

PRIOR DISTRIBUTION AND MODEL CONFIDENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper investigates the impact of training data distribution on the performance of image classification models. By analyzing the embeddings of the training set, we propose a framework to understand the confidence of model predictions on unseen data without the need for retraining. Our approach filters out low-confidence predictions based on their distance from the training distribution in the embedding space, significantly improving classification accuracy. We demonstrate this on the example of several classification models, showing consistent performance gains across architectures. Furthermore, we show that using multiple embedding models to represent the training data enables a more robust estimation of confidence, as different embeddings capture complementary aspects of the data. Combining these embeddings allows for better detection and exclusion of out-of-distribution samples, resulting in further accuracy improvements. The proposed method is model-agnostic and generalizable, with potential applications beyond computer vision, including domains such as Natural Language Processing where prediction reliability is critical.

1 INTRODUCTION

In recent years, we have seen unprecedented development of deep learning models and an extension of the problems that can be effectively solved by them. They have dramatically advanced AI capabilities, but their training paradigm imposes fundamental constraints. Leading companies, such as OpenAI, have utilized nearly all accessible data to build these models, which require immense computational resources and time (Brown et al., 2020). However, this approach presents major limitations, such as the difficulty in optimizing model architectures and the exclusivity of achieving such results to a handful of research institutions and corporations. Consequently, most scientists and engineers are limited to using pre-trained models without the ability to modify or improve them (Bommasani et al., 2021). Shortly after deployment, numerous instances of erroneous outputs, termed ‘hallucinations,’ were identified in these models. These are cases where the response sounds convincing and well-structured but conveys false or unverified information (Huang et al., 2023). A significant challenge with these errors is the difficulty in detecting their root causes through automated means, even though users can often recognize them with ease. This poses a serious risk in industrial environments, where such inaccuracies can trigger unforeseen system behaviors or failures in AI-driven technologies.

Although understanding the factors behind such behavior in AI systems is challenging due to their inherent complexity, models of that type (such as Huawei’s PanGu- Σ) often contain over 1 trillion parameters (Ren et al., 2023). Mathematical structures of this magnitude are extremely difficult to analyze from both theoretical and technical perspectives (Hudson et al., 2024).

The field of computer vision continues to be highly suitable for theoretical exploration. Unlike NLP, where constructing representative datasets for industrial-level models is difficult due to the enormous data requirements, computer vision faces fewer such challenges. As one of the earliest domains enhanced by deep learning, it now benefits from a rich collection of well-curated datasets ideal for research purposes (Wan, 2024).

Consequently, using computer vision models to investigate the nature of ‘hallucinations’ appears more reasonable. As is well-known, training data is a primary factor influencing model performance, as it fundamentally determines the model’s scope and applicability (Chen et al., 2024). Therefore,

054 understanding the distribution of training data is crucial for assessing a model’s predictability (Györi
055 et al., 2022; Dabiri et al., 2023).

056 This paper presents a foundational approach for designing algorithms to detect hallucinations. The
057 core idea is to leverage the latent space representation of training data to assess the model’s pre-
058 dictive capabilities on specific instances. While the current study focuses on deep learning models
059 in computer vision, the proposed methodology is general enough to be extended to other machine
060 learning and deep learning domains, including natural language processing.

062 2 RELATED WORK

063 There are several approaches developed to mitigate “hallucinations” without recreating the model.

064 A practical first step is to adapt the model’s parameters for a given task using tuning strategies like
065 fine-tuning or alignment (Zhang et al., 2023; Wu et al., 2025; Lialin et al., 2023; Prottasha et al.,
066 2024) . However, these techniques often require significant amounts of extra training data and may
067 still lack dependable reliability.

068 Another solution is to mitigate hallucinations by using RAG (Retrieval-Augmented Generation)
069 techniques, which integrate information retrieval with generation. This allows the model to reference
070 real-world information, making the output more accurate and reliable (Bécharde & Ayala, 2024).
071 Another option is to apply statistical and information-theoretic methods to detect hallucinations,
072 aiming to flag outputs that deviate from expected or reliable patterns (Farquhar et al., 2024; Li,
073 2025). While these methods focus on analyzing the output of deep learning networks, they often
074 overlook the distribution of the training data. Yet, the structure and spread of training data across
075 complex feature spaces can critically impact a deep learning model’s performance (Hammoudeh
076 & Lowd, 2024; Lin et al., 2022). Accurately measuring and modeling input data distribution is a
077 crucial strategy for enhancing model performance in machine learning. Numerous methods have
078 been developed to address this challenge effectively (Hammoudeh & Lowd, 2024; Tamang et al.,
079 2025).

080 First, training data is formally described and cataloged using tools like Data Cards, and visualized
081 through techniques such as t-SNE or UMAP applied to CLIP or DINO embeddings. Mislabeled
082 instances can be identified and corrected with Confident Learning (Northcutt et al., 2021), while
083 Dataset Cartography (Swayamdipta et al., 2020) enables categorization of data points into “easy,”
084 “hard,” or “ambiguous” classes. Second, distribution shift and coverage are measured. To detect
085 samples outside the training distribution, OOD (out-of-distribution) methods like ODIN (Liang
086 et al., 2018), Mahalanobis Distance (Lee et al., 2018), Energy-based OOD (Liu et al., 2020), and
087 Open OOD (Yang et al., 2022) are used. As an opposite approach, we can start from the defined
088 distribution of the training data and try to fit the data to it. The Core Set method (Sener & Savarese,
089 2018) selects diverse representatives to cover the input space. To select uncertain and diverse exam-
090 ples, the BADGE method (Ash et al., 2020), based on gradient embeddings, can be used. If you have
091 a set of specific features, the Submodular Selection method (Wei et al., 2015) ensures their cover-
092 age. In the case of class imbalance, the Focal Loss method (Lin et al., 2017) reweights samples from
093 imbalanced classes. Methods like LAION (Schuhmann et al., 2021) and DataComp (Gadre et al.,
094 2023) demonstrate that targeted filtering toward the desired distribution can significantly outperform
095 simply increasing dataset size in terms of accuracy.

096 In addition, several methods focus on evaluating training data by analyzing the contribution of its
097 individual components. Influence Functions (Koh & Liang, 2017) and TracIn (Pruthi et al., 2020)
098 help quantify the impact of specific training samples on model predictions, making it possible to
099 identify non-representative data points that may degrade performance. Other evaluation approaches,
100 such as VQA v2 (Goyal et al., 2017), aim to reduce shortcut learning by rebalancing answer priors,
101 thereby encouraging models to rely on visual information. Similarly, ObjectNet is designed to assess
102 generalization capabilities under varying viewpoints and backgrounds.

