

---

# A Mutual Information Lower Bound for Multimodal Regression Active Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Active learning for continuous regression has lacked an acquisition function that  
2 targets epistemic uncertainty when the predictive distribution is multimodal: vari-  
3 ance misses modal disagreement, and information-theoretic targets like BALD  
4 are designed for discrete outputs. We introduce a Two-Index framework that  
5 makes this separation explicit: one stochastic index selects among competing  
6 model hypotheses (epistemic source), while a second governs within-hypothesis  
7 randomness (aleatoric source). An entropy decomposition within the framework  
8 identifies the mutual information between the output and the epistemic index as a  
9 principled acquisition objective, and we prove this quantity vanishes as the model  
10 is trained on growing datasets, confirming that it captures exactly the uncertainty  
11 data can resolve. Because this mutual information is intractable for continuous  
12 outputs, we derive the Mutual Information Lower Bound (MI-LB) acquisition  
13 function, a closed-form approximation for Mixture Density Network ensembles.  
14 On benchmarks featuring multimodal systems, MI-LB matches or beats every  
15 baseline evaluated and is the only method to do so consistently – geometric and  
16 Fisher-based baselines compete only when the input space already encodes the  
17 multimodality, and collapse otherwise.

## 18 1 Introduction

19 Active learning in continuous, multivariate settings requires acquisition functions that can distinguish  
20 uncertainty the model can reduce by collecting more data (epistemic) from uncertainty intrinsic to  
21 the data-generating process (aleatoric). Existing approaches either operate in restricted settings, such  
22 as discrete classification [1] or unimodal Gaussian outputs [2], or rely on scalar summaries like  
23 predictive variance that conflate the two sources. The problem is worse in multimodal settings: two  
24 model hypotheses may assign the same mean and variance to an output while disagreeing on the  
25 number and location of modes. Variance alone cannot detect this kind of distributional disagreement.

26 We address this gap with a framework that keeps the two sources of uncertainty separate by con-  
27 struction. The framework extends the Epistemic Neural Network formalism [3] by introducing  
28 two independent stochastic indices: an epistemic index  $Z \sim P_Z$  that parameterizes a family of  
29 hypotheses, and an aleatoric index  $\epsilon \sim P_\epsilon$  that governs within-hypothesis stochasticity. A learned  
30 map  $g_\theta(x; z, \epsilon)$  transforms these indices and the input into the output space. This approach provides a  
31 common vocabulary for a broad class of uncertainty quantification methods, facilitating the exchange  
32 of analyses among model families.

33 The two-index structure yields an entropy decomposition that separates predictive uncertainty into an  
34 epistemic term, the mutual information  $I(Y; Z | x)$ , and an aleatoric term, the expected conditional  
35 entropy  $\mathbb{E}_Z[H(Y | x, Z)]$ . Despite being a natural acquisition target, the mutual information requires  
36 computing differential entropies that are intractable for general continuous outputs. When the

37 model family is an ensemble of Mixture Density Networks (MDNs), the predictive and conditional  
38 distributions are both Gaussian mixtures, and known upper and lower bounds on Gaussian-mixture  
39 entropy [4] can be combined to produce a tractable lower bound on  $I(Y; Z | x)$ . We call this the  
40 *Mutual Information Lower Bound* (MI-LB) acquisition function. Across three benchmarks with  
41 controllable multimodality, MI-LB matches or beats every baseline evaluated and is the only method  
42 to do so consistently. Geometric and Fisher-based baselines compete only when the input geometry  
43 already encodes the multimodal structure, and collapse otherwise. To the best of our knowledge,  
44 MI-LB is the first acquisition function specifically designed for multimodal continuous settings.

45 The main contributions of this paper are:

- 46 • The MI-LB acquisition function for active learning in continuous multivariate settings,  
47 with a proof that it lower bounds the true epistemic mutual information  $I(Y; Z | x)$  and a  
48 closed-form expression that is efficient to evaluate under standard MDN assumptions.
- 49 • The Two-Index framework for disentangling epistemic and aleatoric uncertainty, applicable  
50 to a broad class of model families, with a proof that the epistemic term  $I(Y; Z | x)$   
51 vanishes with sufficient data under standard well-specification and consistency conditions  
52 (Appendix D).
- 53 • Experiments on three multimodal benchmarks showing that MI-LB matches or beats every  
54 baseline evaluated (Random, Variance, BAIT, Core-Set), and is the only method to do so  
55 consistently – geometric and Fisher-based baselines compete only when the input geometry  
56 already encodes the multimodal structure, and collapse when the modal disagreement lives  
57 in output space.

## 58 1.1 Related Work

59 Decomposing predictive uncertainty into aleatoric and epistemic components has been a longstanding  
60 goal in the machine learning literature [2, 5]. Foundational approaches include Bayesian Neural  
61 Networks [6], Deep Ensembles [7], and the Epistemic Neural Network framework [3], which  
62 provides a unified interface that does not require explicit Bayesian inference [5, 3]. Despite these  
63 advances, recent benchmarking shows that practical disentanglement remains elusive: modern  
64 evidential, variational, and deterministic methods exhibit rank correlations often exceeding 0.94  
65 between their aleatoric and epistemic estimates [8]. This problem is particularly pronounced in  
66 continuous, multimodal regression tasks, where fitting a single Gaussian to diverse target modes  
67 inflates predictive variance, motivating the use of Mixture Density Networks (MDNs) [9, 10].

68 In the context of active learning, selecting informative data points to minimize labeling costs has  
69 traditionally relied on geometric or gradient-based heuristics [11]. Geometric methods such as  
70 Core-Set minimize L2 distances in embedding space [12], while Fisher-based methods like BAIT  
71 optimize bounds on the maximum likelihood estimator error using network gradients [13]. Both have  
72 important failure modes in multimodal settings: Core-Set ensures coverage in the learned embedding  
73 space and can fail when multimodality manifests in the output distribution rather than in the geometry  
74 of the inputs [12], and BAIT’s last-layer Fisher embedding provides poor candidate rankings when the  
75 predictive likelihood is highly multimodal, since the gradient signal at a single Monte-Carlo sample is  
76 dominated by within-mode variance. In these regimes, black-box acquisition strategies based on the  
77 empirical predictive covariance of ensembles often outperform white-box gradient methods [14, 11].

78 Information-theoretic approaches such as Bayesian Active Learning by Disagreement (BALD) use  
79 mutual information to target epistemic uncertainty directly [1, 15]. BALD is effective in discrete  
80 classification, but extending it to continuous regression requires evaluating differential entropy for  
81 distributions like Gaussian mixtures, which lacks a closed-form solution [4].

## 82 2 The Two-Index Approach for Disentangling Uncertainties

### 83 2.1 The Two-Index Generative Model

84 We build on the Epistemic Neural Network (ENN) framework [3] and the distribution network  
85 formalism of [10] to construct a unified model that explicitly represents both epistemic and aleatoric  
86 uncertainty through two independent stochastic indices.

87 In the ENN framework, a model is specified by a parameterized function  $f_\theta$  and a fixed reference  
 88 distribution  $P_Z$ . An *epistemic index*  $Z \sim P_Z$  is sampled and passed to the network alongside the  
 89 input  $x$ , producing a prediction  $f_\theta(x; Z)$ . Crucially,  $P_Z$  is not updated during training. Instead,  
 90 training optimizes  $\theta$  so that the mapping from  $Z$  to predictions becomes increasingly invariant under  
 91  $P_Z$  as more data is collected. Sensitivity of predictions to  $Z$  reflects epistemic uncertainty: the model  
 92 has not yet identified a unique function consistent with the data.

93 We extend this framework by introducing the *aleatoric index*  $\epsilon \sim P_\epsilon$ , a second fixed reference  
 94 distribution independent of both  $Z$  and  $\theta$ . For a given epistemic index  $Z = z$  and input  $x$ , the model  
 95 defines a conditional distribution over outputs  $Y$  through the generative mapping

$$Y = g_\theta(x; z, \epsilon), \quad Z \sim P_Z, \quad \epsilon \sim P_\epsilon, \quad Z \perp \epsilon. \quad (1)$$

96 Here  $g_\theta$  is a learned transformation that maps the input  $x$  and two independent sources of randomness  
 97 into the output space. The role of  $\epsilon$  is to represent irreducible stochasticity: for any fixed  $z$ , sampling  
 98  $\epsilon \sim P_\epsilon$  produces draws from the aleatoric conditional distribution  $p_\theta(y | x, z)$ . The role of  $Z$  is to  
 99 represent reducible uncertainty: variation in predictions across different realizations of  $Z$  reflects  
 100 hypotheses about the data-generating process that have not yet been eliminated by the available data.  
 101 The network  $g_\theta$  learns a transformation from the tractable distributions  $P_Z$  and  $P_\epsilon$  into a flexible,  
 102 potentially multimodal distribution over the output  $Y$ .

103 **Induced distributions.** The joint generative model in (1) induces a hierarchy of distributions, made  
 104 precise using the push-forward operator. For a measurable space  $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$  and a measurable map  
 105  $h : \mathcal{E} \rightarrow \mathcal{Y}$ , the *push-forward* of a measure  $\mu$  on  $\mathcal{E}$  under  $h$  is the measure  $h_{\#}\mu$  on  $\mathcal{Y}$  defined by

$$(h_{\#}\mu)(A) := \mu(\{\epsilon \in \mathcal{E} : h(\epsilon) \in A\}), \quad A \in \mathcal{B}(\mathcal{Y}). \quad (2)$$

106 For a fixed epistemic realization  $z \in \mathcal{Z}$  and input  $x \in \mathcal{X}$ , the *aleatoric conditional distribution* is  
 107 defined as the push-forward of  $P_\epsilon$  under the map  $g_\theta(x; z, \cdot) : \mathcal{E} \rightarrow \mathcal{Y}$ :

$$p_\theta(\cdot | x, z) := (g_\theta(x; z, \cdot))_{\#}P_\epsilon. \quad (3)$$

108 Marginalizing equation (3) over the epistemic index yields the *predictive distribution*

$$p_\theta(\cdot | x) := \int_{\mathcal{Z}} p_\theta(\cdot | x, z) dP_Z(z). \quad (4)$$

109 The predictive distribution conflates both sources of uncertainty. Disentangling them requires tracking  
 110 the two indices separately, as formalized in the decompositions of Section 2.2.

111 **Interpretation.** Each realization  $Z = z$  instantiates a *hypothesis*  $p_\theta(\cdot | x, z)$  for the true conditional  
 112  $p^*(\cdot | x)$ ; within a fixed hypothesis,  $\epsilon$  governs irreducible stochasticity. Disagreement across  
 113 realizations of  $Z$  is the origin of epistemic uncertainty. The predictive distribution has the same  
 114 mathematical form as Bayesian Model Averaging [16], but the framework does not require a strict  
 115 Bayesian formulation of  $P_Z$ : as illustrated in Appendix B, it accommodates ensembles, BNNs,  
 116 conditional flow matching, and VAEs under a common vocabulary.

## 117 2.2 Quantifying the Two Sources of Uncertainty

118 The two-index structure enables exact decompositions of predictive uncertainty into its epistemic and  
 119 aleatoric components. We provide two such decompositions via variance and entropy.

120 **Variance-Based Decomposition.** In the scalar output ( $Y \in \mathbb{R}$ ) case [5], the law of total variance  
 121 yields an exact additive decomposition. Writing  $Y = g_\theta(x; Z, \epsilon)$  with  $Z \perp \epsilon$ ,

$$\underbrace{\text{Var}_{Z, \epsilon}(Y | x)}_{\text{total variance}} = \underbrace{\mathbb{E}_Z [\text{Var}_\epsilon(Y | Z, x)]}_{\text{aleatoric variance}} + \underbrace{\text{Var}_Z(\mathbb{E}_\epsilon[Y | Z, x])}_{\text{epistemic variance}}. \quad (5)$$

122 The aleatoric term  $\mathbb{E}_Z[\text{Var}_\epsilon(Y | Z, x)]$  averages the output variance within each hypothesis over  
 123 epistemic realizations; it measures the average irreducible spread. The epistemic term  $\text{Var}_Z(\mathbb{E}_\epsilon[Y |$   
 124  $Z, x])$  measures the variance of the conditional mean across hypotheses; it captures disagreement  
 125 between different epistemic realizations of the model.

126 This decomposition is exact and interpretable, but it summarizes uncertainty through second moments  
 127 alone. As noted in the introduction to this section, two hypotheses can share the same conditional  
 128 mean and variance while assigning probability mass to entirely different modes. The variance-based  
 129 epistemic term would report an incomplete summary in such a case. For multimodal or non-Gaussian  
 130 distributions, a richer summary is needed, as further discussed in Appendix C.

131 **Entropy-Based Decomposition.** We complement (equation 5) with an entropy-based decomposi-  
 132 tion that captures distributional disagreement beyond second moments. Define the total predictive  
 133 uncertainty as the differential entropy of  $p_\theta(y | x)$ :

$$H_{Z,\epsilon}(Y | x) = \mathbb{E}_{Z,\epsilon}[-\log p_\theta(Y | x)]. \quad (6)$$

134 Applying the chain rule of mutual information, this decomposes as

$$H_{Z,\epsilon}(Y | x) = \underbrace{\mathbb{E}_Z[H_\epsilon(Y | Z, x)]}_{\text{aleatoric uncertainty}} + \underbrace{I(Y; Z | x)}_{\text{epistemic uncertainty}}, \quad (7)$$

135 where  $I(Y; Z | x)$  is the *mutual information* between the output  $Y$  and the epistemic index  $Z$ ,  
 136 conditioned on the input  $x$ . The aleatoric term  $\mathbb{E}_Z[H_\epsilon(Y | Z, x)]$  is the expected entropy of  
 137 the conditional distribution  $p_\theta(\cdot | x, z)$ , averaged over the epistemic index. The epistemic term  
 138  $I(Y; Z | x)$  quantifies how much knowing  $Z$  reduces uncertainty about  $Y$ : it is large when different  
 139 epistemic indices yield meaningfully different conditional distributions, and vanishes when all  
 140 hypotheses agree. Because mutual information is sensitive to the full shape of each hypothesis  
 141 distribution, it detects the kind of modal disagreement that variance misses.

### 142 2.3 Asymptotic Guarantees

143 The decompositions (5) and (7) have the correct asymptotic behavior by construction. As the size  
 144 of the training dataset grows,  $\theta$  is optimized to produce predictions that are consistent with the data  
 145 under all probable realizations of  $Z$ . In the limit of infinite data and a well-specified model, the  
 146 learned  $\theta$  renders  $g_\theta(x; z, \epsilon)$  approximately independent of  $z$  for  $P_Z$ -almost all  $z$ , so that:

$$I(Y; Z | x, \mathcal{D}_n) \xrightarrow{n \rightarrow \infty} 0, \quad \mathbb{E}_Z[H_\epsilon(Y | Z, x, \mathcal{D}_n)] \xrightarrow{n \rightarrow \infty} H(p^*(\cdot | x)), \quad (8)$$

147 where  $p^*(\cdot | x)$  is the true conditional distribution of the data-generating process. Epistemic  
 148 uncertainty vanishes as the model identifies the correct hypothesis, while aleatoric uncertainty  
 149 converges to the irreducible entropy of the true process. Conversely, at an out-of-distribution input  
 150  $x_{\text{OOD}}$ , training data provides no constraints on  $Z$ , so  $I(Y; Z | x_{\text{OOD}})$  can remain large, providing a  
 151 principled basis for out-of-distribution detection. We note that the convergence of the aleatoric term  
 152 to  $H(p^*(\cdot | x))$  requires that the aleatoric model class — the family of distributions representable  
 153 by  $g_\theta(x; z, \cdot)$  for fixed  $z$  — be sufficiently expressive to capture  $p^*$ . This is a non-trivial condition  
 154 that motivates the use of flexible density estimators, such as Mixture Density Networks or implicit  
 155 generative models, rather than restricting to unimodal Gaussian outputs.

156 This is stated informally as the following theorem, which is made precise and proved in Appendix D:

157 **Theorem 2.1** (Informal: Epistemic Uncertainty Vanishes with Data). *If a model can represent the*  
 158 *true data-generating process  $p^*(\cdot | x)$  and its parameters converge to the correct values as the*  
 159 *dataset grows to infinity, then  $I(Y; Z | x) \rightarrow 0$ . In other words, a well-trained model’s predictions*  
 160 *eventually agree across all epistemic index values, and all remaining uncertainty is aleatoric.*

### 161 2.4 Implications for Data Acquisition

162 The decomposition from (7) and Thm. 2.1 together identify  $I(Y; Z | x)$  as the natural acquisition  
 163 target for active learning. It isolates exactly the component of predictive uncertainty that additional  
 164 data can reduce: epistemic disagreement among hypotheses. The variance-based epistemic term  
 165 from (5) targets the same goal in principle, but summarizes disagreement through second moments  
 166 alone; it can miss regions where hypotheses differ in modal structure rather than spread. The mutual  
 167 information  $I(Y; Z | x)$ , by contrast, is sensitive to the full distribution of each hypothesis, making it  
 168 the right objective when the conditional distribution  $p^*(y | x)$  may be multimodal.

169 The remaining challenge is computational. Evaluating  $I(Y; Z | x)$  requires computing differential  
 170 entropies of the predictive and conditional distributions, which have no closed-form expression for

171 general continuous outputs. In the next section, we show that when the model family is an ensemble  
 172 of Gaussian mixtures, analytical entropy bounds can be used to construct a tractable lower bound on  
 173  $I(Y; Z | x)$  that serves as an effective acquisition function.

