
Heteroscedastic Variational Last Layers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We present a simple, inexpensive, and effective method for heteroscedastic uncertainty
2 quantification in neural networks. We build on Variational Bayesian Last
3 Layers (VBLL), wherein deterministic training objectives are developed for varia-
4 tional inference of the network last layer. In particular, we (1) Introduce t-VBLL
5 layers, which perform variational inference for the aleatoric noise covariance, and
6 (2) Introduce Het-VBLL, a Bayesian last layer scheme to model heteroscedastic
7 noise. These methods are based on novel, analytically tractable evidence lower
8 bounds. We further discuss parameterization and initialization within these models.
9 We show that these novel design elements enable effective uncertainty modeling
10 at minimal additional cost, and substantially improve performance over similar
11 methods such as VBLLs.

12 1 Introduction

13 Effective uncertainty quantification (UQ) in neural networks is crucial for a wide range of downstream
14 applications, such as active learning, exploration, decision making and enhancing model robustness
15 (Neal, 1995; Kendall & Gal, 2017; Wilson & Izmailov, 2020). The scale of modern neural networks,
16 however, necessitates computationally inexpensive UQ methods, particularly when compared to
17 computationally intensive approaches like Markov Chain Monte Carlo (MCMC) (Izmailov et al.,
18 2021) or full-network Bayes-by-Backprop (Blundell et al., 2015). Among these inexpensive methods,
19 last-layer uncertainty methods—which maintain a distribution over the parameters of only the
20 final layer, rather than a point estimate—have emerged as a promising strategy, balancing strong
21 predictive performance with low computational overhead during training and inference (Daxberger
22 et al., 2021; Kristiadi et al., 2021; Harrison et al., 2018, 2024; Liu et al., 2022; Watson et al., 2020,
23 2021). Variational Bayesian Last Layers (VBLLs) in particular offer a lightweight, easy-to-train, and
24 effective approach to UQ (Harrison et al., 2024; Brunzema et al., 2025; Watson et al., 2021).

25 VBLLs replace point estimation of neural network last layers with variational inference, and exploit
26 analytically tractable evidence lower bounds (ELBOs) to yield an easy-to-train, inexpensive approach
27 to Bayesian deep learning. However, the quality of uncertainty quantification in VBLLs—in particular,
28 of epistemic (reducible) uncertainty (Der Kiureghian & Ditlevsen, 2009)—depends on a model of the
29 aleatoric (irreducible) noise. For standard VBLLs, a point estimate of the aleatoric noise variance
30 is computed via MAP estimation, due to the simplicity of this procedure in combination with
31 neural network training. Although the noise term may appear relatively unimportant, it is essential
32 for weighing the predictive accuracy and uncertainty terms, thus playing a vital role in accurate
33 uncertainty quantification.

34 In this work, we investigate and substantially improve the modeling of aleatoric uncertainty in
35 VBLLs. Specifically, we present two methods for more expressive noise modeling. First, we develop
36 t-VBLLs: a variational approach to noise inference that maintains a full variational posterior over
37 the noise variance. In the regression case, this leads to the familiar Student t-distributed predictive
38 distribution. For classification, we introduce a novel sampling-free lower bound on the standard

39 ELBO, enabling effective training with low variance. This approach allows us to characterize the
 40 epistemic uncertainty over the aleatoric noise, leading to a more complete model of uncertainty
 41 in last layer methods. Second, we introduce Het-VBLLs: a method for variational inference of
 42 heteroscedastic (input-dependent) noise (Le et al., 2005; Hayashi, 2011). This builds upon a second
 43 VBLL model to characterize noise and leverages the novel training ELBOs developed in the first part
 44 of the paper.

45 These approaches introduce only a handful of new parameters, making them highly scalable while
 46 substantially improving the calibration and uncertainty quantification of neural networks. In particular,
 47 our contributions are:

- 48 • We introduce t-VBLL and Het-VBLL, two novel last layer variational methods for the
 49 homoscedastic and heteroscedastic settings respectively. These are based on the development
 50 of novel lower bounds on the standard ELBOs which enable sampling-free training.
- 51 • We evaluate these models on a wide variety of supervised regression and classification,
 52 as well as Bayesian optimization settings. We show that our methods consistently result
 53 in better predictive likelihood, calibration, better active learning, and better robustness to
 54 outliers than baseline methods.
- 55 • We release easy-to-use implementations of both t-VBLL and Het-VBLL for the regression
 56 and classification setting (removed for anonymity).

57 2 Preliminaries

58 We parameterize a neural network by separating its last layer weights, W , from the feature extractor
 59 (all preceding layers), denoted by $\phi(\cdot)$. The network output, z , for an input x are modeled as

$$z = W\phi(x) + \varepsilon \quad (1)$$

60 where ε is zero-mean Gaussian noise. For regression tasks, the output is $y = z$. For classification¹,
 61 the class probabilities are given by $p(y | x) = \text{softmax}(z)$. We denote the dimensionalities of the
 62 inputs, outputs, and features as N_x , N_y , and N_ϕ , respectively.

63 Uncertainty quantification in neural networks has been an active area of research since the early
 64 1990s (MacKay, 1992; Nix & Weigend, 1994), gaining renewed attention with the field’s resurgence
 65 in the early 2010s (Blundell et al., 2015; Gal & Ghahramani, 2016; Lakshminarayanan et al., 2017;
 66 Papamarkou et al., 2024). Two prominent approaches have emerged. The first is variance prediction,
 67 where the model directly outputs predictive uncertainty, typically by making the noise term ε
 68 input-dependent (Nix & Weigend, 1994). The second involves quantifying uncertainty over the
 69 network parameters themselves, with Bayesian inference being the most common method (MacKay,
 70 1992; Blundell et al., 2015). Modern Bayesian neural networks typically use variational inference,
 71 optimizing a lower bound on the marginal likelihood by parameterizing distributions over the network
 72 weights.

73 In this work, we build on Variational Bayesian Last Layers (VBLLs), a minimal and low-variance
 74 approach to variational deep learning (Harrison et al., 2024). The VBLL method places a prior $p(W)$
 75 and a variational posterior $q(W)$ only on the weights of the final layer, W , while the weights of the
 76 feature extractor ϕ are learned as point estimates, as in standard neural network training.

77 Harrison et al. (2024) derive a tractable objective by lower-bounding the log marginal likelihood:

$$\log p(Y | X, \Sigma) \geq \sum_{(x,y) \in \mathcal{D}} \mathbb{E}_{q(W)}[\log p(y | x, W, \Sigma)] - \text{KL}(q(W) | p(W)) \quad (2)$$

78 This bound leverages the exact computation of the expected log-likelihood (or a lower bound thereof
 79 for classification), resulting in an inexpensive, sampling-free training objective. The model is trained
 80 by jointly optimizing the feature extractor weights and the parameters of the variational posterior
 81 $q(\cdot)$. Throughout the original work, the noise covariance $\varepsilon \sim \mathcal{N}(0, \Sigma)$ is assumed to have a known,
 82 fixed covariance Σ , or a point estimate for Σ is learned via maximum likelihood or MAP estimation.

83 This paper extends the VBLL framework to perform variational inference directly on the noise
 84 covariance Σ , drawing on techniques from the variance prediction literature. We introduce a factor-
 85 ized variational posterior over the diagonal elements of the noise covariance, $q(\Sigma) = \prod_i^{N_y} q(\Sigma_i)$.

¹Choosing noise identity or zero noise variance recovers standard loss functions for regression and classifica-
 tion, respectively.

86 We explore both homoscedastic posteriors of the form $q(\Sigma)$ and input-dependent heteroscedastic
87 posteriors of the form $q(\Sigma | \mathbf{x})$.

88 3 Related Work

89 While variance prediction is conceptually simple and capable of expressive representation of aleatoric
90 uncertainty, it is limited in its ability to quantify epistemic uncertainty (Der Kiureghian & Ditlevsen,
91 2009; Kendall & Gal, 2017). Similarly, many approaches toward epistemic uncertainty quantification
92 rely on assumed likelihoods that inaccurately model aleatoric uncertainty. Kendall & Gal (2017)
93 made the most direct case for the value of characterizing both aleatoric and epistemic uncertainty,
94 and introduced relatively simple models combining Bayesian dropout for epistemic uncertainty with
95 variance prediction for aleatoric.

96 3.1 Variance Prediction

97 Variance prediction networks aim to, beyond standard point predictions, predict variance terms to
98 capturing (aleatoric) uncertainty (Nix & Weigend, 1994; Bishop & Quazaz, 1996). In regression,
99 this typically corresponds to predicting the mean and variance of a Gaussian likelihood, whilst
100 in classification, this typically corresponds to modeling logit noise (Collier et al., 2020, 2021).
101 Concretely, these networks typically parameterize the variance via

$$\log \Sigma_i = \mathbf{m}_i^\top \phi(\mathbf{x}) \quad (3)$$

102 where \mathbf{m}_i is a (point estimate) last layer. These networks are typically trained via maximum
103 likelihood. Variance prediction networks broadly fall into the set of heteroscedastic models, although
104 there is often a conflation of epistemic and aleatoric uncertainty in the predicted variance (Kendall &
105 Gal, 2017). Skafte et al. (2019) proposed to improve variance prediction by (instead of predicting the
106 log variance) predicting the parameters of an inverse Gamma distribution, and marginalize over these
107 parameters to yield a Student-t predictive distribution. This parameterization has strong similarities to
108 the approaches in this work, and similarly aims to address the shortcomings of aleatoric uncertainty
109 prediction via epistemic noise modeling.

110 3.2 Epistemic Uncertainty and Bayesian Deep Learning

111 Variance prediction networks are limited in their ability to quantify epistemic uncertainty. Like
112 standard networks, predictive behavior far from data is hard to anticipate, and thus predictive
113 uncertainty may be smaller (or larger) than desired. Thus, epistemic uncertainty quantification—
114 even for models that characterize aleatoric uncertainty such as variance prediction models—are
115 necessary. Moreover, the posterior inferred by approximate Bayesian methods relies on the chosen
116 likelihood. Thus, model misspecification in the form of assumed homoscedasticity can substantially
117 harm posterior inference and the predictive performance of the model (Le et al., 2005; Kersting et al.,
118 2007; Lazaro-Gredilla & Titsias, 2011).

119 A wide variety of epistemic uncertainty quantification methods for neural networks have been
120 developed including variational methods (Blundell et al., 2015), Dropout-based methods (Gal &
121 Ghahramani, 2016), ensembling (Lakshminarayanan et al., 2017), randomized priors (Osband et al.,
122 2018, 2023), and many others. However, these methods are typically expensive and scale poorly
123 to larger models. In this work, we build on variational Bayesian last layers (VBLLs) as a scalable
124 and effective Bayesian approach (Harrison et al., 2024). These models fit into a class of similar last
125 layer uncertainty quantification approaches, including SNGP (Liu et al., 2022) and last layer Laplace
126 approximation methods (Daxberger et al., 2021). For these methods, the last layer inference strategy
127 relies on computing a last layer via approximate Bayesian linear regression after each epoch, and
128 thus incorporating variance prediction (and epistemic uncertainty for this term) is challenging.

129 4 Approach

130 In this section we introduce two models: t-VBLLs and Het-VBLLs (see toy experimental results in
131 Fig. 1). We first discuss the noise parameterization that defines each model. We introduce a set of
132 variational bounds used as training objectives, and these general results apply to both homoscedastic
133 and heteroscedastic models. A full discussion of the technical details is provided in the Appendix.

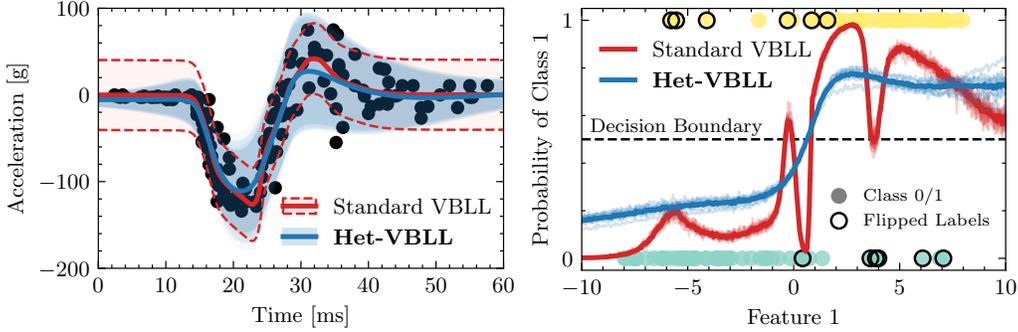


Figure 1: *Left*: Het-VBLLs vs. VBLLs on motorcycle-impact data from Silverman (1985). Het-VBLLs can capture the heteroscedasticity for expressive uncertainty estimates. *Right*: Classification on a toy sigmoid dataset with 100 points and 10% flipped labels. The Het-VBLL models demonstrate greater robustness to outliers.

