Pretrained Visual Uncertainties

Michael Kirchhof¹ Mark Collier² Seong Joon Oh³ Enkelejda Kasneci⁴

Abstract

Accurate uncertainty estimation is vital to trustworthy machine learning, yet uncertainties typically have to be learned for each task anew. This work introduces the first pretrained uncertainty modules for vision models. Similar to standard pretraining this enables the zero-shot transfer of uncertainties learned on a large pretraining dataset to specialized downstream datasets. We enable our large-scale pretraining on ImageNet-21k by solving a gradient conflict in previous uncertainty modules and accelerating the training by up to 180x. We find that the pretrained uncertainties generalize to unseen datasets. In scrutinizing the learned uncertainties, we find that they capture aleatoric uncertainty, disentangled from epistemic components. We demonstrate that this enables safe retrieval and uncertaintyaware dataset visualization. To encourage applications to further problems and domains, we release all pretrained checkpoints and code under https://github.com/mkirchhof/url.

1. Introduction

With every prediction comes the risk of an error. Uncertainty estimates quantify this expected error in order to defer predictions and catch errors before they happen, a key requirement for trustworthy machine learning (Mucsányi et al., 2023). Uncertainty quantification has seen tremendous advances in recent years, bringing principled methods such as Gaussian processes (Liu et al., 2020) and probabilistic embeddings (Oh et al., 2019; Kirchhof et al., 2023a; Kim et al.; Nakamura et al., 2023) to large-scale computer vision (Tran et al., 2022; Dehghani et al., 2023; Collier et al., 2023). Recent benchmarks reveal that they excel at their metrics and are ready for application (Galil et al., 2023a;b). However, there is a lack of widespread adoption of uncertainty

Preliminary work.

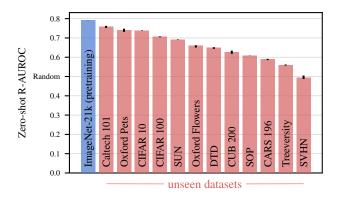


Figure 1. Our pretrained uncertainties generalize to unseen datasets. The R-AUROC measures the quality of uncertainty estimates on zero-shot datasets, see Section 4.

methods by practitioners. The hurdle is that modern uncertainty quantification methods can be complex, making them difficult to implement and increasing the inference costs. So what if uncertainty estimates were as easy to access as being shipped along with every pretrained model?

We seek a method that is simple to implement and train, even on large scales, and most importantly does not interfere with the main objective of a practitioner's model. One group of recent methods stands out on these terms: Feedforward uncertainties (Cui et al., 2023; Kirchhof et al., 2022; Yoo & Kweon, 2019; Oh et al., 2019). They an auxiliary uncertainty head to a deep network that is evaluated alongside every forward pass. Not only is this cheap and simple to implement, it was also recently found that its uncertainty estimates transfer well (Kirchhof et al., 2023b).

Our work solves a remaining gradient conflict of these feed-forward uncertainties to guarantee non-interference with the main objective. We also implement massive caching in our pretraining pipeline, reducing the train time by a factor of up to 180x. This enables us to scale up both the pretraining dataset to ImageNet-21k Winter-2021 (ImageNet-21k-W) (Deng et al., 2009) and the vision backbone to large Vision Transformers (Dosovitskiy et al., 2021).

Figure 1 shows that our pretrained uncertainties now transfer beyond the train dataset to unseen datasets, outperforming previous zero-shot uncertainties (Kirchhof et al., 2023b). We find that the learned uncertainties are generalizable no-

^{*}Equal contribution ¹University of Tübingen, Germany ²Google Research, Switzerland ³University of Tübingen, Tübingen AI Center, Germany ⁴TUM University, Munich, Germany. Correspondence to: Michael Kirchhof <michael dot kirchhof at uni dash tuebingen dot de>.

tions of aleatoric uncertainty disentangled from epistemic uncertainty. This enables multiple use-cases: Beyond error prediction, we showcase novel applications like safe retrieval and uncertainty-aware dataset visualization. To facilitate widespread adoption, we release checkpoints for our pretrained uncertainties along with efficient code to pretrain them for arbitrary model architectures.

In summary, our contributions are:

- We develop a method which learns pretrained uncertainties that transfer zero-shot.
- This is based on fixing a gradient conflict in previous feed-forward uncertainties and speeding up the training by 180x, enabling large-scale pretraining (Section 3.3).
- Our uncertainties represent aleatoric uncertainty, disentangled from epistemic uncertainty (Section 4.5).
- We apply the uncertainties to improve the reliability of retrieval and to aid data visualization (Section 5).

2. Related Work

Large-scale uncertainty quantification. Uncertainty quantification has recently been scaled rapidly. Within one year, the largest vision models capable to perform uncertainty quantification grew from 1B (Tran et al., 2022) to 22B parameters (Dehghani et al., 2023). Benchmarks have increased accordingly. Previous surveys on out-of-distribution detection of 32x32 sized images (Ovadia et al., 2019) have been scaled to real-world images (Galil et al., 2023a) and to several additional tasks (Galil et al., 2023b). One such task is the uncertainty-aware representation learning (URL) benchmark that pretrains an uncertainty estimator and then tests zero-shot uncertainties on unseen datasets (Kirchhof et al., 2023b). Our work sets a new state-of-the-art on this task. In particular, we provide pretrained uncertainty modules for large computer vision models, independent from the classes or specific task of a dataset.