### 174 3 The Mutual Information Lower Bound (MI-LB) Acquisition Function

#### 175 3.1 Entropy Intractability in Continuous Settings

176 Re-arranging the entropy-based decomposition from (7), the mutual information can be written as

$$I(Y; Z | x) = \underbrace{H(Y | x)}_{\text{(i) marginal entropy}} - \underbrace{\mathbb{E}_Z[H(Y | x, Z)]}_{\text{(ii) expected aleatoric entropy}}. \quad (9)$$

177 Computing this quantity requires evaluating two differential entropy terms: term (i), the entropy of  
 178 the predictive distribution  $p_\theta(\cdot | x)$  from equation 4, and term (ii), the entropy of  $Y | x, Z = z$  for a  
 179 fixed  $z$ , averaged over  $P_Z$ . For *implicit* distribution networks [10] — models that produce samples of  
 180  $Y$  via the map  $g_\theta(x; z, \epsilon)$  but do not provide an analytical form for  $p_\theta(\cdot | x, z)$  — both terms must be  
 181 estimated from samples using kernel density estimators, which may scale poorly with the dimension  
 182 of  $\mathcal{Y}$ . For *explicit* distribution networks, the density  $p_\theta(\cdot | x, z)$  is available analytically, enabling  
 183 structured approximation of both terms. We exploit this structure in the case of Mixture Density  
 184 Networks, where each  $p_\theta(\cdot | x, z)$  is a mixture of Gaussians.

#### 185 3.2 MDN Structure and the Marginal Mixture

186 Assume  $Z$  takes values in the discrete index set  $\{1, \dots, n_{\text{ens}}\}$  with  $P_Z(Z = z) = w_z$ , and suppose  
 187 that for each  $z$ , the aleatoric conditional  $p_\theta(\cdot | x, z)$  is a Gaussian mixture with  $K$  components:

$$p_\theta(y | x, z) = \sum_{i=1}^K \alpha_i^{(z)}(x) \mathcal{N}\left(y; \mu_i^{(z)}(x), C_i^{(z)}(x)\right), \quad (10)$$

188 where  $\alpha_i^{(z)}(x) > 0$ ,  $\sum_i \alpha_i^{(z)}(x) = 1$ ,  $\mu_i^{(z)}(x) \in \mathbb{R}^N$ , and  $C_i^{(z)}(x) \in \mathbb{R}^{N \times N}$  is a symmetric positive  
 189 definite covariance matrix. Here  $N$  denotes the dimension of the output space  $\mathcal{Y} \subseteq \mathbb{R}^N$ . Under this  
 190 formulation, the predictive marginal distribution (4) is

$$p_\theta(y | x) = \sum_{z=1}^{n_{\text{ens}}} w_z p_\theta(y | x, z) = \sum_{z=1}^{n_{\text{ens}}} \sum_{i=1}^K \underbrace{w_z \alpha_i^{(z)}(x)}_{=: \beta_{z,i}(x)} \mathcal{N}\left(y; \mu_i^{(z)}(x), C_i^{(z)}(x)\right), \quad (11)$$

191 which is a Gaussian mixture with  $n_{\text{ens}} \cdot K$  components and weights  $\beta_{z,i}(x) = w_z \alpha_i^{(z)}(x)$ . Both  
 192 terms in (9) therefore require computing the differential entropy of a Gaussian mixture, for which no  
 193 closed-form expression exists.

#### 194 3.3 Entropy Bounds for Gaussian Mixtures

195 We make use of two bounds from [4]. Let  $p(y) = \sum_{i=1}^K \pi_i \mathcal{N}(y; \mu_i, C_i)$  be a Gaussian mixture  
 196 with  $K$  components, weights  $\pi_i$ , means  $\mu_i \in \mathbb{R}^N$ , and covariances  $C_i \in \mathbb{R}^{N \times N}$ . [4] shows that the  
 197 differential entropy  $H(p)$  satisfies  $H_{\text{lower}} \leq H(p) \leq H_{\text{upper}}$ , where

$$H_{\text{lower}} = - \sum_{i=1}^K \pi_i \log \left( \sum_{j=1}^K \pi_j \mathcal{N}(\mu_i; \mu_j, C_i + C_j) \right), \quad (12)$$

$$H_{\text{upper}} = \sum_{i=1}^K \pi_i \left( -\log \pi_i + \frac{1}{2} \log((2\pi e)^N |C_i|) \right). \quad (13)$$

198 Both bounds are computable in  $O(K^2)$  operations. The upper bound (13) decomposes as the  
 199 weighted sum of the entropy of individual modes and is tight when the components have negligible  
 200 overlap. The lower bound (12) accounts for inter-component similarity through pairwise Gaussian  
 201 evaluations at the component means. Under the standard MDN assumption of diagonal covariances  
 202  $C_i = \text{diag}(\sigma_{i,1}^2, \dots, \sigma_{i,N}^2)$ , the Gaussian evaluations in (12) and the determinants in (13) reduce to  
 203 sums of univariate quantities, making both bounds particularly efficient to evaluate.

204 **3.4 The MI-LB Acquisition Function**

205 We define the *Mutual Information Lower Bound* (MI-LB) acquisition function as

$$MI-LB(x) := H_{\text{lower}}(Y | x) - \sum_{z=1}^{n_{\text{ens}}} w_z H_{\text{upper}}(Y | x, Z = z), \quad (14)$$

206 where  $H_{\text{lower}}(Y | x)$  is the lower bound (12) applied to the  $n_{\text{ens}} \cdot K$ -component marginal mixture (11)  
 207 with component weights  $\beta_{z,i}(x)$ , means  $\mu_i^{(z)}(x)$ , and covariances  $C_i^{(z)}(x)$ ; and  $H_{\text{upper}}(Y | x, Z = z)$   
 208 is the upper bound (13) applied to the  $K$ -component conditional mixture (10) for each  $z$ . Explicitly,

$$MI-LB(x) = - \sum_{z=1}^{n_{\text{ens}}} \sum_{i=1}^K \beta_{z,i} \log \left( \sum_{l=1}^{n_{\text{ens}}} \sum_{j=1}^K \beta_{l,j} \mathcal{N} \left( \mu_i^{(z)}; \mu_j^{(l)}, C_i^{(z)} + C_j^{(l)} \right) \right) \\ - \sum_{z=1}^{n_{\text{ens}}} w_z \sum_{i=1}^K \alpha_i^{(z)} \left( -\log \alpha_i^{(z)} + \frac{1}{2} \log \left( (2\pi e)^N |C_i^{(z)}| \right) \right), \quad (15)$$

209 where we suppress the dependence on  $x$  for readability. The fundamental property of MI-LB is that  
 210 it constitutes a certified lower bound on the true epistemic uncertainty  $I(Y; Z | x)$ , which we state  
 211 formally as follows, with a proof given in Appendix E.

212 **Theorem 3.1** (MI-LB is a lower bound on mutual information). *Under the model assumptions of*  
 213 *Section 2, with  $Z$  supported on  $\{1, \dots, n_{\text{ens}}\}$  and  $p_{\theta}(\cdot | x, z)$  a Gaussian mixture of the form (10)*  
 214 *for each  $z$ , it holds that*

$$MI-LB(x) \leq I(Y; Z | x) \quad \text{for all } x \in \mathcal{X}. \quad (16)$$

215 The practical consequence of Theorem 3.1 is that maximizing MI-LB over a candidate pool is a  
 216 conservative acquisition strategy: any point selected by MI-LB is guaranteed to have true epistemic  
 217 uncertainty at least as large as the reported score. This mirrors the logic of the Evidence Lower  
 218 Bound (ELBO) in variational inference [17], where optimizing a lower bound on the log-evidence  
 219 provides a principled surrogate without risking false certification of a poor solution. Consequently,  
 220 the MI-LB acquisition rule selects

$$x^* = \arg \max_{x \in \mathcal{X}_{\text{pool}}} MI-LB(x). \quad (17)$$

221 **Connection to BALD.** The Bayesian Active Learning by Disagreement (BALD) acquisition [1]  
 222 is defined identically to (9) in the discrete classification setting, where  $Y$  follows a categorical  
 223 distribution and entropy is Shannon entropy, admitting exact computation. MI-LB can therefore be  
 224 understood as a continuous, multivariate alternative to BALD for regression, with the [4] entropy  
 225 bounds replacing exact entropy to restore tractability in the absence of a finite output alphabet.

226 **4 Experiments**

227 We benchmark MI-LB, the acquisition function proposed in Section 3, against four baselines on  
 228 three multimodal problems: **Random**, **Epistemic Variance**, **BAIT** [13] (last-layer Fisher trace),  
 229 and **Core-Set** [12] ( $k$ -Center-Greedy on the MDN backbone). SBAL [14] and MaxDist [11] batch  
 230 variants of Variance and MI-LB are deferred to the appendix. Among the three benchmarks, §4.2 (the  
 231 coupled double-well) is the discriminating one: input geometry  $(\sigma, \kappa)$  does not encode the modal  
 232 structure, which lives in output space via Kramers escape. §4.1 and §4.3 have phase boundaries in  
 233 the input space and serve to characterise where geometric baselines remain competitive.

234 **4.1 Synthetic Multimodal Conditional Problem**

235 We construct a synthetic conditional distribution  $p^*(y | x)$  as a mixture of  $K=3$  Gaussians:

$$p^*(y | x) = \sum_{k=1}^K \pi_k(x) \mathcal{N}(y; \mu_k(x), \Sigma_k(x)), \quad (18)$$

Table 1: Per-benchmark settings. All used a pool of 50,000 inputs, 100 initial labels, and an ensemble of size 8.

Benchmark	dim $x$	dim $y$	$K_{\text{MDN}}$	rounds $\times$ batch	total labels	NLL*
Multimodal conditional (§4.1)	10	16	5	$20 \times 50$	1,100	22.98
Coupled double-well (§4.2)	7	20	8	$20 \times 50$	1,100	—
Ternary phases (§4.3)	8	1	4	$30 \times 15$	550	1.33

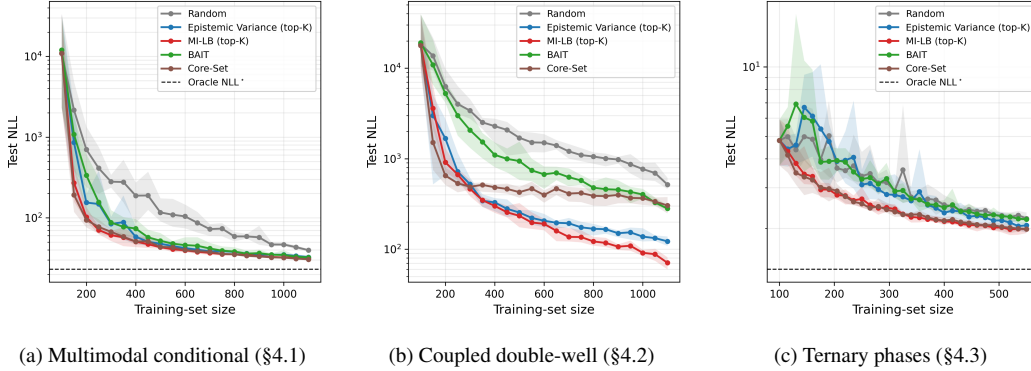


Figure 1: Test NLL vs. training-set size on each benchmark for MI-LB against Random, Epistemic Variance, BAIT, and Core-Set in top- $k$  (5 seeds; bands min-max); dashed = oracle NLL\* where available. MI-LB matches or beats every baseline on every benchmark; geometric methods (Core-Set, BAIT) are competitive only when input geometry encodes the multimodality (a, c) and collapse when it lives in output space (b).

236 where the mixing weights  $\pi_k$ , means  $\mu_k$ , and diagonal covariances  $\Sigma_k$  are input-dependent functions  
 237 (full specification in Appendix F). Inputs  $x \in \mathbb{R}^{10}$  lie on a 4-dimensional manifold embedded via  
 238  $x = \tanh(Al + b)$ ,  $l \sim \mathcal{N}(0, I_4)$ , respecting the manifold hypothesis [18], and outputs  $y \in \mathbb{R}^{16}$ .  
 239 Because  $p^*$  is known in closed form, the oracle NLL\* is computable exactly as a calibration reference.  
 240 A  $K=5$  MDN ensemble trained on the full pool reaches NLL 24.72 vs. oracle 22.98, while a  $K=1$   
 241 MDN attains only 47.76 – the conditional is multimodal and a mixture head is necessary (Fig. 2). The  
 242 mixing weights  $\pi_k(x)$  are gated by the radial coordinate  $r = \|x_{1:L}\|$ : for  $r \lesssim r_0$  a single component  
 243 dominates and  $p^*(y | x)$  is effectively unimodal, while for  $r \gtrsim r_0$  multiple components carry  
 244 comparable mass and  $p^*$  becomes genuinely multimodal, producing a sharp unimodal–multimodal  
 245 transition at  $r \approx r_0$  (full form in App. F.2).

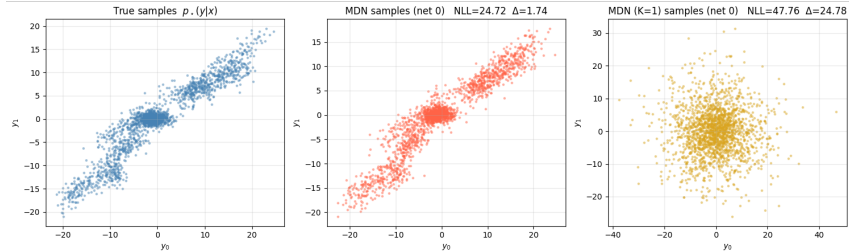


Figure 2: Predicted vs. true samples in the  $(y_0, y_1)$  plane on held-out inputs. **Left:** draws from the oracle  $p^*(y | x)$ , displaying the multimodal structure of the target conditional. **Middle:** MDN ensemble samples; the recovered geometry closely matches the oracle, with calibration gap  $\Delta = 1.74$ . **Right:** single-Gaussian MDN samples collapses to an isotropic blob that cannot represent disjoint modes, yielding  $\Delta = 24.78$ .

246 **Why this is a useful benchmark.** The coexistence of unimodal and multimodal regimes makes  
 247 this problem a natural testbed for the Two-Index framework:  $I(Y; Z | x)$  should be small in the  
 248 interior, where all ensemble members agree on a single-component conditional, and large near the  
 249 transition boundary  $\|x_{1:L}\| \approx r_0$ , where members must resolve the weights and locations of multiple  
 250 components. A well-calibrated acquisition function should therefore concentrate queries on this  
 251 boundary. As a lower bound on  $I(Y; Z | x)$ , MI-LB is designed to detect this distributional disagree-

252 ment, while variance-based scores summarize ensemble disagreement through second moments alone  
 253 and underweight regions where hypotheses differ in modal structure rather than spread.

254 **Results.** Figure 1a compares MI-LB against the four baselines. From the very first acquisition  
 255 rounds MI-LB pulls clearly ahead of Random, Variance, and BAIT and tracks the strongest baseline  
 256 (Core-Set) throughout the budget, ending within seed bands at  $31.1 \pm 0.3$  vs.  $30.6 \pm 0.6$  at  $n =$   
 257 1100. MI-LB matches the best geometric baseline on this benchmark with the tightest seed-to-seed  
 258 std among all methods, tying Core-Set under top-k and beating it under MaxDist batch selection  
 259 (Appendix F.2, Table 5). Core-Set is competitive here for a benchmark-specific reason: the input  
 260 manifold  $x = \tanh(Ax + b)$  is 4-dimensional and the radial coordinate  $\|x_{1:L}\|$  already encodes the  
 261 unimodal–multimodal gate, so k-Center-Greedy in feature space implicitly samples the boundary  
 262 that MI-LB targets explicitly. Sections 4.2–4.3 test this dependence and show that MI-LB retains  
 263 its advantage when geometry no longer encodes multimodality, whereas Core-Set degrades sharply.  
 264 Spatially, MI-LB concentrates on the unimodal–multimodal boundary, Variance is weaker on the  
 265 same region, Core-Set covers the manifold uniformly, Random spreads everywhere (App. F.2, Fig. 6).

## 266 4.2 Coupled Double-Well System

267 We consider a chain of  $P = 5$  particles evolving in coupled one-dimensional double-well potentials  
 268 under overdamped Langevin dynamics,

$$dq_i = \left[ \underbrace{a(q_i - q_i^3)}_{\text{double-well force}} + \underbrace{\kappa \sum_{j \in \text{nn}(i)} (q_j - q_i)}_{\text{nearest-neighbour coupling}} \right] dt + \sigma dW_i, \quad i = 1, \dots, P, \quad (19)$$

269 where  $dW_i$  are independent Wiener increments and we fix  $a = 1$ . Inputs encode the initial particle  
 270 configuration together with  $(\sigma, \kappa)$ ; outputs stack particle positions at four snapshot times. Full  
 271 configuration, integration scheme, and evaluation protocol in Appendix F.3.

272 **Why this is a useful benchmark.** The conditional  $p^*(y | x)$  undergoes a sharp change in mode  
 273 structure controlled by  $\sigma^2/a$ : at low  $\sigma$  each particle stays in its starting well ( $p^*$  unimodal per particle);  
 274 at high  $\sigma$ , noise triggers *Kramers escape* [19], making the single-particle marginal bimodal. Coupling  
 275  $\kappa$  raises the effective barrier for collective flips, coupling single-particle marginals into a joint  
 276 distribution over  $\{-1, +1\}^P$  basins. The Kramers exponent  $a/(2\sigma^2)$  at barrier height  $\Delta V = a/4$   
 277 places the crossover at  $\sigma \gtrsim \sqrt{a/2} \approx 0.71$ . A single-particle sweep (Fig. 3) confirms this scaling, and  
 278 at  $\sigma \gtrsim 1$  noise overwhelms the barrier so that mass spreads beyond both wells into the  $|q| > 1$  tails.

