

448 **A Appendix**

449 **A.1 Extended Derivation for Equation 4**

450 In the main paper, we proposed to decompose the norm and angular similarity into instance-
451 independent and dependent components.

$$\begin{aligned}\|\mathbf{x}\|_2 &= \|\Delta x\|_2 + \mathcal{C}_x \\ |\phi_y| &= |\Delta\phi_y| - |\mathcal{C}_\phi|\end{aligned}$$

452 The $\|\Delta\mathbf{x}\|_2, |\Delta\phi_y|$ are the instance-dependent components and $\mathcal{C}_x, |\mathcal{C}_\phi|$ are the instance-independent
453 components. We can rewrite the pre-softmax logits in Eq. 1 with the decomposed norm and angular
454 similarity.

$$\begin{aligned}\|\mathbf{x}\|_2 \cos \phi_y &= \|\mathbf{x}\|_2 \cos |\phi_y| = (\|\Delta\mathbf{x}\|_2 + \mathcal{C}_x) \cos (|\Delta\phi_y| - |\mathcal{C}_\phi|) \\ &= (\|\Delta\mathbf{x}\|_2 + \mathcal{C}_x) (\cos |\Delta\phi_y| \cos |\mathcal{C}_\phi| + \sin |\Delta\phi_y| \sin |\mathcal{C}_\phi|) \\ &= (\|\Delta\mathbf{x}\|_2 + \mathcal{C}_x) \frac{1}{\cos |\mathcal{C}_\phi|} \left(\cos |\Delta\phi_y| \cos |\mathcal{C}_\phi|^2 + \sin |\Delta\phi_y| \cos |\mathcal{C}_\phi| \sin |\mathcal{C}_\phi| \right) \\ &= (\|\Delta\mathbf{x}\|_2 + \mathcal{C}_x) \frac{1}{\cos |\mathcal{C}_\phi|} \cos |\Delta\phi_y| \left(\cos |\mathcal{C}_\phi|^2 + \cos |\mathcal{C}_\phi| \sin |\mathcal{C}_\phi| \frac{\sin |\Delta\phi_y|}{\cos |\Delta\phi_y|} \right) \\ &= (\|\Delta\mathbf{x}\|_2 + \mathcal{C}_x) \frac{1}{\cos |\mathcal{C}_\phi|} \cos |\Delta\phi_y| \left(\left(1 - \sin |\mathcal{C}_\phi|^2\right) + \cos |\mathcal{C}_\phi| \sin |\mathcal{C}_\phi| \frac{\sin |\Delta\phi_y|}{\cos |\Delta\phi_y|} \right) \\ &= (\|\Delta\mathbf{x}\|_2 + \mathcal{C}_x) \frac{1}{\cos |\mathcal{C}_\phi|} \cos |\Delta\phi_y| \left(1 - \sin |\mathcal{C}_\phi|^2 \left(1 - \frac{\cos |\mathcal{C}_\phi| \sin |\Delta\phi_y|}{\sin |\mathcal{C}_\phi| \cos |\Delta\phi_y|} \right) \right)\end{aligned}$$

455 **A.2 Definitions of Metrics**

456 The problem tackled in this paper is supervised image classification in the face of noise. Assume a data
457 point $X_i \in \mathbf{X}, i \in [1, N]$ each associated with a label $Y \in \mathbf{Y} = \{1, \dots, K\}$. We would like our model
458 M where $M(X_i) = (\hat{Y}_i, \hat{P}_i)$ where \hat{Y}_i is the class prediction and \hat{P}_i is the probability/confidence given
459 by the model to the ground truth distribution $P(Y_i|X_i)$. Ideally \hat{P}_i is well calibrated
460 which means that it represents the likelihood of the true event $\hat{Y}_i = Y_i$. *Perfect calibration* [4] can be
461 defined as:

$$P(\hat{Y}_i = Y_i | \hat{P}_i = P_i) = P_i, \forall P_i \in [0, 1] \quad (11)$$

462 Ways of evaluating Calibration are as follows:

463 **A.2.1 Expected Calibration Error (ECE)**

464 Expected Calibration Error [20] evaluates calibration by calculating the difference in expectation
465 between the confidence and accuracy or:

$$E_{\hat{P}}[|P(\hat{Y} = Y | \hat{P} = p) - p|] \quad (12)$$

466 This can also be computed as the weighted average of bins' accuracy/confidence difference:

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n} |\text{accuracy}(B_m) - \text{confidence}(B_m)| \quad (13)$$

467 where n is the total number of samples. Perfect calibration is achieved bins when confidence equals
468 accuracy and $\text{ECE} = 0$.

