

Benchmarking Uncertainty Disentanglement: Specialized Uncertainties for Specialized Tasks

Bálint Mucsányi¹ Michael Kirchhof¹ Seong Joon Oh^{1,2}

Abstract

Uncertainty quantification, once a singular task, has evolved into a spectrum of tasks, including abstained prediction, out-of-distribution detection, and aleatoric uncertainty quantification. The latest goal is disentanglement: the construction of multiple estimators that are each tailored to one and only one task. Hence, there is a plethora of recent advances with different intentions—that often entirely deviate from practical behavior. This paper conducts a comprehensive evaluation of numerous uncertainty estimators across diverse tasks on ImageNet. We find that, despite promising theoretical endeavors, disentanglement is not yet achieved in practice. Additionally, we reveal which uncertainty estimators excel at which specific tasks, providing insights for practitioners and guiding future research toward task-centric and disentangled uncertainty estimation methods. Our code is available at <https://github.com/bmucsanyi/bud>.

1. Introduction

When uncertainty quantification methods were first pioneered for deep learning (Gal & Ghahramani, 2016; Lakshminarayanan et al., 2017), their task was simple: giving one total uncertainty. The recent demand for trustworthy machine learning (Mucsányi et al., 2023) created new requirements, mostly centering around disentangling the above predictive uncertainty into aleatoric (data-inherent and irreducible) and epistemic (model-centric and reducible) components (Depeweg et al., 2018; Valdenegro-Toro & Mori, 2022; Shaker & Hüllermeier, 2021). They serve different purposes: epistemic uncertainty is widely used for out-of-distribution detection (Van Amersfoort et al., 2020), and two estimators that each estimate one and only one component in a disentangled manner enable tasks like active learning

¹University of Tübingen, Germany ²Tübingen AI Center, Germany. Correspondence to: Bálint Mucsányi <b dot h dot mucsanyi at gmail dot com>.

Preliminary work.

Method ranking	Correctness	Abstinance	Log Prob.	Brier	Aleatoric	ECE	OOD
Deep Ensemble	2	1	1	1	1	4	6
Dropout	1	2	2	2	2	3	2
Baseline	6	3	5	3	3	7	9
SNGP	5	8	4	6	7	1	4
GP	4	6	3	4	6	2	5
Mahalanobis	11	11	–	–	11	–	1
Shallow Ensemble	3	5	6	5	5	5	3
Laplace	7	7	7	7	9	6	11
HET-XL	8	4	8	8	4	8	10
Correctness Pred.	9	9	9	9	8	9	7
Loss Prediction	10	10	–	–	10	–	8

Table 1. Different tasks have different best-performing uncertainty quantification methods on ImageNet-Real. Dropout and deep ensemble are good choices across the board (but expensive, see Appendix K.3). Best, second-best, and third-best method highlighted in gold, silver, and bronze. Beware that differences between ranks can be very small, see the per-task plots for details.

(Lahlou et al., 2023).

One limitation of these recent endeavors is that they are primarily theoretical, supported by toy or small-scale experiments (Shaker & Hüllermeier, 2021; Van Amersfoort et al., 2020; Mukhoti et al., 2023). Larger scale benchmarks often evaluate tasks for only one component and do not test for undesirable side effects on the other component (Galil et al., 2023a; Ovadia et al., 2019). While this approach allows for insights into the performance of a subset of methods on a selection of tasks, there is currently no study that evaluates which component(s) each method captures in practice and which it does not.

Our work establishes an overview of this vast landscape of methods and tasks. We reimplement twelve uncertainty quantification estimators in up to eight ways and evaluate each on seven practically defined tasks on ImageNet-1k (Deng et al., 2009), ranging from abstained prediction to out-of-distribution detection. We further study if recent uncertainty decomposition formulas decompose the estima-

tors into disentangled components as theoretically intended (Wimmer et al., 2023; Pfau, 2013; Depeweg et al., 2018). We find that disentanglement is unachieved in practice since most proposed combinations of estimators are highly internally correlated and fail to unmix aleatoric and epistemic uncertainty into two components (Section 3.1). However, we find that there are groups of approaches specialized on individual tasks, like density-based approaches outperforming the others on out-of-distribution detection (epistemic uncertainty, Section 3.2) but performing close to random on aleatoric and predictive uncertainty tasks (Section 3.3).

These findings emphasize the importance of specifying the particular task one wants to solve and developing uncertainty estimators tailored to it. We anticipate that our insights into the practical workings of uncertainty estimators will drive the field of uncertainty quantification toward developing robust and disentangled uncertainty estimators.

2. Benchmarked Methods

This section provides an overview of current uncertainty estimators and disentanglement formulas we benchmark. We reimplement all methods as plug-and-play modules that will be released after the anonymity period. Details are provided in Appendix A.

2.1. Uncertainty Estimators

We consider a classification setting with a discrete label space \mathcal{Y} of C classes and models $f: \mathcal{X} \rightarrow \Delta^C$ that output a class probability vector on the probability simplex for any input $x \in \mathcal{X}$. The uncertainty estimators can be categorized into two classes: distributional and deterministic methods.

2.1.1. DISTRIBUTIONAL METHODS

Distributional methods output a probability distribution $q(f(x) | x)$ over all possible class probability vectors, abbreviated as $q(f)$. This distribution can, e.g., correspond to a Bayesian hypothesis posterior $p(f | \mathcal{D})$ induced by a parameter posterior $p(\theta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta)$ when training on a dataset \mathcal{D} .

Spectral-normalized Gaussian processes (SNGP) (Liu et al., 2020) obtain these distributions by approximating a Gaussian process over the classifier *output*, aided by spectral normalization on all network parameters. We also benchmark the last-layer Gaussian process without spectral normalization, denoted as GP. The **Laplace approximation** (Daxberger et al., 2021) approximates a Gaussian posterior over the network *parameters* using an efficient Hessian approximation. This is a post-hoc method applied to a point estimate network, allowing to draw multiple outputs per input. **Latent heteroscedastic classifiers (HET-XL)** (Collier et al., 2023) predict a heteroscedastic Gaussian distribution

over the pre-logit *embeddings* and sample multiple embeddings that get turned into class probability vectors.

Dropout (Srivastava et al., 2014) and **deep ensembles** (Lakshminarayanan et al., 2017) do not construct distributions $q(f)$ but directly sample from them, either by M repeated forward passes, or by training M models, respectively. **Shallow ensembles** (Lee et al., 2015) are lightweight approximations of deep ensembles. They use a shared backbone and M output heads (often referred to as “experts”). With a single forward pass, one obtains M logit vectors per input. As a baseline, we also benchmark a **deterministic** network that corresponds to a Dirac posterior in the parameter space.

Practical tasks like threshold-based rejection often need a scalar uncertainty output $u(x) \in \mathbb{R}$ instead of a distribution $q(f)$. To this end, **aggregators** compile the upper distributions into scalar uncertainty estimates $u(x) \in \mathbb{R}$. There are several methods for this aggregation, e.g., calculating the Bayesian Model Average (BMA) $\hat{f}(x) := \mathbb{E}_{q(f)} [f(x)]$ and using its entropy as the uncertainty estimate $u(x)$ or quantifying the variance of $q(f)$. We consider eight aggregators detailed in Appendix B and, unless stated otherwise, use the best-performing one for each distributional method.

2.1.2. DETERMINISTIC METHODS

Deterministic methods (Postels et al., 2022) directly output a scalar uncertainty estimate $u(x)$ instead of modeling a probability distribution over class probability vectors.

Loss prediction (Yoo & Kweon, 2019; Lahlou et al., 2023; Kirchhof et al., 2023b) employs an additional MLP head for $u(x)$ that estimates the loss of the network’s prediction $f(x)$ on each input x , assuming that the loss reflects a notion of (in-)correctness. We also implement a special variant for classification, **correctness prediction**, where $u(x)$ predicts how likely the predicted class $\hat{y} := \arg \max_{c \in \{1, \dots, C\}} f_c(x)$ is to be the correct class y , i.e., $p(\hat{y} = y)$.

Deterministic uncertainty quantification (DUQ) (Van Amersfoort et al., 2020) learns a latent mixture-of-RBF density on the training set and outputs as $u(x)$ how close an input’s embedding is to the mixture means. The **Mahalanobis** method (Lee et al., 2018) builds a similar latent mixture of Gaussians in a post-hoc fashion. It also perturbs the inputs adversarially to train a classifier for separating in-distribution (ID) and out-of-distribution (OOD) samples. This is the only method in our benchmark that requires a validation set for training the logistic regression OOD detector.

2.2. Uncertainty Disentanglement

The previous methods all give one general uncertainty estimate. A second strain of literature outputs not only one estimate but decomposes a posterior $q(f)$ (obtained by any

of the methods above) into multiple estimators, intending to quantify different forms of uncertainty, such as epistemic and aleatoric uncertainty (Hora, 1996). Epistemic uncertainty is caused by a lack of data and can be reduced as one gathers more information. In contrast, aleatoric uncertainty is due to the randomness inherent in the data-generating process itself and is irreducible (Mucsányi et al., 2023). The estimators for each source should be disentangled: the aleatoric estimator should only reflect aleatoric uncertainty, and the epistemic estimator should reflect only epistemic uncertainty. See Appendix E for more details. We benchmark two prominent approaches to obtain such pairs of estimators.

2.2.1. INFORMATION-THEORETICAL DECOMPOSITION

The information-theoretical (IT) decomposition (Depeweg et al., 2018; Shaker & Hüllermeier, 2021; Mukhoti et al., 2021; Wimmer et al., 2023) decompose the entropy of the predictive distribution $p(y | x) = \int p(y | x, f) dq(f)$ into an aleatoric and an epistemic component:

$$\underbrace{\mathbb{H}_{p(y|x)}(y)}_{\text{predictive}} = \underbrace{\mathbb{E}_{q(f)} [\mathbb{H}_{p(y|x,f)}(y)]}_{\text{aleatoric}} + \underbrace{\mathbb{I}_{p(y,f|x)}(y; f)}_{\text{epistemic}}, \quad (1)$$

where $\mathbb{H}_{p(y|x)}(y) \equiv \mathbb{H}(Y|x)$ is the entropy and $\mathbb{I}_{p(y,f|x)}(y; f) \equiv \mathbb{I}(Y; F|x)$ is the mutual information. Intuitively, the aleatoric component gives the spread of the labels that the plausible models in the posterior have on average, whereas the epistemic component only captures the disagreement of the prediction $p(y | x, f)$ between the models f . Since most posteriors $q(f)$ are practically implemented as mixtures of M Diracs $q(f) \approx \{f^{(m)}(x)\}_{m=1}^M$, we show how the aleatoric and epistemic estimates are computed in this special case of Equation (1) in Appendix C.

2.2.2. BREGMAN DECOMPOSITIONS

Bregman decompositions (Pfau, 2013; Gupta et al., 2022; Lahlou et al., 2023; Gruber & Buettner, 2023) use not only the posterior $q(f)$, which is internally computed by each method, but also take the ground-truth generative process $p^{\text{gt}}(x, y)$ into account. Bregman decompositions then decompose the expected loss of a model over all possible training datasets, where the loss D_F is a Bregman divergence like the Euclidean distance or the Kullback-Leibler divergence.

$$\underbrace{\mathbb{E}_{q(f), p^{\text{gt}}(y|x)} [D_F [y \| f(x)]]}_{\text{predictive}} = \underbrace{\mathbb{E}_{p^{\text{gt}}(y|x)} [D_F [y \| f^*(x)]]}_{\text{aleatoric}} + \underbrace{\mathbb{E}_{q(f)} [D_F [\bar{f}(x) \| f(x)]]}_{\text{epistemic}} + \underbrace{D_F [f^*(x) \| \bar{f}(x)]}_{\text{bias}} \quad (2)$$

Since $f^*(x) = \mathbb{E}_{p^{\text{gt}}(y|x)} [y]$ is the Bayes predictor, the

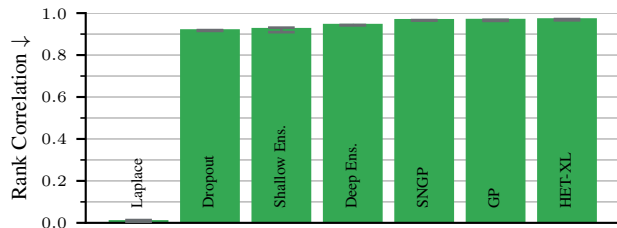


Figure 1. Six out of seven distributional methods exhibit a severely high rank correlation between the information-theoretical aleatoric and epistemic components when evaluated on ImageNet-Real. These methods violate a necessary condition of uncertainty disentanglement.

aleatoric uncertainty is the Bayes risk of the generative process, which is by definition irreducible and independent of the posterior $q(f)$. As this process is unknown in practice, we estimate the aleatoric term by $\mathbb{E}_{q(f)} [\mathbb{H}_{p(y|f,x)}(y)]$. The epistemic uncertainty is, similarly to the IT decomposition, the average distance of the posterior members $f \sim q(f)$ from their centroid $\bar{f}(x) = \arg \min_z \mathbb{E}_{q(f)} [D_F [z \| f(x)]]$. This average is calculated in a dual space, but in certain cases is equal to the BMA (Gupta et al., 2022). To make the decomposition equality complete, the Bregman decomposition has a third term, the bias. This is an uncertainty source that subsumes the uncertainty about the function class (Von Luxburg & Schölkopf, 2011).

The popular DEUP risk decomposition (Lahlou et al., 2023) is a special case of the Bregman decomposition. We provide details and evaluation in Appendix H.

3. Experiments

With these different estimators and pairs of estimators at hand, we now investigate our main research questions: Does any approach give disentangled uncertainty estimators in practice? Furthermore, what type of uncertainty does each estimator capture in terms of practical tasks?

We reimplement and train each approach on a pretrained ResNet-50 for 50 epochs on ImageNet-1k (Deng et al., 2009) with a training pipeline following Tran et al. (2022). Since the DUQ and Mahalanobis methods have memory and stability issues on ImageNet, in Section 3.7, we repeat experiments on CIFAR-10 (Krizhevsky & Hinton, 2009) with the Wide ResNet 28-10 architecture, following Liu et al. (2020). We track the validation performance to conduct early stopping and to choose further hyperparameters of each method. We report mean, minimum, and maximum performance across five seeds. The benchmark took 1 GPU year on RTX 2080 Ti GPUs.

Benchmarking Uncertainty Disentanglement

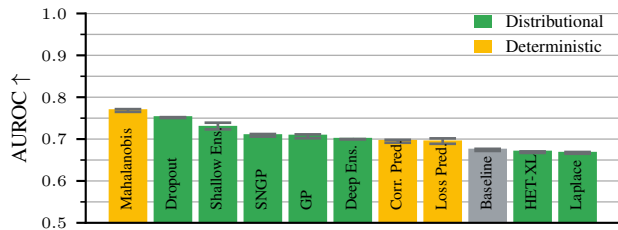


Figure 2. Mahalanobis—a direct OOD detector, dropout, and shallow ensembles distinguish ID and OOD samples considerably better (AUROC ≥ 0.728) than the baseline (AUROC = 0.674). OOD samples are perturbed by ImageNet-C corruptions of severity level two.

3.1. Uncertainty disentanglement often fails

We first study if the IT and Bregman decompositions yield pairs of disentangled estimators. Since they can only decompose the posteriors $q(f)$ of distributional methods, deterministic methods are excluded in this section.

Figure 1 reveals a simple failure: for six of the seven distributional methods, the IT decomposition leads to highly mutually correlated aleatoric and epistemic estimates (rank corr. ≥ 0.92). This correlation remains for Bregman (Appendix H) and does not considerably lower even when we add more epistemic uncertainty into the dataset (Appendix G.1). Therefore, in the majority of cases, disentanglement is violated in practice.

The only exception is the Laplace posterior, whose aleatoric and epistemic estimators are entirely decorrelated. This shows that disentanglement cannot be thought about only on the high level of decompositions but needs to take the exact distributional method into account. However, in the following sections, we will see that Laplace’s decomposed estimators do not sufficiently capture the aleatoric and epistemic ground truth uncertainties.

3.2. OOD-ness is hard to detect

Let us begin by testing which estimators represent epistemic uncertainty. We measure this via an out-of-distribution (OOD) detection task (Gruber & Buettner, 2023; Mukhoti et al., 2021). We use ImageNet-C (Hendrycks & Dietterich, 2019) with corruptions of severity level two as OOD data. This is far enough out-of-distribution to deteriorate the accuracy by 26%, see also Section 3.5. The ground-truth epistemic uncertainty is high on these samples, so we measure via a binary classification AUROC if uncertainty estimators are higher on these OOD samples than on ID samples. From this point forward, we also consider deterministic methods.

As illustrated in Figure 2, Laplace is the least able to detect epistemic uncertainties. Notably, Figure 2 uses the best aggregator available, not just the one that the IT or Bregman

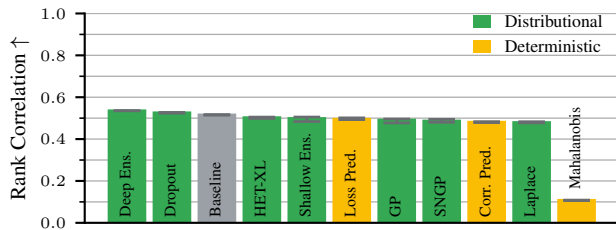


Figure 3. Only deep ensembles and dropout have a higher rank correlation with the ground-truth aleatoric uncertainty than the cross-entropy baseline (rank corr. = 0.516). Methods are evaluated on the ImageNet validation set using the entropy of the ImageNet-Real labels as ground-truth aleatoric uncertainty.

decompositions propose as epistemic estimators. So, while the IT decomposition of Laplace gives decorrelated estimators, they are not good estimators of their corresponding ground-truth uncertainties (Appendix G.2).

The best OOD detection, and thus the highest alignment with epistemic uncertainty, is achieved by the Mahalanobis method. This is a method developed specifically for OOD detection. It is also trained specifically for corruptions of severity two, and its advantage vanishes already when using OOD samples of severity three (Appendix I.3). Interestingly, the second estimator of this latent density type, DUQ, is the worst-performing method on CIFAR-10 (Appendix F.4). This may be because DUQ is developed as a predictive uncertainty estimator, not an OOD detector (Van Amersfoort et al., 2020). This shows that uncertainty estimators must be specifically tailored to the task a practitioner intends to use them for rather than relying on high-level intuitions.

3.3. Aleatoric uncertainty alone is hard to quantify

The previous experiment isolated the epistemic capabilities of uncertainty estimates. Let us now benchmark how well they predict aleatoric uncertainties.

We follow Tran et al. (2022); Kirchhof et al. (2023a;b) and use human annotators as ground-truths for the aleatoric uncertainty, in particular their disagreement: ImageNet-Real (Beyer et al., 2020) (and CIFAR-10H (Peterson et al., 2019)) queries multiple annotators for labels on each image. We showcase some examples in Appendix L. If even humans disagree about the content of an image, it reflects that the image is ambiguous by itself. This means that the entropy of the soft-label distribution per image gives an aleatoric uncertainty ground truth. We then calculate the rank correlation between an uncertainty estimator and these ground-truth values across all images.

Figure 3 shows that most methods perform below the cross-entropy baseline. Deep ensembles are most aligned with human uncertainties on average. This result shows that

Benchmarking Uncertainty Disentanglement

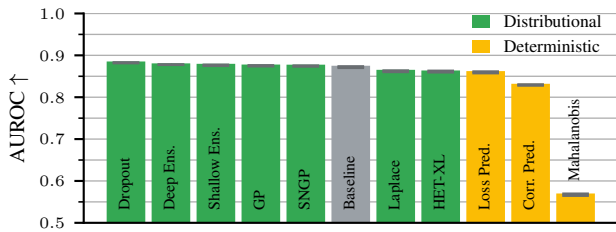


Figure 4. ID, the performance of methods on predicting correctness is saturated, as measured by the AUROC w.r.t. model correctness on the ImageNet validation set. Methods apart from Mahalanobis are within a 0.023 AUROC band. Only dropout achieves consistently better results. The Mahalanobis method is a specialized OOD detector that cannot distinguish ID samples.

even though the ensemble members are initialized to the same pretrained model, the stochasticity induced by independently training them further benefits human alignment. On CIFAR-10, the ensemble members are trained from scratch, resulting in a considerably more performant method.