Feed-forward uncertainties. State-of-the-art models attempt to give such transferable uncertainties by moving away from classifier-layer uncertainties and towards uncertainties in the representation space (Collier et al., 2023). This approach falls under the category of feed-forward or deterministic uncertainties (Postels et al., 2022). They have a specialized uncertainty module that outputs predicted uncertainties during the forward pass of the model at minimal computational costs, which enables scaling. A variational take on this are probabilistic embeddings (Oh et al., 2019; Chun, 2023; Kim et al.; Nakamura et al., 2023) that output a variance estimate to give a distribution of possible representations instead of just one. This has recently been proven to recover the aleatoric uncertainty of the true posterior

(Kirchhof et al., 2023a) and improve retrieval performance (Karpukhin et al., 2022). As opposed to such indirect approaches, a second group of feed-forward approaches makes uncertainty quantification a direct regression task (Yoo & Kweon, 2019; Cui et al., 2023; Lahlou et al., 2023; Laves et al., 2020). In initial experiments, we found this direct approach to scale better. We use it as a starting point to develop our pretrained uncertainties in the next section, overcoming some remaining challenges to enable scaling.

3. Developing Pretrained Uncertainties

In this section we set out the desired properties of our pretrained uncertainties and then extend a popular feed-forward uncertainty method to satisfy these properties in the largescale pretraining case.

3.1. Basic Principles

We develop pretrained uncertainties from basic principles for scalability and ease of use, in order of importance:

- (i) Non-interference with primary task. Adding pretrained uncertainties to a model should not worsen the performance of the pretrained model's primary objective, e.g., its accuracy.
- (ii) Generalization. The predicted uncertainty estimates should reflect general forms of uncertainty that transfer to unseen datasets and tasks beyond the pretraining data and task.
- (iii) **Flexible adjustment**. The uncertainties should be general enough to adapt to new tasks and/or datasets that downstream practitioners might introduce.
- (iv) Minimal overhead. Providing uncertainties add only minimal runtime and memory usage to the main task prediction model.
- (v) Scalable optimization. Training should converge stably to ensure scalability to large pretraining corpora.

In summary, we seek a *download and forget* approach that requires minimal interventions from practitioners.

3.2. Recap: Loss Prediction

We now introduce a simple yet general uncertainty method that we build our method upon. From a decision theory perspective, giving an uncertainty estimate means predicting how wrong one thinks one's estimate is. The key is that any task's level of wrongness is defined by its loss \mathcal{L}_{task} . So in loss prediction (Yoo & Kweon, 2019; Cui et al., 2023; Lahlou et al., 2023; Laves et al., 2020), the model has an additional module u that predicts the model's own loss at

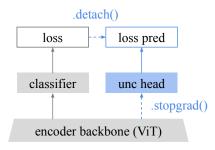


Figure 2. Pretrained uncertainties are returned by an auxiliary head (blue) that is trained to predict the classification loss of each image.

each of its predictions. This is learned via a L_2 loss between u and the (gradient-detached, det.) main task loss $\mathcal{L}_{\text{task}}^{\text{det.}}(y, f(x))$ at every sample (x, y). This is trained along with the main task (Kirchhof et al., 2023b), yielding the combined objective \mathcal{L} :

$$\mathcal{L} = \mathcal{L}_{\text{task}}(y, f(x)) + (u(e(x)) - \mathcal{L}_{\text{task}}^{\text{det.}}(y, f(x)))^{2}. \quad (1)$$

The uncertainty module u is implemented as a small MLP head u(e(x)) on top of the model representations e(x). This makes it cheap to compute during the forward pass, fulfilling the minimal overhead principle (iv). Loss prediction's uncertainties also adapt to any loss, fulfilling the flexibility principle (iii), and transfer well (Kirchhof et al., 2023b), fulfilling the generalization principle (ii).

Yet, loss prediction's implementation has a limitation: Figure 3a depicts a conflict between the uncertainty and the classification task. Their gradients interact negatively with one another, deteriorating the joint backbone and violating the non-interference principle (i). To resolve it, the current implementation stops early, roughly at epoch 12 in the plot. However, this early stopping is at odds with the scalability principle (v). Below, we fix these issues.

3.3. Enhancing Loss Prediction

We introduce four changes to the above loss prediction:

- **1. Stopgrad.** As visualized in Figure 2, we add a stopgrad behind the uncertainty module. This prevents its gradients from flowing to the backbone and interfering with the classification head. This strictly ensures the non-interference principle (i), and, indeed, the training now converges robustly, see Figure 3b. This way, uncertainties can be trained in parallel to the main task, as opposed to only in post-hoc.
- **2. No early stopping.** With the gradient conflict resolved, there is no more need for early stopping. The uncertainty head converges to its maximum at the end of the training in Figure 3b, making the training scalable as per principle (v).
- **3.** Cache everything. Since the classification head and backbone are now independent from the uncertainty head,

we pretrain and then freeze them before training the uncertainty module. Only the uncertainty objective remains:

$$\mathcal{L} = (u(e(x)) - \mathcal{L}_{\text{task}}^{\text{det.}}(y, f(x)))^2$$
 (2)

This can be optimized efficiently: The uncertainty module uses only the representations e(x) as inputs, and, likewise, the task loss depends only on them via f(x) = c(e(x)), where c is the classifier layer. So, we do not need to load the images x or run them through the backbone, but can cache the representations e(x) of the whole training process once (all epochs, including random augmentations). When learning the uncertainty module on top of a pretrained model, this increases the train speed by a factor of 180x and reduces the memory usage so far that we can pretrain uncertainties even for large models on single GPUs (or even CPUs). This paves the way for scalability: After caching the representations once, training the uncertainty module of a ViT-Large for seven ImageNet-21k-W epochs takes 2:26 hours on a single V100 GPU as opposed to 18 days with the standard loss prediction implementation.

4. Scale-free uncertainties. With the current L_2 loss, the uncertainty module is trained to match the scale of the pretraining loss. However, a downstream user might switch to a different loss on a different scale, which would introduce destructive gradients during finetuning. Thus, we switch to the ranking-based objective of Yoo & Kweon (2019):

$$\mathcal{L} = \max(0, \mathbb{1}_{\mathcal{L}} \cdot (u(e(x_1)) - u(e(x_2)) + m)), \tag{3}$$

s.t.
$$\mathbb{1}_{\mathcal{L}} := \begin{cases} +1, & \text{if } \mathcal{L}_{\text{task}}^{\text{det.}}(y_1, f(x_1)) > \mathcal{L}_{\text{task}}^{\text{det.}}(y_2, f(x_2)) \\ -1, & \text{else} \end{cases}$$
(4)

For every pair of images x_1 and x_2 , the indicator function compares which image has the higher primary task loss \mathcal{L}_{task}^{det} . Then, the uncertainty u of that sample is forced to be higher than that of the other sample, by a margin of at least m=0.1. This unties the uncertainty values from the scale of the task loss, improving on the flexibility principle (iii).