104 3 METHODS

105 This manuscript presents a method for understanding the influence of the prior data distribution on
106 the confidence of prediction, independent of the predictive signal of the model.
107

To investigate the effect of data distribution on model confidence, we selected several models trained on the same dataset. Specifically, all models were trained using the ImageNet-1K dataset (Deng et al., 2009), ensuring a consistent distribution of training data between experiments. In order to understand the difference between the performance of our algorithm, we intentionally took several models from different classes with different numbers of parameters. The chosen models are: ResNet type models (Resnet-101 and Resnet-50 (He et al., 2016)), Visual Transformers (Touvron et al., 2021) (DeiT-T, DeiT-S and DeiT-B) and a relatively small model ShuffleNet-V2 (Ma et al., 2018). The brief description of the models is shown in Table 1.

Table 1: Overview of the models used in our experiments. All models were trained on the ImageNet-1K dataset to ensure consistent data distribution.

Model	Training Dataset	Architecture Type	Reference	# Parameters
ResNet-101	ImageNet-1K	Convolutional (CNN)	(He et al., 2016)	44.5M
ResNet-50	ImageNet-1K	Convolutional (CNN)	(He et al., 2016)	25.6M
ShuffleNet-V2	ImageNet-1K	Lightweight CNN	(Ma et al., 2018)	2.3M
DeiT-Tiny	ImageNet-1K	Vision Transformer (ViT)	(Touvron et al., 2021)	5.7M
DeiT-Small	ImageNet-1K	Vision Transformer (ViT)	(Touvron et al., 2021)	22.1M
DeiT-Base	ImageNet-1K	Vision Transformer (ViT)	(Touvron et al., 2021)	86.4M

To understand the data distribution on which the models were trained, we analyzed the structure of the ImageNet-1K training set by computing image embeddings using several pretrained models. Specifically, we selected a diverse set of architectures to generate embeddings for each image in the dataset. These embedding-based representations allowed us to examine the statistical properties of the data in different feature spaces. The analysis provides insight into how various models perceive and encode the visual information present in the training distribution. It is worth mentioning that we took only the "train" part of the set since exactly on this set the models in Table 1 were trained. For the embedding analysis, we chose the following models: The Dino-V1 model (Caron et al., 2021), Dino-V2 models (Oquab et al., 2023) with different embedding sizes (384, 768, 1024) and the MobileNet-V2 model (Sandler et al., 2018). Our choice of embedding models was motivated by two key considerations: the model size and the nature of the data used for pre-training. To assess the sensitivity of our algorithm to embedding dimensionality, we evaluated its performance on closely related models with varying embedding sizes. In addition, we were interested in understanding the effects of the data on which the embedding models were trained. We include a brief description of the embedding models in Table 2.

Table 2: Description of the models and their embedding sizes.

Model	Training Data	Embedding Size	Reference
DINO-V1 ViT-S/16	ImageNet-1K (self-supervised)	384	(Deng et al., 2009)
DINO-V1 ViT-B/16	ImageNet-1K (self-supervised)	768	(Deng et al., 2009)
DINO-V1 ViT-B/8	ImageNet-1K (self-supervised)	768	(Deng et al., 2009)
DINO-V2 ViT-S/14	142M curated dataset (no labels)	384	(Oquab et al., 2023)
DINO-V2 ViT-B/14	142M curated dataset (no labels)	768	(Oquab et al., 2023)
MobileNet-V2	ImageNet-1K (supervised)	1000	(Sandler et al., 2018)

In order to investigate the effect of the prior data distribution on the prediction capabilities of the classification models, we have chosen two distinct datasets: ObjectNet (Barbu et al., 2019) and imageNet-V2 (Recht et al., 2019). We divided the imageNet-V2 into train and internal test sets with a split (75% / 25%) and left the ObjectNet dataset as the external test set. We used only ObjectNet classes that coincide with the classes in the original ImageNet1K dataset (which gave us 10 121 images). We used the train dataset to find the optimal parameter: the number of neighbors N . We divided the test sets (external and internal) into three subsets grouped by the labels. So, images with the same label can only be present in one subset. The goal of this division is to report the variation of the metrics.

For the calculation of embeddings, we used cosine similarity (Manning et al., 2008), which is widely regarded as an effective metric to identify neighbors in high-dimensional spaces. For storing and querying the embeddings, we used ChromaDB (Team, 2023) with ClickHouse (ClickHouse, 2016)

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

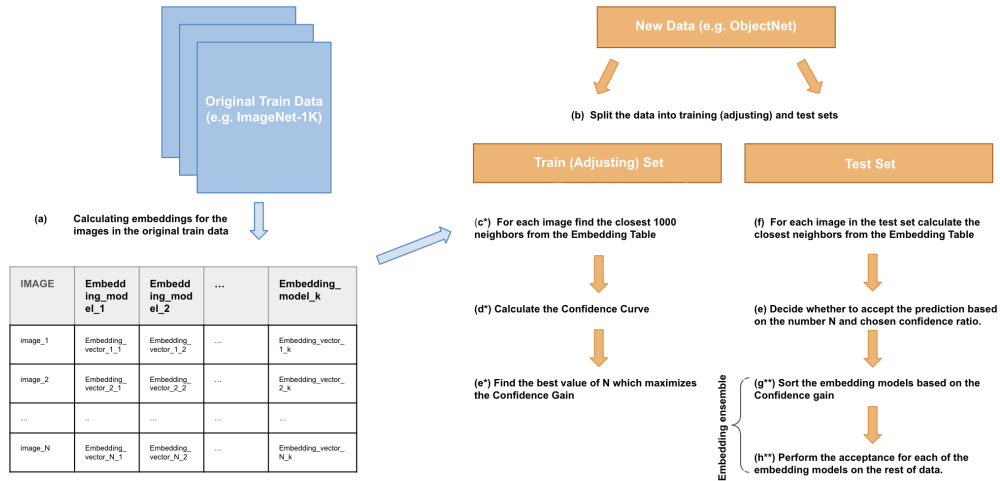


Figure 1: The scheme of the experiment. With (*) the training (adjusting) steps are indicated. With (**) supplementary combination predictions steps are shown