134 4.1 Parameterizing the Noise Distribution

135 We consider independent priors (and variational posteriors) over rows of W and fix a diagonal Σ , and
 136 so it is sufficient to consider each row independently, of the form

$$z_i = \mathbf{w}_i^\top \phi(\mathbf{x}) + \varepsilon_i \quad (4)$$

137 with $\varepsilon_i \sim \mathcal{N}(0, \Sigma_i)$. We fix a prior over the last layer parameters for each row \mathbf{w}_i and the covariance
 138 element Σ_i

$$p(W, \Sigma | \mathbf{x}) = \prod_{i=1}^{N_y} p(\mathbf{w}_i | \Sigma_i, \mathbf{x}) p(\Sigma_i | \mathbf{x}) \quad (5)$$

139 with² $p(\mathbf{w}_i | \Sigma_i, \mathbf{x}) = \mathcal{N}(\bar{\mathbf{w}}_i, \Sigma_i(\mathbf{x})S_i)$. This corresponds to the same choice of prior as Harrison
 140 et al. (2024), although with the covariance parameterization being scaled by Σ_i . This scaling is
 141 standard in Bayesian linear regression to enable recursive updating of sufficient statistics of the
 142 posterior. Here, it yields easy-to-evaluate variational objectives. We similarly structure our variational
 143 posterior as

$$q(W, \Sigma | \mathbf{x}) = \prod_{i=1}^{N_y} q(\mathbf{w}_i | \Sigma_i, \mathbf{x}) q(\Sigma_i | \mathbf{x}) \quad (6)$$

144 following standard results in Bayesian regression, with $q(\mathbf{w}_i | \Sigma_i) = \mathcal{N}(\mathbf{w}_i, \Sigma_i(\mathbf{x})S_i)$.

145 We will consider two variational families for the noise covariance, corresponding to the homoscedastic
 146 case (in which Σ does not depend on \mathbf{x}) and the heteroscedastic case respectively.

147 **Homoscedastic.** We fix an inverse Gamma variational posterior. This is the canonical (conjugate)
 148 prior for the noise covariance in Bayesian linear regression. This choice of prior results in a Student
 149 t-distributed posterior predictive distribution (Box & Tiao, 2011), which we exploit in our prediction,
 150 giving rise to t-VBLLs.

151 **Heteroscedastic.** In the heteroscedastic setting, we parameterize the noise covariance with a VBLL
 152 as

$$\log \Sigma_i = \mathbf{m}_i^\top \phi(\mathbf{x}) \quad (7)$$

153 with $\mathbf{m}_i \sim q(\mathbf{m}_i) = \mathcal{N}(\bar{\mathbf{m}}_i, Z_i)$ (and similarly choose as prior $\mathbf{m}_i \sim \mathcal{N}(\bar{\mathbf{m}}_i, Z_i)$). This yields a
 154 log-Normal variational posterior for Σ . We refer to the approach of using a second VBLL for the
 155 noise covariance as Het-VBLL. While this does not result in a closed-form predictive distribution, it
 156 yields convenient evaluation of the training objective.

157 For both the inverse Gamma and log-Normal variational posteriors, each of $\mathbb{E}[\Sigma_i]$, $\mathbb{E}[\Sigma_i^{-1}]$, and
 158 $\mathbb{E}[\log \Sigma_i]$ (where the expectation is taken with respect to the variational posterior) are analytically
 159 tractable. These results are presented in Appendices B and C. In the next section, we develop lower

²Throughout, we use overbars to denote mean parameters and underbars to denote prior parameters

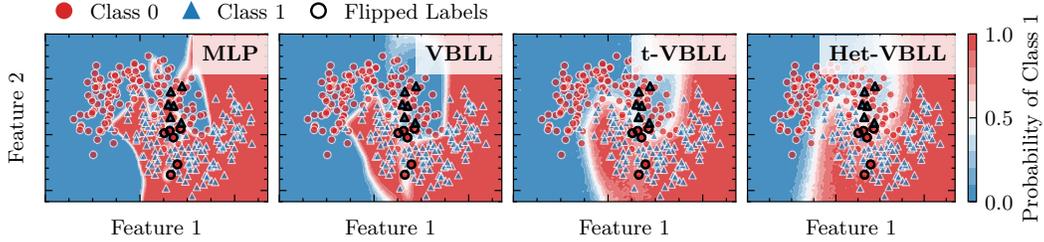


Figure 2: Classification on a noisy two-moon dataset under label noise. A standard MLP struggles on this complex data set. VBLLs show an improved decision boundary but still display slight overfitting. Our models perform best, with Het-VBLL providing the best qualitative results.

160 bounds on the log marginal likelihood for both regression and classification containing (only) these
 161 terms, yielding an analytically tractable training objective.

162 4.2 Variational Lower Bounds

163 The variational lower bound under the variational posterior is

$$\log p(Y | X) \geq \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \mathbb{E}_{q(W, \Sigma | \mathbf{x})} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] - \text{KL}(q(W, \Sigma | \mathbf{x}) || p(W, \Sigma | \mathbf{x})), \quad (8)$$

164 where the likelihood term is

$$\mathbb{E}_{q(W, \Sigma | \mathbf{x})} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] = \sum_{i=1}^{N_y} \mathbb{E}_{q(\Sigma_i | \mathbf{x})} [\mathbb{E}_{q(\mathbf{w}_i | \Sigma_i, \mathbf{x})} [\log p(\mathbf{y}_i | \mathbf{x}, \mathbf{w}_i, \Sigma_i)]]. \quad (9)$$

165 With the chosen variational posterior in (6), we now discuss the bounds that we use as training
 166 objectives in both the homoscedastic and heteroscedastic case. We will focus on the likelihood term
 167 in (8) and discuss both the regression case and the classification case. Full development of results is
 168 presented in Appendix B for the homoscedastic case and Appendix C for heteroscedastic.

169 **Regression.** In the regression case, the inner expectation over \mathbf{w}_i evaluates to

$$\mathbb{E}_{q(W, \Sigma | \mathbf{x})} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] = -\frac{1}{2} \sum_{i=1}^{N_y} \mathbb{E}_{q(\Sigma_i | \mathbf{x})} [\Sigma_i^{-1} (\mathbf{y}_i - \bar{\mathbf{w}}_i^\top \boldsymbol{\phi})^2 + \log \Sigma_i] + \boldsymbol{\phi}^\top S_i \boldsymbol{\phi} \quad (10)$$

170 following results from Harrison et al. (2024). Due to linearity with respect to Σ_i^{-1} and $\log \Sigma_i$, this
 171 expectation is analytically tractable for both the homoscedastic and the heteroscedastic case, yielding
 172 a sampling-free training objective.

173 **Classification.** For classification, we develop the following bound on the likelihood term

$$\begin{aligned} \mathbb{E}_{q(W, \Sigma | \mathbf{x})} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] &\geq \mathbf{y}^\top \bar{W} \boldsymbol{\phi} - \text{LSE}_i(\bar{\mathbf{w}}_i^\top \boldsymbol{\phi} + \alpha_i (\boldsymbol{\phi}^\top S_i \boldsymbol{\phi} + 1)) \\ &\quad - \frac{1}{2} \sum_i \mathbb{E}_{q(\Sigma_i | \mathbf{x})} [4\Sigma_i - \alpha_i + \alpha_i^2 \Sigma_i^{-1}] (\boldsymbol{\phi}^\top S_i \boldsymbol{\phi} + 1) \end{aligned} \quad (11)$$

174 where α_i are variational parameters. This result was derived to yield linearity in Σ_i, Σ_i^{-1} which in
 175 turn yields analytical tractability of this objective. This result builds upon the variational multivariate
 176 logistic regression ELBO developed by Knowles & Minka (2011).

177 4.3 Training and Prediction

178 Both the regression and classification models are trained in the same way as standard neural networks:
 179 the weights of the variational posteriors and the neural network features $\boldsymbol{\phi}$ are jointly trained via
 180 minibatch gradient descent methods. In our experiments we use AdamW (Kingma et al., 2015;
 181 Loshchilov & Hutter, 2017), but other standard optimization algorithms also work well. We find
 182 gradient clipping to substantially improve performance, as noted by Harrison et al. (2024).

183 Following [Blundell et al. \(2015\)](#), we adjust the weight on the KL terms in the training objective to
 184 yield unbiased minibatch estimation. In particular, we divide (8) by the number of data points (which
 185 we write $|\mathcal{D}|$), and thus estimation of the likelihood term can be done by the mean of a minibatch, and
 186 the KL terms have a weight of $1/|\mathcal{D}|$. In the heteroscedastic setting, the KL term in (8) decomposes
 187 into

$$\text{KL}(q(W, \Sigma | \mathbf{x}) \| p(W, \Sigma | \mathbf{x})) = \sum_{i=1}^{N_y} (\text{KL}(q(\mathbf{m}_i) \| p(\mathbf{m}_i))) + \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{q(\Sigma_i | \mathbf{x})} [\text{KL}(q(\mathbf{w}_i | \Sigma_i) \| p(\mathbf{w}_i | \Sigma_i))] \quad (12)$$

188 If we correct for minibatch training, we see that the first term has a $1/|\mathcal{D}|$ weight whereas the second
 189 term has a weight of 1 (this is discussed further in [Appendix C](#)). Intuitively, this means that the
 190 posterior will not concentrate as more data is observed. Therefore, in practice we also use a $1/|\mathcal{D}|$
 191 weight for this term which substantially improves performance.

192 **Training Procedure.** The models presented here can both be used in full model training, or can be
 193 used in a multi-step training procedure. We investigate three approaches to using t- and Het-VBLLs
 194 in this work:

- 195 1. Train a network with a t/Het-VBLL head from scratch.
- 196 2. Use a t/Het-VBLL head on frozen, pretrained features.
- 197 3. Replace the standard last layer of a pretrained network and jointly train the last layer and
 198 fine-tune the features.

199 All three of these methods perform well, although they are suitable to different settings. For smaller
 200 models (such as those we use in regression and Bayesian optimization) we train from scratch. For
 201 very large models, t/Het-VBLL models can require training for as much as 50-100% more total steps³.
 202 Therefore, we generally advocate for using our last layers with pretrained features for large models,
 203 either via freezing features or fine-tuning the model. For the case in which we use t/Het-VBLL
 204 on top of pretrained features, we initialize the mean of the VBLL layer to match the last layer of
 205 the pretrained model. We use methods for improved gradient estimation in heteroscedastic models,
 206 following ([Skafte et al., 2019](#); [Seitzer et al., 2022](#); [Stirn et al., 2023](#)). Details on gradient estimation
 207 and the complexity of training/prediction are provided in the appendix.

208 **Hyperparameters.** Most hyperparameters for t/Het-VBLLs match those from the standard VBLL
 209 model. The primary hyperparameters are the last layer priors, which we chose to be simple zero-mean
 210 isotropic Gaussians. Thus, the only hyperparameter that has to be tuned is the variance scale. For Het-
 211 VBLL models, the use of two VBLL heads leads to two prior variance scale parameters. For t-VBLL
 212 models, the prior for the noise introduces hyperparameters in the form of the inverse Gamma prior
 213 parameters. We share parameters across output dimensions, reducing the number of hyperparams to
 214 two. For classification models, the bound in (11) introduces variational parameters α_i for each output
 215 dimension i . These can be learned jointly with the network weights (which automatically tightens the
 216 bound) or can be fixed and treated as a hyperparameter. Moreover, these parameters can be shared
 217 across output dimensions or treated independently for each dimension.