Table 6: **ResNet18 ECE on CIFAR10/100 Noise**, averaged over 5 seeds

Noise-level	CIFAR10					CIFAR100				
	1	2	3	4	5	1	2	3	4	5
Ensemble	0.051	0.075	0.076	0.118	0.184	0.059	0.076	0.078	0.107	0.149
MCDO	0.076±0.003	0.098±0.004	0.102±0.002	0.164±0.005	0.251±0.008	0.114±0.002	0.147±0.002	0.14±0.006	0.192±0.009	0.255±0.014
ResNet	0.102±0.001	0.141±0.003	0.153±0.007	0.209±0.011	0.293±0.016	0.113±0.004	0.149±0.005	0.152±0.004	0.185±0.005	0.237±0.01
Ours (GS)	0.040±0.002	0.055±0.003	0.060±0.005	0.080±0.01	0.106±0.012	0.067±0.002	0.083±0.002	0.089±0.004	0.116±0.007	0.145±0.013

Table 7: **ResNet18 NLL on CIFAR10/100 Noise**, averaged over 5 seeds

Noise-level	CIFAR10					CIFAR100				
	1	2	3	4	5	1	2	3	4	5
Ensemble	0.544	0.737	0.753	1.055	1.551	1.499	1.927	1.969	2.379	2.99
MCDO	0.667±0.02	0.845±0.029	0.831±0.013	1.262±0.033	1.947±0.054	1.789±0.011	2.225±0.012	2.236±0.043	2.788±0.072	3.585±0.119
ResNet	0.718±0.01	0.979±0.019	1.043±0.046	1.436±0.081	2.052±0.133	1.782±0.018	2.269±0.023	2.326 ± 0.029	2.773±0.027	3.434±0.032
Ours (GS)	0.531±0.006	0.716±0.012	0.785±0.013	1.007±0.018	1.346±0.02	1.786±0.013	2.208±0.013	2.300±0.014	2.676±0.015	3.215±0.028

469 **A.2.2 Negative Log Likelihood (NLL)**

470 A way to measure a model’s probabilistic quality is to use Negative Log Likelihood [18]. Given a
 471 probabilistic model $P(Y|X)$ and N samples it is defined as:

$$\mathbf{L} = - \sum_{i=1}^N \log(\hat{P}(Y_i|X_i)) \quad (14)$$

472 where \hat{P} is the predicted distribution of the ground truth P and Y_i is the true label for input X_i . NLL
 473 belongs to a class of strictly proper scoring rules [30]. A scoring rule is strictly proper if it is uniquely
 474 optimized by only the true distribution. NLL is the negative of the logarithm of the probability of the
 475 true outcome. If the true class is assigned a probability of 1, NLL will be minimum with value 0.

476 **A.2.3 Brier**

477 The Brier score [19] measures accuracy of probabilistic predictions. Across all predicted items N
 478 in a set of predictions, the Brier score measures the mean squared difference between the predicted
 479 probability assigned to possible outcome for $i \in [1, N]$ and the actual outcome.

$$\mathbf{BS} = (1/N) \sum_{t=1}^N \sum_{i=1}^R (f_{ti} - o_{ti})^2 \quad (15)$$

480 Where R is number of possible classes, N is overall number of instances of all classes. f_{ti} is the
 481 approximated probability of the forecast o_{ti} in one hot encoding. Brier score can be intuitively
 482 decomposed into three components: uncertainty, reliability and resolution [31] and it is also a proper
 483 scoring rule.

