

A APPENDIX

- **PseudoCode** (Sec. A.1)
- **Ethics** (Sec. A.2)
- **Reproducibility** (Sec. A.3)
- **Details and Expanded Results for Super-pixel Graph Experiments**(Sec. A.4)
- **Stochastic Centering on the Empirical NTK of Graph Neural Networks** (Sec. A.5)
- **Size-Generalization Dataset Statistics** (Sec. A.6)
- **GOOD Dataset Statistics and Expanded Results** (Sec. A.7)
- **Alternative Anchoring Strategies on GOOD Datasets** (Sec. A.8)
- **Discussion of Post-hoc Calibration Strategies** (Sec. A.9)
- **Details of Generalization Gap Experiments** (Sec. A.10)
- **Expanded Pretraining Results** (Sec. A.12)
- **Runtimes** (Sec. A.13)
- **Mean and Variance of Node Feature Anchoring Distributions** (Sec. A.14)
- **Discussion on Anchoring Design Choices** (Sec. A.15)

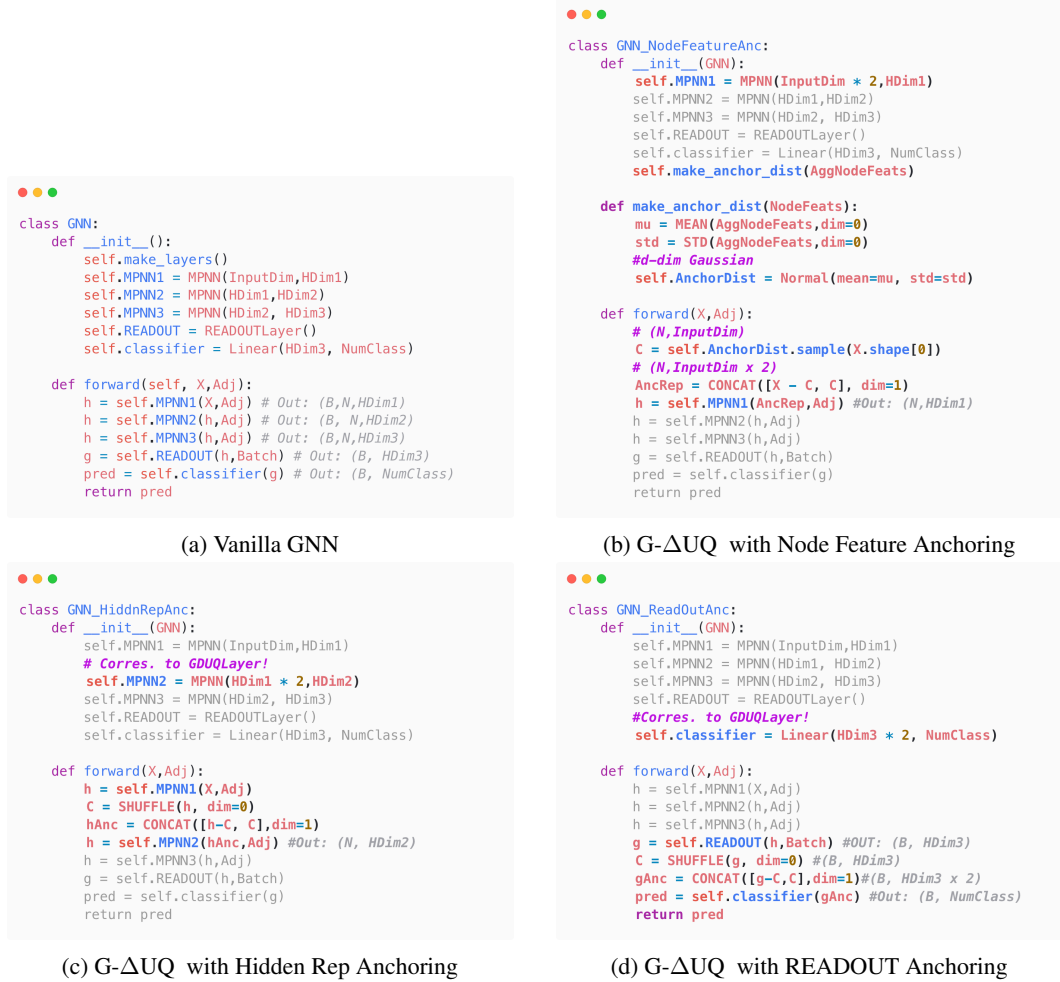
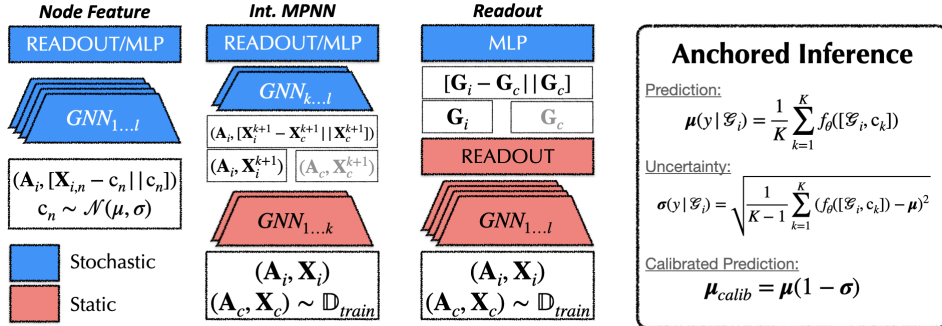
A.1 PSEUDOCODE FOR G- Δ UQ

Figure 5: **PseudoCode for G- Δ UQ**. We provide simplified pseudo-code to demonstrate how anchoring can be performed. We assume PyTorchGeometric style mini-batching. Changes with respect to the vanilla GNN are shown in bold. Unchanged lines are grayed out.



Overview of G- Δ UQ models. Here, we present a conceptual overview of how G- Δ UQ induces partially stochastic models. This figure is complementary to 5.

A.2 ETHICS STATEMENT

This work proposes a method to improve uncertainty estimation in graph neural networks, which has potential broader societal impacts. As graph learning models are increasingly deployed in real-world applications like healthcare, finance, and transportation, it becomes crucial to ensure these models make reliable predictions and know when they may be wrong. Unreliable models can lead to harmful outcomes if deployed carelessly. By improving uncertainty quantification, our work contributes towards trustworthy graph AI systems.

We also consider several additional safety-critical tasks, including generalization gap prediction for graph classification (to the best of our knowledge, we are the first to report results on this task) and OOD detection. We hope our work will encourage further study in these important areas.

However, there are some limitations. Our method requires (modest) additional computation during training and inference, which increases resource usage. Although $G-\Delta UQ$, unlike post-hoc methods, does not need to be fit on a validation dataset, evaluation of its benefits also relies on having some out-of-distribution or shifted data available, which may not always be feasible. We have seen in Table 1 that there are tasks for which $G-\Delta UQ$ fails to improve accuracy and/or calibration of some post-hoc methods, further emphasizing the need to perform appropriate model selection and the risks if shifted validation data is not available. Finally, there are open questions around how much enhancement in uncertainty calibration translates to real-world safety and performance gains.

Looking ahead, we believe improving uncertainty estimates is an important direction for graph neural networks and deep learning more broadly. This will enable the development safe, reliable AI that benefits society. We hope our work inspires more research in the graph domain that focuses on uncertainty quantification and techniques that provide guarantees about model behavior, especially for safety-critical applications. Continued progress will require interdisciplinary collaboration between graph machine learning researchers and domain experts in areas where models are deployed.

A.3 REPRODUCIBILITY

For reproducing our experiments, we have made our code available at this anonymous [repository](#). In the remainder of this appendix (specifically App. A.6, A.7), and A.10), we also provide additional details about the benchmarks and experimental setup.

A.4 DETAILS ON SUPER-PIXEL EXPERIMENTS

We provide an example of the rotated images and corresponding super-pixel graphs in Fig. 6. (Note that classes “6” and “9” may be confused under severe distribution shift, i.e. 90 degrees rotation or more. Hence, to avoid harming class information, our experiments only consider distribution shift from rotation up to 40 degrees.)

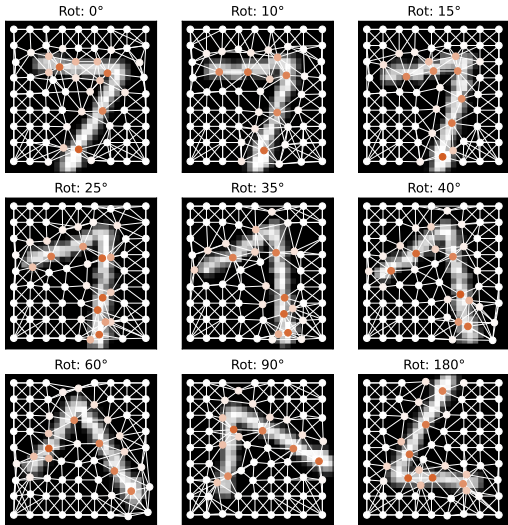


Figure 6: **Rotated Super-pixel MNIST.** Rotating images prior to creating super-pixels to leads to some structural distortion (Ding et al., 2021). However, we can see that the class-discriminative information is preserved, despite rotation. This allows for simulating different levels of graph structure distribution shifts, while still ensuring that samples are valid.

Tables 4 and 5 provided expanded results on the rotated image super-pixel graph classification task, discussed in Sec. 6.1.

In Table 7 we focus on the calibration results on this task for GPS variants alone. Across all levels of distribution shift, the best method is our strategy for applying $G-\Delta UQ$ to a pretrained model—demonstrating that this is not just a practical choice when it is infeasible to retrain a model, but can lead to powerful performance by any measure. Second-best on all datasets is applying $G-\Delta UQ$ during training, further highlighting the benefits of stochastic anchoring.