218 **Prediction.** Prediction for regression in the homoscedastic case is analytically tractable via a
 219 Student t-distributed posterior predictive. For t-VBLL classification, it is necessary to sample noise
 220 realizations via first sampling a realization of Σ from the variational posterior, and then sampling each
 221 $z_i | \mathbf{x}, \Sigma_i \sim \mathcal{N}(\bar{\mathbf{w}}_i^\top \phi(\mathbf{x}), \Sigma_i \phi(\mathbf{x})^\top S_i \phi(\mathbf{x}))$. Given this logit sampling procedure, Monte Carlo
 222 marginalization over the variational posterior can be performed by generating multiple samples. For
 223 the heteroscedastic case, sampling Σ is required for prediction for both regression and classification.
 224 For regression, given a realization of Σ , W can be analytically marginalized (matching standard
 225 VBLLs). For classification, sampling both Σ and logits is required.

226 5 Experiments

227 We investigate the performance of our t-VBLL and Het-VBLL in supervised regression and classifica-
 228 tion, segmentation, and Bayesian optimization. Due to space constraints, we present only a minimal
 229 set of results in the paper body, with experiment details and further results deferred to [Appendix E](#).

³We note that this is similar to SNGP ([Liu et al., 2022](#)) and is substantially faster than many other Bayesian deep learning methods

Table 1: Results for UCI regression tasks.

	BOSTON		CONCRETE		ENERGY	
	NLL (\downarrow)	RMSE (\downarrow)	NLL (\downarrow)	RMSE (\downarrow)	NLL (\downarrow)	RMSE (\downarrow)
t-VBLL	2.51 \pm 0.10	2.87 \pm 0.28	2.95 \pm 0.08	4.70 \pm 0.36	1.19 \pm 0.13	0.85 \pm 0.11
Het-VBLL	2.35 \pm 0.14	3.14 \pm 0.53	3.05 \pm 0.13	5.20 \pm 0.35	1.17 \pm 0.11	1.53 \pm 0.23
VBLL	2.55 \pm 0.06	2.92 \pm 0.12	3.22 \pm 0.07	5.09 \pm 0.13	1.37 \pm 0.08	0.87 \pm 0.04
GBLL	2.90 \pm 0.05	4.19 \pm 0.17	3.09 \pm 0.03	5.01 \pm 0.18	0.69 \pm 0.03	0.46 \pm 0.02
LDGBLL	2.60 \pm 0.04	3.38 \pm 0.18	2.97 \pm 0.03	4.80 \pm 0.18	4.80 \pm 0.18	0.50 \pm 0.02
MAP	2.60 \pm 0.07	3.02 \pm 0.17	3.04 \pm 0.04	4.75 \pm 0.12	1.44 \pm 0.09	0.53 \pm 0.01
RBF GP	2.41 \pm 0.06	2.83 \pm 0.16	3.08 \pm 0.02	5.62 \pm 0.13	0.66 \pm 0.04	0.47 \pm 0.01
Dropout	2.36 \pm 0.04	2.78 \pm 0.16	2.90 \pm 0.02	4.45 \pm 0.11	1.33 \pm 0.00	0.53 \pm 0.01
Ensemble	2.48 \pm 0.09	2.79 \pm 0.17	3.04 \pm 0.08	4.55 \pm 0.12	0.58 \pm 0.07	0.41 \pm 0.02
SWAG	2.64 \pm 0.16	3.08 \pm 0.35	3.19 \pm 0.05	5.50 \pm 0.16	1.23 \pm 0.08	0.93 \pm 0.09
BBB	2.39 \pm 0.04	2.74 \pm 0.16	2.97 \pm 0.03	4.80 \pm 0.13	0.63 \pm 0.05	0.43 \pm 0.01

Table 2: Results on QM9

	MAE (\downarrow)	RMSE (\downarrow)	W/ Hold-Out		W/O Hold-Out	
			NLL (\downarrow)	ECE (\downarrow)	NLL (\downarrow)	ECE (\downarrow)
DNN	6.10 \pm 0.07	14.69 \pm 0.60	-2.78 \pm 0.04	0.10 \pm 0.01	24.35 \pm 6.86	0.15 \pm 0.00
VBLL	5.64 \pm 0.16	13.23 \pm 1.17	-3.01 \pm 0.08	0.09 \pm 0.00	9.28 \pm 2.15	0.13 \pm 0.00
t-VBLL	5.68 \pm 0.08	13.13 \pm 0.57	-3.01 \pm 0.04	0.09 \pm 0.00	-0.85 \pm 0.15	0.08 \pm 0.00
Het-VBLL	6.09 \pm 0.06	14.71 \pm 0.35	-2.79 \pm 0.04	0.11 \pm 0.00	-1.24 \pm 0.00	0.23 \pm 0.00
Var-NN	8.99 \pm 0.12	34.04 \pm 1.30	-3.09 \pm 0.02	0.05 \pm 0.00	-	-

230 5.1 Supervised Regression

231 **Toy Data.** To demonstrate the Het-VBLLs ability to capture heteroscedatic noise, we compare them
 232 against standard VBLLs in Figure 1 on motorcycle-impact data (Silverman, 1985; Kersting et al.,
 233 2007). Het-VBLLs are able to capture the heteroscedasticity in the data whereas standard VBLLs
 234 with an MAP estimate do not.

235 **UCI Datasets.** We also benchmark the t-VBLLs and Het-VBLLs on the standard UCI datasets (Dua
 236 & Graff, 2017) in Tables 1 and 4 (Appendix). Both models exhibit strong performance, matching or
 237 surpassing standard VBLLs (and other baselines) across all tasks. Our experimental procedure and
 238 baselines match Harrison et al. (2024) exactly, and experimental details are provided in Appendix
 239 E.1.2. Both models show strong performance across all tasks.

240 **QM9.** We run experiments on the QM9 property prediction dataset (Wu et al., 2018) (Table 2) using
 241 the zero-point energy U_0 as the target. We use the PaiNN message passing backbone from Schütt
 242 et al. (2021) with a modified readout function that can accommodate last layer approximations with
 243 any predictive distribution without removal of the inductive biases present in the PaiNN network.
 244 See Appendix E.1.3 for more details. For the MAP estimated PaiNN network, we post hoc estimate
 245 homoscedastic output noise based on a hold-out dataset, allowing us to obtain ECE and NLL metrics.
 246 For all variations of VBLL models we utilize standard deviation scaling using the same hold-out
 247 dataset as proposed in Levi et al. (2020) to calibrate the predictive uncertainties of these models (see
 248 *W/ Hold-Out* in Figure 2). Note that this retains rankings of predictive uncertainties, but simply scales
 249 their magnitudes. For some applications (such as BO), hold-out datasets may not be available and we
 250 therefore also present results where temperature scaling and noise estimation only can be done based
 251 on the training set (see *W/O Hold-Out* in Table 2).

252 In Table 2 we note how the VBLL and t-VBLL are the best performing both in terms of MAE and
 253 RMSE and the Var-NNs from Jordahn et al. (2025) have the best NLL and ECE. However, in the case
 254 where no hold-out dataset is used, the Het-VBLLs are the best in terms of NLL and ECE, and have
 255 significantly better precision than the Var-NNs. Further results on calibration of these models are
 256 provided in the Appendix.

257 5.2 Supervised Classification

258 **Toy Data.** We test our models on a noisy version of the two-moon data set in Figure 2. We further
 259 introduce label noise by flipping 20% of the labels for points within a specific range of the first
 260 feature (circled). A standard MLP overfits significantly, and even standard VBLLs struggle with these

Table 3: Results on CIFAR-10

Method	Feature Weights	Accuracy (\uparrow)	ECE (\downarrow)	NLL (\downarrow)	SVHN OOD (\uparrow)	CIFAR-100 OOD (\uparrow)
DNN	Trained	95.8 \pm 0.19	0.028 \pm 0.028	0.183 \pm 0.007	0.946 \pm 0.005	0.893 \pm 0.001
SNGP	Trained	95.7 \pm 0.14	0.017 \pm 0.003	0.149 \pm 0.005	0.960 \pm 0.004	0.902 \pm 0.003
VBLL	Trained	96.4 \pm 0.12	0.022 \pm 0.001	0.160 \pm 0.001	0.969 \pm 0.004	0.900 \pm 0.004
VBLL	Frozen	96.4 \pm 0.01	0.024 \pm 0.000	0.176 \pm 0.000	0.943 \pm 0.002	0.895 \pm 0.000
t-VBLL	Frozen	96.3 \pm 0.03	0.006 \pm 0.000	0.133 \pm 0.001	0.975 \pm 0.000	0.892 \pm 0.001
Het-VBLL	Frozen	96.2 \pm 0.02	0.009 \pm 0.000	0.135 \pm 0.000	0.975 \pm 0.001	0.894 \pm 0.001
LLLA	Frozen	96.3 \pm 0.03	0.010 \pm 0.001	0.133 \pm 0.003	0.965 \pm 0.010	0.898 \pm 0.001
VBLL	Fine-Tuned	95.6 \pm 0.02	0.025 \pm 0.000	0.168 \pm 0.000	0.955 \pm 0.005	0.900 \pm 0.002
t-VBLL	Fine-Tuned	95.6 \pm 0.01	0.014 \pm 0.001	0.175 \pm 0.002	0.950 \pm 0.013	0.834 \pm 0.007
Het-VBLL	Fine-Tuned	95.7 \pm 0.03	0.111 \pm 0.000	0.273 \pm 0.003	0.916 \pm 0.020	0.772 \pm 0.001

261 outliers. In contrast, our proposed t-VBLL and Het-VBLL models demonstrate robustness, achieving
 262 a significantly improved decision boundary.

263 **CIFAR-10/100.** We also evaluate the proposed models on CIFAR-10 and CIFAR-100, comparing t-
 264 VBLL and Het-VBLL against standard DNNs, VBLLs (Harrison et al., 2024), and LLLA (Daxberger
 265 et al., 2021) baselines. On CIFAR-10 (Table 3), both t-VBLL and Het-VBLL achieved competitive
 266 accuracy while significantly reducing Expected Calibration Error and Negative Log-Likelihood,
 267 with improved out-of-distribution detection performance. Similarly, on CIFAR-100 (Table 5), our
 268 models maintained accuracy with a substantial reduction in ECE and NLL. The results indicate that
 269 incorporating heavy-tailed and heteroscedastic modeling improves uncertainty quantification.

270 5.3 Semantic Segmentation

271 Semantic segmentation inherently presents a heteroscedastic problem, with per-pixel uncertainty
 272 varying significantly across diverse scenes (Kendall & Gal, 2017). To validate our models in this
 273 setting, we evaluate on the CityScapes semantic segmentation benchmark validation set (Cordts
 274 et al., 2016). We use the SegFormer model with an MIT-B0 backbone (Xie et al., 2021) as our
 275 base architecture. We compare this baseline to two Bayesian extensions: a standard VBLL (VBLL
 276 SegFormer) and its heteroscedastic counterpart (Het-VBLL SegFormer).

277 The Het-VBLL model achieved the highest mean Intersection-over-Union (mIoU) of 0.769, out-
 278 performing both the VBLL and baseline models, as shown in Table 6 (Appendix). Furthermore,
 279 Het-VBLL improved segmentation performance across 13 of 19 individual classes. Performance
 280 improvements were observed in disambiguating drivable space across classes such as road, sidewalk.
 281 Improvements were also notable for relatively rare classes such as Bus and Truck. These results high-
 282 light the efficacy of explicitly modeling heteroscedastic noise, demonstrating tangible improvements
 283 in semantic segmentation performance. A visualization of the predictive uncertainty is provided
 284 in Figure 3, showing the segmentation for a frame and the aleatoric and epistemic entropy for the
 285 baseline model and the Het-VBLL model.

286 5.4 Bayesian Optimization

287 We evaluate t-VBLL and Het-VBLL networks as surrogates in Bayesian optimization. As baselines,
 288 we use standard VBLLs (Harrison et al., 2024; Brunzema et al., 2025), Gaussian processes (GPs) with
 289 box constraints on the lengthscales, GPs with D-scaled priors on the lengthscales (D-GP) (Hvarfner
 290 et al., 2024), GPs with a Yeo-Johnson outcome transformation (YJ-GP) which is popular for dealing
 291 with heteroscedastic noise (Cowen-Rivers et al., 2022), as well as a variance predictive network
 292 (Var-NN). All experimental details and surrogate configurations are listed in Appendix E.4. We
 293 compare all models using an upper confidence bound acquisition function (Srinivas et al., 2010) with
 294 the same hyperparameter $\beta = 2$ such that the difference in performance is only due to differences in
 295 the model.