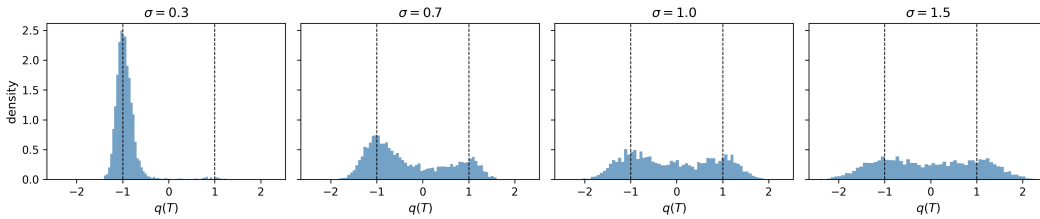


Figure 3: Terminal position  $q(T)$  histograms for an uncoupled particle ( $P = 1, \kappa = 0, q(0) = -0.5$ ) at four noise levels; dashed lines mark  $q = \pm 1$ . At  $\sigma = 0.3$  the particle stays trapped near  $q = -1$ ; at  $\sigma = 0.7 \approx \sqrt{a/2}$  Kramers escape fills both wells; for  $\sigma \geq 1$  noise dominates the barrier, spreading mass into the  $|q| > 1$  tails.

279 With  $\sigma$  and  $\kappa$  as input coordinates, the phase boundary lives directly in input space, marking where  
 280 informative acquisitions should concentrate. A Gaussian head must smear mass across both wells  
 281 in the bimodal regime, inflating test NLL: a single-Gaussian ( $K = 1$ ) head pays a 5.51-nat penalty  
 282 against the  $K = 8$  mixture head used on this benchmark (Appendix F.3, Fig. 8). This penalty is a  
 283 property of the benchmark, not of any acquisition strategy: any method built on a single-Gaussian  
 284 head fails here by construction.

285 **Results.** Figure 1b: MI-LB wins decisively at  $71 \pm 8$  at  $n = 1100$  -  $1.7\times$  better than Variance  
 286 ( $122 \pm 11$ ),  $4\times$  better than Core-Set ( $304 \pm 55$ ) and BAIT ( $281 \pm 15$ ), and  $7\times$  better than Random

287 (518 ± 60). This is the discriminating benchmark: with the bimodal regime’s mode structure living in  
288 *output* space (Kramers escape between wells), two ensemble members can agree on  $\mathbb{E}[Y | x]$  while  
289 disagreeing on whether mass concentrates near  $q = +1$ ,  $q = -1$ , or both – so feature-space coverage  
290 (Core-Set) and last-layer Fisher (BAIT) cannot rank candidates informatively, whereas MI-LB’s  
291 entropy decomposition over the full output mixture targets exactly the disagreement that matters.

### 292 4.3 Synthetic Phase-Competition Benchmark

293 Active learning is a standard tool in materials discovery, where labels are expensive and composition-  
294 process inputs live on low-dimensional manifolds [20, 21, 22]. As in the previous benchmarks,  
295 discrete phases compete across the input space, producing boundary-localised multimodality. We  
296 build a synthetic CALPHAD-style [23] benchmark that reproduces this structure.

297 Each input is a composition on the 3-component simplex (with  $x_C = 1 - x_A - x_B$ ) concatenated  
298 with  $n_{\text{proc}} = 6$  continuous process parameters. Per-phase Gibbs free energies  $G_\phi(x_A, x_B, x_C)$   
299 are random quadratic forms in composition (independent of  $p$ ), and the latent phase distribution is  
300  $\text{softmax}(-G_\phi/\tau_G)$  at temperature  $\tau_G = 0.08$ , producing sharp phase boundaries on the simplex  
301 slice. The scalar response is sampled from a per-phase Gaussian whose mean depends on both  
302 composition and process parameters; full simulator parameters in Appendix F.4.

303 **Why this is a useful benchmark.** The composition  $(x_A, x_B, x_C)$  acts as a categorical latent  
304 selecting a phase, confining multimodality to a thin boundary region on the simplex; the process  
305 subspace contributes only a smooth, phase-conditional shift. This factorisation isolates a clean test of  
306 boundary localisation (realised phase masses and boundary fraction in Appendix F.4).

307 **Results.** The tighter budget than Sections 4.1 and 4.2 reflects the 1-D scalar output, on which NLL  
308 saturates faster than the 16- and 20-dimensional outputs of the other benchmarks. Figure 1c: MI-LB  
309 reaches **1.99 ± 0.04** at  $n = 550$ , tied with Core-Set ( $1.99 \pm 0.04$ ); Variance (2.06), BAIT (2.19), and  
310 Random (2.21) trail. MaxDist on top of MI-LB narrowly wins overall ( $1.95 \pm 0.06$ ). BAIT collapses  
311 to Random-level, the 1-D scalar output makes its single-MC last-layer Fisher estimate too noisy to  
312 rank candidates. Core-Set’s parity again reflects benchmark geometry (phase boundaries lie on the  
313 composition simplex) and Sections 4.2 shows it does not transfer when multimodality is in output  
314 space. MI-LB retains the tightest seed-to-seed spread in the data-scarce regime ( $n \in [115, 250]$ ):  
315  $\sim 2.6 \times$  tighter than Random,  $\sim 4.5 \times$  than Variance (Appendix Table 11).

## 316 5 Conclusion

317 **Summary.** We introduce the Two-Index generative framework to formally disentangle epistemic  
318 and aleatoric uncertainty across diverse model families. By leveraging this formalism, we derive the  
319 MI-LB acquisition function, providing a principled and computationally efficient Mutual Information  
320 lower bound for active learning in continuous, multimodal settings. Across three benchmarks MI-LB  
321 matches or beats every baseline we evaluate (Random, Variance, BAIT, Core-Set) and is the only  
322 method that does so consistently – the geometric baselines are competitive only on benchmarks whose  
323 inputs already encode the multimodality. The discriminating test is the coupled double-well, where  
324 Kramers escape places the modal disagreement in output space: Core-Set and BAIT both collapse  
325 to within  $1.7 \times$  of Random, while MI-LB wins by  $4 \times$  – the empirical signature that distributional  
326 acquisition is doing something feature-space coverage cannot.

327 **Limitations.** MI-LB currently relies on the availability of explicit density components, such as  
328 those in Mixture Density Networks. Its performance is also tied to the quality of the entropy bounds,  
329 which may loosen in some settings. Furthermore, the framework assumes a well-specified model class  
330 capable of recovering the true irreducible stochasticity. Our empirical evaluation is also restricted to  
331 synthetic benchmarks with controllable multimodality, since no standard benchmarks exist for active  
332 learning in continuous multimodal regression; building such benchmarks is left to future work.

333 **Future Work.** Future research may focus on extending the MI-LB objective to implicit generative  
334 models, such as diffusion and flow matching, where densities are not analytically tractable.  
335 Another promising direction is to explore the framework’s utility in safety-critical control tasks where  
336 distinguishing between lack of data and inherent noise is vital for risk-aware planning.

## 337 **Impact Statement**

338 Our contributions enable advances in deep learning and uncertainty quantification. This has the  
339 potential to impact a wide range of downstream applications. While we do not anticipate specific  
340 negative impacts from this work, as with any powerful predictive tool, there is potential for misuse.  
341 We encourage the research community to consider the ethical implications and potential dual-use  
342 scenarios when applying these technologies in sensitive domains and to avoid its application altogether  
343 to weaponry and other military technologies.

## 344 **Declaration of LLM Usage**

345 LLMs were used during the development of this paper, including for editing the writing, developing  
346 ideas, and as code assistants. We believe our use to be within the standard uses of this technology,  
347 and authors have verified the integrity of all material contained in this work.

## 348 **References**

- 349 [1] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning  
350 for classification and preference learning, 2011.
- 351 [2] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for  
352 computer vision? *Advances in neural information processing systems*, 30, 2017.
- 353 [3] Ian Osband, Zheng Wen, Mohammad Asghari, Morteza Ibrahimi, Xiyuan Lu, and Benjamin Van  
354 Roy. Epistemic neural networks. *CoRR*, abs/2107.08924, 2021.
- 355 [4] Marco F. Huber, Tim Bailey, Hugh Durrant-Whyte, and Uwe D. Hanebeck. On entropy  
356 approximation for gaussian mixture random vectors. In *2008 IEEE International Conference  
357 on Multisensor Fusion and Integration for Intelligent Systems*, pages 181–188, 2008.
- 358 [5] Matias Valdenegro-Toro and Daniel Saromo Mori. A deeper look into aleatoric and epistemic  
359 uncertainty disentanglement. In *2022 IEEE/CVF Conference on Computer Vision and Pattern  
360 Recognition Workshops (CVPRW)*, pages 1508–1516. IEEE, 2022.
- 361 [6] Pavel Izmailov, Sharad Vikram, Matthew D Hoffman, and Andrew Gordon Gordon Wilson.  
362 What are bayesian neural network posteriors really like? In *International conference on machine  
363 learning*, pages 4629–4640. PMLR, 2021.
- 364 [7] Andrew Gordon Wilson and Pavel Izmailov. Deep ensembles as approximate bayesian inference.  
365 <https://cims.nyu.edu/~andrewgw/deeensemles/>, 2021.
- 366 [8] Bálint Mucsányi, Michael Kirchhof, and Seong Joon Oh. Benchmarking uncertainty disen-  
367 tanglement: Specialized uncertainties for specialized tasks. *Advances in neural information  
368 processing systems*, 37:50972–51038, 2024.
- 369 [9] Christopher M Bishop. Mixture density networks. *Neural Computing Research Group Report*,  
370 1994.
- 371 [10] Leonardo Ferreira Guilhoto, Akshat Kaushal, and Paris Perdikaris. Multimodal scientific  
372 learning beyond diffusions and flows, 2026.
- 373 [11] David Holzmüller, Viktor Zaverkin, Johannes Kästner, and Ingo Steinwart. A framework and  
374 benchmark for deep batch active learning for regression. *Journal of Machine Learning Research*,  
375 24(164):1–81, 2023.
- 376 [12] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set  
377 approach. In *International Conference on Learning Representations (ICLR)*, 2018.
- 378 [13] Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham Kakade. Gone fishing: Neural  
379 active learning with fisher embeddings. In *Advances in Neural Information Processing Systems  
380 (NeurIPS)*, 2021.

- 381 [14] Andreas Kirsch, Sebastian Farquhar, and Yarin Gal. A simple baseline for batch active learning  
382 with stochastic acquisition functions. *CoRR*, abs/2106.12059, 2021.
- 383 [15] Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch  
384 acquisition for deep bayesian active learning. *Advances in neural information processing*  
385 *systems*, 32, 2019.
- 386 [16] Jennifer A Hoeting, David Madigan, Adrian E Raftery, and Chris T Volinsky. Bayesian model  
387 averaging: a tutorial (with comments by m. clyde, david draper and ei george, and a rejoinder  
388 by the authors. *Statistical science*, 14(4):382–417, 1999.
- 389 [17] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
390 *arXiv:1312.6114*, 2013.
- 391 [18] Lawrence Cayton et al. Algorithms for manifold learning. *Univ. of California at San Diego*  
392 *Tech. Rep.*, 12(1-17):1, 2005.
- 393 [19] H. A. Kramers. Brownian motion in a field of force and the diffusion model of chemical  
394 reactions. *Physica*, 7(4):284–304, 1940.
- 395 [20] Turab Lookman, Prasanna V. Balachandran, Dezhen Xue, and Ruihao Yuan. Active learning in  
396 materials science with emphasis on adaptive sampling using uncertainties for targeted design.  
397 *npj Computational Materials*, 5(1):21, 2019.
- 398 [21] Dezhen Xue, Prasanna V. Balachandran, John Hogden, James Theiler, Deqing Xue, and Turab  
399 Lookman. Accelerated search for materials with targeted properties by adaptive design. *Nature*  
400 *Communications*, 7(1):11241, 2016.
- 401 [22] Eric Stach, Brian DeCost, A. Gilad Kusne, Jason Hatrick-Simpers, Keith A. Brown, Kristofer G.  
402 Reyes, Joshua Schrier, Simon Billinge, Tonio Buonassisi, Ian Foster, Carla P. Gomes, John M.  
403 Gregoire, Apurva Mehta, Joseph Montoya, Elsa Olivetti, Chiwoo Park, Eli Rotenberg, Semion K.  
404 Saikin, Sylvia Smullin, Valentin Stanev, and Benji Maruyama. Autonomous experimentation  
405 systems for materials development: A community perspective. *Matter*, 4(9):2702–2726, 2021.
- 406 [23] Hans Lukas, Suzana G. Fries, and Bo Sundman. *Computational Thermodynamics: The*  
407 *CALPHAD Method*. Cambridge University Press, USA, 1st edition, 2007.
- 408 [24] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using  
409 deep conditional generative models. *Advances in neural information processing systems*, 28,  
410 2015.
- 411 [25] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow  
412 matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 413 [26] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of*  
414 *control, signals and systems*, 2(4):303–314, 1989.
- 415 [27] Yulong Lu and Jianfeng Lu. A universal approximation theorem of deep neural networks  
416 for expressing probability distributions. *Advances in neural information processing systems*,  
417 33:3094–3105, 2020.
- 418 [28] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal  
419 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao  
420 Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- 421 [29] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas  
422 Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2023.
- 423 [30] DeepMind, Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter  
424 Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Antoine Dedieu, Claudio  
425 Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel, Shaobo  
426 Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch, Lena  
427 Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, George Papamakarios, John  
428 Quan, Roman Ring, Francisco Ruiz, Alvaro Sanchez, Laurent Sartran, Rosalia Schneider, Eren  
429 Sezener, Stephen Spencer, Srivatsan Srinivasan, Miloš Stanojević, Wojciech Stokowiec, Luyu  
430 Wang, Guangyao Zhou, and Fabio Viola. The DeepMind JAX Ecosystem, 2020.

431 [31] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*  
432 *arXiv:1606.08415*, 2016.

433 **A Mathematical Notation**

434 Table 2 summarizes the symbols and notation used in this work.

435 For operands that involve expectations, such as expectation  $\mathbb{E}$ , variance  $\text{Var}$  and entropy  $H$ , a sub-  
 436 index indicates what is the random variable for which the expectation is being taken over. When the  
 437 context is clear, this sub-index is often omitted.

Table 2: Summary of the symbols and notation used in this paper.

<b>Symbol</b>	<b>Meaning</b>
$x \in \mathcal{X}$	Input to the model
$Y \in \mathcal{Y}$	Output random variable
$\theta \in \Theta$	Parameters of a neural network
$Z \sim P_Z$	Epistemic index over model hypotheses
$\epsilon \sim P_\epsilon$	Aleatoric index for irreducible stochasticity
$g_\theta(x; z, \epsilon)$	Learned generative map to the output space
$p_\theta(\cdot   x, z)$	Aleatoric conditional distribution for fixed $z$ from eq. (3)
$p_\theta(\cdot   x)$	Predictive distribution, marginal over $P_Z$ from eq. (4)
$p^*(\cdot   x)$	True conditional distribution
$h_{\#}\mu$	Push-forward of measure $\mu$ under map $h$ from eq. (2)
$\mathbb{E}$	Expectation operator
$\text{Var}(\cdot)$	Variance operator
$H(\cdot)$	Differential entropy
$I(Y; Z   x)$	Epistemic mutual information
$\text{KL}(p \  q)$	KL divergence from $q$ to $p$
$\text{TV}(p, q)$	Total variation distance between $p$ and $q$
$\mathcal{N}(\mu, \Sigma)$	Normal distribution with mean $\mu$ and covariance $\Sigma$
$\mathcal{N}(y; \mu, \Sigma)$	PDF of $\mathcal{N}(\mu, \Sigma)$ evaluated at $y$
$n_{\text{ens}}$	Number of ensemble members
$K$	Number of components in a Gaussian mixture
$\alpha_i^{(z)}(x)$	Mixture weight of component $i$ for member $z$
$\mu_i^{(z)}(x)$	Mean of component $i$ for member $z$
$C_i^{(z)}(x)$	Covariance of component $i$ for member $z$
$w_z$	Weight of ensemble member $z$ under $P_Z$
$MI-LB(x)$	Mutual Information Lower Bound acquisition function
$\mathcal{D}_n$	Training dataset of size $n$
$\mathbf{I}_d$	Identity matrix of size $d$

438 **B Concrete Examples of the Two Index Framework**

439 To ground the two-index framework in practice, we show how three families of models already in  
 440 common use fit naturally into the formalism of equation (1). In each case we identify the epistemic

441 index  $Z$ , the aleatoric index  $\epsilon$ , the learned map  $g_\theta$ , and the resulting conditional and predictive  
 442 distributions.

443 **Ensemble of Conditional Variational Autoencoders (C-VAEs).** Consider an ensemble of  $n_{\text{ens}}$   
 444 independently trained C-VAE [24], indexed by  $k \in \{1, \dots, n_{\text{ens}}\}$ . Each model  $k$  consists of a  
 445 decoder  $d_{\theta_k} : \mathcal{X} \times \mathcal{L} \rightarrow \mathcal{Y}$ , where  $\mathcal{L}$  is the latent space. At inference time, a latent code is sampled  
 446 from the prior  $\epsilon \sim P_\epsilon := \mathcal{N}(0, I)$  and passed through the decoder to produce a prediction.

447 In the two-index framework, we set:

$$\begin{aligned} Z &\sim P_Z := \text{Uniform}\{1, \dots, n_{\text{ens}}\}, \\ \epsilon &\sim P_\epsilon := \mathcal{N}(0, I_d), \\ g_\theta(x; z, \epsilon) &:= d_{\theta_z}(x, \epsilon), \end{aligned}$$

448 where  $\theta = (\theta_1, \dots, \theta_{n_{\text{ens}}})$  collects all decoder parameters. The epistemic index  $Z = k$  selects a  
 449 member of the ensemble, capturing uncertainty about which model best represents the data-generating  
 450 process. The aleatoric index  $\epsilon$  is the latent noise injected into the decoder of the selected member,  
 451 capturing the intrinsic stochasticity of the output given a fixed hypothesis. The aleatoric conditional  
 452 from eq. (3) becomes

$$p_\theta(\cdot | x, k) = (d_{\theta_k}(x, \cdot))_{\#} \mathcal{N}(0, I_d),$$

453 which is the generative distribution of the  $k$ -th C-VAE. The predictive distribution equation 4 is the  
 454 equally-weighted mixture of these per-member distributions:

$$p_\theta(\cdot | x) = \frac{1}{n_{\text{ens}}} \sum_{k=1}^{n_{\text{ens}}} p_\theta(\cdot | x, k).$$

455 Epistemic uncertainty, measured by  $I(Y; Z | x)$ , is large when the ensemble members disagree on  
 456 their generative distributions, and vanishes when all decoders produce the same output distribution  
 457 regardless of which member is selected.