484 **A.3 Calibration in the Face of Differing Levels of Noise**

485 We report additional calibration ECE, NLL and Brier results in the face of different levels of
 486 corruption using ResNet18 in Tab. 6, 7 and 8 respectively. CIFAR10 and CIFAR100’s validation set
 487 was corrupted using a library of common corruptions [1] with 5 levels of severity as can be seen in
 488 Fig. 4a, 4b.

Table 8: **ResNet18 Brier on CIFAR10/100 Noise**, averaged over 5 seeds

Noise-level	CIFAR10					CIFAR100				
	1	2	3	4	5	1	2	3	4	5
Ensemble	0.021	0.028	0.03	0.041	0.057	0.005	0.006	0.006	0.007	0.008
MCDO	0.024±0.0	0.03±0.001	0.032±0.0	0.046±0.001	0.065±0.001	0.005±0.0	0.006±0.0	0.006±0.0	0.007±0.0	0.009±0.0
ResNet	0.025±0.0	0.034±0.0	0.038±0.001	0.05±0.002	0.068±0.002	0.005±0.0	0.006±0.0	0.007±0.0	0.007±0.0	0.009±0.0
Ours (GS)	0.022±0.0	0.03±0.0	0.034±0.0	0.043±0.001	0.056±0.001	0.003±0.0	0.005±0.0	0.006±0.0	0.007±0.0	0.008±0.000

	ECE↓	NLL↓	Brier↓	Accuracy↑		ECE↓	NL↓L	Brier↓	Accuracy↑
Ensemble	0.12±0.047	2.973±0.833	0.008±0.002	0.338	Ensemble	0.302	2.397	0.082	0.44
MCDO	0.254±0.005	3.78±0.043	0.009±0.0	0.282±0.002	MCDO	0.373±0.001	3.08±0.025	0.092±0.0	0.401±0.002
ResNet18	0.215±0.007	3.352±0.036	0.009±0.0	0.311±0.003	ResNet18	0.42±0.009	2.941±0.085	0.095±0.001	0.427±0.006
Ours	0.097±0.003	<u>3.189±0.019</u>	0.008±0.0	0.299±0.003	Ours	<u>0.323±0.006</u>	2.211±0.04	<u>0.085±0.001</u>	0.422±0.005

(a) ResNet18 on CIFAR100 Rotate over 5 seeds (b) ResNet18 on CIFAR10 Rotate over 5 seeds

