

BENCHMARKING UNCERTAINTY QUANTIFICATION FOR PROTEIN ENGINEERING

Kevin P. Greenman

Department of Chemical Engineering
Massachusetts Institute of Technology
Cambridge, MA 02139, USA
kpg@mit.edu

Ava P. Soleimany & Kevin K. Yang

Microsoft Research New England
Cambridge, MA 02142, USA
{avasoleimany, yang.kevin}@microsoft.com

ABSTRACT

Machine learning sequence-function models for proteins could enable significant advances in protein engineering, especially when paired with state-of-the-art methods to select new sequences for property optimization and/or model improvement. Such methods (Bayesian optimization and active learning) require calibrated estimations of model uncertainty. While studies have benchmarked a variety of deep learning uncertainty quantification (UQ) methods on standard and molecular machine-learning datasets, it is not clear how well these results extend to protein datasets. In this work, we implement a panel of deep learning UQ methods on the Fitness Landscape Inference for Proteins (FLIP) benchmark regression tasks. We compare results across different degrees of distributional shift using metrics that assess each UQ method’s accuracy, calibration, coverage, width, and rank correlation to provide recommendations for the effective design of biological sequences.

1 INTRODUCTION

Machine learning (ML) has already begun to accelerate the field of protein engineering by providing low-cost predictions of phenomena that require time- and resource-intensive labeling by experiments or physics-based simulations (Yang et al., 2019). It is often necessary to have an estimate of model uncertainty in addition to the property prediction, as the performance of an ML model can be highly dependent on the domain shift between its training and testing data (Kendall & Gal, 2017). Because protein engineering data is often collected in a biased manner (Dallago et al., 2021), tailored ML methods are required to guide the selection of new experiments from a protein landscape. Uncertainty quantification (UQ) can inform the selection of experiments in order to improve a ML model or optimize protein function through active learning or Bayesian optimization.

In chemistry and materials science, several studies have benchmarked common UQ methods against one another on standard datasets and have used or developed appropriate metrics to quantify the quality of these uncertainty estimates (Scalia et al., 2020; Tran et al., 2020; Hirschfeld et al., 2020; Nigam et al., 2021; Soleimany et al., 2021). This work has illustrated that the best choice of UQ method can be dependent on the dataset and other considerations such as scaling.

While some protein engineering work has leveraged uncertainty estimates, these studies have been limited to single UQ methods such as ensembles (Mariet et al., 2020) or Gaussian processes (GPs) (Hie et al., 2020). In this work, we use a group of standardized, public protein datasets to benchmark a panel of UQ methods. Our chosen datasets include splits with varied degrees of domain extrapolation, enabling us to evaluate the methods in a setting similar to what might be experienced while collecting new experimental data. We assess each model using a variety of metrics that capture different aspects of desired performance, including accuracy, calibration, coverage, width, and rank correlation.

2 METHODS

This work uses several dataset splits to evaluate the performance of a panel of uncertainty methods across various evaluation metrics (Figure 1).

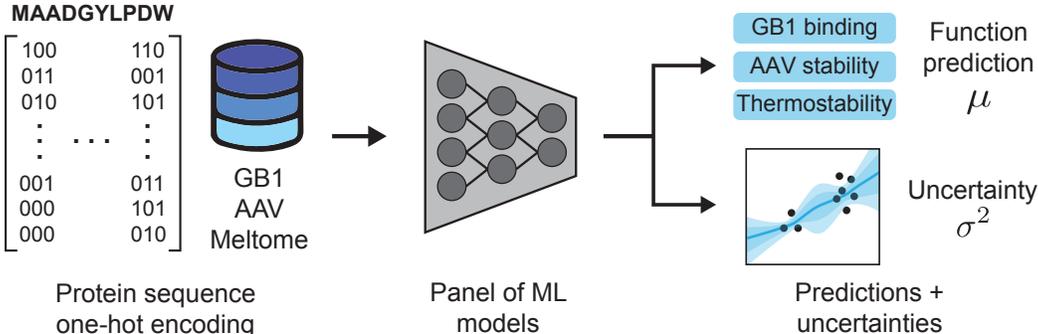


Figure 1: Overall schematic.