In addition to the structural distribution shifts we get by rotating the images before constructing super-pixel graphs, we also simulate feature distribution shifts by adding Gaussian noise with different standard deviations to the pixel value node features in the super-pixel graphs. In Table 8, we report accuracy and calibration results for varying levels of distribution shift (represented by the size of the standard deviation of the Gaussian noise). Across different levels of feature distribution shift, we also see that $G-\Delta UQ$ results in superior calibration, while maintaining competitive or in many cases superior accuracy.

Table 4: **RotMNIST-Accuracy.** Here, we report expanded results (accuracy) on the Rotated MNIST dataset, including a variant that combines G- Δ UQ with Deep Ens. Notably, we see that anchored ensembles outperform basic ensembles in both accuracy and calibration. (Best results for models using Deep Ens. and those not using it marked separately.)

MODEL	G- Δ UQ?	LPE?	Avg. Test (\uparrow)	Acc. (10) (\uparrow)	Acc. (15) (\uparrow)	Acc. (25) (\uparrow)	Acc. (35) (\uparrow)	Acc. (40) (\uparrow)
GatedGCN	×	×	0.947 \pm 0.002	0.918 \pm 0.002	0.904 \pm 0.005	0.828 \pm 0.009	0.738 \pm 0.009	0.679 \pm 0.007
	✓	×	0.933 \pm 0.015	0.894 \pm 0.019	0.878 \pm 0.020	0.794 \pm 0.032	0.698 \pm 0.036	0.636 \pm 0.048
	×	✓	0.949 \pm 0.002	0.917 \pm 0.004	0.904 \pm 0.005	0.829 \pm 0.007	0.744 \pm 0.007	<u>0.685 \pm 0.006</u>
	✓	✓	0.915 \pm 0.032	0.872 \pm 0.038	0.852 \pm 0.0414	0.776 \pm 0.039	0.680 \pm 0.037	0.631 \pm 0.033
GPS	×	✓	0.970 \pm 0.001	0.948 \pm 0.001	0.938 \pm 0.001	0.873 \pm 0.006	0.770 \pm 0.013	0.688 \pm 0.009
	✓	✓	<u>0.969 \pm 0.001</u>	<u>0.946 \pm 0.003</u>	<u>0.937 \pm 0.003</u>	<u>0.869 \pm 0.003</u>	<u>0.769 \pm 0.012</u>	0.679 \pm 0.014
GPS (Pretrained)	✓	✓	0.967 \pm 0.002	0.945 \pm 0.004	0.934 \pm 0.005	0.864 \pm 0.009	0.759 \pm 0.010	0.674 \pm 0.002
GatedGCN-DENS	×	×	0.963 \pm 0.0002	0.943 \pm 0.001	0.933 \pm 0.001	0.874 \pm 0.002	0.794 \pm 0.002	0.731 \pm 0.002
	✓	×	0.949 \pm 0.008	0.922 \pm 0.008	0.907 \pm 0.011	0.828 \pm 0.020	0.733 \pm 0.032	0.662 \pm 0.046
	×	✓	0.965 \pm 0.001	0.943 \pm 0.001	0.933 \pm 0.001	0.873 \pm 0.001	0.792 \pm 0.004	<u>0.736 \pm 0.003</u>
	✓	✓	0.954 \pm 0.005	0.930 \pm 0.010	0.917 \pm 0.011	0.850 \pm 0.023	0.759 \pm 0.025	0.696 \pm 0.032
GPS-DENS	×	✓	0.980 \pm 0.000	0.969 \pm 0.000	0.961 \pm 0.000	0.913 \pm 0.000	0.834 \pm 0.000	0.750 \pm 0.000
	✓	✓	<u>0.978 \pm 0.001</u>	<u>0.963 \pm 0.000</u>	<u>0.953 \pm 0.001</u>	<u>0.905 \pm 0.000</u>	<u>0.822 \pm 0.002</u>	<u>0.736 \pm 0.003</u>

Table 5: **RotMNIST-Calibration.** Here, we report expanded results (calibration) on the Rotated MNIST dataset, including a variant that combines G- Δ UQ with Deep Ens. Notably, we see that anchored ensembles outperform basic ensembles in both accuracy and calibration. (Best results for models using Deep Ens. and those not using it marked separately.)

MODEL	G- Δ UQ	LPE?	Avg.ECE (\downarrow)	ECE (10) (\downarrow)	ECE (15) (\downarrow)	ECE (25) (\downarrow)	ECE (35) (\downarrow)	ECE (40) (\downarrow)
GatedGCN-TS	×	×	0.035 \pm 0.001	0.054 \pm 0.002	0.062 \pm 0.003	0.118 \pm 0.007	0.185 \pm 0.006	0.233 \pm 0.008
	×	✓	0.033 \pm 0.002	0.053 \pm 0.002	0.061 \pm 0.004	0.116 \pm 0.005	0.179 \pm 0.006	0.225 \pm 0.005
GatedGCN	×	×	0.038 \pm 0.001	0.059 \pm 0.001	0.068 \pm 0.008	0.126 \pm 0.008	0.195 \pm 0.012	0.245 \pm 0.011
	✓	×	0.018 \pm 0.008	<u>0.029 \pm 0.013</u>	0.033 \pm 0.164	<u>0.069 \pm 0.033</u>	<u>0.117 \pm 0.048</u>	<u>0.162 \pm 0.067</u>
	×	✓	0.036 \pm 0.003	0.059 \pm 0.002	0.068 \pm 0.008	0.125 \pm 0.006	0.191 \pm 0.007	0.240 \pm 0.008
	✓	✓	0.022 \pm 0.007	0.028 \pm 0.014	<u>0.034 \pm 0.169</u>	0.062 \pm 0.022	0.109 \pm 0.019	0.141 \pm 0.019
GPS-TS	×	✓	0.024 \pm 0.001	0.041 \pm 0.001	0.049 \pm 0.001	0.102 \pm 0.006	0.188 \pm 0.012	0.261 \pm 0.008
GPS	×	✓	0.026 \pm 0.001	0.044 \pm 0.001	0.052 \pm 0.156	0.108 \pm 0.006	0.197 \pm 0.012	0.273 \pm 0.008
	✓	✓	0.022 \pm 0.001	0.037 \pm 0.005	0.044 \pm 0.133	0.091 \pm 0.008	0.165 \pm 0.018	0.239 \pm 0.018
GPS (Pretrained)	✓	✓	<u>0.021 \pm 0.001</u>	0.032 \pm 0.003	0.039 \pm 0.116	0.083 \pm 0.002	0.153 \pm 0.007	0.217 \pm 0.012
GatedGCN-DENS	×	×	0.026 \pm 0.000	0.038 \pm 0.001	0.042 \pm 0.001	0.084 \pm 0.002	0.135 \pm 0.001	0.185 \pm 0.003
	✓	×	0.014 \pm 0.003	0.018 \pm 0.005	0.021 \pm 0.005	<u>0.036 \pm 0.012</u>	<u>0.069 \pm 0.032</u>	<u>0.114 \pm 0.056</u>
	×	✓	0.024 \pm 0.001	0.038 \pm 0.001	0.043 \pm 0.002	0.083 \pm 0.001	0.139 \pm 0.004	0.181 \pm 0.002
	✓	✓	0.017 \pm 0.002	0.024 \pm 0.005	<u>0.027 \pm 0.008</u>	0.030 \pm 0.004	0.036 \pm 0.012	0.059 \pm 0.033
GPS-DENS	×	✓	<u>0.016 \pm 0.001</u>	0.026 \pm 0.002	0.030 \pm 0.000	0.066 \pm 0.000	0.123 \pm 0.000	0.195 \pm 0.000
	✓	✓	0.014 \pm 0.000	<u>0.023 \pm 0.002</u>	<u>0.027 \pm 0.003</u>	0.055 \pm 0.004	0.103 \pm 0.006	0.164 \pm 0.006

Table 6: **Accuracy of GPS Variants on RotatedMNIST.** We focus on the accuracy results for GPS variants on rotated MNIST dataset. Using G- Δ UQ (with or without pretraining) remains close in accuracy to foregoing it, generally within the range of the standard deviation of the results.

MODEL	G- Δ UQ?	Avg. Test (\uparrow)	Acc. (10) (\uparrow)	Acc. (15) (\uparrow)	Acc. (25) (\uparrow)	Acc. (35) (\uparrow)	Acc. (40) (\uparrow)
GPS	×	0.970 \pm 0.001	0.948 \pm 0.001	0.938 \pm 0.001	0.873 \pm 0.006	0.770 \pm 0.013	0.688 \pm 0.009
	✓	<u>0.969 \pm 0.001</u>	<u>0.946 \pm 0.003</u>	<u>0.937 \pm 0.003</u>	<u>0.869 \pm 0.003</u>	<u>0.769 \pm 0.012</u>	<u>0.679 \pm 0.014</u>
GPS (Pretrained)	✓	0.967 \pm 0.002	0.945 \pm 0.004	0.934 \pm 0.005	0.864 \pm 0.009	0.759 \pm 0.010	0.674 \pm 0.002

Table 7: **Calibration of GPS Variants on RotatedMNIST.** We focus on the calibration results for GPS variants on rotated MNIST dataset. Across the board, we see improvements from using G- Δ UQ, with our strategy of applying it to a pretrained model doing best.