4. Experiments

We now study the main objective of pretrained uncertainties, their performance on downstream datasets. Additionally, we investigate which types of uncertainties they represent.

4.1. Experimental Setup

We are interested in how good our uncertainties perform on unseen datasets. This challenging zero-shot transfer is simple with our pretrained uncertainties, our enhanced loss

¹Being scale-free also means being uncalibrated. However, this is not a disadvantage for pretrained uncertainties, because during pretraining the downstream task is unknown, hence it is impossible to be calibrated for the unseen downstream task in the first place.

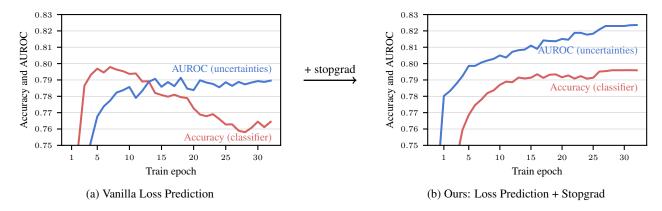


Figure 3. (a) The uncertainty and classification heads of Loss Prediction are in conflict. We solve this in (b) by adding a stopgrad. It ensures that the uncertainty head's gradients do not interfere with those of the classifier head, stabilizing the performance of both. The uncertainty and classifier heads were finetuned on ImageNet-1k on a pretrained (but unfrozen) ViT-Base backbone.

prediction module outputs an uncertainty u(x) for every image x from the downstream task. In order to measure the quality of these uncertainty estimates, we follow URL (Kirchhof et al., 2023b) in using the representation AUROC (R-AUROC) metric. It performs a 1-nearest neighbor classification on the representations e(x) on all images of a downstream dataset. The uncertainties u(x) should then be higher for images that are misclassified, which is quantified by the area under the ROC curve between the predicted uncertainties and whether or not the representation is correct in the sense that it is placed next to another representations of the same class. The R-AUROC can benchmark uncertainties when the classes are unseen during training, but if classes are seen during training, it is highly correlated with a conventional classification AUROC (Kirchhof et al., 2023b). In all experiments we also report the 1-nearest neighbor accuracy (Recall@1) from representation learning to quantify the retrieval performance of the representations and verify the non-interference principle (i) above.

We focus on Vision Transformers (Dosovitskiy et al., 2021) of several sizes and report results for the ViT-Base unless otherwise noted. Their backbone and classifiers were already pretrained by (Steiner et al., 2021) on ImageNet-21k-Winter-2021 (ImageNet-21k-W) (Deng et al., 2009) in timm (Wightman, 2019), so we only train the uncertainty module. As we show in Section 4.4, our approach is robust to architecture and optimizer hyperparameters, so we use the default values reported in Appendix A, inter alia a lightweight 2-layer MLP with width 512 for the uncertainty head. We train each model on five seeds and report the median as well as the distance to the maximum or minimum, whichever is larger, to provide an interpretable means to judge the variation.

As datasets, we use ImageNet-21k-W for pretraining and twelve datasets that span a variety of natural image domains for zero-shot transfer. Three are used in the URL benchmark, namely CUB-200-2011 (Wah et al., 2011), CARS196 (Krause et al., 2013), and Stanford Online Products (SOP) (Song et al., 2016). Seven are the natural images datasets of the Visual Task Adaption Benchmark (VTAB) (Zhai et al., 2020), namely Caltech101 (Fei-Fei et al., 2004), Oxford IIIT Pets (Parkhi et al., 2012), CIFAR100 (Krizhevsky, 2009), Scene Understanding 397 (SUN) (Xiao et al., 2010), Oxford Flowers 102 (Nilsback & Zisserman, 2008), Describable Textures (DTD) (Cimpoi et al., 2014), and Street View House Numbers (SVHN) (Netzer et al., 2011). The remaining two are CIFAR10 (Krizhevsky, 2009) and Treeversity#1 (Schmarje et al., 2022).

4.2. Pretrained Uncertainties Generalize

We first test the generalization principle (ii) on the twelve unseen datasets. Figure 1 shows that our pretrained uncertainties generalize well on eleven of the twelve datasets. The best R-AUROCs (Caltech101: 0.758 ± 0.006 , Oxford Pets: 0.740 ± 0.008 , and CIFAR10: 0.739 ± 0.002) are close to that of the pretraining dataset (0.791 \pm 0.001). In Appendix B we find that the zero-shot R-AUROC is higher on datasets that are closer to the domain spanned by ImageNet-21k-W, as one would expect from a pretrained model. This implies that further scaling the pretraining corpus, which is possible with our efficient training, may further benefit performance. The performance also depends on the granulatity of the zero-shot dataset. SVHN for example demands fine-grained house number disambiguation. It is harder to assign a pretrained uncertainty to such specialized tasks without knowing them in advance.

4.3. A New State-of-the-art

How do these results compare to the transfer performances of other methods in the field? The URL benchmark (Kirchhof et al., 2023b) has recently tested the R-AUROC of eleven

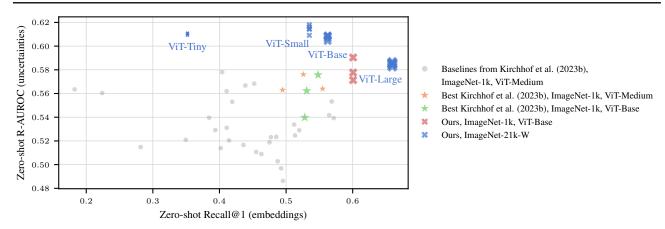


Figure 4. Our pretrained uncertainties outperform the approaches in the URL benchmark (Kirchhof et al., 2023b). The URL benchmark trained ViT-Mediums on ImageNet-1k. We reimplement its best approach (orange) on ViT-Base (green), then enhance it with our changes (red), and finally scale the training of ours to ImageNet-21k with various ViT sizes (blue). Each dot is one seed.