On the other side of the spectrum, the Mahalanobis method is almost uncorrelated with aleatoric uncertainty. Coupled with an estimator that represents only aleatoric without epistemic uncertainties, this could give rise to disentangled uncertainty estimators in future research. This result is surprising since latent density methods are intended to capture aleatoric uncertainty by placing aleatorically uncertain samples in between class centroids where density and thus uncertainty is high (Van Amersfoort et al., 2020).

3.4. Correctness prediction works across the board

Let us now broaden the view beyond disentanglement to benchmark how well uncertainty estimators solve other practically relevant tasks. We start with correctness prediction, where the AUROC quantifies whether wrong predictions generally have higher uncertainties than correct predictions.

Figure 4 shows that most uncertainty estimators perform within ± 0.014 of the cross-entropy baseline when predicting correctness. Modern methods like HET-XL do not outperform older methods like deep ensembles or dropout. We see a similar saturation when slightly altering the correctness metric to account for soft labels in Appendix I.1.

There is a related task called abstained prediction, where the predictions for the $x\%$ most uncertain examples are excluded and we measure the accuracy on the remainder. The area under this curve shows how many errors removing the supposedly uncertain samples prevents. Figure 5 shows that the saturation is just as pronounced on the abstained prediction task. All uncertainty methods apart from Mahalanobis obtain an AUC score greater than 0.92. Practically, this means that one can obtain a high classification accuracy by

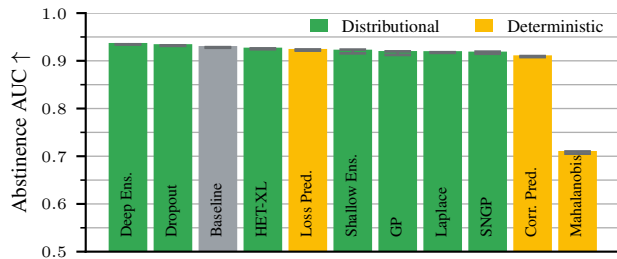


Figure 5. ID, all benchmarked methods apart from a specialized OOD detector perform very well on the abstinance task (AUC ≥ 0.9) based on the AUC of the cumulative abstinance accuracy curve. However, only dropout and deep ensemble surpass the baseline, and the methods are saturated, with almost all of them being within a 0.03 AUC band. Evaluation performed on the ImageNet validation dataset.

utilizing any of the current uncertainty quantifiers to abstain from prediction on a tiny set of uncertain samples. While computationally expensive distributional methods have a slight edge, the cheaper deterministic methods also give considerable performance.

We would like to highlight the poor performance of the Mahalanobis method. Being a specialized OOD detector, it aligns with our expectations that it cannot tell the correctness of only in-distribution samples apart. In Appendix F.2, we show that DUQ, which likewise models the data density, also falls behind the baseline on the correctness prediction task. This suggests that the predictive uncertainty DUQ is intended to achieve is less aligned with the notion of correctness benchmarked here.

3.5. Uncertainties can generalize well to OOD settings

A necessary condition for the reliable deployment of uncertainty quantification methods is that their estimates should stay performant when facing uncertain inputs. We test this by checking if their previous abstinance and correctness performances remain high longer than the model’s accuracy when increasing the OOD perturbation level. Only then can we trust them and base, e.g., the abstinance from prediction on these uncertainty estimates.

Figure 6 shows the correctness prediction AUROC, abstinance AUC, and model accuracy as we increasingly perturb the samples and go OOD. The results show an almost constant correctness prediction performance as we go more OOD, whereas the accuracy degrades considerably. Abstained prediction performance degrades together with accuracy, which is a fundamental property of the metric itself since the area under the accuracy curve depends on the baseline accuracy. This is maintained even when we normalize the metrics (solid lines) according to their random predictive performance using the formula $(\text{metric} - \text{rnd}) / (1 - \text{rnd})$ for

Benchmarking Uncertainty Disentanglement

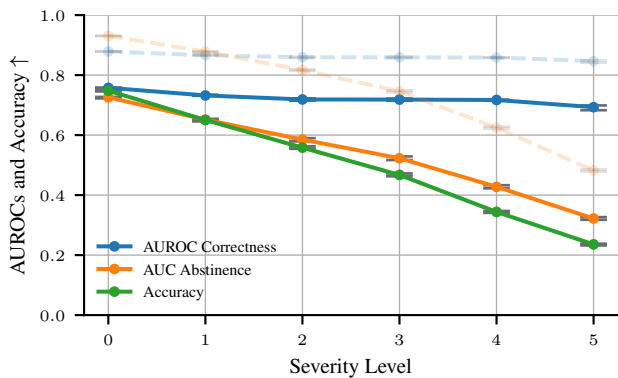


Figure 6. The predictive performance of uncertainty methods degrades much slower on the correctness prediction task than the accuracy of the model as the samples become more OOD by corrupting the ImageNet validation images. The displayed method is dropout, whose results are representative of all other methods (except Mahalanobis). Solid lines correspond to metrics normalized to the $[0, 1]$ range w.r.t. the random predictor on the particular task. Dashed lines correspond to unnormalized values.

direct comparability, where rnd is the base value that a random predictor achieves on that metric (0.5 for AUROC, $1/C$ for classification accuracy). This holds across all methods (except Mahalanobis), see Appendix I.2. This observation underlines the trustworthiness of existing uncertainty quantification methods on OOD correctness prediction.

3.6. Different tasks require different estimators

In the previous sections, we have hinted at the fact that the performance across methods is very similar on some tasks and dissimilar on others. In this section, we investigate the correlation among the previous practical tasks and further popular metrics using a correlation matrix. To construct the matrix, we consider all benchmarked methods with all uncertainty aggregators that can be benchmarked on the considered metrics (see Appendix B) and calculate the rank correlation on different pairs of metrics.

Figure 7 shows two clusters of metrics. The Brier score and log probability proper scoring rules, coupled with the ECE metric, are all recognized as *predictive uncertainty* metrics in the literature (Mucsányi et al., 2023). As such, they are also tightly connected to the correctness prediction task. The high rank correlation among these metrics (rank corr. $\in [0.707, 0.948]$) evidences this claim empirically. The second cluster is the accuracy, abstinance, aleatoric uncertainty triad: interestingly, methods that are good at capturing human uncertainty are also the most accurate, even though models are trained without access to those human uncertainty soft labels. There is no notable connection between the two clusters: Although proper scoring rules evaluate both the uncertainty estimates and the correctness

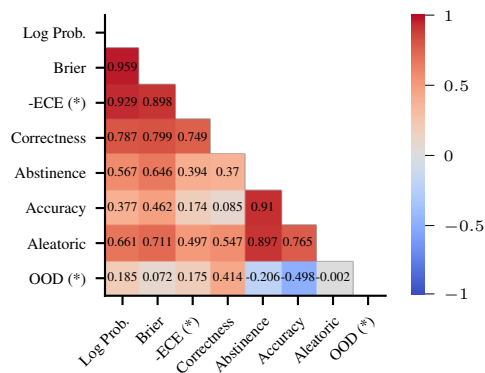


Figure 7. Only some of the considered metrics have a very high correlation among methods on the ImageNet validation dataset: most capture different aspects of uncertainty methods. Rank correlation of metric pairs across all methods and aggregators. (*) OOD’s and ECE’s correlation with the other methods is sensitive to the choice of aggregator, see Appendix J.

of prediction in one, they are not correlated with accuracy. OOD detection does not belong to any of these clusters, underlining that it benchmarks a different type of uncertainty, namely, epistemic uncertainty.

As there are different groups of tasks, there is no one-fits-all uncertainty estimator. Table 1 demonstrates this by ranking all methods on all tasks. An uncertainty estimator has to be chosen or developed for the specific task a practitioner is interested in. If the task is unknown, dropout and deep ensembles offer a good compromise, but even the baseline is a good starting point if the runtime costs of deep ensembles are too high (Appendix K.3).

3.7. Conclusions do not always transfer among datasets

We conclude our experiments with a word of caution: Appendix F repeats all above experiments on CIFAR-10, which is widely used in the uncertainty quantification literature (Van Amersfoort et al., 2020; Mukhoti et al., 2021; 2023; Gruber & Buettner, 2023) but can sometimes lead to conclusions not in line with the larger scale ImageNet. We share some key differences below.

Disentanglement. Most methods show an almost perfect rank correlation between the components of the IT decomposition (see Figure 1 for ImageNet and Figure 8 for CIFAR-10). However, *which* of the methods shows a promising decorrelation of the components is seemingly random across the datasets: on ImageNet, Laplace leads to uncorrelated components, but it showcases perfect correlation on CIFAR-10. The opposite holds for SNGP and its GP variant.

Aleatoric uncertainty. On CIFAR-10, all methods are consistently less aligned with human uncertainties (best

Benchmarking Uncertainty Disentanglement

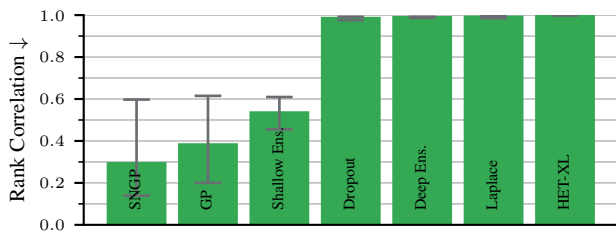


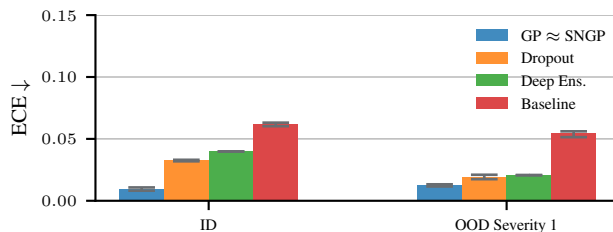
Figure 8. Four out of seven distributional methods exhibit an almost perfect rank correlation (≥ 0.986) between the IT aleatoric and epistemic components when tested on CIFAR-10. The Laplace method that disentangles the IT aleatoric and epistemic components on ImageNet shows an extreme correlation.

rank corr. 0.40 vs 0.54 on ImageNet, Appendix F.5). This is peculiar, as the vast number of classes on ImageNet could introduce more noise into the human soft labels. For example, human labelers might unanimously recognize that the object in question is a dog, yet there may be variations in opinion regarding its specific breed. This indicates that the way soft labels are obtained, which is different in ImageNet-Real and CIFAR-10H, has a significant impact on the results.

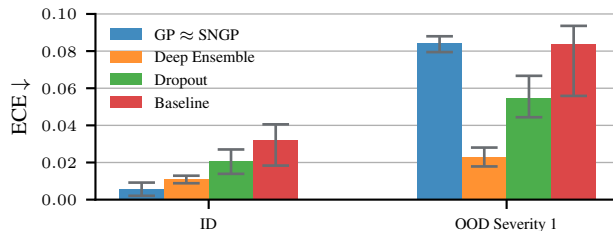
Robustness. We observe that correctness predictors are much more robust on ImageNet than on CIFAR-10, even though the drop in accuracy is very similar. Unlike on ImageNet, where the uncertainty estimators maintain a close to constant performance in predicting correctness as we go more and more OOD (Figure 6), on CIFAR-10, almost all correctness estimators deteriorate together with the model’s accuracy (Figure 10). While robustness would appear as a striking problem on CIFAR-10, it gets resolved by simply switching to a larger-scale dataset.

Calibration. SNGP shows a highly different performance depending on the dataset and task it is trained for. On CIFAR-10, Figure 9b shows that SNGP (and GP) provide the best-calibrated uncertainties ID, but already at the lowest OOD perturbation level drop to the baseline level, with the ECE jumping from 0.005 to 0.084. In Appendix I.4, we show that the best way to aggregate SNGP’s posterior $q(f)$ into an uncertainty score $u(x)$ is different when optimizing for the ECE versus correctness prediction, showing that subtle design choices greatly affect performance. On ImageNet, Figure 9a evidences that the exact opposite happens: most methods become *more calibrated* as we go slightly OOD, and SNGP variants retain their calibration.

Method rankings. These fundamentally different behaviors also change the ranking of the approaches in some tasks. Table 2 shows the correlation of rankings on CIFAR-10 and ImageNet. 6 out of 7 metrics have substantially different



(a) ImageNet calibration results. Methods preserve their rankings as the dataset becomes more OOD via ImageNet-C corruptions, and most of them even become *more calibrated*.



(b) On CIFAR-10, methods become less calibrated, do not preserve their rankings, and SNGP breaks down to the baseline level at severity level one already.

Figure 9. Methods display drastically different behavior on ImageNet and CIFAR-10 regarding the robustness of their calibration.

Metric	Rank Corr. CIFAR-10 vs ImageNet
Correctness	0.286
Abstinance	0.546
Log Prob.	0.358
Brier	0.373
Aleatoric	0.316
ECE	0.330
OOD	-0.289

Table 2. The rankings of approaches are considerably different on CIFAR-10 and ImageNet. Correlations of method rankings on different metrics for all combinations of methods and aggregators.

rankings (rank. corr. ≤ 0.373). This is also reflected in the best-performing methods: In four of the seven tasks, it is not the same on ImageNet as it is on CIFAR-10. This indicates that performance on CIFAR-10 should not be taken as an estimate for performance on larger-scale datasets.

These experiments underline that methods might show substantially different behaviors on large-scale datasets. We encourage, as best practice, to first scale the approaches to the final deployment domain (and define a precise task) instead of making fundamental design choices on toy datasets.

4. Related Works

This section discusses previous quantitative uncertainty studies and how their findings connect to ours.

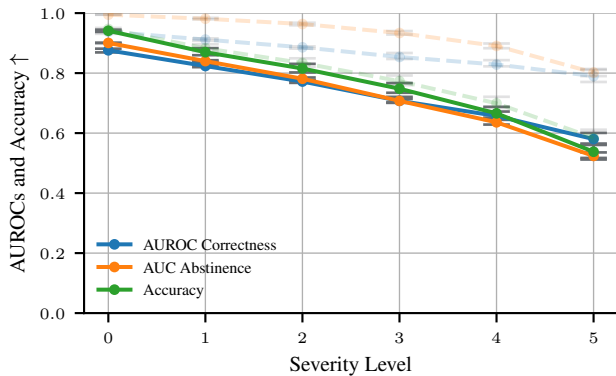


Figure 10. On CIFAR-10, model accuracy and the performance of the uncertainty method degrade together as the samples become more OOD by corrupting the images. The displayed method is dropout, whose results are representative of all other methods (except Mahalanobis). Solid lines correspond to metrics normalized to the $[0, 1]$ range w.r.t. the random predictor; dashed lines correspond to the unnormalized values.

Disentanglement of aleatoric and epistemic uncertainties via decompositions (Pfau, 2013; Depeweg et al., 2018) has recently been shown to have failure cases (Wimmer et al., 2023; Bengs et al., 2023; Gruber et al., 2023; Valdenegro-Toro & Mori, 2022). These papers make their claims theoretically or via toy problems like binary classification or 1D regression. Our results support this discussion with a practical and quantitative perspective. We find that uncertainty decompositions do not work in general and, if at all, depend greatly on their practical implementation. Our findings thus encourage us to take a more holistic view of disentanglement, including the decomposition formula, method, and implementation. For example, according to our results, it is promising to combine separate methods, such as loss prediction and the Mahalanobis distance, where each method handles a specific type of uncertainty, similar to Mukhoti et al. (2021).

Sensitivity to implementation. The finding that the performance of uncertainty estimators depends greatly on their implementation and that different design choices are better for different metrics is in line with recent benchmarks (Galil et al., 2023b; Kirchof et al., 2023b). Our benchmark further shows that the aggregator function of distributional methods is a crucial component and that simply averaging outputs is often inferior to, e.g., averaging in the dual space (Gupta et al., 2022).

Robustness. Recent benchmarks on OOD detection and robustness (Nado et al., 2021; Ovidia et al., 2019; Postels et al., 2022; Galil et al., 2023a) have first highlighted robustness issues of uncertainty estimates. Our benchmark

supports these findings on CIFAR-10, especially in the region where the OOD-ness is only slight yet already causes degradation of both the main task and the uncertainty estimator. The latter implies that uncertainty estimators either need to become more robust to distribution shifts (Kirchof et al., 2023b) or be better able to detect subtle epistemic uncertainties. However, our experiments on ImageNet do not show robustness issues, highlighting the importance of the used dataset.

Aleatoric uncertainty, as opposed to epistemic uncertainty with the OOD detection proxy task, still lacks a standardized testing protocol. The current approaches seem to converge to soft labels, but nuances in how they are collected still need discussion (compare, for example, CIFAR-10H (Peterson et al., 2019) to CIFAR-10S (Collins et al., 2022) and CIFAR-10N (Wei et al., 2022)). An increasing number of uncertainty quantification approaches compare to such human ground-truth notions of aleatoric uncertainty (Tran et al., 2022; Kirchof et al., 2023a;b), indicating the interest in the field. Our benchmark shows that no method can yet reliably give aleatoric uncertainty estimates, stressing the need for benchmarks and methods to develop along.

Predictive uncertainty tasks and calibration, on the other hand, start to become saturated and ready for application according to our experiments. This corroborates recent findings by Galil et al. (2023b). In comparison to this benchmark that compared model architectures, we compared different approaches on the same backbone.

A limitation of our study is the focus on classification. Future research should aim to expand the scope of the investigation to include a wider array of tasks like uncertainties for regression (Upadhyay et al., 2023) or unsupervised learning (Kirchof et al., 2023a). This expansion would help in generalizing our findings and understanding the dynamics of uncertainty quantification in diverse real-world scenarios.