296 As a benchmark experiments, we focus on the classic Ackley objective ($N_x = 5$), as well as pest
 297 control with $N_x = 25$ (Oh et al., 2019) and Lunar Lander Eriksson et al. (2019), where automatic
 298 feature learning in Bayesian neural networks has been shown to outperform standard GPs (Li et al.,
 299 2024; Brunzema et al., 2025). To test our models, we add heavy-tailed noise from a zero-mean
 300 Laplace distribution to the outcome of an experiment, controlled by an outlier probability. In Fig. 4,
 301 we show the performance of all surrogates over this outlier probability for 10 random seeds each.

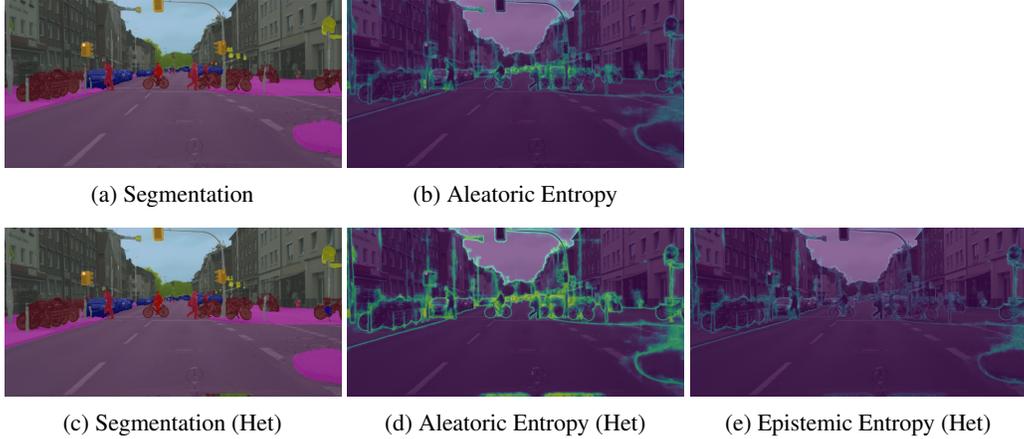


Figure 3: Comparison of the base model (top row) and the proposed Het-VBLL model (bottom row). The models are evaluated on semantic segmentation and uncertainty estimation. Note that the baseline model does not yield an epistemic uncertainty prediction.

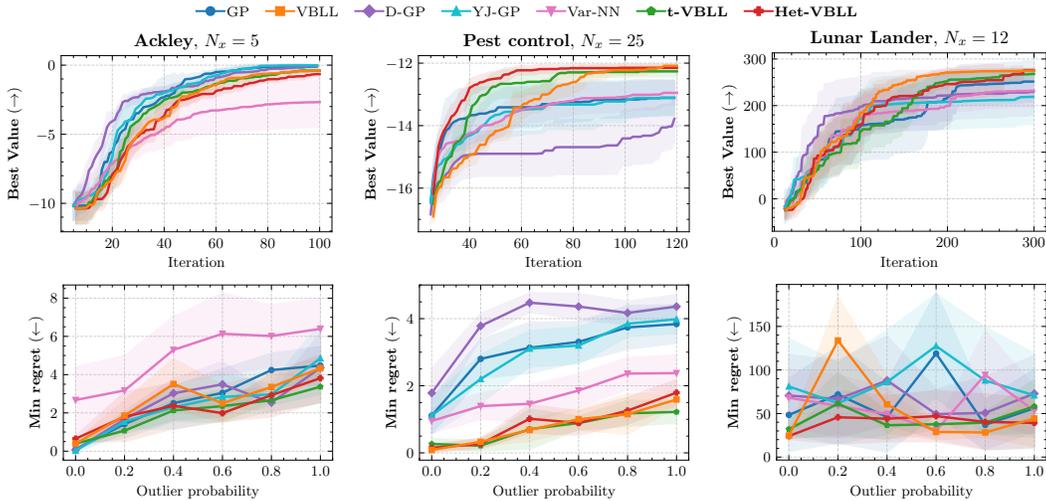


Figure 4: *Top row*: Best value without outliers. *Bottom row*: Minimum regret obtained under heavy-tail distributed outliers on Ackley, Pest Control, and Lunar Lander. Our models demonstrate more robustness to outliers as well as competitive performance regarding best value for no outliers.

302 VBLL-based models clearly outperform GPs. We can further see, that as the outlier probability
 303 increases also the minimum regret obtained by a surrogate increases. Our proposed models perform
 304 on par with VBLLs for an outlier probability of zero are more robust to outliers.

305 6 Discussion and Conclusion

306 In this paper we have introduced two novel methods for scalable Bayesian deep learning, each
 307 applicable to both regression and classification. Development of these methods relied on the design of
 308 novel training objectives and model architectures that enabled sampling-free computation of a lower
 309 bound on the (standard) ELBO. Both t-VBLL and Het-VBLL show strong performance relatively to
 310 both similar last layer methods and considerably more expensive methods. Several questions remain
 311 with these models, however. First, we have investigated basic approaches described in the literature
 312 to stabilize the heteroscedastic models, but they still show inconsistent performance (as in the QM9
 313 experiments). Second, training of VBLL models is still relatively slow compared to vanilla models,
 314 and optimization schemes to accelerate training are likely possible to accelerate training.

315 **References**

- 316 Balandat, M., Karrer, B., Jiang, D. R., Daulton, S., Letham, B., Wilson, A. G., and Bakshy, E.
317 BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Neural Information*
318 *Processing Systems (NeurIPS)*, 2020.
- 319 Bishop, C. and Quazaz, C. Regression with input-dependent noise: A bayesian treatment. *Neural*
320 *Information Processing Systems (NeurIPS)*, 9, 1996.
- 321 Blei, D. M. and Lafferty, J. D. A correlated topic model of science. *The annals of applied statistics*,
322 2007.
- 323 Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. Weight uncertainty in neural network.
324 In *International Conference on Machine Learning (ICML)*, 2015.
- 325 Box, G. E. and Tiao, G. C. *Bayesian inference in statistical analysis*, volume 40. John Wiley & Sons,
326 2011.
- 327 Brunzema, P., Jordahn, M., Willes, J., Trimpe, S., Snoek, J., and Harrison, J. Bayesian optimization
328 via continual variational last layer training. *International Conference on Learning Representations*
329 *(ICLR)*, 2025.
- 330 Collier, M., Mustafa, B., Kokiopoulou, E., Jenatton, R., and Berent, J. A simple probabilistic method
331 for deep classification under input-dependent label noise. *arXiv:2003.06778*, 2020.
- 332 Collier, M., Mustafa, B., Kokiopoulou, E., Jenatton, R., and Berent, J. Correlated input-dependent
333 label noise in large-scale image classification. In *IEEE Conference on Computer Vision and Pattern*
334 *Recognition (CVPR)*, 2021.
- 335 Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth,
336 S., and Schiele, B. The cityscapes dataset for semantic urban scene understanding. In *IEEE*
337 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 338 Cowen-Rivers, A. I., Lyu, W., Tutunov, R., Wang, Z., Grosnit, A., Griffiths, R. R., Maraval, A. M.,
339 Jianye, H., Wang, J., Peters, J., et al. Hebo: Pushing the limits of sample-efficient hyper-parameter
340 optimisation. *Journal of Artificial Intelligence Research (JAIR)*, 74:1269–1349, 2022.
- 341 Daxberger, E., Kristiadi, A., Immer, A., Eschenhagen, R., Bauer, M., and Hennig, P. Laplace
342 redux-effortless Bayesian deep learning. *Neural Information Processing Systems (NeurIPS)*, 2021.
- 343 Der Kiureghian, A. and Ditlevsen, O. Aleatory or epistemic? does it matter? *Structural safety*, 2009.
- 344 Dua, D. and Graff, C. UCI machine learning repository, 2017. URL [http://archive.ics.uci.](http://archive.ics.uci.edu/ml)
345 [edu/ml](http://archive.ics.uci.edu/ml).
- 346 Eriksson, D., Pearce, M., Gardner, J., Turner, R. D., and Poloczek, M. Scalable global optimization
347 via local Bayesian optimization. In *Neural Information Processing Systems (NeurIPS)*, 2019.
- 348 Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty
349 in deep learning. In *International Conference on Machine Learning (ICML)*, 2016.
- 350 Gardner, J. R., Pleiss, G., Bindel, D., Weinberger, K. Q., and Wilson, A. G. Gpytorch: Blackbox
351 matrix-matrix Gaussian process inference with gpu acceleration. In *Neural Information Processing*
352 *Systems (NeurIPS)*, 2018.
- 353 Harrison, J., Sharma, A., and Pavone, M. Meta-learning priors for efficient online Bayesian regression.
354 *Workshop on the Algorithmic Foundations of Robotics (WAFR)*, 2018.
- 355 Harrison, J., Willes, J., and Snoek, J. Variational Bayesian last layers. *International Conference on*
356 *Learning Representations (ICLR)*, 2024.
- 357 Hayashi, F. *Econometrics*. Princeton University Press, 2011.
- 358 Hvarfner, C., Hellsten, E. O., and Nardi, L. Vanilla Bayesian optimization performs great in high
359 dimensions. In *International Conference on Machine Learning (ICML)*, 2024.

- 360 Izmailov, P., Vikram, S., Hoffman, M. D., and Wilson, A. G. G. What are bayesian neural network
361 posteriors really like? In *International Conference on Machine Learning (ICML)*, 2021.
- 362 Jordahn, M., Jensen, J. V., Schmidt, M. N., and Andersen, M. R. On local posterior structure in deep
363 ensembles, 2025. URL <https://arxiv.org/abs/2503.13296>.
- 364 Kendall, A. and Gal, Y. What uncertainties do we need in bayesian deep learning for computer
365 vision? *Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- 366 Kersting, K., Plagemann, C., Pfaff, P., and Burgard, W. Most likely heteroscedastic gaussian process
367 regression. In *International Conference on Machine Learning (ICML)*, 2007.
- 368 Kingma, D. P., Salimans, T., and Welling, M. Variational dropout and the local reparameterization
369 trick. *Neural Information Processing Systems (NeurIPS)*, 28, 2015.
- 370 Knowles, D. and Minka, T. Non-conjugate variational message passing for multinomial and binary
371 regression. In *Neural Information Processing Systems (NeurIPS)*, 2011.
- 372 Kristiadi, A., Hein, M., and Hennig, P. Learnable uncertainty under laplace approximations. In
373 *Uncertainty in Artificial Intelligence (UAI)*, 2021.
- 374 Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple and scalable predictive uncertainty
375 estimation using deep ensembles. *Neural Information Processing Systems (NeurIPS)*, 2017.
- 376 Lazaro-Gredilla, M. and Titsias, M. K. Variational heteroscedastic gaussian process regression. In
377 *International Conference on Machine Learning (ICML)*, 2011.
- 378 Le, Q. V., Smola, A. J., and Canu, S. Heteroscedastic gaussian process regression. In *International
379 Conference on Machine Learning (ICML)*, pp. 489–496, 2005.
- 380 Levi, D., Gispan, L., Giladi, N., and Fetaya, E. Evaluating and calibrating uncertainty prediction in
381 regression tasks, 2020. URL <https://openreview.net/forum?id=ryg8wpEtvB>.
- 382 Li, Y. L., Rudner, T. G. J., and Wilson, A. G. A study of Bayesian neural network surrogates for
383 Bayesian optimization. In *International Conference on Learning Representations (ICLR)*, 2024.
- 384 Liu, J., Padhy, S., Ren, J., Lin, Z., Wen, Y., Jerfel, G., Nado, Z., Snoek, J., Tran, D., and Lakshmi-
385 narayanan, B. A simple approach to improve single-model deep uncertainty via distance-awareness.
386 *Journal of Machine Learning Research*, 2022.
- 387 Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv:1711.05101*, 2017.
- 388 MacKay, D. J. A practical bayesian framework for backpropagation networks. *Neural Computation*,
389 1992.
- 390 MMSegmentation. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark.
391 2020.
- 392 Neal, R. M. *Bayesian Learning for Neural Networks*. PhD thesis, University of Toronto, 1995.
- 393 Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images
394 with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised
395 Feature Learning*, 2011.
- 396 Nix, D. A. and Weigend, A. S. Estimating the mean and variance of the target probability distribution.
397 In *International conference on neural networks (ICNN)*, 1994.
- 398 Oh, C., Tomczak, J., Gavves, E., and Welling, M. Combinatorial Bayesian optimization using the
399 graph cartesian product. *Neural Information Processing Systems (NeurIPS)*, 32, 2019.
- 400 Osband, I., Aslanides, J., and Cassirer, A. Randomized prior functions for deep reinforcement
401 learning. *Neural Information Processing Systems (NeurIPS)*, 2018.
- 402 Osband, I., Wen, Z., Asghari, S. M., Dwaracherla, V., Ibrahimi, M., Lu, X., and Van Roy, B. Epistemic
403 neural networks. In *Neural Information Processing Systems (NeurIPS)*, 2023.