Figure 5. Pretrained uncertainties separate clear from ambiguous images on Stanford Online Products, a zero-shot dataset.

approaches on the CUB, CARS, and SOP datasets, averaged. URL tested the methods on ViT-Medium, a niche architecture that does not have ImageNet-21k checkpoints available. Thus, we switch to ViT-Base and reimplement URL's best performing method, the original loss prediction from Section 3.2. We reuse their codebase for compatibility.

Figure 4 shows the results, both in terms of R-AUROC and Recall@1. First, we find that our ViT-Base loss prediction reimplementation (red stars) achieves comparable performance to the original ViT-Medium backbone (orange stars). We then enhance the original loss prediction with our changes from Section 3.3 (green crosses). We find that the Recall@1 increases by 0.065 ± 0.012 because the backbone is no longer deteriorated by the uncertainty module gradients. The Recall@1 is now constant because we only train the uncertainty head anymore, keeping the pretrained backbone frozen. This does not restrict the R-AUROC of the uncertainty head, it in fact increases slightly by 0.021 ± 0.030 . Finally, we scale from ImageNet-1k to ImageNet-21k-W pretraining (blue crosses). This increases the R-AUROC again by 0.028 ± 0.014 on the ViT-Base. Because we now

use a different ImageNet-21k-W pretrained checkpoint, its Recall@1 changes. Upon training uncertainty modules for various ViT sizes, we find that they all have higher R-AUROC than the previous state-of-the-art. This demonstrates the generality of our approach.

4.4. Negative Results: Simple Beats Complex

Before we continue, we share some negative results. In Appendix C, we experiment with several techniques to further improve our method, including softening the loss function, uncertainty-induced training data augmentations, architecture changes, optimizer modifications, and initialization schemes. However, none of them significantly improve the uncertainties beyond the method presented in Section 3.3. Thus, to avoid adding unnecessary complexity, we decide to keep our approach as clean and simple as it is.

4.5. Pretrained Uncertainties \approx Aleatoric Uncertainty

If pretrained uncertainties work on unseen datasets, then which uncertainties do they capture? In this section, we

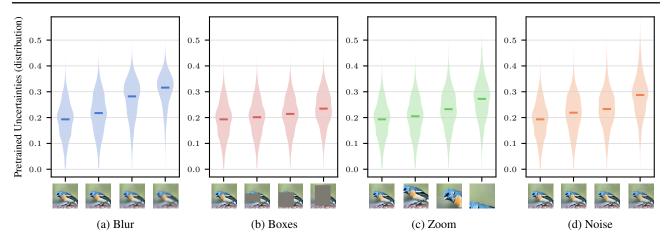


Figure 6. Pretrained uncertainties grow as images are deteriorated. Distributions and medians over unseen datasets (CUB, CARS, SOP). Note that except zooming, the pretrained model is not exposed to these data augmentations during training.

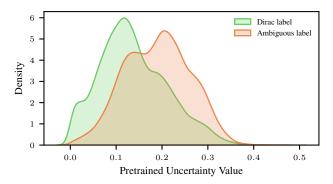


Figure 7. Pretrained uncertainties are systematically higher for ImageNet ReaL-H images with multiple possible labels.

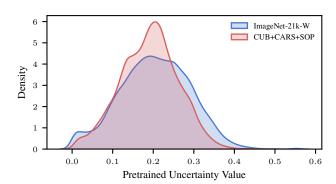


Figure 8. Pretrained uncertainties are consistent between train and unseen datasets, indicating the absence of epistemic uncertainty.

find that they model primarily aleatoric uncertainty and are mostly invariant to epistemic uncertainty. We conduct all analyses on unseen datasets and the ViT-Base model, unless otherwise noted.

To form a working hypothesis, we give some randomly selected examples with low and high predicted uncertainty in Figure 5. Although the network was not trained on this task, eBay product images, it correctly gives low uncertainties to clear and high uncertainties to ambiguous images. Not even a human expert or a Bayes classifier will be able to reduce the ambiguous images' uncertainty to zero, their ambiguity is intrinsic. This is known as aleatoric uncertainty. We hypothesize that pretrained uncertainties represent this form of uncertainty. Below, we investigate this hypothesis.

Human ambiguities. Aleatoric uncertainty is what is left even when an expert makes a prediction, for example a human annotator. So, we compare the model uncertainties to those of human annotators. ImageNet-1k ReaL-H (Beyer et al., 2020) re-collected labels for the 50,000 images in the

ImageNet-1k validation set.² While clear images kept their original Dirac label, annotators gave multiple or no labels to an image if it was ambiguous. Figure 7 shows that pretrained uncertainties are systematically higher on images the human annotators considered ambiguous (AUROC 0.701). This reinforces the aleatoric uncertainty hypothesis.

Interventional study. Second, we run an interventional experiment to induce aleatoric uncertainty. We deteriorate the images of the unseen datasets by blurring, overlaying with grey boxes, zooming in strongly, and adding Gaussian noise. Except zooming, these transformations were not applied during pretraining. Figure 6 shows that each of these transformations increases the pretrained uncertainties, the more strongly we deteriorate the images. This is additional evidence for the aleatoric uncertainty hypothesis.

No sign of epistemic uncertainty. We now consider the opposite hypothesis: Besides aleatoric uncertainty, do pre-

²Our pretraining dataset ImageNet-21k-W covers the classes of ImageNet-1k but neither its validation images nor any soft labels.