5. Conclusion

We study how current uncertainty estimators and disentanglement formulas perform on a wide array of uncertainty quantification tasks. In summary, our findings encourage a pragmatic reassessment of uncertainty quantification research. There is no general uncertainty; instead, uncertainty quantification covers a spectrum of tasks where the definition of the exact task heavily influences the optimal method and performance. Such a precise definition of tasks per estimator would also benefit constructing disentangled uncertainties. This shift could lead to the alignment of theoretical development and intuitive descriptions about what particular types of uncertainty a method aims to capture, with tangible improvements on the benchmark tasks we consider.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

References

- Bengs, V., Hüllermeier, E., and Waegeman, W. On second-order scoring rules for epistemic uncertainty quantification. In *International Conference on Machine Learning (ICML)*, 2023.
- 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.
- Biewald, L. Experiment tracking with weights and biases. 2020. URL <https://www.wandb.com/>. Software available from wandb.com.
- Collier, M., Jenatton, R., Mustafa, B., Houlsby, N., Berent, J., and Kokiopoulou, E. Massively scaling heteroscedastic classifiers. In *The Eleventh International Conference on Learning Representations, 2023*. URL <https://openreview.net/forum?id=sIoED-yPK9l>.
- Collins, K. M., Bhatt, U., and Weller, A. Eliciting and learning with soft labels from every annotator. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing (HCOMP)*, volume 10, 2022.
- Daxberger, E., Kristiadi, A., Immer, A., Eschenhagen, R., Bauer, M., and Hennig, P. Laplace redux—effortless bayesian deep learning. *Advances in Neural Information Processing Systems*, 34:20089–20103, 2021.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Depeweg, S., Hernandez-Lobato, J.-M., Doshi-Velez, F., and Udluft, S. Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning. In *International Conference on Machine Learning*, pp. 1184–1193. PMLR, 2018.
- Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059. PMLR, 2016.
- 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 *International Conference on Learning Representations (ICLR)*, 2023b.
- Gruber, C., Schenk, P. O., Schierholz, M., Kreuter, F., and Kauermann, G. Sources of uncertainty in machine learning—a statisticians’ view. *arXiv preprint arXiv:2305.16703*, 2023.
- Gruber, S. and Buettner, F. Uncertainty estimates of predictions via a general bias-variance decomposition. In *International Conference on Artificial Intelligence and Statistics*, pp. 11331–11354. PMLR, 2023.
- Gupta, N., Smith, J., Adlam, B., and Mariet, Z. Ensembling over classifiers: a bias-variance perspective. *arXiv preprint arXiv:2206.10566*, 2022.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- Hora, S. C. Aleatory and epistemic uncertainty in probability elicitation with an example from hazardous waste management. *Reliability Engineering & System Safety*, 54(2-3):217–223, 1996.
- 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. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023b. URL <https://openreview.net/forum?id=e9n4JjkmXZ>.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. *Master’s thesis, Department of Computer Science, University of Toronto*, 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*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=eGLdVRvfvfQ>. Expert Certification.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple and scalable predictive uncertainty estimation using deep

- ensembles. *Advances in neural information processing systems*, 30, 2017.
- Lee, K., Lee, K., Lee, H., and Shin, J. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018.
- Lee, S., Purushwalkam, S., Cogswell, M., Crandall, D., and Batra, D. Why m heads are better than one: Training a diverse ensemble of deep networks. *arXiv preprint arXiv:1511.06314*, 2015.
- 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*, 33:7498–7512, 2020.
- Loshchilov, I. and Hutter, F. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Skq89Scxx>.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Mucsányi, B., Kirchhof, M., Nguyen, E., Rubinstein, A., and Oh, S. J. Trustworthy machine learning. *arXiv preprint arXiv:2310.08215*, 2023.
- Mukhoti, J., Kirsch, A., van Amersfoort, J., Torr, P. H., and Gal, Y. Deep deterministic uncertainty: A simple baseline. *arXiv preprint arXiv:2102.11582*, 2021.
- Mukhoti, J., Kirsch, A., van Amersfoort, J., Torr, P. H., and Gal, Y. Deep deterministic uncertainty: A new simple baseline. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24384–24394, 2023.
- Nado, Z., Band, N., Collier, M., Djolonga, J., Dusenberry, M., Farquhar, S., Filos, A., Havasi, M., Jenatton, R., Jerfel, G., Liu, J., Mariet, Z., Nixon, J., Padhy, S., Ren, J., Rudner, T., Wen, Y., Wenzel, F., Murphy, K., Sculley, D., Lakshminarayanan, B., Snoek, J., Gal, Y., and Tran, D. Uncertainty Baselines: Benchmarks for uncertainty & robustness in deep learning. *arXiv preprint arXiv:2106.04015*, 2021.
- 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.
- Peterson, J. C., Battleday, R. M., Griffiths, T. L., and Rusakovsky, O. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9617–9626, 2019.
- Pfau, D. A generalized bias-variance decomposition for bregman divergences. *Unpublished Manuscript*, 2013.
- 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.
- Shaker, M. H. and Hüllermeier, E. Ensemble-based uncertainty quantification: Bayesian versus credal inference. In *PROCEEDINGS 31. WORKSHOP COMPUTATIONAL INTELLIGENCE*, volume 25, pp. 63, 2021.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Tran, D., Liu, J. Z., Dusenberry, M. W., Phan, D., Collier, M., Ren, J., Han, K., Wang, Z., Mariet, Z. E., Hu, H., Band, N., Rudner, T. G. J., Nado, Z., van Amersfoort, J., Kirsch, A., Jenatton, R., Thain, N., Buchanan, E. K., Murphy, K. P., Sculley, D., Gal, Y., Ghahramani, Z., Snoek, J., and Lakshminarayanan, B. Plex: Towards reliability using pretrained large model extensions. In *First Workshop on Pre-training: Perspectives, Pitfalls, and Paths Forward at ICML 2022*, 2022. URL <https://openreview.net/forum?id=6x0gB9gOHFg>.
- Upadhyay, U., Kim, J. M., Schmidt, C., Schölkopf, B., and Akata, Z. Likelihood annealing: Fast calibrated uncertainty for regression. *arXiv preprint arXiv:2302.11012*, 2023.
- 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 Amersfoort, J., Smith, L., Teh, Y. W., and Gal, Y. Uncertainty estimation using a single deep deterministic neural network. In *International conference on machine learning*, pp. 9690–9700. PMLR, 2020.
- Von Luxburg, U. and Schölkopf, B. Statistical learning theory: Models, concepts, and results. In *Handbook of the History of Logic*, volume 10, pp. 651–706. Elsevier, 2011.
- Wei, J., Zhu, Z., Cheng, H., Liu, T., Niu, G., and Liu, Y. Learning with noisy labels revisited: A study using real-world human annotations. In *International Conference on Learning Representations (ICLR)*, 2022.

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*, pp. 2282–2292. PMLR, 2023.

Yoo, D. and Kweon, I. S. Learning loss for active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

Zagoruyko, S. and Komodakis, N. Wide residual networks. In *Proceedings of the British Machine Vision Conference 2016*. British Machine Vision Association, 2016.

A. Benchmarked Methods

We consider a classification setting with discrete label space $\{1, \dots, C\}$ of C classes and models that output a probability vector $f(x) \in \Delta^{C-1}$ for input $x \in \mathcal{X}$. The (pre-softmax) logits of the models are denoted by $\log \hat{f}(x)$.

We evaluate two classes of methods: direct prediction methods and distributional methods.

A.1. Direct Prediction Methods

Direct prediction methods output an uncertainty estimate $u(x)$ for input x , such as $u(x) \approx p(f(x) \text{ is a correct prediction})$.

A.1.1. RISK PREDICTION

Risk prediction (Upadhyay et al., 2023; Lahlou et al., 2023; Kirchhof et al., 2023b) employs an additional output head u^{RP} connected to the pre-logit layer that predicts the risk of the network’s prediction on input $x \in \mathcal{X}$. The risk predictor head is trained in a supervised fashion by making $u^{\text{RP}}(x)$, the predicted risk, closer to the actual loss $\ell(f(x), y) = -\log \hat{f}_y(x)$. Precisely, we use the objective

$$\mathcal{L} = \sum_{i=1}^n -\log \hat{f}_{y_i}(x_i) + \lambda \left(u^{\text{RP}}(x_i) + \log \hat{f}_{y_i}(x_i) \right)^2, \quad (3)$$

where the risk predictor loss (squared Euclidean distance) is traded off with the label predictor loss (cross-entropy) with a hyperparameter $\lambda \in \mathbb{R}^+$.

Note that y is a random variable in the presence of aleatoric uncertainty. In expectation, this encourages $u^{\text{RP}}(x)$ to approximate the true *risk* $\mathcal{R}(f, x) = \mathbb{E}_{p^{\text{st}}(y|x)} [\ell(f(x), y)]$ at each input x .

A.1.2. CORRECTNESS PREDICTION

Correctness prediction is a variant of risk prediction that, instead of aiming to predict the risk of the network on input x , predicts the probability of correctness $p \left(\arg \max_{c \in \{1, \dots, C\}} f_c(x) = y \mid x \right)$ on input x . This is achieved by using a sigmoid correctness predictor head h and using the objective

$$\mathcal{L} = \sum_{i=1}^n -\log \hat{f}_{y_i}(x_i) - \lambda (l_i \log h(x_i) + (1 - l_i) \log(1 - h(x_i))), \quad (4)$$

where $l_i = \mathbf{I} \left[\arg \max_{c \in \{1, \dots, C\}} \log \hat{f}_c(x_i) = y_i \right] \forall i \in \{1, \dots, n\}$, and the correctness predictor loss (binary cross-entropy) is traded off with the label predictor loss (cross-entropy) with a hyperparameter $\lambda \in \mathbb{R}^+$. The uncertainty estimate is $u^{\text{CP}}(x) = 1 - h(x)$ (i.e., the probability of making an error).

A.1.3. DETERMINISTIC UNCERTAINTY QUANTIFICATION

The deterministic uncertainty quantification (DUQ) method of Van Amersfoort et al. (2020) learns a latent mixture-of-RBF density on the training set with a strictly proper scoring rule to capture the uncertainty in the prediction based on the Euclidean distance of the input’s embedding to the mixture means. The training objective is

$$\mathcal{L} = - \sum_{i=1}^n \sum_{c=1}^C y_{ic} \log K_c(x_i) + (1 - y_{ic}) \log(1 - K_c(x_i)), \quad (5)$$

where $K_c(x) = \exp \left(-\frac{1}{2\gamma} \left\| \log \hat{f}(x) - m_c \right\|^2 \right)$ is the RBF value corresponding to class $c \in \{1, \dots, C\}$ identified by its mean vector m_c in the latent space. To facilitate minibatch training, Van Amersfoort et al. (2020) employ an exponential

moving average (EMA) to learn the mean vector using the following update rules:

$$n_c \leftarrow \gamma \cdot n_c + (1 - \gamma)|\mathcal{B}_c| \quad (6)$$

$$M_c \leftarrow \gamma \cdot M_c + (1 - \gamma) \sum_{(x,y) \in \mathcal{B}_c} W_c \log \hat{f}(x) \quad (7)$$

$$m_c \leftarrow \frac{M_c}{n_c}, \quad (8)$$

where \mathcal{B} is a minibatch of samples and $\mathcal{B}_c = \{(x, y) \in \mathcal{B} \mid y = c\} \forall c \in \{1, \dots, C\}$. γ is the EMA parameter and W_c characterizes a linear mapping of the logits for each class.

To regularize the latent density and prevent feature collapse, Van Amersfoort et al. (2020) use the following gradient penalty added to $\nabla_{\theta} \mathcal{L}$:

$$\lambda \cdot \left(\left\| \nabla_x \sum_{c=1}^C K_c \right\|_2^2 - 1 \right)^2 \quad (9)$$

Each RBF component in the latent space corresponds to one class. The confidence output of the method is the maximal RBF value of the input over all classes. Therefore, the *uncertainty* estimate can be calculated as $u^{\text{dug}}(x) = 1 - \max_{c \in \{1, \dots, C\}} K_c(x)$.

The predicted class of the trained network is $\arg \max_{c \in \{1, \dots, C\}} K_c(x)$.

A.1.4. MAHALANOBIS

The Mahalanobis method (Lee et al., 2018) builds a post-hoc latent density for the training set in the latent space by calculating per-class means and covariances, and using the induced mixture-of-Gaussians as the latent density estimate. Such latent densities are estimated in multiple layers of the network. One layer’s confidence estimate is the maximal Mahalanobis score (Gaussian log-likelihood) $K_{\ell}(x)$ over all classes:

$$K_{\ell,c}(x) = -(f_{\ell}(x) - \mu_{\ell,c})^{\top} \Sigma_{\ell}^{-1} (f_{\ell}(x) - \mu_{\ell,c}) \quad (10)$$

$$K_{\ell}(x) = \max_{c \in \{1, \dots, C\}} K_{\ell,c}(x), \quad (11)$$

where f_{ℓ} is the ℓ -th layer’s output,

$$\mu_{\ell,c} = \frac{1}{n_c} \sum_{i=1}^n \mathbf{I}[y_i = c] f_{\ell}(x_i) \quad (12)$$

is the centroid of the Gaussian for class $c \in \{1, \dots, C\}$ in layer $\ell \in \{1, \dots, L\}$, n_c is the number of samples with label c , and

$$\Sigma_{\ell} = \frac{1}{n} \sum_{c=1}^C \sum_{i=1}^n \mathbf{I}[y_i = c] (f_{\ell}(x) - \mu_{\ell,c}) (f_{\ell}(x) - \mu_{\ell,c})^{\top} \quad (13)$$

is the tied covariance matrix used for all classes in layer $\ell \in \{1, \dots, L\}$.

To make the differences of latent embeddings of ID and OOD samples more pronounced, all samples are adversarially perturbed w.r.t. the maximal Mahalanobis score for each layer’s confidence score:

$$\hat{x}^{(\ell)} = x - \epsilon \operatorname{sgn}(\nabla_x - K_{\ell}(x)). \quad (14)$$

This perturbed sample is used to compute $K_{\ell}(\hat{x}^{(\ell)})$. Finally, a logistic regression OOD detector is learned on a held-out validation set of a balanced mix of ID and OOD samples to learn weights w_{ℓ} for each layer $\ell \in \{1, \dots, L\}$ using the L -dimensional inputs $[K_1(\hat{x}^{(1)}), \dots, K_L(\hat{x}^{(L)})]^{\top}$. The final *uncertainty* estimate becomes $u^{\text{Mah}}(x) = \sum_{\ell=1}^L w_{\ell} K_{\ell}(x)$.

This is the only method in our benchmark that requires a validation set for training the logistic regression OOD detector and not just the hyperparameters.

A.2. Distributional Methods

Distributional methods output a conditional probability distribution *over probability vectors* $q(f(x) | x)$, abbreviated as $q(f)$.

A.2.1. SPECTRAL NORMALIZED GAUSSIAN PROCESS

Spectral normalized Gaussian processes (SNGP) (Liu et al., 2020) give an approximate Bayesian treatment to obtain uncertainty estimates using spectral normalization of the parameter tensors and a last-layer Gaussian process approximated by Fourier features. For an input x , it predicts a multivariate Gaussian distribution

$$\mathcal{N}\left(\beta\phi(x), \phi(x)^\top (\Phi^\top \Phi + I)^{-1} \phi(x)I\right), \quad (15)$$

where β is a learned parameter matrix that maps from the pre-logits to the logits, and $\phi(x) = \cos(Wh(x) + b)$ is a random feature embedding of the input x with $h(x)$ being a pre-logit embedding, W a fixed semi-orthogonal random matrix, and b a fixed random vector sampled from $\text{Uniform}(0, 2\pi)$. $\Phi^\top \Phi$ is the empirical covariance matrix of the pre-logits of the training set. This is calculated during the last epoch. The multivariate Gaussian presented above can be Monte-Carlo sampled to obtain M logit vectors.

The method also applies spectral normalization to the hidden weights in each layer in order to satisfy input distance awareness. We treat whether to apply spectral normalization through the network and whether to use layer normalization in the last layer as hyperparameters.

We benchmark both SNGPs and their non-spectral-normalized variants (simply denoted by GP).

A.2.2. LATENT HETEROSCEDASTIC CLASSIFIER

Latent heteroscedastic classifiers (HET-XL) (Collier et al., 2023) construct a heteroscedastic Gaussian distribution in the pre-logit layer to model per-input uncertainties: $\mathcal{N}(\phi(x), \Sigma(x))$, where

$$\Sigma(x) = V(x)^\top V(x) + \text{diag}(d(x)) \quad (16)$$

is an input-conditional full-rank covariance matrix. Both the low-rank term's $V(x)$ and the diagonal term's $d(x)$ are calculated as a linear function of the layer's output before the pre-logit layer.

One can Monte-Carlo sample the pre-logits from the above Gaussian distribution and obtain a set of logits by transforming each using the last linear layer of the network. During training, this set is used to calculate the Bayesian Model Average (BMA) whose argmax is the final prediction.

HET-XL uses a temperature parameter to scale the logits before calculating the BMA. This is chosen using a validation set.

A.2.3. LAPLACE APPROXIMATION

The Laplace approximation (Daxberger et al., 2021) approximates a Gaussian posterior $q(\theta | \mathcal{D})$ over the network parameters for a Gaussian prior $p(\theta)$ and likelihood defined by the network architecture. It uses the maximum a posteriori (MAP) estimate as the mean and the inverse Hessian of the loss evaluated at the MAP as the covariance matrix:

$$\mathcal{N}\left(\theta_{\text{MAP}}, \left(\frac{\partial^2 \mathcal{L}(\mathcal{D}; \theta)}{\partial \theta_i \partial \theta_j} \Big|_{\theta_{\text{MAP}}}\right)^{-1}\right). \quad (17)$$

This is a post-hoc method applied to a point estimate network. Following the recommendation of Daxberger et al. (2021), we employ a last-layer KFAC Laplace approximation and find the prior variance using evidence optimization.

A.2.4. DROPOUT

Dropout (Srivastava et al., 2014) has been shown to be a variational approximation to a deep Gaussian process (Gal & Ghahramani, 2016). Dropout in the realm of uncertainty quantification remains active during inference, and is used to sample M logits by performing M forward passes. Therefore, it directly samples from $q(f)$ without characterizing it.

A.2.5. DEEP ENSEMBLE

Deep ensembles (Lakshminarayanan et al., 2017) are approximate model distributions that give rise to a mixture of Dirac deltas in parameter space: $q(\theta) = \frac{1}{M} \sum_{i=1}^M \delta(\theta - \theta^{(i)})$. Predominantly used to reduce the variance in the predictions and improve model accuracy, deep ensembles can also be used as approximators to the true distribution $p(\theta)$ induced by the randomness over datasets $\mathcal{D} := \{(x_i, y_i) | i \in \{1, \dots, n\}, x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ in the generative process $p(x, y)$.

We obtain a set of logits by performing a forward pass over all models. Similarly to dropout, deep ensembles do not explicitly parameterize the distribution over the predictions, they only sample from it. We ensemble five models trained with cross-entropy.

A.2.6. SHALLOW ENSEMBLE

Shallow ensembles (Lee et al., 2015) are lightweight approximations of deep ensembles. They use a shared backbone and M output heads (often referred to as “experts”). With a single forward pass, one obtains M logit vectors per input.

A.2.7. DETERMINISTIC BASELINE

As a baseline, we also benchmark a deterministic network trained with the cross-entropy loss that corresponds to a Dirac delta in parameter space: $q(\theta') = \delta(\theta - \theta')$.

B. Aggregators

In practical applications, distributional methods output a discrete set of probability vectors $\{f^{(m)}(x)\}_{m=1}^M$ per input x . This set can be aggregated in several ways to construct an uncertainty estimate $u(x)$.

The Bayesian Model Average (BMA) is given by $\tilde{f}(x) = \frac{1}{M} \sum_{m=1}^M f^{(m)}(x)$. We consider two aggregators for the BMA:

- Its entropy: $u(x) = \mathbb{H}(\tilde{f}(x))$.
- One minus its maximum probability entry: $u(x) = 1 - \max_{c \in \{1, \dots, C\}} \tilde{f}_c(x)$.

The $\tilde{f}(x)$ value in the Bregman decomposition (Section 2.2.2) averages in the log probability space instead of the probability simplex: $\tilde{f}(x) = \text{softmax}\left(\frac{1}{M} \sum_{m=1}^M \log f^{(m)}(x)\right)$. Similarly to the BMA, we can take the entropy of $\tilde{f}(x)$ or one minus its maximum probability entry as an uncertainty estimate.

One can also calculate the entropy of the individual probability vectors and then average them, leading to

$$u(x) = \frac{1}{M} \sum_{m=1}^M \mathbb{H}(f^{(m)}(x)). \quad (18)$$

Similarly, one can determine the expected maximum probabilities and construct the estimator as

$$u(x) = 1 - \frac{1}{M} \sum_{m=1}^M \max_{c \in \{1, \dots, C\}} f_c^{(m)}(x). \quad (19)$$

Finally, one can directly use the epistemic components of the Bregman and IT decompositions as they do not require a ground truth. In particular, one can use

$$u(x) = \mathbb{H}(\tilde{f}(x)) - \frac{1}{M} \sum_{m=1}^M \mathbb{H}(f^{(m)}(x)), \quad f^{(m)} \sim q(f) \forall m \in \{1, \dots, M\}, \quad (20)$$

the (discretized) epistemic part of the IT decomposition (see Appendix C), or

$$u(x) = \frac{1}{M} \sum_{m=1}^M \left[D_{\text{KL}}(\tilde{f}(x) \parallel f^{(m)}(x)) \right], \quad f^{(m)} \sim q(f) \forall m \in \{1, \dots, M\}, \quad (21)$$

the (discretized) epistemic part of the Bregman decomposition (see Appendix D).