MODEL	G- Δ UQ	Avg.ECE (\downarrow)	ECE (10) (\downarrow)	ECE (15) (\downarrow)	ECE (25) (\downarrow)	ECE (35) (\downarrow)	ECE (40) (\downarrow)
GPS-TS	×	0.024 \pm 0.001	0.041 \pm 0.001	0.049 \pm 0.001	0.102 \pm 0.006	0.188 \pm 0.012	0.261 \pm 0.008
GPS	×	0.026 \pm 0.001	0.044 \pm 0.001	0.052 \pm 0.156	0.108 \pm 0.006	0.197 \pm 0.012	0.273 \pm 0.008
	✓	<u>0.022 \pm 0.001</u>	<u>0.037 \pm 0.005</u>	<u>0.044 \pm 0.133</u>	<u>0.091 \pm 0.008</u>	<u>0.165 \pm 0.018</u>	<u>0.239 \pm 0.018</u>
GPS (Pretrained)	✓	0.021 \pm 0.001	0.032 \pm 0.003	0.039 \pm 0.116	0.083 \pm 0.002	0.153 \pm 0.007	0.217 \pm 0.012

Table 8: **MNIST Feature Shifts**. G- Δ UQ improves calibration and maintains competitive or even improved accuracy across varying levels of feature distribution shift.

MODEL	LPE?	G- Δ UQ?	Calibration	STD = 0.1		STD = 0.2		STD = 0.3		STD = 0.4	
				Accuracy (\uparrow)	ECE (\downarrow)	Accuracy (\uparrow)	ECE (\downarrow)	Accuracy (\uparrow)	ECE (\downarrow)	Accuracy (\uparrow)	ECE (\downarrow)
GatedGCN	×	×	×	0.742 \pm 0.005	0.186 \pm 0.018	0.481 \pm 0.015	0.414 \pm 0.092	0.293 \pm 0.074	0.606 \pm 0.147	0.197 \pm 0.092	0.71 \pm 0.178
	×	✓	×	0.773\pm0.053	0.075\pm0.032	<u>0.536\pm0.010</u>	0.160\pm0.087	0.356\pm0.101	<u>0.422\pm0.083</u>	0.249\pm0.074	0.529\pm0.047
	✓	×	×	<u>0.751\pm0.02</u>	0.176 \pm 0.014	0.519 \pm 0.004	0.348 \pm 0.03	<u>0.345\pm0.032</u>	0.485 \pm 0.096	0.233 \pm 0.043	0.581 \pm 0.142
	✓	✓	×	0.745 \pm 0.026	<u>0.100\pm0.036</u>	0.541\pm0.040	<u>0.235\pm0.067</u>	0.355\pm0.062	0.408\pm0.116	<u>0.242\pm0.063</u>	<u>0.539\pm0.139</u>

A.5 STOCHASTIC CENTERING ON THE EMPIRICAL NTK OF GRAPH NEURAL NETWORKS

Using a simple grid-graph dataset and 4 layer GIN model, we compute the Fourier spectrum of the NTK. As shown in Fig. 7, we find that shifts to the node features can induce systematic changes to the spectrum.

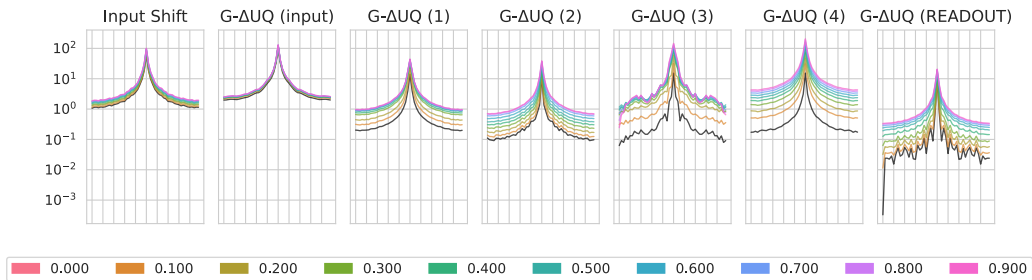


Figure 7: **Stochastic Centering with the empirical GNN NTK.** We find that performing constant shifts at intermediate layers introduces changes to a GNN’s NTK. We include a vanilla GNN NTK in black for reference. Further, note the shape of the spectrum should not be compared across subplots as each subplot was created with a different random initialization.

A.6 SIZE-GENERALIZATION DATASET STATISTICS

The statistics for the size generalization experiments (see Sec. 5.1) are provided below in Table 9.

Table 9: **Size Generalization Dataset Statistics:** This table is directly reproduced from (Buffelli et al., 2022), who in turn used statistics from (Yehudai et al., 2021; Bevilacqua et al., 2021).

	NCII			NCI109		
	ALL	SMALLEST 50%	LARGEST 10%	ALL	SMALLEST 50%	LARGEST 10%
CLASS A	49.95%	62.30%	19.17%	49.62%	62.04%	21.37%
CLASS B	50.04%	37.69%	80.82%	50.37%	37.95%	78.62%
# OF GRAPHS	4110	2157	412	4127	2079	421
AVG GRAPH SIZE	29	20	61	29	20	61

	PROTEINS			DD		
	ALL	SMALLEST 50%	LARGEST 10%	ALL	SMALLEST 50%	LARGEST 10%
CLASS A	59.56%	41.97%	90.17%	58.65%	35.47%	79.66%
CLASS B	40.43%	58.02%	9.82%	41.34%	64.52%	20.33%
# OF GRAPHS	1113	567	112	1178	592	118
AVG GRAPH SIZE	39	15	138	284	144	746

A.7 GOOD BENCHMARK EXPERIMENTAL DETAILS

For our experiments in Sec. 5.2, we utilize the in/out-of-distribution covariate and concept splits provided by Gui et al. (2022). Furthermore, we use the suggested models and architectures provided by their package. In brief, we use GIN models with virtual nodes (except for GOODMotif) for training, and average scores over 3 seeds. When performing stochastic anchoring at a particular layer, we double the hidden representation size for that layer. Subsequent layers retain the original size of the vanilla model.

When performing stochastic anchoring, we use 10 fixed anchors randomly drawn from the in-distribution validation dataset. We also train the G- Δ UQ for an additional 50 epochs to ensure that models are able to converge. Please see our code repository for the full details.

We also include results on additional node classification benchmarks featuring distribution shift in Table 12. In Table 14, we compare models without G- Δ UQ to the use of G- Δ UQ with randomly sampled anchors at the first or second layer.

Dataset	Shift	Train	ID validation	ID test	OOD validation	OOD test	Train	OOD validation	ID validation	ID test	OOD test
Length											
GOOD-SST2	covariate	24744	5301	5301	17206	17490					
	concept	27270	5843	5843	15142	15944					
Color											
GOOD-CMNIST	covariate	42000	7000	7000	7000	7000					
	concept	29400	6300	6300	14000	14000					
	no shift	42000	14000	14000	-	-					
Base						Size					
GOOD-Motif	covariate	18000	3000	3000	3000	3000	18000	3000	3000	3000	3000
	concept	12600	2700	2700	6000	6000	12600	2700	2700	6000	6000
Word						Degree					
GOOD-Cora	covariate	9378	1979	1979	3003	3454	8213	1979	1979	3841	3781
	concept	7273	1558	1558	3807	5597	7281	1560	1560	3706	5686
University											
GOOD-WebKB	covariate	244	61	61	125	126					
	concept	282	60	60	106	109					
Color											
GOOD-CBAS	covariate	420	70	70	70	70					
	concept	140	140	140	140	140					

Table 10: Number of Graphs/Nodes per dataset.

Dataset	Model	# Model layers	Batch Size	# Max Epochs	# Iterations per epoch	Initial LR	Node Feature Dim
GOOD-SST2	GIN-Virtual	3	32	200/100	-	1e-3	768
GOOD-CMNIST	GIN-Virtual	5	128	500	-	1e-3	3
GOOD-Motif	GIN	3	32	200	-	1e-3	4
GOOD-Cora	GCN	3	4096	100	10	1e-3	8710
GOOD-WebKB	GCN	3	4096	100	10	1e-3/5e-3	1703
GOOD-CBAS	GCN	3	1000	200	10	3e-3	8

Table 11: Model and hyperparameters for GOOD datasets.

A.8 GOOD DATASET ADDITIONAL RESULTS

We also include results on additional node classification benchmarks featuring distribution shift in Table 12. In Table 14, we compare models without G- Δ UQ to the use of G- Δ UQ with randomly sampled anchors at the first or second layer.

Table 12: **Additional Node Classification Benchmarks.** Here, we compare accuracy and calibration error of G- Δ UQ and "no G- Δ UQ" (vanilla) models on 4 node classification benchmarks across concept and covariate shifts. First, we note that across all our evaluations, without any posthoc calibration, G- Δ UQ is superior to the vanilla model on nearly every benchmark for better or same accuracy (8/8 benchmarks) and better calibration error (7/8), often with a significant gain in calibration performance. However, due to the challenging nature of these shifts, achieving state-of-the-art calibration performance often requires the use of post-hoc calibration methods – so we also evaluate how these posthoc methods can be elevated when combined with G- Δ UQ (versus the vanilla variant). When combined with popular posthoc methods, we highlight that performance improves across the board, when combined with G- Δ UQ (including in WebKB and CBAS-Concept). For example, on WebKB – across the 9 calibration methods considered, "G- Δ UQ + calibration method" improves or maintains the calibration performance of the analogous "no G- Δ UQ + calibration method" in 7/9 (concept) and 6/9 (covariate). In CBAS, calibration is improved or maintained as the no-G- Δ UQ version on 5/9 (concept) and 9/9 (covariate). In all cases, this is achieved with little or no compromise on classification accuracy (often improving over "no G- Δ UQ" variant). We also emphasize that, across all the 8 evaluation sets (4 datasets x 2 shift types) in Table 10, the best performance is almost always obtained with a GDUQ variant: (accuracy: 8/8) as well as best calibration (6/8) or second best (2/8).