458 **Ensemble of Conditional Flow Matching Models.** Flow matching models [25] learn a deter-  
 459 ministic vector field that transports an initial noise sample  $\epsilon \sim P_\epsilon$  to a target distribution over  $\mathcal{Y}$ ,  
 460 conditioned on input  $x$ . At inference time, the output is obtained by integrating the learned vector  
 461 field from time  $t = 0$  to  $t = 1$ , starting from  $\epsilon$ . An ensemble of  $n_{\text{ens}}$  such models captures epistemic  
 462 uncertainty through disagreement across members.

463 The two-index instantiation mirrors the C-VAE case, with a key structural difference: the aleatoric  
 464 index  $\epsilon$  is not a latent code injected at an intermediate layer but rather the *initial noise* supplied to the  
 465 flow at inference time. Letting  $\Phi_{\theta_k}(\cdot; x) : \mathcal{Y} \rightarrow \mathcal{Y}$  denote the learned flow map of member  $k$  (i.e.,  
 466 the solution operator of the ODE from  $t = 0$  to  $t = 1$ ), we set:

$$\begin{aligned} Z &\sim P_Z := \text{Uniform}\{1, \dots, n_{\text{ens}}\}, \\ \epsilon &\sim P_\epsilon := \mathcal{N}(0, I_d), \\ g_\theta(x; z, \epsilon) &:= \Phi_{\theta_z}(\epsilon; x). \end{aligned}$$

467 The aleatoric conditional is

$$p_\theta(\cdot | x, k) = (\Phi_{\theta_k}(\cdot; x))_{\#} \mathcal{N}(0, I_d),$$

468 which, for an exactly learned flow, equals the target conditional distribution of member  $k$ .

469 **Bayesian Neural Network for Classification.** Bayesian Neural Networks (BNNs) are a special  
 470 case of ENNs [3], with the epistemic index  $Z$  playing the role of a sample from the weight posterior  
 471  $p(\theta | \mathcal{D})$ . We instantiate this within the two-index framework for a classification setting with  $n_{\text{cls}}$   
 472 classes.

473 Let  $f_w : \mathcal{X} \rightarrow \Delta^{n_{\text{cls}}-1}$  be a neural network with weights  $w$ , mapping inputs to the probability  
 474 simplex. In a BNN, weights are treated as random variables with posterior  $p(w | \mathcal{D})$  approximated by  
 475 a distribution  $Q_\theta(w)$  with learnable parameters  $\theta$  (e.g., a mean-field Gaussian). The epistemic index  
 476  $Z$  is a draw from this approximate posterior. The aleatoric index  $\epsilon$  serves as the source of randomness

477 for sampling a class label from the predicted categorical distribution. Concretely:

$$\begin{aligned} Z &\sim P_Z := Q_\theta(w), \\ \epsilon &\sim P_\epsilon := \text{Uniform}(0, 1), \\ g_\theta(x; z, \epsilon) &:= \min\{c \in \{1, \dots, n_{\text{cls}}\} : (F_z(x))_c \geq \epsilon\}, \end{aligned}$$

478 where  $(F_z(x))_c := \sum_{j=1}^c [f_z(x)]_j$  is the  $c$ -th entry of the cumulative sum of the predicted class  
479 probabilities under weights  $z$ . This is the inverse CDF (quantile) transform:  $g_\theta(x; z, \epsilon)$  samples a  
480 class label from the categorical distribution  $\text{Cat}(f_z(x))$  using  $\epsilon$  as the uniform source of randomness.  
481 The aleatoric conditional is then

$$p_\theta(\cdot | x, z) = \text{Cat}(f_z(x)),$$

482 and the predictive distribution is equal to the Bayesian Model Average (BMA):

$$p_\theta(\cdot | x) = \mathbb{E}_{Z \sim Q_\theta} [\text{Cat}(f_Z(x))].$$

483 The epistemic term  $I(Y; Z | x)$  measures disagreement among weight posterior samples about  
484 the predicted class probabilities, recovering the standard notion of epistemic uncertainty in BNN  
485 classification [3]. The aleatoric term  $\mathbb{E}_Z[H(f_Z(x))]$  is the expected entropy of each categorical  
486 prediction, averaged over the weight posterior. It measures the label ambiguity that remains under a  
487 fixed weight configuration.

## 488 C Limitations of Variance for Multimodal Distributions

489 The Epistemic Variance acquisition function (Section 2.2), also known as Uncertainty Sampling, is  
490 the standard choice for active learning with real-valued predictions. There are settings, however,  
491 where it assigns identical scores to distributions with very different uncertainty profiles. We give a  
492 concrete example on the unit sphere that illustrates this failure mode.

493 For a distribution  $p$  over  $\mathbb{R}^m$  with mean  $\bar{y} = \mathbb{E}_p[Y]$ , we define the *multivariate variance* as the trace  
494 of the covariance matrix:

$$\text{Var}(p) := \text{tr}(\text{Cov}_p(Y)) = \mathbb{E}_p[\|Y - \bar{y}\|_2^2]. \quad (20)$$

495 This is the expected squared  $\ell_2$  distance from a sample to the mean, and it is the quantity that  
496 variance-based acquisition functions rank candidates by.

497 Consider the unit sphere  $S^{m-1} := \{y \in \mathbb{R}^m : \|y\|_2 = 1\}$  and two inputs  $x, x' \in \mathbb{R}^n$ . For  $x$ , suppose  
498  $p_\theta(y | x)$  is the *uniform distribution* on  $S^{m-1}$ . For  $x'$ , let  $C_\delta(e_1)$  and  $C_\delta(-e_1)$  denote the spherical  
499 caps of geodesic radius  $\delta > 0$  centered at  $e_1$  and  $-e_1$  respectively, and define the *polar distribution*  
500  $p_\theta(y | x')$  as the distribution that selects one of the two caps with equal probability and then draws  
501 uniformly within it. Both distributions are continuous with respect to the surface measure on  $S^{m-1}$ ,  
502 both have mean  $\bar{y} = 0$  by symmetry, and every sample lies on  $S^{m-1}$ . Since  $\|Y\|_2^2 = 1$  almost surely  
503 and  $\bar{y} = 0$  for both distributions, equation (20) gives  $\text{Var}(p) = 1$  in both cases, exactly and for all  $m$   
504 (see Figure 4 for an illustration). Any acquisition function that ranks inputs by variance assigns the  
505 same score to  $x$  and  $x'$ .

506 The two distributions are not equally uncertain, however. The polar distribution confines the output to  
507 two small patches of the sphere, while the uniform distribution spreads mass over the entire surface.  
508 The entropy-based decomposition of Section 2.2 captures this gap. Writing  $A_m = \text{Area}(S^{m-1})$  for  
509 the total surface area and  $A_\delta = \text{Area}(C_\delta)$  for the area of each cap, the differential entropies with  
510 respect to the surface measure are

$$h(p_\theta(\cdot | x)) = \log A_m, \quad h(p_\theta(\cdot | x')) = \log 2 + \log A_\delta. \quad (21)$$

511 The entropy gap is  $\log(A_m/(2A_\delta))$ , which grows without bound as  $m \rightarrow \infty$  for any fixed cap radius  
512  $\delta$ . This follows from concentration of measure on the sphere: the fraction  $A_\delta/A_m$  of the sphere  
513 covered by a cap of fixed geodesic radius  $\delta < \pi/2$  vanishes exponentially in  $m$ , because the mass  
514 of the integrand  $\sin^{m-2}(\phi)$  in the surface area element concentrates near the equator  $\phi = \pi/2$ .  
515 Entropy-based acquisition functions detect this growing gap, while variance-based scores cannot.

## Two Probability Distributions with Equal Variance

Both distributions have the same variance despite their different shapes

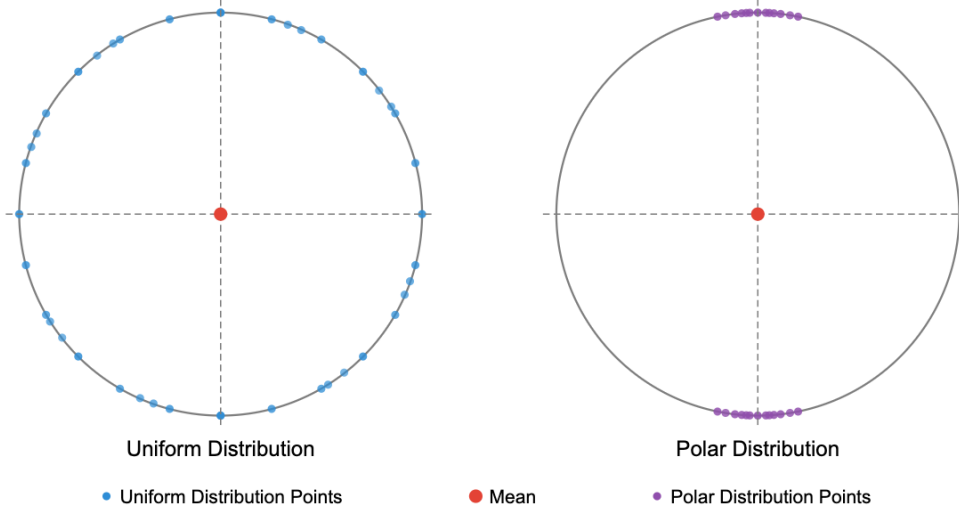


Figure 4: Two distributions on the unit circle with identical variance but different entropy. **Left:** the uniform distribution spreads samples evenly across the circle. **Right:** the polar distribution concentrates samples near two antipodal points. Both distributions have mean zero and multivariate variance equal to 1, yet the uniform distribution has strictly higher entropy.

## 516 D Proof of Asymptotic Guarantees

517 In what follows, we formalize the insights from Section 2.3 and theorem 2.1 in order to provide  
 518 guarantees on the asymptotic behavior under infinite data. First, we prove that epistemic uncertainty  
 519 measured as the mutual information  $I(Y; Z|x)$  converges to zero. Then, we show that the predicted  
 520 aleatoric uncertainty  $\mathbb{E}_Z[H_\epsilon(Y | Z, x, \mathcal{D}_n)]$  converges to the true quantity  $H(p^*(\cdot | x))$ . As a  
 521 corollary of these two results, we also observe that the total predictive entropy also converges to the  
 522 true entropy of the data-generating process.

### 523 D.1 Mutual Information Collapse

524 We begin by stating and proving a lemma that will be helpful in proving the main theorem.

525 **Lemma D.1** (KL divergence–total variation bound under bounded likelihood ratio). *Let  $p$  and  $q$  be*  
 526 *probability measures on a measurable space  $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$  with  $p \ll q$ . Suppose there exists a constant*  
 527  *$C \geq 1$  such that  $dp/dq \leq C$  holds  $q$ -almost everywhere. Then*

$$\text{KL}(p \| q) \leq 2M_C \cdot \text{TV}(p, q),$$

528 where

$$M_C := \max\left(1, \frac{C \log C - C + 1}{C - 1}\right) \quad (22)$$

529 for  $C > 1$ , and  $M_1 := 1$ .

530 *Proof.* Let  $h := dp/dq$ , so that  $0 \leq h \leq C$  holds  $q$ -a.e. Since  $\int h dq = 1$ , the KL divergence can be  
 531 written as

$$\text{KL}(p \| q) = \int h \log h dq = \int \phi(h) dq, \quad (23)$$

532 where  $\phi(t) := t \log t - t + 1$ . The function  $\phi$  is non-negative and convex on  $[0, \infty)$  with  $\phi(1) = 0$   
 533 and  $\phi'(1) = 0$ .

534 We claim that  $\phi(t) \leq M_C |t - 1|$  for all  $t \in [0, C]$ . On the interval  $[0, 1]$ : since  $t \log t \leq 0$  for  
 535  $t \in [0, 1]$ , we have  $\phi(t) = t \log t + (1 - t) \leq 1 - t = |t - 1|$ . On the interval  $[1, C]$ : convexity of  $\phi$

536 together with  $\phi(1) = 0$  gives  $\phi(t) \leq \phi(C) \cdot \frac{t-1}{C-1}$ , where  $\phi(C) = C \log C - C + 1$ . Combining both  
 537 cases yields  $\phi(t) \leq M_C |t - 1|$  on  $[0, C]$ .

538 Integrating against  $q$  and applying the bound gives

$$\begin{aligned} \text{KL}(p \parallel q) &= \int \phi(h) dq \leq M_C \int |h - 1| dq \\ &= 2M_C \cdot \text{TV}(p, q), \end{aligned} \quad (24)$$

539 where the last equality uses  $\text{TV}(p, q) = \frac{1}{2} \int |dp - dq| = \frac{1}{2} \int |h - 1| dq$ .  $\square$

540 Equipped with this lemma, we now state the first part of the theorem precisely.

541 **Theorem D.2** (Mutual Information Collapse). *Let  $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$  be a measurable space, let  $(\mathcal{Z}, P_Z)$   
 542 and  $(\mathcal{E}, P_\epsilon)$  be probability spaces, let  $\Theta \subseteq \mathbb{R}^d$ , and let  $g_\theta : \mathcal{X} \times \mathcal{Z} \times \mathcal{E} \rightarrow \mathcal{Y}$  be measurable for  
 543 each  $\theta \in \Theta$ . Fix  $x \in \mathcal{X}$ , define the aleatoric conditional  $p_\theta(\cdot \mid x, z) := (g_\theta(x; z, \cdot))_{\#} P_\epsilon$  and the  
 544 predictive marginal  $p_\theta(\cdot \mid x) := \int p_\theta(\cdot \mid x, z) dP_Z(z)$ . Suppose the following conditions hold.*

545 **A1 Well-specification.** *There exists  $\theta^* \in \Theta$  such that  $p_{\theta^*}(\cdot \mid x, z) = p^*(\cdot \mid x)$  for  $P_Z$ -almost  
 546 every  $z$ .*

547 **A2 Consistency.**  $\hat{\theta}_n \rightarrow \theta^*$  in probability as  $n \rightarrow \infty$ .

548 **A3 Uniform continuity in total variation.** *The map  $\theta \mapsto p_\theta(\cdot \mid x, z)$  is continuous at  $\theta^*$  in total  
 549 variation, uniformly in  $z$ : for every  $\epsilon > 0$  there exists  $\delta > 0$  such that  $\|\theta - \theta^*\| < \delta$  implies  
 550  $\sup_{z \in \mathcal{Z}} \text{TV}(p_\theta(\cdot \mid x, z), p_{\theta^*}(\cdot \mid x, z)) < \epsilon$ .*

551 **A4 Bounded likelihood ratio.** *There exists  $C \geq 1$  such that for all  $\theta$  in a neighborhood of  $\theta^*$   
 552 and  $P_Z$ -almost every  $z$ , the aleatoric conditional is absolutely continuous with respect to  
 553 the predictive marginal and*

$$\frac{dp_\theta(\cdot \mid x, z)}{dp_\theta(\cdot \mid x)} \leq C \quad p_\theta(\cdot \mid x)\text{-a.e.}$$

554 *Then*

$$I_{\hat{\theta}_n}(Y; Z \mid x) := \mathbb{E}_Z \left[ \text{KL} \left( p_{\hat{\theta}_n}(\cdot \mid x, Z) \parallel p_{\hat{\theta}_n}(\cdot \mid x) \right) \right] \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty.$$

555 Before proving this theorem, we describe what assumptions A1–A4 intuitively mean. Assumption A1  
 556 relates to *universality* of the chosen model class: the neural network architecture should be expressive  
 557 enough to contain a parameter  $\theta^*$  that captures the desired map, and is a standard assumption in  
 558 the machine learning literature, having been proved for many choices of architectures via universal  
 559 approximation theorems [26, 27]. Assumption A2 relates to well-posedness of training: under  
 560 infinite data, it should be possible to recover the true optimum of the problem. Assumption A3  
 561 is a regularity condition requiring that the aleatoric conditionals respond smoothly to parameter  
 562 perturbations, uniformly across the epistemic index. Assumption A4 requires a uniform upper bound  
 563 on the likelihood ratio  $dp_\theta(\cdot \mid x, z)/dp_\theta(\cdot \mid x)$ ; for finite ensembles with prior weights  $w_z > 0$ , this  
 564 is automatically satisfied with  $C = 1/\min_z w_z$ , since the marginal  $p_\theta(\cdot \mid x) = \sum_z w_z p_\theta(\cdot \mid x, z) \geq$   
 565  $w_z p_\theta(\cdot \mid x, z)$  dominates each conditional. For continuous epistemic indices, A4 is a mild regularity  
 566 condition excluding pathological concentration of the conditional relative to the marginal.

567 *Proof.* Fix  $x \in \mathcal{X}$  and write  $I_\theta := I_\theta(Y; Z \mid x)$ ,  $p_\theta(z) := p_\theta(\cdot \mid x, z)$ , and  $p_\theta := p_\theta(\cdot \mid x)$  for  
 568 brevity. The strategy is to show that  $\theta \mapsto I_\theta$  is continuous at  $\theta^*$  with  $I_{\theta^*} = 0$ , after which the  
 569 conclusion follows from A2 by the continuous mapping theorem.

570 By A1,  $p_{\theta^*}(z) = p^*(\cdot \mid x)$  for  $P_Z$ -a.e.  $z$ , and marginalizing gives

$$\begin{aligned} p_{\theta^*} &= \int p_{\theta^*}(z) dP_Z(z) = \int p^*(\cdot \mid x) dP_Z(z) \\ &= p^*(\cdot \mid x), \end{aligned} \quad (25)$$

571 so that  $p_{\theta^*}(z) = p_{\theta^*}$  for  $P_Z$ -a.e.  $z$ . Consequently  $\text{KL}(p_{\theta^*}(z) \parallel p_{\theta^*}) = 0$  for  $P_Z$ -a.e.  $z$ , and integrat-  
 572 ing against  $P_Z$  yields  $I_{\theta^*} = 0$ .