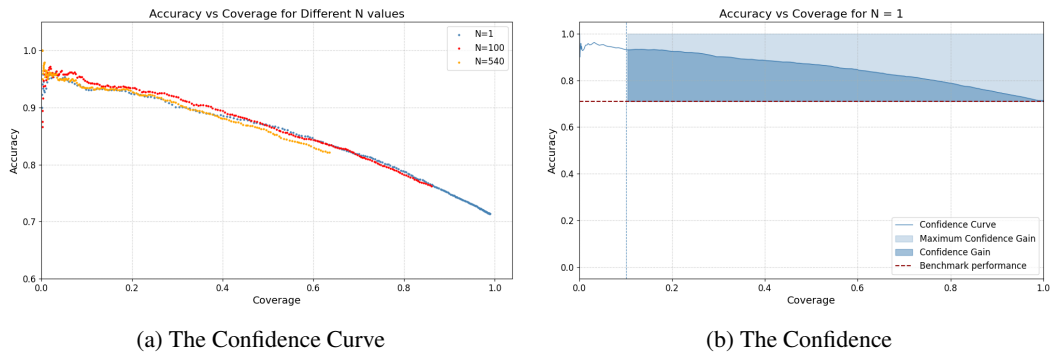


Figure 2: The Confidence Curves and Confidence Gain for the ResNet-101 & DINO-V2 ViT-B/14 .

as the backend server to enable efficient retrieval. All calculations were performed on a MacBook with an Apple M1 chip.

4 EXPERIMENTS

Confidence estimation was conducted on both the training and test sets. In the training set, we adjusted the parameter N to optimize performance. The procedure for adjusting confidence in the training set involves the following steps. An illustration of both the training and test pipelines is provided in Figure 1.

1. For each image in the original training dataset ImageNet-1K (let us call it Base Image Set), we calculate the embeddings for embedding models presented in Table 2 using cosine similarity. We store these embeddings in the Base Embedding Database. Although the Base Image Set serves as a training set for the models we investigated, we do not fine-tune its parameters. Therefore, the word "training" could be misleading.
2. For the training dataset (a subset of ImageNet-V2), we compute the embeddings for each image and identify 1000 closest neighbors from a base embedding database. This allows us to incorporate information about the local distribution of embeddings. We choose this

particular number to ensure that we fully describe the neighbors of a particular image. We refer to this collection as the Train Embedding Set. The term "train" emphasizes that this dataset is used to tune the parameter N , which is then kept fixed during the evaluation on the test sets.

3. For each image in the train image set, we calculate and store the predictions of the classification model.
4. For each image in the train set, we calculate the following statistics. We use the parameters $L_{threshold}$ and N as the threshold value for the distance of the neighbors and the necessary number of neighbors within the proximity $L_{threshold}$. Should the quantity of neighboring elements for the embedding of the image exceed N within the $L_{threshold}$ vicinity, we accept the model's prediction as valid. Conversely, if this threshold is not met, we conclude that the model lacks sufficient confidence to make a prediction. We calculate the accuracy for all images where the confidence is enough. Let us call the percentage of approved images the *Coverage*. The Confidence Curve would be the dependence of the accuracy upon the coverage. The examples of Confident Curves for ResNet-101 & DINO-V2 ViT-B/14 for three different N are presented in Figure 2(a). In Figure 3(b) *Maximum Confidence Gain* and *Confidence Gain* are shown for the same embedding / model pair, and $N = 1$. For more details on how the values are calculated, please see Appendix 1.
5. We find the values of N that maximize the *Confidence Gain* (the ratio between the *Maximum Confidence Gain* and the *Confidence Gain*). We use the value of neighbors for the test sets.
6. We repeated steps 3,4,5 for each of the classification models.

In Figure 4, the results of the adjustment (or training) are shown. On the x-axis, the number of neighbors N is presented, and on the y-axis, the calculated *Normalized Confidence Gain* is depicted.

Table 3: N and best *Normalized Confidence Gain* for each classification model across embedding models.

Classification Model	DINO-V1 ViT-B/16	DINO-V1 ViT-B/8	DINO-V1 ViT-S/16	DINO-V2 ViT-B/14	DINO-V2 ViT-S/14	MobileNet-V2
deit_base_patch16_224	1 / 0.398	1 / 0.441	1 / 0.368	24 / 0.480	4 / 0.419	1 / 0.194
deit_small_patch16_224	1 / 0.408	1 / 0.449	1 / 0.379	7 / 0.472	2 / 0.427	1 / 0.199
deit_tiny_patch16_224	1 / 0.402	1 / 0.433	1 / 0.379	23 / 0.442	4 / 0.406	1 / 0.218
ResNet-101	1 / 0.394	1 / 0.434	1 / 0.364	60 / 0.477	4 / 0.417	1 / 0.195
ResNet-50	1 / 0.407	1 / 0.448	1 / 0.380	24 / 0.487	2 / 0.426	1 / 0.209
shufflenet_v2_x1_0	1 / 0.418	1 / 0.436	1 / 0.401	7 / 0.417	4 / 0.395	1 / 0.244

Table 4: The Normalized Confidence Gain for different pairs of embedding/classification models and for the combination of embedding models. Benchmark accuracy is shown for reference. All values are mean \pm std across three label subsets for internal and external test sets.