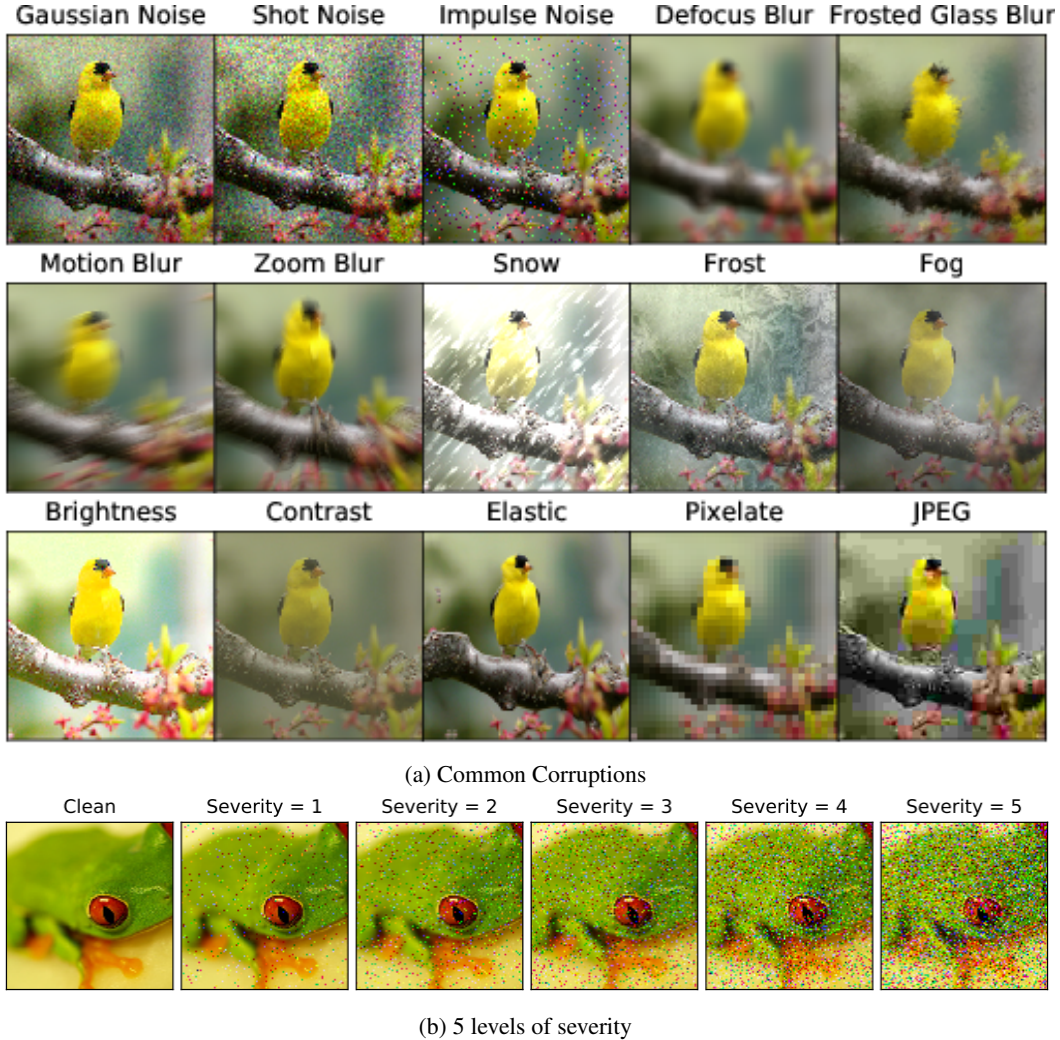


Figure 4: Examples of common corruptions and the 5 levels of severity from [1].

489 In Tab. 6, 7 and 8 we show how differing levels of common corruptions effect the calibration of
 490 models. Across all levels of corruption our model consistently had the stronger Brier score in
 491 CIFAR100 and much strong ECE and NLL on CIFAR10.

492 A.4 Calibration in the Face of Rotation

493 In Tab. 9b, 9a we rotated CIFAR10 and CIFAR100 validation data set by [0, 350] degrees with 10
 494 degree steps in between, the calibration metrics and accuracy were then averaged. For each model 5
 495 seeds were trained, for MCDO 5 passes were done on each model for inference with a dropout rate of
 496 50% as suggested in the original paper and 5 models were ensembled for Deep Ensemble. β' for our
 497 models were 4 on CIFAR10 and 10 for CIFAR100.