Dataset	Domain	Calibration	Shift: Concept				Shift: Covariate			
			Accuracy (\uparrow)		ECE (\downarrow)		Accuracy (\uparrow)		ECE (\downarrow)	
			No G- Δ UQ	G- Δ UQ	No G- Δ UQ	G- Δ UQ	No G- Δ UQ	G- Δ UQ	No G- Δ UQ	G- Δ UQ
WebKB	University	X	0.253 \pm 0.003	0.281\pm0.009	0.67 \pm 0.061	0.593 \pm 0.025	0.122 \pm 0.029	0.115 \pm 0.041	0.599 \pm 0.091	0.525 \pm 0.033
		CAGCN	0.253 \pm 0.005	0.268 \pm 0.008	0.452 \pm 0.14	0.473 \pm 0.12	0.122 \pm 0.018	0.092 \pm 0.161	0.355 \pm 0.227	0.396 \pm 0.161
		Dirichlet	0.229 \pm 0.018	0.22 \pm 0.022	0.472 \pm 0.06	0.472 \pm 0.03	0.244\pm0.105	0.295\pm0.044	0.299\pm0.092	0.328\pm0.044
		ETS	0.253 \pm 0.005	0.273 \pm 0.012	0.64 \pm 0.06	0.575 \pm 0.019	0.121 \pm 0.021	0.084 \pm 0.027	0.539 \pm 0.112	0.499 \pm 0.027
		GATS	0.253 \pm 0.005	0.273 \pm 0.01	0.608 \pm 0.008	0.485 \pm 0.02	0.122 \pm 0.018	0.079 \pm 0.029	0.455 \pm 0.057	0.376 \pm 0.029
		IRM	0.251 \pm 0.005	0.266 \pm 0.011	0.342\pm0.017	0.349\pm0.006	0.097 \pm 0.04	0.046 \pm 0.013	0.352 \pm 0.037	0.422 \pm 0.013
		Orderinvariant	0.253 \pm 0.005	0.27 \pm 0.01	0.628 \pm 0.026	0.564 \pm 0.024	0.122 \pm 0.018	0.106 \pm 0.065	0.545 \pm 0.079	0.47 \pm 0.065
		Spline	0.237 \pm 0.012	0.257 \pm 0.023	0.436 \pm 0.029	0.386 \pm 0.034	0.122 \pm 0.013	0.171 \pm 0.056	0.472 \pm 0.031	0.39 \pm 0.056
		VS	0.253 \pm 0.005	0.275\pm0.011	0.67 \pm 0.009	0.588 \pm 0.011	0.122 \pm 0.018	0.095 \pm 0.014	0.602 \pm 0.044	0.507 \pm 0.014
Cora	Degree	X	0.581 \pm 0.003	0.595 \pm 0.003	0.307 \pm 0.009	0.13 \pm 0.011	0.47 \pm 0.002	0.518 \pm 0.014	0.348 \pm 0.032	0.141 \pm 0.008
		CAGCN	0.581 \pm 0.003	0.597\pm0.002	0.135 \pm 0.009	0.128 \pm 0.025	0.47 \pm 0.002	0.522\pm0.025	0.256 \pm 0.08	0.231 \pm 0.025
		Dirichlet	0.534 \pm 0.007	0.551 \pm 0.004	0.12 \pm 0.004	0.196 \pm 0.003	0.414 \pm 0.007	0.449 \pm 0.01	0.163 \pm 0.002	0.356 \pm 0.01
		ETS	0.581 \pm 0.003	0.596\pm0.004	0.301 \pm 0.009	0.116 \pm 0.018	0.47 \pm 0.002	0.523\pm0.003	0.31 \pm 0.077	0.141 \pm 0.003
		GATS	0.581 \pm 0.003	0.596\pm0.004	0.185 \pm 0.018	0.229 \pm 0.039	0.47 \pm 0.002	0.521 \pm 0.011	0.211 \pm 0.004	0.308 \pm 0.011
		IRM	0.582 \pm 0.002	0.597 \pm 0.002	0.125 \pm 0.001	0.102 \pm 0.002	0.469 \pm 0.001	0.522\pm0.004	0.194 \pm 0.005	0.13\pm0.004
		Orderinvariant	0.581 \pm 0.003	0.592 \pm 0.002	0.226 \pm 0.024	0.213 \pm 0.049	0.47 \pm 0.002	0.498 \pm 0.027	0.318 \pm 0.042	0.196 \pm 0.027
		Spline	0.571 \pm 0.003	0.595 \pm 0.003	0.080\pm0.004	0.068\pm0.004	0.459 \pm 0.003	0.52 \pm 0.004	0.158 \pm 0.01	0.098\pm0.004
		VS	0.581 \pm 0.003	0.596\pm0.004	0.306 \pm 0.004	0.127 \pm 0.002	0.47 \pm 0.001	0.522\pm0.005	0.345 \pm 0.005	0.146 \pm 0.005
Cora	Word	X	0.607 \pm 0.003	0.628 \pm 0.001	0.284 \pm 0.009	0.111 \pm 0.013	0.603 \pm 0.004	0.633 \pm 0.031	0.263 \pm 0.004	0.118 \pm 0.019
		CAGCN	0.607 \pm 0.002	0.628 \pm 0.002	0.138 \pm 0.011	0.236 \pm 0.019	0.603 \pm 0.004	0.634 \pm 0.035	0.129 \pm 0.009	0.253 \pm 0.035
		Dirichlet	0.579 \pm 0.007	0.588 \pm 0.006	0.105 \pm 0.011	0.168 \pm 0.005	0.562 \pm 0.007	0.578 \pm 0.007	0.095 \pm 0.006	0.269 \pm 0.007
		ETS	0.607 \pm 0.002	0.628 \pm 0.002	0.282 \pm 0.002	0.11 \pm 0.003	0.603 \pm 0.004	0.634 \pm 0.013	0.243 \pm 0.023	0.106 \pm 0.013
		GATS	0.607 \pm 0.002	0.628 \pm 0.002	0.166 \pm 0.009	0.261 \pm 0.028	0.603 \pm 0.004	0.635\pm0.037	0.16 \pm 0.015	0.293 \pm 0.037
		IRM	0.608 \pm 0.001	0.63 \pm 0.002	0.115 \pm 0.002	0.088 \pm 0.003	0.602 \pm 0.003	0.635\pm0.004	0.106 \pm 0.002	0.098 \pm 0.004
		Orderinvariant	0.607 \pm 0.002	0.624 \pm 0.002	0.174 \pm 0.024	0.201 \pm 0.061	0.603 \pm 0.004	0.621 \pm 0.076	0.154 \pm 0.022	0.202 \pm 0.076
		Spline	0.598 \pm 0.005	0.629\pm0.002	0.073\pm0.002	0.062\pm0.005	0.591 \pm 0.002	0.635\pm0.004	0.063\pm0.006	0.053\pm0.004
		VS	0.607 \pm 0.001	0.63\pm0.002	0.283 \pm 0.003	0.111 \pm 0.003	0.603 \pm 0.004	0.636\pm0.003	0.261 \pm 0.005	0.119 \pm 0.003
CBAS	Color	X	0.83 \pm 0.014	0.829 \pm 0.011	0.169 \pm 0.013	0.151 \pm 0.014	0.703 \pm 0.015	0.746 \pm 0.027	0.266 \pm 0.02	0.169 \pm 0.018
		CAGCN	0.83 \pm 0.013	0.83 \pm 0.013	0.137\pm0.011	0.143 \pm 0.022	0.703 \pm 0.019	0.749 \pm 0.033	0.25 \pm 0.021	0.186 \pm 0.017
		Dirichlet	0.801 \pm 0.02	0.806 \pm 0.008	0.161 \pm 0.012	0.17 \pm 0.01	0.671 \pm 0.018	0.771 \pm 0.03	0.241 \pm 0.029	0.217 \pm 0.017
		ETS	0.83 \pm 0.013	0.827 \pm 0.014	0.146 \pm 0.013	0.164 \pm 0.007	0.703 \pm 0.019	0.76 \pm 0.037	0.28 \pm 0.023	0.176 \pm 0.019
		GATS	0.83 \pm 0.013	0.83 \pm 0.021	0.16 \pm 0.009	0.173 \pm 0.021	0.703 \pm 0.019	0.751 \pm 0.016	0.236 \pm 0.039	0.16\pm0.015
		IRM	0.829 \pm 0.013	0.839\pm0.015	0.142 \pm 0.009	0.133\pm0.006	0.72 \pm 0.019	0.803\pm0.04	0.207 \pm 0.035	0.158\pm0.017
		Orderinvariant	0.83 \pm 0.013	0.803 \pm 0.008	0.174 \pm 0.006	0.173 \pm 0.009	0.703 \pm 0.019	0.766 \pm 0.045	0.261 \pm 0.017	0.194 \pm 0.031
		Spline	0.82 \pm 0.016	0.824 \pm 0.011	0.159 \pm 0.009	0.16 \pm 0.014	0.683 \pm 0.019	0.786 \pm 0.038	0.225 \pm 0.034	0.179 \pm 0.035
		VS	0.829 \pm 0.012	0.840\pm0.011	0.166 \pm 0.011	0.146 \pm 0.012	0.717 \pm 0.019	0.809\pm0.008	0.242 \pm 0.019	0.182 \pm 0.014

Table 15: **Alternative Anchoring Strategies.** Here, we consider an alternative anchoring formulation for graph classification. Namely, instead of shuffling features across the batch (denoted Batch in the table), we perform READOUT anchoring by fitting a normal distribution over the hidden representations. We then randomly sample from this distribution to create anchors. Conceptually, this is similar to the node feature anchoring strategy. One potential direction of future work that is permitted by this formulation is to optimize the parameters of this distribution given a signal from an appropriate auxiliary task or loss. For example, we could perform an alternating optimization where the GNN is trained to minimize the loss, and the mean and variance of the anchoring distribution are optimized to minimize the expected calibration error on a separate calibration dataset. While a rigorous formulation is left to future work, we emphasize that the potential for improving the anchoring distribution, and thus controlling corresponding hypothesis diversity, is in fact a unique benefit of G- Δ UQ.

