
Supplementary material to

A Reproduction of Ensemble Distribution Distillation

Anonymous Author(s)

Affiliation

Address

email

1	Contents	
2	A The EnD² Algorithm	2
3	B Experiments on Synthetic Data	3
4	B.1 Methodology	3
5	B.1.1 Dataset	3
6	B.1.2 Model description and hyperparameters	3
7	B.1.3 Experimental setup and code	3
8	B.1.4 Computational requirements	4
9	B.2 Results	4
10	B.2.1 Classification accuracy	4
11	B.2.2 Visualization of uncertainty	4
12	C Computational requirements for reproduction	4
13	D Histograms	6
14	E Relative performance of EnD² compared to ensemble and original article	6

15 A The EnD² Algorithm

16 The original paper features an excellent description of the mathematical formulation of the EnD² model, but we did not
 17 find it immediately obvious how to translate this into an implementation in a modern deep learning framework. For this
 18 reason, we will now briefly describe it from an algorithmic-centred perspective using pseudocode and plain English.

19 The process of training an EnD² model is described in Algorithm 1. In practice, the optimization in line 7 can easily be
 20 achieved using the standard "fit" method of frameworks such as Keras and PyTorch, by constructing an intermediate
 21 dataset and using a custom loss function with a callback for annealing the temperature.

22 The intermediate dataset is constructed by first adding any auxiliary images to the training images, and then passing the
 23 extended image set as input to the ensemble. The ensemble should output an array of logits as described in line 5 of
 24 Algorithm 1. The new dataset is then formed by matching each image to its corresponding ensemble logits, using the
 25 latter as the target.

26 The custom loss function is described in Algorithm 2. This formulation includes temperature annealing. This loss
 27 function is the only modification necessary to adapt a general classification model into an EnD² model, providing it is
 28 then trained on an intermediate dataset as described in the previous paragraph. Note that this formulation assumes that
 29 the model outputs logits. This output can be converted into Dirichlet probabilities by applying the standard softmax
 30 operation.

Algorithm 1: Training algorithm for EnD² given an ensemble

Input : Ensemble En outputting logits, training data X (same as the ensemble is trained on), (optional) Out of
 distribution data X_{OOD}

Output : Trained EnD² model

```

1 if  $X_{OOD}$  not None then
2   |  $X = [X, X_{OOD}]$  // append OOD data to training set
3 end
4  $\phi = En.predict(X)$  //  $\exp(\phi)$  are the labels for EnD2
5 //  $\phi$  is a tensor of logits corresponding to the true distribution, each row corresponds to a model and each column a
   class. Each matrix corresponds to one image
6  $model_\theta \leftarrow classifier$  // create a new classifier model with weights  $\theta$ , with logits as output
7  $EnD^2 = argmin_\theta \{ Loss_{EnD^2}(\phi, model_\theta(X)) \}$  // train model backpropagation
8 return EnD2

```

Algorithm 2: loss for EnD²

Input : Ensemble logits: ϕ , predicted logits: z , temperature: $T = T(t)$, annealing

Output : cost: C

```

1  $\epsilon = 10^{-8}$  // Smoothing factor
2  $\delta = 1 - 10^{-3}$  // Central smoothing factor
3  $\alpha = e^{z/T(t)}$  // elementwise exponential
4  $M = \#models$ 
5  $N = \#classes$ 
6 for  $i \leftarrow 1$  to  $M$  do
7   |  $\alpha_{0i} = \sum_j \alpha_{i,j}$  // sum over the classes to produce the precision factor
8 end
9  $P_{En} = softmax(\phi/T(t))$  // softmax over classes
10  $P_{En} = \delta(P_{En} - \frac{1}{N}) + \frac{1}{N}$ 
11  $TIT = \sum_i^N (log(\Gamma(\alpha_i + \epsilon)) - log(\Gamma(\alpha_0 + \epsilon)))$  // target independent term, where  $log(\Gamma(x)) = log((x-1)!)$ 
12  $A = \frac{1}{M} \sum_i^M (log(P_{En_i} + \epsilon))$  // mean over ensemble
13  $TDT = - \sum_i^N ((\alpha_i - 1)A_i)$  // target dependent term, sum over classes
14 return  $(TDT + TIT)T(t)^2$ 

```



Figure 1: The synthetic, spiral dataset.