2.1 DATASETS AND TASKS

The landscapes used in this work were taken from the Fitness Landscape Inference for Proteins (FLIP) benchmark (Dallago et al., 2021). These included the binding domain of an immunoglobulin binding protein (GB1), adeno-associated virus stability (AAV), and thermostability (Meltome), which cover a large sequence space and a broad range of protein families. The FLIP benchmark includes several train-test splits, or tasks, for each landscape; most of these tasks are designed to mimic common, real-world data collection scenarios and are thus a more realistic assessment of generalizability than random train-test splits. However, random splits are also included as a point of reference. We chose 8 of the 15 FLIP tasks to benchmark the panel of uncertainty methods. We selected these tasks to be representative of several regimes of domain shift; they include random sampling with no domain shift (AAV/sampled, Meltome/mixed, and GB1/sampled); the highest (and most relevant) domain-shift regimes (AAV/Mut-Des, AAV/7 vs. Many, and GB1/1 vs. Rest); and less aggressive domain shifts (GB1/2 vs. Rest and GB1/3 vs. Rest).

2.2 EVALUATION METRICS

To give a comprehensive report of model accuracy, we compute the following metrics on the test sets: root mean square error (RMSE), mean absolute error (MAE), coefficient of determination (R^2), and Spearman rank correlation (ρ). RMSE is more sensitive to outliers than MAE, so while both are informative independently, the combination of the two gives additional information about the distribution of errors. R^2 and ρ are both unitless and are thus more easily interpreted and compared across datasets.

We evaluated the quality of the uncertainty estimates using four metrics. First, ρ_{unc} is the Spearman rank correlation between uncertainty and absolute prediction error. Following Kompa et al. (2021), we measure the coverage as the percentage of true values that fall within the 95% confidence interval ($\pm 2\sigma$) of each prediction. Kompa et al. (2021) define the width as the size of the 95% confidence region (4σ), but we calculate the width relative to the range of the training set as $4\sigma/R$ to make these values more interpretable across datasets. Finally, the miscalibration area (also called the area under the calibration error curve or AUCE) quantifies the absolute difference between the calibration plot and perfect calibration in a single number (Gustafsson et al., 2020).

2.3 UNCERTAINTY METHODS

We implemented eight uncertainty methods for this benchmark: linear Bayesian ridge regression (BRR), Gaussian processes (GPs), and six methods using variations on a convolutional neural network (CNN) architecture. We used a `scikit-learn` implementation of BRR (Pedregosa et al.,

2011) and a GPyTorch implementation of continuous-kernel GPs (Gardner et al., 2018). The CNN implementation from FLIP provided the core architecture used by our dropout (Gal & Ghahramani, 2016), ensemble (Lakshminarayanan et al., 2017), evidential (Amini et al., 2020), mean-variance estimation (MVE) (Nix & Weigend, 1994), and last-layer stochastic variational inference (SVI) (Hoffman et al., 2013) models. For each dropout model, we tested dropout fractions of 0.1, 0.2, 0.3, 0.4, and 0.5 and reported the model with the lowest miscalibration area. The output of the evidential uncertainty can be divided into an epistemic and aleatoric component following the analysis of Amini et al. (2020), so we report these uncertainties separately in Tables 1-4 along with their sum as the total uncertainty in Figures 2-4. For all models and landscapes, the sequences were featurized as one-hot encodings.

3 RESULTS AND DISCUSSION

We trained each of the seven models described in Section 2.3 on each of the eight tasks described in Section 2.1 and evaluated their performance on the test set using the metrics described in Section 2.2. We compare model calibration and accuracy in Figure 2 and the percent coverage vs. average width relative to range in Figure 3.

As expected, the splits with the least required domain extrapolation tend to have more accurate models (lower RMSE). However, the relationship between miscalibration area and extrapolation is less clear; some models were highly calibrated on the most rigorous (highest domain shift) splits, while others were poorly calibrated even on random splits. There is no single method that performs consistently well across splits and landscapes, but some trends can be observed. For example, the CNN ensemble is often one of the highest accuracy models, but also one of the most poorly calibrated.