573 For continuity at  $\theta^*$ , fix  $\varepsilon > 0$ . Since  $p_\theta(\cdot | x)$  is a mixture over  $P_Z$  with  $p_\theta(\cdot | x, z) \ll p_\theta(\cdot | x)$ ,  
 574 Lemma D.1 together with A4 gives, for every  $\theta$  in the neighborhood specified by A4 and  $P_Z$ -a.e.  $z$ ,

$$\text{KL}(p_\theta(z) \| p_\theta) \leq 2M_C \cdot \text{TV}(p_\theta(z), p_\theta), \quad (26)$$

575 where  $M_C$  is the constant from equation 22. It therefore suffices to bound  $\mathbb{E}_Z[\text{TV}(p_\theta(Z), p_\theta)]$ . The  
 576 triangle inequality for total variation, applied with  $p_{\theta^*}(z)$  and  $p_{\theta^*}$  as intermediate points, gives

$$\begin{aligned} \text{TV}(p_\theta(z), p_\theta) &\leq \text{TV}(p_\theta(z), p_{\theta^*}(z)) \\ &\quad + \text{TV}(p_{\theta^*}(z), p_{\theta^*}) \\ &\quad + \text{TV}(p_{\theta^*}, p_\theta), \end{aligned} \quad (27)$$

577 in which the middle term vanishes for  $P_Z$ -a.e.  $z$  by equation 25. By A3 there exists  $\delta > 0$  such that  
 578  $\|\theta - \theta^*\| < \delta$  implies

$$\sup_{z \in \mathcal{Z}} \text{TV}(p_\theta(z), p_{\theta^*}(z)) < \frac{\varepsilon}{4M_C}, \quad (28)$$

579 controlling the first term in equation 27 uniformly in  $z$ . For the third term, joint convexity of total  
 580 variation together with the representation  $p_\theta = \int p_\theta(z) dP_Z(z)$  yields

$$\begin{aligned} \text{TV}(p_{\theta^*}, p_\theta) &\leq \int \text{TV}(p_{\theta^*}(z), p_\theta(z)) dP_Z(z) \\ &< \frac{\varepsilon}{4M_C}, \end{aligned} \quad (29)$$

581 where the strict bound again follows from equation 28 at the same  $\delta$ . Integrating equation 27 against  
 582  $P_Z$  and substituting equation 28 and equation 29 gives

$$\begin{aligned} \mathbb{E}_Z[\text{TV}(p_\theta(Z), p_\theta)] &< \frac{\varepsilon}{4M_C} + 0 + \frac{\varepsilon}{4M_C} \\ &= \frac{\varepsilon}{2M_C}, \end{aligned} \quad (30)$$

583 which combined with equation 26 yields

$$I_\theta \leq 2M_C \cdot \mathbb{E}_Z[\text{TV}(p_\theta(Z), p_\theta)] < 2M_C \cdot \frac{\varepsilon}{2M_C} = \varepsilon. \quad (31)$$

584 Hence  $\theta \mapsto I_\theta$  is continuous at  $\theta^*$ .

585 Since  $\hat{\theta}_n \xrightarrow{P} \theta^*$  by A2 and  $I_\theta$  is continuous at  $\theta^*$  with  $I_{\theta^*} = 0$ , the continuous mapping theorem  
 586 gives  $I_{\hat{\theta}_n} \xrightarrow{P} I_{\theta^*} = 0$ .  $\square$

## 587 D.2 Aleatoric Entropy Convergence

588 Theorem D.2 establishes that the epistemic mutual information collapses to zero. We now show that  
 589 the aleatoric uncertainty estimator converges to the entropy of the true conditional. The argument  
 590 follows the same three-step template as Theorem D.2: point evaluation at  $\theta^*$ , continuity at  $\theta^*$ , and  
 591 the continuous mapping theorem. The new ingredient is a regularity condition that does not follow  
 592 from A3 alone.

593 Assumption A3 controls  $\theta \mapsto p_\theta(\cdot | x, z)$  in total variation, but differential entropy is not continuous  
 594 in TV: distributions can be arbitrarily TV-close with entropies that differ by any prescribed amount.  
 595 We therefore add the following assumption.

596 **A5 Uniform entropy continuity.** The map  $\theta \mapsto H(p_\theta(\cdot | x, z))$  is continuous at  $\theta^*$ , uni-  
 597 formly in  $z$ : for every  $\varepsilon > 0$  there exists  $\delta > 0$  such that  $\|\theta - \theta^*\| < \delta$  implies  
 598  $\sup_{z \in \mathcal{Z}} |H(p_\theta(\cdot | x, z)) - H(p_{\theta^*}(\cdot | x, z))| < \varepsilon$ .

599 For Gaussian-mixture conditionals of the type used in Section 4.1, A5 follows from continuous  
 600 parametrisation of the mixture weights, means, and covariances together with a uniform lower bound  
 601 on component covariance eigenvalues in a neighborhood of  $\theta^*$ . For finite ensembles, the supremum  
 602 over  $z$  reduces to a maximum over finitely many continuous functions and is automatic.

603 **Theorem D.3** (Aleatoric Entropy Consistency). *Adopt the setup of Theorem D.2 and suppose A1, A2,*  
604 *and A5 hold. Then*

$$A_{\hat{\theta}_n}(x) := \mathbb{E}_Z \left[ H_\epsilon \left( Y \mid Z, x; \hat{\theta}_n \right) \right] \xrightarrow{P} H(p^*(\cdot \mid x)) \quad \text{as } n \rightarrow \infty.$$

605 *Proof.* Fix  $x \in \mathcal{X}$  and write  $A_\theta := \mathbb{E}_Z [H(p_\theta(\cdot \mid x, Z))]$  and  $p_\theta(z) := p_\theta(\cdot \mid x, z)$  for brevity. The  
606 strategy mirrors Theorem D.2: show that  $\theta \mapsto A_\theta$  is continuous at  $\theta^*$  with  $A_{\theta^*} = H(p^*(\cdot \mid x))$ , then  
607 apply the continuous mapping theorem under A2.

608 By A1,  $p_{\theta^*}(z) = p^*(\cdot \mid x)$  for  $P_Z$ -a.e.  $z$ , so  $H(p_{\theta^*}(z)) = H(p^*(\cdot \mid x))$  for  $P_Z$ -a.e.  $z$ . Integrating  
609 against  $P_Z$  gives

$$A_{\theta^*} = H(p^*(\cdot \mid x)). \quad (32)$$

610 For continuity at  $\theta^*$ , fix  $\epsilon > 0$ . By A5 there exists  $\delta > 0$  such that  $\|\theta - \theta^*\| < \delta$  implies

$$\sup_{z \in \mathcal{Z}} |H(p_\theta(z)) - H(p_{\theta^*}(z))| < \epsilon. \quad (33)$$

611 Applying the triangle inequality for integrals followed by equation 33 yields

$$\begin{aligned} |A_\theta - A_{\theta^*}| &\leq \int |H(p_\theta(z)) - H(p_{\theta^*}(z))| dP_Z(z) \\ &\leq \sup_{z \in \mathcal{Z}} |H(p_\theta(z)) - H(p_{\theta^*}(z))| \\ &< \epsilon, \end{aligned} \quad (34)$$

612 so  $\theta \mapsto A_\theta$  is continuous at  $\theta^*$ .

613 Since  $\hat{\theta}_n \xrightarrow{P} \theta^*$  by A2 and  $\theta \mapsto A_\theta$  is continuous at  $\theta^*$ , the continuous mapping theorem combined  
614 with equation 32 gives  $A_{\hat{\theta}_n} \xrightarrow{P} H(p^*(\cdot \mid x))$ .  $\square$

615 Theorems D.2 and D.3 together imply that the total predictive entropy converges to the entropy of the  
616 true conditional, which gives a complete consistency picture for the entropy decomposition equation 7.

617 **Corollary D.4** (Total Predictive Entropy Consistency). *Under A1–A5,*

$$H_{Z,\epsilon}(Y \mid x; \hat{\theta}_n) \xrightarrow{P} H(p^*(\cdot \mid x)) \quad \text{as } n \rightarrow \infty.$$

618 *Proof.* By the entropy decomposition equation 7,

$$H_{Z,\epsilon}(Y \mid x; \hat{\theta}_n) = A_{\hat{\theta}_n}(x) + I_{\hat{\theta}_n}(Y; Z \mid x).$$

619 Theorem D.3 gives  $A_{\hat{\theta}_n}(x) \xrightarrow{P} H(p^*(\cdot \mid x))$  and Theorem D.2 gives  $I_{\hat{\theta}_n}(Y; Z \mid x) \xrightarrow{P} 0$ . The  
620 continuous mapping theorem applied to the sum yields the result.  $\square$

## 621 E Proof of Theorem 3.1

622 *Proof.* We wish to show that  $MI-LB(x) \leq I(Y; Z \mid x)$  for all  $x \in \mathcal{X}$ . Recalling the decomposition  
623 from (9),

$$I(Y; Z \mid x) = H(Y \mid x) - \mathbb{E}_Z [H(Y \mid x, Z)]. \quad (35)$$

624 We bound the two terms on the right-hand side in opposite directions using the results of [4].

625 **Step 1: Lower bounding  $H(Y \mid x)$ .**

626 The marginal predictive distribution  $p_\theta(\cdot \mid x)$  is the Gaussian mixture with  $n_{\text{ens}} \cdot K$  components  
627 given in equation 11. By the lower bound of [4],

$$H(Y \mid x) \geq H_{\text{lower}}(Y \mid x), \quad (36)$$

628 where  $H_{\text{lower}}(Y \mid x)$  is equation 12 applied to the marginal mixture with weights  $\beta_{z,i}$ , means  $\mu_i^{(z)}$ ,  
629 and covariances  $C_i^{(z)}$ .

630 **Step 2: Upper bounding**  $\mathbb{E}_Z[H(Y | x, Z)]$ .

631 For each fixed  $z \in \{1, \dots, n_{\text{ens}}\}$ , the aleatoric conditional  $p_\theta(\cdot | x, z)$  is a Gaussian mixture with  $K$   
 632 components as in equation 10. By the upper bound of [4],

$$H(Y | x, Z = z) \leq H_{\text{upper}}(Y | x, Z = z) \quad \text{for each } z. \quad (37)$$

633 Multiplying by  $w_z \geq 0$  and summing over  $z$ ,

$$\mathbb{E}_Z[H(Y | x, Z)] = \sum_{z=1}^{n_{\text{ens}}} w_z H(Y | x, Z = z) \leq \sum_{z=1}^{n_{\text{ens}}} w_z H_{\text{upper}}(Y | x, Z = z). \quad (38)$$

634 **Step 3: Combining the bounds.**

635 Substituting equation 36 and equation 38 into equation 35, and using the fact that  $f(a, b) = a - b$  is  
 636 non-decreasing in  $a$  and non-increasing in  $b$ ,

$$\begin{aligned} I(Y; Z | x) &= H(Y | x) - \mathbb{E}_Z[H(Y | x, Z)] \\ &\geq H_{\text{lower}}(Y | x) - \sum_{z=1}^{n_{\text{ens}}} w_z H_{\text{upper}}(Y | x, Z = z) = MI-LB(x). \quad \square \end{aligned}$$

## 637 F Experimental Details

### 638 F.1 Software

639 Our code is implemented in JAX [28] using the Flax [29] and Optax [30] libraries to define and train  
 640 our neural networks.

### 641 F.2 Synthetic Multimodal Conditional Problem

642 This appendix specifies the full data-generating process, model, trainer, and active-learning protocol  
 643 used in Section 4.1. All hyperparameters listed here are the values used to produce Figures 1a and 6  
 644 and are exposed as defaults in `examples/multimodal_conditional/experiment_config.py`.

645 **Input manifold.** Latent codes  $l \in \mathbb{R}^L$  are drawn  $l \sim \mathcal{N}(0, I_L)$  and embedded into the input space  
 646 via

$$x = \tanh(A l + b_m), \quad A \in \mathbb{R}^{D \times L}, \quad A_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/L), \quad b_m \sim \mathcal{N}(0, I_D).$$

647 Both  $A$  and  $b_m$  are drawn once using `manifold_seed = 1` and held fixed across the benchmark; the  
 648 `tanh` bounds  $x$  to  $[-1, 1]^D$ .

649 **Random Fourier feature map.** All mixture parameters are input-dependent through a fixed random  
 650 Fourier feature map

$$h(x) = \cos(\Omega x + \varphi) \in \mathbb{R}^P, \quad \Omega_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/D), \quad \varphi_i \stackrel{\text{iid}}{\sim} \text{Uniform}(0, 2\pi).$$

651 The parameters  $(\Omega, \varphi)$  are drawn once with `dist_seed = 42` and held fixed.

652 **Per-component means and variances.** For each of the  $K$  mixture components, uncentered means  
 653 and log-variances are

$$\tilde{\mu}_k(x) = B_k h(x) + c_k, \quad \log \Sigma_k(x) = C_k h(x) + b_k^{\text{var}},$$

654 with  $B_k, C_k \in \mathbb{R}^{M \times P}$  having i.i.d.  $\mathcal{N}(0, 1/P)$  entries,  $c_k \sim \mathcal{N}(0, c_{\text{scale}}^2 I_M)$ , and  $b_k^{\text{var}} = 0$ . The  
 655 per-component offsets  $c_k$  control how far apart the modes sit in output space. Component means are  
 656 then centered to enforce  $\mathbb{E}[y | x] = 0$ :

$$\mu_k(x) = \tilde{\mu}_k(x) - \sum_{j=1}^K \pi_j(x) \tilde{\mu}_j(x),$$

657 and the covariance is diagonal,  $\Sigma_k(x) = \text{diag}(\exp \log \Sigma_k(x))$ . The variance-coupling coefficient  $\alpha$   
 658 in the code is set to 0, so log-variances do not depend on  $\pi_k$ .

659 **Structured mixing weights.** To create a clean unimodal/multimodal phase boundary, we use the  
 660 structured mixing mode. Letting  $r = \|x_{1:L}\|$ , define a radial gate and  $K-1$  angular scores,

$$g(x) = \frac{1}{2}(1 + \tanh(\beta(r - r_0))),$$

$$s(x) = \text{softmax}(\gamma V x_{1:L}) \in \Delta^{K-2}, \quad V \in \mathbb{R}^{(K-1) \times L}, \quad V_{k,:}/\|V_{k,:}\| = V_{k,:},$$

661 where the rows of  $V$  are unit vectors drawn from  $\mathcal{N}(0, I_L)$  and normalised, and then assemble the  
 662 logits

$$\ell_0(x) = \text{scale} \cdot (1 - g(x)), \quad \ell_k(x) = \text{scale} \cdot g(x) \cdot s_k(x) \quad (k = 1, \dots, K-1),$$

663 and set  $\pi(x) = \text{softmax}(\ell(x))$ . Inside the radius ( $r < r_0$ ) the mass concentrates on component 0,  
 664 making  $p^*(y | x)$  effectively unimodal; outside the radius the angular scores activate components  
 665  $1, \dots, K-1$  in sectors.

666 **Hyperparameters.** Table 3 lists the distribution hyperparameters and Table 4 the model and  
 667 active-learning hyperparameters. All values match the defaults in `experiment_config.py`.

Table 3: Distribution hyperparameters for the multimodal conditional problem.

Symbol	Value	Meaning
$D$	10	Input dimension
$M$	16	Output dimension
$L$	4	Latent manifold dimension
$K$	3	Number of true mixture components
$P$	128	Random Fourier features
$c_{\text{scale}}$	10.0	Std. of per-component offsets $c_k$
$\alpha$	0.0	Variance-coupling coefficient (disabled)
$\tau$	1.0	Softmax temperature (unused in structured mode)
$\beta$	8.0	Transition sharpness
$r_0$	1.3	Transition radius
$\gamma$	2.0	Angular sharpness
scale	3.0	Logit magnitude
dist_seed	42	PRNG seed for distribution parameters
manifold_seed	1	PRNG seed for manifold map $(A, b_m)$

668 **Oracle NLL.** Because  $p^*(y | x)$  is known in closed form, the oracle negative log-  
 669 likelihood  $\text{NLL}^* = -\mathbb{E}_{(x,y) \sim p^*}[\log p^*(y | x)]$  can be evaluated exactly on the test set.  
 670 Under the settings above it equals  $\text{NLL}^* = 22.98$  (computed by `compute_true_nll` in  
 671 `examples/multimodal_conditional/utisls.py`).

672 **Model architecture.** Each ensemble member is a Mixture Density Network with shared backbone:  
 673 a 2-layer MLP with 128 hidden units per layer and GELU activations [31]. The MLP output is  
 674 linearly projected to  $K_{\text{MDN}} = 5$  mixture components, each with a mean in  $\mathbb{R}^M$  and a diagonal  
 675 covariance; mixture weights come from a softmax head. Ensemble size is  $n_{\text{ens}} = 8$ , with members  
 676 initialised from different PRNG seeds and trained independently on the same labeled set (the standard  
 677 deep-ensemble construction used to instantiate  $P_Z = \text{Uniform}\{1, \dots, n_{\text{ens}}\}$ ).

678 **Training schedule.** All members are trained with AdamW (peak learning rate  $5 \times 10^{-4}$ , weight  
 679 decay  $10^{-2}$ ) preceded by adaptive gradient clipping at 0.1. The learning rate follows a linear warmup  
 680 over  $\min(500, n_{\text{iter}}/5)$  steps to the peak value, after which it decays exponentially at rate 0.9 per  
 681 2,000 steps. The number of gradient steps per active learning round is set adaptively to

$$n_{\text{iter}}(n_{\text{lab}}) = \min(10,000, 10 \cdot n_{\text{lab}}),$$

682 where  $n_{\text{lab}}$  is the current labeled-set size; this prevents overfitting on small labeled sets in early rounds  
 683 while allowing full optimisation once  $n_{\text{lab}} \geq 1000$ . The training mini-batch size is 128 (full-batch  
 684 when  $n_{\text{lab}} < 128$ ). Optimiser, architecture, and active-learning hyperparameters are fixed across all  
 685 acquisition functions; we do not perform per-method tuning.