Internal Test								
Classification Model	Benchmark Accuracy	DINO-V1 ViT-B/16	DINO-V1 ViT-B/8	DINO-V1 ViT-S/16	DINO-V2 ViT-B/14	DINO-V2 ViT-S/14	MobileNet-V2	Combination
deit_base_patch16_224	0.696 \pm 0.014	0.390 \pm 0.052	0.414 \pm 0.048	0.357 \pm 0.035	0.449 \pm 0.077	0.404 \pm 0.046	0.194 \pm 0.046	0.361 \pm 0.058
deit_small_patch16_224	0.666 \pm 0.019	0.376 \pm 0.068	0.403 \pm 0.071	0.343 \pm 0.062	0.431 \pm 0.076	0.376 \pm 0.058	0.189 \pm 0.058	0.320 \pm 0.060
deit_tiny_patch16_224	0.582 \pm 0.037	0.376 \pm 0.076	0.388 \pm 0.079	0.353 \pm 0.078	0.408 \pm 0.070	0.380 \pm 0.067	0.202 \pm 0.069	0.366 \pm 0.067
ResNet-101	0.687 \pm 0.025	0.345 \pm 0.074	0.369 \pm 0.066	0.308 \pm 0.072	0.410 \pm 0.106	0.366 \pm 0.058	0.169 \pm 0.065	0.316 \pm 0.071
ResNet-50	0.676 \pm 0.014	0.377 \pm 0.050	0.408 \pm 0.049	0.347 \pm 0.050	0.446 \pm 0.075	0.398 \pm 0.044	0.189 \pm 0.068	0.344 \pm 0.061
shufflenet_v2_x1_0	0.538 \pm 0.022	0.363 \pm 0.056	0.374 \pm 0.063	0.345 \pm 0.056	0.379 \pm 0.053	0.357 \pm 0.046	0.220 \pm 0.050	0.351 \pm 0.046
External Test								
Classification Model	Benchmark Accuracy	DINO-V1 ViT-B/16	DINO-V1 ViT-B/8	DINO-V1 ViT-S/16	DINO-V2 ViT-B/14	DINO-V2 ViT-S/14	MobileNet-V2	Combinations
deit_base_patch16_224	0.269 \pm 0.059	0.079 \pm 0.109	0.132 \pm 0.120	0.059 \pm 0.107	0.214 \pm 0.153	0.189 \pm 0.149	0.044 \pm 0.084	0.141 \pm 0.130
deit_small_patch16_224	0.229 \pm 0.045	0.072 \pm 0.092	0.120 \pm 0.099	0.055 \pm 0.082	0.188 \pm 0.125	0.170 \pm 0.119	0.037 \pm 0.063	0.086 \pm 0.087
deit_tiny_patch16_224	0.152 \pm 0.041	0.073 \pm 0.077	0.102 \pm 0.085	0.058 \pm 0.072	0.119 \pm 0.092	0.119 \pm 0.093	0.040 \pm 0.057	0.106 \pm 0.090
ResNet-101	0.304 \pm 0.057	0.069 \pm 0.112	0.131 \pm 0.112	0.049 \pm 0.101	0.238 \pm 0.144	0.211 \pm 0.144	0.036 \pm 0.080	0.147 \pm 0.123
ResNet-50	0.293 \pm 0.056	0.069 \pm 0.102	0.127 \pm 0.106	0.047 \pm 0.101	0.243 \pm 0.130	0.211 \pm 0.135	0.035 \pm 0.082	0.136 \pm 0.114
shufflenet_v2_x1_0	0.117 \pm 0.027	0.062 \pm 0.064	0.079 \pm 0.065	0.047 \pm 0.061	0.096 \pm 0.069	0.090 \pm 0.068	0.038 \pm 0.048	0.083 \pm 0.066

We can see that the highest *Normalized Confidence Gain* value is achieved for all models in besides the ShuffleNet-V2 by the DINO-V2 ViT-B / 14 embedding model, which is so far the strongest model among embedding models in terms of the number of parameters. Interestingly, the second version of the DINO models outperforms the first version. In table 3, the best parameter N and the corresponding *Normalized Confidence Gain* are presented.

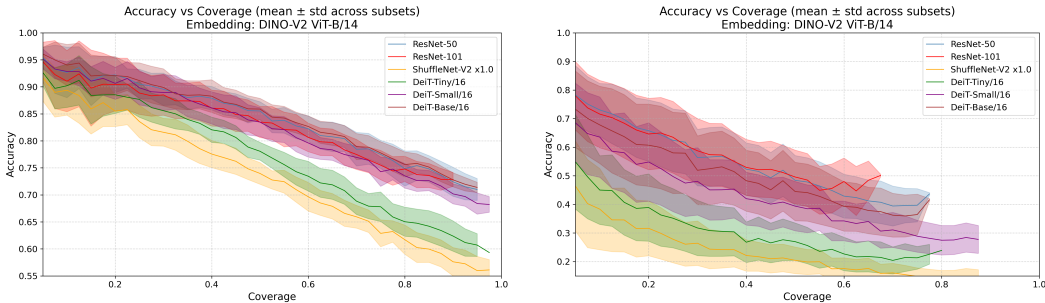
For the test pipeline, we use the number of parameters N , which maximizes the *Normalized Confidence Gain* for the specific pair of the classification model and the embedding model. We use this value to calculate the metric for the Internal Test and External Test sets. The metrics and their standard deviations (for 3 subsets of each of the External and Internal Tests) are reported in Table 4. We can observe that the highest value of *Normalized Confidence Gain* is reached almost for all cases by Dino-v2 ViT-B14. As mentioned earlier, this is quite logical and can be explained by the model’s higher capacity and the higher resolution of the input images.

The next logical step would be to combine several embedding models (or use an approach similar to the ensemble of models, as in (Dietterich, 2000)), which has been done using greedy heuristics. We fix the same coverage for each embedding model. We take the model with the highest *Normalized Confidence Gain* and make predictions. We accept a portion of the predictions and leave the rest to other models. It should be clear that the total coverage for all embedding models is slightly higher than the parameter *coverage* for individual models, but is not fixed. The results of the experiments (accuracy vs. coverage) are presented in Figure 5. The *Normalized Confidence Gain* for the combination of embedding models is shown in Table 4 in the column "Combination". The observed *Confidence Gain* for the combination of embedding models is lower than that of the best-performing individual model. In particular, at very low coverage values for the internal dataset, the accuracy of all models converges to similar levels. This suggests that the underlying data distribution plays a more decisive role than the intrinsic performance of the models. A plausible explanation is that, within the test datasets, certain images are located in close proximity (in the embedding space) to images from the training set. Consequently, under such conditions, the specific choice of the classification model becomes less critical. For the external dataset, the same conclusions cannot be drawn. This is probably related to the fact that the images in the external dataset (ObjectNet) exhibit a greater distributional shift relative to the original training set (ImageNet-1K) than those in the internal test set.

5 DISCUSSION

From the experiment chapter, we can infer that a more sophisticated embedding model is more likely to have a better Confidence Curve (or higher *Normalized Confidence Gain*). With "sophisticated" we mean a larger number of dimensions, larger size of the entering image, and a richer dataset on which the model was trained. In particular, the dimensions of the images and the capacity of the embedding model significantly determine the performance.

We showed that combining embeddings from multiple models can be a promising strategy for integrating diverse representation spaces. However, the *Normalized Confidence Gain* obtained from



(a) Confidence Curves for Internal Dataset.

(b) Confidence Curves for External Dataset.