B Experiments on Synthetic Data

B.1 Methodology

The goal with these experiments is to provide qualitative justification for [Claim 5](#) and illustrate the inner workings of EnD². We also provide some new experiments on temperature annealing and the size of the auxiliary dataset, to visualize their effect.

B.1.1 Dataset

To illustrate the model, Malinin et. al. use a synthetic dataset in \mathbb{R}^2 . Our rendering of this dataset can be seen in [Figure 1](#). This is advantageous since it enables plotting both knowledge and data uncertainty over the entire data manifold, giving a qualitative understanding of whether the algorithm works or not, in contrast to higher dimensional data (images, etc.) that cannot be plotted. The dataset itself looks like a spiral, divided into three classes shaped as spiralling arms of increasing radius. The spirals are centred and almost symmetric around the origin. Furthermore, they have increased noise and overlap with radius, which leads us to believe that uncertainty should vary as well. In addition to the spiral data an OOD data-set, referred to as the AUX data-set is also used, which takes the form of a ring slightly outside the spiral.

For the experiments, 1000 samples per ID class are used, both for training and test. The number of AUX samples was also 1000. This is the same setting as the original paper. The generation of the data uses the original paper’s code, but the hyperparameters were not specified. Our hyperparameters can be found in our code. We manually searched for hyperparameters, so that our plot would look as close to theirs, but the exact correspondence is probably not achieved.

B.1.2 Model description and hyperparameters

The original paper does not specify what type of neural network was used for classification. We were also unable to find it in the (unofficial) code. Instead, we chose to use a simple DNN with four hidden layers, each of width 64 with ReLu-activation functions, trained by minimizing the categorical cross-entropy using the Adam-optimizer, all with standard `tf.keras` settings, for 85 epochs. EnD and EnD² used the same base model but was instead trained for 500 epochs.

B.1.3 Experimental setup and code

On the output of an ensemble of 100 models, all differently randomly initialized, we train EnD and EnD² both with and without auxiliary data, using an initial temperature of 1, as in the paper. Doing this, we observed that the training diverges for many initialisations, mainly for EnD²_{AUX}. Thus, we also used an initial temperature of $T = 2.5$, both with and without annealing. The annealing schedule was $T = 2.5$ between epoch 0 and 200, linearly decreasing to 1 between

Table 1: Classification error on Spiral Dataset, compared with [1]. Error bars are 95%-confidence intervals assuming normal distribution. Note that our results likely use a different base model and training procedure than the original paper, since it was not specified there.

ERR↓	IND	ENSM	EnD	EnD ²	EnD _{+AUX}	EnD ² _{+AUX}	EnD ² _{+AUX, ANN}	EnD ² _{+AUX, T=2.5}	EnD ² _{+AUX20}
Our results	8.20±0.67	2.3±NA	3.90±0.65	3.86±0.70	2.61±0.11	4.67±3.26	3.30±0.59	3.45±0.96	5.0±1.54
Paper [1]	13.21	12.37	12.39	12.47	12.41	12.40	-	-	-

epoch 200 and 400 and 1 between epoch 400 and 500. Additionally, we also trained a model EnD_{+AUX20}^2 with only 20 samples from the auxiliary dataset.

All 7 models were trained 20 times, with different random initialisations. To make sure they converged, the test error was calculated. In cases test error was above 10%, it was deemed as non-convergence, and not taken into account. Among the converged ones, the mean error and the 95%-confidence interval around the mean is calculated, assuming a normal distribution. This means that for cases with fewer samples, the confidence interval is larger.

The main goal of this experiment is to visually show the total uncertainty, the data uncertainty and the knowledge uncertainty. They were calculated as specified in [1] and [2], for the grid $[-2000, 2000] \times [-2000, 2000]$ at all coordinates divisible by four, for a total of 10^6 points.

The full code is available at <https://anonymous.4open.science/r/4ee2c9ef-295f-44e2-8214-f0818b932817/>.