686 **Active-learning protocol.** For each of 5 seeds we instantiate a fresh pool of 50,000 candidate  
687 inputs, a held-out test set of 2,000 inputs, and an initial labeled set of 100 inputs. Acquisition scores  
688 are evaluated on the remaining pool in chunks of 256. Each active-learning run performs 20 rounds,  
689 acquiring 50 queries per round, so the final labeled set contains  $100 + 20 \cdot 50 = 1100$  examples, i.e.  
690 2.2% of the pool. Test NLL is evaluated on the held-out set at the end of every round.

Table 4: Model and active-learning hyperparameters for the multimodal conditional problem.

Symbol / option	Value	Meaning
$n_{\text{ens}}$	8	Ensemble size
$K_{\text{MDN}}$	5	MDN mixture components per member
hidden features	128	MLP width
depth	2	Number of MLP hidden layers
activation	GELU	Nonlinearity
training batch size	128	Mini-batch size for AdamW
peak learning rate	$5 \times 10^{-4}$	AdamW peak LR (after warmup)
weight decay	$10^{-2}$	AdamW weight decay
gradient clip	0.1	Adaptive gradient clipping threshold
$n_{\text{iter}}$ (cap)	10,000	Max gradient steps per round
iter-per-sample	10	Slope of adaptive schedule
candidate pool size	50,000	Unlabeled pool $ \mathcal{X}_{\text{pool}} $
test set size	2,000	Held-out evaluation set
initial labelled	100	Initial labeled budget
AL rounds	20	Number of acquisition rounds
query batch size	50	Queries acquired per round
acquisition batch size	256	Chunk size for scoring the pool
seeds	$\{0, 1, 2, 3, 4\}$	Data / training seeds for reported curves

691 **Acquisition functions.** Five scoring functions are evaluated. **Random** draws scores i.i.d.  
692 Uniform(0, 1). **Epistemic Variance** uses the trace of the covariance of conditional means across  
693 ensemble members. **MI-LB** is equation 15. **BAIT** [13] uses the paper’s last-layer mean-head Fisher  
694 of ensemble member 0 with one MC sample  $y \sim p_{\theta}(\cdot | x)$  per pool point. **Core-Set** [12] runs  
695 k-Center-Greedy on the shared MDN backbone activations of ensemble member 0 (the “final FC  
696 layer” recipe of §4.4). The latter two produce no scalar score and bypass the score-to-batch step.

697 **Implementation of BAIT and Core-Set in the MDN setting.** Both baselines are adapted from  
698 their original classification specifications. **BAIT** requires per-input Fisher embeddings  $G(x)$  such that  
699  $F(x) = G(x)^{\top} G(x)$  approximates the per-sample Fisher information; following the last-layer recipe  
700 of [13] we restrict the Fisher to the mean-head weights of ensemble member 0. Drawing one Monte-  
701 Carlo sample  $y \sim p_{\theta}(\cdot | x)$  and using the closed-form MDN log-likelihood gradient  $\nabla_{W_{\mu}} \log p(y |$   
702  $x; \theta) = \gamma_k(y, x) (y - \mu_k(x)) / \sigma_k^2(x) \otimes z(x)$ , where  $\gamma_k$  are posterior mixture responsibilities and  $z(x)$   
703 is the shared backbone activation, the embedding is the flattened outer product,  $G(x) \in \mathbb{R}^{1 \times hKd}$ .  
704 Selection minimises  $\text{tr}((F_{\text{train}} + \lambda I)^{-1} F_{\text{cand}})$  via the forward+backward greedy of [13] with ridge  
705  $\lambda = 10^{-3}$  and the labelled-set Fisher as the burden matrix. This single-ensemble, single-MC choice  
706 is a compute trade-off rather than a capacity claim: extending the Fisher embedding to all  $n_{\text{ens}} = 8$   
707 members (to capture cross-member disagreement and the mixture-weight signal the mean-head Fisher  
708 omits) or to multiple MC samples per pool point would multiply an already dominant selection  
709 overhead — BAIT alone consumes  $\sim 12$  of the  $\sim 25$  GPU-h total compute budget (Table 13), pushing  
710 a richer embedding into a 100+ GPU-h regime no other acquisition here requires. We therefore  
711 report the canonical last-layer recipe of [13] and flag the resulting Fisher-estimate variance as a  
712 known limitation in the empirical discussion below. **Core-Set** uses the activations of the shared MDN  
713 backbone of ensemble member 0 (the layer immediately before the mixture head) as  $h$ -dimensional  
714 features, then runs k-Center-Greedy [12] on Euclidean distances initialised at the labelled set. Both  
715 methods consume the full pool without subsampling.

716 **Selection strategies.** Three batch-selection rules convert per-point scores into a query batch of size  
717 50:

- 718 • **Top- $k$  (greedy)**. Select the 50 points with the highest acquisition scores. Used for all  
719 main-text results.
- 720 • **SBAL [14]**. We use the *softmax* variant of SBAL (as opposed to the softmax and power  
721 variants in Kirsch et al.), which handles negative scores without modification: sample  
722 50 points without replacement from  $\text{softmax}(\text{score}/T)$  via the Gumbel-top- $k$  trick. The  
723 temperature  $T$  interpolates between exploitation ( $T \rightarrow 0$ , recovers top- $k$ ) and uniform  
724 exploration ( $T \rightarrow \infty$ , recovers Random). All SBAL runs reported here use  $T = 1.0$ .
- 725 • **MaxDist [11]**. Acquisition-weighted farthest-point sampling in feature space (LCMD-TP  
726 variant). Given input features  $x$  standardised per dimension and normalised scores  $\tilde{s} \in [0, 1]$ ,  
727 greedily pick the point maximising  $d_{\min}(i) \cdot (1 + w \tilde{s}(i))$ , where  $d_{\min}(i)$  is the squared  
728 distance from point  $i$  to its nearest already-selected or already-training point. All MaxDist  
729 runs use score-weight  $w = 1$ . At  $w = 0$  this reduces to pure farthest-point sampling, and as  
730  $w \rightarrow \infty$  it approaches top- $k$  on  $\tilde{s}$ .

731 SBAL and MaxDist only apply to the deterministic acquisitions (Epistemic Variance and MI-LB);  
732 combining either with Random is undefined (Random has no meaningful score ranking or signal  
733 to weight against). Any acquisition of the form `sbal_X` or `maxdist_X` in the experiment log  
734 corresponds to applying the named selection strategy on top of base score  $X$ . BAIT and Core-Set  
735 bypass the score-to-batch step entirely (Fisher-trace forward+backward greedy and k-Center-Greedy  
736 respectively) and have no SBAL/MaxDist variants.

737 **Results for SBAL and MaxDist variants.** Figure 5 and Table 5 report final test NLL at  $n = 1100$   
738 across 5 seeds. At  $T = 1.0$ , MI-LB (SBAL) sits roughly halfway between MI-LB (top- $k$ ) and  
739 Random (33.99 vs 31.12 and 39.73), because softmax sampling at  $T = 1.0$  places non-trivial mass  
740 outside the high-score region near the radial gate; lowering  $T$  tightens the distribution back onto  
741 MI-LB (top- $k$ ) and recovers it exactly as  $T \rightarrow 0$ . MI-LB (MaxDist) at  $w = 1$  matches MI-LB (top- $k$ )  
742 (30.49 vs 31.12), adding batch diversity without sacrificing score quality; larger  $w$  moves it toward  
743 top- $k$  on the score, smaller  $w$  toward pure farthest-point sampling.

744 **BAIT and Core-Set.** Core-Set ties MI-LB within seed-to-seed noise ( $30.57 \pm 0.56$  vs  $31.12 \pm 0.34$ );  
745 MI-LB retains the tighter spread. The tie reflects benchmark geometry: pool inputs lie on a 4-D tanh  
746 manifold, so k-Center-Greedy spreads queries across that manifold and ends up sampling the gate  
747 region that MI-LB targets through entropy disagreement. BAIT effectively ties Variance ( $32.26 \pm 1.55$   
748 vs  $32.56 \pm 1.14$ ), but with the largest seed-to-seed spread of any acquisition we evaluate — the  
749 variance of the single-MC Fisher estimate at small budgets. Both baselines bypass the Two-Index  
750 decomposition; the next two appendices test whether the apparent tie generalises beyond a benchmark  
751 in which input geometry already encodes the multimodality.

Table 5: Final test NLL at  $n = 1100$ , mean  $\pm$  std across 5 seeds (min–max in brackets). SBAL uses temperature  $T = 1.0$ ; MaxDist uses score weight  $w = 1$ . Oracle NLL $^* = 22.98$ .

Acquisition	Mean $\pm$ std	Min – Max
Random	$39.73 \pm 0.90$	38.89 – 40.97
Epistemic Variance (top- $k$ )	$32.56 \pm 1.14$	31.41 – 34.14
Epistemic Variance (SBAL)	$31.01 \pm 0.57$	30.43 – 31.61
MI-LB (top- $k$ )	$31.12 \pm 0.34$	30.58 – 31.37
MI-LB (SBAL)	$33.99 \pm 0.74$	33.21 – 35.01
MI-LB (MaxDist)	$30.49 \pm 0.46$	29.84 – 30.91
BAIT	$32.26 \pm 1.55$	30.40 – 34.34
Core-Set	<b><math>30.57 \pm 0.56</math></b>	29.97 – 31.31

### 752 F.3 Coupled Double-Well System

753 This appendix specifies the full data-generating process, model, trainer, and active-learning protocol  
754 used in Section 4.2.

755 **Simulator.**  $P = 5$  particles evolve according to the overdamped Langevin SDE equation 19 under  
756 open boundary conditions: particle 1 couples only to particle 2, particle  $P$  only to particle  $P-1$ , and

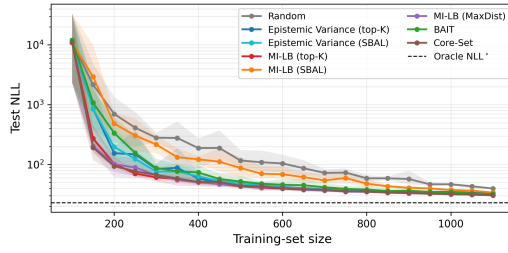


Figure 5: Learning curves on the multimodal benchmark for SBAL ( $T = 1.0$ ) and MaxDist ( $w = 1$ ) variants of Variance and MI-LB, plus the BAIT and Core-Set baselines, alongside the three top- $k$  curves from Fig. 1a (5 seeds; bands min-max).

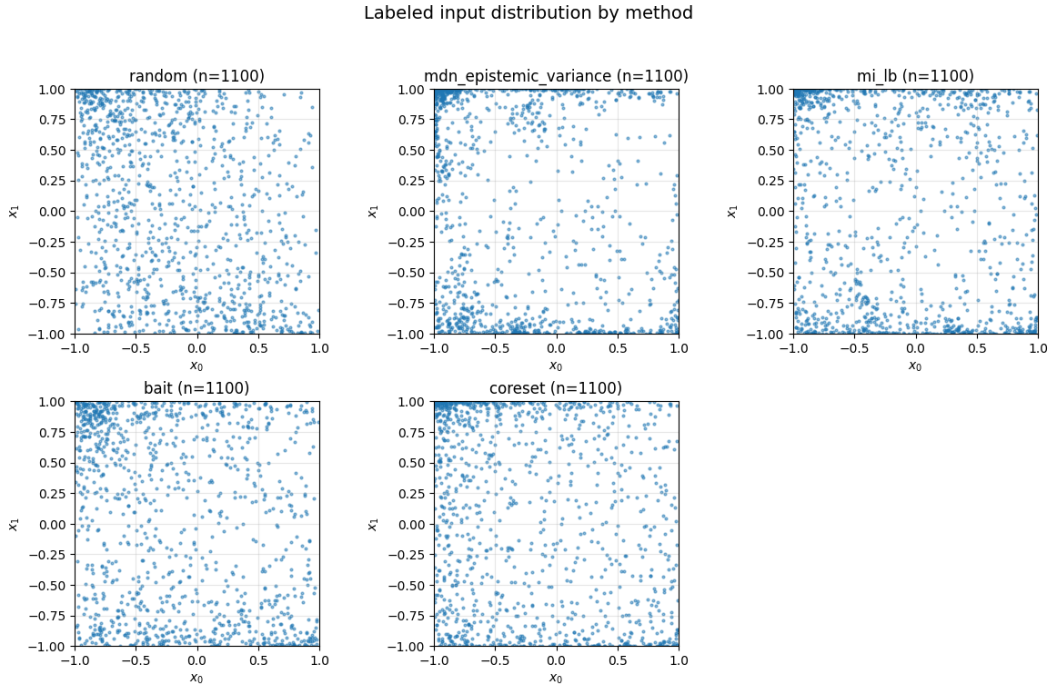


Figure 6: Spatial distribution of all labeled inputs (initial + acquired) at  $n = 1100$  projected onto  $(x_0, x_1)$  for the five base acquisitions: Random, Epistemic Variance, and MI-LB (**top row**), and the BAIT and Core-Set baselines (**bottom row**). By the end of the budget the four informative acquisitions look broadly similar — all concentrate mass along the multimodal edge of the pool ( $x_0 \rightarrow -1$ ) and the upper/lower boundary — while Random remains the clear outlier, spreading queries uniformly across the support. All five panels inherit the same off-center pool support from the fixed bias  $b$  in  $x = \tanh(AI + b)$  (drawn once at benchmark instantiation; see *Input manifold* above).

757 interior particles to both neighbours. Barrier heights are uniform,  $a = 1$ , giving per-particle potential  
 758  $V(q) = q^4/4 - q^2/2$  with minima at  $q = \pm 1$  and barrier height  $a/4 = 0.25$ . We integrate with the  
 759 Euler–Maruyama scheme

$$q_i^{(t+dt)} = q_i^{(t)} + [a(q_i^{(t)} - (q_i^{(t)})^3) + \kappa \sum_{j \in \text{nn}(i)} (q_j^{(t)} - q_i^{(t)})] dt + \sigma \sqrt{dt} \eta_i^{(t)}, \quad \eta_i^{(t)} \sim \mathcal{N}(0, 1),$$

760 with  $dt = 0.005$  and terminal time  $T = 5.0$  (i.e.  $n_{\text{steps}} = 1000$ ).

761 **Inputs and outputs.** Each trajectory is summarised by the 7-dimensional input  $x =$   
 762  $(q_1(0), \dots, q_5(0), \sigma, \kappa) \in \mathbb{R}^{P+2}$ , drawn component-wise from

$$q_i(0) \sim \text{Uniform}(-1.5, 1.5), \quad \sigma \sim \text{Uniform}(0.3, 2.0), \quad \kappa \sim \text{Uniform}(0, 3.0).$$

763 The output is the concatenation of  $n_{\text{snap}} = 4$  snapshots of the particle positions evenly spaced in  
 764  $[T/n_{\text{snap}}, T]$  (i.e. at  $t = 1.25, 2.5, 3.75, 5.0$ ), flattened into  $y \in \mathbb{R}^{n_{\text{snap}} \cdot P} = \mathbb{R}^{20}$ . Multi-snapshot  
 765 outputs give the model direct access to trajectory-level structure and substantially improve training  
 766 stability compared to single-endpoint outputs, at no additional simulation cost.

767 **Hyperparameters.** Table 6 lists the simulator and data hyperparameters and Table 7 the model and  
 768 active-learning hyperparameters.

Table 6: Simulator and data hyperparameters for the coupled double-well benchmark.

Symbol / option	Value	Meaning
$P$	5	Number of coupled particles
$a$	1 (uniform)	Barrier-height coefficient
$T$	5.0	Terminal integration time
$dt$	0.005	Euler–Maruyama step (1000 steps)
$n_{\text{snap}}$	4	Snapshots recorded per trajectory
boundary	open	Chain ends couple only to one neighbour
$q_i(0)$ range	$[-1.5, 1.5]$	Initial-configuration prior
$\sigma$ range	$[0.3, 2.0]$	Noise-intensity prior
$\kappa$ range	$[0, 3.0]$	Coupling-strength prior

769 **Model architecture.** Each ensemble member is an MDN with a 3-hidden-layer MLP backbone,  
 770 128 hidden units per layer, GELU activations, and a final linear projection to  $K_{\text{MDN}} = 8$  mixture  
 771 components with diagonal covariances over the 20-dimensional output. The ensemble size is  $n_{\text{ens}} = 8$ ,  
 772 with members initialised from different PRNG seeds and trained independently on the same labeled  
 773 set.

774 **Training schedule.** All members are trained with AdamW using the same optimiser configuration  
 775 as in Appendix F.2: peak learning rate  $5 \times 10^{-4}$ , weight decay  $10^{-2}$ , linear warmup followed by  
 776 exponential decay at rate 0.9 per 2,000 steps, and adaptive gradient clipping at 0.1. The number of  
 777 gradient steps per round is set adaptively to  $n_{\text{iter}}(n_{\text{lab}}) = \min(10,000, 10 \cdot n_{\text{lab}})$ , with mini-batch  
 778 size 128.

779 **Active-learning protocol.** For each of 5 seeds we instantiate a fresh pool of 50,000 candidate  
 780 inputs, a held-out test set of 2,000 inputs, and an initial labeled set of 100 inputs. Acquisition scores  
 781 are evaluated on the remaining pool in chunks of 256. Each run performs 20 rounds of 50 queries  
 782 each, so the final labeled set contains  $100 + 20 \cdot 50 = 1,100$  trajectories (2.2% of the pool). Test  
 783 NLL is evaluated at the end of every round. Seeds  $\{0, 1, 2, 3, 4\}$  are used for all reported curves.