- 404 Papamarkou, T., Skoularidou, M., Palla, K., Aitchison, L., Arbel, J., Dunson, D., Filippone, M.,
405 Fortuin, V., Hennig, P., Hernández-Lobato, J. M., et al. Position: Bayesian deep learning is needed
406 in the age of large-scale ai. In *International Conference on Machine Learning (ICML)*, 2024.
- 407 Schütt, K., Unke, O., and Gastegger, M. Equivariant message passing for the prediction of tensorial
408 properties and molecular spectra. In *International Conference on Machine Learning*, pp. 9377–
409 9388. PMLR, 2021.
- 410 Seitzer, M., Tavakoli, A., Antic, D., and Martius, G. On the pitfalls of heteroscedastic uncer-
411 tainty estimation with probabilistic neural networks. In *International Conference on Learning
412 Representations (ICLR)*, 2022.
- 413 Silverman, B. W. Some aspects of the spline smoothing approach to non-parametric regression curve
414 fitting. *Journal of the Royal Statistical Society: Series B (Methodological)*, 47, 1985.
- 415 Skafte, N., Jorgensen, M., and Hauberg, S. Reliable training and estimation of variance networks.
416 *Neural Information Processing Systems (NeurIPS)*, 2019.
- 417 Srinivas, N., Krause, A., Kakade, S., and Seeger, M. Gaussian process optimization in the bandit
418 setting: No regret and experimental design. In *International Conference on Machine Learning
419 (ICML)*, 2010.
- 420 Stirn, A., Wessels, H., Schertzer, M., Pereira, L., Sanjana, N., and Knowles, D. Faithful heteroscedas-
421 tic regression with neural networks. In *Artificial Intelligence and Statistics (AISTATS)*, 2023.
- 422 Watson, J., Lin, J. A., Klink, P., and Peters, J. Neural linear models with functional gaussian process
423 priors. In *Advances in Approximate Bayesian Inference (AABI)*, 2020.
- 424 Watson, J., Lin, J. A., Klink, P., Pajarinen, J., and Peters, J. Latent derivative Bayesian last layer
425 networks. In *Artificial Intelligence and Statistics (AISTATS)*, 2021.
- 426 Wilson, A. G. and Izmailov, P. Bayesian deep learning and a probabilistic perspective of generalization.
427 In *Neural Information Processing Systems (NeurIPS)*, volume 33, 2020.
- 428 Wu, Z., Ramsundar, B., Feinberg, E. N., Gomes, J., Geniesse, C., Pappu, A. S., Leswing, K., and
429 Pande, V. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):
430 513–530, 2018.
- 431 Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., and Luo, P. Segformer: Simple and
432 efficient design for semantic segmentation with transformers. In *Neural Information Processing
433 Systems (NeurIPS)*, 2021.

434	A Background	14
435	A.1 Inverse Gamma	14
436	A.2 Log-Normal	14
437	B Variational Inference for Homoscedastic Noise	14
438	B.1 Variational Lower Bound and Approach	15
439	B.2 Regression	15
440	B.2.1 Training Objective	15
441	B.2.2 Prediction	16
442	B.3 Classification	16
443	B.3.1 Training Objective	16
444	B.3.2 Prediction	17
445	B.4 Complexity and Parameterization	17
446	C Heteroscedastic Noise Modeling	17
447	C.1 Bayesian Last Layer Methods for Heteroscedastic Noise	18
448	C.2 Variational Lower Bound for VBLL-Variance Networks	18
449	C.2.1 Regression	19
450	C.2.2 Classification	19
451	C.3 Prediction	19
452	C.4 Complexity and Parameterization	19
453	D Training and Algorithmic Details	19
454	D.1 Training Methodology	19
455	D.2 Improving Gradient Estimation in Heteroscedastic Models	20
456	E Experiment Details and Further Results	20
457	E.1 Supervised Regression	20
458	E.1.1 Motorcycle Dataset	20
459	E.1.2 UCI Datasets	20
460	E.1.3 QM9 Property Prediction	20
461	E.2 Supervised Classification	22
462	E.2.1 Half Moon Dataset	22
463	E.2.2 CIFAR 10 and CIFAR 100	22
464	E.3 Semantic Segmentation	22
465	E.4 Bayesian Optimization	22

466 **A Background**

467 In this section we discuss the two variational posteriors used in this paper for modeling noise
 468 covariance.

469 **A.1 Inverse Gamma**

470 Our first approach is an inverse Gamma distribution,

$$q(\Sigma_i) = \mathcal{IG}(\nu_i, \Psi_i) \tag{13}$$

471 where $\nu_i > 1$, $\Psi_i > 0$ are shape and scale parameters. This distribution is equivalent to a Gamma
 472 distribution for the inverse covariance. This distribution has several desirable properties. First, it
 473 is conjugate: in the Bayesian linear regression (BLR) model (with appropriately chosen priors), an
 474 inverse Gamma prior yields an inverse Gamma posterior. Additionally, within the BLR model, an
 475 inverse-Gamma posterior yields a t-distributed posterior predictive. While this approach is desirable
 476 for several reasons in the classical BLR setting, it is less suited to the heteroscedastic setting as we
 477 discuss in Section C. The (inverse) Gamma posterior also has a straightforward generalization for
 478 non-diagonal covariances in the (inverse) Wishart distribution. We note a few useful identities, which
 479 will be useful in developing our main training objectives:

$$\mathbb{E}[\Sigma_i^{-1}] = \nu_i \Psi_i^{-1} \tag{14}$$

$$\mathbb{E}[\log \Sigma_i] = \log(\Psi_i) - \psi(\nu_i) \tag{15}$$

480 where $\psi(\cdot)$ denotes the digamma function.

481 For this parameterization, we write our prior as $p(\Sigma_i) = \mathcal{IG}(\underline{\nu}_i, \underline{\Psi}_i)$. The KL divergence between
 482 inverse Gammas is a standard result, with

$$\text{KL}(q(\Sigma_i) \| p(\Sigma_i)) = \underline{\nu}_i \log \frac{\underline{\Psi}_i}{\underline{\Psi}_i} - \log \frac{\Gamma(\underline{\nu}_i)}{\Gamma(\underline{\nu}_i)} + (\nu_i - \underline{\nu}_i) \psi(\nu_i) - (\Psi_i - \underline{\Psi}_i) \frac{\nu_i}{\underline{\Psi}_i} \tag{16}$$

483 where $\Gamma(\cdot)$ is the gamma function.

484 **A.2 Log-Normal**

485 A second approach to the variational posterior is a log-Normal⁴

$$q(\Sigma_i) = \log \mathcal{N}(\mu_i, C_i) \tag{17}$$

486 with mean μ_i and variance $C_i > 0$. This choice of variational posterior may initially seem like an
 487 odd one: we lose the favorable conjugacy properties of the inverse Gamma posterior, and gain little
 488 in return. However, as we see in Section C, the log-Normal approach has substantial benefits for
 489 heteroscedastic modeling. We have the same identities,

$$\mathbb{E}[\Sigma_i^{-1}] = \exp(-\mu_i + \frac{1}{2}C_i) \tag{18}$$

$$\mathbb{E}[\log \Sigma_i] = \mu_i. \tag{19}$$

490 For this posterior specification, the prior is also log-Normal and written as $p(\Sigma_i) = \log \mathcal{N}(\underline{\mu}_i, \underline{C}_i)$.
 491 The KL divergence between log-Normals is equivalent to the KL divergence between their corre-
 492 sponding Normal distributions.

493 **B Variational Inference for Homoscedastic Noise**

494 In this section, we derive a variational last layer objective with noise inference in the homoscedastic
 495 case. First, we will lay out the structure of the variational lower bounds we develop throughout the
 496 paper. We then describe two variational posterior design options, and prior choices. We will then
 497 define the training objectives and the resulting posterior predictive distribution.

⁴Note that if $z \sim \mathcal{N}(\mu, \sigma^2)$ then $\exp(z) \sim \log \mathcal{N}(\mu, \sigma^2)$.

498 **B.1 Variational Lower Bound and Approach**

499 We begin by writing the lower bound on the marginal likelihood for arbitrarily specified noise
500 covariance Σ , and then discuss outcomes for inverse Gamma priors/variational posteriors. We
501 exclude dependence on the features ϕ , which we assume fixed in the derivation. Point estimates
502 for the feature weights are learned via stochastic gradient descent on the ELBO. Generally, we will
503 choose noise priors with distributions that match the variational posterior for tractability, although
504 this is not strictly necessary.

505 We structure our variational posterior as

$$q(W, \Sigma) = \prod_{i=1}^{N_y} q(\mathbf{w}_i | \Sigma_i) q(\Sigma_i) \quad (20)$$

506 with $q(\mathbf{w}_i | \Sigma_i) = \mathcal{N}(\bar{\mathbf{w}}_i, S_i \Sigma_i)$, following standard results in Bayesian regression. The variational
507 lower bound is

$$\log p(Y | X) \geq \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \mathbb{E}_{q(W, \Sigma)} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] - \text{KL}(q(W, \Sigma) \| p(W, \Sigma)) \quad (21)$$

508 By independence assumptions in the prior and variational posterior, we have

$$\text{KL}(q(W, \Sigma) \| p(W, \Sigma)) = \sum_{i=1}^{N_y} (\text{KL}(q(\Sigma_i) \| p(\Sigma_i)) + \mathbb{E}_{q(\Sigma_i)} [\text{KL}(q(\mathbf{w}_i | \Sigma_i) \| p(\mathbf{w}_i | \Sigma_i))]). \quad (22)$$

509 The expectation of the first KL term is straightforward, and is

$$\mathbb{E}_{q(\Sigma_i)} [\text{KL}(q(\mathbf{w}_i | \Sigma_i) \| p(\mathbf{w}_i | \Sigma_i))] = \frac{1}{2} (\log \det S_i - \log \det S_i - N_\phi + \text{tr}(S_i^{-1} S_i)) \quad (23)$$

$$+ \frac{1}{2} (\mathbb{E}_{q(\Sigma_i)} [\Sigma_i^{-1}] (\mathbf{w}_i - \bar{\mathbf{w}}_i) S_i^{-1} (\mathbf{w}_i - \bar{\mathbf{w}}_i)) \quad (24)$$

510 The expectation of Σ_i^{-1} can be computed via the identities in Section 2. The KL divergence between
511 prior and variational posteriors for Σ_i can also be computed via the identities in Section 2. Thus, the
512 unresolved aspect in evaluating the variational lower bound is computing the expected likelihood
513 $\mathbb{E}_{q(W, \Sigma)} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)]$.

514 **B.2 Regression**

515 We first discuss the likelihood term and prediction in the regression setting.

516 **B.2.1 Training Objective**

517 The predictive likelihood term in the ELBO is

$$\mathbb{E}_{q(W, \Sigma)} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] = \sum_{i=1}^{N_y} \mathbb{E}_{q(\Sigma_i)} [\mathbb{E}_{q(\mathbf{w}_i | \Sigma_i)} [\log p(\mathbf{y}_i | \mathbf{x}, \mathbf{w}_i, \Sigma_i)]] \quad (25)$$

518 where the inner expectation is computed as in [Harrison et al. \(2024\)](#) as

$$\mathbb{E}_{q(\mathbf{w}_i | \Sigma_i)} [\log p(\mathbf{y}_i | \mathbf{x}, \mathbf{w}_i, \Sigma_i)] = \log \mathcal{N}(\mathbf{y}_i | \bar{\mathbf{w}}_i^\top \phi, \Sigma_i) - \frac{1}{2} \phi^\top S_i \phi \quad (26)$$

$$= -\frac{1}{2} (\Sigma_i^{-1} (\mathbf{y}_i - \bar{\mathbf{w}}_i^\top \phi)^2 + \log \Sigma_i + \phi^\top S_i \phi). \quad (27)$$

519 To evaluate the outer expectation,

$$\mathbb{E}_{q(W, \Sigma)} [\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] = -\frac{1}{2} \sum_{i=1}^{N_y} \mathbb{E}_{q(\Sigma)} [\Sigma_i^{-1} (\mathbf{y}_i - \bar{\mathbf{w}}_i^\top \phi)^2 + \log \Sigma_i + \phi^\top S_i \phi] \quad (28)$$

520 we leverage the identities for $\mathbb{E}[\Sigma_i^{-1}]$ and $\mathbb{E}[\log \Sigma_i]$ described in the previous section.