Shift Type	Method	Test Acc			Test Cal			OOD Acc			OOD Cal		
		MPNN	Batch	Random	MPNN	Batch	Random	MPNN	Batch	Random	MPNN	Batch	Random
GoodMoif_basis_concept	Dirichlet	0.995 ± 0.0007	0.994 ± 0.0002	0.996 ± 0.0009	0.040 ± 0.0037	0.036 ± 0.0016	0.035 ± 0.0054	0.924 ± 0.0069	0.923 ± 0.0117	0.942 ± 0.0034	0.080 ± 0.0153	0.102 ± 0.0071	0.063 ± 0.0086
	ETS	0.995 ± 0.0007	0.995 ± 0.0005	0.996 ± 0.0007	0.035 ± 0.0034	0.036 ± 0.0101	0.032 ± 0.0052	0.925 ± 0.0095	0.926 ± 0.009	0.935 ± 0.0068	0.095 ± 0.0098	0.096 ± 0.0128	0.087 ± 0.0141
	IRM	0.995 ± 0.0007	0.9957 ± 0.0009	0.996 ± 0.0004	0.0198 ± 0.0089	0.0229 ± 0.0105	0.0225 ± 0.0038	0.925 ± 0.0096	0.9301 ± 0.0123	0.946 ± 0.0024	0.0873 ± 0.0170	0.0966 ± 0.0103	0.0907 ± 0.0276
	OderInvariant	0.995 ± 0.0007	0.995 ± 0.0005	0.995 ± 0.0005	0.033 ± 0.0094	0.028 ± 0.0047	0.032 ± 0.0009	0.925 ± 0.0095	0.928 ± 0.0104	0.935 ± 0.0027	0.090 ± 0.0092	0.093 ± 0.0070	0.075 ± 0.0029
	Spline	0.995 ± 0.0007	0.995 ± 0.0007	0.996 ± 0.0005	0.034 ± 0.0002	0.035 ± 0.0090	0.032 ± 0.0048	0.924 ± 0.0098	0.926 ± 0.0092	0.937 ± 0.0030	0.091 ± 0.0084	0.089 ± 0.0123	0.082 ± 0.0065
	VS	0.995 ± 0.0007	0.995 ± 0.0005	0.996 ± 0.000	0.035 ± 0.0034	0.036 ± 0.0087	0.033 ± 0.0068	0.925 ± 0.0094	0.928 ± 0.0095	0.936 ± 0.0053	0.094 ± 0.0066	0.095 ± 0.0133	0.082 ± 0.009
GoodMoif_basis_covariate	Dirichlet	0.999 ± 0.0003	0.999 ± 0.0004	0.999 ± 0.0002	0.017 ± 0.0054	0.017 ± 0.0019	0.014 ± 0.0004	0.685 ± 0.0504	0.650 ± 0.0450	0.698 ± 0.0139	0.336 ± 0.0667	0.371 ± 0.0474	0.320 ± 0.0140
	ETS	0.997 ± 0.0004	0.999 ± 0.0005	0.999 ± 0.0002	0.0095 ± 0.0091	0.017 ± 0.0064	0.017 ± 0.0056	0.690 ± 0.0434	0.649 ± 0.0476	0.686 ± 0.0226	0.313 ± 0.0413	0.379 ± 0.0185	0.334 ± 0.0167
	IRM	0.997 ± 0.0004	0.999 ± 0.0006	0.999 ± 0.0003	0.0085 ± 0.0012	0.010 ± 0.0032	0.014 ± 0.0042	0.690 ± 0.0434	0.647 ± 0.0472	0.692 ± 0.0226	0.315 ± 0.0595	0.354 ± 0.0450	0.328 ± 0.0211
	OderInvariant	0.997 ± 0.0004	0.999 ± 0.0005	0.999 ± 0.0003	0.014 ± 0.0028	0.020 ± 0.0090	0.013 ± 0.0081	0.690 ± 0.0434	0.649 ± 0.0450	0.689 ± 0.0170	0.320 ± 0.0501	0.358 ± 0.0410	0.328 ± 0.0218
	Spline	0.997 ± 0.0004	0.999 ± 0.0005	0.999 ± 0.0003	0.016 ± 0.0049	0.017 ± 0.0053	0.017 ± 0.0052	0.690 ± 0.0434	0.649 ± 0.0476	0.692 ± 0.0199	0.324 ± 0.0548	0.373 ± 0.0507	0.327 ± 0.0165
	VS	0.999 ± 0.0001	0.999 ± 0.0003	0.999 ± 0.0002	0.011 ± 0.0033	0.014 ± 0.0034	0.012 ± 0.0016	0.682 ± 0.0561	0.650 ± 0.0546	0.687 ± 0.0251	0.325 ± 0.0568	0.374 ± 0.0591	0.337 ± 0.0264
GOODSST2_length_concept	Dirichlet	0.938 ± 0.0019	0.939 ± 0.0056	0.942 ± 0.00180	0.189 ± 0.01989	0.165 ± 0.0179	0.187 ± 0.0256	0.684 ± 0.0193	0.693 ± 0.0020	0.687 ± 0.0027	0.146 ± 0.0196	0.133 ± 0.015	0.169 ± 0.0168
	ETS	0.938 ± 0.0020	0.939 ± 0.0060	0.941 ± 0.0017	0.189 ± 0.0018	0.180 ± 0.0022	0.183 ± 0.0007	0.694 ± 0.0193	0.692 ± 0.0019	0.687 ± 0.0034	0.216 ± 0.0033	0.220 ± 0.0057	
	IRM	0.939 ± 0.0016	0.939 ± 0.0058	0.941 ± 0.0018	0.126 ± 0.0011	0.126 ± 0.0013	0.127 ± 0.0017	0.697 ± 0.0185	0.692 ± 0.0028	0.685 ± 0.0026	0.240 ± 0.0017	0.232 ± 0.0050	0.242 ± 0.0053
	OderInvariant	0.938 ± 0.0020	0.939 ± 0.0060	0.941 ± 0.0022	0.114 ± 0.0014	0.115 ± 0.0029	0.115 ± 0.0012	0.694 ± 0.0193	0.692 ± 0.0019	0.687 ± 0.0033	0.224 ± 0.0010	0.222 ± 0.0030	0.225 ± 0.0054
	Spline	0.938 ± 0.0026	0.938 ± 0.0044	0.941 ± 0.0010	0.129 ± 0.0021	0.129 ± 0.0019	0.128 ± 0.0012	0.692 ± 0.0190	0.692 ± 0.0022	0.687 ± 0.0035	0.234 ± 0.0052	0.231 ± 0.0044	0.243 ± 0.0034
	VS	0.938 ± 0.0027	0.939 ± 0.0057	0.941 ± 0.0018	0.290 ± 0.0099	0.484 ± 0.0008	0.487 ± 0.0007	0.693 ± 0.0184	0.693 ± 0.0018	0.687 ± 0.0031	0.331 ± 0.0484	0.375 ± 0.0222	0.382 ± 0.0048
GOODSST2_length_covariate	Dirichlet	0.896 ± 0.0029	0.893 ± 0.0009	0.895 ± 0.00095	0.196 ± 0.0155	0.172 ± 0.0091	0.197 ± 0.0109	0.825 ± 0.0037	0.827 ± 0.0066	0.805 ± 0.0150	0.163 ± 0.0198	0.141 ± 0.0087	0.142 ± 0.0122
	ETS	0.896 ± 0.0023	0.894 ± 0.0011	0.894 ± 0.0006	0.157 ± 0.0013	0.159 ± 0.0004	0.162 ± 0.0019	0.826 ± 0.0036	0.828 ± 0.0065	0.806 ± 0.0117	0.309 ± 0.0050	0.314 ± 0.0076	0.300 ± 0.0070
	IRM	0.895 ± 0.0019	0.893 ± 0.0003	0.894 ± 0.0007	0.107 ± 0.0004	0.107 ± 0.0003	0.106 ± 0.0020	0.826 ± 0.0040	0.828 ± 0.0065	0.809 ± 0.0152	0.276 ± 0.0046	0.277 ± 0.0061	0.265 ± 0.0078
	OderInvariant	0.896 ± 0.0023	0.894 ± 0.0011	0.894 ± 0.0008	0.288 ± 0.0008	0.285 ± 0.0008	0.284 ± 0.0013	0.826 ± 0.0036	0.828 ± 0.0065	0.806 ± 0.0106	0.241 ± 0.0022	0.241 ± 0.0063	0.254 ± 0.0054
	Spline	0.894 ± 0.0016	0.890 ± 0.0009	0.892 ± 0.0010	0.309 ± 0.0024	0.307 ± 0.0009	0.307 ± 0.0022	0.822 ± 0.0026	0.822 ± 0.0092	0.801 ± 0.0110	0.275 ± 0.0043	0.276 ± 0.0063	0.264 ± 0.0063
	VS	0.898 ± 0.0028	0.893 ± 0.0008	0.894 ± 0.0007	0.291 ± 0.0183	0.460 ± 0.0011	0.465 ± 0.0010	0.821 ± 0.0053	0.827 ± 0.0071	0.806 ± 0.0119	0.299 ± 0.0188	0.431 ± 0.0061	0.429 ± 0.0054