Figure 3 illustrates that many methods perform relatively well in either coverage or width (corresponding to the the top and left limits of the plot, respectively), but few methods perform well in both. Among those that do, the GP model is the only method that is among the best across all landscapes and most splits. Similarly to Figure 2, there is some observable trend that more challenging splits are further from the optimal part (upper left) of the plot; this trend is more clear for the GB1 splits than for the AAV splits.

Figure 4 compares the ranking performance of each method in terms of predictions relative to true values and uncertainty estimates relative to true errors. The splits are ordered according to domain shift within their respective landscapes, and the rank correlation of the predictions to the true labels generally decreases moving from left to right within a landscape. Ensembling was often among the best performance in ρ , and the GP and BRR models performed better on ρ on the splits with the highest domain shift. Performance on ρ_{unc} is generally much worse than that on ρ , with a large number of results showing negative correlation. Dropout has one of the highest and most consistent ρ_{unc} across splits. All methods have ρ_{unc} near zero for the most challenging splits. MVE performs particularly poorly in regimes of high domain shift, which is consistent with its intended use as an estimator of aleatoric (data-dependent) uncertainty.

4 CONCLUSION

Calibrated uncertainty estimations are necessary for effective property optimization or model improvement using Bayesian optimization or active learning. In this work, we have benchmarked a panel of uncertainty quantification (UQ) methods on protein datasets, including on train-test splits that are representative of real-world data collection practices. After evaluating each method based on accuracy, calibration, coverage, width, and rank correlation, there is no method that performs consistently well across all metrics or all landscapes and splits. However, this study only examines models using one-hot-encoding representations of sequences; further research is needed to understand how more informative and generalizable representations such as the pretrained embeddings from (Dallago et al., 2021) may impact performance, particularly on splits with higher domain shifts. Additionally, while the evaluation metrics used herein provide valuable information, they are ultimately only a proxy for expected performance in Bayesian optimization and active learning. Retrospective studies using holdouts of the training sets would be useful for evaluating this performance more directly. A more thorough understanding of how to best apply UQ to sequence-function models will ultimately enable more effective protein engineering.