Figure 3: Accuracy vs Coverage (Confidence Curve) for the best Embedding Model (DINO-V2 ViT-B/14).

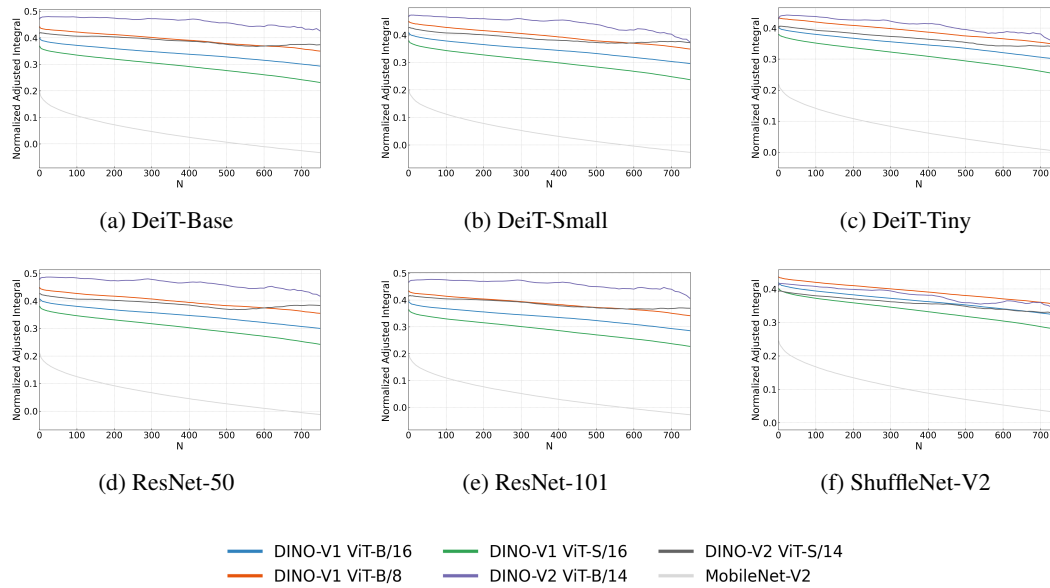
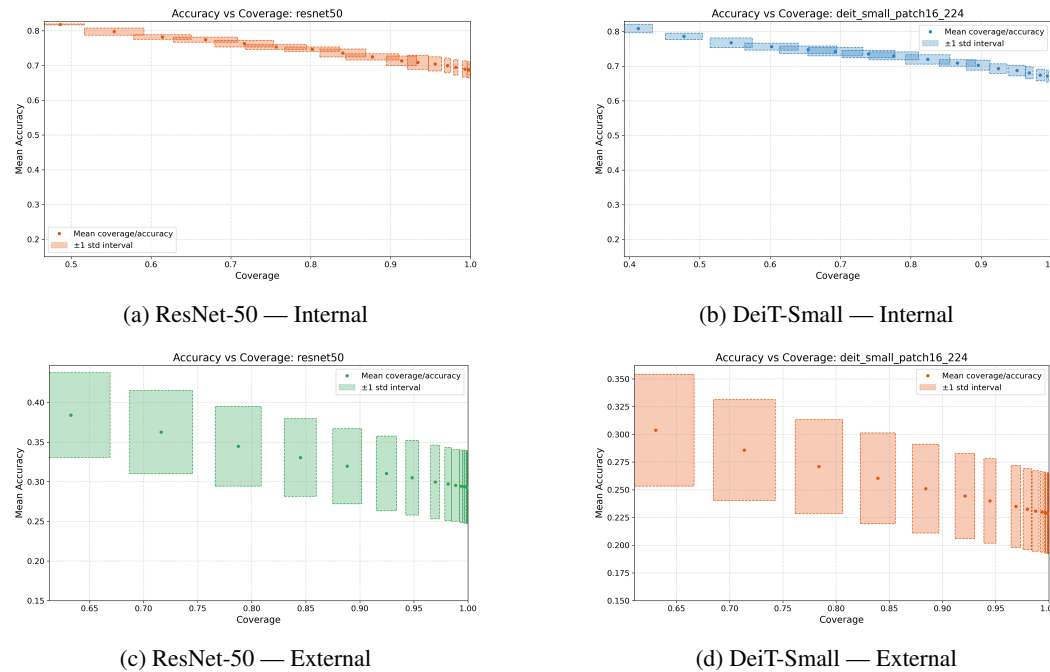
Figure 4: Normalized Confidence Gain vs number of neighbors N for different classification models.

Figure 5: The Confidence Curve (Accuracy vs Total coverage) for the internal (a,b) and external (c,d) datasets for models ResNet-101 (a,c) and DeiT-Small (b,d). using the combinations of the embedding models.

such combinations was lower than that of the best performing individual embedding model. Unlike classical ensemble methods, where aggregating weaker models often leads to improved performance, we did not observe a performance gain in this setting. We identify two possible reasons for this. First, our combination strategy was a simple greedy approach and more sophisticated methods could yield better results. For example, a meta-classifier could be trained to select the most suitable embedding model for confidence estimation on a per-sample basis. Second, most of the models in our study were trained on only two distinct datasets. We expect that training on more diverse datasets

could enrich the combined confidence estimates. However, achieving such diversity is challenging, as off-the-shelf models are typically trained on large, overlapping datasets, limiting the potential for complementary knowledge.

We can see that after training, most of the models achieved the highest *Normalized Confidence Gain* with $N = 1$. Only the best-performing embedding models achieved the maximum value at significantly larger values of N ($N > 1$). We believe that the more complicated the model, the larger the optimal value of N .

Our experiments were aimed solely at the image classification model. However, we suggest that a similar approach can work for other areas of Machine Learning, particularly for Natural Language Processing. But there is a caveat: we do not know how to partition the data efficiently. Images are distinct entities, and it is logical to compute the embeddings for each image. How exactly to divide the text corpora is unclear. There are two main questions in this investigation. What size of window to use and what jump step to use? Perhaps we need to have embedding collections of different window sizes, which could enrich the understanding of the neighborhood locality. Regarding the size of the jump, it would be interesting to calculate the embeddings with a step of one token. However, it would be computationally very heavy. In fact, English Wikipedia contains approximately 13.5 billion token instances across all articles, according to the TokTrack dataset (Flöck et al., 2017) (as of 2017). The training corpora used for large language models (LLMs) are estimated to be two to three orders of magnitude larger, although the exact size is unknown. For example, whereas GPT-3 was trained on approximately 499 billion tokens (Brown et al., 2020). In comparison, ImageNet-1K (used in this paper) contains approximately 1 million images and, correspondingly, 1 million embedding vectors. Moreover, it remains unclear how much efficiency is lost in neighborhood locality calculations when the step size exceeds a single token.