Unless stated otherwise, we use the best-performing alternative for each distributional method in the benchmarks. For these methods, the model’s prediction is always the most confident class of the BMA. For direct prediction methods, we use their “canonical” uncertainty estimator introduced in Appendix A.1.

C. Special Form on the Information-Theoretical Decomposition for Discrete Posteriors

Below, we show that the information-theoretical (IT) decomposition (Depeweg et al., 2018) separates the entropy of the BMA into an expected entropy term and a Jensen-Shannon divergence term when considering discrete uniform distributions $q(f(x) | x) = \frac{1}{M} \sum_{m=1}^M \delta(f(x) - f^{(m)}(x))$ abbreviated as $q(f)$ in the main paper.

The IT decomposition treats the entropy of the predictive distribution $p(y | x) = \int p(y | f(x)) dq(f(x) | x)$ as the predictive uncertainty metric and decomposes it into

$$\underbrace{\mathbb{H}_{p(y|x)}(y)}_{\text{predictive}} = \underbrace{\mathbb{E}_{q(f(x)|x)} [\mathbb{H}_{p(y|x,f)}(y)]}_{\text{aleatoric}} + \underbrace{\mathbb{I}_{p(y,f(x)|x)}(y; f(x))}_{\text{epistemic}}, \quad (22)$$

where \mathbb{H} is the entropy and \mathbb{I} is the mutual information.

Under a discrete uniform approximate posterior $q(f)$, the predictive uncertainty is still the entropy of the BMA and the aleatoric uncertainty also stays the expected entropy of the probability vectors of non-zero measure. We only have to show that the mutual information takes the convenient form of the Jensen-Shannon divergence under such an approximate posterior. Using $p(y, f(x) | x) = p(y | f(x))q(f(x) | x)$, we have

$$\mathbb{I}_{p(y,f(x)|x)}(y; f(x)) = \sum_{y=1}^C \int \log \frac{p(y, f(x) | x)}{p(y | x)q(f(x) | x)} dp(y, f(x) | x) \quad (23)$$

$$= \frac{1}{M} \sum_{m=1}^M \sum_{y=1}^C p(y | f^{(m)}(x)) \log \frac{p(y | f^{(m)}(x))}{p(y | x)} \quad (24)$$

$$= -\frac{1}{M} \sum_{m=1}^M \mathbb{H}(f^{(m)}(x)) - \sum_{y=1}^C \frac{1}{M} \sum_{m=1}^M p(y | f^{(m)}(x)) \log p(y | x) \quad (25)$$

$$= \mathbb{H}\left(\frac{1}{M} \sum_{m=1}^M f^{(m)}(x)\right) - \frac{1}{M} \sum_{m=1}^M \mathbb{H}(f^{(m)}(x)) \quad (26)$$

which is the Jensen-shannon divergence of the distributions $p(y | f^{(m)}(x))$, $m \in \{1, \dots, M\}$.

D. Special Form of the Bregman Decomposition for the Kullback-Leibler Divergence

Considering a Bregman divergence induced by the strictly convex function F as the loss function \mathcal{L} , the Bregman decomposition disentangles it into

$$\underbrace{\mathbb{E}_{q(f), p^{\text{st}}(y|x)} [D_F[y \| f(x)]]}_{\text{predictive}} = \underbrace{\mathbb{E}_{p^{\text{st}}(y|x)} [D_F[y \| f^*(x)]]}_{\text{aleatoric}} + \underbrace{\mathbb{E}_{q(f)} [D_F[\bar{f}(x) \| f(x)]]}_{\text{epistemic}} + \underbrace{D_F[f^*(x) \| \bar{f}(x)]}_{\text{bias}}, \quad (27)$$

where $f^*(x) = \mathbb{E}_{p^{\text{st}}(y|x)}[y]$ is the Bayes predictor and $\bar{f}(x) = \arg \min_{z \in \Delta^{C-1}} \mathbb{E}_{q(f)} [D_F[z \| f(x)]]$ is the central predictor.

When choosing $F(\cdot) = -\mathbb{H}(\cdot)$, we obtain $D_F[\cdot \| \cdot] = D_{\text{KL}}(\cdot \| \cdot)$. Consider the predictive uncertainty term. Note that we now treat samples from $p^{\text{st}}(y | x)$ as one-hot vectors. Therefore, a sample’s entropy is zero. Consequently, the predictive uncertainty becomes $\mathbb{E}_{q(f), p^{\text{st}}(y|x)} [\text{CE}(y, f(x))]$. The aleatoric term takes a convenient form:

$$\mathbb{E}_{p^{\text{st}}(y|x)} [D_{\text{KL}}(y \| f^*(x))] = \mathbb{E}_{p^{\text{st}}(y|x)} \left[\sum_{i=1}^C y_i \log \frac{y_i}{f_i^*(x)} \right] = - \sum_{i=1}^C f_i^*(x) \log f_i^*(x) = \mathbb{H}(f_i^*(x)). \quad (28)$$

On datasets with multiple labels per input, this quantity is precisely the entropy of the (normalized) label distribution corresponding to the labeler votes.

To calculate $\bar{f}(x)$, we can proceed as follows.

$$\bar{f}(x) = \arg \min_{z \in \Delta^{C-1}} \mathbb{E}_{q(f)} [D_{\text{KL}}(z \parallel f(x))] \quad (29)$$

$$= \arg \min_{z \in \Delta^{C-1}} \sum_{i=1}^C z_i \log z_i - \sum_{i=1}^C z_i \log (\exp (\mathbb{E}_{q(f)} [\log f_i(x)])) \quad (30)$$

$$= \arg \min_{z \in \Delta^{C-1}} \sum_{i=1}^C z_i \log z_i - \sum_{i=1}^C z_i \log (\exp (\mathbb{E}_{q(f)} [\log f_i(x)])) + \sum_{i=1}^C z_i \log \left(\sum_{j=1}^C \exp (\mathbb{E}_{q(f)} [\log f_j(x)]) \right) \quad (31)$$

$$= \arg \min_{z \in \Delta^{C-1}} \sum_{i=1}^C z_i \log z_i - \sum_{i=1}^C z_i \log \underbrace{\frac{\exp (\mathbb{E}_{q(f)} [\log f_i(x)])}{\sum_{j=1}^C \exp (\mathbb{E}_{q(f)} [\log f_j(x)])}}_{p_i:=} \quad (32)$$

$$= \arg \min_{z \in \Delta^{C-1}} D_{\text{KL}}(z \parallel p) \quad (33)$$

$$= p. \quad (34)$$

Therefore, $\bar{f}(x) = \text{softmax} (\mathbb{E}_{q(f)} [\log f(x)])$, where \log is applied elementwise.

D.1. DEUP is a Special Case of Bregman

As mentioned in Section 2.2.2 of the main paper, a closely related formula to Bregman is the risk decomposition of [Lahlou et al. \(2023\)](#) where the predictive uncertainty is directly equated to the risk of a deterministic predictor $f: \mathcal{X} \rightarrow \mathcal{Y}$, not an expectation of risks over datasets or hypothesis distributions:

$$\underbrace{\mathcal{R}(f, x)}_{\text{predictive}} = \underbrace{\mathcal{R}(f^*, x)}_{\text{aleatoric}} + \underbrace{\mathcal{R}(f, x) - \mathcal{R}(f^*, x)}_{\text{bias}} \quad (35)$$

where $\mathcal{R}(f, x) = \mathbb{E}_{p(y|x)} [\mathcal{L}(f(x), y)]$ is the pointwise risk of f on $x \in \mathcal{X}$. When choosing \mathcal{L} to be the squared Euclidean distance or the Kullback-Leibler divergence, Equation 35 becomes a special case of Equation 2 for a Dirac distribution $q(f') = \delta(f' - f)$ at an arbitrary predictor f . This formulation is desired when one wants the predictive uncertainty to be aligned with the risk of one particular predictor and not the expected risk over a hypothesis distribution.

E. Goals of Disentanglement

What does it mean to have disentangled uncertainty estimators? Consider two estimators $u^{(a)}(x_i), u^{(e)}(x_i)$ and *ground-truth* aleatoric and epistemic uncertainties $U^{(a)}(x_i), U^{(e)}(x_i)$ for each input x_i . The estimators $u^{(a)}$ and $u^{(e)}$ are disentangled if

1. $u^{(a)}$ has low rank correlation with $U^{(e)}$ and
2. $u^{(e)}$ has low rank correlation with $U^{(a)}$.

Importantly, $u^{(a)}$ and $u^{(e)}$ having a severely high rank correlation prohibits disentanglement. Further, they are well-performing if

3. $u^{(a)}$ has high rank correlation with $U^{(a)}$ and
4. $u^{(e)}$ has high rank correlation with $U^{(e)}$.

Inspired by generalized bias-variance decompositions ([Pfau, 2013](#); [Gruber & Buettner, 2023](#)), one may treat the training dataset \mathcal{D} as a random variable sampled from the generative process $p(x, y)$ and record the variability of the trained predictor under dataset change. Following the Bregman decomposition, one may then define

$$U^{(e)}(x) := \mathbb{E}_{p(\mathcal{D})} [D_F [\bar{f}(x) \parallel f_{\mathcal{D}}(x)]] \quad (36)$$

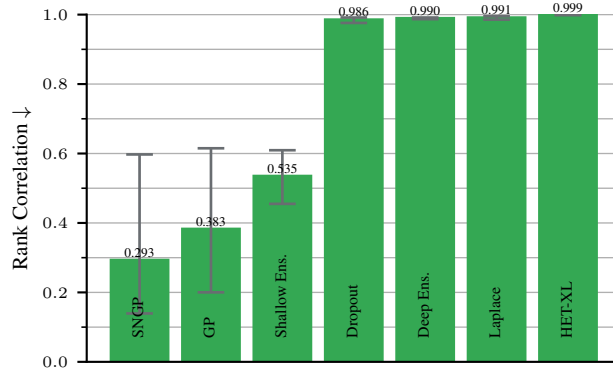


Figure 11. Four out of seven distributional methods exhibit an almost perfect rank correlation between the IT aleatoric and epistemic components on the CIFAR-10 dataset. These methods violate a necessary condition of uncertainty disentanglement.

with the corresponding central predictor $\bar{f}(x) = \arg \min_z \mathbb{E}_{p(\mathcal{D})} [D_F [z \| f_{\mathcal{D}}(x)]]$. As this is impossible to obtain in practical setups (or is too noisy to Monte Carlo estimate), we instead consider the proxy task of OOD detection for evaluating the disentanglement of aleatoric and epistemic uncertainties. This is not needed for, e.g., the evaluation of the Bregman bias and aleatoric components’ disentanglement.

F. CIFAR-10 Experiments

This section mirrors Section 3 of the main paper on the CIFAR-10 dataset. We want to understand the behavior of current uncertainty quantification methods. We first benchmark disentanglement formulas and then individual estimators.

F.1. Disentanglement formulas often fail

Figure 11 reveals a surprisingly simple failure of disentanglement formulas: In four of the seven distributional methods, the IT decomposition leads to aleatoric and epistemic estimates that are highly mutually correlated (Rank. Corr ≥ 0.99). This violates the first two necessary conditions for disentanglement. The correlation remains even when we add more epistemic uncertainty into the dataset in Appendix G.1 or switch to Bregman estimates in Appendix H. Three posterior estimators, SNGP, GP, and shallow ensembles (but not deep ensembles) show a lower rank correlation between their aleatoric and epistemic components. From this, it follows that disentanglement cannot be thought about only on the level of decomposition formulas but needs to take the explicit implementation into account.

To find out whether the remaining three estimator pairs are not only disentangled but also well-performing, we compare them to the above-defined aleatoric and epistemic ground truths in Appendix G.2. We find that they perform better than random but worse than specialized deterministic estimators for both tasks. In other words, combining two specialized estimators may practically perform better than deriving them by decomposing a single posterior.

F.2. Correctness prediction and abstained prediction work across the board

Figure 12 shows that most uncertainty estimators perform within ± 0.015 of the baseline when predicting correctness ID, and modern methods like HET-XL do not outperform older methods like deep ensembles. We see a similar saturation when slightly altering the correctness metric to account for soft labels in Appendix I.1.

The saturation is even more pronounced on the abstained prediction task. All uncertainty methods apart from Mahalanobis obtain an AUC score greater than 0.99 in Figure 13. Practically, this means that one can obtain a close-to-perfect classification accuracy by abstaining from prediction on a tiny set of samples. In both tasks, the more computationally expensive distributional methods have a slight edge over deterministic methods.

We would like to highlight the poor predictive performance of the Mahalanobis method. Being a specialized OOD detector, it aligns with our expectations that it cannot tell the correctness of only in-distribution samples apart. DUQ, which also models the data density, also falls behind the baseline on the correctness prediction task. This suggests that the notion of

Benchmarking Uncertainty Disentanglement

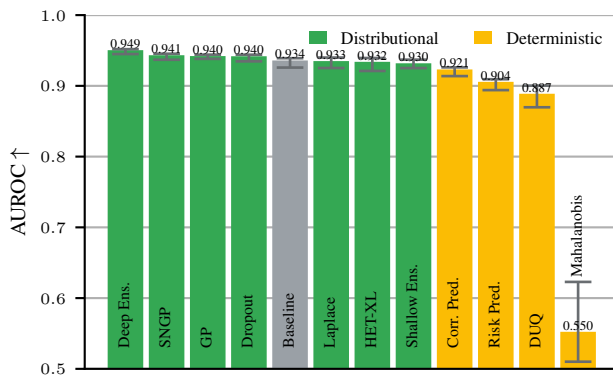


Figure 12. On ID CIFAR-10 samples, the performance of methods on predicting correctness is saturated, with only deep ensembles achieving consistently better results. The Mahalanobis method is a specialized OOD detector that cannot distinguish ID samples.

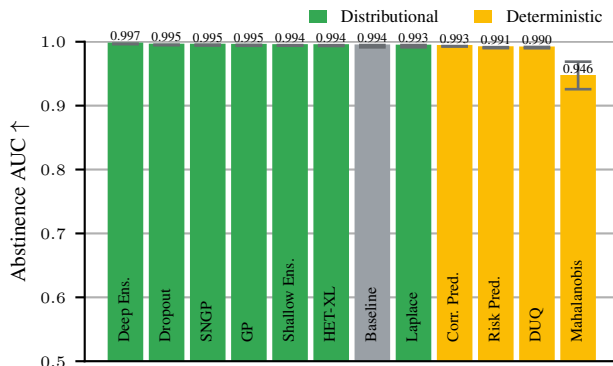


Figure 13. On ID CIFAR-10 samples, most methods solve the abstinence task almost perfectly.

predictive uncertainty imposed by the latent density of DUQ is less aligned with the notion of correctness.

F.3. Uncertainties are (only) as robust as accuracy

A necessary condition for the reliable deployment of uncertainty quantification methods is that their estimates should stay performant longer than the model’s predictive performance when going OOD. Only then can we trust them and base, e.g., the abstinence of the model on their predictions. At first glance, when plotting the raw performance values as dashed lines in Figure 10, this is the case. However, AUROCs have a random performance value of 0.5 and the random predictor’s expected accuracy is the inverse of the number of classes. Normalizing these values out by FORMULA reveals that both correctness prediction and abstained prediction degrade consistently with the model’s accuracy as the data becomes more and more OOD (solid lines). Results per method are reported in Appendix I.2. This observation shows that all benchmarked methods are incapable of generalizing better than the models, making their correctness predictions not trustworthy at higher perturbation levels.

F.4. Subtle OOD-ness is hard to detect

If both accuracy and uncertainty estimates deteriorate OOD, can current methods reliably detect OOD samples?

We use balanced mixtures of ID and OOD datasets. OOD samples are perturbed ID samples with severity level two where the models’ accuracy already severely deteriorates according to Figure 10. The uncertainty estimators are tasked to predict which sample is OOD, i.e., OOD inputs should have higher uncertainty estimates. As shown in Figure 14, the Mahalanobis method, which in the previous tasks was far off, is by far the most performant on average in telling apart clean

Benchmarking Uncertainty Disentanglement

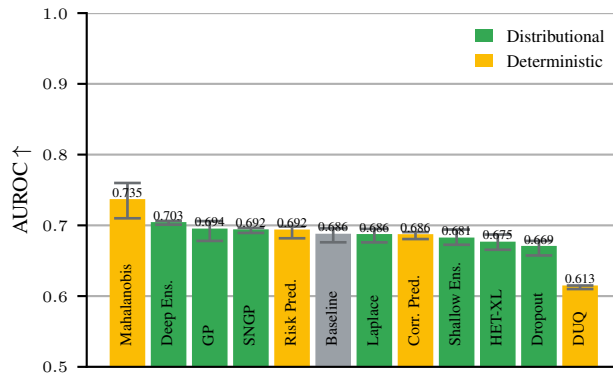


Figure 14. Only deep ensembles and Mahalanobis, a direct OOD detector, can distinguish ID and OOD samples considerably better than the baseline on CIFAR-10. OOD samples are perturbed by CIFAR-10C corruptions of severity level two.

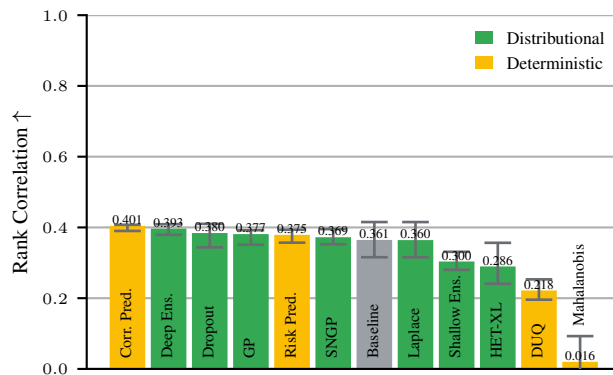


Figure 15. On CIFAR-10, none of the methods have a considerably higher rank correlation with the ground-truth aleatoric uncertainty than the baseline.

CIFAR-10 samples from perturbed ones. This shows that the OOD task indeed benchmarks a notably different uncertainty quantification capability than the previous predictive uncertainty tasks, namely epistemic uncertainty. We further investigate the correlations between the tasks in Appendix F.6.

One may expect that latent-density-based methods are naturally good at predicting OOD-ness since, for a well-behaved latent density, the embeddings of OOD inputs will lie far away from the class centroids. Interestingly, the worst-performing method also builds a latent density estimate: we cannot establish clear categories of methods that work and that do not. We hypothesize that the main power of the Mahalanobis method lies in its adversarial perturbation of the inputs and that adding this to DUQ would make it as good as Mahalanobis. We leave this experiment to future work.

The uncertainty methods show a steady increase in OOD detection performance as we increase the severity of the perturbed half of the dataset (see Appendix I.3). The advantage of the Mahalanobis method vanishes as the severity increases. At severity 5, most methods become saturated.

One may expect that the risk/correctness prediction methods' OOD detection performance is worse than distributional methods as they are trained only on ID samples and cannot utilize the disagreement among models to modulate their prediction vectors. Surprisingly, this is not the case. To resolve this paradox, we note that it does not matter what these methods' exact behavior on OOD inputs is. As long as it is sufficiently different from their ID behavior, they can be employed to detect OOD-ness.

Benchmarking Uncertainty Disentanglement

Method ranking	Correctness	Abstinance	Log Prob.	Brier	ECE	Aleatoric	OOD
Deep Ensemble	1	1	1	3	1	2	2
Dropout	4	2	2	4	2	3	11
Baseline	5	7	7	7	9	7	6
SNGP	2	3	4	1	5	6	4
GP	3	4	8	2	7	4	3
Mahalanobis	12	12	-	-	-	12	1
Shallow Ensemble	8	5	6	8	6	9	9
Laplace	6	8	3	5	8	8	7
HET-XL	7	6	5	9	4	10	10
Correctness Pred.	9	9	10	10	10	1	8
Loss Prediction	10	10	-	-	-	5	5
DUQ	11	11	9	3	6	11	12

Table 3. Different tasks have different best-performing methods on CIFAR-10. Deep ensemble is a good (but expensive, see Appendix K.3) choice across the board. Best, second-best, and third-best method highlighted in gold, silver, and bronze. Note that differences between ranks can be very small, see the plots per task for details.