784 **Acquisition functions and selection strategies.** Five base acquisitions, defined in Appendix F.2  
 785 unchanged: Random, Variance, MI-LB, BAIT [13], Core-Set [12]. Selection strategies: top- $k$   
 786 Random / Variance / MI-LB; SBAL ( $T = 1.0$ ) for Variance and MI-LB; MaxDist ( $w = 1$ ) for MI-LB.  
 787 BAIT and Core-Set bypass the score-to-batch step entirely.

Table 7: Model and active-learning hyperparameters for the coupled double-well benchmark.

Symbol / option	Value	Meaning
$n_{\text{ens}}$	8	Ensemble size
$K_{\text{MDN}}$	8	MDN mixture components per member
hidden features	128	MLP width
depth	3	Number of MLP hidden layers
activation	GELU	Nonlinearity
training batch size	128	Mini-batch size for AdamW
peak learning rate	$5 \times 10^{-4}$	AdamW peak LR (after warmup)
weight decay	$10^{-2}$	AdamW weight decay
gradient clip	0.1	Adaptive gradient clipping threshold
$n_{\text{iter}}$ (cap)	10,000	Max gradient steps per round
iter-per-sample	10	Slope of adaptive schedule
candidate pool size	50,000	Unlabeled pool $ \mathcal{X}_{\text{pool}} $
test set size	2,000	Held-out evaluation set
initial labelled	100	Initial labeled budget
AL rounds	20	Number of acquisition rounds
query batch size	50	Queries acquired per round
acquisition batch size	256	Chunk size for scoring the pool
seeds	$\{0, 1, 2, 3, 4\}$	Data / training seeds

788 **Results for SBAL and MaxDist variants.** Figure 7 and Table 8 report the full learning curves and  
789 the final test NLL at  $n = 1100$  across 5 seeds. Two patterns are worth highlighting. First, MaxDist  
790 pairs well with MI-LB: MI-LB (MaxDist) at  $w = 1$  finishes at  $131 \pm 15$ , within a factor of two of  
791 MI-LB (top- $k$ ) and comparable to Epistemic Variance (top- $k$ ). Acquisition-weighted farthest-point  
792 sampling adds batch diversity without discarding the informative ranking. Second, SBAL severely  
793 degrades both base scores on this benchmark. MI-LB (SBAL) ends at  $553 \pm 92$  and Epistemic  
794 Variance (SBAL) at  $457 \pm 46$ , both within striking distance of the Random baseline ( $518 \pm 60$ ) and  
795 roughly 4–8 $\times$  worse than their respective top- $k$  counterparts (122 and 71). This is the opposite of  
796 the multimodal-conditional result in Appendix F.2, where SBAL was competitive or helpful, and  
797 it reflects a quantitative rather than a qualitative difference: the total budget (1,100 points) is small  
798 relative to the pool, so softmax sampling at  $T = 1.0$  spreads enough mass onto low-score points  
799 that the effective fraction of informative queries collapses toward the Random baseline; annealing  $T$   
800 toward zero would tighten the distribution onto the top- $k$  result and recover it in the limit  $T \rightarrow 0$ . This  
801 is consistent with the original observations of Kirsch et al. [14]: stochastic batch acquisition helps  
802 when the top of the score ranking is corrupted by approximation noise and diverse sampling hedges  
803 against that noise; it hurts when the top ranking is already approximately correct and exploration  
804 trades informative points for uninformative ones.

805 **BAIT and Core-Set.** Both geometric baselines fail decisively on this benchmark: Core-Set finishes  
806 at  $304.3 \pm 54.8$  and BAIT at  $281.3 \pm 15.1$ ,  $\sim 4\times$  worse than MI-LB ( $70.8 \pm 8.0$ ) and within  $1.7\times$   
807 of Random ( $518.1 \pm 60.1$ ). Mode structure here lives in *output* space (Kramers escape between  
808  $q = \pm 1$ ): k-Center-Greedy spreads queries uniformly over  $(\sigma, \kappa)$  regardless of where the bimodal  
809 regime sits, and the single-MC last-layer Fisher embedding is dominated by within-mode noise.  
810 Their joint collapse is the strongest evidence in the paper that MI-LB’s advantage is not subsumed by  
811 feature-space coverage or last-layer information geometry.

812 **Why a mixture head:  $K = 1$  vs  $K = 8$ .** To justify the choice of a mixture head ( $K_{\text{MDN}} = 8$ )  
813 over a single-Gaussian head ( $K = 1$ ), we train both with the identical architecture used in the  
814 AL experiments (depth-3 MLP, 128 hidden units,  $n_{\text{ens}} = 8$ ) on a larger offline dataset of 200,100  
815 trajectories drawn from the same input distribution as the active-learning experiments, and evaluate  
816 both on a held-out test set of 5,000 trajectories. The ensemble-averaged test NLL is **15.90** for  $K = 8$   
817 versus **21.41** for  $K = 1$ , a gap of 5.51 nats on the 20-dimensional output (roughly 0.28 nats per  
818 output dimension). Absolute NLL values here are below those in the AL experiments because training  
819 uses 200,100 examples rather than the 1,100 queried by the AL loop; the quantity of interest is the  
820  $K = 8$ -vs- $K = 1$  gap, which is driven by the head structure, not the training-set size. The  $K = 1$

Table 8: Final test NLL at  $n = 1100$ , mean  $\pm$  std across 5 seeds (min–max in brackets) for the coupled double-well benchmark. All SBAL runs use  $T = 1.0$ ; MaxDist uses score weight  $w = 1$ .

Acquisition	Mean $\pm$ std	Min – Max
Random	518.1 $\pm$ 60.1	462.6 – 616.2
Epistemic Variance (top- $k$ )	122.1 $\pm$ 11.3	109.1 – 138.5
Epistemic Variance (SBAL)	456.8 $\pm$ 45.7	392.9 – 515.3
MI-LB (top- $k$ )	<b>70.8 <math>\pm</math> 8.0</b>	59.9 – 81.6
MI-LB (SBAL)	553.0 $\pm$ 91.9	423.5 – 630.2
MI-LB (MaxDist)	130.7 $\pm$ 14.9	113.5 – 147.0
BAIT	281.3 $\pm$ 15.1	268.8 – 305.3
Core-Set	304.3 $\pm$ 54.8	249.3 – 377.5

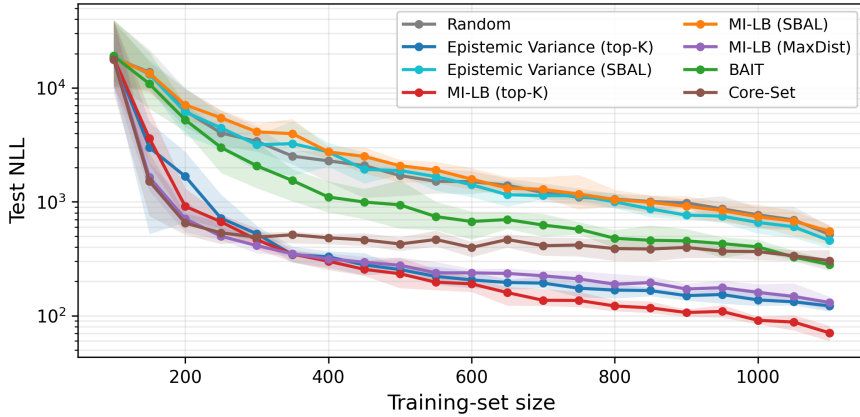


Figure 7: Learning curves for all eight (acquisition, selection-strategy) combinations on the coupled double-well benchmark including the BAIT and Core-Set baselines (5 seeds; bands min–max). MaxDist-MI-LB stays within  $\sim 2\times$  of top- $k$  MI-LB; both SBAL variants and the two geometric baselines (BAIT, Core-Set) collapse toward the Random regime.

821 model must place a single Gaussian per input, so wherever the conditional is bimodal it is forced to  
 822 smear mass between the wells and pay an unavoidable log-likelihood penalty; the  $K = 8$  mixture  
 823 can allocate one component per well. Figure 8 shows the qualitative picture on the test distribution:  
 824 ground-truth samples of  $(q_0(T), q_1(T))$  concentrate near  $(\pm 1, \pm 1)$  in the bimodal regime and the  
 825  $K = 8$  MDN recovers the same multi-well envelope, while the  $K = 1$  model collapses to a single  
 826 broad ellipse near the origin. This is the failure mode anticipated in Section 4.2 and the quantitative  
 827 reason the active-learning experiments on this benchmark use  $K_{\text{MDN}} = 8$  throughout.

#### 828 F.4 Material Science Application: Alloy Phase Competition

829 This appendix details the synthetic phase-competition simulator and the active-learning protocol  
 830 summarised in Section 4.3, and reports the SBAL and MaxDist batch-diversity variants deferred from  
 831 the main text.

832 **Phase model.** The simulator implements a softmin-over-quadratic-free-energies model inspired by  
 833 the CALPHAD framework [23]. For each of  $N_\phi = 4$  nominal phases  $\phi$  we draw a positive-definite  
 834  $3 \times 3$  Hessian as  $H_\phi = R_\phi R_\phi^\top + 0.3 \cdot I_3$  with  $R_\phi$  a Gaussian matrix ( $R_{\phi,ij} \sim \mathcal{N}(0, 0.5^2)$ ), and a linear  
 835 bias  $b_\phi \in \mathbb{R}^3$  with  $b_{\phi,i} \sim \mathcal{N}(0, 2^2)$ . Writing  $x_3 = (x_A, x_B, x_C)$  for the composition coordinates  
 836 extended onto the 3-simplex, the Gibbs free energy of phase  $\phi$  is the quadratic form

$$G_\phi(x_3) = \frac{1}{2} x_3^\top H_\phi x_3 + b_\phi^\top x_3,$$

837 and the phase posterior at temperature  $\tau_G$  is  $\pi(\phi | x_3) \propto \exp(-G_\phi(x_3)/\tau_G)$ . This posterior does  
 838 not depend on the process parameters  $p \in \mathbb{R}^{n_{\text{proc}}}$ ; phase assignment is a function of composition  
 839 alone.

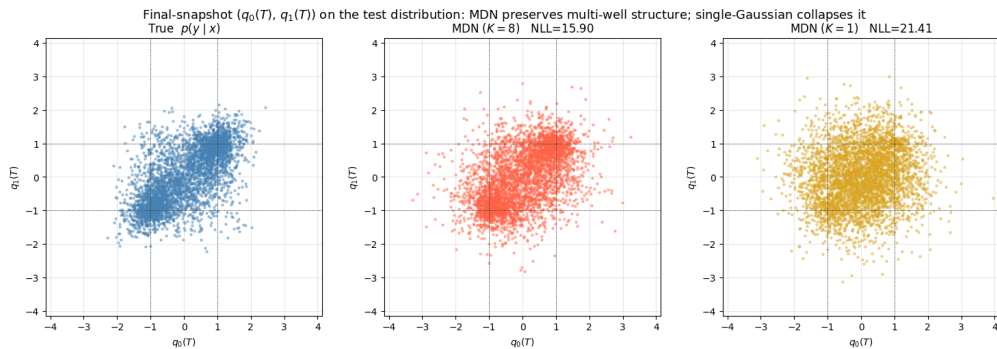


Figure 8: Coupled double-well benchmark: final-snapshot positions  $(q_0(T), q_1(T))$  for 4,000 held-out test inputs drawn from the same joint distribution as the AL test set. **Left:** ground-truth simulator samples; mass concentrates near  $(\pm 1, \pm 1)$ , the signature of the per-particle bimodal regime at  $\sigma \gtrsim \sqrt{a/2} \approx 0.71$ , with a diagonal enhancement that reflects neighbour alignment under non-zero  $\kappa$ . **Middle:** samples from the  $K = 8$  MDN ensemble (net 0); the mixture head preserves the same multi-well envelope (test NLL 15.90). **Right:** samples from a  $K = 1$  single-Gaussian MDN trained with the same architecture and budget; per-input unimodality forces mass into a single broad ellipse centred near the origin (test NLL 21.41; +5.51 nats over  $K = 8$ ). Dotted lines mark the well locations  $q = \pm 1$ .

840 The conditional response mean for phase  $\phi$  is

$$\mu_\phi(x_3, p) = c_\phi^\top x_3 + d_\phi + \frac{1}{2} \sin(\omega_\phi^\top x_3) + W_\phi^\top p,$$

841 with  $c_\phi \in \mathbb{R}^3$ ,  $d_\phi \in \mathbb{R}$ ,  $\omega_\phi \in \mathbb{R}^3$ , and  $W_\phi \in \mathbb{R}^{n_{\text{proc}}}$  all drawn from independent Gaussians (entry-  
 842 wise scales  $c_{\mu\text{-scale}} = 6$  for  $c_\phi$ , 2 for  $d_\phi$  and  $\omega_\phi$ , 1 for  $W_\phi$ ). The per-phase log-variance is the affine  
 843 function  $\log \sigma_\phi^2(x_3) = e_\phi^\top x_3 + f_\phi$  with  $e_{\phi,i} \sim \mathcal{N}(0, 0.5^2)$  and  $f_\phi \sim \mathcal{N}(-1, 0.3^2)$ , so per-phase  
 844 noise scales vary smoothly with composition around a typical  $\sigma_\phi \approx 0.6$ .

845 The output is sampled by first drawing  $\phi \sim \pi(\cdot | x_3)$ , then drawing  $y | \phi, x_3, p \sim$   
 846  $\mathcal{N}(\mu_\phi(x_3, p), \sigma_\phi^2(x_3))$ ; the ground-truth conditional  $p^*(y | x)$  is therefore a Gaussian mixture  
 847 in  $y$  with composition-dependent mixing weights, composition- and-process-dependent component  
 848 means, and composition-dependent component variances. The number of components with apprecia-  
 849 ble weight at any given  $x$  is determined by the realised phase structure of the chosen system seed,  
 850 summarised below.

851 **Inputs and outputs.** The input is the 8-dimensional vector  $x = (x_A, x_B, p_1, \dots, p_6)$  with  
 852  $(x_A, x_B, x_C)$  drawn from a symmetric Dirichlet( $\alpha = (1, 1, 1)$ ) on the 3-simplex and  $p_i \sim$   
 853 Uniform( $-1, 1$ ) component-wise. The output is the scalar response  $y \in \mathbb{R}$ .

854 **Realised phase structure at the system seed.** At the system seed used throughout  
 855 (system\_seed = 12,  $N_\phi = 4$ ,  $\tau_G = 0.08$ ), all four nominal phases carry appreciable mass on  
 856 the simplex: a grid-evaluation of  $\pi(\phi | x_3)$  on 20,301 uniform simplex points gives marginal phase  
 857 masses 48.18%, 25.60%, 18.84%, 7.38% (sorted, summing to 100%). The boundary region defined  
 858 by  $\max_\phi \pi(\phi | x_3) < 0.7$  occupies 11.40% of the simplex slice; tightening to  $\max_\phi \pi < 0.5$  con-  
 859 tracts it to 1.14%, and relaxing to  $\max_\phi \pi < 0.8$  expands it to 17.60%. The realised problem is thus  
 860 a genuine 4-phase competition with curved boundaries between phase domains in the 8-dimensional  
 861 input space. Inside any single phase domain the property response collapses to a single Gaussian  
 862 whose mean varies smoothly with the 6-dimensional process subspace, while in the boundary region  
 863 the conditional is a genuine multi-component mixture.

864 **Model architecture.** Each ensemble member is an MDN with a 2-hidden-layer MLP backbone, 64  
 865 hidden units per layer, GELU activations [31], and a final linear projection to  $K_{\text{MDN}} = 4$  mixture  
 866 components with diagonal covariances over the scalar property output. The ensemble size is  $n_{\text{ens}} = 8$ ,  
 867 with members initialised from different PRNG seeds and trained independently on the same labeled  
 868 set.

869 **Training schedule.** All members are trained with AdamW (peak learning rate  $2 \times 10^{-4}$ , weight de-  
870 cay  $5 \times 10^{-2}$ ) preceded by adaptive gradient clipping at 0.1, with the same warmup-then-exponential-  
871 decay schedule as in Appendix F.2. The number of gradient steps per round is set adaptively to  
872  $n_{\text{iter}}(n_{\text{lab}}) = \min(40,000, 200 \cdot n_{\text{lab}})$  with a minimum of 2,000 steps per round, and mini-batch  
873 size 64. The smaller peak LR and stronger weight decay (relative to the multimodal and double-well  
874 benchmarks) reflect the larger per-round step budget; as in the other benchmarks, optimiser, architec-  
875 ture, and active-learning hyperparameters are fixed across all acquisition functions and are not tuned  
876 per method.

877 **Hyperparameters.** Tables 9 and 10 list the simulator/data and model/AL hyperparameters used for  
878 all 30 runs reported below.

Table 9: Simulator and data hyperparameters for the ternary phase-competition benchmark.

Symbol / option	Value	Meaning
$N_\phi$	4	Nominal phases (all four active at the chosen seed)
$\tau_G$	0.08	Free-energy softmin temperature
$n_{\text{proc}}$	6	Process parameters per input
$c_{\mu\text{-scale}}$	6	Scale of per-phase composition coupling
$H_\phi$ raw scale	0.5	$R_{\phi,ij} \sim \mathcal{N}(0, 0.5^2)$
$H_\phi$ regulariser	$0.3 \cdot I_3$	PD floor on Hessian
$b_\phi$ scale	2	Linear-bias scale on composition
$d_\phi$ scale	2	Per-phase mean offset scale
$\omega_\phi$ scale	2	Sinusoidal-modulation frequency scale
$W_\phi$ scale	1	Per-phase process-coupling scale
$e_\phi$ scale	0.5	Composition-dependence of $\log \sigma_\phi^2$
$f_\phi$ mean / std	-1 / 0.3	Per-phase log-variance offset
$(x_A, x_B, x_C)$ prior	Dirichlet(1, 1, 1)	Composition prior on 3-simplex
$p_i$ prior	Uniform(-1, 1)	Process-parameter prior
system_seed	12	PRNG seed for $H_\phi, b_\phi, c_\phi, \dots$
candidate pool size	50,000	Unlabeled pool $ \mathcal{X}_{\text{pool}} $
test set size	2,000	Held-out evaluation set
initial labelled	100	Initial labeled budget

Table 10: Model and active-learning hyperparameters for the ternary phase-competition benchmark.