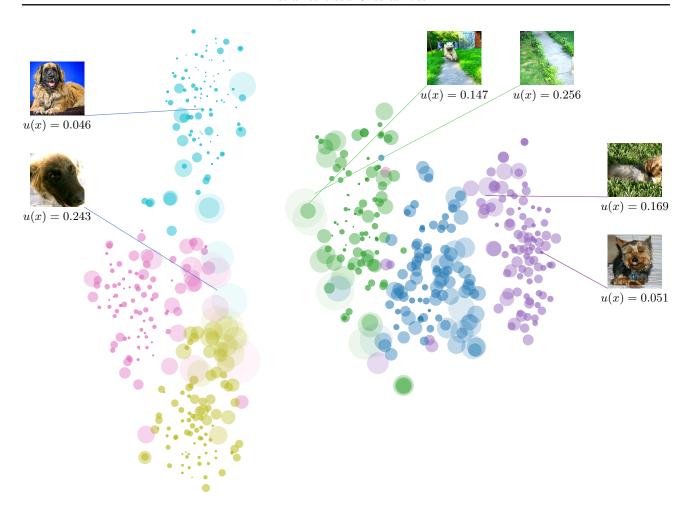


Figure 9. Visualizing pretrained uncertainties makes it easy to identify outliers in a tSNE plot. Images with high uncertainty u(x) are larger and transparent. Six classes of the zero-shot Oxford Pets dataset.

trained uncertainties comprise epistemic uncertainty? This uncertainty arises when a model has not seen an input before. Note that this would be detrimental to a pretrained uncertainty model, since it is intended to be used (exclusively) on unseen datasets where every image would be highly epistemically uncertain, drowning out the aleatoric signal. To test for epistemic uncertainty, we compare the pretrained uncertainties of in-distribution pretraining images to those of out-of-distribution images from unseen datasets. Figure 8 shows that the uncertainties on the pretraining data are similarly distributed to the unseen dataset (pairwise AUROC 0.503 ± 0.004). This suggests they are primarily capturing aleatoric uncertainty.

In summary, our results suggest that pretrained uncertainties quantify the amount of aleatoric uncertainty in an image, both in- and out-of-distribution, without being confounded by epistemic uncertainty. This constitutes significant progress for the ongoing efforts to disentangle epis-

temic from aleatoric uncertainty (Wimmer et al., 2023; Valdenegro-Toro & Mori, 2022).

5. Application Examples

In this section we showcase two applications that are unlocked by our new pretrained uncertainties.

5.1. Uncertainty-aware tSNE

Pretrained representations are often used to visualize datasets using methods like tSNE (van der Maaten & Hinton, 2008) or UMAP (McInnes et al., 2018). With pretrained uncertainties, we can now communicate inhererent ambiguities and explain outliers in these plots.

Figure 9 shows a tSNE visualization of six dog breeds in Oxford Pets, an unseen dataset. If an image has a high pretrained uncertainty, it is plotted as a larger and increasingly

transparent circle. The core region of each class cluster, typified by many small opaque circles, is now visually distinct from border regions and outliers, which are typified by collections of large transparent circles. By inspecting the images that underlie each representation, we verify that the core regions with low uncertainties comprise prototypical images whereas more uncertain images are often cropped out or in a camera angle that makes the exact dog breed ambiguous. We can also see that images lying in other classes' regions are often highly uncertain. Such misclassifications can be prevented by using our uncertainty-enhanced tSNE plots by allowing practitioners to understand and adjust the data preprocessing and filtering.

5.2. Safe Retrieval

This outlier identification can also be automated and utilized to make image retrieval more robust, enabling safe retrieval.

Consider again the Oxford Pets dataset in Figure 9. If we add a new dog image and search for its nearest neighbor, existing next-neighbor retrieval systems (Douze et al., 2024) may match it to an ambiguous image since these tend to lay at border regions. Similarly, if our new image itself is ambiguous, it is likely misplaced and existing systems will match it an arbitrary class. We can utilize pretrained uncertainties to tackle both of these problems by

- 1. Rejecting queries that are uncertain, and
- 2. Removing ambiguous images in the existing dataset, making it impossible to match to them.

As an example, we reject and/or clean the 10% most uncertain images per class in Oxford Pets. Table 1 shows that a typical cosine-distance based next-neighbour search achieves a Recall@1 of 0.772 ± 0.000 , or in other words a rate of 0.228 ± 0.000 wrong retrievals. Refusing to retrieve images when the input query is uncertain reduces this error rate by 14% to 0.196 ± 0.003 and cleaning the dataset from ambiguous images as potential retrieval partners reduces it by an additional 17%. These improvements are not only observed on unseen dataset but also on data that the retrieval system is familiar with. When using the ImageNet-1k validation set, whose classes were seen during ImageNet-21k-W pretraining (but whose validation images are unseen), deferring ambiguous queries reduces the error rate by 10% and cleaning the database reduces another 10%. All of these improvements are obtained automatically and fully unsupervised - we do not need to know the ground truth label of either the input or the database, since pretrained uncertainties can be computed for any input.

These are only first demonstrations of the opportunities that pretrained uncertainties offer. We anticipate further applications building up on pretrained uncertainties. For example,

1-NN error	Oxford Pets	ImageNet-1k
Full datasets	0.228 ± 0.000	0.382 ± 0.000
+ clean queries	0.196 ± 0.003	0.343 ± 0.001
+ clean database	0.163 ± 0.003	0.307 ± 0.001

Table 1. Pretrained uncertainties reduce the error rate of nextneighbour retrieval systems by rejecting ambiguous queries and/or removing ambiguous images from the database.

recent literature proposes retrieving a set of potential neighbours that is close to an ambiguous input with respect to its representation uncertainty (Kirchhof et al., 2023a). This can be implemented with pretrained uncertainty since they give uncertainties about representations. Similarly, conformal prediction (Angelopoulos & Bates, 2022) can view our pretrained uncertainties as a scoring function and calibrate its uncertainty predictions to downstream datasets. To facilitate future research, we provide all pretrained uncertainty checkpoints and code under the link in the abstract.