The current framework does not directly improve the performance of the classification model (or potentially other models). However, it enables the identification of samples that are out-of-distribution with respect to the original training set, thus making the model more confident in certain predictions.

6 CONCLUSION

In this paper, we showed that knowledge of the training dataset and its embeddings can be beneficial for determining a model’s confidence in its predictions on newly seen data. The presented framework provides an alternative way to estimate prediction confidence. In particular, we significantly improved the accuracy of classification models without additional retraining—only by adjusting the number of neighbors—while skipping a minimal number of images that lie far from the training distribution. We demonstrated that using a single embedding model to filter out low-confidence samples can improve overall accuracy. Finally, we propose that the framework is not limited to object recognition and computer vision tasks, but could also be applied to other sensitive areas, such as Natural Language Processing, where hallucinations are common. To the best of our knowledge, this work does not present any ethical issues.

REPRODUCIBILITY STATEMENT

All data used in this study are publicly available: ImageNet-1K, ImageNet-V2, and ObjectNet. The embedding models used in our experiments are also publicly available and were not retrained. The classification models evaluated for confidence and accuracy are likewise publicly available and were not retrained.

We provide an anonymized `.zip` archive in the supplementary materials containing: (i) all source code for the embedding extraction, confidence-curve calculation, and evaluation, (ii) instructions and configuration files to reproduce every experiment, and (iii) scripts to regenerate all tables and figures in the paper.

Running the provided scripts on the stated datasets reproduces every reported metric and plot.

REFERENCES

- 432
433
434 Jordan Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep
435 batch active learning by diverse, uncertain gradient lower bounds. In *International Conference on*
436 *Learning Representations (ICLR)*, 2020. URL <https://arxiv.org/abs/1906.03671>.
- 437 Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Joshua
438 Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the
439 limits of object recognition models. In *Advances in Neural Information Processing Systems*
440 *(NeurIPS)*, pp. 9448–9458, 2019. URL [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2019/hash/97af07a14caca681feacf301271f014-Abstract.html)
441 [2019/hash/97af07a14caca681feacf301271f014-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/97af07a14caca681feacf301271f014-Abstract.html).
- 442 P. B  chard and O. M. Ayala. Reducing hallucination in structured outputs via retrieval-augmented
443 generation. *arXiv preprint arXiv:2404.08189*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2404.08189)
444 [2404.08189](https://arxiv.org/abs/2404.08189).
- 445 Rishi Bommasani et al. On the opportunities and risks of foundation models. *arXiv preprint*
446 *arXiv:2108.07258*, 2021. URL <https://arxiv.org/abs/2108.07258>.
- 447 Tom Brown et al. Language models are few-shot learners. In *Advances in Neural Information*
448 *Processing Systems*, volume 33, pp. 1877–1901, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2005.14165)
449 [2005.14165](https://arxiv.org/abs/2005.14165).
- 450 Mathilde Caron, Hugo Touvron, Ishan Misra, Herv   J  gou, Julien Mairal, Piotr Bojanowski, and
451 Armand Joulin. Emerging properties in self-supervised vision transformers. *arXiv preprint*
452 *arXiv:2104.14294*, 2021. URL <https://arxiv.org/abs/2104.14294>.
- 453 X. Chen, Z. Ma, X. Zhang, S. Xu, S. Qian, J. Yang, D. F. Fouhey, and J. Chai. Multi-object
454 hallucination in vision-language models. *arXiv preprint arXiv:2407.06192*, 2024. doi: 10.48550/
455 [arXiv.2407.06192](https://doi.org/10.48550/arXiv.2407.06192). URL <https://doi.org/10.48550/arXiv.2407.06192>.
- 456 Inc. ClickHouse. Clickhouse: Open-source column-oriented database management system.
457 <https://clickhouse.com>, 2016. Accessed: 2025-07-27.
- 458 S. Dabiri, V. Lioutas, B. Zwartsenberg, Y. Liu, M. Niedoba, X. Liang, D. Green, J. Sefas, J. W.
459 Lavington, F. Wood, and A. Scibior. Realistically distributing object placements in synthetic
460 training data improves the performance of vision-based object detection models. *arXiv preprint*
461 *arXiv:2305.14621*, 2023. URL <https://arxiv.org/abs/2305.14621>.
- 462 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
463 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
464 pp. 248–255. IEEE, 2009. doi: 10.1109/CVPR.2009.5206848.
- 465 Thomas G. Dietterich. Ensemble methods in machine learning. In *Multiple Classifier Systems:*
466 *First International Workshop, MCS 2000, Cagliari, Italy, June 21–23, 2000, Proceedings*, pp.
467 1–15. Springer, 2000. doi: 10.1007/3-540-45014-9_1.
- 468 S. Farquhar, J. Kossen, L. Kuhn, and Y. Gal. Principled detection of hallucinations
469 in large language models using entropy-based uncertainty estimators. *Nature*, 2024.
470 doi: 10.1038/s41586-024-07421-0. URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41586-024-07421-0)
471 [s41586-024-07421-0](https://www.nature.com/articles/s41586-024-07421-0).
- 472 Fabian Fl  ck, Kenan Erdogan, and Maribel Acosta. Toktrack: A complete token provenance and
473 change tracking dataset for the english wikipedia. *arXiv preprint arXiv:1703.08244*, 2017. Token
474 count: approximately 13.545 billion.
- 475 Syed Yousuf Gadre, Gabriel Ilharco, Mitchell Wortsman, Ludwig Schmidt, and Hannaneh Ha-
476 jishirzi. Datacomp: In search of the next generation of multimodal datasets. In *Advances in*
477 *Neural Information Processing Systems (NeurIPS)*, volume 36, pp. 13271–13295, 2023. URL
478 <https://arxiv.org/abs/2304.14108>.
- 479 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V
480 in VQA matter: Elevating the role of image understanding in visual question answering. In
481 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.
482 6904–6913. IEEE, 2017. URL <https://arxiv.org/abs/1612.00837>.