F.5. Aleatoric uncertainty alone is hard to quantify

The previous experiment isolated the epistemic capabilities of uncertainty estimates. Let us now benchmark how much they predict aleatoric uncertainties. Since we use the entropy of human annotator label distributions as ground truths, this could also be considered the alignment with human uncertainties. Figure 15 shows that most methods perform within the error bar of the baseline. Correctness prediction is most aligned *on average* with a notably small min-max error bar. This is reasonable since ID, the aleatoric uncertainty of the sample determines the network’s correctness the most.¹ The latent density methods DUQ and Mahalanobis are particular outliers. Even though they are intended to capture aleatoric uncertainty by placing aleatorically uncertain samples in between class centroids, resulting in a low density (Van Amersfoort et al., 2020), they perform worse than the remaining methods. Mahalanobis even performs randomly. This reinforces that it is an epistemic, not an aleatoric uncertainty estimator. We obtain similar results for the rank correlation with the Bregman bias component: no method can significantly surpass the baseline (see Appendix H.2). We also investigate the correlation of methods with different sources of uncertainty on OOD datasets in Appendix H.2.

F.6. Different tasks require different estimators

In the previous sections, we have hinted at the fact that the performance across methods is very similar on some tasks and dissimilar on others. In this section, we investigate the correlation among the previous practical tasks using a correlation matrix displayed in Figure 16. To construct the matrix, we consider all benchmarked methods with all uncertainty aggregators (see Appendix B) and calculate the correlation of their rankings on different metrics. We find that methods good in predicting correctness are good in abstaining from prediction and vice versa (rank corr. 0.894), and both tasks are also correlated with the log-likelihood. These methods form a cluster that captures predictive uncertainty capabilities. The log-likelihood, a proper scoring rule, is also highly correlated with the Brier score, another proper scoring rule. Both metrics evaluate both the uncertainty estimates and the correctness of prediction in one, which also explains why the Brier score is highly correlated with accuracy. Aleatoric uncertainty, OOD detection, and calibration (via ECE) form clusters on their own, underlining that they benchmark different forms of uncertainty.

As there are different groups of tasks, there is no one-fits-all uncertainty estimator. Table 3 demonstrates this by ranking all methods on all tasks. An uncertainty estimator has to be chosen or developed for the specific task a practitioner is interested in. If the task is unknown, deep ensembles offer a good compromise, but even the baseline is a good starting point if the runtime costs of deep ensembles are too high.

¹Note that we train with only one label per input.

Benchmarking Uncertainty Disentanglement

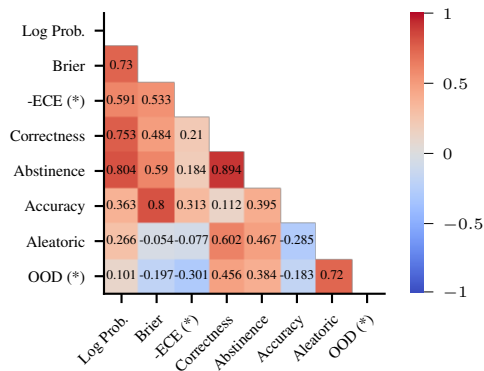


Figure 16. The rank correlation of metrics on CIFAR-10 is notably different from that of ImageNet (Figure 7). Rank correlation of metric pairs calculated over all (method, aggregator) pairs. (*) OOD’s and ECE’s correlation with the other methods is sensitive to the choice of aggregator, see Appendix J.

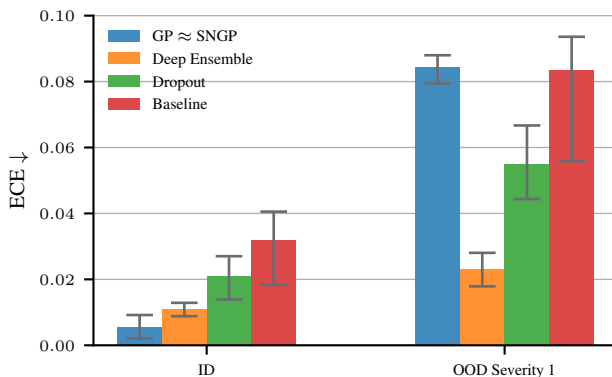


Figure 17. While most methods preserve their rankings as the dataset becomes more and more OOD via CIFAR-10C corruptions, SNGP variants break down to the baseline level already for severity level one on the ECE metric.

E.7. Performance depends on implementation details

We conclude with an example where defining the specific task is of particular importance. SNGP shows a highly different performance depending on the dataset and task it is trained for. Figure 17 shows that SNGP (and GP) provide the best-calibrated uncertainties ID, but already at the lowest OOD perturbation level drop to the baseline level, with the ECE jumping from 0.005 to 0.084. In Appendix I.4, we show that the best way to aggregate SNGP’s posterior $q(f)$ into an uncertainty score $u(x)$ is different when optimizing for the ECE versus for correctness prediction. This is not a bug; we implemented SNGP three times from scratch with the same results. It rather goes to show that subtle design choices greatly affect performance. Hence, we encourage, as best practice, to first define the uncertainty quantification task at hand and then develop a specialized uncertainty estimator.

G. Further Results on the Information-Theoretical Decomposition

G.1. OOD Generalization Performance

The ID correlation of the IT decomposition’s components using different distributional methods is discussed in Appendix F.1 and Section 3.1 of the main paper. In this section, we focus on how this correlation changes as we go more and more OOD.

G.1.1. IMAGENET

Figure 18 shows the generalization performance of benchmarked distributional methods using the IT decomposition at severity levels two and five. Only SNGP and GP show a reasonably low rank correlation between the IT aleatoric and epistemic components across all severities. These sources generally become less correlated as we go more OOD. Balanced mixtures of OOD samples lead to higher correlations, but the ranking of methods remains unchanged.

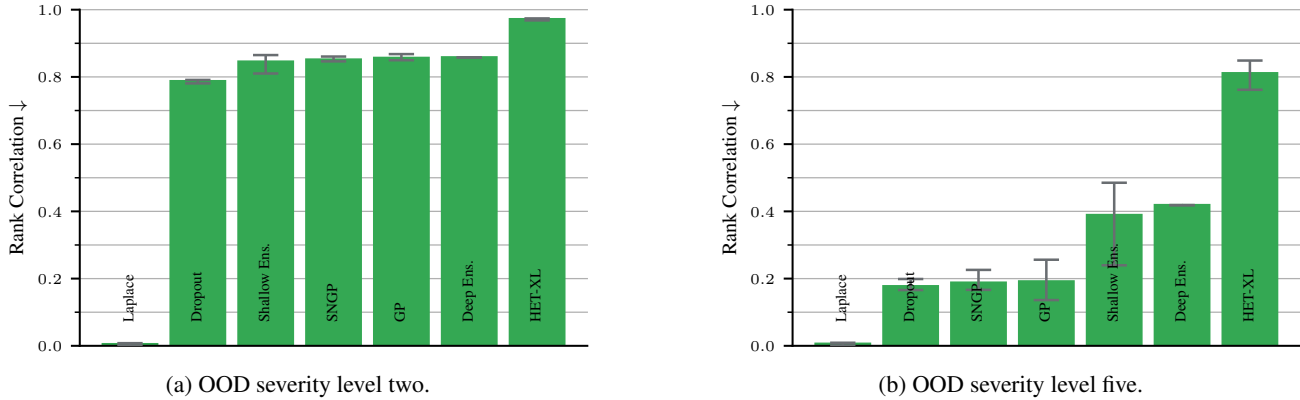


Figure 18. Laplace has uncorrelated epistemic and aleatoric IT estimates across different OOD levels on the ImageNet validation set. Other methods only achieve a low rank correlation at the highest severity levels. OOD rank correlation results of the IT decomposition using different distributional methods.

G.1.2. CIFAR-10

Results for severity levels two and five are shown in Figure 19. Only SNGP and GP show a reasonably low rank correlation between the IT aleatoric and epistemic components across all severities. These sources become slightly less correlated as we go more OOD—much less so than on ImageNet (Figure 18). Balanced mixtures of OOD samples lead to higher correlations, but the ranking of methods remains unchanged.

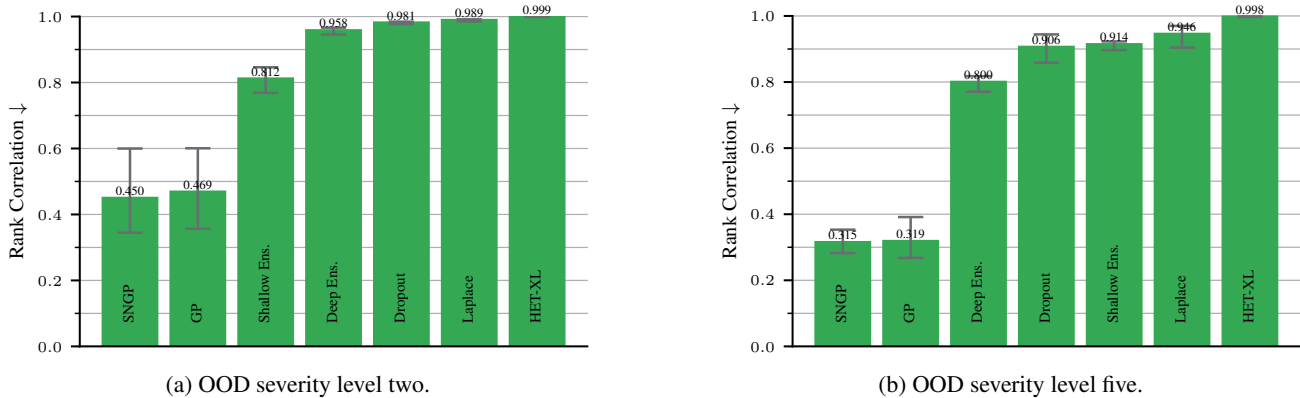


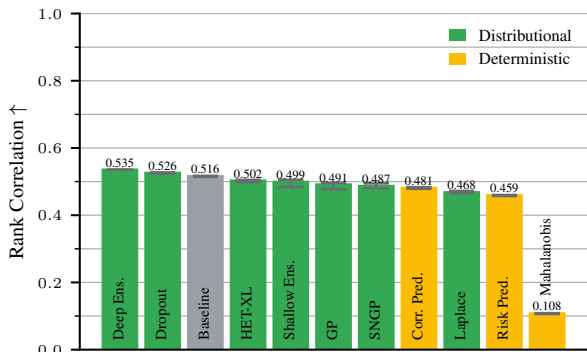
Figure 19. On CIFAR-10, SNGPs show a stably low rank correlation between the Bregman epistemic and aleatoric components across different OOD levels. OOD rank correlation results of the IT decomposition’s epistemic and aleatoric components using different distributional methods.

G.2. Performance of Decorrelated Methods using the Information-Theoretical Components

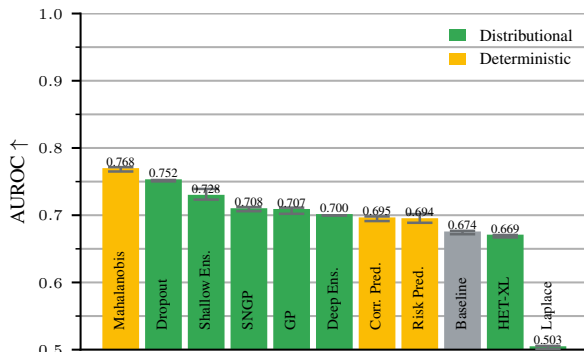
G.2.1. IMAGENET

In Section 3.1 of the main paper, we show that using Laplace, the aleatoric and epistemic components of the IT decomposition can become uncorrelated. In this section, we show that they do not perform well on the tasks they are made for. Figure 20

shows that Laplace does not match the baseline with its best estimator when Laplace uses the estimators of the IT decomposition—in fact, it even performs randomly on the OOD detection task. This limitation prohibits its practical use, and we cannot benefit from the less correlated components.



(a) Rank correlation with the ground-truth human uncertainties.

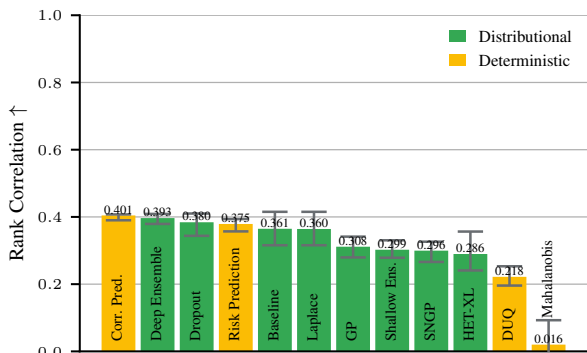


(b) AUROC of the OOD detection task using severity level two.

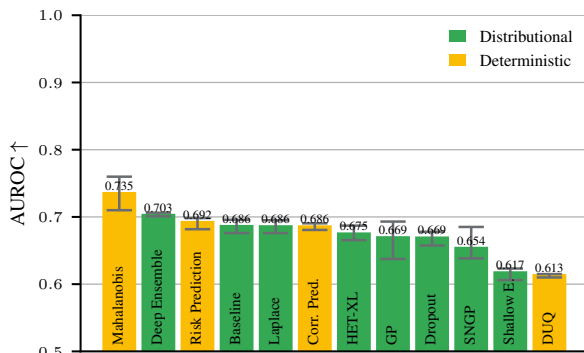
Figure 20. On ImageNet, Laplace cannot match the baseline when using the estimators of the IT decomposition. All other methods are equipped with their best-performing estimator for the respective tasks, showing that specialized estimators work better.

G.2.2. CIFAR-10

In Appendix F.1, we discuss that when using SNGP variants or shallow ensembles on CIFAR-10, the aleatoric and epistemic components of the IT decomposition can become quite uncorrelated. In this section, we show that while these components might be uncorrelated, they are not performant on the tasks they are made for. Figure 21 shows that both the SNGP variants and the shallow ensemble method underperform the baseline with its best estimator when the SNGPs and the shallow ensemble use the estimators of the IT decomposition.



(a) Rank correlation between methods and the Bregman aleatoric component.



(b) AUROC of OOD detection performance of methods using perturbations of severity level two.

Figure 21. SNGP, GP, and shallow ensemble underperform the baseline on CIFAR-10 when using the estimators of the IT decomposition. All other methods are equipped with their best-performing estimator for the respective tasks, showing that the IT decomposition is not practically beneficial.

H. Further Results on the Bregman Decomposition

H.1. Correlation of Components and Limitations

Let us carry out the same experiments for Bregman as we did for the IT decomposition in Appendix F.1. As the Bregman and risk decompositions (Equations 2 and 35) consider the *ground-truth* label distribution as the aleatoric component, we

use the IT aleatoric uncertainty as an estimator of it.

H.1.1. IMAGENET

Figure 22 shows that there is a considerable rank correlation between the Bregman ground-truth aleatoric and bias components but is not severe enough such that it prevents the theoretical possibility of disentangling them via estimators. On the right, we can also see that the IT and Bregman correlation results are very similar.

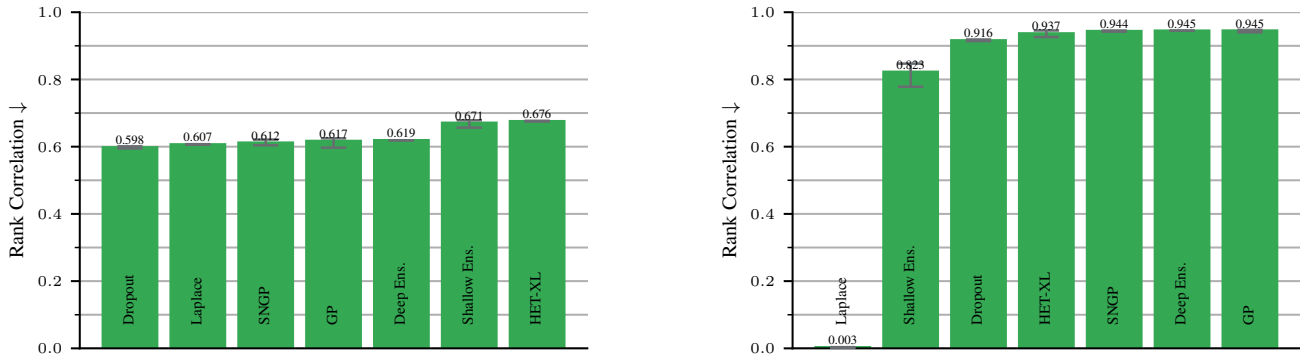


Figure 22. *Left.* On ImageNet-Real, the rank correlation of the Bregman aleatoric and bias terms is between 0.6 and 0.7 for all distributional methods we benchmark. *Right.* The Bregman decomposition shows similar rank correlation results to the IT decomposition between the estimated aleatoric uncertainty and the epistemic component on the ImageNet validation dataset.

H.1.2. CIFAR-10

We see in Figure 23 that the results using Bregman are virtually the same: no distributional method can provide decorrelated uncertainty estimates that align with their ground truth. Given the high ground-truth rank correlation of the aleatoric and bias components shown in Figure 25, there also seems to be a fundamental limitation in disentangling them. Note that none of the benchmarked distributional methods truly separate the aleatoric and bias terms.

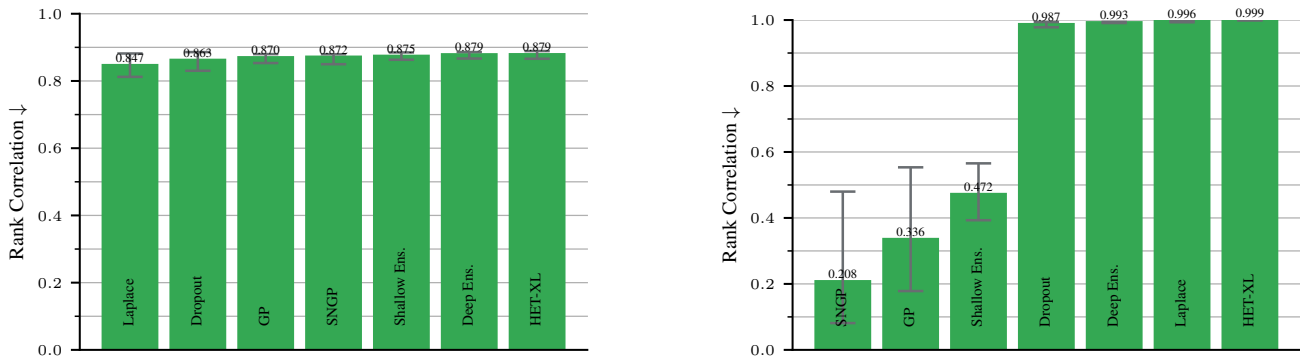


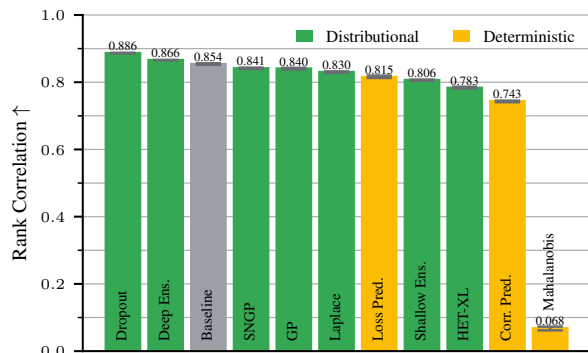
Figure 23. *Left.* The rank correlation of the Bregman aleatoric and bias terms is above 0.84 for all distributional methods we benchmark on CIFAR-10. *Right.* On CIFAR-10, the Bregman decomposition shows similar rank correlation results to the IT decomposition between the estimated aleatoric uncertainty and the epistemic component.