B.1.4 Computational requirements

The experiments were run on the CPU of a normal laptop (2.7 GHz Dual-Core i5). The total time to reproduce the ensemble of 100 models and all 20 repetitions of all 7 tested distillation methods, is around 5 to 6 hours.

B.2 Results

B.2.1 Classification accuracy

In Table 1, the classification accuracy from our experiment and the original paper is reported. We see that

- the ensemble outperforms the individual models, and that all distillation methods perform closer to the ensemble, than an individual model.
- the best performance is achieved by EnD with auxiliary data.
- using annealing or not when starting at $T = 2.5$ does not affect the final classification accuracy.

B.2.2 Visualization of uncertainty

The total, data and knowledge uncertainty is plotted in Figure 2 for a grid of 10^6 points. In contrast to the original paper, we fix the scale of the colour bar for better comparability between plots.

We observe that

- EnD^2 is not able to emulate the uncertainty landscape of the ensemble, but EnD_{+AUX}^2 can approximate it fairly well.
- Starting at a higher temperature ($T = 2.5$) and using annealing produces similar results as starting at temperature 1, but starting at temperature 2.5 and keeping it there for the entire training duration does not capture the true uncertainty.
- Using a smaller auxiliary dataset gives a worse approximation of the ensemble’s uncertainty landscape.

C Computational requirements for reproduction

In this section, we report the computational resources used for this reproduction. The running time of the major experiments on CIFAR-10 is expressed in time on an RTX 2070. For easier comparison, we also report the equivalent cost when running on a V100 GPU on Google Cloud for \$2.48 per hour, given a relative performance of 2.89 versus an