A.9 POST-HOC CALIBRATION STRATEGIES

Several post hoc strategies have been developed for calibrating the predictions of a model. These have the advantage of flexibility, as they operate only on the outputs of a model and do not require that any changes be made to the model itself. Some methods include:

- **Temperature scaling (TS)** (Guo et al., 2017) simply scales the logits by a temperature parameter $T > 1$ to smooth the predictions. The scaling parameter T can be tuned on a validation set.
- **Ensemble temperature scaling (ETS)** (Zhang et al., 2020) learns an ensemble of temperature-scaled predictions with uncalibrated predictions ($T = 1$) and uniform probabilistic outputs ($T = \infty$).
- **Vector scaling (VS)** (Guo et al., 2017) scales the entire output vector of class probabilities, rather than just the logits.
- **Multi-class isotonic regression (IRM)** (Zhang et al., 2020) is a multiclass generalization of the famous isotonic regression method (Zadrozny & Elkan, 2002): it ensembles predictions and labels, then learns a monotonically increasing function to map transformed predictions to labels.
- **Order-invariant calibration** (Rahimi et al., 2020) uses a neural network to learn an intra-order-preserving calibration function that can preserve a model’s top-k predictions.
- **Spline calibration** instead uses splines to fit the calibration function (Gupta et al., 2021).
- **Dirichlet calibration** (Kull et al., 2019) models the distribution of outputs using a Dirichlet distribution, using simple log-transformation of the uncalibrated probabilities which are then passed to a regularized fully connected neural network layer with softmax activation.

For node classification, some graph-specific post-hoc calibration methods have been proposed. **CaGCN** (Wang et al., 2021) uses the graph structure and an additional GCN to produce node-wise temperatures. **GATS** (Hsu et al., 2022) extends this idea by using graph attention to model the influence of neighbors’ temperatures when learning node-wise temperatures. We use the post hoc calibration baselines provided by Hsu et al. in our experiments.

All of the above methods, and others, may be applied to the output of any model including one using G- Δ UQ. As we have shown, applying such post hoc methods to the outputs of the calibrated models may improve uncertainty estimates even more. Notably, calibrated models are expected to produce confidence estimates that match the true probabilities of the classes being predicted (Naeini et al., 2015; Guo et al., 2017; Ovadia et al., 2019). While poorly calibrated CIs are over/under confident in their predictions, calibrated CIs are more trustworthy and can also improve performance on other safety-critical tasks which implicitly require reliable prediction probabilities (see Sec. 5). We report the top-1 label expected calibration error (ECE) (Kumar et al., 2019; Detlefsen et al., 2022). Formally, let p_i be the top-1 probability, c_i be the predicted confidence, b_i a uniformly sized bin in $[0, 1]$. Then,

$$ECE := \sum_i^N b_i \| (p_i - c_i) \|$$

A.10 DETAILS ON GENERALIZATION GAP PREDICTION

Accurate estimation of the expected generalization error on unlabeled datasets allows models with unacceptable performance to be pulled from production. To this end, generalization error predictors (GEPs) (Garg et al., 2022; Ng et al., 2022; Jiang et al., 2019; Trivedi et al., 2023a; Guillory et al., 2021) which assign sample-level scores, $S(x_i)$ which are then aggregated into dataset-level error estimates, have become popular. We use maximum softmax probability and a simple thresholding mechanism as the GEP (since we are interested in understanding the behavior of confidence indicators), and report the error between the predicted and true target dataset accuracy: $GEPError := ||Acc_{target} - \frac{1}{|X|} \sum_i \mathbb{I}(S(\bar{x}_i; F) > \tau)||$ where τ is tuned by minimizing GEP error on the validation dataset. We use the confidences obtained by the different baselines as sample-level scores, $S(x_i)$ corresponding to the model’s expectation that a sample is correct. The MAE between the estimated error and true error is reported on both in- and out-of -distribution test splits provided by the GOOD benchmark.

A.11 RESULTS ON GENERALIZATION ERROR PREDICTION

GEP Experimental Setup. GEPs (Garg et al., 2022; Ng et al., 2022; Jiang et al., 2019; Trivedi et al., 2023a; Guillory et al., 2021) aggregate sample-level scores capturing a model’s uncertainty about the correctness of a prediction into dataset-level error estimates. Here, we use maximum softmax probability for scores and a thresholding mechanism as the GEP. (See Appendix A.10 for more details.) We consider READOUT anchoring with both pretrained and end-to-end training, and report the mean absolute error between the predicted and true target dataset accuracy on the OOD test split.

GEP Results. As shown in Table 16, both pretrained and end-to-end G- Δ UQ outperform the vanilla model on 7/8 datasets. Notably, we see that pretrained G- Δ UQ is particularly effective as it obtains the best performance across 6/8 datasets. This not only highlights its utility as a flexible, light-weight strategy for improving uncertainty estimates without sacrificing accuracy, but also emphasizes that importance of structure, in lieu of full stochasticity, when estimating GNN uncertainties.

Table 16: **GOOD-Datasets, Generalization Error Prediction Performance.** The MAE between the predicted and true test error on the OOD test split is reported. G- Δ UQ variants outperform vanilla models on 7/8 datasets (GOODMotif(Basis,Covariate) being the exception). **Pretrained G- Δ UQ is particularly effective at this task as it achieves the best performance overall on 6/8 datasets.** Promisingly, we see that regular G- Δ UQ improves performance over the vanilla model on 6/8 datasets (even if it is not the best overall). We further observe that performing generalization error prediction is more challenging under covariate shift than concept shift on the GOODCMNIST, GOODMotif(Basis) and GOODMotif(Size) datasets. On these datasets, the MAE is almost twice as large than their respective concept shift counterparts, across methods. GOODSST2 is the exception, where concept shift is in fact more challenging. To the best of our knowledge, we are the first to investigate generalization error prediction on GNN-based tasks under distribution shift. Understanding this behavior further is an interesting direction of future work.

Method	CMNIST (Color)		MotifLPE (Basis)		MotifLPE (Size)		SST2	
	Concept(\downarrow)	Covariate (\downarrow)	Concept(\downarrow)	Covariate(\downarrow)	Concept(\downarrow)	Covariate(\downarrow)	Concept(\downarrow)	Covariate(\downarrow)
Vanilla	0.200 \pm 0.009	0.510 \pm 0.089	0.045 \pm 0.003	0.570 \pm 0.012	0.324 \pm 0.018	0.537 \pm 0.146	0.117 \pm 0.006	0.056 \pm 0.044
G- Δ UQ	0.190 \pm 0.010	0.493 \pm 0.072	0.023 \pm 0.003	0.572 \pm 0.019	0.317 \pm 0.007	0.528 \pm 0.189	0.124 \pm 0.016	0.054 \pm 0.043
Pretr. G- Δ UQ	0.192 \pm 0.005	0.387 \pm 0.048	0.018 \pm 0.012	0.573 \pm 0.004	0.307 \pm 0.016	0.356 \pm 0.143	0.114 \pm 0.004	0.030 \pm 0.026

A.12 ADDITIONAL STUDY ON PRETRAINED G- Δ UQ

For the datasets and data shifts on which we reported out-of-distribution calibration error of pretrained vs. in-training G- Δ UQ earlier in Fig. 4, we now report additional results for in-distribution and out-of-distribution accuracy as well as calibration error. We also include results for the additional GOODMotif-basis benchmark for completeness, noting that the methods provided by the original benchmark Gui et al. (2022) generalized poorly to this split (which may be related to why G- Δ UQ methods offer little improvement over the vanilla model.) Fig. 8 shows these extended results. By these additional metrics, we again see the competitiveness of applying G- Δ UQ to a pretrained model versus using it in end-to-end training.

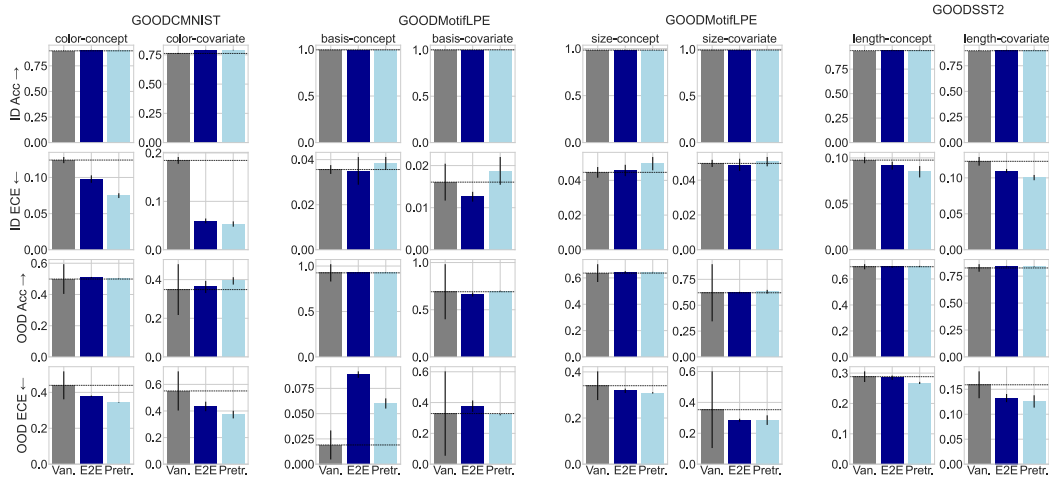


Figure 8: Evaluating Pretrained G- Δ UQ. Here, we report the performance of pretrained G- Δ UQ models vs. end-to-end and vanilla models with respect to in-distribution and out-of-distribution accuracy as well as expected calibration error. With the exception of the GOODMotif (basis) dataset, pretrained G- Δ UQ improves the OOD ECE over both the vanilla model and end-to-end G- Δ UQ at comparable or improved OOD accuracy on 7/8 datasets. Furthermore, pretrained G- Δ UQ also improves the ID ECE on all but the GOODMotif (size) datasets (6/8), where it performs comparably to the vanilla model, and maintains the ID accuracy. (We note that all methods are comparably better calibrated on the GOODMotif ID data than GOODCMNIST/GOODSST2 ID data; we suspect this is because there may exist simple shortcuts available in the GOODMotif dataset that can be used on the ID test set effectively.) Overall, these results clearly demonstrate that pretrained G- Δ UQ does offer some performance advantages over end-to-end G- Δ UQ and does so at reduced training times (see Table. A.13). For example, on GOODCMNIST (covariate shift), pretrained G- Δ UQ is not only 50% faster than end-to-end G- Δ UQ, it also improves OOD accuracy and OOD ECE over both the vanilla and end-to-end G- Δ UQ models.