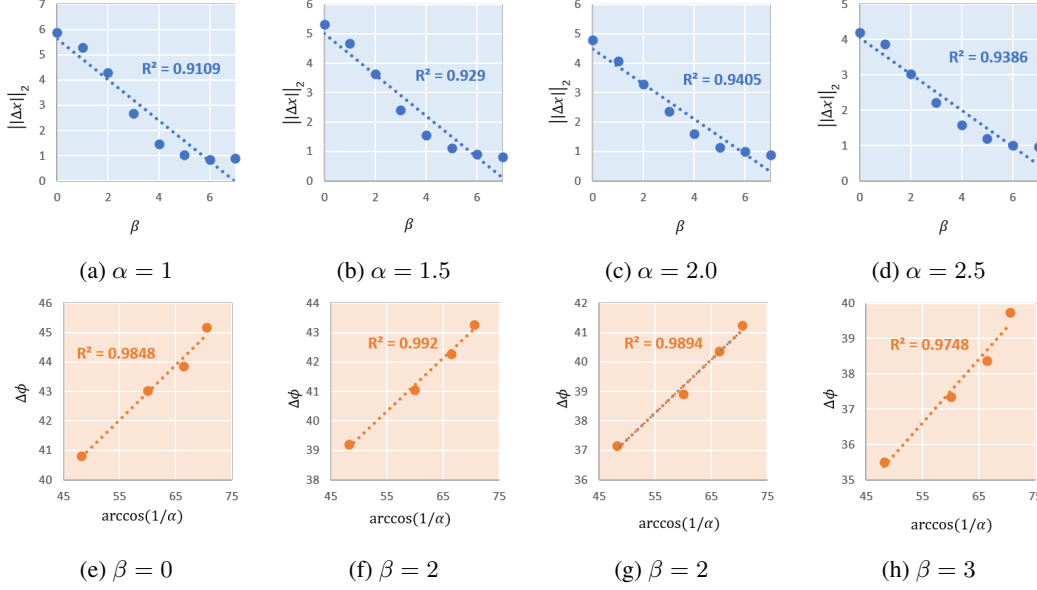


Figure 5: **Properties of $\|\Delta x\|_2$ and $\Delta\phi$.** (a) - (b): $\|\Delta x\|_2$ decreases linearly with β for fixed α reflecting Eq. 2 and 6. (e) - (h) $\Delta\phi$ increases linearly with $\arccos(1/\alpha)$ for fixed β reflecting Eq. 3 and 6. Note that all plots include R-squared values to indicate goodness-of-fit of the linear relationship.

Table 10: Model Requirements

Model	Loss Function	Hyperparameters	Output	Multi-pass Infer
Ensemble	CE	Number of Models	LL	True
MCDO	CE	Dropout %	LL	True
SNGP	CE	Spectral Norm Bound, GP scale & bias & discount factor & covariance factor & field factor & ridge penalty	GP	False
DUQ	Multi-BCE	Gradient penalty, RBF sigma, embedding gamma	RBF	False
Ours	CE	β^* , error (default = 0.1)	Decomposed LL	False

LL: Linear Layer. **CE:** Cross-Entropy, **BCE:** Binary Cross-Entropy, **GP:** Gaussian Process, **RBF:** Radial Basis Function

498 A.5 Empirical Support for Disentangled Training: Figures

499 We examine how the norm, and angle change with different $\alpha - \beta$ configurations. From Fig. 5a
500 - 5d, we observe that the norm decreases linearly with β for fixed α . From Fig. 5e - 5h, we
501 observe that the angle increases linearly with $\arccos(1/\alpha)$. The observations are consistent with the
502 original geometric motivation. β encodes an instance-independent portion, \mathcal{C}_x , of the norm. As β
503 increases the instance-independent portion increases and therefore the magnitude of the dependent
504 component, $\|\Delta x\|_2$ decreases linearly. α encodes the inverse of the cosine of a relaxation angle, \mathcal{C}_ϕ .
505 As $\arccos(1/\alpha)$ increases, the resulting angle, $\Delta\phi$ increases linearly due to the increased relaxation
506 angle encoded by α .

507 A.6 Qualitative Comparison: Extended Discussion

508 **GSD vs. Single Pass Models** The current state-of-the-art single pass models for inference on OOD
509 data, without training on OOD data, are SNGP [9] and DUQ [8]. The primary disadvantages
510 of these models is: **1) Hyperparameter Combinatorics:** Both DUQ and SNGP require many
511 hyperparameters as shown in Tab. 10. SNGP requires the most hyperparameters out of all the single
512 pass models. The large combinatoric scale, in addition to the fact that these hyperparameters must
513 be tuned via pre-training grid search, make these methods costly as a full training procedure with
514 multiple epochs are required before evaluating calibration. Our model only has *one* hyperparameter
515 that is tuned post-training with 1 epoch on validation set. **2) Extended Training Time:** DUQ
516 requires a centroid embedding update every epoch, while SNGP requires sampling potentially high
517 dimensional embeddings of training points for generating the covariance matrix as well as updates to

Table 11: **Extended Generalizability Experiments** We benchmark our method against the vanilla models using 12 different backbones and 4 different dOPTasets. **Grid Searched (GS):** β' grid searched on validOPTion ECE, **Optimized (OPT):** β' optimized via SGD on validOPTion NLL .