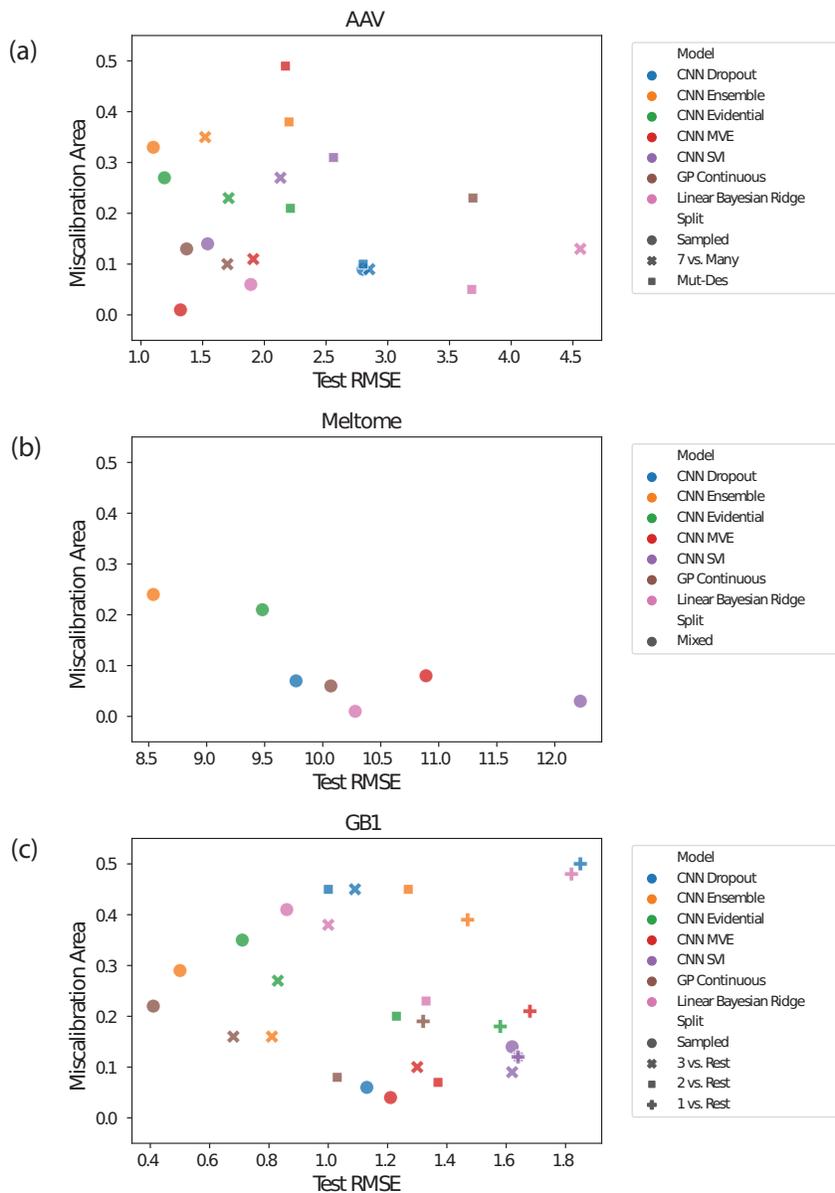


Figure 2: Miscalibration area vs. root mean square error (RMSE) for the (a) AAV, (b) Meltome, and (c) GB1 landscapes. Miscalibration area (also called the area under the calibration error curve or AUCE) quantifies the absolute difference between the calibration plot and perfect calibration. It is desirable to have a model that is both accurate and well-calibrated, so the best performing points are those closest to the lower left corner of the plots.

AUTHOR CONTRIBUTIONS

A.P.S. and K.K.Y. conceived the project. K. P. G. wrote the computer code, analyzed the data, and wrote the first manuscript draft. A.P.S. and K.K.Y. supervised the research and edited the manuscript.

ACKNOWLEDGMENTS

K.P.G. was supported by a Microsoft Research micro-internship and by the National Science Foundation Graduate Research Fellowship Program under Grant No. 1745302. We acknowledge the MIT Lincoln Laboratory Supercloud cluster (Reuther et al., 2018) at the Massachusetts Green High