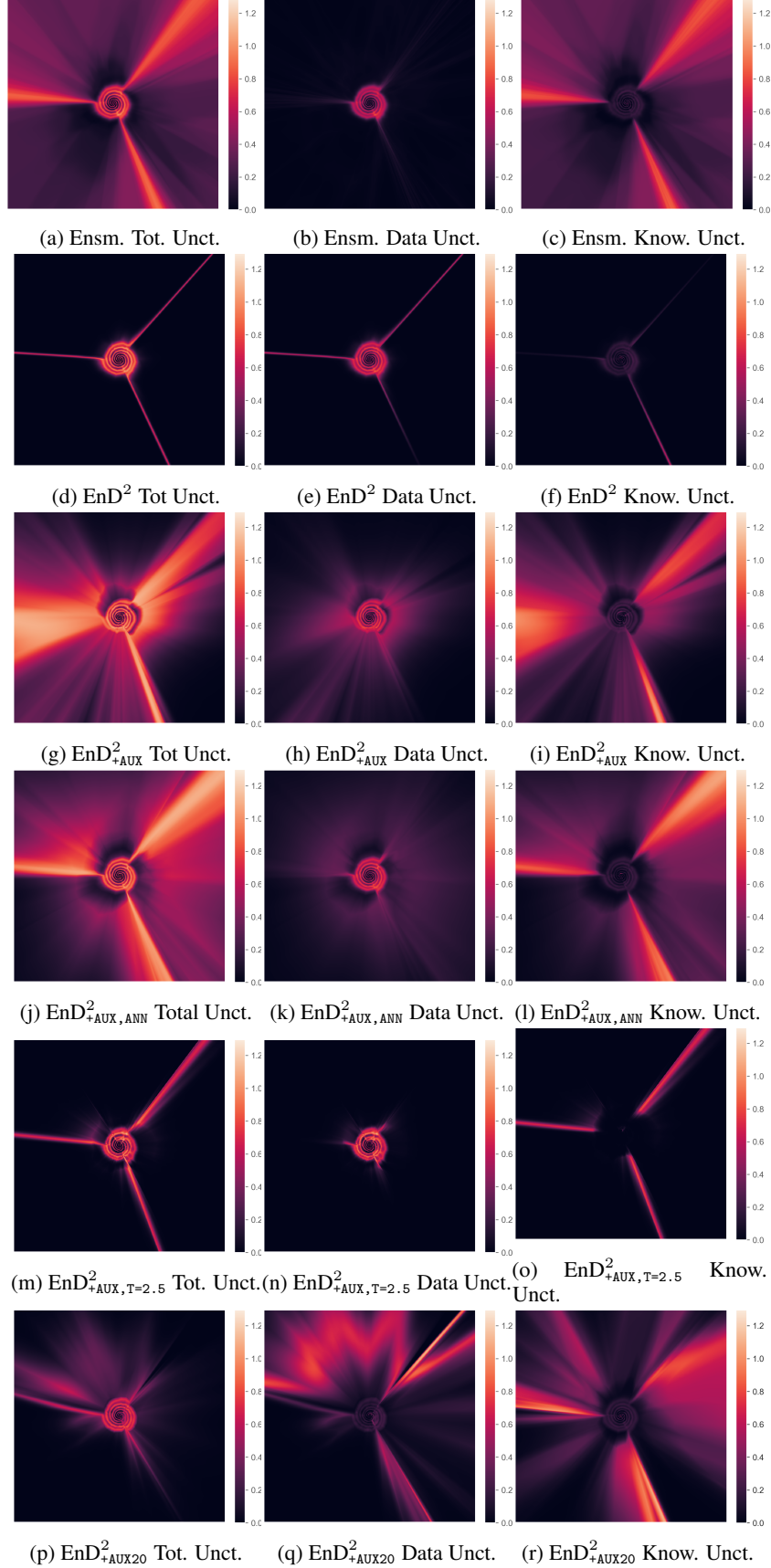


Figure 2: Recreation of Figure 3 in [1], showing uncertainties over entire data manifold.

Table 2: Computation requirements for major experiments, and which claims they test. GPU time refers to time on an NVIDIA GeForce RTX 2070. Equivalent cost represents the cost if run on a V100 on Google cloud, for \$2.48 per hour.

Experiment	Models	GPU minutes per model	GPU days	Equivalent cost in \$
Ensemble, training	400	16	4.44	91.53
Ensemble, labeling	400	0.45	0.13	2.57
Ensemble, inference	400	0.23	0.06	1.32
Evaluation, claim 1 and 2	15	51	0.53	10.94
Size ablation, training, claim 3	112	51	3.97	81.69
Temperature ablation, training, claim 4	27	51	0.96	19.69
3-class ensemble, training, claim 5	100	5.25	0.36	7.51
Total			11.413	235.06

Table 3: OOD ROC-AUC \uparrow on CIFAR-10 (in) and LSUN (out), normalized to ensemble results. Error bounds signify two standard deviations, taken over three models.

Unc.	IND	ENSM	EnD	EnD ²	EnD _{+AUX}	EnD _{+AUX} ²	PN _{+AUX}
Tot. our	0.96 \pm 0.00	1.00 \pm NA	1.00 \pm 0.01	0.98 \pm 0.00	1.01 \pm 0.00	1.00 \pm 0.00	1.02 \pm 0.01
Tot. paper	0.97 \pm 0.01	1.00 \pm NA	0.94 \pm 0.01	0.97 \pm 0.01	0.94 \pm 0.01	1.00 \pm 0.01	1.01 \pm 0.01
Know., our	-	1.00 \pm NA	-	0.95 \pm 0.01	-	0.99 \pm 0.01	1.02 \pm 0.00
Know., paper	-	1.00 \pm NA	-	0.98 \pm 0.01	-	0.99 \pm 0.01	1.01 \pm 0.01

RTX 2070¹. Note that these figures represent the time to reproduce only the final experiments. We estimate that the total GPU time used for this reproduction, including experimentation and bug-hunting, to be 3 to 5 times as long. The full data can be seen in Table 2.

D Histograms

To compare ensembles, EnD² and EnD_{+AUX}² on the CIFAR-10 and 3-class CIFAR-10 datasets, we provide histograms of data and knowledge uncertainty for in- and out-of-domain-distribution, in Figure 3 and 4.

E Relative performance of EnD² compared to ensemble and original article

In Tables 3 and 4 of the main report we report several measures for the 7 different models tested. For better comparability, we here also provide the values normalized to the ensembles' performance, both for our experiments, and for the original paper, in Table 4 and 3.