model	dataset	Clean				Corrupt/Rotate			
		accuracy \uparrow	ECE \downarrow	NLL \downarrow	Brier \downarrow	accuracy \uparrow	ECE \downarrow	NLL \downarrow	Brier \downarrow
LeNet5	Mnist	96.16%	0.01	0.132	0.006	33.95%	0.43	4.533	0.104
GSD LeNet5 GS	Mnist	96.86%	0.005	0.103	0.005	35.73%	0.42	4.405	0.101
GSD LeNet5 OPT	Mnist	96.86%	0.012	0.106	0.005	35.73%	0.406	4.173	0.01
DenseNet	SVHN	41.72%	0.051	1.71	0.072	14.31%	0.301	3.844	0.107
GSD DenseNet GS	SVHN	41.7%	0.027	1.62	0.069	14.41%	0.287	3.134	0.106
GSD DenseNet OPT	SVHN	41.7%	0.04	1.62	0.069	14.41%	0.277	3.25	0.105
ResNet	CIFAR10	95.9%	0.007	0.149	0.006	76.54%	0.178	1.28	0.603
GSD ResNet34 GS	CIFAR10	95.9%	0.005	0.148	0.006	76.54%	0.088	0.882	0.037
GSD ResNet34 OPT	CIFAR10	95.9%	0.011	0.162	0.007	76.54%	0.054	0.813	0.035
ResNet50	CIFAR10	95.32%	0.03	0.203	0.008	76.32%	0.17	1.23	0.039
GSD ResNet50 GS	CIFAR10	95.82%	0.008	0.147	0.007	76.23%	0.057	0.766	0.033
GSD ResNet50 OPT	CIFAR10	95.82%	0.01	0.158	0.007	76.32%	0.115	0.928	0.038
ResNet101	CIFAR10	95.61%	0.028	0.197	0.007	77.59%	0.154	1.118	0.037
GSD ResNet101 GS	CIFAR10	95.62%	0.007	0.158	0.007	77.21%	0.075	0.852	0.036
GSD ResNet101 OPT	CIFAR10	95.62%	0.007	0.155	0.007	77.21%	0.086	0.788	0.033
ResNet152	CIFAR10	95.7%	0.028	0.196	0.007	75.2%	0.179	1.337	0.041
GSD ResNet152 GS	CIFAR10	95.63%	0.007	0.151	0.007	76.58%	0.058	0.765	0.033
GSD Resnet152 OPT	CIFAR10	95.63%	0.01	0.154	0.007	76.58%	0.043	0.756	0.032
ResNet34	CIFAR100	78.81%	0.071	0.868	0.003	51.16%	0.19	2.387	0.007
GSD ResNet34 GS	CIFAR100	78.02%	0.037	0.938	0.003	49.27%	0.098	2.361	0.007
GSD ResNet34 OPT	CIFAR100	78.02%	0.043	0.93	0.003	49.27%	0.112	2.372	0.007
ResNet50	CIFAR100	79.28%	0.0746	0.861	0.003	49.71%	0.213	2.477	0.007
GSD ResNet50 GS	CIFAR100	78.97%	0.0326	0.879	0.003	50.12%	0.08	2.264	0.006
GSD ResNet50 OPT	CIFAR100	78.97%	0.041	0.856	0.003	50.12%	0.110	2.28	0.007
ResNet101	CIFAR100	79.21%	0.725	2.98	0.009	51.34%	0.470	3.62	0.009
GSD ResNet101 GS	CIFAR100	79.82%	0.034	0.834	0.003	53.14%	0.082	2.11	0.006
GSD ResNet101 OPT	CIFAR100	79.82%	0.038	0.829	0.003	53.14%	0.092	2.114	0.006
ResNet152	CIFAR100	80.71%	0.0895	0.815	0.003	54.2%	0.233	2.45	0.007
GSD ResNet152 GS	CIFAR100	79.85%	0.0364	0.827	0.003	53%	0.078	2.12	0.006
GSD ResNet152 OPT	CIFAR100	79.85%	0.0397	0.821	0.003	53%	0.087	2.12	0.006

518 the bounded spectral norm on each training step, thus increasing training time while our model trains
519 in the same amount of time as the model it is applied to.