- 486 N. G. Györi et al. Training data distribution significantly impacts the accuracy and precision of
487 parameter estimates. *Quantitative Imaging in Medicine and Surgery*, 2022. URL [https://](https://pubmed.ncbi.nlm.nih.gov/34545955/)
488 pubmed.ncbi.nlm.nih.gov/34545955/.
489
- 490 Zayd Hammoudeh and Daniel Lowd. Training data influence analysis and estimation: a survey.
491 *Machine Learning*, 113:2351–2403, 2024. doi: 10.1007/s10994-023-06495-7. URL [https://](https://link.springer.com/article/10.1007/s10994-023-06495-7)
492 link.springer.com/article/10.1007/s10994-023-06495-7.
- 493 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
494 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*
495 *(CVPR)*, pp. 770–778, 2016.
- 496
- 497 L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and
498 T. Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges,
499 and open questions. *arXiv preprint arXiv:2311.05232*, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2311.05232)
500 [abs/2311.05232](https://arxiv.org/abs/2311.05232).
- 501 N. Hudson, J. G. Pauloski, M. Baughman, A. Kamatar, M. Sakarvadia, L. Ward, R. Chard, A. Bauer,
502 M. Levental, W. Wang, W. Engler, O. P. Skelly, B. Blaiszik, R. Stevens, K. Chard, and I. Foster.
503 Trillion parameter ai serving infrastructure for scientific discovery: A survey and vision. *arXiv*
504 *preprint arXiv:2402.03480*, 2024. URL <https://arxiv.org/abs/2402.03480>.
- 505
- 506 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
507 *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70 of
508 *Proceedings of Machine Learning Research*, pp. 1885–1894. PMLR, 2017. URL [https://](https://arxiv.org/abs/1703.04730)
509 arxiv.org/abs/1703.04730.
- 510 Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. A simple unified framework for detecting
511 out-of-distribution samples and adversarial attacks. In *Advances in Neural Information Process-*
512 *ing Systems (NeurIPS)*, volume 31, 2018. URL <https://arxiv.org/abs/1807.03888>.
- 513
- 514 Jiawei Li. Principled detection of hallucinations in llms via multiple-testing framework. *arXiv*
515 *preprint arXiv:2508.18473*, 2025. URL <https://arxiv.org/abs/2508.18473>.
- 516
- 517 V. Lialin, V. Deshpande, and A. Rumshisky. Scaling down to scale up: A guide to parameter-
518 efficient fine-tuning. *arXiv preprint arXiv:2303.15647*, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2303.15647)
[abs/2303.15647](https://arxiv.org/abs/2303.15647).
- 519
- 520 Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image
521 detection in neural networks. In *International Conference on Learning Representations (ICLR)*,
522 2018. URL <https://arxiv.org/abs/1706.02690>.
- 523
- 524 J. Lin et al. Measuring the effect of training data on deep learning: A query-wise perspec-
525 tive. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162
526 of *Proceedings of Machine Learning Research*, pp. 13253–13271. PMLR, 2022. URL [https://](https://proceedings.mlr.press/v162/lin22h/lin22h.pdf)
proceedings.mlr.press/v162/lin22h/lin22h.pdf.
- 527
- 528 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object
529 detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp.
530 2980–2988. IEEE, 2017. URL <https://arxiv.org/abs/1708.02002>.
- 531
- 532 Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detec-
533 tion. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pp. 21464–
21475, 2020. URL <https://arxiv.org/abs/2010.03759>.
- 534
- 535 Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines
536 for efficient cnn architecture design. In *Proceedings of the European Conference on Computer*
537 *Vision (ECCV)*, pp. 116–131, 2018.
- 538
- 539 Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information*
Retrieval. Cambridge University Press, 2008. ISBN 9780521865715. URL [https://nlp.](https://nlp.stanford.edu/IR-book/)
[stanford.edu/IR-book/](https://nlp.stanford.edu/IR-book/).

- 540 Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. Confident learning: Estimating uncertainty in
541 dataset labels. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pp.
542 16921–16940, 2021. URL <https://arxiv.org/abs/1911.00068>.
- 543
544 Maxime Oquab, Thomas Darcet, Theo Moutakanni, Maciej Szafraniec, Vladislav Khalidov, Pierre
545 Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Wojciech Galuba, et al. Di-
546 nov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*,
547 2023. URL <https://arxiv.org/abs/2304.07193>.
- 548 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
549 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
550 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
551 12:2825–2830, 2011.
- 552 N. J. Prottasha, A. Mahmud, M. S. I. Sobuj, P. Bhat, M. Kowsher, N. Yousefi, and Ö. Ö. Garibay.
553 Parameter-efficient fine-tuning of large language models using semantic knowledge tuning. *Sci-
554 entific Reports*, 14:30667, 2024. doi: 10.1038/s41598-024-75599-4. URL [https://doi.
555 org/10.1038/s41598-024-75599-4](https://doi.org/10.1038/s41598-024-75599-4).
- 556
557 Gaurav Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data
558 influence by tracing gradient descent. In *Advances in Neural Information Processing Systems
559 (NeurIPS)*, volume 33, pp. 20013–20024, 2020. URL [https://arxiv.org/abs/2002.
560 08484](https://arxiv.org/abs/2002.08484).
- 561 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
562 generalize to imagenet? In *Proceedings of the 36th International Conference on Machine Learn-
563 ing (ICML)*, pp. 5389–5400. PMLR, 2019. URL [https://proceedings.mlr.press/
564 v97/recht19a.html](https://proceedings.mlr.press/v97/recht19a.html).
- 565
566 X. Ren, P. Zhou, X. Meng, X. Huang, Y. Wang, W. Wang, P. Li, X. Zhang, A. Podolskiy, G. Ar-
567 shinov, A. Bout, I. Piontkovskaya, J. Wei, X. Jiang, T. Su, Q. Liu, and J. Yao. Pangu- σ : To-
568 wards trillion-parameter language model with sparse heterogeneous computing. *arXiv preprint
569 arXiv:2303.10845*, 2023. URL <https://arxiv.org/abs/2303.10845>.
- 570 Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh
571 Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the
572 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4510–
573 4520, 2018. URL [https://openaccess.thecvf.com/content_cvpr_2018/
574 html/Sandler_MobileNetV2_Inverted_Residuals_CVPR_2018_paper.html](https://openaccess.thecvf.com/content_cvpr_2018/html/Sandler_MobileNetV2_Inverted_Residuals_CVPR_2018_paper.html).
- 575 Christoph Schuhmann, Romain Vencu, Romain Beaumont, Cade Gordon, Ross Wightman, Mehdi
576 Cherti, Thea Coombes, Ashish Katta, Charlie Mullis, Mitchell Wortsman, Patrick Schramowski,
577 Shubhendu Kundurthy, Katherine Crowson, Ludwig Schmidt, Robin Kaczmarczyk, and Jenia
578 Jitsev. LAION-400M: Open dataset of CLIP-filtered 400 million image-text pairs. *arXiv preprint
579 arXiv:2111.02114*, 2021. URL <https://arxiv.org/abs/2111.02114>.
- 580
581 Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set
582 approach. In *International Conference on Learning Representations (ICLR)*, 2018. URL <https://arxiv.org/abs/1708.00489>.
- 583
584 Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi,
585 Noah A. Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with
586 training dynamics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Lan-
587 guage Processing (EMNLP)*, pp. 9275–9293. Association for Computational Linguistics, 2020.
588 URL <https://aclanthology.org/2020.emnlp-main.746>.
- 589
590 Lakpa Tamang, Mohamed Reda Bouadjeneq, Richard Dazeley, and Sunil Aryal. Handling out-of-
591 distribution data: A survey. *arXiv preprint arXiv:2507.21160*, 2025. URL [https://arxiv.
592 org/abs/2507.21160](https://arxiv.org/abs/2507.21160).
- 593 Chroma Team. Chroma: The ai-native open-source embedding database. [https://github.
com/chroma-core/chroma](https://github.com/chroma-core/chroma), 2023. Accessed: 2025-07-27.