6. Conclusion

This work introduces a pretrained uncertainty module for computer vision models that is simple, cheap and scalable. We demonstrate its scalability by pretraining on ImageNet-21k-W. Our pretrained uncertainties method gives state-of-the-art zero shot uncertainty estimates on unseen datasets i.e. without finetuning. In future work, we anticipate scaling to even larger pretraining datasets as well as extending the method to pretraining objectives beyond classification and beyond the vision domain. We expect our fixes to the vanilla loss prediction method, that eliminate interference between uncertainty prediction and the main task, to also help in other feed-forward uncertainty quantifiers. By providing pretrained checkpoints, we intend to support applications similar to enhanced visualization and safe retrieval.

Impact Statement

Our uncertainties are intended to capture errors before they happen and reveal uncertainties that would remain undetected when images are solely expressed as representation vectors. This makes models more safe and trustworthy by allowing them to fulfill their tasks with less errors. We enable an easier access to these uncertainties by our ease-of-use principles and providing plug-and-play checkpoints. We see this as a positive impact on both the community and society. As with all general-purpose machine learning advancements, this assumes that a practitioner does not develop a model with a harmful task, which is beyond our sphere of influence. Additionally, we encourage researchers to follow our example of finding ways to significantly reduce training costs for a lower energy consumption during training. We provide our code standalone to help start these efforts.

References

- Angelopoulos, A. N. and Bates, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv* preprint arXiv:2107.07511, 2022.
- Beyer, L., Hénaff, O. J., Kolesnikov, A., Zhai, X., and Oord, A. v. d. Are we done with ImageNet? *arXiv preprint arXiv:2006.07159*, 2020.
- Chen, X., Liang, C., Huang, D., Real, E., Wang, K., Liu, Y., Pham, H., Dong, X., Luong, T., Hsieh, C.-J., et al. Symbolic discovery of optimization algorithms. *arXiv* preprint arXiv:2302.06675, 2023.
- Chun, S. Improved probabilistic image-text representations. *arXiv preprint arXiv:2305.18171*, 2023.
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and Vedaldi, A. Describing textures in the wild. In *Computer Vision and Pattern Recognition (CVPR)*, 2014.
- Collier, M., Jenatton, R., Mustafa, B., Houlsby, N., Berent, J., and Kokiopoulou, E. Massively scaling heteroscedastic classifiers. *arXiv preprint arXiv:2301.12860*, 2023.
- Cui, P., Zhang, D., Deng, Z., Dong, Y., and Zhu, J. Learning sample difficulty from pre-trained models for reliable prediction. In *Neural Information Processing Systems* (NeurIPS), 2023.
- Dehghani, M., Djolonga, J., Mustafa, B., Padlewski, P., Heek, J., Gilmer, J., Steiner, A. P., Caron, M., Geirhos, R., Alabdulmohsin, I., et al. Scaling vision transformers to 22 billion parameters. In *International Conference on Machine Learning (ICML)*, 2023.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. ImageNet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition (CVPR)*, 2009.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn,
 D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer,
 M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby,
 N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvasy, G., Mazaré, P.-E., Lomeli, M., Hosseini, L., and Jégou, H. The Faiss library. *arXiv preprint arXiv:2401.08281*, 2024.
- Fei-Fei, L., Fergus, R., and Perona, P. Learning generative visual models from few training examples: An incremental Bayesian approach tested on 101 object categories. Computer Vision and Pattern Recognition Workshop, 2004.

- Galil, I., Dabbah, M., and El-Yaniv, R. A framework for benchmarking class-out-of-distribution detection and its application to ImageNet. In *The Eleventh International* Conference on Learning Representations (ICLR), 2023a.
- Galil, I., Dabbah, M., and El-Yaniv, R. What can we learn from the selective prediction and uncertainty estimation performance of 523 ImageNet classifiers? In *Interna*tional Conference on Learning Representations (ICLR), 2023b.
- Hendrycks, D., Mu, N., Cubuk, E. D., Zoph, B., Gilmer, J., and Lakshminarayanan, B. AugMix: A simple method to improve robustness and uncertainty under data shift. In *International Conference on Learning Representations* (ICLR), 2020.
- Karpukhin, I., Dereka, S., and Kolesnikov, S. Probabilistic embeddings revisited. *arXiv preprint arXiv:2202.06768*, 2022.
- Kim, E., Jung, D., Park, S., Kim, S., and Yoon, S. Probabilistic concept bottleneck models. In *International Conference on Machine Learning (ICML)*.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*, 2015.
- Kirchhof, M., Roth, K., Akata, Z., and Kasneci, E. A non-isotropic probabilistic take on proxy-based deep metric learning. In *European Conference on Computer Vision* (*ECCV*), 2022.
- Kirchhof, M., Kasneci, E., and Oh, S. J. Probabilistic contrastive learning recovers the correct aleatoric uncertainty of ambiguous inputs. *International Conference on Machine Learning (ICML)*, 2023a.
- Kirchhof, M., Mucsányi, B., Oh, S. J., and Kasneci, E. Url: A representation learning benchmark for transferable uncertainty estimates. *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, 2023b.
- Krause, J., Stark, M., Deng, J., and Fei-Fei, L. 3D object representations for fine-grained categorization. In *Conference on Computer Vision and Pattern Recognition* (CVPR) Workshop, 2013.
- Krizhevsky, A. Learning multiple layers of features from tiny images. Technical report, 2009.
- Lahlou, S., Jain, M., Nekoei, H., Butoi, V. I., Bertin, P., Rector-Brooks, J., Korablyov, M., and Bengio, Y. DEUP: Direct epistemic uncertainty prediction. *Transactions on Machine Learning Research (TMLR)*, 2023. ISSN 2835-8856.