A.13 RUNTIME TABLE

Table 17: **Runtimes.** We include the runtimes of both training per epoch (in seconds) and performing calibration. Reducing stochasticity can help reduce computation (L1 \rightarrow L3). Cost can also be reduced by using a pretrained model.

Dataset	GOODCMNIST		GOODSST2		GOODMotifLPE	
	Training (S)	Inference (S)	Training (S)	Inference (S)	Training (S)	Inference (S)
Vanilla	18.5	25.8	10.8	18.5	3.8	4.5
Temp. Scaling	18.5	23.5	10.8	13.4	3.8	5.3
DEns (Ens Size=3)	18.456 x Ens Size	59.4	10.795 x Ens Size	29.0	3.8 x Ens Size	11.8
G- Δ UQ (L1, 10 anchors)	22.1	181.5	15.9	17.1	5.8	15.5
G- Δ UQ (L2, 10)	22.4	148.6	12.7	15.5	5.8	11.8
G- Δ UQ (HiddenRep, 10)	18.5	28.0	13.8	19.6	3.9	6.5
G- Δ UQ (Pretr. HiddenRep, 10)	8.6	27.8	6.8	16.0	2.5	6.4

A.14 MEAN AND VARIANCE OF NODE FEATURE GAUSSIANS

Table 18: **Mean and Variance of Node Feature Anchoring Gaussians.** We report the mean and variance of the Gaussian distributions fitted to the input node features. Because the input node features vary in size, we report aggregate statistics over the mean and variance corresponding to each dimension. For example, Min(Mu) indicates that we are reports the minimum mean over the d -dim set of means.

Dataset	Domain	Shift	Min (Mu)	Max (Mu)	Mean (Mu)	Std (Mu)	Min (Std)	Max (Std)	Mean (Std)	Std (Std)
GOODSST2	length	concept	-4.563	0.69	-0.011	0.278	0.163	0.803	0.242	0.049
GOODSST2	length	covariate	-4.902	0.684	-0.01	0.3	0.175	0.838	0.255	0.05
GOODCMNIST	color	concept	0.117	0.133	0.127	0.008	0.092	0.097	0.095	0.003
GOODCMNIST	color	covariate	0.087	0.131	0.102	0.025	0.108	0.109	0.108	0
GOODMotifLPE	size	covariate	0.003	0.021	0.011	0.008	0.835	1.728	1.248	0.377
GOODMotifLPE	size	concept	-0.006	0	-0.002	0.003	0.542	1.114	0.783	0.242
GOODMotifLPE	basis	concept	-0.011	0.015	0.001	0.011	0.721	1.464	1.09	0.304
GOODMotifLPE	basis	covariate	-0.007	-0.002	-0.004	0.002	0.808	1.913	1.251	0.469
GOODWebKB	university	concept	0	0.95	0.049	0.099	0.001	0.5	0.168	0.095
GOODWebKB	university	covariate	0	0.934	0.05	0.104	0.001	0.5	0.164	0.098
GOODCora	degree	concept	0	0.507	0.007	0.017	0.001	0.5	0.061	0.051
GOODCora	degree	covariate	0	0.518	0.007	0.017	0.001	0.5	0.061	0.052
GOODCBAS	color	covariate	0.394	0.591	0.471	0.093	0.142	0.492	0.403	0.174
GOODCBAS	color	concept	0.23	0.569	0.4	0.144	0.168	0.495	0.39	0.152

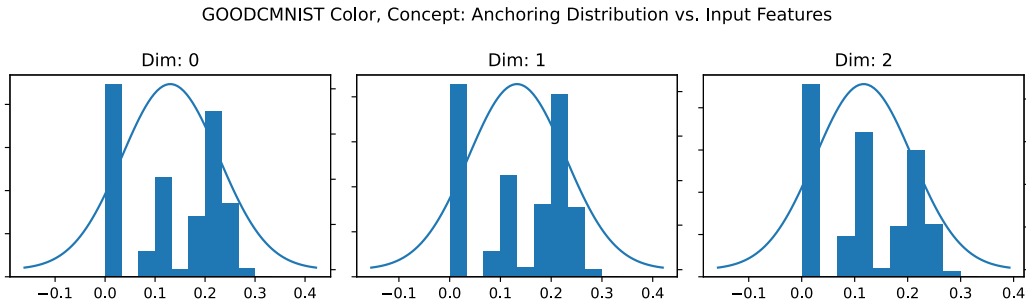


Figure 9: **GOODCMNIST, Concept, Anchoring Distribution.** We plot the mean and variance of the fitted anchoring distribution vs. the true feature distribution for each input dimension. We observe there is a mismatch between the empirical distribution and the fitted Gaussian. However, we did not find this mismatch to harm the effectiveness of G- Δ UQ.

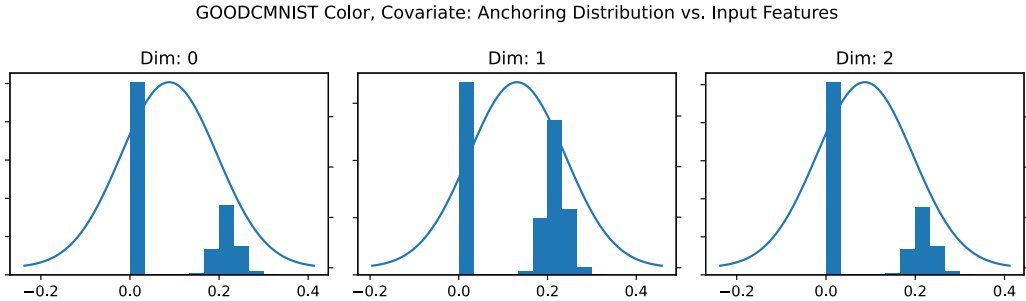


Figure 10: **GOODCMNIST, Covariate, Anchoring Distribution.** We plot the mean and variance of the fitted anchoring distribution vs. the true feature distribution for each input dimension. We observe there is a mismatch between the empirical distribution and the fitted Gaussian. However, we did not find this mismatch to harm the effectiveness of G- Δ UQ.

Table 19: **Number of Parameters per Model.** We provide the number of parameters in the vanilla and modified parameter as follows. Note, that the change in parameters is architecture and input dimension dependent. For example, GOODCMNIST, and GOODSST2 use GIN MPNN layers. Therefore, when changing the layer dimension, we are changing the dimension of its internal MLP. It is not an error that intermediate layer G- Δ UQ have the same number of parameters, this is due to the architecture: these layers are the same size in the vanilla model. Likewise, GOODCora’s input features have dimension is 8701, so doubling the input layer’s dimension appears to add a significant number of parameters. We do not believe this

Dataset	GOODCMNIST	GOODMotif	GOODSST2	GOODCORA	GOODWebKB	GOODCBAS
Baseline	2001310	911403	1732201	2816770	695105	185104
G- Δ UQ(NFA)	2003110	913803	2193001	5429770	1206005	186304
G- Δ UQ(L1)	2360110	1633203	2091001			
G- Δ UQ(L2)	2360110	1633203	2091001			
G- Δ UQ(L3)	2360110					
G- Δ UQ(L4)	2360110					
G- Δ UQ(Readout)	2004310	912303	1732501			

A.15 EXPANDED DISCUSSION ON ANCHORING DESIGN CHOICES

Below, we expand upon some of the design choices for the proposed anchoring strategies.

When performing node featuring anchoring, how does fitting a Gaussian distribution to the input node features help manage the combinatorial stochasticity induced by message passing?

Without loss of generality, consider a node classification setting, where every sample is assigned a unique anchor. Then, due to message passing, after l hops, a given node's representation will have aggregated information from its l hop neighborhood. However, since each node in this neighborhood has a unique anchor, we see that any given node's representation is not only stochastic due to its own anchor but also that of its neighbors. For example, if any of its neighbors are assigned a different anchor, then the given node's representation will change, even if its own anchor did not. Since this behavior holds true for all nodes and each of their respective neighborhoods, we loosely refer to this phenomenon having combinatorial complexity, as effectively marginalizing out the anchoring distribution would require handling any and all changes to all l -hop neighbors. In contrast, when performing anchored image classification, the representation of a sample is only dependent on its unique, corresponding anchor, and is not influenced by the anchors of other samples. To this end, using the fitted Gaussian distribution helps manage this complexity, since changes to the anchors of a node's l -hop neighborhood are simpler to model as they require only learning to marginalize out a Gaussian distribution (instead of the training distribution). Indeed, for example, if we were to assume a simplified model where message passing only summed node neighbors, the anchoring distribution would remain gaussian after l rounds of message passing since the sum of gaussians is still gaussian (the exact parameters of the distribution would depend on the normalization used however).