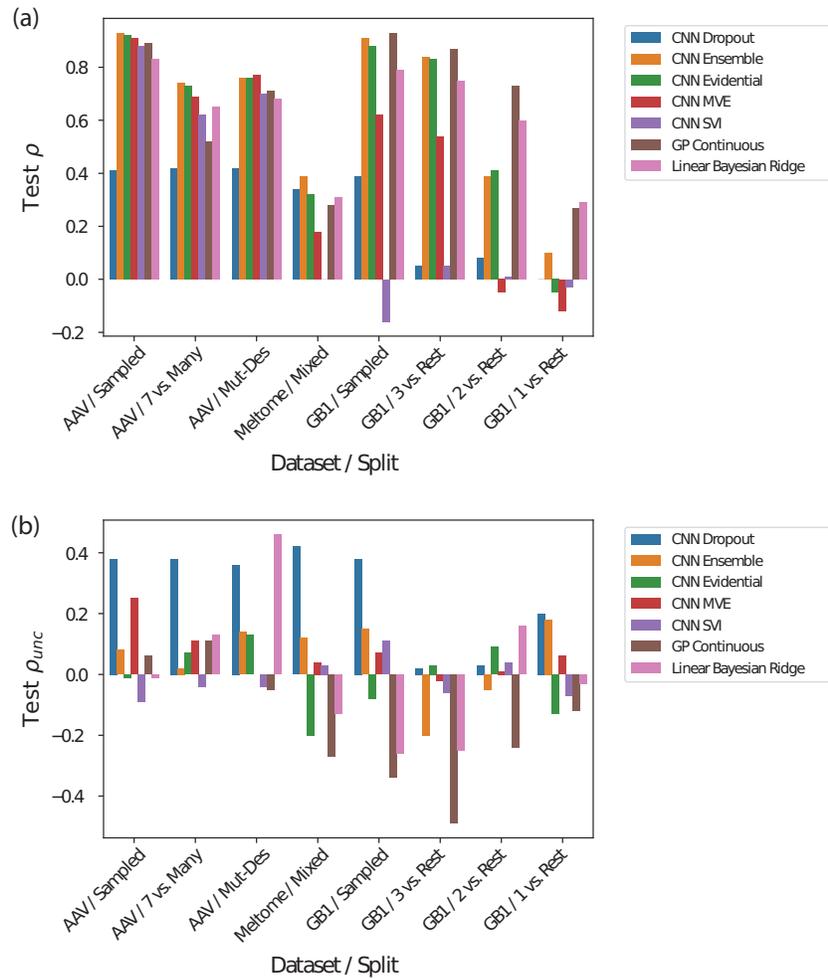


Figure 4: Spearman rank correlations of (a) predictions (ρ) and (b) uncertainties (ρ_{unc}) vs. extrapolation. Within each landscape (AAV, Meltome, and GB1), splits are qualitatively ordered by the amount of domain shift between train and test sets, with the lowest domain shift on the left and highest shift on the right.

REFERENCES

- Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. Deep evidential regression. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 14927–14937. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/aab085461de182608ee9f607f3f7d18f-Paper.pdf>.
- Christian Dallago, Jody Mou, Kadina E Johnston, Bruce Wittmann, Nick Bhattacharya, Samuel Goldman, Ali Madani, and Kevin K Yang. FLIP: Benchmark tasks in fitness landscape inference for proteins. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL <https://openreview.net/forum?id=p2dMLEwL8tF>.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/gall16.html>.
- Jacob Gardner, Geoff Pleiss, Kilian Q Weinberger, David Bindel, and Andrew G Wilson. Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration. *Advances in neural information processing systems*, 31, 2018.
- Fredrik K Gustafsson, Martin Danelljan, and Thomas B Schon. Evaluating scalable bayesian deep learning methods for robust computer vision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 318–319, 2020.
- Brian Hie, Bryan D. Bryson, and Bonnie Berger. Leveraging uncertainty in machine learning accelerates biological discovery and design. *Cell Systems*, 11(5):461–477.e9, 2020. ISSN 2405-4712. doi: <https://doi.org/10.1016/j.cels.2020.09.007>. URL <https://www.sciencedirect.com/science/article/pii/S2405471220303641>.
- Lior Hirschfeld, Kyle Swanson, Kevin Yang, Regina Barzilay, and Connor W Coley. Uncertainty quantification using neural networks for molecular property prediction. *Journal of Chemical Information and Modeling*, 60(8):3770–3780, 2020.
- Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 2013.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in Neural Information Processing Systems*, 30, 2017.
- Benjamin Kompa, Jasper Snoek, and Andrew L. Beam. Empirical frequentist coverage of deep learning uncertainty quantification procedures. *Entropy*, 23(12), 2021. ISSN 1099-4300. doi: 10.3390/e23121608. URL <https://www.mdpi.com/1099-4300/23/12/1608>.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf>.
- Zelda Mariet, Ghassen Jerfel, Zi Wang, Christof Angermüller, David Belanger, Suhani Vora, Maxwell Bileschi, Lucy Colwell, D Sculley, Dustin Tran, et al. Deep uncertainty and the search for proteins. In *Workshop: Machine Learning for Molecules*, 2020.
- AkshatKumar Nigam, Robert Pollice, Matthew FD Hurley, Riley J Hickman, Matteo Aldeghi, Naruki Yoshikawa, Seyone Chithrananda, Vincent A Voelz, and Alán Aspuru-Guzik. Assigning confidence to molecular property prediction. *Expert opinion on drug discovery*, 16(9):1009–1023, 2021.