Symbol / option	Value	Meaning
$n_{\text{ens}}$	8	Ensemble size
$K_{\text{MDN}}$	4	MDN mixture components per member
hidden features	64	MLP width
depth	2	Number of MLP hidden layers
activation	GELU	Nonlinearity [31]
training batch size	64	Mini-batch size for AdamW
peak learning rate	$2 \times 10^{-4}$	AdamW peak LR (after warmup)
weight decay	$5 \times 10^{-2}$	AdamW weight decay
gradient clip	0.1	Adaptive gradient clipping threshold
min iter per round	2,000	Floor on adaptive iteration count
$n_{\text{iter}}$ (cap)	40,000	Max gradient steps per round
iter-per-sample	200	Slope of adaptive schedule
AL rounds	30	Number of acquisition rounds
query batch size	15	Queries acquired per round
acquisition batch size	256	Chunk size for scoring the pool
seeds	{0, 1, 2, 3, 4}	Data / training seeds

879 **Active-learning protocol.** For each of 5 seeds we instantiate a fresh pool of 50,000 candidate  
880 inputs, a held-out test set of 2,000 inputs, and an initial labeled set of 100 inputs. Acquisition scores

881 are evaluated on the remaining pool in chunks of 256. Each run performs 30 rounds of 15 queries  
 882 each, so the final labeled set contains  $100 + 30 \cdot 15 = 550$  inputs (1.1% of the pool). Test NLL is  
 883 evaluated at the end of every round.

884 **Acquisition functions and selection strategies.** We evaluate the same five base acquisitions as  
 885 in Appendix F.2 (Random, Epistemic Variance, MI-LB, BAIT [13], Core-Set [12]), paired with  
 886 three selection strategies for the score-based ones: top- $k$  for Random / Variance / MI-LB, SBAL for  
 887 Variance and MI-LB, and MaxDist for MI-LB. SBAL temperature is  $T = 1.0$ ; MaxDist score weight  
 888 is  $w = 1$ . BAIT and Core-Set bypass the score-to-batch step entirely (Fisher-trace forward+backward  
 889 greedy and k-Center-Greedy on the MDN backbone respectively). Definitions follow Appendix F.2  
 890 unchanged.

891 **Seed-to-seed spread of the top- $k$  base acquisitions.** Section 4.3 reports that MI-LB is markedly  
 892 more consistent across seeds than Random or Epistemic Variance in the data-scarce regime, with the  
 893 gap closing by  $n \gtrsim 400$ . Table 11 gives the underlying numbers: the seed-to-seed min–max range  
 894 of test NLL averaged over each budget band. In the early band  $n \in [115, 250]$ , MI-LB’s range is  
 895  $\sim 2.6\times$  tighter than Random’s and  $\sim 4.5\times$  tighter than Epistemic Variance’s; the inflated Variance  
 896 spread comes from seeds in which the second-moment signal selects unhelpful queries before the  
 897 ensemble has accumulated enough data to disagree informatively over modal structure. By the late  
 898 band ( $n \gtrsim 400$ ) all three ranges sit within a factor of  $\sim 2.4$  of each other.

Table 11: Seed-to-seed min–max range of test NLL (in nats), averaged over each budget band, for the top- $k$  base acquisitions on the synthetic phase-competition benchmark (`system_seed = 12`, 5 seeds).

Acquisition (top- $k$ )	Early ( $n \in [115, 250]$ )	Late ( $n \gtrsim 400$ )
Random	1.82	0.24
Epistemic Variance	3.13	0.31
MI-LB	<b>0.69</b>	<b>0.13</b>

899 **Results for SBAL and MaxDist variants.** Table 12 reports the final test NLL at  $n = 550$  across 5  
 900 seeds for all combinations, and Figure 9 shows the corresponding learning curves. Three observations.

901 First, MaxDist is the mean-best acquisition on this benchmark: MI-LB (MaxDist) at  $w = 1$  finishes  
 902 at  $1.945 \pm 0.061$ , narrowly ahead of MI-LB (top- $k$ ) at  $1.985 \pm 0.043$ . The seed-to-seed spread is  
 903 large enough that the two are not separated cleanly—per-seed, MaxDist wins on 2 seeds, ties on 1,  
 904 and is within 0.01 nat of top- $k$  on the other 2—so the safest summary is that MaxDist matches MI-LB  
 905 (top- $k$ ) and offers a small mean improvement, the same qualitative conclusion as on the coupled  
 906 double-well benchmark of Appendix F.3.

907 Second, SBAL hurts MI-LB: MI-LB (SBAL) at  $T = 1.0$  degrades from  $1.985 \pm 0.043$  to  $2.109 \pm$   
 908  $0.054$ , with MI-LB (top- $k$ ) winning on all 5 seeds. This is consistent with the small-budget, narrow-  
 909 high-information-manifold argument of Appendix F.3: at  $T = 1.0$  enough mass leaks onto low-score  
 910 points that the effective fraction of informative queries drops, and the realisable improvement from  
 911 batch diversity is bounded above by the MaxDist result that preserves the informative ranking exactly.

912 Third, SBAL modestly helps Epistemic Variance:  **$2.029 \pm 0.050$**  for SBAL versus  $2.062 \pm 0.087$   
 913 for top- $k$ , with the SBAL variant winning on 4 of 5 seeds and producing a tighter seed-to-seed  
 914 spread. Because the Variance-based score is itself already a noisier ranking signal than MI-LB on  
 915 this benchmark, its top- $k$  ordering is less informative, and stochastic relaxation does not cost what it  
 916 does for MI-LB, at this temperature it actively helps by hedging against unhelpful early queries.

917 **BAIT and Core-Set.** Core-Set ties MI-LB within seed-to-seed noise ( $1.989 \pm 0.043$  vs  $1.985 \pm$   
 918  $0.043$ ); the tie reflects benchmark geometry, with phase boundaries lying directly on the 2-D composi-  
 919 tion simplex so that k-Center-Greedy in feature space hits the same boundary MI-LB targets through  
 920 entropy disagreement. BAIT is the only acquisition whose mean NLL falls within seed-to-seed noise  
 921 of Random ( $2.190 \pm 0.090$  vs  $2.209 \pm 0.085$ ): the 0.019 improvement is an order of magnitude  
 922 smaller than the  $\sim 0.09$  standard deviation of either run, whereas every other acquisition in Table 12  
 923 beats Random by at least 0.10 in mean NLL — well outside the noise band. The single-MC last-layer  
 924 Fisher embedding produces a noisy candidate ranking on the 1-D scalar output, and the  $30 \times 15$   
 925 small-budget protocol does not give the forward+backward greedy enough rounds to recover — the

926 same failure mode anticipated in Appendix F.2 (where BAIT had the largest seed-to-seed spread of  
 927 any reported method) but more pronounced here because of the lower-dimensional output.

Table 12: Final test NLL at  $n = 550$ , mean  $\pm$  std across 5 seeds  $\{0, 1, 2, 3, 4\}$  (min–max in brackets) for the synthetic phase-competition benchmark at `system_seed = 12`. SBAL temperature is  $T = 1.0$  for both SBAL runs; MaxDist uses score weight  $w = 1$ .

Acquisition	Mean $\pm$ std	Min – Max
Random	$2.209 \pm 0.085$	2.136 – 2.346
Epistemic Variance (top- $k$ )	$2.062 \pm 0.087$	1.918 – 2.122
Epistemic Variance (SBAL)	$2.029 \pm 0.050$	1.977 – 2.107
MI-LB (top- $k$ )	$1.985 \pm 0.043$	1.911 – 2.017
MI-LB (SBAL)	$2.109 \pm 0.054$	2.057 – 2.174
MI-LB (MaxDist)	<b><math>1.945 \pm 0.061</math></b>	1.870 – 2.023
BAIT	$2.190 \pm 0.090$	2.087 – 2.265
Core-Set	$1.989 \pm 0.043$	1.951 – 2.058

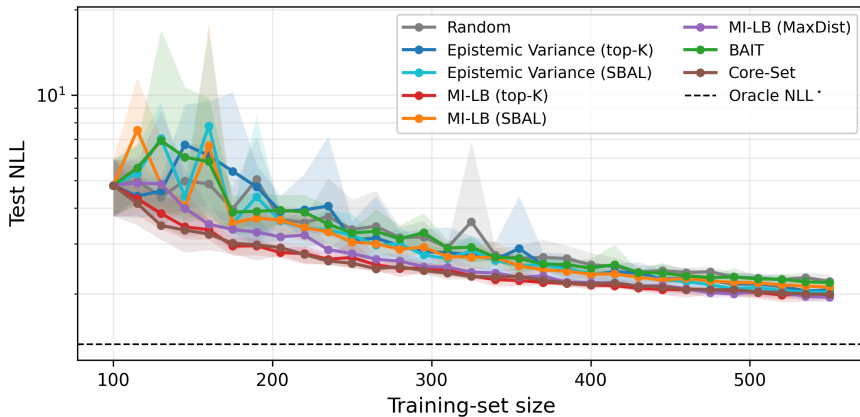


Figure 9: Synthetic phase-competition benchmark: test NLL vs. training-set size for all eight (acquisition, selection-strategy) combinations including the BAIT and Core-Set baselines (5 seeds; shaded bands min–max). SBAL  $T = 1.0$ , MaxDist  $w = 1$ . The dotted line marks the oracle floor  $NLL^* \approx 1.33$ . The  $y$ -axis is clipped to  $[1.2, 7.0]$  to suppress early-iteration training-instability spikes for the SBAL variants and BAIT.

928 **Remark on per-acquisition simplex visualisation.** A per-acquisition visualisation of selected  
 929 points projected onto the  $(x_A, x_B)$  simplex slice would complement the boundary-seeking narrative  
 930 of Section 4.3. The production runs reported in Table 12 did not log query indices to the experiment  
 931 tracker, so we do not include such a figure here; the example notebook accompanying this benchmark  
 932 provides a single-seed reproduction that materialises the simplex projection from local state.

933 **Why a mixture head:  $K = 1$  vs  $K = 4$ .** To justify the choice of a mixture head ( $K_{MDN} = 4$ )  
 934 over a single-Gaussian head ( $K = 1$ ), we train both with the identical architecture used in the  
 935 AL experiments (depth-2 MLP, 64 hidden units,  $n_{ens} = 8$ ) on a larger offline dataset of 100,000  
 936 samples drawn from the same input distribution as the active-learning experiments, and evaluate  
 937 both on a held-out test set of 10,000 samples. The ensemble-averaged test NLL is **1.34** for  $K = 4$   
 938 versus **1.88** for  $K = 1$ , against an analytically computed oracle floor of  $NLL^* = 1.33$ . The  $K = 4$   
 939 ensemble therefore sits within 0.01 nat of the oracle—consistent with the fact that the simulator is  
 940 itself a 4-component Gaussian mixture per input, so the  $K = 4$  MDN model class contains the true  
 941 distribution. The  $K = 1$  model pays a 0.55-nat penalty at the same budget: forced to place a single  
 942 Gaussian per input, it cannot resolve phase competition in the boundary region and must smear mass  
 943 across phase means that differ by several nats (per-phase median means  $+2.9, +6.9, +3.7, -4.0$  on  
 944 this simulator instance). The penalty is smaller than on the coupled double-well benchmark (5.51  
 945 nats, Appendix F.3) because the per-phase response is a unimodal Gaussian with bounded variance,  
 946 so the smearing cost is bounded; it is still clearly nonzero, and the active-learning experiments on this  
 947 benchmark therefore use  $K_{MDN} = 4$  throughout.

948 Figure 10 shows the qualitative picture on a 120-point simplex grid with process parameters pinned to  
 949 zero. The  $K = 4$  ensemble’s predicted conditional mean  $\hat{\mathbb{E}}[Y | x]$  tracks the ground-truth  $\mathbb{E}^*[Y | x]$   
 950 to within roughly  $\pm 1$  nat, with residuals smoothly distributed across the simplex. The  $K = 1$   
 951 residual is visibly structured and larger (up to  $\pm 4$ ) at the intersections of the phase-boundary network,  
 952 where the single-Gaussian head cannot simultaneously fit the correct mean and an appropriate  
 953 variance—exactly the failure mode anticipated in Section 4.3.

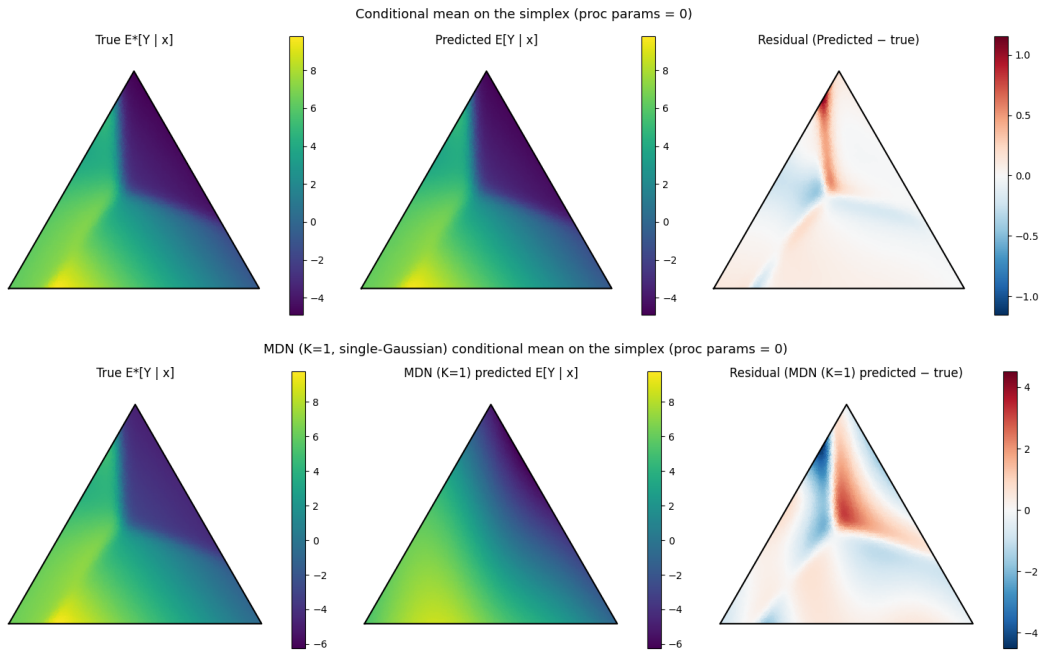


Figure 10: Synthetic phase-competition benchmark (`system_seed = 12`): predicted conditional mean on a 120-point simplex grid (process parameters pinned to zero) for the  $K = 4$  MDN ensemble (**top row**) and a  $K = 1$  single-Gaussian MDN (**bottom row**), both trained offline on 100,000 samples with the architecture used throughout the AL experiments. **Left column:** ground-truth  $\mathbb{E}^*[Y | x]$ . **Middle column:** ensemble-predicted  $\hat{\mathbb{E}}[Y | x]$ . **Right column:** residual (predicted  $-$  true). The  $K = 4$  residuals sit within  $\sim \pm 1$  and are smoothly distributed; the  $K = 1$  residuals reach  $\pm 4$  and concentrate at the intersections of the phase-boundary network, where the single-Gaussian head cannot represent the competing phase components. The ensemble-averaged test NLL on a held-out set of 10,000 samples is 1.34 ( $K = 4$ ) vs 1.88 ( $K = 1$ ), against an oracle floor of 1.33.

## 954 F.5 Compute Resources

955 All experiments were run on a workstation with  $7 \times$  NVIDIA RTX A6000 GPUs (48 GB VRAM  
 956 each), 128 CPU cores, and approximately 504 GiB of system memory. Each active-learning run  
 957 uses a single A6000 GPU; peak GPU memory fits comfortably within the 48 GB budget for all  
 958 configurations (small MLP backbones, ensemble size 8), so GPU memory is not a binding constraint.  
 959 Per benchmark we evaluate eight (acquisition, selection-strategy) combinations across 5 seeds, giving  
 960 40 runs per benchmark and 120 runs in total across the three benchmarks. Per-run wall-clock times,  
 961 extracted from the W&B logs and reported as the median over 5 seeds, are summarised in Table 13.  
 962 Total compute across the three benchmarks is approximately 25 GPU-hours, of which  $\sim 12$  h (roughly  
 963 half) is the 15 BAIT runs, BAIT is the only acquisition with non-trivial selection overhead, since its  
 964 forward+backward greedy on the per-input Fisher embeddings of the full 50,000-point pool scales  
 965 with both pool size and query batch. All other acquisitions, including Core-Set, complete in roughly  
 966 the same wall time as Random per benchmark.

Table 13: Per-run wall-clock time on a single NVIDIA RTX A6000, median over 5 seeds, for the three paper benchmarks. Each entry is the end-to-end runtime of one 20-round (multimodal, double-well) or 30-round (ternary) active-learning run. Total compute across all 120 runs is  $\sim 25$  GPU-hours, of which  $\sim 12$  h is BAIT ( $\sim 8.6$  h on the double-well benchmark alone).

Acquisition (selection)	Multimodal	Double-well	Ternary
Random (top- $k$ )	2.7 min	3.1 min	12.8 min
Epistemic Variance (top- $k$ )	3.2 min	3.4 min	13.6 min
MI-LB (top- $k$ )	3.1 min	3.4 min	13.1 min
BAIT (top- $k$ )	22.4 min	1.35 h	15.0 min
Core-Set (top- $k$ )	3.5 min	4.4 min	14.2 min
Epistemic Variance (SBAL)	3.1 min	3.4 min	13.2 min
MI-LB (SBAL)	3.1 min	3.4 min	13.2 min
MI-LB (MaxDist)	3.7 min	3.9 min	13.6 min