H.2. Alignment of Methods with the Bregman Bias

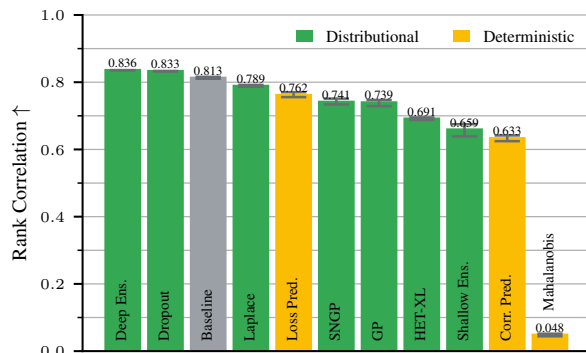
H.2.1. IMAGENET

The rank correlation of benchmarked methods with the bias component of the Bregman decomposition is shown in Figure 24 for ID and OOD with severity two. Only dropout and deep ensembles are notably more correlated with the Bregman bias component than the baseline, but most methods exhibit a high rank correlation (≥ 0.8) (unlike CIFAR-10 below).

This suggests that uncertainty estimators are most aligned with the bias component of Bregman out of the three. Considering an OOD dataset with severity-two perturbations, all methods become less correlated with bias, unlike CIFAR-10 below.



(a) ID rank correlation of methods with the Bregman decomposition's bias component.

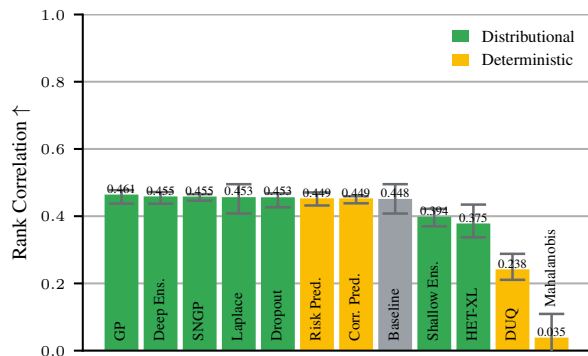


(b) OOD rank correlation of methods with the Bregman bias using severity-two perturbations.

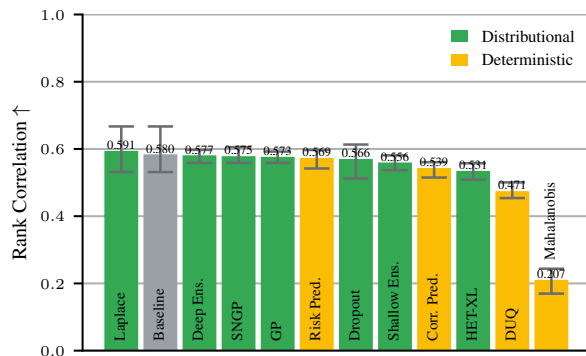
Figure 24. Only dropout and deep ensembles are notably more correlated with the Bregman bias component than the baseline on ImageNet. Most methods exhibit a high rank correlation (≥ 0.8). When going more OOD, all methods become less correlated with bias.

H.2.2. CIFAR-10

The rank correlation of benchmarked methods with the bias component of the Bregman decomposition is shown in Figure 25 for ID and OOD with severity two. None of the benchmarked methods are notably more correlated with the Bregman bias component than the baseline. Considering an OOD dataset with severity-two perturbations, all methods become better correlated with bias.



(a) ID rank correlation of methods with the Bregman decomposition's bias component.



(b) OOD rank correlation of methods with the Bregman bias using severity-two perturbations.

Figure 25. None of the benchmarked methods are significantly more correlated with the Bregman bias component than the baseline on CIFAR-10. When going more OOD, all methods become better correlated with bias.

I. Further Practical Results

I.1. Correctness Prediction

I.1.1. IMAGENET

We show the correctness prediction performance of methods on OOD and mixed ID + OOD datasets in Figure 26. OOD, no method has an edge over the baseline. We observe a consistent but not significant degradation of performance across

methods on both dataset types. Interestingly, the performance of Mahalanobis does not increase on mixed datasets, even though models perform worse on OOD images than on ID ones, and it is a suitable OOD detector.

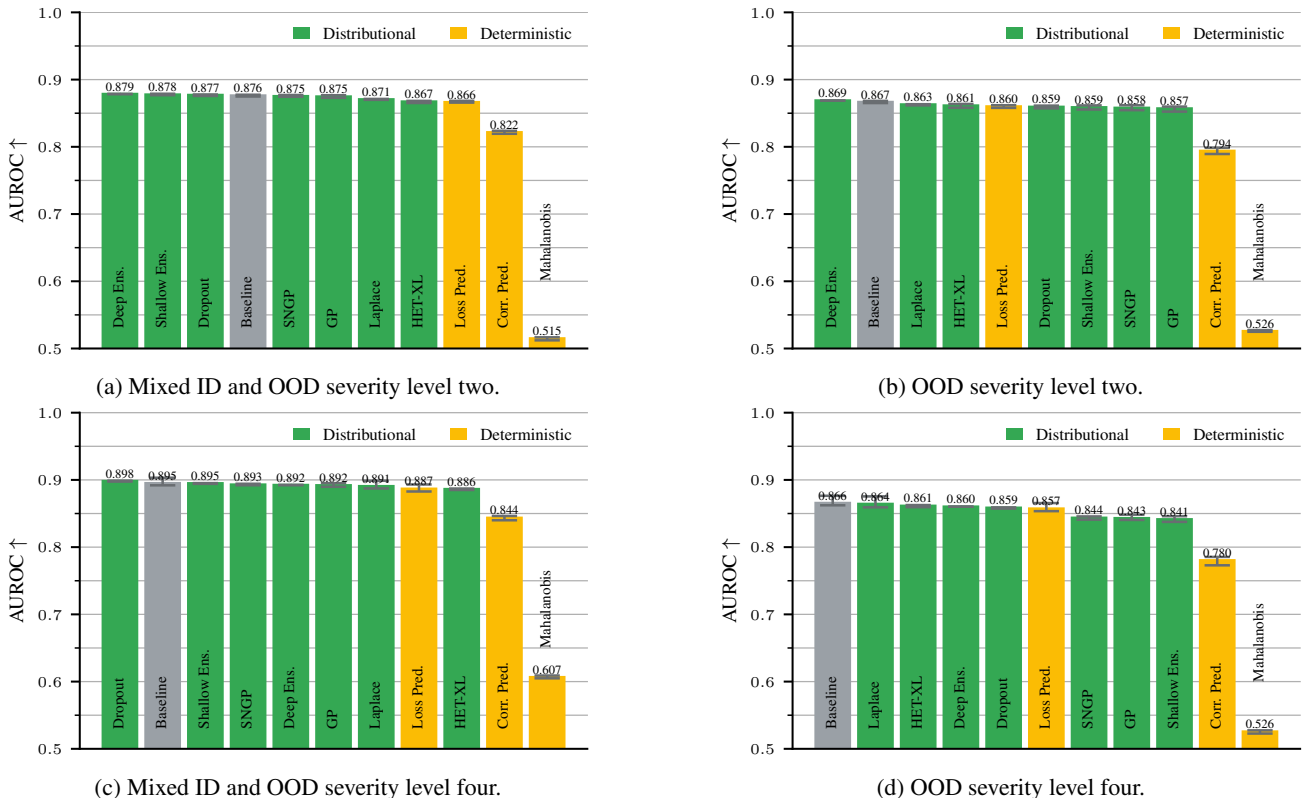


Figure 26. On ImageNet, the correctness prediction performance of the methods consistently drops both on completely OOD datasets (top right and bottom right) and on balanced mixtures of ID and OOD datasets (top left and bottom left) as we increase the severity. The performance of Mahalanobis increases on mixed datasets, as models perform worse on OOD images than on ID ones, and it is a suitable OOD detector.

As we have access to validation sets with multiple labels per input, the notion of a “correct” prediction becomes unclear. In the main paper and the previous plots, we focus on the canonical notion of correctness: whether the model predicts the most likely class. A related notion of correctness is soft correctness: here, the model does not receive a binary reward for its prediction but rather a continuous number $c \in [0, 1]$: this gives the ground-truth probability of the predicted class being the correct one. We can calculate a similar AUROC correctness score as before, but as the AUROC requires binary labels, one has to unroll the continuous correctness value: e.g., if $c = 0.8$ for an input x , one can represent this with 8 correct and 2 incorrect predictions on the same input x . The resulting AUROC has a theoretical limit strictly less than one when $c < 1$. The corresponding results are shown in Figure 27 that follows the same structure as Figure 26. The saturation of methods is unchanged; however, the ordering is affected.

I.1.2. CIFAR-10

We show the correctness prediction performance of methods on OOD and mixed ID + OOD datasets in Figure 28. We observe a consistent degradation of performance across methods on both dataset types. The only exception is the Mahalanobis method: here, we observe an increase in performance. For the mixed datasets, it can be explained by the fact that the models perform worse on OOD images than on ID ones, making the Mahalanobis method a suitable estimator of correctness. However, the result is unexpected for solely OOD samples.

As we have access to validation sets with multiple labels per input, the notion of a “correct” prediction becomes unclear. In the main paper and the previous plots, we focus on the canonical notion of correctness: whether the model predicts the most likely class. A related notion of correctness is soft correctness: here, the model does not receive a binary reward for its prediction but rather a continuous number $c \in [0, 1]$: this gives the ground-truth probability of the predicted class being the

correct one. We can calculate a similar AUROC correctness score as before, but as the AUROC requires binary labels, one has to unroll the continuous correctness value: e.g., if $c = 0.8$ for an input x , one can represent this with 8 correct and 2 incorrect predictions on the same input x . The resulting AUROC has a theoretical limit strictly less than one when $c < 1$. The corresponding results are shown in Figure 29 that follows the same structure as Figure 28. The saturation of methods is unchanged; however, the ordering is affected.

I.2. Performance Tendency for Increasing Severity

In Section 3.5 of the main paper and Appendix F.3, we show the performance of dropout when going OOD, claiming that it is prototypical for other methods. Figure 30 and Figure 31 show for CIFAR-10 and ImageNet, respectively, that dropout is not an outlier and other methods show very similar generalization capabilities.

I.3. OOD Detection

I.3.1. IMAGENET

In Appendix F.4, we hint at the fact that nearly all methods show a steady increase in OOD detection performance as we increase the severity of the perturbed half of the dataset. Figure 32 shows how the performance of each method changes as we increase the severity level. We can see a steady increase in OOD detection performance for all methods. However, the specialized OOD detector, Mahalanobis, benefits less than the other methods. In particular, at severity level three, dropout becomes best on average, and shallow ensembles also overtake the Mahalanobis method. At severity level five, Mahalanobis becomes the *worst* OOD detector out of the benchmarked methods. This may be because Mahalanobis was trained to detect samples at severity level two and cannot generalize as well to higher severity levels as the other methods.

I.3.2. CIFAR-10

Figure 33 shows how the OOD detection performance of each method changes as we increase the severity level. We can see a steady increase in OOD detection performance for all methods. However, the specialized OOD detector, Mahalanobis, benefits less than the other methods. In particular, at severity level four, deep ensembles become best on average, and SNGP variants also overtake the Mahalanobis method at severity level five.

I.4. Sensitivity to the Choice of Aggregator

This section demonstrates that the choice of aggregator is of crucial importance for specific methods and tasks, showing results on the CIFAR-10 dataset. Some tasks, such as correctness prediction and abstained prediction tasks (ID), have highly correlated performances across methods. This is an intuitive result, as both tasks fundamentally require telling apart correct and incorrect samples. However, this is not always the case. OOD detection and the ECE score are two such tasks/metrics.

Considering that the ECE is closely connected to correctness but on a much more fine-grained scale than the binary correctness prediction task, one would expect that a high ECE score translates over to a high correctness AUROC score. Surprisingly, this is often not the case for distributional methods, as shown in Figure 34. Importantly, the choice of aggregator (see Appendix B) per distributional method has a high influence on predictive power and the correlation of performances: SNGP variants, Laplace, deep ensemble, and HET-XL all have different optimal aggregators for correctness and ECE, and depending on which task we optimize the estimator for, we can obtain both positive and negative rank correlation among the task performances. SNGP variants are particular outliers: they have the largest trade-off between the best possible performances in the two tasks.

I.5. ECE OOD Generalization

In Figure 17 of the appendix, we show that SNGP variants are the most calibrated ID on the CIFAR-10 dataset, but they also break down to the baseline level from severity level one. The four methods that bring considerable improvements compared to the baseline ID are SNGP, GP, deep ensemble, and dropout, in the order of increasing ECE. Figure 35 shows that as we increase the severity level, only deep ensemble and dropout are capable of performing considerably better ($\geq .05$ ECE improvement) than the baseline on average.

I.6. Log Probability Proper Scoring Rule

In Figure 36, we present the methods’ results on the log probability proper scoring rule considering the CIFAR-10 dataset. We find that deep ensemble, dropout, and Laplace are the only methods that consistently outperform the baseline on average, both ID and OOD for all severity levels. Still, their performance advantage is within the error bars of the baseline.

J. Sensitivity of Correlation Matrix

Figure 37 shows rank correlation results across metrics using different estimators on the ImageNet dataset. Notably, as foreshadowed in the caption of Figure 7, the OOD and ECE metrics exhibit different rank correlation scores depending on the estimator we choose.

K. Training and Implementation Details

For both datasets, we train and evaluate on an NVIDIA GeForce RTX 2080 Ti. We only use an NVIDIA A100 Tensor Core GPU for the construction of the Laplace approximation on ImageNet, owing to the VRAM requirements of this method.

K.1. CIFAR-10

For CIFAR-10, we follow the augmentations and training schedules of the [uncertainty_baselines](#) GitHub repository. In particular, we train a Wide ResNet 28-10 (Zagoruyko & Komodakis, 2016) for 200 epochs with a step decay schedule at [60, 120, 160] epochs with decay rate 0.2. We use stochastic gradient descent with momentum 0.9 and a batch size of 128. Our training augmentation comprises a random crop using padding 2 and a random flip on the vertical axis with probability 0.5. The learning rate and weight decay hyperparameters are chosen by the Bayesian optimization scheme of Weights and Biases (Biewald, 2020). The additional hyperparameters of benchmarked methods are determined by either using values suggested by the original authors or including these in the hyperparameter sweep.

K.2. ImageNet

On ImageNet, we fine-tune a pretrained ResNet 50 (He et al., 2016) using the `resnet50.a1_in1k` parameters from the `timm` library as initialization. We fine-tune for 50 epochs following a cosine learning rate schedule (Loshchilov & Hutter, 2017) using the AdamW optimizer (Loshchilov & Hutter, 2019) and a learning rate warmup period of 5 epochs. We use a batch size of 128 with 16 accumulation steps, resulting in an effective batch size of 1024. The hyperparameters are chosen identically to those on CIFAR-10 (see Appendix K.1).

K.3. Runtime

Table 4 and Table 5 show statistics of the per-epoch runtime for each method on ImageNet and CIFAR-10, respectively. As Laplace, Mahalanobis, and deep ensemble are post-hoc methods, their reported time comprises the construction of the method and its evaluation.

L. Visualization of Images and Label Distributions

This section displays both easy (low human uncertainty) ImageNet samples in Figure 38 and hard (high human uncertainty) ones in Figure 39 using the ImageNet-ReaL labels and ImageNet-C perturbations.

Figure 40 and Figure 41 give summary statistics of the label distributions of ImageNet-ReaL and CIFAR-10H, respectively.

Method	Mean (s)	Min (s)	Max (s)	Std Dev (s)
Deterministic	90.1829	83.7423	122.5249	6.0592
GP	91.2036	83.9509	236.0735	9.5304
Correctness Prediction	91.8449	82.6371	133.5970	6.4087
Shallow Ensemble	91.9881	83.0224	119.2658	5.1190
Risk Prediction	94.7748	83.2490	123.5470	8.8957
SNGP	96.8125	88.4623	129.1145	6.4576
HET-XL	99.0390	89.4465	161.1186	8.7700
Dropout	134.0979	126.2741	188.3551	7.0563
DUQ	148.6540	137.8707	197.6619	7.9503
Laplace	273.2982	250.8563	307.6892	22.9335
Mahalanobis	370.4277	360.0912	376.3072	5.7191
Deep Ensemble	865.0582	835.1872	904.2822	22.8003

Table 4. Summary of per-epoch times for the benchmarked methods on CIFAR-10. As Laplace, Mahalanobis, and Deep Ensemble are post-hoc methods, their reported time comprises the construction of the method and its evaluation. Methods are sorted by increasing mean per-epoch runtime, separately for trained and post-hoc methods.

Method	Mean (s)	Min (s)	Max (s)	Std Dev (s)
Deterministic	2646.6512	2566.9703	2899.1143	44.7298
Risk Prediction	2692.1484	2637.5598	2908.5478	33.4657
Correctness Prediction	2732.9235	2562.6328	3284.8763	230.3123
Shallow Ensemble	2803.4469	2671.9823	3268.5573	181.0415
GP	3059.7266	2645.5470	3880.5432	399.7528
Dropout	3145.4667	3034.7784	3307.8100	80.4335
SNGP	3233.9454	3081.7184	3742.6995	145.4325
HET-XL	4018.7616	3915.7693	4214.7061	55.4737
Mahalanobis	33929.2063	32972.1129	35235.9114	956.6842
Laplace	52836.5020	52298.8008	53949.9782	588.1313
Deep Ensemble	161492.2153	161492.2153	161492.2153	0.0000

Table 5. Summary of per-epoch times for the benchmarked methods on ImageNet. As Laplace, Mahalanobis, and Deep Ensemble are post-hoc methods, their reported time comprises the construction of the method and its evaluation. Methods are sorted by increasing mean per-epoch runtime separately for trained and post-hoc methods.