- 594 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
 595 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *Pro-*
 596 *ceedings of the 38th International Conference on Machine Learning (ICML)*, pp. 10347–10357.
 597 PMLR, 2021. URL <https://proceedings.mlr.press/v139/touvron21a.html>.
- 598 Z. Wan. A survey of dataset refinement for problems in computer vision. *ACM Computing Sur-*
 599 *veys*, 2024. doi: 10.1145/3627157. URL [https://dl.acm.org/doi/abs/10.1145/](https://dl.acm.org/doi/abs/10.1145/3627157)
 600 [3627157](https://dl.acm.org/doi/abs/10.1145/3627157).
- 601 Kai Wei, Rishabh Iyer, and Jeff Bilmes. Submodularity in data subset selection and active learning.
 602 In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, volume 37 of
 603 *Proceedings of Machine Learning Research*. PMLR, 2015. URL [https://proceedings.](https://proceedings.mlr.press/v37/weil5.html)
 604 [mlr.press/v37/weil5.html](https://proceedings.mlr.press/v37/weil5.html).
- 605 Y. Wu, J. Li, Z. Guo, and L. Li. Learning like humans: Resource-efficient federated fine-tuning
 606 through cognitive developmental stages. *arXiv preprint arXiv:2508.00041*, 2025. URL [https:](https://arxiv.org/abs/2508.00041)
 607 [//arxiv.org/abs/2508.00041](https://arxiv.org/abs/2508.00041).
- 608 Jingkang Yang, Yiyao Li, Juntao Li, Chixiang Xu, Yuxuan Wang, Chenyang Zhang, Bing Xu, and
 609 Ziwei Lin. Openood: Benchmarking generalized out-of-distribution detection. In *Advances in*
 610 *Neural Information Processing Systems (NeurIPS)*, volume 35, pp. 29983–29997, 2022. URL
 611 <https://arxiv.org/abs/2210.07242>.
- 612 S. Zhang, L. Dong, X. Li, S. Zhang, X. Sun, S. Wang, J. Li, R. Hu, T. Zhang, F. Wu, and G. Wang.
 613 Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023.
 614 URL <https://arxiv.org/abs/2308.10792>.

618 A CONFIDENCE GAIN DETAILS

619 In this section, we explain in detail how the *Confidence Gain* is computed, alongside its key prop-
 620 erties. Figure 1(a) illustrates the confidence curve for the [ResNet-101 & DINO-V2 ViT-B/14]
 621 model–embedding pair; however, the method applies to all model–embedding combinations.

622 We can observe that the distribution of points along the curve is not uniform: for large coverage
 623 values, there are almost no points, whereas for small coverage values, the points are densely packed.
 624 This is because a small coverage corresponds to more confident predictions, which occur more fre-
 625 quently. To handle the lack of points near full coverage (that is, near *coverage* = 1), we linearly
 626 extrapolated the confidence curve using the 20 points with the highest value of *coverage*. Further-
 627 more, we noticed that for small values (below values of 0.1), the confidence curve tends to behave
 628 chaotically. Therefore, to ensure robustness, we computed the *Confidence Gain* by integrating over
 629 the interval *coverage* $\in [0.1, 1]$. Empirically, we observed that for *coverage* > 0.1, the accuracy
 630 behaves monotonically with the coverage and is less sensitive to noise.

631 We define the *Confidence Gain* as the area between the confidence curve (including extrapolation
 632 when necessary), the vertical line *coverage* = 0.1, and the horizontal line *accuracy* = acc_b . The
 633 acc_b is the accuracy of the classification model when using 100% of the data, i.e., without any
 634 coverage-based filtering. This value is intrinsic to the classification model and the dataset and does
 635 not depend on the embedding model. As *coverage* $\rightarrow 1$, the accuracy along the confidence curve
 636 naturally approaches acc_b . We calculate the area using the Scikit-Learn library (Pedregosa
 637 et al., 2011).

638 We define the *Maximum Confidence Gain* as the area bounded by the four lines:

$$\begin{aligned} 639 &x = 0.1, \\ 640 &x = 1.0, \\ 641 &y = 1.0, \\ 642 &y = acc_b. \end{aligned}$$

643 The value of the *Confidence Gain* is always less than the *Maximum Confidence Gain*. In other
 644 words, the Maximum Confidence Gain is a theoretical upper bound for the *Confidence Gain*. There-
 645 fore, it is logical to introduce the *Normalized Confidence Gain* as a ratio between these two values.
 646 The *Normalized Confidence Gain* would be determined in the range $[0, 1]$, where the value 1 corre-
 647 sponds to the perfect Confidence Gain, and the value 0 corresponds to no gain at all.

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

B LLM USAGE

This manuscript’s text was polished using the large language model ChatGPT (GPT-4, May 2025 release). The model was used only for grammar and style improvements. All scientific ideas, experiments, analyses, and conclusions were conceived and verified by the authors, who take full responsibility for the final content.