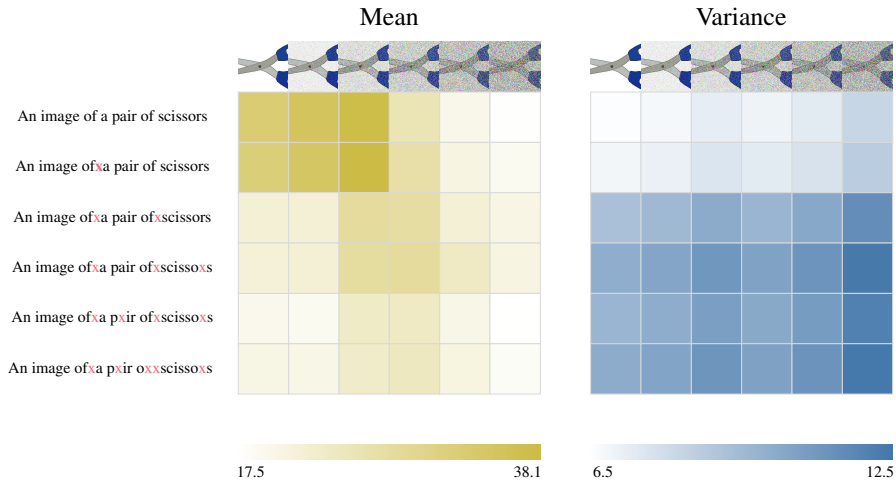


Figure 10: Illustration of the distribution over cosine similarities, depicting mean and variance, for varying image and text perturbations. We can observe that the mean cosine similarity decreases with increasing perturbation, while the variance increases, indicating that the distribution over cosine similarities captures model uncertainties in out-of-distribution settings.

### Logit Covariance

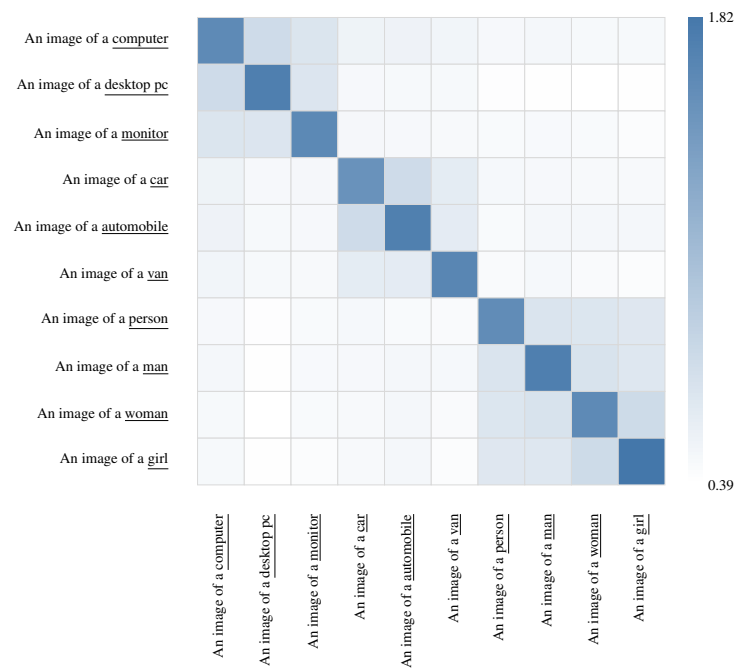


Figure 11: Illustration of the cosine similarity covariance between different text encodings. We find that the covariance captures correlations between semantically similar descriptions/class labels.