521 **B.2.2 Prediction**

522 We now discuss computing the predictive distribution

$$p(\mathbf{y}_i | \mathbf{x}) = \mathbb{E}_{q(\mathbf{w}_i, \Sigma_i)}[p(\mathbf{y}_i | \mathbf{x}, \mathbf{w}_i, \Sigma_i)]. \quad (29)$$

523 **Inverse Gamma.** For the inverse Gamma variational posterior, we can exploit standard conjugacy
524 results. In particular, the posterior predictive for each row i is multivariate t -distributed,

$$p(\mathbf{y}_i | \mathbf{x}) = t_{2\nu_i}(\bar{\mathbf{w}}_i^\top \phi, \frac{\Psi}{\nu_i}(1 + \phi^\top S_i \phi)) \quad (30)$$

525 **Log-Normal.** For the log-Normal variational posterior, we must turn instead to a Monte Carlo
526 approximation. We will sample

$$\hat{\Sigma}_i \sim q(\Sigma_i) \quad (31)$$

527 for all i , and for which sampling is straight-forward by simply sampling from a normal and exponen-
528 tiating. Given this realized sample, we can marginalize over the last layer yielding predictive

$$p(\mathbf{y}_i | \mathbf{x}) = \mathcal{N}(\bar{\mathbf{w}}_i^\top \phi, \hat{\Sigma}_i(\phi^\top S_i \phi + 1)) \quad (32)$$

529 **B.3 Classification**

530 We consider a classification model of the form

$$p(\mathbf{y} | \mathbf{x}) = \text{softmax}(W\phi(\mathbf{x}) + \varepsilon) \quad (33)$$

531 where the addition of the aleatoric noise ε is optional. Our model exactly matches the regression
532 model, with the only difference being the softmax. Critically, we again assume Σ is diagonal, with
533 inverse-Gamma distributed diagonal entries.

534 **B.3.1 Training Objective**

535 We write the likelihood

$$\mathbb{E}_{q(W, \Sigma)}[\log p(\mathbf{y} | \mathbf{x}, W, \Sigma)] = \mathbf{y}^\top \bar{W} \phi - \mathbb{E}_{q(\Sigma)} \mathbb{E}_{q(W|\Sigma)}[\text{LSE}_i(\mathbf{w}_i^\top \phi + \varepsilon_i)] \quad (34)$$

536 where we assume \mathbf{y} is a one-hot encoding of the class labels, and lower bound the $\text{LSE}(\cdot)$ term. It
537 may be tempting to exploit the relatively standard (Blei & Lafferty, 2007) approach for variational
538 multinomial logistic regression to construct a lower bound on the log-sum-exp term for the inner
539 expectation (as is used in Harrison et al. (2024)), yielding

$$-\mathbb{E}_{q(W|\Sigma)}[\text{LSE}_i(\mathbf{w}_i^\top \phi + \varepsilon_i)] \geq -\mathbb{E}_{q(\Sigma)}[\log \sum_i \exp(\bar{\mathbf{w}}_i^\top \phi + \frac{\Sigma_i}{2}(1 + \phi^\top S_i \phi))]. \quad (35)$$

540 It is possible to further exchange the expectation and the negative log/sum terms, yielding

$$-\mathbb{E}_{q(\Sigma)}[\text{LSE}_i(\bar{\mathbf{w}}_i^\top \phi + \frac{\Sigma_i}{2}(1 + \phi^\top S_i \phi))] \geq -\log \sum_i \mathbb{E}_{q(\Sigma)}[\exp(\bar{\mathbf{w}}_i^\top \phi + \frac{\Sigma_i}{2}(1 + \phi^\top S_i \phi))]. \quad (36)$$

541 However, the expectation on the RHS generally does not exist⁵. Thus, we propose two possible
542 approaches, each of which we discuss below.

543 **Semi-Monte Carlo.** First, we may turn to sampling. We sample (via reparameterization trick,
544 available in standard automatic differentiation packages) gamma random variables to compute a
545 Monte Carlo approximation to the expectation in (35).

546 **Reduced Knowles-Minka.** Instead of directly using the bound on the log-sum-exp used in Harrison
547 et al. (2024), we can instead use the main result from Knowles & Minka (2011), where (applying
548 their result to our chosen parameterization)

$$-\mathbb{E}_{q(W|\Sigma)}[\text{LSE}_i(\mathbf{w}_i^\top \phi + \varepsilon)] \geq -\frac{1}{2} \sum_i a_i^2 (\phi^\top S_i \phi + 1) \Sigma_i - \text{LSE}_i(\bar{\mathbf{w}}_i^\top \phi + (\frac{1}{2} - a_i)(\phi^\top S_i \phi + 1) \Sigma_i) \quad (37)$$

⁵This can be seen by the non-existence of the moment generating function of the inverse gamma distribution.

549 where a_i are variational parameters that are typically optimized to be maximally tight. Note that
 550 choosing $a_i = 0$ for all i exactly recovers (35). To yield a tractable bound for the outer expectation
 551 with respect to Σ , it is necessary to remove it from inside the log-sum-exp term. So, we choose

$$a_i = \frac{1}{2} - \frac{\alpha_i}{\Sigma_i} \quad (38)$$

552 and plugging in yields

$$-\mathbb{E}[\text{LSE}_i(\mathbf{w}_i^\top \phi + \varepsilon)] \geq -\frac{1}{2} \sum_i \left(\frac{1}{2} - \frac{\alpha_i}{\Sigma_i}\right)^2 (\phi^\top S_i \phi + 1) \Sigma_i - \text{LSE}_i(\bar{\mathbf{w}}_i^\top \phi + \alpha_i (\phi^\top S_i \phi + 1)) \quad (39)$$

$$= -\frac{1}{2} \sum_i \left(\frac{\Sigma_i}{4} - \alpha_i + \frac{\alpha_i^2}{\Sigma_i}\right) (\phi^\top S_i \phi + 1) - \text{LSE}_i(\bar{\mathbf{w}}_i^\top \phi + \alpha_i (\phi^\top S_i \phi + 1)) \quad (40)$$

553 which is analytically tractable for the outer expectation over Σ , yielding (for the inverse Gamma
 554 variational posterior)

$$\mathbb{E}_{q(W, \Sigma)}[\log p(\mathbf{y} \mid \mathbf{x}, W, \Sigma)] \geq \mathbf{y}^\top \bar{W} \phi - \text{LSE}_i(\bar{\mathbf{w}}_i^\top \phi + \alpha_i (\phi^\top S_i \phi + 1)) \quad (41)$$

$$- \frac{1}{2} \sum_i \left(\frac{\Psi_i}{4(\nu_i - 1)} - \alpha_i + \frac{\alpha_i^2 \nu_i}{\Psi_i} \right) (\phi^\top S_i \phi + 1). \quad (42)$$

555 We have several options with α_i ; setting $\alpha_i = \frac{1}{2}$ results in the first two terms exactly matching the
 556 standard VBLL objective (Harrison et al., 2024). Choosing $\alpha_i = 0$ yields the remarkably simple
 557 overall bound

$$-\mathbb{E}_{q(W \mid \Sigma)}[\text{LSE}_i(\mathbf{w}_i^\top \phi + \varepsilon)] \geq -\frac{1}{8} \sum_i (\phi^\top S_i \phi + 1) \Sigma_i - \text{LSE}_i(\bar{\mathbf{w}}_i^\top \phi). \quad (43)$$

558 Other options can be chosen for learning α_i 's. We can treat them as hyperparameters, in which it is
 559 convenient to set $\alpha_i = \alpha$ to reduced the number of parameters, and it can be swept over. Alternatively,
 560 because the bound holds for any α_i , they can be learned as model parameters together with the other
 561 model parameters. Knowles & Minka (2011) propose an iterative update that exploits convexity with
 562 respect to α_i 's—while better optimization schemes exploiting convexity are possible, we will not
 563 investigate them in this paper.

564 B.3.2 Prediction

565 For prediction, we again turn to Monte Carlo approximation within the hierarchical model, combined
 566 with local reparameterization, and sample $\hat{\Sigma}_i$ from the variational posterior, and compute

$$\hat{z}_i \sim \mathcal{N}(\bar{\mathbf{w}}_i^\top \phi, \hat{\Sigma}_i (1 + \phi^\top S_i \phi)) \quad (44)$$

567 for each logit element. Variance reduction schemes are possible for prediction, but they are beyond
 568 the scope of this paper.

569 B.4 Complexity and Parameterization

570 We follow the parameterization presented in Harrison et al. (2024), which proposes to parameterize
 571 the variational posterior $q(\bar{\mathbf{w}}_i, \Sigma_i S_i)$ via a simple unconstrained tensor for $\bar{\mathbf{w}}_i$ and with either a
 572 strictly positive diagonal or Cholesky-decomposed representation for S_i . Recall the inverse Gamma
 573 variational posterior is $q(\Sigma_i) = \mathcal{IG}(\nu_i, \Psi_i)$ with $\nu_i > 1, \Psi_i > 0$. Enforcing strict positivity is done
 574 via exponentiating tensors (that are unconstrained). Lower bounds on parameter values are similarly
 575 accomplished by adding the offset.

576 The addition of variational inference in the homoscedastic case adds minimal additional complexity.
 577 The only additional tensors added are those of the variational posterior. Thus, we add $2N_y$ parameters
 578 but do not estimate a point estimate for Σ_i , and thus add only N_y parameters.

579 C Heteroscedastic Noise Modeling

580 In this section we discuss approaches to modeling heteroscedastic noise that also quantifies the
 581 epistemic uncertainty associated with aleatoric noise prediction.

582 **C.1 Bayesian Last Layer Methods for Heteroscedastic Noise**

583 Our approach to aleatoric noise modeling in this work builds on the standard VBLL model. In
 584 particular, we use a VBLL variational posterior for the last layer in a covariance predictive model of
 585 the form

$$\log \Sigma_i = \mathbf{m}_i^\top \phi(\mathbf{x}) \quad (45)$$

586 with $\mathbf{m}_i \sim q(\mathbf{m}_i) = \mathcal{N}(\bar{\mathbf{m}}, Z)$ (and similarly choose a Normal prior $\mathbf{m}_i \sim \mathcal{N}(\bar{\mathbf{m}}, Z)$).

587 Given this structure, $\log \Sigma_i$ is Normally distributed and Σ_i is log-Normally distributed, as

$$\Sigma_i \sim \log \mathcal{N}(\bar{\mathbf{m}}_i^\top \phi(\mathbf{x}), \phi(\mathbf{x})^\top Z \phi(\mathbf{x})) \quad (46)$$

588 We choose a prior of the same structure over \mathbf{m} , which we write $p(\mathbf{m})$. We established in Section
 589 2 that the log-Normal covariance distribution is a reasonable one. While it allows for tractable
 590 variational objectives, it does not allow analytical marginalization for the predictive distribution. We
 591 will note identities which are critical in our development, which build upon those presented in Section
 592 2:

$$\mathbb{E}_{q(\mathbf{m}_i)}[\Sigma_i] = \exp(\bar{\mathbf{m}}_i^\top \phi(\mathbf{x}) + \frac{1}{2} \phi(\mathbf{x})^\top Z_i \phi(\mathbf{x})) \quad (47)$$

$$\mathbb{E}_{q(\mathbf{m}_i)}[\Sigma_i^{-1}] = \exp(-\bar{\mathbf{m}}_i^\top \phi(\mathbf{x}) + \frac{1}{2} \phi(\mathbf{x})^\top Z_i \phi(\mathbf{x})) \quad (48)$$

$$\mathbb{E}_{q(\mathbf{m}_i)}[\log \Sigma_i] = \bar{\mathbf{m}}_i \quad (49)$$

593 We can compare this modeling approach to learning point estimates of Σ in the standard VBLL model,
 594 or to learning a standard variance prediction network. If we set $N_\phi = 1$ and set $\phi = 1$, we exactly
 595 recover the VBLL noise estimation scheme under a log-Normal prior. Thus, this heteroscedastic noise
 596 scheme represents a strict generalization of the standard VBLL model. If we replace the variational
 597 posterior over \mathbf{m}_i with a point estimate, we recover a standard variance prediction model.