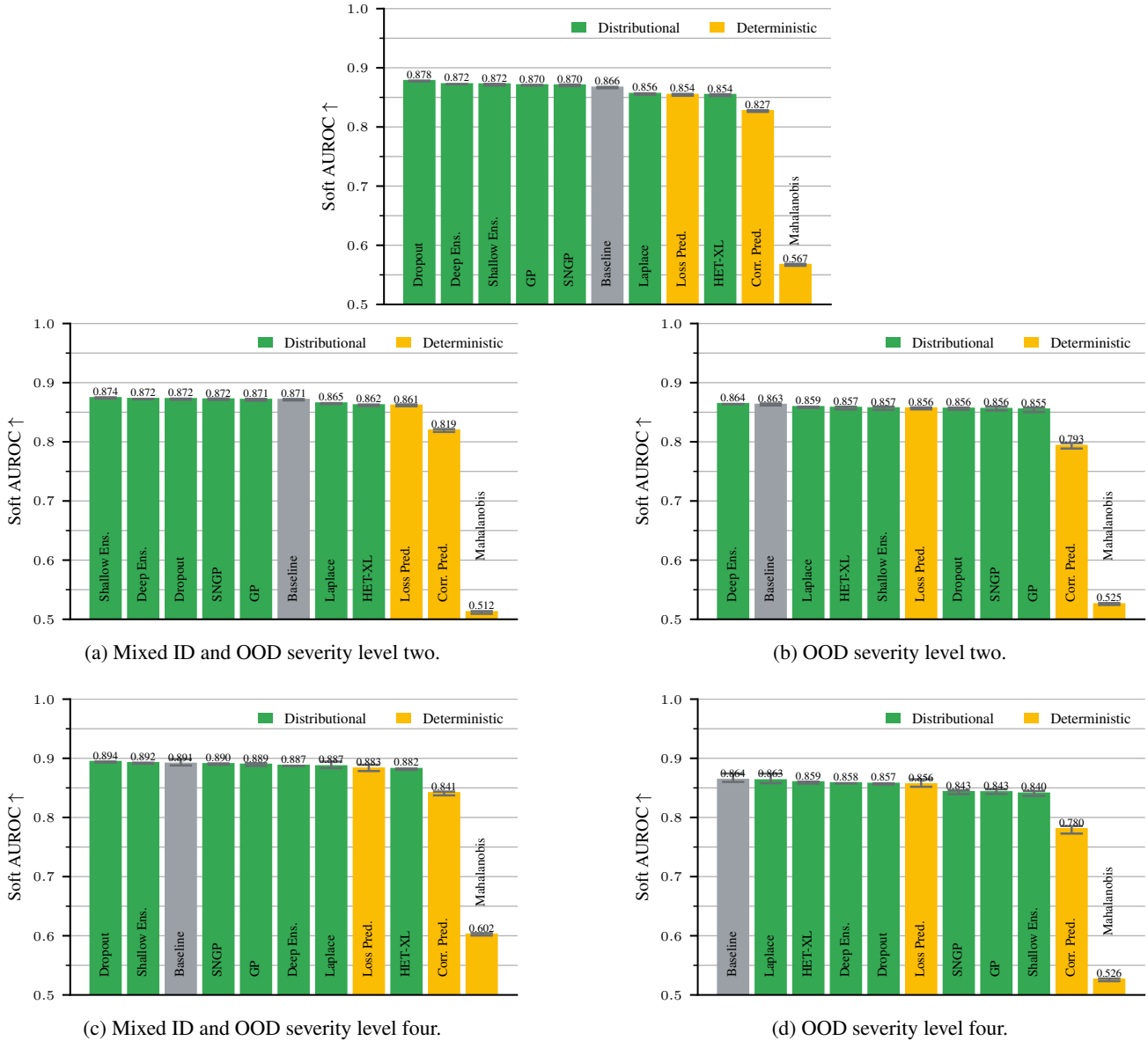
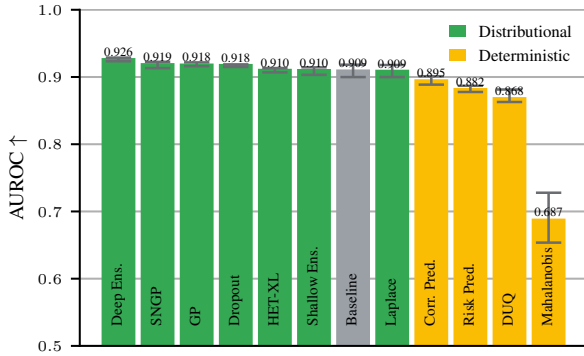
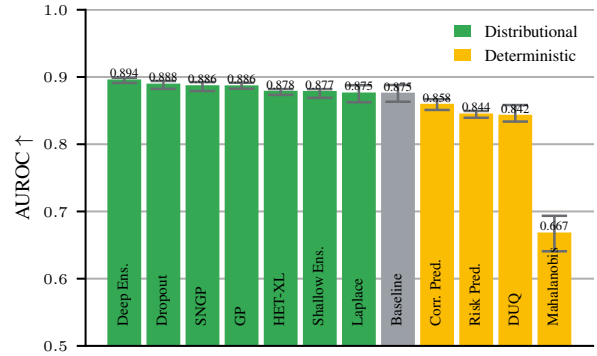


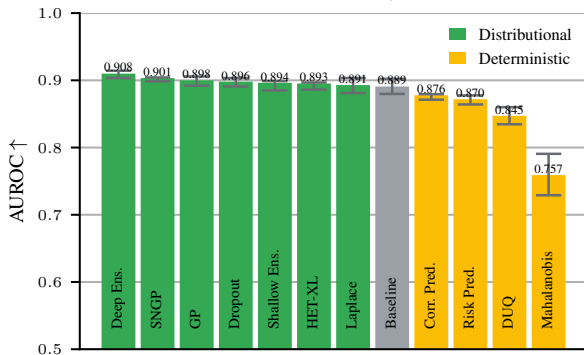
Figure 27. Slightly altering the evaluation criterion has an influence on the ranking of methods on the ImageNet validation set, evidencing that the methods are saturated on the correctness prediction task. Variant of 26 where correctness is calculated w.r.t. the soft labels of ImageNet-Real. ID results are added on top.



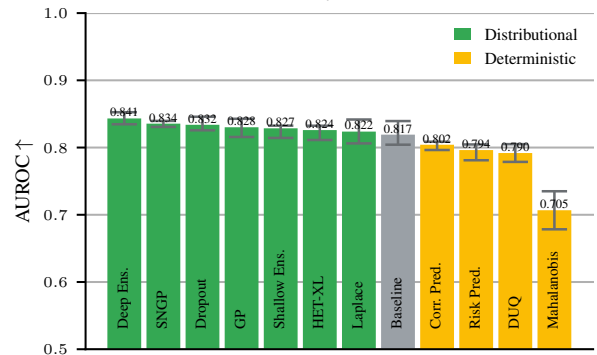
(a) Mixed ID and OOD severity level two.



(b) OOD severity level two.



(c) Mixed ID and OOD severity level four.



(d) OOD severity level four.

Figure 28. On CIFAR-10, the performance of the methods consistently drops on the correctness prediction task, both on completely OOD datasets (top right and bottom right) and on balanced mixtures of ID and OOD datasets (top left and bottom left). A notable exception is the Mahalanobis method: here, we observe an increase in performance as we increase the level of severity. For the mixed datasets, this can be explained by the fact that the models perform worse on OOD images than on ID ones, making the Mahalanobis method a suitable estimator of correctness. However, the result is unexpected for solely OOD samples.

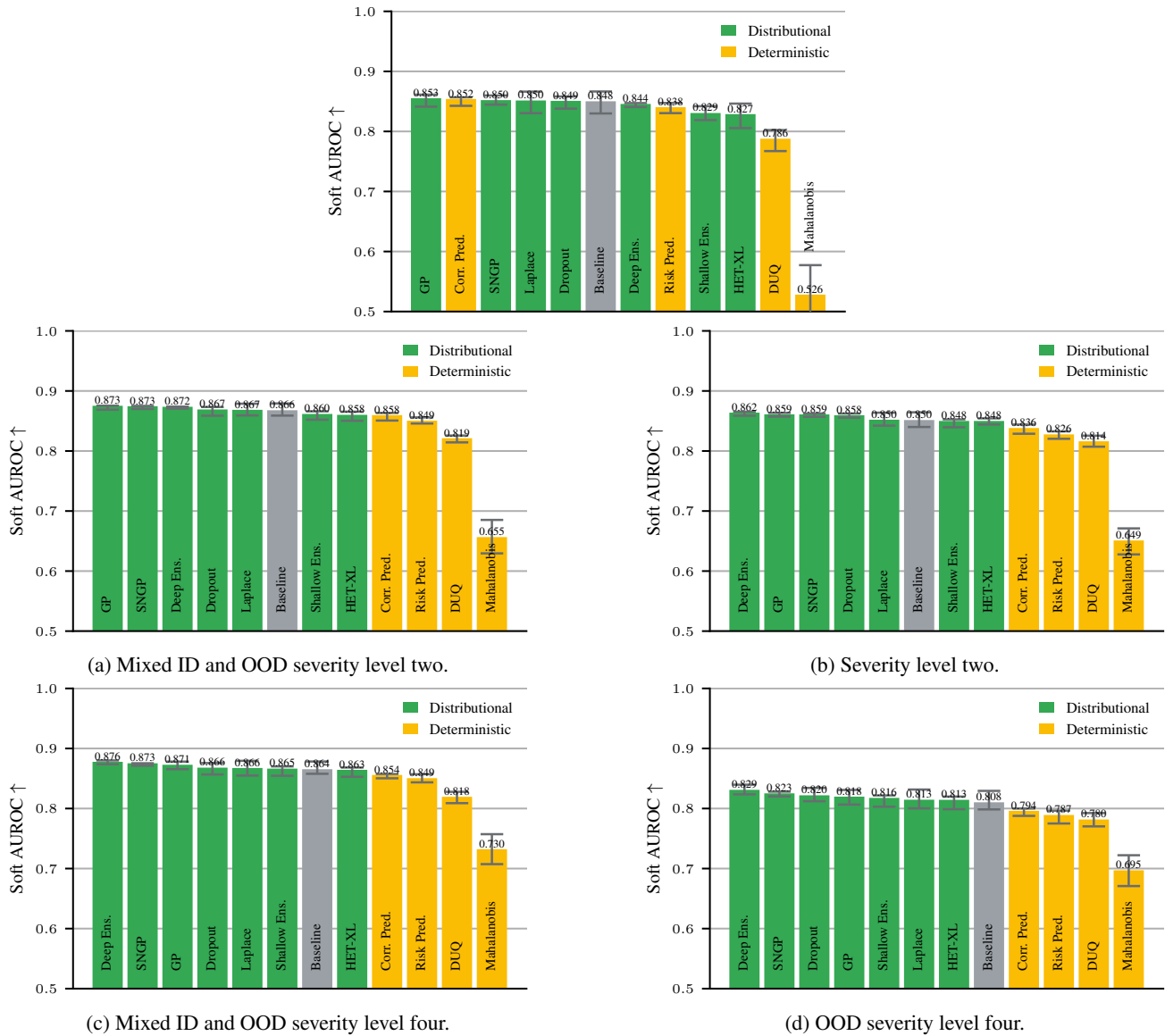


Figure 29. Slightly altering the evaluation criterion on the correctness prediction task changes the ranking of methods considerably, evidencing that the methods are saturated on CIFAR-10. Variant of Figure 28 w.r.t. soft label correctness with the ID results added on top.

Benchmarking Uncertainty Disentanglement

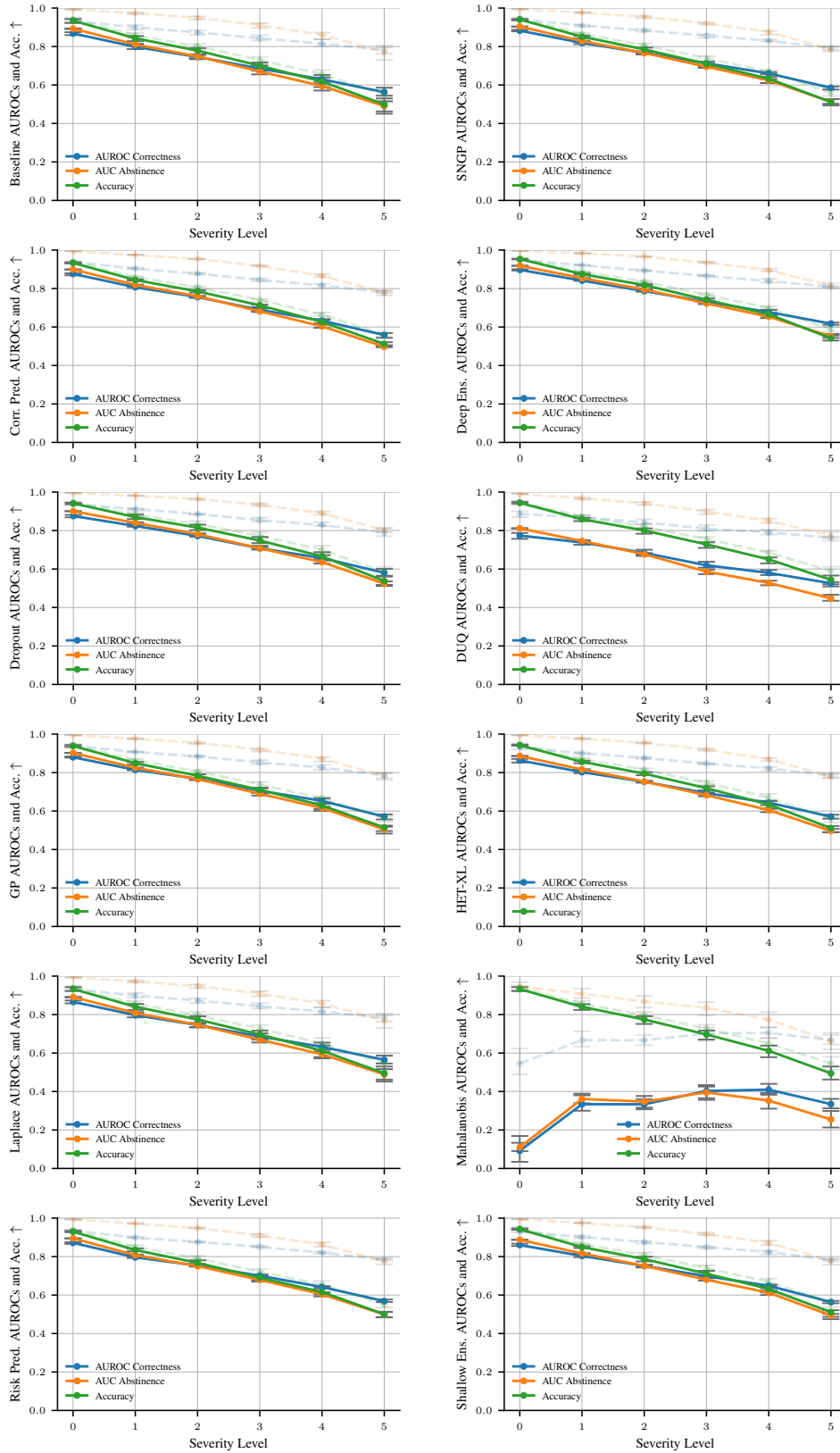


Figure 30. On CIFAR-10, all methods’ performance deteriorates at the same rate as the model’s accuracy on the correctness and abstinance tasks. The only exception is Mahalanobis, which is a specialized OOD detector.

Benchmarking Uncertainty Disentanglement

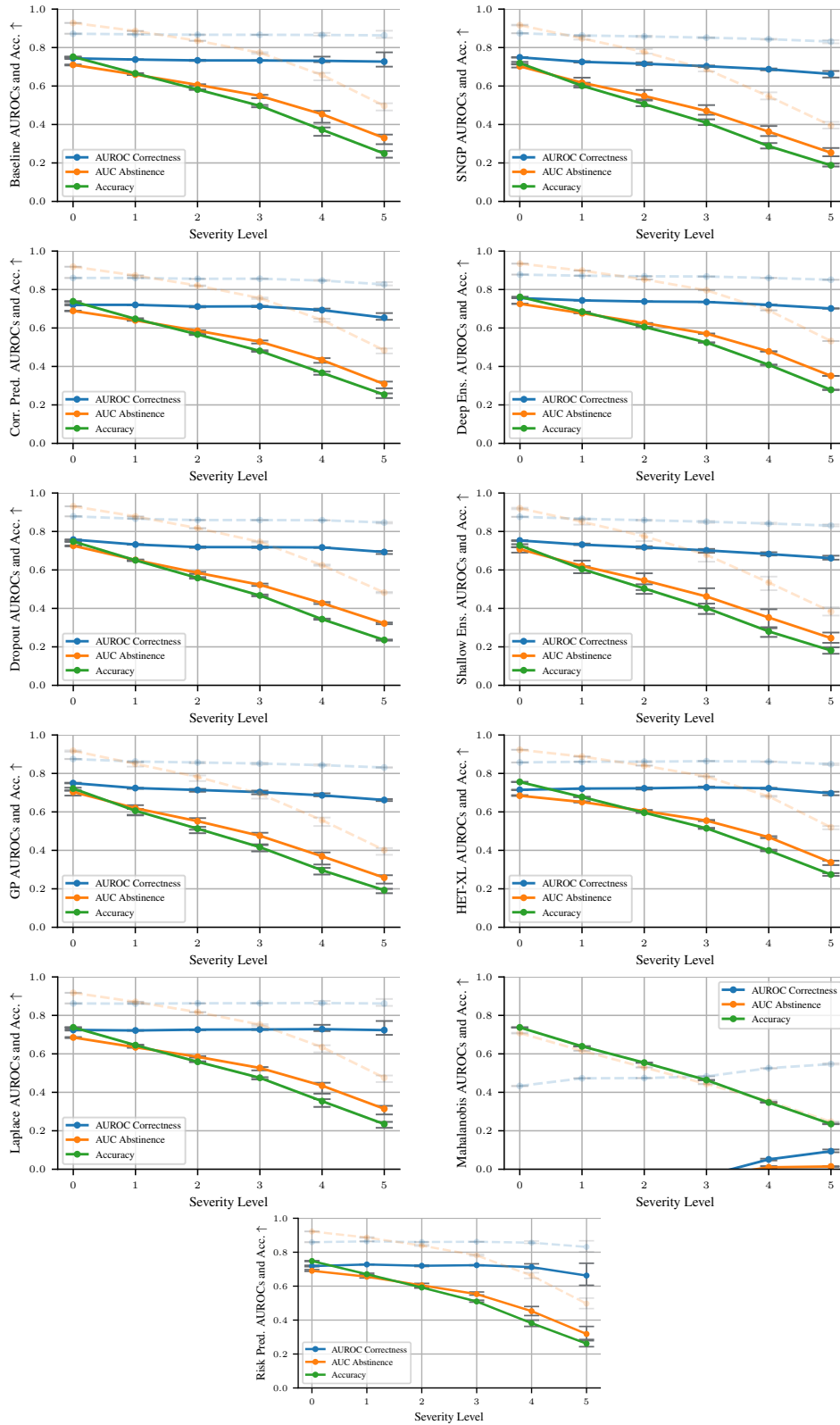
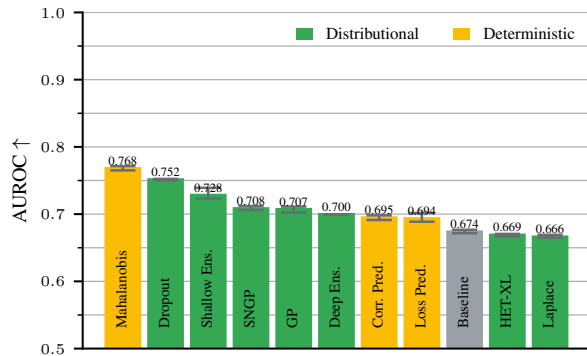
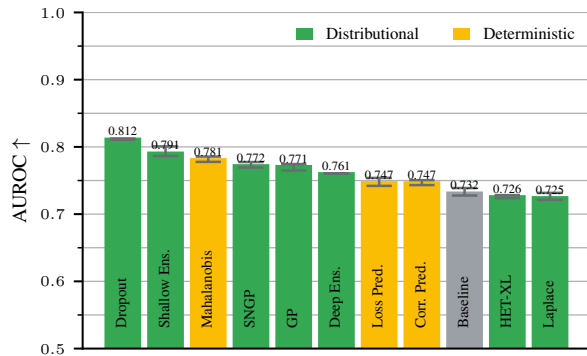


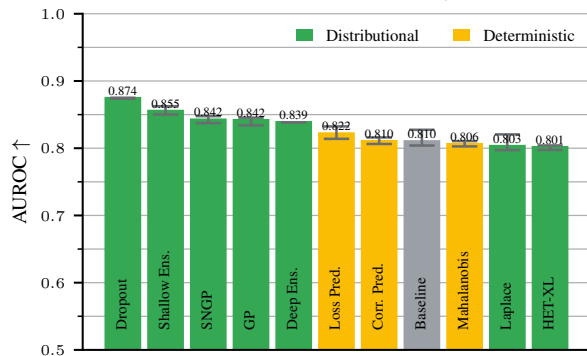
Figure 31. On ImageNet, the estimate for predictive correctness is much more robust to OOD perturbations than the model’s accuracy for all methods except Mahalanobis (which is a specialized OOD detector). The AUC abstinance score deteriorates at the same rate as the model’s accuracy.



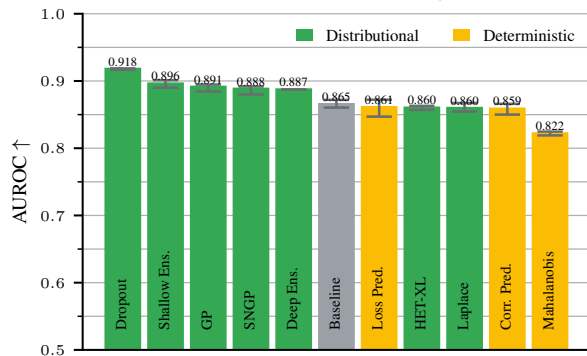
(a) AUROC OOD-ness with OOD severity level two.