- David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability distribution. In *Proceedings of 1994 IEEE international conference on neural networks (ICNN'94)*, volume 1, pp. 55–60. IEEE, 1994.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Albert Reuther, Jeremy Kepner, Chansup Byun, Siddharth Samsi, William Arcand, David Bestor, Bill Bergeron, Vijay Gadepally, Michael Houle, Matthew Hubbell, Michael Jones, Anna Klein, Lauren Milechin, Julia Mullen, Andrew Prout, Antonio Rosa, Charles Yee, and Peter Michaleas. Interactive supercomputing on 40,000 cores for machine learning and data analysis. In *2018 IEEE High Performance Extreme Computing Conference (HPEC)*, pp. 1–6. IEEE, 2018.
- Gabriele Scalia, Colin A Grambow, Barbara Pernici, Yi-Pei Li, and William H Green. Evaluating scalable uncertainty estimation methods for deep learning-based molecular property prediction. *Journal of chemical information and modeling*, 60(6):2697–2717, 2020.
- Ava P Soleimany, Alexander Amini, Samuel Goldman, Daniela Rus, Sangeeta N Bhatia, and Connor W Coley. Evidential deep learning for guided molecular property prediction and discovery. *ACS central science*, 7(8):1356–1367, 2021.
- Kevin Tran, Willie Neiswanger, Junwoong Yoon, Qingyang Zhang, Eric Xing, and Zachary W Ulissi. Methods for comparing uncertainty quantifications for material property predictions. *Machine Learning: Science and Technology*, 1(2):025006, may 2020. doi: 10.1088/2632-2153/ab7e1a. URL <https://doi.org/10.1088/2632-2153/ab7e1a>.
- Kevin K Yang, Zachary Wu, and Frances H Arnold. Machine-learning-guided directed evolution for protein engineering. *Nature methods*, 16(8):687–694, 2019.

A APPENDIX

Table 1: RMSE (\downarrow) / MAE (\downarrow) / R^2 (\uparrow)

Dataset (Split)	BRR	Ens.	MVE	Dropout	GP	Evi. (a)	Evi. (e)	SVI
AAV (sampled)	1.89	1.10	1.32	2.80	1.37	1.19	1.19	1.54
	1.50	0.83	1.01	2.37	1.03	0.89	0.89	1.18
	0.62	0.87	0.82	0.18	0.80	0.85	0.85	0.75
AAV (7-vs-many)	4.56	1.52	1.91	2.85	1.70	1.71	1.71	2.13
	3.30	1.19	1.53	2.40	1.33	1.33	1.33	1.76
	-4.02	0.44	0.12	0.17	0.30	0.29	0.29	-0.10
AAV (mut-des)	3.68	2.20	2.17	2.80	3.69	2.21	2.21	2.56
	2.58	1.68	1.68	2.38	3.15	1.69	1.69	2.11
	-0.12	0.60	0.61	0.17	-0.13	0.60	0.60	0.46
Meltome (mixed)	10.28	8.54	10.89	9.77	10.07	9.48	9.48	12.22
	7.88	6.55	7.88	6.89	7.49	7.12	7.12	8.79
	0.22	0.46	0.12	0.15	0.25	0.34	0.34	-0.11
GB1 (sampled)	0.86	0.50	1.21	1.13	0.41	0.71	0.71	1.62
	0.65	0.34	0.92	0.84	0.27	0.43	0.43	1.15
	0.50	0.83	0.00	0.15	0.88	0.65	0.65	-0.80
GB1 (3-vs-rest)	1.00	0.81	1.30	1.09	0.68	0.83	0.83	1.62
	0.77	0.58	1.06	0.83	0.45	0.54	0.54	1.08
	0.38	0.60	-0.03	0.00	0.71	0.57	0.57	-0.61
GB1 (2-vs-rest)	1.33	1.27	1.37	1.00	1.03	1.23	1.23	1.64
	1.04	1.05	0.94	0.75	0.68	0.93	0.93	1.13
	-0.17	-0.07	-0.24	0.00	0.30	0.01	0.01	-0.78
GB1 (1-vs-rest)	1.82	1.47	1.68	1.85	1.32	1.58	1.58	1.64
	1.53	1.26	1.16	1.61	1.13	1.07	1.07	1.15
	-1.23	-0.17	-0.90	-3.10	-0.17	-0.67	-0.67	-0.80

Table 2: ρ (\uparrow) / ρ_{unc} (\uparrow)