¹Benchmark taken from <https://timdettmers.com/2020/09/07/which-gpu-for-deep-learning/>

Table 4: Classification metrics on CIFAR-10, normalized to ensemble results. Error bounds signify two standard deviations, taken over three models.

Crit.	IND	ENSM	EnD	EnD ²	EnD _{+AUX}	EnD _{+AUX} ²	PN _{+AUX}
ERR \downarrow , our	1.12 \pm 0.08	1.00 \pm NA	0.99 \pm 0.06	1.13 \pm 0.02	1.13 \pm 0.02	1.16 \pm 0.01	1.14 \pm 0.04
ERR \downarrow , paper	1.29 \pm 0.06	1.00 \pm NA	1.08 \pm 0.05	1.18 \pm 0.03	1.08 \pm 0.03	1.11 \pm 0.06	1.21 \pm 0.10
PRR \uparrow , our	0.87 \pm 0.02	1.00 \pm NA	0.98 \pm 0.00	0.96 \pm 0.01	0.98 \pm 0.02	0.96 \pm 0.01	0.70 \pm 0.12
PRR \uparrow , paper	0.97 \pm 0.01	1.00 \pm NA	0.98 \pm 0.01	0.98 \pm 0.01	0.98 \pm 0.00	0.99 \pm 0.00	0.94 \pm 0.02
ECE \downarrow , our	41.37 \pm 0.35	1.00 \pm NA	0.94 \pm 0.05	1.45 \pm 0.13	1.08 \pm 0.19	1.85 \pm 0.29	5.69 \pm 0.37
ECE \downarrow , paper	1.69 \pm 0.31	1.00 \pm NA	2.00 \pm 0.15	0.77 \pm 0.15	2.00 \pm 0.46	1.69 \pm 0.31	9.23 \pm 0.54
NLL \downarrow , our	6.38 \pm 0.04	1.00 \pm NA	1.06 \pm 0.04	1.35 \pm 0.02	1.19 \pm 0.01	1.38 \pm 0.01	1.86 \pm 0.04
NLL \downarrow , paper	1.32 \pm 0.05	1.00 \pm NA	1.16 \pm 0.05	1.32 \pm 0.05	1.16 \pm 0.05	1.26 \pm 0.00	2.00 \pm 0.05