(b) AUROC OOD-ness with OOD severity level three.



(c) AUROC OOD-ness with OOD severity level four.



(d) AUROC OOD-ness with OOD severity level five.

Figure 32. The OOD detection performance of all methods increases steadily as we increase the severity of the perturbed half of the mixed dataset on the ImageNet validation dataset. However, the specialized OOD detector, Mahalanobis, generalizes worse than the other methods.

Benchmarking Uncertainty Disentanglement

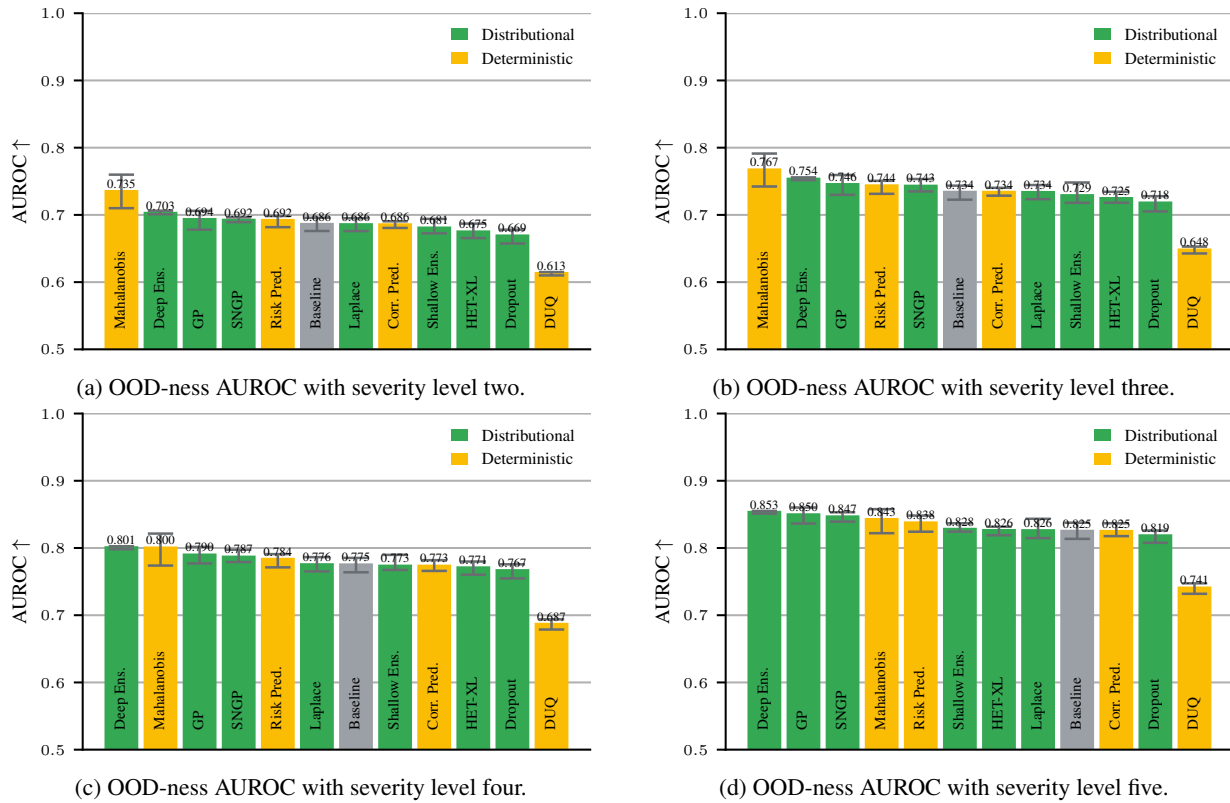


Figure 33. On CIFAR-10, the OOD detection performance of all methods increases steadily as we increase the severity of the perturbed half of the mixed dataset. Mahalanobis generalizes worse than the other methods.

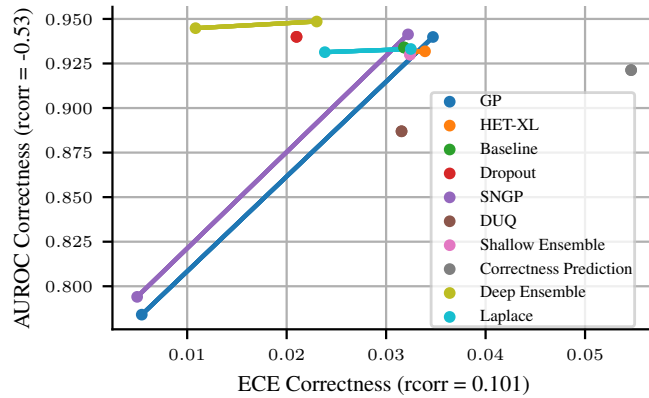


Figure 34. Depending on what task we optimize the aggregator for, we obtain notably different results for SNGP variants, deep ensembles, and Laplace networks on CIFAR-10. Each color corresponds to one method. The point pairs per method show the performance of the method optimized for ECE and that optimized for the correctness AUROC. The rank correlation values on the axes are with respect to all five seeds for all methods when optimized for ECE or the correctness AUROC. The sign of the correlation can flip based on what aggregator we choose.

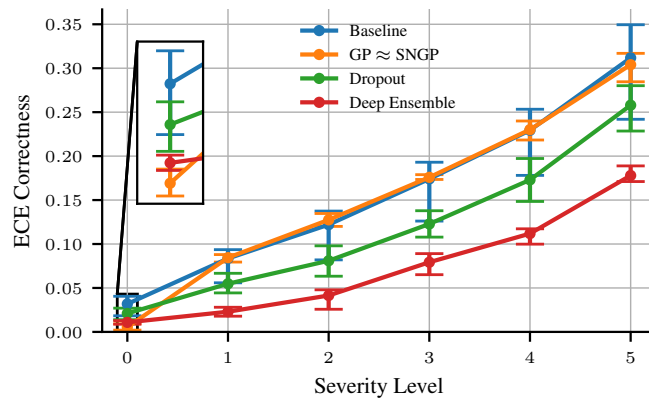


Figure 35. On CIFAR-10, as we increase the severity level, only deep ensemble and dropout are capable of performing considerably better than the baseline on average. SNGP is omitted as it behaves identically to GP. The baseline is included for reference.

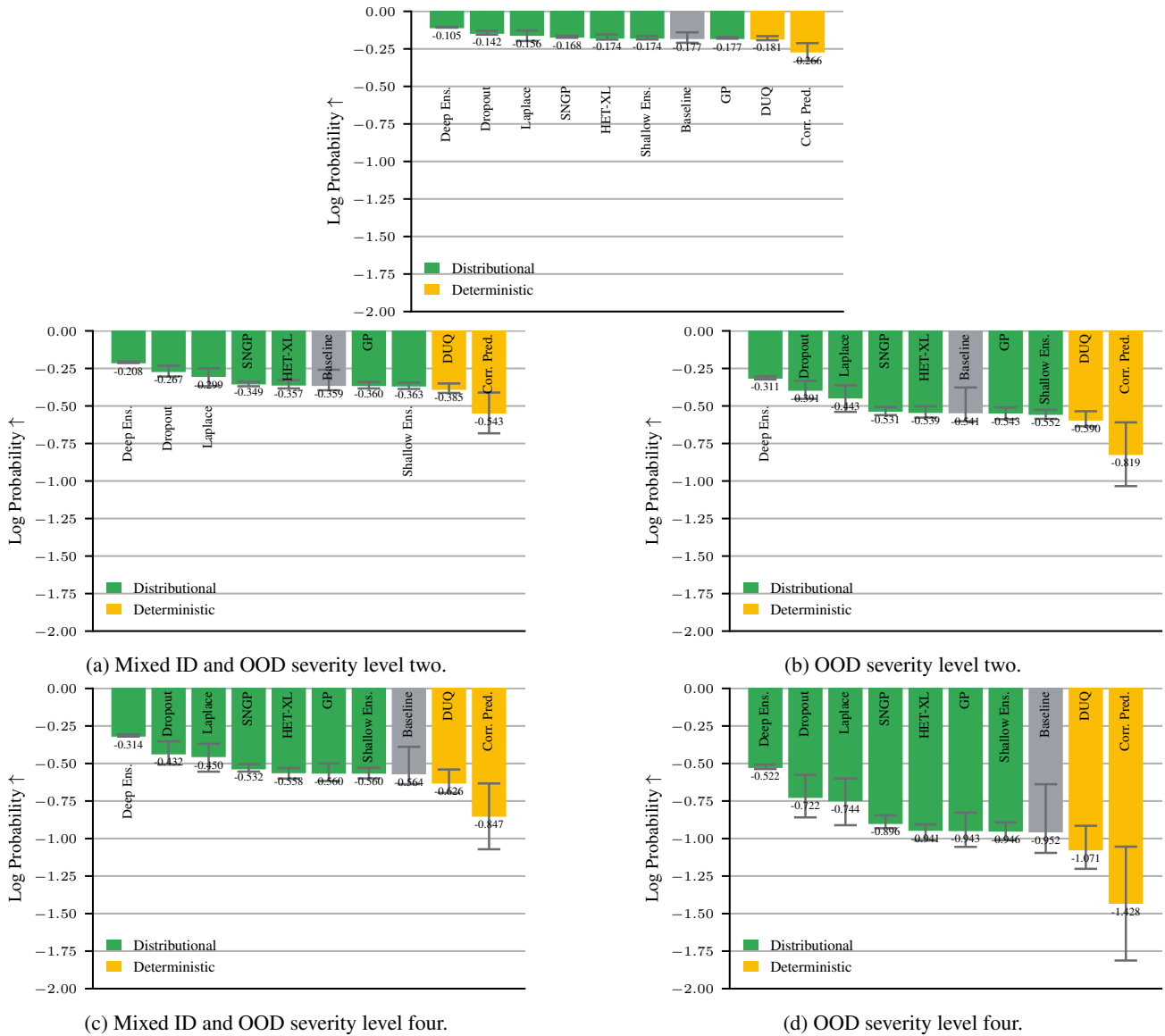
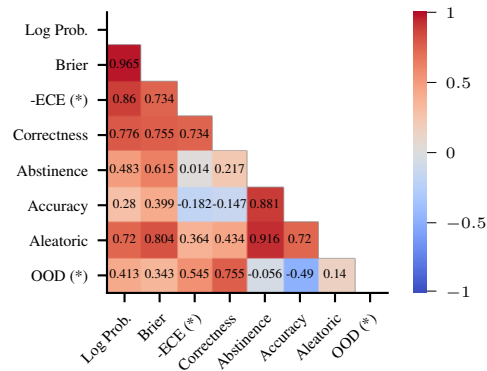
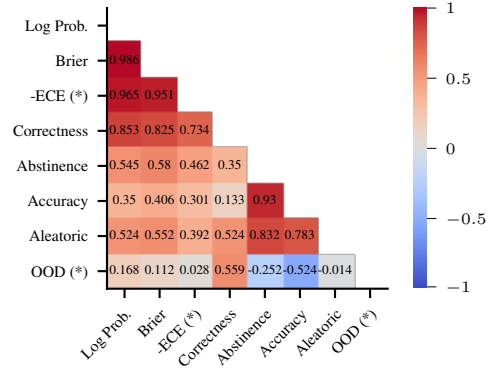


Figure 36. On CIFAR-10, deep ensemble, dropout, and Laplace are the only methods that consistently outperform the baseline *on average*, both ID and OOD for all severity levels when evaluating on the log probability proper scoring rule.

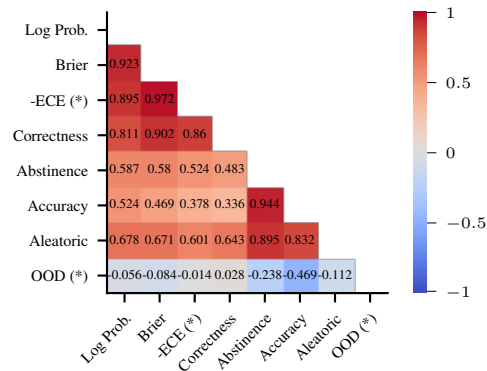
Benchmarking Uncertainty Disentanglement



(a) Results using one minus the maximum probability of \tilde{f} as the uncertainty estimator.



(b) Results using one minus the maximum probability of \bar{f} as the uncertainty estimator.



(c) Results using one minus the expected maximum probability as the uncertainty estimator.

Figure 37. Rank correlation results across metrics using different estimators on ImageNet. The OOD and ECE metrics exhibit highly different rank correlation scores depending on the estimator we choose.

Original Samples

Perturbed Samples



Figure 38. Easy ImageNet-ReaL cases with no human disagreement on the labels. OOD samples are of severity two.

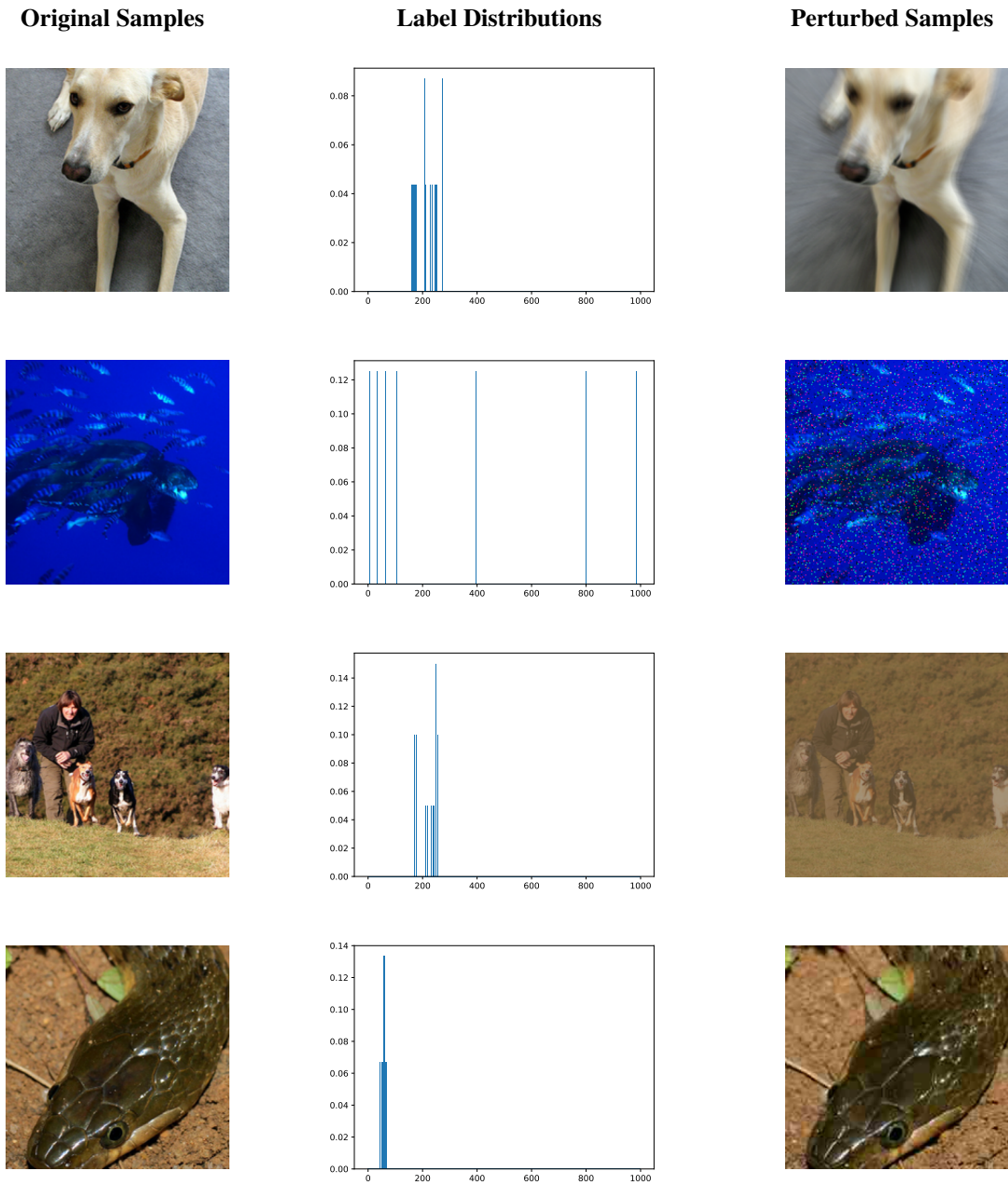


Figure 39. Hard ImageNet-ReaL cases with high human uncertainty (i.e., high disagreement among annotators on the correct label). OOD samples are of severity two.

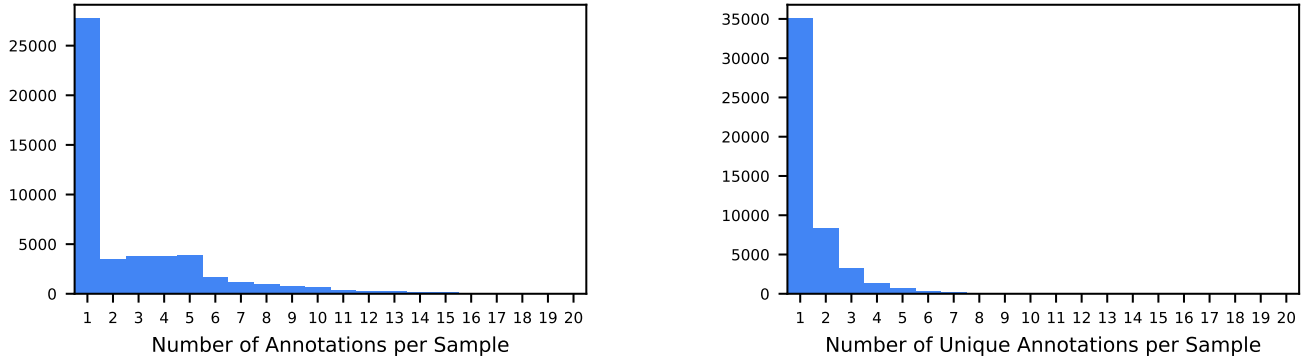


Figure 40. Histograms of the label distributions of the ImageNet-Real validation set.

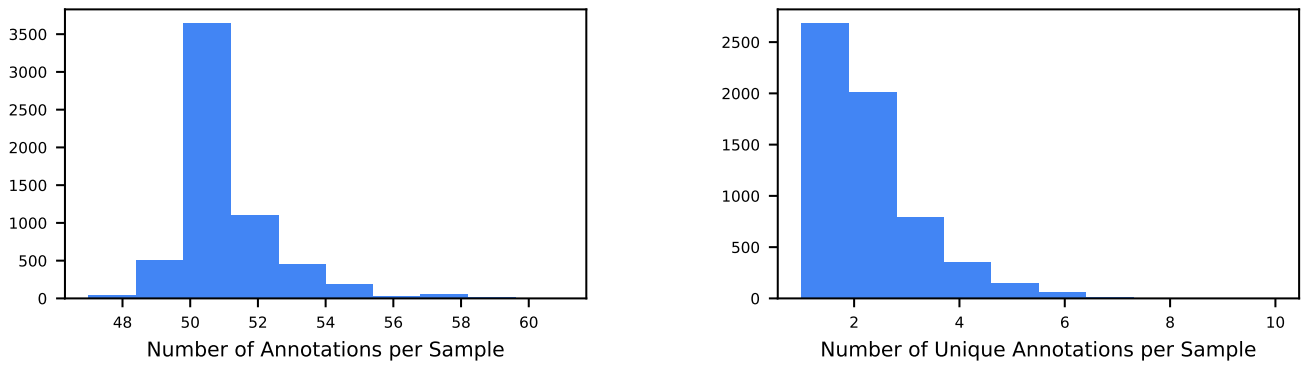


Figure 41. Histograms of the label distributions of the CIFAR-10H validation set.