- Laves, M.-H., Ihler, S., Fast, J. F., Kahrs, L. A., and Ortmaier, T. Well-calibrated regression uncertainty in medical imaging with deep learning. In *Medical Imaging with Deep Learning*, pp. 393–412. PMLR, 2020.
- Liu, J., Lin, Z., Padhy, S., Tran, D., Bedrax Weiss, T., and Lakshminarayanan, B. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. Advances in Neural Information Processing Systems (NeurIPS), 2020.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv* preprint arXiv:1711.05101, 2017.
- McInnes, L., Healy, J., Saul, N., and Großberger, L. UMAP: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29):861, 2018. doi: 10.21105/joss.00861. URL https://doi.org/10.21105/joss.00861.
- Mucsányi, B., Kirchhof, M., Nguyen, E., Rubinstein, A., and Oh, S. J. Trustworthy machine learning, 2023.
- Nakamura, H., Okada, M., and Taniguchi, T. Representation uncertainty in self-supervised learning as variational inference. In *International Conference on Computer Vision* (ICCV), 2023.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.
- Nilsback, M.-E. and Zisserman, A. Automated flower classification over a large number of classes. In *Proceedings* of the Indian Conference on Computer Vision, Graphics and Image Processing, 2008.
- Oh, S. J., Gallagher, A. C., Murphy, K. P., Schroff, F., Pan, J., and Roth, J. Modeling uncertainty with hedged instance embeddings. In *International Conference on Learning Representations (ICLR)*, 2019.
- Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., Dillon, J., Lakshminarayanan, B., and Snoek, J. Can you trust your model's uncertainty? Evaluating predictive uncertainty under dataset shift. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019.
- Parkhi, O. M., Vedaldi, A., Zisserman, A., and Jawahar,C. V. Cats and dogs. In *Computer Vision and Pattern Recognition (CVPR)*, 2012.
- Postels, J., Segù, M., Sun, T., Sieber, L. D., Van Gool, L., Yu, F., and Tombari, F. On the practicality of deterministic epistemic uncertainty. In *International Conference on Machine Learning (ICML)*, 2022.

- Schmarje, L., Grossmann, V., Zelenka, C., Dippel, S., Kiko, R., Oszust, M., Pastell, M., Stracke, J., Valros, A., Volkmann, N., et al. Is one annotation enough? A data-centric image classification benchmark for noisy and ambiguous label estimation. *arXiv preprint arXiv:2207.06214*, 2022.
- Song, H. O., Xiang, Y., Jegelka, S., and Savarese, S. Deep metric learning via lifted structured feature embedding. In *Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Steiner, A., Kolesnikov, A., , Zhai, X., Wightman, R., Uszkoreit, J., and Beyer, L. How to train your ViT? Data, augmentation, and regularization in vision transformers. *arXiv* preprint arXiv:2106.10270, 2021.
- TorchVision. Torchvision: Pytorch's computer vision library. *GitHub repository:* https://github.com/pytorch/vision, 2016.
- Tran, D., Liu, J., Dusenberry, M. W., Phan, D., Collier, M., Ren, J., Han, K., Wang, Z., Mariet, Z., Hu, H., et al. Plex: Towards reliability using pretrained large model extensions. *arXiv preprint arXiv:2207.07411*, 2022.
- Valdenegro-Toro, M. and Mori, D. S. A deeper look into aleatoric and epistemic uncertainty disentanglement. In *Computer Vision and Pattern Recognition Workshops* (CVPRW), 2022.
- van der Maaten, L. and Hinton, G. Visualizing data using t-SNE. *Journal of Machine Learning Research (JMLR)*, 9(86):2579–2605, 2008.
- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The Caltech-UCSD birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
- Wightman, R. PyTorch image models, 2019.
- Wimmer, L., Sale, Y., Hofman, P., Bischl, B., and Hüllermeier, E. Quantifying aleatoric and epistemic uncertainty in machine learning: Are conditional entropy and mutual information appropriate measures? In *Uncertainty in Artificial Intelligence (UAI)*, 2023.
- Xiao, J., Hays, J., Ehinger, K. A., Oliva, A., and Torralba, A. Sun database: Large-scale scene recognition from abbey to zoo. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- Yoo, D. and Kweon, I. S. Learning loss for active learning. In Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. Cutmix: Regularization strategy to train strong classifiers

- with localizable features. In *International Conference on Computer Vision (ICCV)*, 2019.
- Zhai, X., Puigcerver, J., Kolesnikov, A., Ruyssen, P., Riquelme, C., Lucic, M., Djolonga, J., Pinto, A. S., Neumann, M., Dosovitskiy, A., Beyer, L., Bachem, O., Tschannen, M., Michalski, M., Bousquet, O., Gelly, S., and Houlsby, N. A large-scale study of representation learning with the visual task adaptation benchmark, 2020.
- Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. Mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations (ICLR)*, 2018.

A. Training details

Architecture. With d_e denoting the dimensionality of the (flattened) embeddings e(x) of each ViT size, our pretrained uncertainty module has the following size across all ViT sizes: Linear $(d_e$, 512), LeakyReLU (negative slope 0.01), Linear (512, 512), LeakyReLU (negative slope 0.01), Linear (512, 1), Softplus $(\beta=1)$, threshold=20). The softplus in the end is to ensure that all uncertainties are strictly positive. This could be dropped since our uncertainties are scale free, but we added it for convenience of interpretation.

Optimizer. We train on ImageNet-21k-W for 460 episodes of 200,000 images, corresponding to roughly seven full epochs. We use a cosine learning rate scheduler that warms up the learning rate from 0.0001 to 0.0028 for 25 episodes and then decays it down to 1e-8 for the remaining episodes. We use an AdamW (Loshchilov & Hutter, 2017) optimizer with $\beta_1 = 0.8$ and $\beta_2 = 0.95$. We apply weight decay of strength 0.0001. These settings are constant for all experiments, without any hyperparameter tuning.