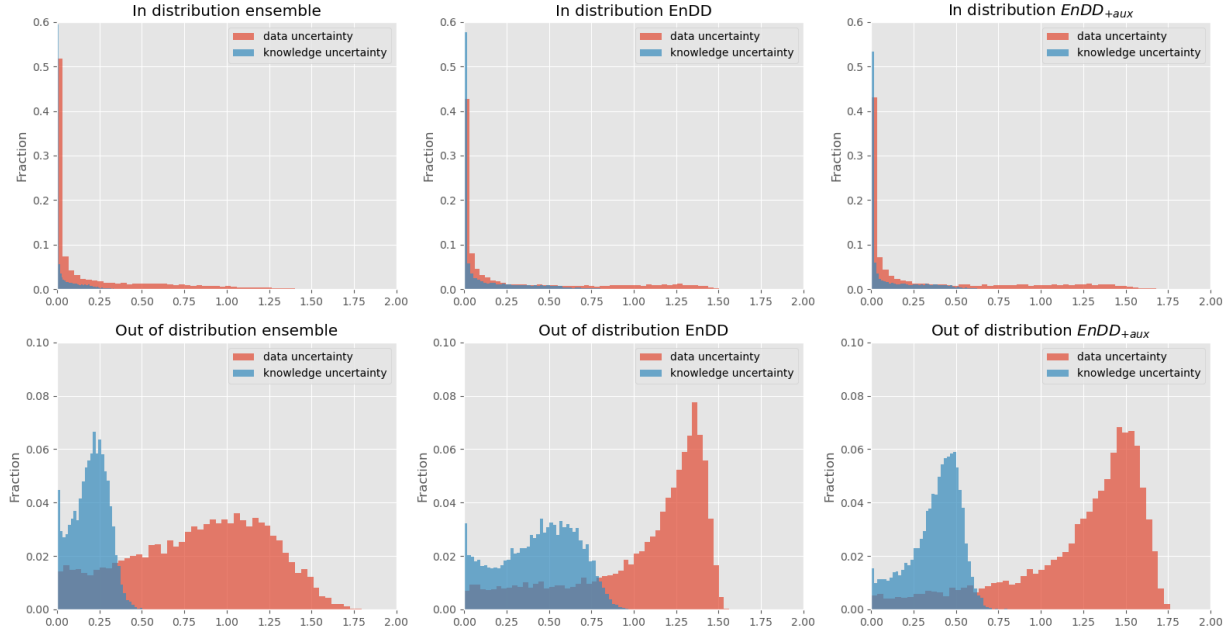


Figure 3: Data/knowledge uncertainty-distributions for ensemble, EnD^2 and $\text{EnD}^2_{\text{AUX}}$.

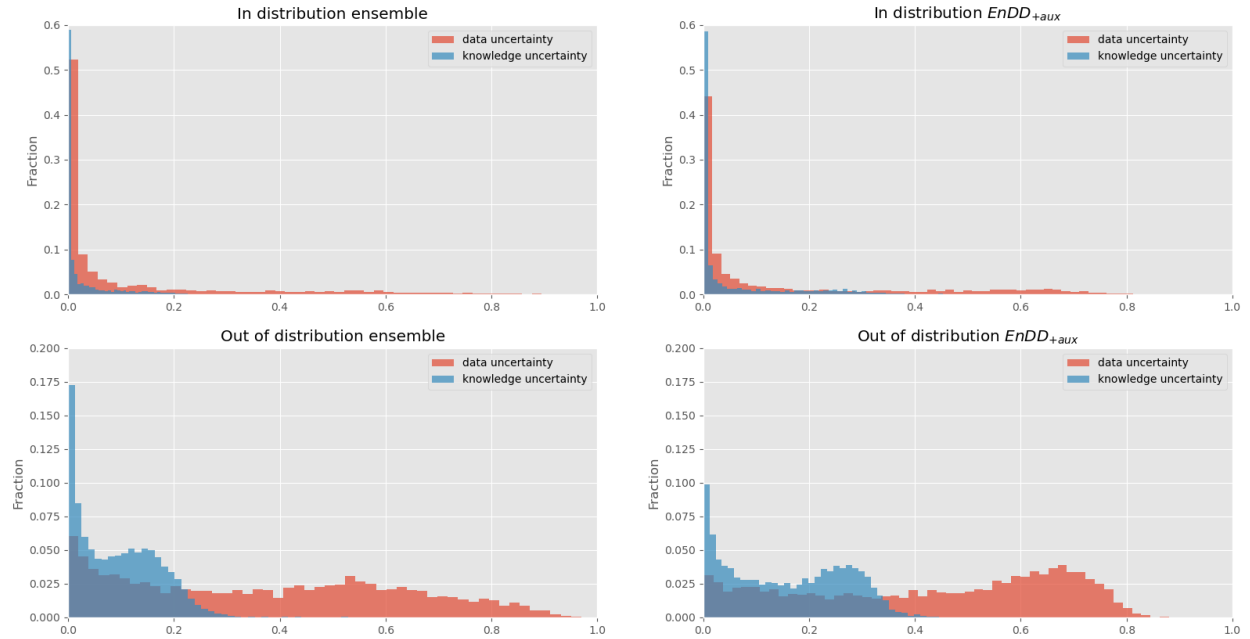


Figure 4: Data/knowledge uncertainty-distributions for ensemble and $\text{EnD}^2_{\text{AUX}}$ on the 3-class CIFAR10 dataset

References

- [1] Andrey Malinin, Bruno Mlodozienec, and Mark Gales. “Ensemble Distribution Distillation”. In: *International Conference on Learning Representations (ICLR)*. 2020.
- [2] Andrey Malinin and Mark Gales. “Predictive Uncertainty Estimation via Prior Networks”. In: *Advances in Neural Information Processing Systems 31*. 2018.