598 **C.2 Variational Lower Bound for VBLL-Variance Networks**

599 We can obtain tractable variational objectives by combining our variance parameterization with the
 600 variational objectives obtained in the last section. Note that each input \mathbf{x} induces a Σ and W in our
 601 generative model. For T training examples, our variational posterior is

$$q(W_{1:T}, M | X) = q(M) \prod_{t=1}^T q(W_t | M, \mathbf{x}_t) \quad (50)$$

602 where t indexes training data. Following the previous section and indexing rows with i (and dropping
 603 data indexing), the terms in this factorized variational posterior can be written

$$q(W | M, \mathbf{x}) = \prod_{i=1}^{N_y} q(\mathbf{w}_i | \mathbf{m}_i, \mathbf{x}) \quad (51)$$

604 with

$$q(\mathbf{w}_i | \mathbf{m}_i, \mathbf{x}) = \mathcal{N}(\bar{\mathbf{w}}_i, \Sigma_i(\mathbf{m}_i, \mathbf{x}) S_i) \quad (52)$$

605 and

$$q(M) = \prod_{i=1}^{N_y} q(\mathbf{m}_i). \quad (53)$$

606 With this variational posterior, we have

$$\log p(Y | X) \geq \mathbb{E}_{q(M)}[\log p(Y | X, M)] - \text{KL}(q(M) \| p(M)) \quad (54)$$

$$\begin{aligned} &\geq \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \mathbb{E}_{q(W|M, \mathbf{x})q(M)}[\log p(\mathbf{y} | \mathbf{x}, W, M)] - \text{KL}(q(M) \| p(M)) \quad (55) \\ &\quad - \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{q(M)}[\text{KL}(q(W | M, \mathbf{x}) \| p(W | M, \mathbf{x}))]. \end{aligned}$$

607 There are two KL terms in (55). The term for $q(M)$ is a straightforward KL between Gaussians, and
 608 factorizes over the dimensionality of \mathbf{y} . The second KL term is

$$\sum_{\mathbf{x} \in \mathcal{D}} \sum_{i=1}^{N_y} \frac{1}{2} (\mathbb{E}_{q(m_i)}[\Sigma_i^{-1}] (\bar{\mathbf{w}} - \bar{\mathbf{w}})^\top S_i^{-1} (\bar{\mathbf{w}} - \bar{\mathbf{w}}) + \text{tr}(S_i^{-1} S_i) - \log \frac{\det S_i}{\det \underline{S}_i} - N_\phi) \quad (56)$$

609 which is tractable via (48).

610 Practically, the relative weight of the second KL term is much larger as it is the sum of T terms. In
 611 practice, we will scale this second KL term by $1/T$ (to match the relative weight of the first KL term),
 612 which improves performance. This larger weighting factor results from sharing the parameters of the
 613 variational posterior for each W (specifically \bar{W} , S).

614 C.2.1 Regression

615 To compute the likelihood term in (55) in the regression case, we build upon the objective as written
 616 in (28), and apply the developed identities for $\mathbb{E}[\Sigma_i^{-1}]$ and $\mathbb{E}[\log \Sigma_i]$.

617 C.2.2 Classification

618 To compute the likelihood term for the classification case, we use the previously developed objective in
 619 (40) and apply the expressions for $\mathbb{E}[\Sigma_i^{-1}]$ and $\mathbb{E}[\Sigma_i]$, yielding an analytically tractable heteroscedastic
 620 classification model.

621 C.3 Prediction

622 Because the variational posterior for the noise covariance is log-Normal, we lose conjugacy for
 623 prediction in the regression case. Thus, for both regression and classification, we turn to sampling.
 624 For the regression case, we sample Σ from the variational posterior and then marginalize \mathbf{w} as in
 625 standard VBLL regression models. In classification, we are forced to sample the noise covariance Σ ,
 626 and then sample from the conditional distribution over logits under the variational posterior.

627 C.4 Complexity and Parameterization

628 The approach to heteroscedastic modeling fundamentally relies on parameterizing two VBLL heads:
 629 one for the mean, and one for the noise covariance. Both are parameterized as with standard VBLLs,
 630 which can have a dense, diagonal, or low-rank covariance structure, and thus the complexity of
 631 this layer is twice the complexity of a standard VBLL layer. However, we note that VBLLs (with
 632 appropriately chosen covariance parameterizations) have comparable complexity to standard neural
 633 network layers, and thus the added computational cost of adding heteroscedasticity is comparable to
 634 adding an additional layer to a neural network.

635 D Training and Algorithmic Details

636 We train our t-VBLL and Het-VBLL models with standard neural network optimization strategies.
 637 Concretely, we jointly train the parameters of the variational posterior together with the the neural
 638 network features ϕ . In this work we only train feature point estimates, but it is also possible to train
 639 variational posteriors for features via Bayes-by-backprop (as discussed in Harrison et al. (2024)), or
 640 train ensembles of models with t-VBLL or Het-VBLL heads.

641 D.1 Training Methodology

642 We train via minibatch optimization, with the sums over the data replaced with (re-weighted)
 643 expectations over minibatches as in Blundell et al. (2015); Harrison et al. (2024). Note that this
 644 implies that, in the heteroscedastic model, the KL term for $q(M)$ should be weighted by one over the
 645 dataset size, whereas the other KL term and the likelihood term are averaged over the dataset.

646 The models presented in this paper can be used for full model training, or for a phase of training in
 647 more complex network training pipelines. For example, a t- or Het-VBLL head can be trained as a
 648 linear probe on pre-trained features, or may be used as a head for a fine-tuned model. Practically,
 649 VBLL-based models train slightly more slowly than standard neural networks due to (typically) being
 650 more heavily regularized, and thus training VBLL-based heads in a second phase of training often
 651 accelerates training versus training from scratch with VBLL heads.

Table 4: Further results for UCI regression tasks.

	POWER		WINE		YACHT	
	NLL (\downarrow)	RMSE (\downarrow)	NLL (\downarrow)	RMSE (\downarrow)	NLL (\downarrow)	RMSE (\downarrow)
t-VBLL	2.75 ± 0.02	3.83 ± 0.07	0.91 ± 0.05	0.62 ± 0.03	0.99 ± 0.34	0.87 ± 0.21
Het-VBLL	2.73 ± 0.03	3.75 ± 0.09	0.92 ± 0.07	0.61 ± 0.03	0.74 ± 0.50	1.94 ± 1.03
VBLL	2.73 ± 0.01	3.68 ± 0.03	1.02 ± 0.03	0.65 ± 0.01	1.29 ± 0.17	0.86 ± 0.17
GBLL	2.77 ± 0.01	3.85 ± 0.03	1.02 ± 0.01	0.64 ± 0.01	1.67 ± 0.11	1.09 ± 0.09
LDGBLL	2.77 ± 0.01	3.85 ± 0.04	1.02 ± 0.01	0.64 ± 0.01	1.13 ± 0.06	0.75 ± 0.10
MAP	2.77 ± 0.01	3.81 ± 0.04	0.96 ± 0.01	0.63 ± 0.01	5.14 ± 1.62	0.94 ± 0.09
RBF GP	2.76 ± 0.01	3.72 ± 0.04	0.45 ± 0.01	0.56 ± 0.05	0.17 ± 0.03	0.40 ± 0.03
Dropout	2.80 ± 0.01	3.90 ± 0.04	0.93 ± 0.01	0.61 ± 0.01	1.82 ± 0.01	1.21 ± 0.13
Ensemble	2.70 ± 0.01	3.59 ± 0.04	0.95 ± 0.01	0.63 ± 0.01	0.35 ± 0.07	0.83 ± 0.08
SWAG	2.77 ± 0.02	3.85 ± 0.05	0.96 ± 0.03	0.63 ± 0.01	1.11 ± 0.05	1.13 ± 0.20
BBB	2.77 ± 0.01	3.86 ± 0.04	0.95 ± 0.01	0.63 ± 0.01	1.43 ± 0.17	1.10 ± 0.11

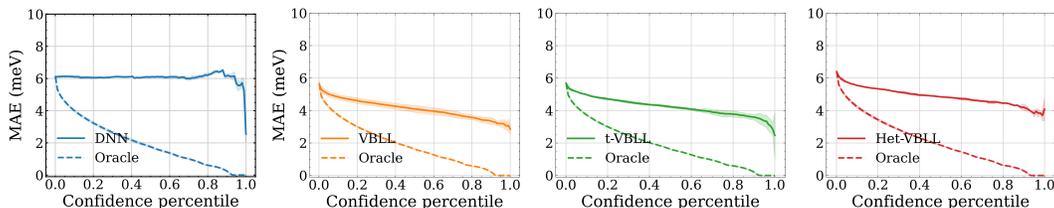


Figure 5: Calibration curves from models trained on QM9 where MAE is plotted as a function of confidence threshold. Oracle lines indicate the theoretically optimal ranking of predictive uncertainties.

652 D.2 Improving Gradient Estimation in Heteroscedastic Models

653 The faithful heteroscedastic regression method proposed in [Stirn et al. \(2023\)](#) has been reported to
 654 improve heteroscedastic variance output networks and can easily be applied on top of the proposed
 655 Het-VBLL models. We implement the faithful regression method by detaching the gradients from
 656 the variance output, such that they do not affect the learning of the features $\phi(\mathbf{x})$, and scaling the
 657 gradient accordingly to obtain a Newton step rather than a gradient step. In regression, the detaching
 658 can be done specifically by detaching the features $\phi(\mathbf{x})$ before computing the log-noise covariance
 659 term. Generally we observe minimal improvements when using faithful regression.

660 E Experiment Details and Further Results

661 E.1 Supervised Regression

662 E.1.1 Motorcycle Dataset

663 For the experiment in Figure 1 (left) on the data set provided by [Silverman \(1985\)](#), we use for both
 664 the VBLL and three Het-VBLL three hidden layers with a feature dimension of $N_\phi = 64$, ELU
 665 activation functions, and a prior scale of one. As optimizer, we use AdamW ([Loshchilov & Hutter, 2017](#)).
 666 For the Het-VBLL, we choose a learning rate of $3e-4$ and a weight decay of $1e-5$. For
 667 the standard VBLLs, we choose a learning rate of $3e-3$ and a weight decay of 0. For both, we use
 668 gradient clipping at 1. We train both surrogates for 2000 epochs with a batch size of 32.

669 E.1.2 UCI Datasets

670 Results on UCI datasets ([Dua & Graff, 2017](#)) are shown in Tables 1 and 4. These experiments match
 671 the setting used in [Watson et al. \(2021\)](#) and [Harrison et al. \(2024\)](#), and we compare against their
 672 baselines. In particular, we use two hidden layer MLPs of width 50, and use a batch size of 32 (except
 673 for POWER in which we use batch size 256). We use AdamW ([Loshchilov & Hutter, 2017](#)) with a
 674 learning rate of $1e-3$ and weight decay of $1e-2$ on the hidden layers. We ran 10 seeds for each
 675 dataset. For more details on the experimental setting, we refer the reader to [Harrison et al. \(2024\)](#).

676 E.1.3 QM9 Property Prediction

677 We also run experiments on the QM9 property prediction dataset ([Wu et al., 2018](#)) using the zero-point
 678 energy U_0 as the target. We use the PaiNN message passing backbone from [Schütt et al. \(2021\)](#)
 679 with a modified readout function. The readout function has been modified such that full molecule

680 energy is predicted based on *aggregate features* from all atoms in the molecule, rather than as a sum
 681 of energies from individual atoms in the molecule. This modification allows us to directly apply
 682 VBLL layers with any predictive distribution, which the original read-out function would not. To
 683 retain the inductive bias that total molecule energy scales linearly in the number of atoms in the
 684 molecule, we use summation as our aggregation function. See Figure 6 for illustration of original and
 685 modified readout functions. Finally, we use the reference atom energy scaling as provided in QM9,
 686 and perform a *per atom* standardization and scaling of energy values thus retaining the inductive bias
 687 of energy scaling with number of atoms in the molecule. We refer to the code for details.