Augmentations. We use the torchvision (TorchVision, 2016) augmentations that timm applies by default. Inter alia, all images are cropped to 224x224 pixels. We first apply a RandomResizedCropAndInterpolation (size=(224, 224), scale=(0.08, 1.0), ratio=(0.75, 1.3333), interpolation=bilinear bicubic), and then randomly add RandomHorizontalFlip (p=0.5) and with p = 0.4 a ColorJitter (brightness=[0.6, 1.4], contrast=[0.6, 1.4], saturation=[0.6, 1.4], hue=None).

B. Transfer analysis

In this section, we analyze which datasets our pretrained uncertainties transfer to. We hypothesize that they behave similarly to the classifier head for ImageNet-21k. To test this, we use the entropy of the 21k class predictions of the classifier head as uncertainties and test its R-AUROC. We find that our pretrained uncertainties achieve a similar performance on all datasets. This indicates that pretrained uncertainties works on datasets similar enough to ImageNet-21k that its classifier is also informative. Note that this classifier method is not applicable to provide pretrained uncertainties in practice since it requires maintaining a heavy classifier head (17M parameters for ViT-Base), violating principle (iv), is not scalable to datasets with more classes, violating principle (v), and is not available outside ImageNet-21k classification, violating principle (iii).

Dataset	Pretrained Uncertainties	21k Classifier Entropy
ImageNet-21k	0.791 ± 0.001	0.798
Caltech 101	0.758 ± 0.006	0.808
Oxford Pets	0.740 ± 0.008	0.724
CIFAR 10	0.739 ± 0.002	0.716
CIFAR 100	0.706 ± 0.002	0.696
SUN	0.691 ± 0.002	0.697
Oxford Flowers	0.659 ± 0.009	0.679
Describable Textures	0.649 ± 0.006	0.610
CUB 200	0.626 ± 0.008	0.608
Stanford Online Products	0.607 ± 0.001	0.591
CARS 196	0.589 ± 0.003	0.554
Treeversity	0.560 ± 0.003	0.565
SVHN	0.495 ± 0.005	0.524

Table 2. Pretrained uncertainties perform similarly to using the entropy of the ImageNet-21k classifier head as uncertainty estimate (note that this is impractical due to its size and violating the first principles in Section 3.1). This implies that pretrained uncertainties cover roughly the classes that ImageNet-21k also covers.

C. Simple beats Complex: Negative Results

We test multiple adjustments to our loss, architecture and optimizer in Table 3. We report the average R-AUROC on the unseen datasets CUB, CARS, and SOP as in the URL protocol, evaluated for five seeds on a ViT-Base pretrained on ImageNet-21k-W.

The first change regards the loss function. Currently, when comparing two images, it always requires one image to have a pretrained uncertainty of at least 0.1 larger than the other one. In the formula below, we add an 'approximately equal'

category, such that images whose ground-truth loss is within a leeway l are not required to have different loss values:

$$\mathbb{1}_{\mathcal{L}} := \begin{cases}
+1 &, \text{ if } \mathcal{L}_{\text{task}}^{\text{det.}}(y_1, f(x_1)) > l + \mathcal{L}_{\text{task}}^{\text{det.}}(y_2, f(x_2)) \\
-1 &, \text{ if } \mathcal{L}_{\text{task}}^{\text{det.}}(y_1, f(x_1)) + l < \mathcal{L}_{\text{task}}^{\text{det.}}(y_2, f(x_2)) \\
0 &, \text{ else}
\end{cases} \tag{5}$$

However, at several values of the allowed leeway l, this does not change the performance outside the margin of error of the baseline method (0.608 ± 0.004) .

Second, we change the size of the uncertainty head MLP. By default it has 2 hidden layers of width 512. We either shrink it to 1 hidden layer with width 256 or enlargen it to 3 hidden layers of width 1024. While there is a slight trend favoring smaller heads, it does not exceed the margin of chance.

Third, we briefly experimented with initializing the uncertainty module with zero-weights. However, this failed to train at all, which is theoretically expected.

Fourth, we add strong augmentations that add different types of aleatoric uncertainty to half of the train dataset. None of these increase the performance, with some even deteriorating it. While this might seem counterintuitive, we presume such artificial sources of uncertainty do not reflect the uncertainties occuring on real images.

Last, we experiment with optimizers other than our default AdamW with cosine learning rate scheduler. While Lion collapses after less than one epoch, SGD performs slightly better than the baseline. However, it might be a false positive, especially taking multiple testing into account. Indeed, the reason we did not select SGD for the main paper is that it did not systematically outperform AdamW during our preliminary experiments on the validation splits with less seeds. The test splits were held strictly secret until the writing of the paper. We suggest future researchers to experiment with replacing their advanced optimizers by SGD.

Method	R-AUROC
Default	0.608 ± 0.004
Softened loss $(l = 0.001)$	0.607 ± 0.006
Softened loss ($l = 0.01$)	$\boldsymbol{0.607} \pm 0.005$
Softened loss $(l = 0.1)$	0.609 ± 0.005
Smaller uncertainty module	0.611 ± 0.002
Larger uncertainty module	0.606 ± 0.004
Initialize uncertainty module with zero	0.500 ± 0.000
AugMix (Hendrycks et al., 2020) for 50% of train data	0.606 ± 0.005
CutMix (Yun et al., 2019) for 50% of train data	0.575 ± 0.002
MixUp (Zhang et al., 2018) for 50% of train data	0.552 ± 0.009
Blurred images for 50% of train data	0.609 ± 0.003
Small crops for 50% of train data	0.610 ± 0.004
Box overlay for 50% of train data	0.606 ± 0.001
Adam optimizer (Kingma & Ba, 2015)	0.609 ± 0.002
Lion optimizer (Chen et al., 2023)	0.500 ± 0.000
SGD optimizer	0.614 ± 0.005
Step learning rate scheduler	0.607 ± 0.003

Table 3. No change to loss, architecture, optimizer, or data augmentation improves the performance. R-AUROC averaged across CUB, CARS, and SOP as in the URL protocol, for five seeds on a ViT-Base.