REFERENCES

- Uri Alon and Eran Yahav. On the bottleneck of graph neural networks and its practical implications. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2021.
- Rushil Anirudh and Jayaraman J. Thiagarajan. Out of distribution detection via neural network anchoring. In *Asian Conference on Machine Learning, ACML 2022, 12-14 December 2022, Hyderabad, India, 2022*.
- Beatrice Bevilacqua, Yangze Zhou, and Bruno Ribeiro. Size-invariant graph representations for graph classification extrapolations. In *Proc. Int. Conf. on Machine Learning (ICML)*, 2021.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *Proc. Int. Conf. on Machine Learning (ICML)*, 2015.
- Davide Buffelli, Pietro Liò, and Fabio Vandin. Sizeshiftreg: a regularization method for improving size-generalization in graph neural networks. In *Proc. Adv. in Neural Information Processing Systems (NeurIPS)*, 2022.
- Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Velickovic. Principal neighbourhood aggregation for graph nets. In *NeurIPS*, 2020.
- Nicki Skaftø Detlefsen, Jiri Borovec, Justus Schock, Ananya Harsh, Teddy Koker, Luca Di Liello, Daniel Stancl, Changsheng Quan, Maxim Grechkin, and William Falcon. Torchmetrics - measuring reproducibility in pytorch, 2022. URL <https://github.com/Lightning-AI/torchmetrics>.
- Mucong Ding, Kezhi Kong, Jiuhai Chen, John Kirchenbauer, Micah Goldblum, David Wipf, Furong Huang, and Tom Goldstein. A closer look at distribution shifts and out-of-distribution generalization on graphs. In *NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications*, 2021.
- Vijay Prakash Dwivedi, Chaitanya K. Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks. *CoRR*, 2020.
- Vijay Prakash Dwivedi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Graph neural networks with learnable structural and positional representations. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2022a.
- Vijay Prakash Dwivedi, Ladislav Rampásek, Michael Galkin, Ali Parviz, Guy Wolf, Anh Tuan Luu, and Dominique Beaini. Long range graph benchmark. In *Proc. Adv. in Neural Information Processing Systems NeurIPS, Datasets and Benchmark Track*, 2022b.
- Saurabh Garg, Sivaraman Balakrishnan, Zachary C. Lipton, Behnam Neyshabur, and Hanie Sedghi. Leveraging unlabeled data to predict out-of-distribution performance. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2022.
- Thomas Gaudelot, Ben Day, Arian R. Jamasb, Jyothish Soman, Cristian Regep, Gertrude Liu, Jeremy B. R. Hayter, Richard Vickers, Charles Roberts, Jian Tang, David Roblin, Tom L. Blundell, Michael M. Bronstein, and Jake P. Taylor-King. Utilising graph machine learning within drug discovery and development. *CoRR*, abs/2012.05716, 2020.
- Shurui Gui, Xiner Li, Limei Wang, and Shuiwang Ji. GOOD: A graph out-of-distribution benchmark. In *Proc. Adv. in Neural Information Processing Systems (NeurIPS), Benchmark Track*, 2022.
- Devin Guillory, Vaishaal Shankar, Sayna Ebrahimi, Trevor Darrell, and Ludwig Schmidt. Predicting with confidence on unseen distributions. In *ICCV*, 2021.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *Proc. of the Int. Conf. on Machine Learning, (ICML)*, 2017.
- Kartik Gupta, Amir Rahimi, Thalaisyasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu, and Richard Hartley. Calibration of neural networks using splines. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2021.

- Arman Hasanzadeh, Ehsan Hajiramezani, Shahin Boluki, Mingyuan Zhou, Nick Duffield, Krishna Narayanan, and Xiaoning Qian. Bayesian graph neural networks with adaptive connection sampling. In *ICML*, 2020.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2017.
- Dan Hendrycks, Mantas Mazeika, and Thomas G. Dietterich. Deep anomaly detection with outlier exposure. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2019.
- Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in ML safety. *CoRR*, abs/2109.13916, 2021.
- Dan Hendrycks, Steven Basart, Mantas Mazeika, Andy Zou, Joseph Kwon, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Song. Scaling out-of-distribution detection for real-world settings. In *Proc. Int. Conf. on Machine Learning (ICML)*, 2022a.
- Dan Hendrycks, Andy Zou, Mantas Mazeika, Leonard Tang, Bo Li, Dawn Song, and Jacob Steinhardt. Pixmix: Dreamlike pictures comprehensively improve safety measures. In *Proc. Int. Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022b.
- José Miguel Hernández-Lobato and Ryan P. Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *Proc. Int. Conf. on Machine Learning (ICML)*, 2015.
- Hans Hao-Hsun Hsu, Yuesong Shen, Christian Tomani, and Daniel Cremers. What makes graph neural networks miscalibrated? In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2022.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: an introduction to concepts and methods. *Mach. Learn.*, 110(3), 2021.
- Yiding Jiang, Dilip Krishnan, Hossein Mobahi, and Samy Bengio. Predicting the generalization gap in deep networks with margin distributions. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- Konstantin Kirchheim, Marco Filax, and Frank Ortmeier. Pytorch-ood: A library for out-of-distribution detection based on pytorch. In *Workshop at the Proc. Int. Conf. on Computer Vision and Pattern Recognition CVPR*, 2022.
- Boris Knyazev, Graham W. Taylor, and Mohamed R. Amer. Understanding attention and generalization in graph neural networks. In *Proc. Adv. in Neural Information Processing Systems (NeurIPS)*, 2019.
- Ranganath Krishnan and Omesh Tickoo. Improving model calibration with accuracy versus uncertainty optimization. 2020.
- Meelis Kull, Miquel Perelló-Nieto, Markus Kängsepp, Telmo de Menezes e Silva Filho, Hao Song, and Peter A. Flach. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with dirichlet calibration. In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2019.
- Ananya Kumar, Percy Liang, and Tengyu Ma. Verified uncertainty calibration. In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2019.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Proc. Adv. in Neural Information Processing Systems (NeurIPS)*, 2017.
- Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2018.

- Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2020.
- Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. In *Proc. Adv. in Neural Information Processing Systems (NeurIPS)*, 2021.
- Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In *ICML 2020 Workshop on Graph Representation Learning and Beyond (GRL+ 2020)*, 2020. URL www.graphlearning.io.
- Luis Müller, Mikhail Galkin, Christopher Morris, and Ladislav Rampásek. Attending to graph transformers. *CoRR*, abs/2302.04181, 2023.
- Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proc. Conf. on Adv. of Artificial Intelligence (AAAI)*, 2015.
- Aviv Netanyahu, Abhishek Gupta, Max Simchowitz, Kaiqing Zhang, and Pulkit Agrawal. Learning to extrapolate: A transductive approach. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2023.
- Nathan Ng, Neha Hulkund, Kyunghyun Cho, and Marzyeh Ghassemi. Predicting out-of-domain generalization with local manifold smoothness. *CoRR*, abs/2207.02093, 2022.
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2019.
- Amir Rahimi, Amirreza Shaban, Ching-An Cheng, Richard Hartley, and Byron Boots. Intra order-preserving functions for calibration of multi-class neural networks. *Advances in Neural Information Processing Systems*, 33:13456–13467, 2020.
- Mrinank Sharma, Sebastian Farquhar, Eric Nalisnick, and Tom Rainforth. Do bayesian neural networks need to be fully stochastic? In *AISTATS*, 2023.
- Jayaraman J. Thiagarajan, Rushil Anirudh, Vivek Narayanaswamy, and Peer-Timo Bremer. Single model uncertainty estimation via stochastic data centering. In *Proc. Adv. in Neural Information Processing Systems (NeurIPS)*, 2022.
- Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M. Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature. In *Proc. Int. Conf. on Learning Representations ICLR*, 2022.
- Puja Trivedi, Danai Koutra, and Jayaraman J. Thiagarajan. A closer look at scoring functions and generalization prediction. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023a.
- Puja Trivedi, Danai Koutra, and Jayaraman J. Thiagarajan. A closer look at model adaptation using feature distortion and simplicity bias. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2023b.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018.
- Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-logit matching. In *Proc. Int. Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022a.
- Haorui Wang, Haoteng Yin, Muhan Zhang, and Pan Li. Equivariant and stable positional encoding for more powerful graph neural networks. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2022b.

- Xiao Wang, Hongrui Liu, Chuan Shi, and Cheng Yang. Be confident! towards trustworthy graph neural networks via confidence calibration. In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2021.
- Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre-Alvise Rebuffi, Ira Ktena, Krishnamurthy Dj Dvijotham, and Ali Taylan Cemgil. A Fine-Grained Analysis on Distribution Shift. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2022.
- Andrew Gordon Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. In *Proc. Adv. in Neural Information Processing Systems NeurIPS*, 2020.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *ICLR*, 2019.
- Gilad Yehudai, Ethan Fetaya, Eli Meiron, Gal Chechik, and Haggai Maron. From local structures to size generalization in graph neural networks. In *International Conference on Machine Learning*, pp. 11975–11986. PMLR, 2021.
- Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 694–699, 2002.
- Jize Zhang, Bhavya Kailkhura, and Thomas Yong-Jin Han. Mix-n-match : Ensemble and compositional methods for uncertainty calibration in deep learning. In *Proc. Int. Conf. on Machine Learning (ICML)*, 2020.
- Yingxue Zhang, Soumyasundar Pal, Mark Coates, and Deniz Üstebay. Bayesian graph convolutional neural networks for semi-supervised classification. In *AAAI*, 2019.
- Yanqiao Zhu, Yuanqi Du, Yinkai Wang, Yichen Xu, Jieyu Zhang, Qiang Liu, and Shu Wu. A survey on deep graph generation: Methods and applications. In *Learning on Graphs Conference (LoG)*, 2022.