520 **GSD vs. Multi-Pass Models** Bayesian MCDO [7] and Deep Ensemble [14] are considered the
521 current state-of-the-art methods for multi-pass calibration. Bayesian MCDO requires multiple passes
522 with dropout during training and inference in order to achieve stronger calibration. Deep Ensembles
523 requires N times the number of parameters as the single model it is ensembling where N is the
524 number of models ensembled. The obvious disadvantage to Deep Ensembles is that it requires N
525 times as long to train and run inference as its base model. While no model currently beats Deep
526 Ensemble in accuracy on both clean data and corrupted data, we have shown that our model has
527 stronger calibration in the face of certain levels of severity of corruption Tab. 1 and 2. Bayesian
528 MCDO has shown to have stronger calibration than the same model not trained with dropout, but
529 tends to suffer large accuracy drops as well as not being as strong as other single pass models or Deep
530 Ensemble in calibration, even with many passes. Our model empirically suffers minimal accuracy
531 drops when compared to its backbone and in some conditions led to stronger accuracy on corrupted
532 data (Tab. 1 and 2).

533 A.7 Generalizability: Extended Table

534 **Generalizability** We explored how generalizable our method is by applying it to 12 different models
535 and 4 different datasets in Tab. 11. We report results for both variants of our model: **Grid Searched:**
536 grid search β' on the validation set to minimize ECE and **Optimized:** optimize β' on the validation
537 set via gradient decent to minimize NLL for 10 epochs, similar to temperature scaling. We can

538 see consistently that our model had stronger calibration across all models and metrics, including
 539 models known to be well calibrated like LeNet [22]. All models were tested on CIFAR10C and
 540 CIFAR100C datasets offered by [1] where the original CIFAR10 and CIFAR100 were pre-corrupted;
 541 these were used for consistent corruption benchmarking across all models. All non-CIFAR datasets
 542 were corrupted via rotation from angles [0,350] with 10 step angles in between and the average
 543 calibration and accuracy was taken across all degrees of rotation. Our models included: DenseNet [23],
 544 LeNet [22] and 6 varying sizes of ResNet, which are described in [24]. The datasets we experimented
 545 on CIFAR10 [25], CIFAR100 [26], MNIST [27] and SVHN [28], CIFAR10C [1], CIFAR100C [1].

546 A.8 Training Parameters and Dataset License

547 We train all our models using stochastic gradient descent for 200 epochs and a batch size of 128 on
 548 RTX 2080 GPUs. We use a starting learning rate of 0.1 and a weight decay of $5.0e - 4$. For ResNet18
 549 experiments, we use a cosine scheduler for learning rate. For Wide ResNet-20-10 experiments, we
 550 use a step scheduler which multiplies the learning rate at epoch 60, 120 and 160 by 0.2.

551 The CIFAR10/100 datasets [25, 26] are released under MIT license. The CIFAR10/100C datasets [1]
 552 are released under Creative Commons Public license.

553 A.9 Introduction to Temperature Scaling

554 Temperature scaling is a simple form of Platt scaling [32]. Temperature scaling uses a scalar T to
 555 adjust the confidence of the softmax probability in a classification model. Following the notation from
 556 the main paper, let l denotes the logits. The temperature scalar is applied to all classes as following:

$$P(y|x) = \frac{\exp \frac{1}{T} l_y}{\sum_{j=1}^c \exp \frac{1}{T} l_j} = \frac{\exp (\|\mathbf{w}_y\|_2 \frac{1}{T} \|\mathbf{x}\|_2 \cos \phi_y)}{\sum_{j=1}^c \exp (\|\mathbf{w}_j\|_2 \frac{1}{T} \|\mathbf{x}\|_2 \cos \phi_j)} \quad (16)$$

557 As described in Fig. 1a, the temperature effectively changes the slope of $\|\mathbf{x}\|_2$ from 1 to $\frac{1}{T}$. The
 558 temperature parameter is optimized by minimizing negative log likelihood on a validation set while
 559 freezing all the other model parameters [4]. Temperature scaling calibrates a model’s confidence on
 560 IND data and does not change accuracy. However, it does not provide any mechanism to improve
 561 calibration on shifted distribution and is inferior to other uncertainty estimation methods in terms of
 562 calibration [5].