Dataset (Split)	BRR	Ens.	MVE	Dropout	GP	Evi. (a)	Evi. (e)	SVI
AAV (sampled)	0.83	0.93	0.91	0.41	0.89	0.92	0.92	0.88
	-0.01	0.08	0.25	0.38	0.06	-0.01	-0.01	-0.09
AAV (7-vs-many)	0.65	0.74	0.69	0.42	0.52	0.73	0.73	0.62
	0.13	0.02	0.11	0.38	0.11	0.07	0.07	-0.04
AAV (mut-des)	0.68	0.76	0.77	0.42	0.71	0.76	0.76	0.70
	0.46	0.14	-	0.36	-0.05	0.14	0.13	-0.04
Meltome (mixed)	0.31	0.39	0.18	0.34	0.28	0.32	0.32	0.00
	-0.13	0.12	0.04	0.42	-0.27	-0.20	-0.20	0.03
GB1 (sampled)	0.79	0.91	0.62	0.39	0.93	0.88	0.88	-0.16
	-0.26	0.15	0.07	0.38	-0.34	-0.06	-0.08	0.11
GB1 (3-vs-rest)	0.75	0.84	0.54	0.05	0.87	0.83	0.83	0.05
	-0.25	-0.20	-0.02	0.02	-0.49	0.04	0.03	-0.06
GB1 (2-vs-rest)	0.60	0.39	-0.05	0.08	0.73	0.41	0.41	0.01
	0.16	-0.05	0.01	0.03	-0.24	0.09	0.08	0.04
GB1 (1-vs-rest)	0.29	0.10	-0.12	0.00	0.27	-0.05	-0.05	-0.03
	-0.03	0.18	0.06	0.20	-0.12	-0.13	-0.13	-0.07

Table 3: Coverage (\uparrow) / ($4\sigma/R$) (\downarrow)

Dataset (Split)	BRR	Ens.	MVE	Dropout	GP	Evi. (a)	Evi. (e)	SVI
AAV (sampled)	0.98 0.02	0.40 0.00	0.94 0.01	0.95 0.03	0.98 0.02	0.88 0.18	0.88 0.18	0.74 0.01
AAV (7-vs-many)	0.83 0.03	0.35 0.00	0.83 0.01	0.95 0.02	0.99 0.02	0.83 0.14	0.83 0.14	0.56 0.01
AAV (mut-des)	0.85 0.02	0.30 0.00	0.69 0.01	0.94 0.02	0.59 0.02	0.79 0.15	0.79 0.15	0.46 0.01
Meltome (mixed)	0.92 0.01	0.61 0.00	0.84 0.01	0.86 0.01	0.94 0.01	0.74 0.89	0.73 0.89	0.91 0.01
GB1 (sampled)	1.00 0.30	0.49 0.01	0.94 0.06	0.89 0.04	0.98 0.04	0.93 0.27	0.93 0.27	0.71 0.04
GB1 (3-vs-rest)	1.00 0.53	0.64 0.04	0.93 0.11	0.11 0.01	0.96 0.09	0.87 0.30	0.87 0.32	0.74 0.08
GB1 (2-vs-rest)	1.00 0.39	0.13 0.02	0.88 0.15	0.11 0.01	0.89 0.12	0.93 0.19	0.93 0.19	0.71 0.10
GB1 (1-vs-rest)	0.05 0.03	0.49 0.09	0.58 0.19	0.00 0.02	0.86 0.32	0.85 0.41	0.85 0.41	0.71 0.28

Table 4: Miscalibration Area (\downarrow)

Dataset (Split)	BRR	Ens.	MVE	Dropout	GP	Evi. (a)	Evi. (e)	SVI
AAV (sampled)	0.06	0.33	0.01	0.09	0.13	0.19	0.19	0.14
AAV (7-vs-many)	0.13	0.35	0.11	0.09	0.10	0.17	0.17	0.27
AAV (mut-des)	0.05	0.38	0.49	0.10	0.23	0.17	0.17	0.31
Meltome (mixed)	0.01	0.24	0.08	0.07	0.06	0.20	0.21	0.03
GB1 (sampled)	0.41	0.29	0.04	0.06	0.22	0.24	0.24	0.14
GB1 (3-vs-rest)	0.38	0.16	0.10	0.45	0.16	0.15	0.15	0.09
GB1 (2-vs-rest)	0.23	0.45	0.07	0.45	0.08	0.04	0.04	0.12
GB1 (1-vs-rest)	0.48	0.39	0.21	0.50	0.19	0.07	0.07	0.12