688 In Figure 5 we highlight the calibrated ranking of the predictive uncertainty of the VBLL models by
 689 plotting confidence graphs that illustrate that the uncertainty estimates of the models generally have
 690 good rankings. The oracle lines indicate the theoretically optimal rankings. We illustrate that the
 691 predictive uncertainty rankings of the VBLL models, where we note that the t-VBLL models have
 692 more predictions with high predictive uncertainties, which we hypothesize is an affect of the heavier
 693 tailed predictive distributions. Finally, we find that even though the Het-VBLL models appear better
 694 calibrated without post-hoc calibration, their predictive uncertainty rankings are very similar to the
 695 other VBLL models.

696 We train 3 seeds for all configurations and report mean and ± 2 standard error of the mean in all QM9
 697 metrics and plots. For the DNN models we use all the same training parameters as in Schütt et al.
 698 (2021) except we set our weight decay to $3e - 3$. In the VBLL models, we use a dense covariance for
 699 the last layer, a Wishart scale of $1e2$ and a prior scale of 1. For the t-VBLL models, we use a dense
 700 covariance, with degrees of freedom prior of 1, Wishart scale of 0.5 and prior scale of 2.0. For the
 701 Het-VBLL models we use a prior scale of 1, noise prior scale of 0.1 and grad correction of 1. We
 702 also use gradient clipping in the Het-VBLL models for stability. All models are trained with AdamW
 703 optimizer - as in Schütt et al. (2021) we use a reduce on plateau learning rate scheduler for the MAP
 704 estimated PaiNN network but use cosine annealing with VBLL, t-VBLL and Het-VBLL due to the
 705 validation loss being more volatile in the VBLL models.

706 We train the DNN, VBLL and t-VBLL models for 650 epochs and the Het-VBLL models for 1000
 707 epochs, within which they all have converged on MAE, which is the accuracy metric of interest for
 708 molecular property prediction datasets. The models were each trained (with random allocation) on a
 709 single GPU using an internal CUDA HPC cluster. The hardware contained in the cluster includes:
 710 RTX A4000, RTX 2080 Ti, RTX A5000, GTX 1080 Ti, RTX 4070 Ti SUPER, Titan Xp, Titan
 711 V, L40S. Although training time varies with allocated GPU, running a single seed for any of the
 712 aforementioned methods takes approximately 15 hours. Based on this, reproducing the full QM9
 713 results would take approximately 180 GPU hours. We do however note that tuning prior values and
 714 weight decay parameters cost additional compute.

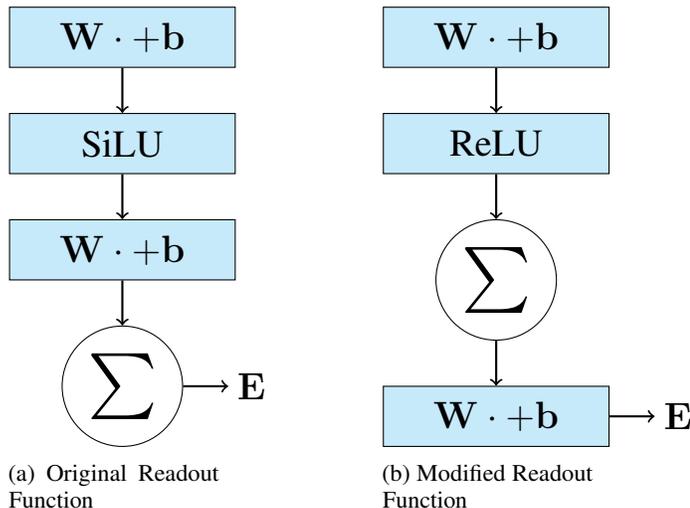


Figure 6: Original (a) and modified (b) readout functions for PaiNN network

Table 5: Results on CIFAR-100

Method	Feature Weights	Accuracy (\uparrow)	ECE (\downarrow)	NLL (\downarrow)	SVHN OOD (\uparrow)	CIFAR-10 OOD (\uparrow)
DNN	Trained	80.4 \pm 0.29	0.107 \pm 0.004	0.941 \pm 0.016	0.799 \pm 0.020	0.795 \pm 0.001
SNGP	Trained	80.3 \pm 0.23	0.030 \pm 0.004	0.761 \pm 0.007	0.846 \pm 0.019	0.798 \pm 0.001
VBLL	Trained	80.7 \pm 0.03	0.040 \pm 0.002	0.913 \pm 0.011	0.849 \pm 0.006	0.791 \pm 0.003
VBLL	Frozen	80.7 \pm 0.02	0.063 \pm 0.000	0.831 \pm 0.005	0.843 \pm 0.001	0.804 \pm 0.001
t-VBLL	Frozen	81.3 \pm 0.10	0.039 \pm 0.001	0.782 \pm 0.004	0.681 \pm 0.017	0.811 \pm 0.003
Het-VBLL	Frozen	81.2 \pm 0.07	0.039 \pm 0.001	0.777 \pm 0.005	0.685 \pm 0.001	0.811 \pm 0.001
LLLA	Frozen	80.4 \pm 0.29	0.210 \pm 0.018	1.048 \pm 0.014	0.834 \pm 0.014	0.811 \pm 0.002
VBLL	Fine-Tuned	80.8 \pm 0.11	0.047 \pm 0.001	0.783 \pm 0.002	0.779 \pm 0.008	0.794 \pm 0.003
t-VBLL	Fine-Tuned	80.2 \pm 0.25	0.035 \pm 0.001	0.810 \pm 0.006	0.815 \pm 0.032	0.791 \pm 0.003
Het-VBLL	Fine-Tuned	80.2 \pm 0.12	0.048 \pm 0.002	0.823 \pm 0.003	0.838 \pm 0.006	0.805 \pm 0.003

715 E.2 Supervised Classification

716 E.2.1 Half Moon Dataset

717 For the two-moon data set in Figure 2, we use `sklearn` to generate the data and set the noise level to
718 0.25. We introduce input-dependent label noise by flipping 20% of the labels between $[0.5, 1]$ for
719 feature 1. For all baselines, we use the same backbone configuration consisting of two hidden layers
720 with 128 neurons ($N_\phi = 128$) and ELU activations. Further, we choose a learning rate of $1e - 3$
721 and a weight decay of $1e - 4$ and train all models for 1000 epochs. For the standard VBLL and the
722 t-VBLL, we use a prior scale of 1 and a Wishart scale of 1. For the Het-VBLL, we choose a noise
723 prior scale of 0.1 and a prior scale of 1.

724 E.2.2 CIFAR 10 and CIFAR 100

725 For the CIFAR 10 and CIFAR 100 datasets, we compare our models against other last-layer methods,
726 swapping out the classification head of a pretrained and frozen Wide ResNet-28-10 network following
727 Liu et al. (2022). These experiments match the settings used in Harrison et al. (2024), and we compare
728 against their reported baselines. In particular we compare against standard VBLL (Harrison et al.,
729 2024) and last-layer Laplace (Daxberger et al., 2021) methods, as well as the results for the standard
730 (vanilla) network. Performance is evaluated through accuracy, calibration error (ECE), negative log-
731 likelihood (NLL) and out-of-distribution AUROC for both near and far OOD detection capabilities.
732 For CIFAR-10, we assess near OOD performance using CIFAR-100 as the out-of-distribution dataset
733 and vice versa for the CIFAR-100 OOD evaluation. In both cases, we utilize Street View House
734 Numbers (SVHN) (Netzer et al., 2011) as a far OOD dataset. For t-VBLL we set a prior scale of 1
735 and an inverse Gamma scale parameter of 10. For Het-VBLL we select a noise prior scale of 0.01 and
736 prior scale of 10. All models are trained using the AdamW optimizer and a learning rate of $1e-3$ with
737 a linear warmup. During fine-tuning, the classification head is first trained for two warmup epochs,
738 after which the encoder unfrozen the entire network is jointly optimized. We train 3 seeds for all
739 configurations and report mean and ± 1 standard error of the mean for all metrics.

740 E.3 Semantic Segmentation

741 VBLL SegFormer and Het-VBLL SegFormer were trained and implemented using the MMSegmen-
742 tation (MMSegmentation, 2020) framework. We employed an AdamW optimizer with a learning rate
743 of $6e-5$ for 160,000 iterations. A linear learning rate warm up period was followed by polynomial
744 learning rate decay schedule, according to the standard training protocol for this architecture on
745 CityScapes implemented by MMSegmentation. For Het-VBLL Segformer we select a noise prior
746 scale of 0.01 and prior scale of 1.0. VBLL Segformer utilizes a prior scale of 1.0 and a Wishart scale
747 of 10.0.

748 E.4 Bayesian Optimization

749 All models are implemented using GPyTorch (Gardner et al., 2018) and BoTorch (Balandat et al.,
750 2020). For the Student-t predictive of the t-VBLLs, we directly use the standard deviation of the
751 predictive in UCB. For the Het-VBLL model, we construct this standard deviation by sampling ten
752 Σ from the variational posterior as discussed in Sec. C.3 and then use their mean in the acquisition
753 function. For the outliers, we add heavy-tailed noise from a zero mean Laplace distribution with a
754 standard deviation of 0.5 (Pest Control), 1 (Ackley), or 5 (Lunar Lander) to the output of an experiment

Table 6: Results on CityScapes Semantic Segmentation - Validation

IoU	SegFormer-B0	VBLL SegFormer-B0	Het-VBLL SegFormer-B0
Road	0.980	0.980	0.982
Sidewalk	0.840	0.843	0.847
Building	0.922	0.922	0.923
Wall	0.588	0.628	0.587
Fence	0.568	0.564	0.552
Pole	0.626	0.613	0.630
Traffic Light	0.698	0.685	0.698
Traffic Sign	0.777	0.764	0.782
Vegetation	0.926	0.924	0.926
Terrain	0.641	0.642	0.639
Sky	0.950	0.948	0.950
Person	0.807	0.799	0.808
Rider	0.574	0.574	0.587
Car	0.943	0.942	0.945
Truck	0.694	0.778	0.719
Bus	0.839	0.841	0.856
Train	0.757	0.759	0.762
Motorcycle	0.649	0.631	0.650
Bicycle	0.763	0.754	0.761
mIoU	0.765	0.768	0.769

755 according to the outlier probability of the experiment. To calculate the minimum instantiations regret,
756 we approximate the optimal value of Pest Control as 12 and for Lunar Lander as 300. For Ackley, the
757 known optimal value is 0.

758 **VBLL Baselines** For all VBLL baselines, we use 3 hidden layers with 64 neurons and ELU
759 activations and set the maximum training epochs to 10000. For all models but the Het-VBLL, we use
760 a patience of 100 (Brunzema et al., 2025). For the Het-VBLL, we further set the gradient correction
761 scale to 1 as in Appendix E.1.3. All models are trained using AdamW (Loshchilov & Hutter, 2017)
762 and we use a gradient clipping of 1.

763 **Variance Predictive Baseline** For the variance predictive baseline, we use the same backbone as
764 for the VBLL baselines, i.e., 3 hidden layers with 64 neurons and ELU activations, and further adopt
765 the same optimization routine above.

766 **GP Baselines** For all the GP baselines, we choose a Matérn kernel with $\nu = 2.5$ and use individual
767 lengthscales ℓ_i for all input dimensions. Our standard GP baseline, GP, optimizes these individual
768 lengthscales within the box constraints $\ell_i \in [0.005, 4]$ as in (Eriksson et al., 2019). Furthermore, we
769 add a GP with D -scaled priors (Hvarfner et al., 2024), D-GP, which have demonstrated increased
770 sample efficiency on some tasks in previous work. Lastly, we add a GP with a Yeo-Johnson (YJ)
771 outcome transformation to cope with heteroskedasticity as proposed in (Cowen-Rivers et al., 2022),
772 YJ-GP, as a baseline. We refit the YJ transformation at each time step and use the same lengthscales
773 configuration as for the base GP baseline. For all GP baselines, we learn the hyperparameter of the
774 GP at each iteration through maximum likelihood estimation.

775 **Computational Resources** All BO simulations were performed on an HPC cluster with Intel Xeon
776 8468 Sapphire at 2.1 GHz using 4 cores. The average wall clock time for Lunar Lander was for
777 GP 1.85 hours, for VBLL 9.02 hours, for t-VBLL 7.2 hours, and for het-VBLL 20 hours. It is
778 however important to note that surrogate fit time is not of high importance in the context of BO as the
779 experiment is usually assumed to be the time consuming bottleneck.